

## Molecular Phylogenetics of the Lizard Genus *Microlophus* (Squamata:Tropiduridae): Aligning and Retrieving Indel Signal from Nuclear Introns

EDGAR BENAVIDES,<sup>1</sup> REBECCA BAUM,<sup>3</sup> DAVID MCCLELLAN,<sup>1</sup> AND JACK W. SITES, JR.<sup>1,2</sup>

<sup>1</sup>Department of Integrative Biology, and

<sup>2</sup>M.L. Bean Life Science Museum, Brigham Young University, Provo, UT 84602, USA; E-mail: eb235@email.byu.edu (E.B.)

<sup>3</sup>Department of Chemistry, Brigham Young University, Provo, UT 84602, USA

**Abstract.**— We use a multigene data set (the mitochondrial locus and nine nuclear gene regions) to test phylogenetic relationships in the South American “lava lizards” (genus *Microlophus*) and describe a strategy for aligning noncoding sequences that accounts for differences in tempo and class of mutational events. We focus on seven nuclear introns that vary in size and frequency of multibase length mutations (i.e., indels) and present a manual alignment strategy that incorporates insertions and deletions (indels) for each intron. Our method is based on mechanistic explanations of intron evolution that does not require a guide tree. We also use a progressive alignment algorithm (Probabilistic Alignment Kit; PRANK) and distinguishes insertions from deletions and avoids the “gapcost” conundrum. We describe an approach to selecting a guide tree purged of ambiguously aligned regions and use this to refine PRANK performance. We show that although manual alignment is successful in finding repeat motifs and the most obvious indels, some regions can only be subjectively aligned, and there are limits to the size and complexity of a data matrix for which this approach can be taken. PRANK alignments identified more parsimony-informative indels while simultaneously increasing nucleotide identity in conserved sequence blocks flanking the indel regions. When comparing manual and PRANK with two widely used methods (CLUSTAL, MUSCLE) for the alignment of the most length-variable intron, only PRANK recovered a tree congruent at deeper nodes with the combined data tree inferred from all nuclear gene regions. We take this concordance as an objective function of alignment quality and present a strongly supported phylogenetic hypothesis for *Microlophus* relationships. From this hypothesis we show that (1) a coded indel data partition derived from the PRANK alignment contributed significantly to nodal support and (2) the indel data set permitted detection of significant conflict between mitochondrial and nuclear data partitions, which we hypothesize arose from secondary contact of distantly related taxa, followed by hybridization and mtDNA introgression. [Indels; length-variable introns; *Microlophus*; mitochondrial-nuclear conflict; phylogenetics; progressive alignment; secondary contact.]

Sequence homology statements are key to formulating sound phylogenetic hypotheses, and multiple alignment programs based on either progressive (Thompson et al., 1994 [CLUSTAL]; Edgar, 2004 [MUSCLE]) or consistency-based scoring alignments (Notredame et al., 2000 [T-Coffe]), are widely used by systematists. Improvements such as iterative optimizations to correct errors introduced early in the alignment process (Wallace et al., 2005) have greatly increased the accuracy and sensitivity of these methods (Thompson et al., 1999; Pollard et al., 2004; Wallace et al., 2005; Lunter et al., 2005a), but the majority of systematic studies commonly apply a single heuristic alignment algorithm to data sets with different mutation mechanisms and mutation patterns (i.e., ribosomal versus protein-coding sequences), implying that all gene regions can be aligned under the same set of assumptions. For example, the distribution of ribosomal gene substitutions is tightly constrained by secondary structure, whereas protein gene substitutions are constrained by reading frame and codon conservation (Li, 1997). In contrast, point mutations plus short and long length mutations characterize the evolution of intron regions (Belshaw and Bensasson, 2006), and the indiscriminate application of a single heuristic alignment method (e.g., CLUSTAL) to protein-coding, noncoding, and ribosomal sequences is an inappropriate oversimplification.

The increased use of novel nuclear gene regions in phylogenetic studies, alone or combined with more traditional markers such as ribosomal and mitochondrial gene regions, raises new theoretical and empirical issues

(e.g., parameter and tree optimization [Pagel and Meade, 2004], data partitioning [Castoe et al., 2004; Brandley et al., 2005], model choice [Sullivan and Joyce, 2005], and topological incongruence [Gatesy and Baker, 2005; Phillips et al., 2004]). Likewise, investigators have compared the performance of different alignment methods (Whiting et al., 2006; Kjer et al., 2007), but to our knowledge, a clear distinction of issues affecting alignment procedures across genes that do vary in tempo and class of mutational event (point versus small indels versus large length mutations) has not been presented. In this study, we empirically differentiate the alignment of (a) protein coding and (b) ribosomal gene regions from (c) nuclear introns characterized by length mutations of variable sizes and frequencies and compare phylogenetic hypotheses generated from alignments that make these distinctions with those from two widely used heuristic methods that do not. Our analyses benefit from recent insights for the global alignment of complex indel substitution patterns (Loytinoja and Goldman, 2005).

### *The South American “Lava Lizards” (Genus Microlophus)*

The “lava lizards” (genus *Microlophus*; Tropiduridae) display an unusual geographic distribution among terrestrial vertebrates; the 21 recognized species include nine taxa endemic to the Galápagos Archipelago and 12 species mostly confined to a linear strip of 5000 km along rain-shadowed western coastal deserts of South America

(Figs. 1 and 2). Monophyly of the genus is well supported (Frost, 1992; Harvey and Gutberlet, 2000; Frost et al., 2001), but the currently recognized *Occipitalis* (Dixon and Wright, 1975; Frost, 1992) and *Peruvianus* groups (Van Denburgh and Slevin, 1913; Frost, 1992; Heise, 1998) are less supported. The 12 species of the *Occipitalis* group include the aforementioned nine Galápagos endemics, which represent two independent colonization events from the continent (Lopez et al., 1992; Wright, 1983; Heise, 1998; Kizirian et al., 2004). The remaining three species of *Occipitalis* plus the nine species of the *Peruvianus* group are confined to the continent. This study is the first to include all recognized species and multiple localities of widespread species and thus improves on earlier studies by increasing these two critical variables of phylogenetic sampling design (Rokas et al., 2005). Our objectives are to (1) test a novel alignment strategy to the specific problem of alignment of length-variable mutations typical of nuclear introns and (2) obtain a well-supported phylogeny for the genus *Microlophus*, as a foundation for ongoing studies of interisland colonization patterns in the Galápagos Archipelago, and patterns

of speciation within the *Peruvianus* group in mainland South America.

## MATERIALS AND METHODS

### Taxon Sampling

This study includes 73 terminals, of which 70 are ingroup samples from one or more localities for each of the 21 species of *Microlophus*; online Appendix 1 (<http://www.systematicbiology.org>) summarizes locality information for each terminal and provides details on museum vouchers. All localities are shown in Figures 1 and 2.

### Sampling and Laboratory Methods

Genomic DNA was extracted from muscle tissue preserved in 100% ethanol using either a slightly modified version of the procedure of Fetzner (1999) or the Qiagen extraction kit (Qiagen, Valencia, CA). Amplifications were performed (under varying profiles) in 20- $\mu$ L reaction volumes using TaKaRa hotstart *Taq*

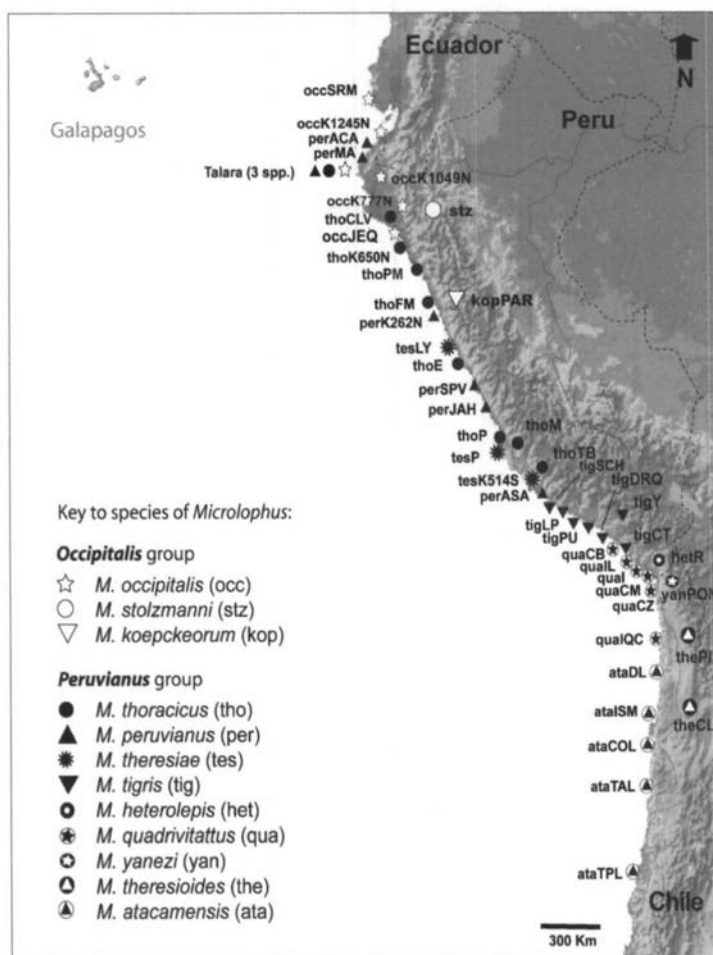


FIGURE 1. South American continental localities sampled for the 12 recognized species of mainland *Microlophus* used in this study. Locality and voucher details are summarized in Appendix 1; open and solid symbols identify species of the *Occipitalis* and *Peruvianus* groups (species acronyms in parentheses).

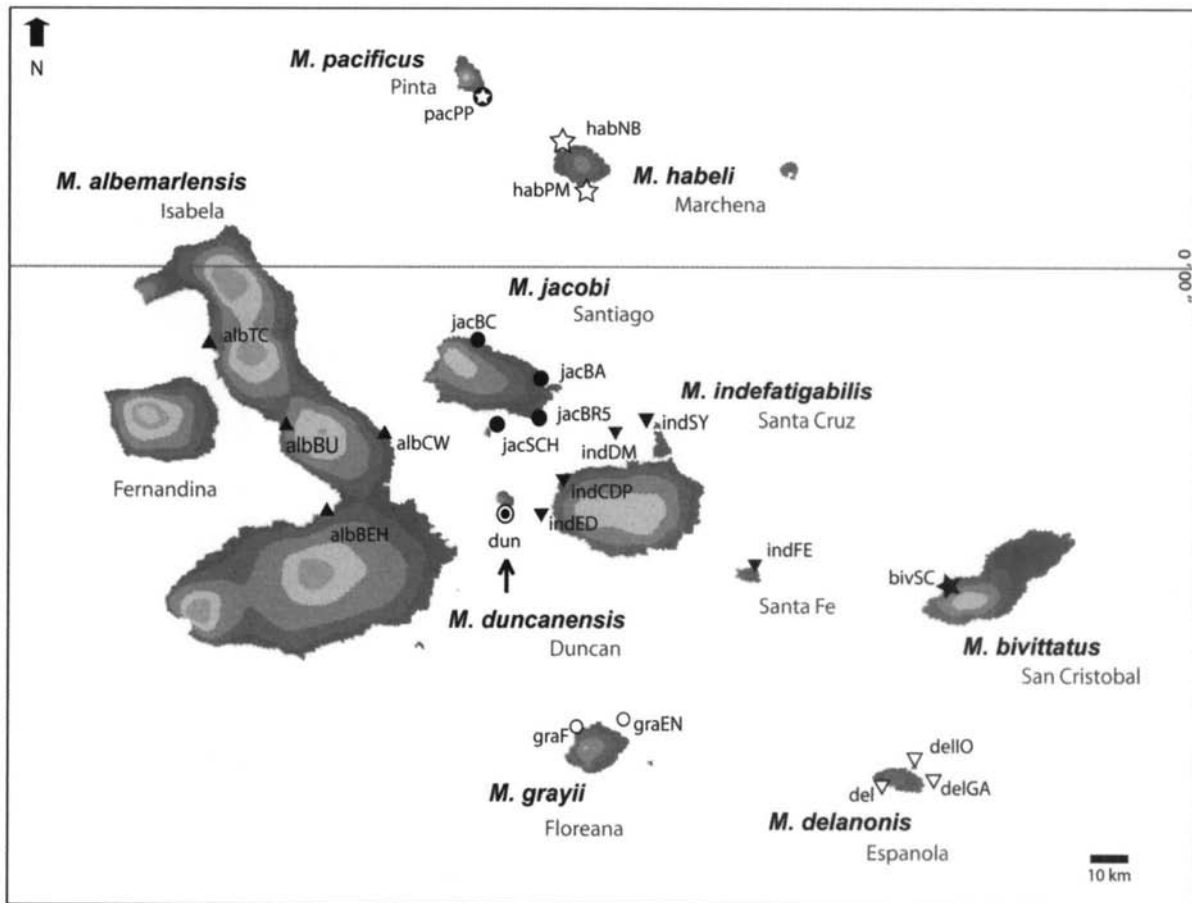


FIGURE 2. Localities sampled for the nine species of *Microlophus* endemic to the Galápagos Archipelago. Locality details are summarized in Appendix 1; solid and open stars identify the two species of the Eastern Galápagos Radiation (*M. bivittatus* and *M. habeli*).

DNA polymerase and 10× reaction buffer (100 mM Tris-HCl [pH 8.3], 500 mM KCl, 15 mM MgCl<sub>2</sub>). Five mitochondrial (cyt-*b*, ND4, tRNA<sup>his+ser</sup>, 12S, and 16S) and nine nuclear (recombination activating gene 1 [Rag-1], nuclear oocyte maturation factor gene [Cmos],  $\alpha$ -enolase [Enol], glyceraldehyde-phosphate dehydrogenase [Gapdh], creatin kinase [CK], ribosomal protein 40 [RP40],  $\beta$ -crystallin [Cryba],  $\alpha$ -tropomyosin [Atrp], and an anonymous [Anon] region) gene regions were sampled for all terminals. Among the nuclear regions, only the Rag-1 and Cmos genes are coding genes, whereas the remaining seven are noncoding introns. The primers and PCR conditions are summarized in online Appendix 2 (<http://www.systematicbiology.org>). A BLAST search did not match the Anon locus to any known gene, and it displayed no conflict when used in phylogenetic analyses. It was therefore presumed to contain phylogenetic signal and is included in all analyses (see Shaw [2002] and Dolman and Moritz [2006] for examples).

Purified double-stranded products were used directly in 1/4 volume dideoxy-termination sequencing reactions using BigDye Terminator v3.1 (Applied Biosystems). Both strands were sequenced for all PCR products, and sequences were edited and proofread with Sequencher v. 4.1 (Gene Codes). All sequenced

gene regions were queried in BLAST searches of GenBank to confirm homology. Complete sequences have been deposited in GenBank (online Appendix 3; <http://www.systematicbiology.org>), whereas aligned sequences have been stored in TreeBase (project no. SN3395).

