

Phylogenetic Utility of Different Types of Molecular Data Used to Infer Evolutionary Relationships among Stalk-Eyed Flies (Diopsidae)

RICHARD H. BAKER,^{1,2,4} GERALD S. WILKINSON,³ AND ROB DE SALLE¹

¹*Department of Entomology, American Museum of Natural History, 79th Street at Central Park West, New York, NY 10024, USA; E-mail: desalle@amnh.org*

²*Department of Biology, Yale University, New Haven, CT 06520, USA*

³*Department of Zoology, University of Maryland, College Park, MD 21742, USA; E-mail: gw10@umail.umd.edu*

Abstract.—A phylogenetic hypothesis of relationships among 33 species of stalk-eyed flies was generated from a molecular data set comprising three mitochondrial and three nuclear gene regions. A combined analysis of all the data equally weighted produced a single most-parsimonious cladogram with relatively strong support at the majority of nodes. The phylogenetic utility of different classes of molecular data was also examined. In particular, using a number of different measures of utility in both a combined and separate analysis framework, we focused on the distinction between mitochondrial and nuclear genes and between faster-evolving characters and slower-evolving characters. For the first comparison, by nearly any measure of utility, the nuclear genes are substantially more informative for resolving diopsid relationships than are the mitochondrial genes. The nuclear genes exhibit less homoplasy, are less incongruent with one another and with the combined data, and contribute more support to the combined analysis topology than do the mitochondrial genes. Results from the second comparison, however, provide little evidence of a clear difference in utility. Despite indications of rapid divergence and saturation, faster-evolving characters in both the nuclear and mitochondrial data sets still provide substantial phylogenetic signal. In general, inclusion of the more rapidly evolving data consistently improves the congruence among partitions. [Diopsidae; incongruence; nuclear genes; partitioned Bremer support; saturation.]

Flies in the family Diopsidae are characterized by the elongation of the head into long stalks, with the eyes and antennae laterally displaced at the ends of these stalks. In several species, this elongation is so extreme that the length of their eye stalks exceeds the length of their body. Several dipteran families possess hypercephalic species (Grimaldi and Fenster, 1989; Wilkinson and Dodson, 1997), but the Diopsidae are unique in that both males and females of all the species within the family have some degree of head modification. Recently, experiments examining the importance of eye stalks in the mating system of diopsids have provided considerable information about their adaptive significance (Burkhardt and de la Motte, 1985, 1988; Lorch et al., 1993; Wilkinson, 1993; Burkhardt et al., 1994; Wilkinson and Reillo, 1994; Wilkinson and Dodson, 1997). For the majority of diopsid species, males have markedly larger eye stalks than females, and several studies have demonstrated that this increased eye span functions as an ornament for both male combat (Burkhardt and de la Motte, 1983, 1987; Lorch et al., 1993) and

female mate choice (Burkhardt and de la Motte, 1988; Burkhardt et al., 1994; Wilkinson and Reillo, 1994). There is also considerable variation among species in both the size of the eye stalks and the extent of the sexual dimorphism. Given this interspecific variation, understanding of the forces influencing the evolution of eye stalks would be enhanced by comparative information concerning the pattern of morphological change within the family. This type of analysis requires a robust phylogenetic hypothesis of the evolutionary relationships among these flies. Therefore, we conducted a systematic study of several diopsid species, using molecular sequence data from six different gene regions. Various comparative analyses utilizing this phylogenetic framework will be presented elsewhere (Presgraves et al., 1999; Baker and Wilkinson, in prep.).