#### Alignment of Protein-Coding and Ribosomal Gene Regions

Protein-coding regions (Cyt-*b* and ND4 [mtDNA] and Cmos and Rag-1 [nuclear]) were aligned using MUSCLE (Edgar, 2004) and open reading frames checked with SeAl (Rambaut, 1996). Ribosomal (12S, 16S, and tRNAs) gene regions were also aligned with MUSCLE under default parameters and then the "refine" command was used twice to further improve the existing alignment. Final mtRNA positional homology was derived by visual adjustment to secondary structure models developed by Van de Peer et al. (1994 [for the 12S gene]); Gutell and Fox (1988 [for the 16S gene], but see Wiens and Reeder, 1997; and Cannone et al., 2002), and Kumazawa and Nishida (1993 [for the tRNAs]). In some cases, single nucleotides were shifted by a few base pairs to conserve base identity with blocks inferred to be stem motifs (from the models), thus constraining small indels to appear in loop regions.

TABLE 1. Summary of sequence and indel variability in the seven length-variable nuclear introns inferred after manual (M) and PRANK (PR) alignments.

Gene region	No. of indels and parsimony-informative indels (in parentheses)		No. of characters after alignment		Sequence size range (bp)	No. of parsimony informative sites		No. of invariant sites	
	M	PR	M	PR		M	PR	M	PR
Enol	6(5)	6(5)	284	284	147–284	36	36	224	224
Gapdh	13(4)	16(6)	341	342	289–341	53	55	259	258
Atrp	10(10)	13(5)	258	260	237–248	42	41	197	203
Anon	12(10)	11(11)	415	415	360–414	97	97	309	310
RP40	14(7)	14(7)	806	806	327–765	87	87	665	665
Cryba	51(39)	56(38)	841	897	185–839	204	178	583	657
CK	11(9)	15(13)	358	363	310–355	84	76	248	254
Total	117(84)	131(85)	3303	3367	1855–3246	603	570	2485	2571

The alignment of nuclear indel regions incorporated information from both point and length mutations that characterize regions presumed to experience limited or no selective pressures (Thorne et al., 1992; but see Belshaw and Bensasson, 2006; Roy and Gilbert, 2006). All seven introns displayed length variation among ingroup terminals (Table 1), and we took two approaches to their alignment.

#### Manual Alignment of Introns

Manual alignments were implemented under two important assumptions. First, all indels represent single mutational events (Graham et al., 2000; Kelchner, 2000), and second, we define alignments as most parsimonious when indels are placed to preserve blocks of sequence integrity (i.e., indel-free regions with maximum base-pair identity) in the indel-flanking regions (Morgenstern, 1999; Britten et al., 2003; Brudno et al., 2003; Ogden and Rosenberg, 2006). The “single mutation event” assumption represents a fundamental difference from any algorithm-based optimization of positional homology, and we justify this assumption based on recent descriptions of mutational hot spots, secondary structure configurations, and repeat motifs of variable complexity (Graham et al., 2000; Kelchner, 2000). These studies suggest that some introns evolve under structural constraints in a nonrandom, nonindependent fashion in which length differences between sequences are best explained as single mutational events (Lohne and Borsch, 2005; Andolfatto, 2005). Therefore, we manually aligned the seven nuclear introns individually, on the basis of these two basic assumptions, using the “4-step” procedure outlined in Figure 3.

The upper panel of Figure 3 represents an unaligned intron (Cryba) for which length polymorphisms may be caused by both insertions and deletions. In order to manually align this intron, we first grouped conspecific ingroup terminals, one under the other, which easily permitted identification of blocks of sequence base-pair identity in some regions. In other parts of the sequence for which such blocks were not obvious among conspecifics, we shifted sequences to create indels having identical positional extension—indels of the

same length whose insertion into a sequence created additional blocks of nucleotide identity on either side of the indel. These are the entire indels described by Graham et al. (2000), and their placement reduced or eliminated the number of base substitutions in the flanking sequences (Fig. 3, step 1).

After completing step 1 for conspecific sequences, we compared heterospecific blocks side by side (= step 2; done independently for ingroup and outgroup terminals—see details in Fig. 3 legend), to identify either the same indels, or those for which positional extension differed slightly between different taxa (“overlapping” indels [Graham et al., 2000]). Complex arrangements were usually located after the first step and always inferred between heterospecific terminals, never between conspecifics. The iterative collapsing of identical alignments identified in steps 1 and 2 is analogous to the population aggregation analysis described by Davis and Nixon (1992), which is used to identify diagnostic character differences between species. At a coarse scale, different taxonomic units are evident as consistent blocks of nucleotide colored columns in the unaligned Cryba intron depicted in the upper panel of Figure 3.

The third step required a number of alternative options to deal with more complex rearrangements identified in step 2. At this step, indel placement was further improved (i.e., the overall alignment made more stringent with respect to nucleotide identity in conserved blocks) by comparing alignments made in steps 1 and 2 to specific intron microstructural changes first described by Golenberg et al. (1993) and Gu and Li (1995) and later employed by others (summarized in Lohne and Borsch [2005]) for the manual alignment of intron regions. The use of an a priori set of rules further reduces subjectivity of the alignment process and enhances its repeatability (Sanchis et al., 2001). Step 3 (“subroutines” and details are given in Fig. 3) emphasizes the identification of simple sequence repeat (SSR) motifs and/or possible inversions, as they are evidence of slipped-strand mispairing (Levinson and Gutman, 1987) or simple hairpin structures (Kelchner, 2000; Lohne and Borsch, 2005).

In the final step, we concatenated all length-variable aligned introns and the two protein-coding nuclear

regions (Cmos and Rag-1) to create the final nuclear data matrix (Fig. 3, step 4). Although time-consuming, our approach allows visualization of blocks of sequence integrity within and across multiple taxa and gives an alignment based on the premise of primary homology assessment (De Pinna, 1991). The matrix provides a reasonable comparative framework for input tree-based computer-generated alignments (Sanchis et al., 2001; Creer et al., 2006) and maximizes nucleotide identity in flanking gapless intron sequences (Ogden and Rosenberg, 2006; Siddharthan, 2006) as a proxy for nucleotide homology hypotheses.

#### *Progressive Alignment of Introns*

We employed a modification of the current progressive algorithms (Edgar, 2004; Notredame et al., 2000; Katoh et al., 2000; Keightley and Johnson, 2004; Do et al., 2005) as an alternative to the manual alignment. The Probabilistic Alignment Kit (PRANK; Loytynoja and Goldman, 2005) algorithm implements Markov models and probabilistic score schemes to handle multinucleotide indel events and distinguishes insertions from deletions, a step that is fundamental in the context of intron sequence alignment. In most other programs, indels are penalized relative to nucleotide changes, and arbitrarily chosen penalties might produce either highly fragmented sequences with multiple indels and few nucleotide differences (false negatives; Morrison 2006:512) or few indels coupled with many nucleotide differences (false positives or overalignment; Cline et al., 2002:309).

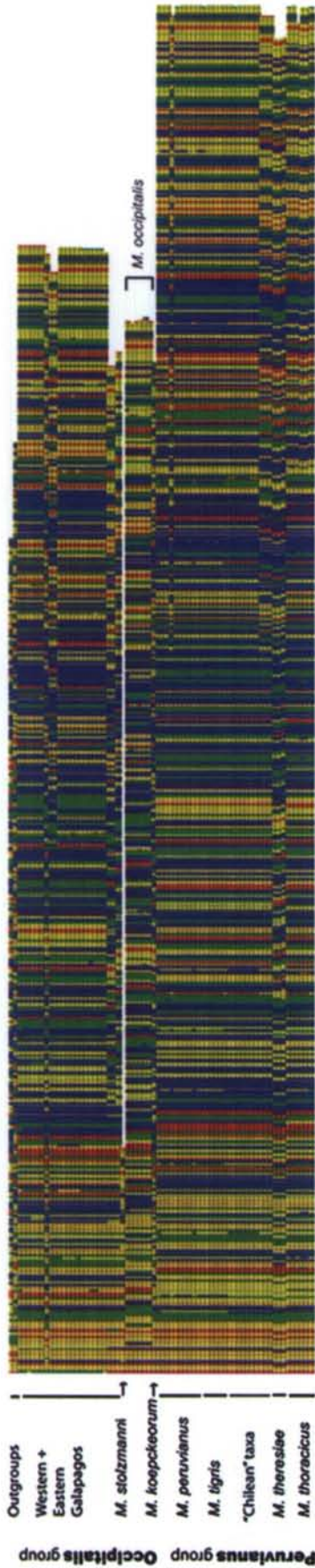
Multiple alignment programs penalize deletions only where they occur, whereas single insertion events are penalized multiple times in each alignment iteration between their original occurrence and the root of the guide tree (Loytynoja and Goldman, 2005). PRANK uses outgroup rooting to explicitly distinguish deletions from insertions and avoids repeated penalization of insertions through the use of storage matrices that allows for distinct subsolutions to be compared through a match/gap-scoring scheme that simultaneously models substitution processes. Because this method considers indels as phylogenetic information, the algorithm may be sensitive to the order in which sequences are added (i.e., the initial guide tree topology). Further, the effects of guide tree bias on phylogenetic inference are likely to be amplified in large data sets, and such errors cannot be disregarded (Redelings and Suchard, 2005; Kumar and Filipowski, 2007:132). We generated a starting guide tree by concatenating intron fragments that were shown to be unambiguously aligned by multiple alignment programs, and we describe a 5-step protocol to construct an "unbiased" guide tree (one not influenced by indel placement) for the PRANK alignment (Fig. 4).

First, we used parameter perturbation and multiple alignment programs to identify and then purge the original intron matrices of positions with uncertain homology. In step 1, we used the program SOAP (Loytynoja and Milinkovitch, 2001) to align each intron under different parameters/algorithms and to purge the provi-

sional alignment from unstable blocks (i.e., those sensitive to parameter perturbation). Three programs were used to generate provisional alignments of each intron. Parameter perturbation alignments were produced by CLUSTAL W with gap-opening penalties ranging from 1 to 20 and gap-extension penalties ranging from 0.1 to 0.5 ( $4 \times 2$  combinations). Two additional alignments (for a total of 10 aligned matrices per intron) were included in the alignment pool, including (1) a tree-based partitioning algorithm coupled with multiple iterations (MUSCLE; Edgar, 2004); and (2) a Bayesian probabilistic sequence alignment (ProAlign; Loytynoja and Milinkovitch, 2003). MUSCLE and ProAlign were run with default parameters. All blocks or positions supported by less than 95% in the set of 10 provisional alignments were excluded from each intron to produce a single "purged" alignment. In step 2, the seven purged matrices were concatenated to each other and to the Cmos and Rag-1 coding sequences. The final matrix has a length of 2568 bp and includes roughly 54% of the original nuclear region data matrix of 4691 bp (Fig. 4). In step 3, the purged matrix was used to construct the guide tree through a Bayesian analysis (10 million generations) based on a GTR+I+G substitution model. This topology was then used as the single input tree for the progressive alignment of each of the seven introns with PRANK (Fig. 4, step 4). Unless indicated, all PRANK alignments were made with default parameters and the HKY substitution model. Once alignments for all indels were completed, the seven introns were concatenated with the two protein-coding genes to produce the final PRANK-aligned nuclear matrix (Fig. 4, step 5).

To further establish an a priori "baseline" against which to compare both manual and PRANK-aligned sequences, we constructed trees for the Cryba intron alone based on CLUSTAL and MUSCLE alignments; all of Cryba gene trees are then compared for congruence to the combined data tree (hereinafter the "combined" tree; see below). Our goal is not to provide an intensively "bench-marked" baseline on optimal alignment parameters (see instead, Terry and Whiting, 2005; Smythe et al., 2007), but rather to qualitatively evaluate tree topology and nodal support in the manual and PRANK-aligned sequences. We chose the Cryba intron because its sequence complexity features: (1) "orphan" sequences (sequences with no close similarity to the remaining species; e.g., *M. stolzmanni*); (2) distinct sequence subgroups (sequences with high similarity within but not between subgroups); (3) long- and short-repeat motifs, and (4) the highest frequency of indels among the sampled introns. All of these factors are likely to drastically reduce the accuracy of alignment algorithms (Pollard et al., 2004; Morrison, 2006).

We used topological congruence to the combined tree as the single criterion to compare accuracy of alignment results for the Cryba intron (Creer et al., 2006) and explicitly avoid the use of scoring functions (sum-of-pairs or column scores) because algorithms that depend on a pattern-matching reference alignment model sequence patterns rather than historical processes



## Manual alignment; Motif recognition and gap placement

### Step 1. Conspecifics

lulique CTCCCACCTTTGGCTTTGGCTTCAGC CTCCCACCTTTGGCTTTGGCTTCAGC  
 Talita CTCCCACCTTTGGCTTTGGCTTCAGC CTCCCACCTTTGGCTTTGGCTTCAGC  
 Calama CTCCCCTTTGGCTTCAGC CTCCCCTTTGGCTTCAGC  
 Pika CTCCCCTTTGGCTTCAGC CTCCCCTTTGGCTTCAGC

### Step 2. Heterospecifics

pacificus TGTTCGTGGCTTTGGCTAGCGGCTGGGCTCT TGTTCGTGGCTTTGGCTAGCGGCTGGGCTCT  
 abnormalemsis TGTTCGTGGCTTTGGCTAGCGGCTGGGCTCT TGTTCGTGGCTTTGGCTAGCGGCTGGGCTCT  
 habell CGTTCGTGCTGCG CGTTCCT CGTTCCT  
 habell CGTTCGTGCTGCG CGTTCCT CGTTCCT  
 bhvittatus CGTTCGTGGCTTTGGCTGCG CGTTCCT CGTTCCT  
 occipitalis CGTTCGTGGCTTTGGCTCT CGTTCCT CGTTCCT

### Step 3. Simple sequence repeats multiple alternatives:

3.1. alternative (1)  
 GCTGGGGGGGGGCTTT GCTGGGGGGGGGGCTTT  
 ACTGGGGG-----TCCT ACTGGGGGGGGG---TCCT  
 ACTGGGGGGGGG---TCCT ACTGGGGGGGGG---TCCT

3.2. alternative (1); inconsistent  
 GGTAT-----GAATTAADA GGTAT-----TGAATTAACA  
 GGTAT-----AAATTAACA GGTAT-----AAATTAACA  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA

alternative (2); inconsistent  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA

alternative (3); consistent  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA  
 GGTATATAATAATAATAACA GGTATATAATAATAATAACA

### 3.3.

alternative (1); inconsistent  
 CAGATTGATGATATATATACGATTATGC CAGATTGATGATATATATACGATTATGC  
 CAGATT-----gatlatgc CAGATTgatt-----ATGC  
 CAGATTgatt-----ATGC CAGATTgatt-----ATGC  
 CAGATT-----gat-----ATGC CAGATTgatt-----ATGC

alternative (2); consistent  
 CAGATTGATGATATATATACGATTATGC CAGATTGATGATATATATACGATTATGC  
 CAGATT-----gatlatgc CAGATTgatt-----ATGC  
 CAGATTgatt-----ATGC CAGATTgatt-----ATGC  
 CAGATT-----gat-----ATGC CAGATTgatt-----ATGC

### 3.4.

alternative (1); inconsistent  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC

alternative (2); inconsistent  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC

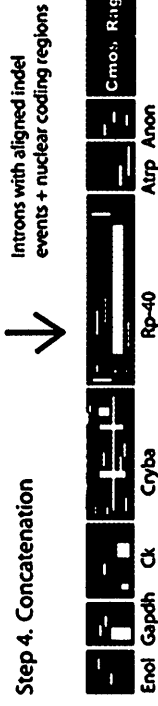
alternative (3); consistent  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC  
 TCCTCCCTCCCTCCCTCACC TCCTCCCTCCCTCCCTCACC

### 3.5.

alternative (1); inconsistent  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC  
 ATGA-----GATAGTGC ATGA-----GATAGTGC  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC  
 ATGA-----GATAGTGC ATGA-----GATAGTGC  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC

alternative (2); consistent  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC  
 ATGA-----GATAGTGC ATGA-----GATAGTGC  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC  
 ATGA-----GATAGTGC ATGA-----GATAGTGC  
 ATGACGAGATAGTAGTGC ATGACGAGATAGTAGTGC

### Step 4. Concatenation



Introns with aligned indel events + nuclear coding regions

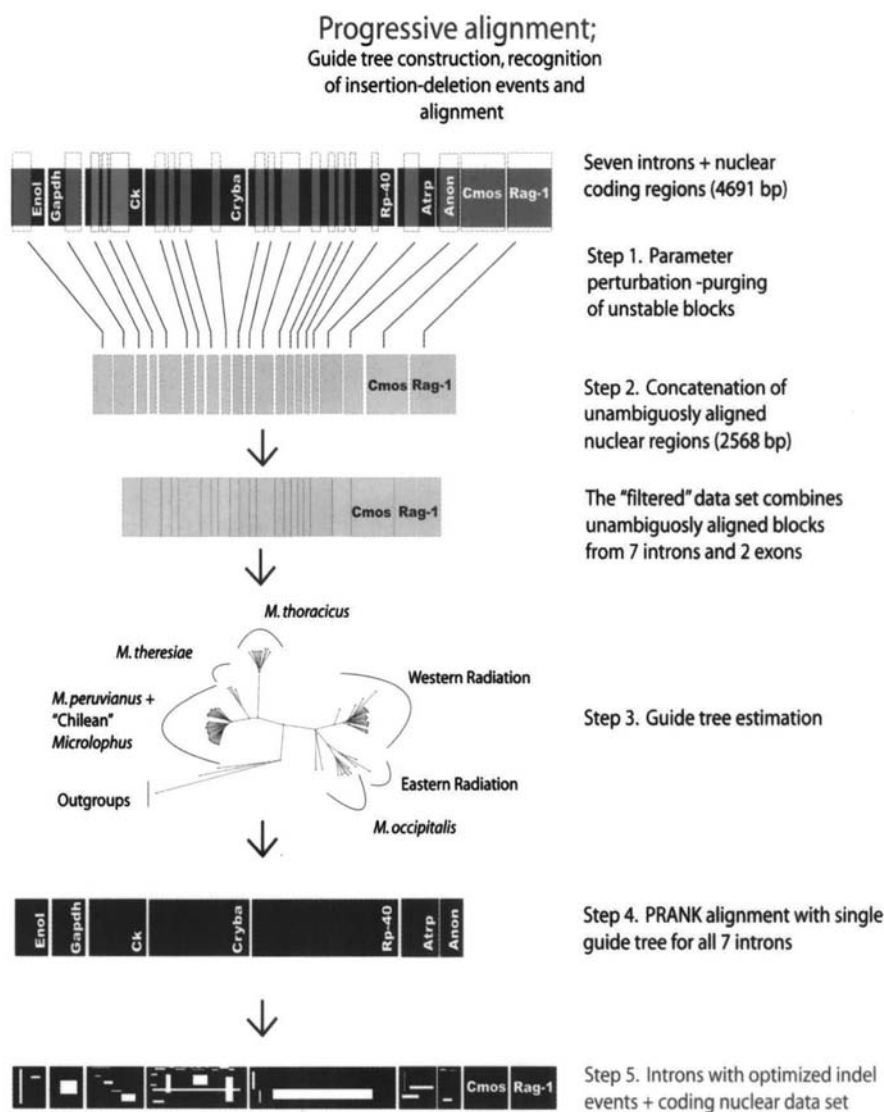


FIGURE 4. The software-assisted alignment to obtain the input tree for PRANK analyses. Step 1: Selection of unambiguously aligned blocks (in gray) through parameter perturbation of multiple alignments of each intron using SOAP and a 95% cutoff level (see text for details). Step 2: Concatenation of unambiguously aligned intron blocks plus nuclear coding genes (Cmos, Rag-1). Step 3: Building of a guide tree based on concatenated sequences from step 2. Step 4: Use of PRANK with the input tree obtained in step 3, for alignment of each of the seven introns. Step 5: Concatenation of PRANK-aligned introns plus coding regions (Cmos and Rag-1) to produce final nuclear sequence matrix.

FIGURE 3. Illustration of the protocol followed for the manual alignment of nuclear intron sequences in the lizard genus *Microlophus*. The color graph shows a gapless length-variable nuclear intron (Cryba) for all terminals of *Microlophus* and *Tropidurus* (outgroup species). Four colors identify different DNA nucleotides. Intron length, as measured in base pairs found after the amplification with a single pair of primers, varied from 185 bp in *M. stolzmanni* to 839 bp in *M. peruvianus*. The manual alignment steps are based on mechanistic explanations of intron evolution and given here in basic outline, and the same protocol was followed for each of the seven length-variable introns. In all frames below the color graph, flanking regions showing sequence integrity are shown in uppercase letters, whereas bold-font letters identify aligned motifs, underlining identifies simple sequence repeat (SSR) motifs, and base substitutions are shown with a "\*" Step 1: Conspecifics. Conspecific terminals from multiple localities that show intron length variation and three base changes are shown in the left panel, and the right panel shows subsequent indel placement in Calama and Pica sequences, which eliminates base substitutions and preserves flanking sequence integrity. Step 2: Heterospecifics. Heterospecific terminals showing original indel positional differences in left panel (overlapping indels), and subsequent placement of indels to improve sequence base-pair identify of flanking regions across five species (*M. pacificus*, *M. albemarlensis*, *M. habeli*, *M. bivittatus*, and *M. occipitalis*). Alternative step 3.1: Identification of length-variable mononucleotide strings of uncertain positional homology (in the third and fourth sequences); these were eliminated from the matrix (see text for details). Alternative step 3.2: Indel placements and motif recognition (alignments 1 and 2); only the third alternative reflects a simple sequence repeat (SSR) event that is consistent with hypothesized mechanisms. Alternative step 3.3: Sequence composition adjacent to an indel that could not resolve the position of a putative indel. In alternative 1 the indel is placed in three different blocks and sequence similarity is preserved; in alternative 2 the indel is placed arbitrarily but it ensures inference of a single mutational event across all ingroup taxa (maintaining the single mutation). Alternative step 3.4: Length mutations may overlap with one another to create a series of overlapping step indels, but here only the third alternative reduces the total number of possible mutations (alternatives 1 and 2 require 3 events and alternative 3 only 2). Alternative step 3.5: In this example two indels (alternative 2) instead of one (alternative 1) are inserted to eliminate the two inferred base substitutions in the flanking sequence required by inference of a single indel. Step 4: Concatenation. Manually aligned introns are concatenated to each other and the two coding regions (Cmos and Rag-1), to produce the nuclear gene matrix.

(Morrison, 2006). In contrast, PRANK uses an evolutionary scoring function that is expected to emphasize the biological correctness of both point and length mutation alignments, allowing, for example, insertions to be kept unaligned (Edgar and Sjölander, 2003; Lassmann and Sonnhammer, 2005).

#### *Phylogenetic Methods*

All introns were concatenated with both the nuclear protein-coding and the mitochondrial data sets for phylogenetic analyses. Indels were coded as an additional partition (see below) and their effect on the final topology was recorded through partitioned and overall measures of clade support.

Phylogeny reconstruction was carried out separately for the mtDNA locus, the nuclear data set, and then the combined data set. We used minimum and maximum numbers of a priori partitions for the three main data sets: (A) mtDNA partitions—2, 5, and 12 partitions; (B) nuclear DNA partitions—3, 10, and 14 partitions; and (C) combined data set—12, 16, and 26 partitions (see online Appendices 4 and 8 for further details on partitions; <http://www.systematicbiology.org>). We derived likelihood scores for each partition format of each data set, and used log likelihood-ratio tests to evaluate alternative partitions under the null hypothesis that adding more partitions (i.e., parameters) does not significantly alter likelihood scores. Failure to reject the null hypothesis means that an increase in parameter number does not improve the phylogenetic result, whereas significant differences in likelihood scores suggest that more complex models add phylogenetic information (Sullivan and Joyce, 2005).

Maximum parsimony (MP) analyses were run under equal character weighted heuristic searches with 1000 replicates of random addition and tree bisection and reconnection branch swapping (TBR). A max trees limit of 1000 and 10,000 was used for single and concatenated genes, respectively. MP was implemented in PAUP\* 4.0b5 (Swofford, 2002), with branch support estimated from 500 bootstrap pseudoreplicates (summarizing five independent runs of 100 replicates with 10 random addition replicates each). Maximum likelihood (ML) bootstrap values were obtained through 1000 pseudoreplicates using a fast algorithm (PHYML; Guindon and Gascuel, 2003).

Replicated Bayesian analyses (started from independent random trees) coupled with Markov chain Monte Carlo (MCMC) simulations were run for 10 million generations using four incrementally heated Markov chains sampled every 1000 generations. Model selection for all partitions was based on DT-ModSel (Minin et al., 2003), because it typically selects simpler (less parameter rich) models. Simpler models have been shown to estimate branch lengths and tree topologies with less error and yet the same accuracy as more complex models selected by hLRTs, AIC, or BIC approaches (see Abdo et al., 2005). Mixing of phylogenetic parameters and stationarity of likelihood scores were assessed using the program Tracer v1.2 (Rambaut and Drummond, 2003).

All analyses were run using the parallel version of Mr. Bayes (v3.04b; Ronquist and Huelsenbeck, 2003) on the Beowulf cluster housed in the Department of Integrative Biology at Brigham Young University. Nodal support was arbitrarily considered strong with bootstrap values >70% (but see caveats in Hillis and Bull, 1993) and posterior probability values >95% (Alfaro et al., 2003; see Lewis et al., 2005 for caveats), and gene trees recovering strongly supported conflicting nodes are interpreted as real conflict (Wiens, 1998). Conflict between gene trees might be due to differences in coalescent histories, recombination, nonorthology, pseudogene amplification, or human error, and the absence of conflict by these criteria is taken here to imply that these potentially confounding factors have minimal influence on the genealogies of the nuclear regions.

#### *Coding Length Mutations for Phylogenetic Analyses*

Single length mutational events inferred from manual and PRANK alignments were coded as individual binary characters following Simmons and Ochoterena (2000). Indel information was incorporated in both MP and Bayesian tree searches, with the Bayesian analyses of indel partitions based on the maximum likelihood model (Markov k) developed for morphological characters (Lewis, 2001). In all cases we assumed that the indel partitions had unequal rates among characters and so we incorporated a gamma distribution (Mk+G). We assessed the effect of including the indel partition with Bremer partitioned support values (PBS; using TreeRot [Sorenson, 1999] under the heuristic search parameters: addseq = random nreps = 1000 swap = tbr hold = 10) over the preferred Bayesian topologies. We are aware of the limitations of Bremer support values (see DeBry, 2001) and do not interpret any of these values as relative support for a given node. Below we show that individual nuclear gene trees are topologically congruent among themselves but conflict with the mtDNA locus at some nodes, and we use PBS to identify character incongruence localized to specific nodes. This conflict has significant evolutionary ramifications in our study.

## RESULTS

### *Patterns of Sequence Variability*

Ten genomic regions (the mtDNA locus, nine nuclear gene regions) were collected for all terminals, and after coding for length mutations, an eleventh partition was added. Both mitochondrial and nuclear genes show moderate to high levels of variation; the ND4 and the tRNAs regions were the most variable mitochondrial regions, whereas Cryba was the most variable among nuclear introns (estimated by maximum pairwise uncorrected distance values; online Appendix 4; <http://www.systematicbiology.org>). BLAST searches did not reveal any match to known genes for the Anon locus, and patterns of variation showed that the nuclear intron regions evolve roughly at the same rate as ribosomal genes (12S and 16S) and two to four times more rapidly than the two nuclear exons (Cmos and Rag-1).



### Protein-Coding and Ribosomal Gene Alignments

Alignment of mitochondrial protein and nuclear protein-coding genes (Cyt-*b* and ND4; Cmos and Rag-1, respectively) was facilitated by conservation of the codon reading frames. Alignment of the ribosomal 12S, 16S, and tRNA<sup>his</sup> + tRNA<sup>ser</sup> regions first identified conserved blocks that, upon comparison with secondary structure models, allowed the recognition of stem and loop partitions. Regions of questionable homology—commonly found in loops—were excluded; 9 of 10 excluded regions in the 12S gene corresponded to loop regions (80 bp in total), and 4 of the 5 excluded regions in the 16S gene were located in regions identified as loops (32 bp in total). Only four tRNA loop base positions were deemed ambiguously aligned and removed from phylogenetic analyses.

### Manual Indel Alignments

We inferred a total of 117 length mutation events across all 73 terminals in the seven noncoding gene regions (Table 1), with numbers of indels/parsimony informative indels ranging from lows in Enol (6/5) and Gapdh (13/4) to a high in Cryba (51/39). Online Appendix 5 (<http://www.systematicbiology.org>) lists these events and provides a detailed description of the 117 manually aligned length mutations coded in the indel matrix. On a single gene basis, we inferred 51 indels in Cryba (the gapless length of this intron ranges from 185 to 839 bp; Table 1), and we derived these as 45 deletions, five insertions, and a single simple sequence repeat (SSR) insertion. Large length modifications are also apparent in the RP40 intron (327 to 755 bp; Table 1); here the original intron length has been increased by an insertion that characterizes terminals of *M. occipitalis* (mainland; northern Peru) and *M. habeli* and *M. bivittatus* (Eastern Galápagos; Marchena and San Cristobal Islands).