In addition, the quality of the molecular data set allows for a detailed examination of the phylogenetic utility of different classes of molecular data. As the quantity and diversity of molecular data have expanded in systematic studies, there has been increased emphasis on discriminating among different types of character information and accommodating these differences into the tree-building process (Fitch and Ye, 1991; Knight

⁴Present address (and address for correspondence): The Galton Laboratory, University College London, 4 Stephenson Way, London NW1 2HE, U.K.

and Mindell, 1993, 1995; Chippendale and Wiens, 1994; Collins et al., 1994; Yang, 1996). This accommodation ranges from down-weighting characters to separating heterogeneous data, but, in all cases, follows from the principle that some sources of information are more reliable than others and that the more reliable information can be identified as such. In practice, the reliability or utility of systematic data is either assessed directly from the data at hand or assumed to be reasonable, given the results of previous studies. The methods for quantifying utility are varied and include (1) examining properties of the data before phylogenetic analysis (e.g., transition/transversion ratios, saturation curves, and compositional bias) (Friedlander et al., 1994; Graybeal, 1994), (2) examining properties of the data specific to a given phylogenetic analysis (e.g., resolution, homoplasy measures, and support measures) (Milinkovitch et al., 1996; Naylor and Brown, 1998), (3) comparing the results with those for a "known" or well-supported phylogeny (e.g., simulation studies and higher-level vertebrate phylogenies) (Friedlander et al., 1994; Graybeal, 1994; Russo et al., 1996; Cunningham, 1997a, 1997b), and (4) assessing the congruence among comparable data types from different sources (e.g., transversion data from all of the mitochondrial genes) (Miyamoto et al., 1994; Friedlander et al., 1998; Allard et al., 1999). Despite the ubiquitous use of these measures, little attention has been paid to whether or not, when compared with each other, they provide consistent or conflicting information about the nature of data (but see Davis et al., 1998, for a recent example). In this study, by calculating and comparing several different measures, we have comprehensively examined both the utility of different data types and the equivocity with which this utility can be identified. In general, we focused on two character data dichotomies: nuclear genes versus mitochondrial genes and faster-evolving characters versus slower-evolving characters.

The majority of molecular sequence data for animal systematic studies has traditionally come from either mitochondrial genes or nuclear ribosomal genes. This bias stems primarily from the ease with which these genes can be amplified relative to the nuclear protein-coding genes. Recently, however, the use of nuclear protein genes in systematic studies has received increased

emphasis (Friedlander et al., 1994, 1996, 1998; Gatesy et al., 1996; Fang et al., 1997; Ramos-Onsins et al., 1998). As this type of molecular information becomes more prominent, its utility relative to other types of data must be examined. The data matrix constructed for diopsid flies contains three mitochondrial genes and three nuclear protein-coding genes and therefore provides an excellent opportunity for this kind of analysis.

Concern about the phylogenetic effects of rapidly evolving nucleotide sites has become one of the most prominent issues in molecular systematics. Since Felsenstein's (1978) seminal paper on the consistency of molecular data, numerous studies have focused on the relationship between the rate of molecular evolution and the accuracy of phylogenetic reconstruction (Fitch and Ye, 1991; Hillis, 1991; Bull et al., 1993; Mindell and Thacker, 1996; Yang, 1996). These studies have generally concluded that, because multiple substitutions potentially mask historical patterns of grouping, rapidly evolving characters have little phylogenetic utility, such that their inclusion in an equally weighted combined matrix will likely confound the results. The majority of empirical studies make at least some attempt to identify the more variable characters, and reducing their influence on the final hypothesis has become commonplace (Mindell et al., 1996; Murphy and Collier, 1997; Bloomer and Crowe, 1998; Danforth and Ji, 1998; Martin and Bermingham, 1998). Although some methods of differential weighting discriminate among characters individually (e.g., successive approximation [Farris, 1969; Carpenter, 1988] and Goloboff weighting [Goloboff, 1993]), most studies downweight an entire class of data or type of transformation (e.g., third positions or transitions) that is deemed to be too fast. Some authors (Nixon and Carpenter, 1996; Kluge, 1997; Miller et al., 1997; Allard et al., 1999; Wenzel and Siddall, 1999) have questioned the validity of this approach, arguing that the justification for choosing among the array of possible a priori differential weightings is rarely strong and that a combined analysis of all the characters equally weighted represents the most severe test of the hypothesis. In addition, empirical work (Philippe et al., 1996; Baker and DeSalle, 1997; Olmstead et al., 1998; Bjorkland, 1999; Källersjö et al., 1999) has

demonstrated that, despite indications of high divergences and saturation, rapidly changing molecular data still contain substantial phylogenetic signal. Overall, the utility of rapidly evolving data types needs to be examined in greater detail with a wide array of measures to determine whether a consistent pattern of unreliability emerges. The diversity of the diopsid data set provides several instances for this type of comparison between faster- and slower-evolving characters.