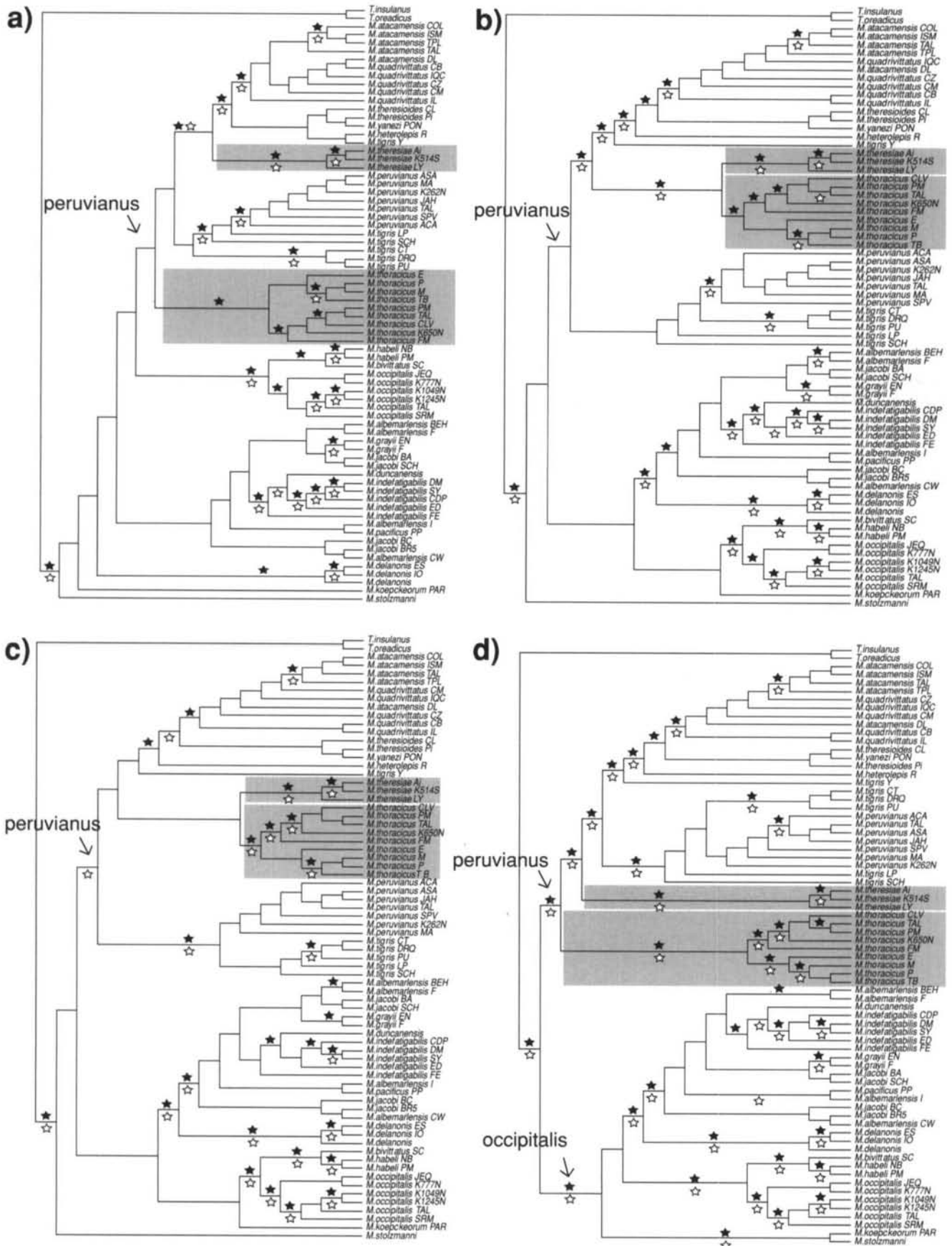
We were unable to discern nucleotide positional homology for parts of two gene regions (CK and Cryba). In the first case, CK presents a SSR trinucleotide repeat (TCC) between positions 76 to 95 in the *Peruvianus* group (PAUP file in online Appendix 6A; <http://www.systematicbiology.org>). The CK alignment shows two continuous homonucleotide strings (cytosine and guanine between positions 267 to 288) in all taxa except *Tropidurus* and *M. thoracicus* terminals that seem to have independently lost this repeat. Homonucleotide strings can be considered as stepwise indels (Lohne and Borsch, 2005), but we ignored this region (it was coded as missing data) because we could not reject the possibility it may have resulted from inaccurate *Taq* enzymatic activity during PCR (Kelchner, 2000; several amplification/sequence attempts failed to show a repeatable number of mononucleotide repeats). In the second case, the Cryba intron shows an insertion of a nonrepeat 8-bp segment in the *Occipitalis* group (positions 1001 to 1007; Appendix 6A) that has no apparent positional homology with remaining ingroup or outgroup taxa, and it was also coded as missing data.

In some cases extensive variation between the ingroup and outgroup terminals forced indel placements that produced blocks of putatively homologous sites that were not continuous with blocks that showed nucleotide identity. Tentatively aligned segments may result due to high divergence between ingroup and outgroup taxa (e.g., positions 762 to 773 and 991 to 997; Appendix 6A). In a second case, a particular segment of outgroup sequence is included within a larger indel in the ingroup (e.g., positions 790 to 811, 825 to 829, 837 to 859, 871 to 880, and 974 to 983; Appendix 6A) and can only be shifted within the boundaries of the longest indel. Under this constraint, the limited number of possible homology statements forces the outgroup blocks to align with the only available, but not necessarily similar, ingroup positions within the larger indel. This constitutes over-alignment (i.e., overlapping of nonhomologous sequences; Cline et al., 2002; Morrison, 2006), which by the criterion of maximizing blocks of sequence identity is a suboptimal solution.

### PRANK Alignments

The alignment of all intron sequences was optimized under the same guide tree obtained from the filtered unambiguous alignment of the nine nuclear regions (2568 bp in total; see Material and Methods). The guide tree (not shown) successfully recovered nodes at deep and intermediate levels of divergence, but shallower nodes were not fully resolved. Because PRANK ignores unresolved nodes, we used a neighbor-joining tree based on uncorrected distances of this same data set. The NJ tree (not shown) is fully resolved and completely congruent with the topology resolved in the filtered tree. PRANK allows the use of simpler models of substitution (Jukes-Cantor [JC] and Hasegawa-Kishino-Yano [HKY]) so we approximated the models selected by DT-ModSel for the manually aligned partitions (Appendix 4) by selecting the HKY option in PRANK. Overall, results of the progressive alignment contain more indels than the manual solution (up to 131 coded gaps; Table 1 and online Appendix 6B [<http://www.systematicbiology.org>]). More specifically, the number of indel characters after the pairwise alignment remained the same for two genes (Enol and RP40), decreased for one (Anon), increased slightly in three (Gapdh, Atrp, and CK), and was substantially higher in the most variable nuclear intron (Cryba; Table 1).

Overall, PRANK found one third of the indels we inferred manually (38 of 117; Appendix 5) and it produced alternative solutions that accommodated hard-to-align regions in more parsimonious ways than our manual alignments (e.g., positions 185 to 196 in Appendices 6A and 6B). PRANK alignments normally increased the number of identical base pairs (2482 to 2571) and decreased the number of parsimony informative sites in the blocks of sequences flanking the indels (Table 1). Examining CK as an example, PRANK increased the number of identical base pairs from 248 to 254 and the



number of indels from 11 to 15 (9 and 13 parsimony informative, respectively), whereas decreasing the number of parsimony informative sites from 84 to 76 (Table 1). Additional differences were the rearrangement of the trinucleotide simple sequence repeat (SSR) of positions 76 to 95 (Appendices 6A and 6B), and PRANK inferred two single-base pair insertions restricted to the outgroup taxa (positions 191 and 275; Appendix 6B).

#### *The Cryba Intron as an Example*

The alignment of Cryba underscored the need for selection of a substitution model and correction parameters implemented in PRANK. Our preferred alignment (i.e., the one consistent with the combined evidence tree) was obtained when using the HKY model ( $\kappa = 1.830$  and empirical base frequencies 0.229/0.276/0.230/0.263 for A/C/G/T) and forcing insertions to be always skipped. Compared to the manual solution, the PRANK alignment significantly reduced the number of parsimony informative sites (from 204 to 178; Table 1), increased the number of identical positions (from 583 to 657; implying fewer point mutations after indel optimization), and obtained higher likelihood scores (see Fig. 5). PRANK opened a number of gaps and insertions already present in the manual alignment but refined this alignment by inferring more parsimonious solutions for the location of indels (see Appendix 6B for additional details). Most importantly, PRANK distributed the reduced sequence of *M. stolzmanni* (185 versus 839 bp in the outgroup) in a solution that maximized nucleotide identity in small blocks, without modifying the default values for the gap extension (gaprate) and gap-opening (gaptxt) probabilities.

To compare the level of phylogenetic congruence of the PRANK versus other alignments, we realigned the Cryba intron with CLUSTAL and MUSCLE (using default parameters in both) and then generated Cryba gene trees for these plus the manual and PRANK alignments. Figure 5 illustrates four gene trees and shows that the PRANK tree recovers more of the deep nodes present in the combined tree, and with higher support, than any of the others. It is the only topology that, for example, recovers the *Occipitalis* and *Peruvianus* clades with strong support and, within the latter, the *M. thoracicus* and *M. theresiae* clades at the base of the *Peruvianus* group (all with strong support). Gene trees obtained from CLUSTAL, MUSCLE, and manual alignments recovered “wrong” topologies in conflict with the combined tree and having fewer well-supported clades. Notably, the PRANK alignment that modeled the Cryba transition/transversion ratio (using the HKY model) recovered a similar topology but 126.2 likelihood units better and 18 steps shorter

(TL = 312) than the PRANK Jukes and Cantor default alignment.

#### *Phylogenetic Analyses—Mitochondrial DNA*

A single-alignment hypothesis based on codon conservation and secondary structure was used for phylogenetic analyses of mitochondrial partitions. Appendix 4 summarizes substitution models for three partitions of the mtDNA locus, and likelihood scores of 2 ( $\ln L = -24556.524$ ), 5 ( $\ln L = -24519.738$ ), and 12 ( $\ln L = -24257.128$ ) partitions showed an increase of  $\sim 200$  log likelihood units by fully partitioning the mtDNA locus. Log likelihood-ratio tests showed that these partitions differed significantly ( $P < 0.001$ ) from each other (details available from EB upon request), so we illustrate the 12 partition tree. Figure 6 shows the Bayesian consensus topology, which recovers all members of the *Peruvianus* group (node 3) except *M. thoracicus* (node 58) and resolves “northern” (node 14) and “southern” groups (node 4) nested within node 3. At more nested levels, mtDNA data do not resolve the polytomy at node 15, nor do they resolve exclusive species within these groups (*M. atacamensis*, *M. tigris*).

In contrast, the monophyly of the *Occipitalis* group is recovered with strong support (node 30), and within this clade, the “Western Galápagos Radiation” proposed earlier (Wright, 1983; Lopez et al., 1992; Heise, 1998; Kizirian et al., 2004) is recovered with strong support (node 31). *Microlophus delanonis* is strongly supported as the sister species to the remaining taxa in this clade, and although there is strong support for monophyly of conspecific terminals, the polytomy of node 32 prevents a full understanding of the interspecific relationships in this radiation. The two species comprising the “Eastern Galápagos Radiation” (endemics from Marchena [*M. habeli*] and San Cristobal [*M. bivittatus*] islands are strongly supported as the sister group of the mainland *M. occipitalis*) (node 50). Phylogenetic placement of the mainland *Occipitalis* group species *M. koepckeorum* and *M. stolzmanni* are equivocal (nodes 48 and 49).

#### *Phylogenetic Analyses—Nuclear DNA*

Online Appendices 7A and 7B (<http://www.systematicbiology.org>) show individual nuclear gene trees and their nodal support, respectively, and there was no evidence of strongly supported conflict among any of these trees. Appendix 4 summarizes the substitution models for the three partition strategies for the nuclear data, whereas online Appendix 8 (<http://www.systematicbiology.org>) shows likelihood values for individual gene trees for manual or PRANK

FIGURE 5. Topologies from phylogenetic analyses of four alignment hypotheses of the Cryba intron (MP topology illustrated). (a) CLUSTAL ( $\ln L = -3730.23$  [Bayesian];  $\ln L = -3659.377$  [ML]; TL = 563 [MP]). (b) MUSCLE ( $\ln L = -3672.755$ ;  $\ln L = -3326.900$ ; TL = 393). (c) Manual ( $\ln L = -3568.648$ ;  $\ln L = -3239.158$ ; TL = 359). (d) PRANK (default parameters;  $\ln L = -3530.668$ ;  $\ln L = -3209.069$ ; TL = 330). Shaded clades show the relative positions of *M. thoracicus* and *M. theresiae* in each topology, whereas the *Occipitalis* and *Peruvianus* groups are identified only if recovered. Solid and open stars identify nodes supported by Bayesian PP  $\geq 0.95$ , and MP bootstrap proportions  $\geq 70$ , respectively.

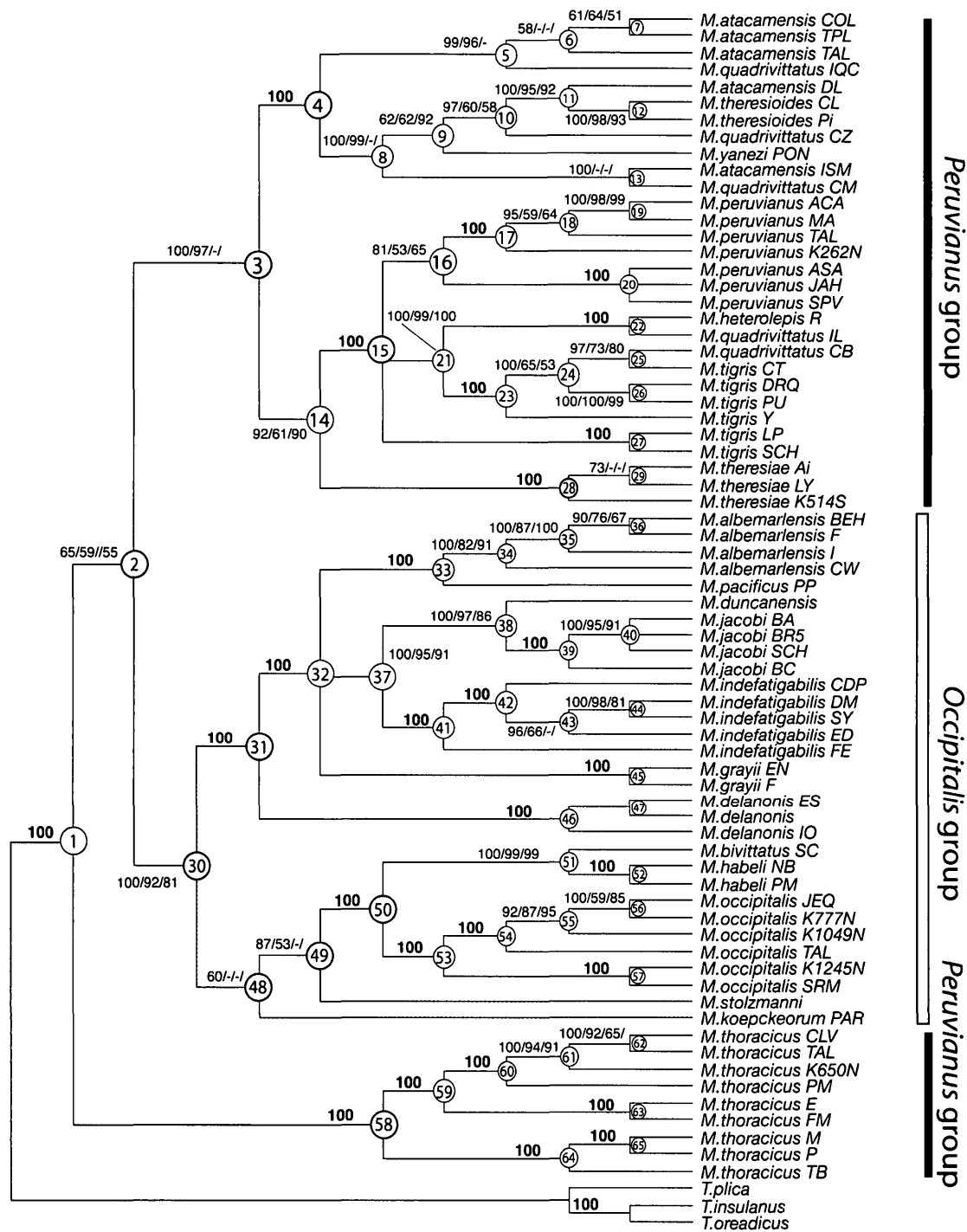


FIGURE 6. The 50% majority-rule consensus of Bayesian MCMC trees for the mtDNA locus. Bolded circles identify the nodes discussed in the text. Numbers above nodes indicate Bayesian/ML/MP support values (ML and MP shown only if > 50); a single bold 100 identifies nodes for which all support indices are 100.

alignments. Log likelihood ratio tests of among 3 ( $\ln L = -16,052.580$ ), 10 ( $\ln L = -16,070.398$ ), and 14 ( $\ln L = -16,113.739$ ) partitions revealed a significant difference between all paired combinations ( $P < 0.001$ ), so the 14-partition tree is preferred. Figure 7 shows the consensus Bayesian topology, and again both the *Occipitalis* (node 35; supported by eight

partitions) and the *Peruvianus* group (node 2; supported by four partitions) are recovered with strong support.

Within the *Occipitalis* group there is some support for *M. stolzmanni* as the sister clade of the two Galápagos radiations (node 36), and seven partitions recovered *M. occipitalis* as the sister group of the Eastern Galápagos

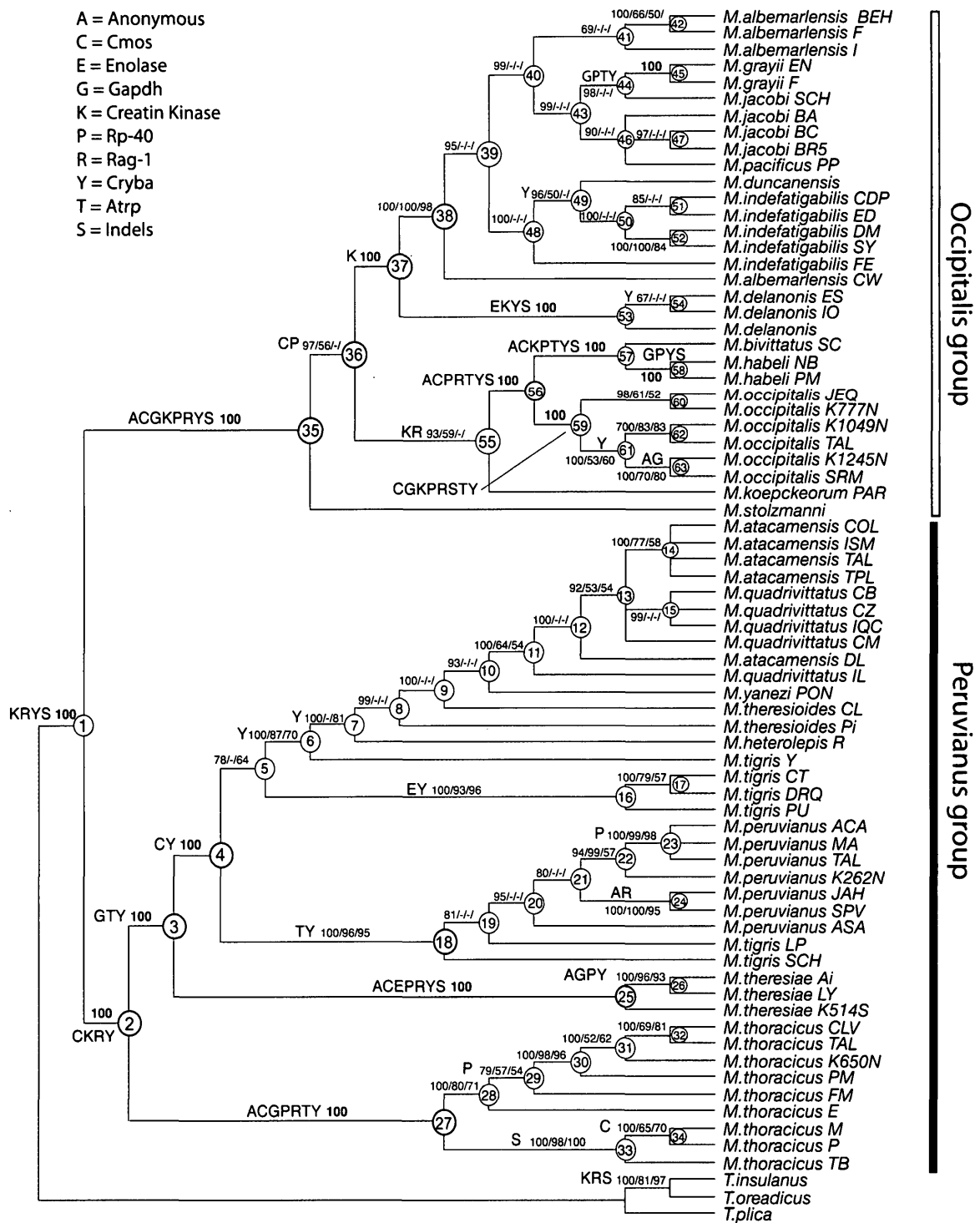


FIGURE 7. The 50% majority-rule consensus of Bayesian MCMC trees from combined nDNA gene regions; all symbols are as in Figure 6. Letters identify individual gene regions that provide significant independent support ( $PP \geq 0.95$ ) for resolved clades.

Radiation (*M. bivittatus* and *M. habeli*; node 56). Resolution of a monophyletic Western Galápagos Radiation is strongly supported (node 37), as is *M. delanonis* as the sister clade to all others in this radiation (nodes 37 and 38). Within this radiation, some heterospecific clades are evident due to nonmonophyly of *M. albemarlensis* and *M. jacobi*. Within the *Peruvianus* group, seven genes recover

the *M. thoracicus* clade (node 27), three genes recover *M. theresiae* (node 28) as the sister taxon of the remaining *Peruvianus* group species (node 3), and two genes recover populations of *M. tigris* in two distinct clades (nodes 4/18). Heterospecific clades are also evident at more nested levels due to nonmonophyly of *M. atacamensis*, *M. quadrivittatus*, and *M. theresioides*.

### Combined Data and a Phylogenetic Hypothesis for *Microlophus*

Partitioned Bayesian analysis of the combined mtDNA and nuclear DNA (with coded indels) data sets produced a 50% majority rule consensus tree with an average  $\ln L = -42,033.714$  and an average of  $\ln L = -42,091.843$  for the manual and PRANK-aligned introns, respectively (Appendix 5). There are no differences in the topology between these trees and Figure 8 displays the topology recovered with PRANK-aligned introns. For this data set, likelihood values extracted from Bayesian analyses with 12 and 26 partitions were not significantly different ( $\chi^2 = 48.627$ ;  $df = 104$ ;  $P = 1.000$ ), thus we preferred the partition format that estimated the optimal topology with fewer parameters. Table 2 summarizes different measures of clade support for key nodes shown in Figure 8 (details for all nodes are given in online Appendix 9; <http://www.systematicbiology.org>).

The combined data topology recovers the majority of the clades previously recovered with either the mtDNA or nuclear data sets (see exceptions below). At the deepest levels of divergence, these include the *Peruvianus* and *Occipitalis* groups (nodes 2 and 40, respectively), each of which is strongly supported by all support indices (Bayesian, ML, and MP), as well as multiple data partitions. Within the *Occipitalis* group, there is strong support for the Western and Eastern Galápagos clades (nodes 42 and 62, respectively). Within the *Peruvianus* group, there is strong support for the placement of *M. thoracicus* (nodes 2) and *M. theresiae* (nodes 3) as second and first outgroups, respectively, to a strongly supported “peruvianus” plus “Chilean *Microlophus*” clade (node 4). Each of these two latter clades also received strong support (nodes 15 and 5, respectively), and the general pattern is strong support of the deepest clades by the nuclear partitions and shallower clades by the mtDNA (Table 2).

We note some topological discordance between mitochondrial and nuclear genes stemming from (1) suboptimal node resolution due to markers resolving different temporal scales, coupled with probable inflated PP values at short nodes; and (2) conflict due to significant incongruence between the mitochondrial and nuclear data sets (Wiens, 1998). An example of suboptimal node resolution is the ambiguous placement of *M. stolzmanni* and *M. koepckeorum* within the *Occipitalis* group (nodes 41 and 61; Fig. 8); this appears to be due to a limited number of characters and weak conflict among data sets, as both species are equivocally placed by either genome.

The most obvious example of mtDNA versus nDNA conflict is found at clades nested within the *Peruvianus* group; compare the symmetrical mtDNA topology (nodes 4 and 14 in Fig. 6) with the pectinate nDNA topology (nodes 5 to 13 in Fig. 7) versus the combined data that differs slightly from both (nodes 5 and 15 in Fig. 8). As an example, terminals in the strongly supported Chilean *Microlophus* clade in the combined tree (node 5; Fig. 8) are grouped in the mtDNA topology into an ar-

angement that does not correlate with the geographic distribution of populations of the putative species *M. atacamensis* and *M. quadrivittatus*. Nodes within this clade are short (Fig. 8) and weakly supported except for some high PP values, suggesting inflated Bayesian support (Lewis et al., 2005). We nevertheless accept this topology as the best working hypothesis for *Microlophus*, although as a species tree it suggests that both the Chilean *Microlophus* and “peruvianus” clades (nodes 5 and 15) need revisionary work.

## DISCUSSION

### *Alignment of Noncoding Regions: Manual and Progressive Approaches*

Most of the indels inferred by PRANK and by our manual assessments are phylogenetically consistent (i.e., they recover the same synapomorphies for the same nodes), but PRANK found additional indels that we could not find by manual adjustments in ambiguous regions (Table 1). In our experience, trying to resolve length-variable nucleotide regions flanked by conserved blocks is a daunting task that requires sliding and partitioning sequence blocks in a highly subjective manner. In the Cryba alignment, for example, the CLUSTAL matrix (not shown) was characterized by outgroup aligned sequences that are more or less in the same positions as our manual alignment, and these positions were overaligned (i.e., nonsimilar regions that should not be aligned; see positions 804 to 828; Appendix 6A). Applying different gap costs do not solve the issue; gap costs set too high result in few or no gaps in the matrix (and hence the overaligned sequences), or if too low the alignment is extensively fragmented into many small gaps with little or no phylogenetic consistency among sequences (i.e., they recover the same gaps for terminals that are not related). For the Cryba intron the degree of congruence to the combined tree, likelihood scores, and tree lengths based on four alignment strategies is directly proportional to the number of invariant positions (CLUSTAL [418], MUSCLE [550], manual [583], PRANK [657]) and inversely proportional to the number of parsimony informative sites. Using congruence to the combined data tree as a proxy of alignment accuracy, the PRANK alignment suggests that an increase in the number of invariant bases (i.e., more statements of sequence homology) increases the phylogenetic signal in the remaining variable characters. This observation corroborates a recent simulation study (Ogden and Rosenberg, 2006) showing that sequence homology accuracy correlates positively with accuracy of inferred phylogenies across different tree shapes and phylogenetic methods.

It can be argued that the score functions of the CLUSTAL and MUSCLE algorithms applied to the alignment of an indel-rich intron are not accurate because they do not incorporate, for instance, penalties for overalignment (Cline et al., 2002; Edgar and Sjölander, 2004; Blackshields et al., 2006). However, PRANK bypasses this objective function problem by using outgroups for “weighing” indels, and the HMM solution and

TABLE 2. Summary of support indices (MLBS [PHYML], MPBS [PAUP\*], PP [MrBayes]) for selected nodes discussed in the text (highlighted in Fig. 8); the first and second rows show values obtained for manual and PRANK aligned intron regions, respectively. Fifteen partitions were considered for the calculation of Bremer supports, including: five mtDNA genes (2991 bp), nine nuclear regions (4691 and 4755 bp by manual and PRANK alignments), and one indel partition ( $n = 117$  and 131 characters [manual and PRANK alignments]). The full record for all nodes can be viewed in Appendix 9.

No.	Node description	Partitioned Bremer support																	
		MLBS	MPBS	PP	12S	16S	ND4	tRNA	Cyt <i>b</i>	CK	Cryba	RP40	Enol	Gapdh	Anon	Atrp	Cmos	Rag-1	Indels
2	<i>Peruvianus</i> group	100	100	1.0	0	1	1	-1	-1	5	1	1	1	2	4	1	4	3	-1
3	Node 4 + node 30	100	100	1.0	0	1	1	-2	6	4	1	2	1	1	6	2	5	4	2
4	Node 5 + node 15	100	100	1.0	3	-1	-4	-1	5	12	2	4	3	4	2	1	1	1	1
5	Chilean <i>Microlophus</i>	98	87	1.0	-3	-1	3	-3	7	12	2	4	3	4	3	3	1	1	1
15	"Peruvianus"	100	100	1.0	-5	-1	3	0	-1	4	1	0	0	1	4	1	2	2	2
30	<i>M. theresiae</i>	100	100	1.0	25	10	29	4	29	1	3	1	0	0	0	0	0	0	1
32	<i>M. thoracicus</i>	100	100	1.0	32	14	34	6	38	1	2	1	0	1	0	0	0	0	1
40	<i>Occipitalis</i> group	100	100	1.0	16	7	16	-0	12	-2	-3	-1	0	-2	-1	0	-1	0	-1
41	<i>Occipitalis</i> group— <i>M. stolzmanni</i>	100	100	1.0	21	10	21	0	17	-3	-5	-1	0	-2	-1	0	-1	0	-1
42	Western Galápagos Radiation	55	51	1.0	8	4	19	1	23	2	5	3	2	0	2	0	3	2	4
59	<i>M. delanonis</i>	100	100	1.0	11	6	25	1	32	3	5	4	3	1	4	0	4	2	7
61	<i>M. koeppckeorum</i> + node 62	100	100	1.0	16	5	13	1	11	4	5	11	1	5	5	5	4	4	2
62	E. Galápagos + <i>M. occipitalis</i>	100	100	1.0	21	8	17	1	17	6	10	15	2	8	8	7	6	6	4
		100	100	1.0	3	-1	1	1	-1	3	0	3	2	2	1	-1	1	2	1
		100	100	1.0	6	0	1	1	-5	6	0	5	5	6	5	-1	3	5	3
		100	100	1.0	0	-1	-1	1	3	0	-1	2	0	0	-2	-1	0	2	-1
		100	100	1.0	-3	-1	-3	2	9	0	-2	3	0	0	-3	-1	1	3	-2
		100	100	1.0	7	-3	7	2	7	2	3	3	1	1	2	0	3	2	1
		100	100	1.0	13	-6	12	2	13	7	6	3	3	3	5	1	6	3	3
		100	100	1.0	10	8	5	4	15	2	3	0	1	1	0	1	0	1	2
		100	100	1.0	12	10	7	6	23	2	4	-1	2	1	0	1	0	1	2
		54	100	.87	-2	0	3	0	3	0	-1	1	0	0	-1	-1	0	1	-1
		100	100	1.0	-4	-1	3	0	8	0	-1	0	0	0	-1	0	0	1	-2
		100	100	1.0	-1	2	3	2	4	2	0	5	0	0	-1	1	1	3	1
		100	100	1.0	-2	5	7	3	7	4	0	10	1	0	0	3	2	5	5