MATERIAL AND METHODS

Taxa

Currently ~150 described diopsid species are grouped in 13 different genera (Steyskal, 1972; Feijen, 1984, 1989), although 7 of these genera contain only one or a few species. Estimates of undescribed taxa suggest that the eventual number of species within the family will exceed 200 (Feijen, 1989). For our analysis, molecular sequence data were collected for 33 diopsid species and two outgroup taxa (Table 1). The choice of taxa was limited by our ability to obtain fresh specimens. Representatives from each of the six major genera except *Diopsina* were included, although sampling is skewed toward *Diasemopsis*: of the ~40 species currently recognized in *Diasemopsis* (Feijen, 1989), 15 were included in this study, whereas *Diopsis* contains >60 described species, of which only 4 were sampled. Two *Teloglabus* species from the family Centroncidae were used as outgroup taxa. Originally included within the family Diopsidae (Shillito, 1950; Hennig, 1958), these flies were placed in a separate family by Feijen (1983, 1989). In his analysis, a clade that included both the Centroncidae and Syringogastridae was designated as the sister taxon of the Diopsidae. Recent character data from egg morphology (Meier and Hilger, 2000), however, contradict the monophyly of Centroncidae and Syringogastridae and instead place the Centroncidae alone as sister to the Diopsidae.

DNA Extraction and Sequencing

Molecular character information was generated from six different gene regions that included fragments from three mitochondrial genes—cytochrome oxidase II (CO II), 12S ribosomal RNA (12S), and 16S ribosomal RNA (16S)—and three nuclear genes—elongation

factor-1 α (EF-1 α), *wingless*, and *white* (Table 2). After alignment, these data comprised a total of 3236 characters, of which 966 were phylogenetically informative.

Genomic DNA was extracted from single flies according to the preparation outlined by Vogler et al. (1993) but with an additional phenol purification and only a single ethanol-ammonium acetate precipitation. Primers for the various gene regions are listed in Table 2. For EF-1 α and *wingless*, several diopsid-specific internal sequencing primers were made. Polymerase chain reaction (PCR) protocols differed for the various gene fragments. The protocol for the three mitochondrial genes was generally 94°C for 1 min, 50°C for 1 min, and 72°C for 1 min, for a total of 33 cycles. The conditions for EF-1 α were 94°C for 1 min, 55°C for 1.5 min, and 72°C for 2 min, for a total of 35 cycles; for *wingless* and *white* these were 94°C for 1 min, 47°C for 1 min, and 72°C for 1.5 min, for 35 cycles. PCR products were cleaned with GeneClean kits (BIO 101) and were sequenced by using either a manual dideoxy double-stranded protocol with ³⁵S or automated sequencing with a fluorescent dideoxy terminator mix on the ABI 373 and 377 automated sequencers. Automated sequence outputs were imported into Sequencher (Gene Codes Corp., Ann Arbor, MI) for visual inspection of the chromatographs and alignment of the various sequencing reactions for a given taxon and a given gene into a single sequence. Because of the limited quantity and poor quality of the *Teloglabus milleri* genomic DNA, we were unable to obtain sequence data from the *wingless* and *white* gene for this outgroup taxon. Therefore, for several of the nuclear gene analyses, this taxon was excluded.

Data Analysis

Alignment.—The ribosomal sequences were aligned by using MALIGN 2.7 (Wheeler and Gladstein, 1994) and several gap:change costs (2:1, 4:1, 6:1, 8:1, 10:1) to identify regions of alignment stability and instability (Gatesy et al., 1994). The search parameters used included: build, score 4, alignswap, treeswap, keepaligns 100, keeptrees 100, and extragaps 1. Regions of alignment ambiguity among the different alignments were then deleted from the final matrix by using the cull option

TABLE 1. List of 35 taxa used in the study.