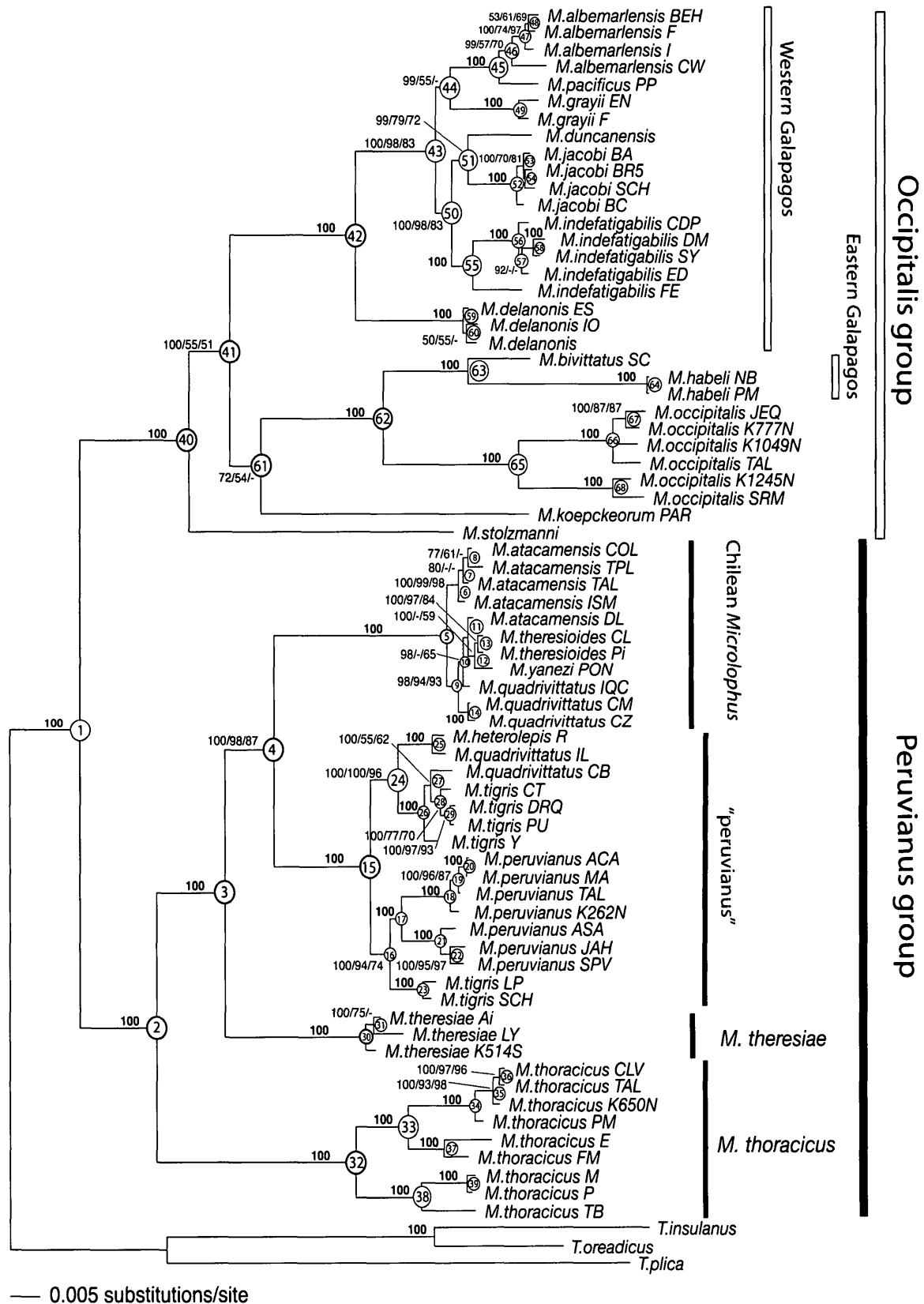


FIGURE 8. The 50% majority rule consensus phylogram of 9001 Bayesian MCMC trees (10 million generations) from 12 data partitions for the combined data set; nodes and support values are numbered as in Figures 6 and 7. Topologies obtained from manual (117 indels; 7682 characters total) versus PRANK (131 indels; 7902 characters) alignments are identical; Table 2 summarizes additional measures of support for key nodes under both alignments.



evolutionary models result in “staggered alignments” (Higgins et al., 2005; Morrison, 2006) that recovered a higher number of indels in all seven introns.

A practical difference among alignments refers to sensitivity (ability to detect all residues that should be aligned), selectivity (ability to align only those residues that should be aligned; Lambert et al., 2003; Morrison, 2006), and the solution found by Loytynoja and Goldman (2005) implements this distinction successfully. Still, in some situations, manual alignment may be more sensitive and selective than PRANK. For example, in the alignment of simple sequence repeats (SSR) of positions 70 to 110 (Appendices 6A and 6B), the manual approach provides a more parsimonious solution of two versus four evolutionary (indel) events.

As larger and more heterogeneous data sets become routine in phylogenetic analyses, intron-rich data sets will become common because noncoding sequences comprise the vast majority of metazoan genomes, particularly in vertebrates (Roy and Gilbert, 2006). Accordingly, both theoretical and empirical studies focusing on alignment issues for these types of markers are critical (Metzler, 2003; Miklos et al., 2004; Keightley and Johnson 2004; Holmes, 2005; Creer et al., 2006; Siddharthan, 2006), but in contrast to protein and ribosomal gene regions, there are no sources of independent evidence (i.e., secondary structure or codon conservation) to corroborate the correct alignment of noncoding DNA (Pollard et al., 2004). Furthermore, heuristic algorithms are undermined by high-frequency length mutations and repeat motifs of variable complexity, and the sheer number of possible mutation patterns of noncoding DNA defies standard procedures identified for other gene classes.

In practice, common heuristic methods search for unambiguously aligned regions first; they build an approximate guide tree and from this assemble an alignment that in the case of introns should optimize gap placement. Gaps are opened under the constraint of penalties that define the extension and cost of each according to an objective function that balances gaps and matches. Interestingly, most programs that rely on a guide tree regard insertions introduced early in the matching of two similar sequences as a nuisance because, unlike deletions, they are not penalized only at the place where they occur but rather have to be penalized repeatedly as the algorithm progresses down to the root of the tree (Loytynoja and Goldman, 2005). Heuristic algorithms (i.e., implemented in the CLUSTAL, MUSCLE, and T-COFFE programs, among others) do find approximate solutions to the alignment of both point and length mutations, but with the caveat that insertion placements may be suboptimal as long as heuristic programs consider insertions equivalent to deletions.

The complexity of inferring indels is also increased because progressive methods that rely on the construction of an initial guide tree assume that the best information on gap placement will be found among the most similar sequences, when in reality there may be better information in the alignment of all sequences considered in the group of interest. No computer approach is fully

capable of correcting a particular alignment based on a global solution, but the first step in this direction should be a program that handles insertions and deletions differently, and this is the advantage of the PRANK hidden Markov model (HMM) implementation. The method developed by Loytynoja and Goldman (2005) distinguishes insertions and deletions as separate events (via outgroup rooting) and then corrects early mistakes in the placement of gaps (indels) in a progressive alignment context. A key innovation of this development is that PRANK does not simply maximize a similarity function score to assess the accuracy of its alignment; rather it implements a phylogenetic scoring function that is based on simple evolutionary models (Loytynoja and Goldman, 2005). Normally, score or objective functions (e.g., sum-of-pairs or column scores) are calculated by maximizing the sum of similarities for pairs of sequences and a reference alignment (Blackshields et al., 2006; Edgar and Batzoglou, 2006), and their utilization represents the core of the disagreement between biological correctness and algorithm optimization (Lassmann and Sonnhammer, 2005; Konagurthu and Stuckey, 2006). The use of different objective functions is pivotal to understand why different algorithms can produce different alignments (Edgar and Sjölander, 2003; Morrison, 2006; Ogden and Rosenberg, 2006). The PRANK scoring function uses dynamic programming to move among storage matrices and select the optimal subalignment based on substitution models, thereby avoiding entrapment on a local alignment optimum.

#### *Aligning Length-Variable Introns and Guide Trees*

A key component of tree-based progressive alignments is the order in which pairwise alignments are made (Thompson et al., 1994). The modification of the hierarchy of these successive alignments towards the root of the tree has a large effect in the overall alignment (Redelings and Suchard, 2005). The degree to which guide trees affect downstream phylogenetic analysis in a multigene approach has not been studied and deserves increased attention (Kumar and Filipowski, 2007). We are also unaware of studies on the effect of length polymorphisms on the construction of guide trees, and a recent genomic approach showing different starting trees produced distinct but equally well supported topologies (Kumar and Filipowski, 2007). In our data set, Cryba shows “orphan” sequences and extensive length polymorphisms that produce “incorrect” phylogenies with CLUSTAL or MUSCLE alignments (Fig. 5a, b).

Another issue is that coding and noncoding DNA may be found as a mosaic of homologous and nonhomologous regions for which differentiation in an alignment context will be particularly difficult for noncoding DNA without independent evidence to corroborate primary homology (Morrison, 2006). Recombination and lateral transfer are frequently mentioned for tree discrepancies in genomic studies (Raymond et al., 2002), implying that the hierarchy of the “Tree of Life” might

be harder to obtain (Doolittle, 1999; Lebrun et al., 2006). Very few alignment programs model recombination and they do so in a protein alignment context (Lee et al., 2002; Blanchette et al., 2004). We suggest that some of these concerns may be obviated by building a guide tree based on unambiguously aligned fragments. Its application may be contested on grounds of circularity (see de Queiroz, 1995, for contextual details), but it also represents a compromise option for progressive alignments that require guide trees.

#### Phylogenetic Information Content of Intron Length Substitutions

Several important observations emerged from the inclusion of the length mutations from the seven nuclear introns considered in this study. First, inclusion of manual and PRANK-aligned gap-coded partitions had no effect on the final topology; these are identical for the nuclear data set with and without gap partitions (Fig. 7). The inclusion of gap partitions resulted in lowered likelihood scores and the associated need for more intense searches because a new partition was added, but in the absence of topological differences, the likelihood score of the nuclear tree inferred using PRANK-aligned sequences is higher than the manually aligned tree (Appendix 8). Gap-coded partitions contributed positive support in 21 nodes in the combined data set, no support in 44 nodes, and showed conflict in 6 nodes (Appendix 9). Comparing manual and PRANK-aligned partitions, these ratios (positive/zero/negative support) are 21/44/6 and 21/40/10, respectively. These contributions were independent of the differences in the number of parsimony informative sites in both data sets. The total number of parsimony informative characters gained after coding the indel characters from manual and PRANK alignments is 82 and 94, respectively, suggesting that the overall contribution of the gap partitions in the context of the remaining nuclear and the mitochondrial data is surprisingly important (Table 2). With eventual modeling of these kinds of mutations (which has begun; Holmes, 2005), coded indel partitions should also improve estimates of branch lengths as noncoding sequences become more common in phylogenetic studies.

#### Phylogenetic and Evolutionary Implications

Within the genus *Microlophus* both the *Occipitalis* (Dixon and Wright, 1975; Frost, 1992) and *Peruvianus* groups (Van Denburgh and Slevin, 1913; Frost, 1992; Heise, 1998) are only moderately supported by morphological data, and previous molecular studies have not included all species, thereby precluding rigorous independent tests of the monophyly of each clade. In this study, nuclear and the combined data sets always recover the *Occipitalis* and *Peruvianus* groups with strong statistical support, and this support is consistently spread across many data partitions. The mtDNA partition corroborated the nDNA data in strong support of the *Occipitalis* group, but it failed to recover the *Peruvianus* group; one explanation for this result is that the mtDNA data set

may show base compositional bias, so some terminals would be attracted to branches sharing similar base frequencies rather than ancestry (Wiens and Hollingsworth, 2000). We reevaluated the mtDNA locus using a log-determinant method to accommodate this bias (Jordan and Hewitt, 2004) and recovered a well-supported monophyletic *Peruvianus* group that places *M. thoracicus* as the sister species to all others in the group (tree not shown). This result suggests that compositional bias partially accounts for failure of the mtDNA locus to corroborate the nuclear genes in support of the *Peruvianus* group.

Elsewhere, the resolution of relationships within *Peruvianus* and *Occipitalis* groups depends on the combined effect of all markers. Resolution of relationships among the 12 species recognized within the *Occipitalis* group corroborates earlier studies supporting two colonizations of the Galápagos (Lopez et al., 1992; Wright, 1983; Heise, 1998; Kizirian et al., 2004). Specifically, there is strong support for a small Eastern Galápagos Radiation consisting of the sister species *M. bivittatus* and *M. habeli* endemic to the islands of San Cristobal and Marchena, respectively, and sharing a sister clade relationship to the mainland *M. occipitalis*. Our results also confirm monophyly of the “Western Galápagos Radiation” that includes the remaining seven species, and the absence of a close relationship of this radiation to any individual mainland species. The mtDNA data set shows conflict between genes for all nodes that resolve interspecific (interisland) relationships of the Western Radiation, and the source of this conflict must reside in homoplasy, insufficient sampling, or other confounding factors that are not fully accommodated by our methods. We accept the combined data topology (node 42; Fig. 8) as the best available working hypothesis for colonization of the Galápagos Archipelago and are conducting studies to test alternative interisland colonization hypotheses for the Western Radiation.

#### Nuclear-Cytoplasmic Conflict: Secondary Contact and Hybridization?

The combined tree resolves many nodes within the *Peruvianus* group with strong support, and the nDNA recover a pectinate topology for many terminals in the Chilean clade (see node 5; Fig. 7), albeit some supported only by high PP values. The mtDNA gene tree is symmetrical for these taxa, and although not all nodes are in strong conflict by the criterion of Wiens (1998), the nuclear tree is qualitatively improved in the sense that its topology better agrees with both currently recognized species for the *Peruvianus* group and their distributions (neither of which holds for the mtDNA tree). The nuclear gene tree match to species boundaries is not exact, and several processes likely have contributed to the discordance. In Figure 9 we have juxtaposed the relevant clades from the mtDNA and nDNA trees and mapped the distributions of relevant terminals on this section of the coastal desert of western South America.

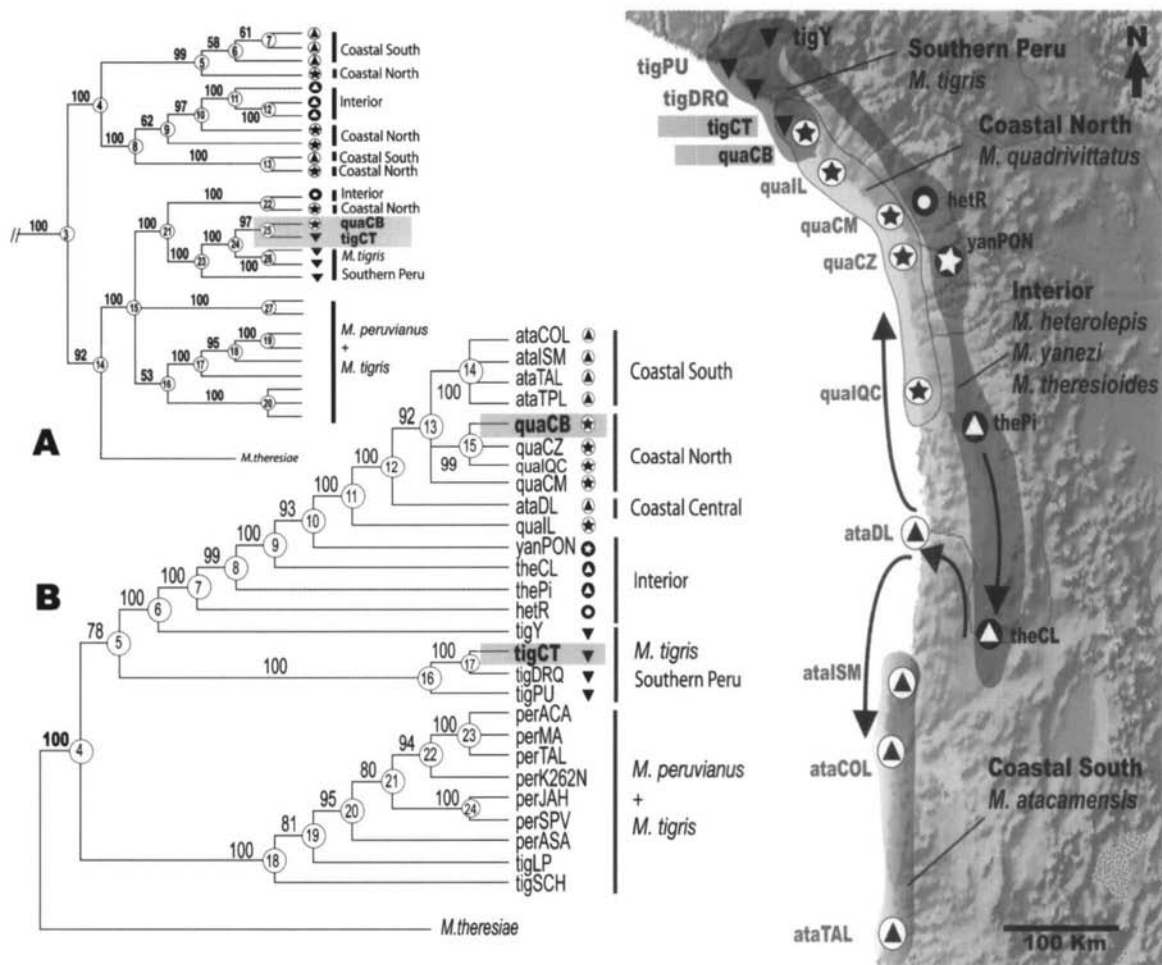


FIGURE 9. Conflicting mitochondrial (tree A depicting node 3 in Fig. 6) and nuclear DNA (tree B depicting node 4 in Fig. 7) topologies (node numbers and PP values as in the original figures). The map enlarges the southern Peru–northern Chile area and shows approximate distributions of putative species in this region (symbols are as in Fig. 1). Arrows and shaded areas illustrate a geographically based hypothesis of evolutionary relationships among interior and coastal terminals implied by the nDNA topology (with caveats). Note the sister-group relationship between populations of *M. quadrivittatus* and *M. tigris* at the shaded node (25) in the mtDNA topology (A) versus support for an alternative placement of the same terminals in the nDNA topology (shaded nodes 15 and 17 in B); the inferred secondary contact zone is shown by the juxtaposed circled star (quaCB) and solid inverted triangle (tigCT) in southern Peru. Symbols are as in Figure 1.

Examination of these topologies suggests that one cause for this conflict could be mtDNA introgression between *M. quadrivittatus* from Caleta Ballenita and *M. tigris* from Caleta Meca (terminals quaCB and tigCT in tree B [nuclear] in Fig. 9). These species are sympatric and marginally segregated by habitat (specimens were collected < 0.5 km apart), and secondary contact followed by mtDNA introgression could explain why *M. quadrivittatus* is recovered with a geographically overlapping terminal of *M. tigris*, yet remains a clearly distinct clade that more closely matches current species boundaries in the nDNA tree. Three nuclear markers (Cryba, CK, Gapdh) support a phylogenetic pattern of differentiation of the Chilean *Microlophus* (nodes 12 to 15 in Fig. 9) that occurred via derivation from an inland ancestor shared with the paraphyletic *M. tigris*–*M. heterolepis*–*M. yanezi*–*M. theresioides* group (nodes 6 to 11). The nDNA tree further implies that *M. theresioides* diverged via dispersal through the Rio Loa valley to the coast and then contin-

ued south as *M. atacamensis* and north as *M. quadrivittatus* until overlapping and hybridizing with *M. tigris* in southern Peru.

We suggest that the combined data topology of the Chilean *Microlophus* clade is both misleading due to secondary contact and mtDNA introgression, and that the match between the nDNA topology and current taxonomy offers the best working hypothesis for this region of the tree. A recent study by Leaché and McGuire (2006) reached a similar conclusion (mtDNA tree was misleading due to introgression) for the lizard genus *Phrynosoma* in North America. The combined data support the mtDNA topology, but some nodes (9 and 10) show relatively high PP values associated with low bootstrap values and short branch lengths, so we cannot rule out the possibility that the high PP values are overestimates (Lewis et al., 2005). However, collapse of the weakly supported nodes in the nuclear tree would not significantly alter the

phylogenetic position of *M. quadrivittatus* relative to *M. tigris*. Paraphyly of some taxa in the nDNA tree suggests that recent speciation/incomplete lineage sorting, poorly defined species boundaries/hybridization, or some combination of these processes has contributed to the geographic complexity implied by this topology, and investigations to resolve these issues are in progress.

### CONCLUSION

There are challenges in any alignments of large multi-gene data sets, but we suggest that a clear distinction be made among the classes of gene regions used; sequence alignments should reflect mutational patterns, and in the case of introns successful incorporation of length mutations in the alignment process needs to be based on methods tailored to their unique features (Morrison, 2006). Although manual alignments are helpful, we suggest that the PRANK algorithm provides an efficient alternative that effectively distinguishes insertions from deletions and recovers maximum phylogenetic signal from both. We are encouraged by results of this study and are optimistic that the current attention focused on model-based coestimation of alignments and phylogeny (Redelings and Suchard, 2005; Fleissner et al., 2005; Lunter et al., 2005b), and developments in explicit models of length mutations (Metzler, 2003; Miklos et al., 2004; Keightley and Johnson, 2004; Holmes, 2005), will allow full use of this class of markers in multilocus phylogenetic studies.

### ACKNOWLEDGEMENTS

We thank H. M. and H. L. Snell and the staff of the Charles Darwin Research Station and the Galápagos National Park service for logistical support in Galápagos; L. Coloma, M.T. Rodrigues, H. L. Snell, J. C. Ortiz, A. Catenazzi, S. Kelez, and J. Cordoba for tissue loans; F. Torrez, P. Victoriano, R. Grams, and M. Vidal for field work in Chile and Peru; J. Wells, P. Alsbury, and I. Stehmeier for lab assistance at BYU; and A. Loytynoja for useful recommendations for implementing PRANK. Lizards were collected in the field under permits issued by INRENA-Peru (TUPA 0113-2002-AG), Parque Nacional Galápagos Ecuador (PT 7.5 FR 28 to H. L. Snell), and Servicio Agrícola y Ganadero in Chile (Permit 3112 to J. C. Ortiz and M. Vidal). E.B. was supported by various units of BYU, including the Department of Integrative Biology and the M. L. Bean Life Science Museum, mentoring and research fellowships from the Office of Graduate Studies, a Larsen Scholarship, and a B. F. Harrison Scholarship. Additional research support was provided by a graduate research award from the Society of Systematic Biologists, and NSF awards DEB 0132227 (to J. W. Sites, Jr., and D. McClellan), DEB 0309111 (to J.W. Sites, Jr., and E. Benavides), and EF 0334966 (to T. Reeder, M. Kearney, J. Wiens, and J.W. Sites, Jr.). Earlier drafts of this paper were critiqued by L. Johnson, D. Mulcahy, and M. Whiting.

### REFERENCES

Abdo, Z., V. Minin, P. Joyce, and J. Sullivan. 2005. Accounting for uncertainty in the tree topology has little effect on the decision theoretic approach to model selection in phylogeny estimation. *Mol. Biol. Evol.* 22:691–703.

Alfaro, M. E., S. Zoller, and F. Lutzoni. 2003. Bayes or bootstrap? A simulation study comparing the performance of Bayesian Markov chain Monte Carlo sampling and bootstrapping in assessing phylogenetic confidence. *Mol. Biol. Evol.* 20:255–266.

Andolfatto, P. 2005. Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* 437:1149–1152.

Arévalo, E., S. K. Davis, and J. W. Sites, Jr. 1994. Mitochondrial DNA sequence divergence and phylogenetic relationships among eight chromosome races of the *Sceloporus grammicus* complex (Sauria, Phrynosomatidae) in central Mexico. *Syst. Biol.* 43:387–418.

Belshaw, R., and D. Bensasson. 2006. The rise and fall of introns. *Heredity* 96:208–213.

Blackshields, G., I. M. Wallace, M. Larkin, and D. G. Higgins. 2006. Analysis and comparison of benchmarks for multiple sequence alignments. *In Silico Biol.* 6:0030. ([/isb/2006/06/0030/](http://isb/2006/06/0030/)).

Blanchette, M., W. J. Kent, C. Riemer, L. Elnitski, A. F. A. Smit, K. M. Roskin, R. Baertsch, K. Rosenbloom, H. Clawson, E. D. Green, D. Haussler, and W. Miller. 2004. Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res.* 14:708–715.

Brandley, M. C., A. Schmitz, and T. W. Reeder. 2005. Partitioned Bayesian analyses, partition choice, and the phylogenetic relationships of scincid lizards. *Syst. Biol.* 54:373–390.

Britten, R. J., L. Rowen, J. Williams, and R. A. Cameron. 2003. Majority of divergence between closely related DNA samples is due to indels. *Proc. Natl. Acad. Sci. USA* 100:4661–4665.

Brudno, M., M. Chapman, B. Göttens, S. Batzoglou, and B. Morgenstern. 2003. Fast and sensitive alignment of large genomic sequences. *BMC Bioinformatics* 4:66.

Cannone, J. J., S. Subramanian, N. M. Schnare, J. R. Collett L. M. D'Souza, Y. Du, B. Feng, N. Lin, L. V. Madabusi, K. M. Muller, N. Pande, Z. Shang, N. Yu, and R. R. Gutell. 2002. The comparative RNA web (CRW) site: An online database of comparative sequence and structure information for ribosomal, intron, and other RNAs. *BMC Bioinformatics* 3:15.

Castoe, T.A., T. M. Doan, and C. L. Parkinson. 2004. Data partitions and complex models in Bayesian analysis: The phylogeny of gymnophthalmid lizards. *Syst. Biol.* 53:448–459.

Cline, M., R. Hugheyook, and K. Karplus. 2002. Predicting reliable regions in protein sequence alignments. *Bioinformatics* 18:306–314.

Creer, S., C. E. Pook, A. Malhotra, and R. S. Thorpe. 2006. Optimal intron analyses in the *Trimeresurus* radiation of Asian pitvipers. *Syst. Biol.* 55:57–72.

Davis, J. I., and K. C. Nixon. 1992. Populations, genetic variation, and the delimitation of phylogenetic species. *Syst. Biol.* 41:421–435.

DeBry, R. W. 2001. Improving the interpretation of the decay index for DNA sequence data. *Syst. Biol.* 50:742–752.

De Pinna, M. C. C. 1991. Concepts and tests of homology in the cladistic paradigm. *Cladistics* 7:367–394.

De Queiroz, K. 1996. Including the characters of interest during tree reconstruction and the problems of circularity and bias in studies of character evolution. *Am. Nat.* 148:700–708.

Dixon, J. R., and J. W. Wright. 1975. A review of the lizards of the iguanid genus *Tropidurus* in Peru. Contributions in Science, Natural History Museum of Los Angeles County 271:1–39.

Do, C. B., M. S. P., Mahabhashyam, M. Brudno, and S. Batzoglou. 2005. PROBCONS: Probabilistic consistency-based multiple sequence alignment. *Genome Res.* 15:330–340.

Dolman, G., and B. Phillips. 2004. Single copy nuclear DNA markers characterized for comparative phylogeography in Australian wet tropics rainforest skinks. *Mol. Ecol. Notes* 4:185–187.

Dolman, G., and C. Moritz. 2006. A multi-locus perspective on refugial isolation and divergence in rainforest skinks (*Carlia*). *Evolution* 60:573–582.

Doolittle, W. F. 1999. Phylogenetic classification and the universal tree. *Science* 284:2124–2129.

Edgar, R. C. 2004. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32:1792–1797.

Edgar, R. C., and K. Sjölander. 2003. SATCHMO: Sequence alignment and tree construction using hidden Markov models. *Bioinformatics* 19:1404–1411.

Edgar, R. C., and S. Batzoglou. 2006. Multiple sequence alignment. *Curr. Opin. Struct. Biol.* 16:368–373.

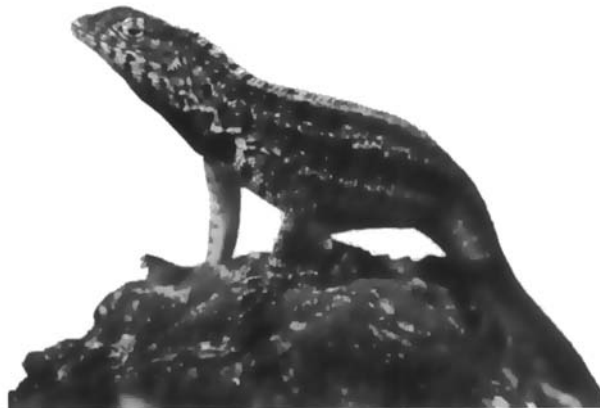
Fetzner, J. 1999. Extracting high-quality DNA from shed reptiles skins: A simplified method. *BioTechniques* 26:1052–1054.

Fleissner, R., D. Metzler, and A. von Haeseler. 2005. Simultaneous statistical multiple alignment and phylogeny reconstruction. *Syst. Biol.* 54:548–561.

- Friesen, V. L., B. C. Congdon, M. G. Kidd, and T. P. Bird. 1999. Polymerase chain reaction (PCR) primers for the amplification of five nuclear introns in vertebrates. *Mol. Ecol.* 8:2147–2149.
- Friesen V. L., B. C. Congdon, H. E. Walsh, and T. P. Birt. 1997. Intron variation in marbled murrelets detected using analyses of single-stranded conformational polymorphisms. *Mol. Ecol.* 6:1047–1058.
- Frost, D. 1992. Phylogenetic analysis and taxonomy of iguanian lizards (Reptilia: Squamata). *Am. Mus. Novitates* 3033:1–68.
- Frost, D. R., M. T. Rodrigues, T. Grant, and T. A. Titus. 2001. Phylogenetics of the lizard genus *Tropidurus* (Squamata: Tropiduridae: Tropidurinae): Direct optimization, descriptive efficiency, and sensitivity analysis of congruence between molecular data and morphology. *Mol. Phylogenet. Evol.* 21:352–371.
- Gatesy, J., and R. H. Baker. 2005. Hidden likelihood support in genomic data: Can forty-five wrongs make a right? *Syst. Biol.* 54:483–492.
- Golenberg, E. M., M. T. Clegg, M. L. Durbin, J. Doebly, and D. P. Ma. 1993. Evolution of a non-coding region of the chloroplast genome. *Mol. Phylogenet. Evol.* 2:52–64.
- Graham, S. W., P. A. Reeves, A. C. E. Burns, and R. G. Olmstead. 2000. Microstructural changes in noncoding chloroplasts DNA: Interpretation evolution and utility of indels and inversions in basal angiosperm phylogenetic inference. *Int. J. Plant. Sci.* 161:583–596.
- Gu, X., and W. -H. Li. 1995. The size distribution of insertions and deletions in human and rodent pseudogenes suggests the logarithmic gap penalty for sequence alignment. *J. Mol. Evol.* 40:464–473.
- Guindon, S., and O. Gascuel. 2003. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst. Biol.* 52:696–704.
- Gutell, R. R., and G. E. Fox. 1988. A compilation of large subunit RNA sequences presented in structural format. *Nucleic Acids Res.* 16:r175–r269.
- Harvey, M. B., and R. L. Gutberlet, Jr. 2000. A phylogenetic analysis of tropidurine lizards (Squamata: Tropiduridae), including new characters of squamation and epidermal microstructure. *Zool. J. Linn. Soc.* 128:189–233.
- Heise, P. J. 1998. Phylogeny and biogeography of Galápagos lava lizards (*Microlophus*) inferred from nucleotide sequence variation in mitochondrial DNA. A dissertation presented for the Doctor of Philosophy degree. The University of Tennessee, Knoxville, Tennessee.
- Higgins, D. G., G. Blackshields, and I. M. Wallace. 2005. Mind the gaps: Progress in progressive alignment. *Proc. Natl. Acad. Sci. USA* 102:10411–10412.
- Hillis, D. M., and J. J. Bull. 1993. An empirical testing of bootstrapping as a method for assessing confidence in phylogenetic analysis. *Syst. Biol.* 42:182–192.
- Holmes, I. 2005. Using evolutionary expectation maximization to estimate indel rates. *Bioinformatics* 21:2294–2300.
- Janke, A., D. Erpenbeck, M. Nilsson, and U. Arnason. 2001. The mitochondrial genomes of the iguana (*Iguana iguana*) and the caiman (*Caiman crocodylus*): Implications for amniote phylogeny. *Proc. Biol. Sci.* 268:623–631.
- Jordan, B. H., and G. M. Hewitt. 2004. The origin and radiation of Macaronesian beetles breeding in *Euphorbia*: The relative importance of multiple data partitions and population sampling. *Syst. Biol.* 53:711–734.
- Katoh, K., K. Misawa, K. Kuma, and T. Miyata. 2002. MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30:3059–3066.
- Keightley, P. D., and T. Johnson. 2004. MCALIGN: Stochastic alignment of non-coding DNA sequences based on an evolutionary model of sequence evolution. *Genome Res.* 14:442–450.
- Kelchner, S. A. 2000. The evolution of non-coding chloroplast DNA and its application in plant systematics. *Ann. Missouri Bot. Gard.* 87:482–498.
- Kjer, K. M., Gillespie, J. J., and K. A. Ober. 2007. Opinions on multiple sequence alignment, and an empirical comparison of repeatability and accuracy between POY and structural alignment. *Syst. Biol.* 56:133–146.
- Kizirian, D., A. Trager, M. A. Donnelly, and J. W. Wright. 2004. Evolution of Galápagos Island lizards (Iguania: Tropiduridae: *Microlophus*). *Mol. Phylogenet. Evol.* 32:761–769.
- Konagurthu, A. S., and P. J. Stuckey. 2006. Optimal sum-of-pairs multiple sequence alignment using incremental Carrillo and Lipman bounds. *J. Comput. Biol.* 13:668–685.
- Kumar, S., and A. Filipinski. 2007. Multiple sequence alignment: In pursuit of homologous DNA positions. *Genome Res.* 17:127–135.
- Kumazawa, Y., and M. Nishida. 1993. Sequence evolution of mitochondrial tRNA genes and deep branch animal phylogenetics. *J. Mol. Evol.* 37:380–398.
- Lambert, C., J. M. Campenhout, X. DeBolle, and E. Depiereux. 2003. Review of common sequence alignment methods: Clues to enhance reliability. *Curr. Genomics* 4:131–146.
- Lassmann, T., and E. L. L. Sonnhammer. 2005. Automatic alignment assessment of alignment quality. *Nucleic Acids Res.* 33:7120–7128.
- Leaché, A. D., and J. A. McGuire. 2006. Phylogenetic relationships of horned lizards (*Phrynosoma*) based on nuclear and mitochondrial data: Evidence for a misleading mitochondrial gene tree. *Mol. Phylogenet. Evol.* 39:628–644.
- Lebrun, E., J. M. Santini, M. Brugna, A. -L. Ducluzeau, S. Ouchane, B. Scoepp-Cothenet, F. Baymann, and W. Nitschke. 2006. The Rieske protein: A case study on the pitfalls of multiple sequence alignments and phylogenetic reconstruction. *Mol. Biol. Evol.* 23:1180–1190.
- Lee, C., C. Grasso, and M. F. Sharlow. Multiple sequence alignment using partial order graphs. *Bioinformatics* 18:425–464.
- Levinson, G., and G. A. Gutman. 1987. Slipped-strand mispairing: A major mechanism for DNA sequence evolution. *Mol. Bio. Evol.* 4:203–221.
- Lewis, P. O. 2001. A likelihood approach to estimating phylogeny from discrete morphological character data. *Syst. Biol.* 50:913–925.
- Lewis, P. O., M. T. Holder, and K. E. Holsinger. 2005. Polytomies and Bayesian and phylogenetic inference. *Syst. Biol.* 54:241–253.
- Li, W.-H. 1997. Molecular evolution. Sinauer Associates, Sunderland, Massachusetts.
- Lohne, C., and T. Borsch. 2005. Molecular evolution and phylogenetic utility of the petD group II intron: A case study in basal angiosperms. *Mol. Biol. Evol.* 22:317–322.
- Lopez, T. J., E. D. Hauselman, L. R. Maxson, and J. W. Wright. 1992. Preliminary analysis of phylogenetic relationships among Galápagos Island lizards of the genus *Tropidurus*. *Amphibia-Reptilia* 13:327–339.
- Loytynoja, A., and N. Goldman. 2005. An algorithm for progressive multiple alignment of sequences with insertions. *Proc. Natl. Acad. Sci. USA* 102:10557–10562.
- Loytynoja, A., and M. C. Milinkovitch. 2001. SOAP. Cleaning multiple alignments from unstable blocks. *Bioinformatics* 17:277–286.
- Loytynoja, A., and M. C. Milinkovitch. 2003. ProAlign, a probabilistic multiple alignment. *Bioinformatics* 19:1505–1513.
- Lunter, G., A. J. Drummond, I. Miklos, and J. Hein. 2005a. Statistical alignment: Recent progress, new applications, and challenges. Pages 382–411 in *Statistical methods in molecular evolution* (R. Nielsen, ed.). Springer, New York.
- Lunter, G., I. M. Miklos, A. Drummond, J. L. Jensen, and J. Hein. 2005b. Bayesian coestimation of phylogeny and sequence alignment. *BMC Bioinformatics* 6:83.
- Metzler, D. 2003. Statistical alignment based on fragment insertion and deletion models. *Bioinformatics* 19:490–499.
- Miklos, I., G. A. Lunter, and I. Holmes. 2004. A “long-indel” model for evolutionary sequence alignment. *Mol. Biol. Evol.* 21:529–540.
- Minin, V., Z. Abdo, P. Joyce, and J. Sullivan. 2003. Performance-based selection of likelihood models for phylogeny estimation. *Syst. Biol.* 52:674–683.
- Morgenstern, B. 1999. DIALIGN 2: Improvement of the segment-to-segment approach to multiple sequence alignment. *Bioinformatics* 15:211–218.
- Morrison, D. A. 2006. Multiple sequence alignment for phylogenetic purposes. *Austral. Syst. Bot.* 19:479–539.
- Notredame, C., D. G. Higgins, and J. Heringa. 2000. T-coffee: A novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol.* 302:205–217.
- Ogden, T. H., and M. S. Rosenberg. 2006. Multiple sequence alignment accuracy and phylogenetic inference. *Syst. Biol.* 55:314–328.
- Pagel, M., and A. Meade. 2004. A phylogenetic mixture model for detecting pattern heterogeneity in gene sequence or character state data. *Syst. Biol.* 53:571–581.

- Palumbi, S. R. 1996. Nucleic acids I: The polymerase chain reaction. Pages 205–247 in *Molecular systematics* (D. M. Hillis, C. Moritz, and B. K. Mable, eds.). Sinauer Associates, Sunderland, Massachusetts.
- Phillips, M. J., F. Delsuc, and D. Penny. 2004. Genomescale phylogeny and the detection of systematic biases. *Mol. Biol. Evol.* 21:1455–1458.
- Pollard, D. A., C. M. Bergman, J. Stoye, S. E. Celniker, and M. B. Eisen. 2004. Benchmarking tools for the alignment of functional noncoding DNA. *BMC Bioinformatics* 5:6.
- Rambaut, A. 1996. Se-AL: Sequence alignment editor. Available at <http://evolve.zoo.ox.ac.uk/>.
- Rambaut, A., and A. J. Drummond. 2003. Tracer v1.2, <http://evolve.zoo.ox.ac.uk/>.
- Raymond, J., O. Zhaxybayeva, J. P. Gogarten, S. Y. Gerdes, and R. E. Blankenship. 2002. Whole-genome analysis of photosynthetic prokaryotes. *Science* 298:1614–1619.
- Redelings, B. D., and M. A. Suchard. 2005. Joint Bayesian estimation of alignment and phylogeny. *Syst. Biol.* 54:401–418.
- Rokas, A., D. Kruger, and S. B. Carroll. 2005. Animal evolution and the molecular signature of radiations compressed in time. *Science* 310:1933–1938.
- Ronquist, F., and J. P. Huelsenbeck. 2003. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* 19:1572–1574.
- Roy, S. W., and W. Gilbert. 2006. The evolution of spliceosomal introns: Patterns, puzzles, and progress. *Nat. Rev. Genet.* 7:211–221.
- Saint, K. M., C. C. Austin, S. C. Donnellann, and M. N. Hutchinson. 1998. C-mos, a nuclear marker useful for squamate phylogenetic analysis. *Mol. Phylogenet. Evol.* 3:240–247.
- Sanchis, A., J. M. Michelena, A. Latorre, D. L. J. Quicke, U. Gardenfors, and R. Belshaw. 2001. The phylogenetic analysis of variable-length sequence data: Elongation factor-1 $\alpha$  introns in European populations of the parasitoid wasp genus *Pauesia* (Hymenoptera: Braconidae: Aphidiinae). *Mol. Biol. Evol.* 18:1117–1131.
- Shaw, K. L. 2002. Conflict between nuclear and mitochondrial DNA phylogenies of a recent species radiation: What mtDNA reveals and conceals about modes of speciation in Hawaiian crickets. *Proc. Natl. Acad. Sci. USA* 99:16122–16127.
- Siddharthan, R. 2006. Sigma: Multiple alignment of weakly-conserved non-coding DNA sequence. *BMC Bioinformatics* 7:143.
- Simmons, M. P., and H. Ochoterena. 2000. Gaps as characters in sequence-based phylogenetic analyses. *Syst. Biol.* 49:369–381.
- Smythe, A. B., M. J. Sanderson, and S. A. Nadler. 2007. Nematode small subunit phylogeny correlates with alignment parameters. *Syst. Biol.* 55:972–992.
- Sorenson, M. D. 1999. TreeRot, version 2. Boston University, Boston, Massachusetts.
- Sullivan, J., and P. Joyce. 2005. Model selection in phylogenetics. *Ann. Rev. Ecol. Syst.* 36:445–466.
- Swofford, D. L. 2002. PAUP\*: Phylogenetic analysis using parsimony (\*and other methods), beta version 4.0.b5b. Sinauer, Sunderland, Massachusetts.
- Terry, M. D., and M. F. Whiting. 2005. Mantophasmatodea and phylogeny of the lower neopterous insects. *Cladistics* 21:240–257.
- Thompson, J. D., D. G. Higgins, and T. J. Gibson. 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, J. L., H. Kishino, and J. Felsenstein. 1992. Inching towards reality: An improved likelihood model of sequence evolution. *J. Mol. Evol.* 34:3–16.
- Van Denburgh, J., and J. R. Slevin. 1913. Expedition of the California Academy of Sciences to the Galápagos Islands, 1905–1906. IX. The Galapagoan lizards of the genus *Tropidurus* with notes on iguanas of the genera *Conolophus* and *Amblyrhynchus*. *Proc. Calif. Acad. Sci. Ser.* 42:132–202.
- Van de Peer, Y., I. Van den Broek, P. de Rijk, and R. de Wachter. 1994. Database on the structure of small ribosomal subunit RNA. *Nucleic Acids Res.* 22:3488–3494.
- Wallace, I. A., O. O'Sullivan, and D. G. Higgins. 2005. Evaluation of iterative alignment algorithms for multiple alignments. *Bioinformatics* 21:1408–1414.
- Whiting, A. S., A. M. Bauer, and J. W. Sites, Jr. 2003. Phylogenetic relationships and limb loss in sub-Saharan South African scincine lizards (Squamata: Scincidae). *Mol. Phylogenet. Evol.* 3:582–592.
- Whiting, A. S., J. W. Sites, Jr., K. C. M. Pellegrino, and M. T. Rodrigues. 2006. Comparing alignment methods for inferring the history of the New World lizards genus *Mabuya* (Squamata: Scincidae). *Mol. Phylogenet. Evol.* 38:719–730.
- Wiens, J. J. 1998. Combining data sets with different phylogenetic histories. *Syst. Biol.* 47:568–581.
- Wiens, J. J., and B. D. Hollingsworth. 2000. War of the iguanas: Conflicting molecular and morphological phylogenies and long-branch attraction. *Syst. Biol.* 49:143–159.
- Wiens, J. J., and T. W. Reeder. 1997. Phylogeny of the spiny lizards (*Sceloporus*) based on molecular and morphological evidence. *Herpetol. Monogr.* 11:1–101.
- Wiens, J. J., T. W. Reeder, and A. Nieto Montes de Oca. 1999. Molecular phylogenetics and evolution of sexual dichromatism among populations of the Yarrow's spiny lizard (*Sceloporus jarrovi*). *Evolution* 53:1884–1897.
- Wright, J. W. 1983. The evolution and biogeography of the lizards of the Galápagos Archipelago: Evolutionary genetics of *Phyllodactylus* and *Tropidurus* populations. Pages 123–155 in *Patterns of evolution in Galápagos organisms* (R. I. Bowman, M. Berson, and A. E. Levinton, eds.). AAAS Symposium Volume, San Francisco, California.

First submitted 9 June 2006; reviews returned 26 August 2006;  
final acceptance 5 July 2007  
Associate Editor: Karl Kjer



*Microlophus indefatigabilis* from the islet of Plaza Sur near the Island of Santa Cruz. *M. indefatigabilis* is part of the Western Radiation that represents one of two independent colonizations of the Galapagos Islands from the South American mainland (Photo taken by Heidi Snell). ©Heidi Snell/Visual Escapes