Species	Locality ^a	GenBank accession no.					
		COII	12S	16S	EF1-1 α	Wingless	White
<i>Chaetodiopsis meigenii</i> (Westwood)	Nelspruit, South Africa	AF304752	AF304682	AF304717	AF303182	AF304787	AF304821
<i>Cyrtodiopsis curranii</i> (Shillito)	Chaing Mai, Thailand	AF304781	AF304711	AF304746	AF303211	AF304815	AF304849
<i>Cyrtodiopsis dahmani</i> (Wiedemann)	Ulu Gombak, Malaysia	AF304782	AF304712	AF304747	AF303212	AF324429 ^c	AF304850
<i>Cyrtodiopsis quinqueguttata</i> (Walker)	Ulu Gombak, Malaysia	AF304783	AF304713	AF304748	AF303213	AF304817	AF304851
<i>Cyrtodiopsis whitiei</i> (Curran)	Ulu Gombak, Malaysia	AF304784	AF304714	AF304749	AF303214	AF304818	AF304852
<i>Diasemopsis aethiopia</i> (Rondani)	Pietermaritzburg, South Africa	AF304771	AF304701	AF304736	AF303201	AF304805	AF304839
<i>Diasemopsis albifacies</i> (Curran)	Lope, Gabon	AF304758	AF304688	AF304723	AF303188	AF304793	AF304827
<i>Diasemopsis conjuncta</i> (Curran)	Bonjongo, Cameroon	AF304773	AF304703	AF304738	AF303203	AF304807	AF304841
<i>Diasemopsis dubia</i> (Bigot)	Pietermaritzburg, South Africa	AF304750	AF304680	AF304715	AF303180	AF304785	AF304819
<i>Diasemopsis elongata</i> (Curran)	Lamborene, Gabon	AF304775	AF304705	AF304740	AF303205	AF304809	AF304843
<i>Diasemopsis fasciata</i> (Gray)	Limbe, Cameroon	AF304754	AF304684	AF304719	AF303184	AF304789	AF304823
<i>Diasemopsis hirsuta</i> (Curran)	Bonjongo, Cameroon	AF304760	AF304690	AF304725	AF303190	AF304795	AF304829
<i>Diasemopsis longipeditunculata</i> (Curran)	Bomossa, Republic of Congo	AF304776	AF304706	AF304741	AF303206	AF304810	AF304844
<i>Diasemopsis munroi</i> (Curran)	Gillits, South Africa	AF304757	AF304687	AF304722	AF303187	AF304792	AF304826
<i>Diasemopsis nebulosa</i> (Curran)	Mabeta, Cameroon	AF304774	AF304704	AF304739	AF303204	AF304788	AF304822
<i>Diasemopsis obstans</i> (Walker)	Pietermaritzburg, South Africa	AF304753	AF304683	AF304718	AF303183	AF304790	AF304824
<i>Diasemopsis signata</i> (Dalman)	Limbe, Cameroon	AF304755	AF304685	AF304720	AF303185	AF304790	AF304824
<i>Diasemopsis silvatica</i> (Eggers)	Pietermaritzburg, South Africa	AF304751	AF304681	AF304716	AF303181	AF304786	AF304820
<i>Diasemopsis sp.M</i>	Pietermaritzburg, South Africa	AF304756	AF304686	AF304721	AF303186	AF304791	AF304825
<i>Diasemopsis sp.W</i>	Bomossa, Republic of Congo	AF304761	AF304691	AF304726	AF303191	AF304796	AF304830
<i>Diopsis apicalis</i> (Dalman)	Pietermaritzburg, South Africa	AF304777	AF304707	AF304742	AF303207	AF304811	AF304845
<i>Diopsis fumipennis</i> (Westwood)	Pietermaritzburg, South Africa	AF304778	AF304708	AF304743	AF303208	AF304812	AF304846
<i>Diopsis gnu</i> (Hendel)	Mzingazi, South Africa	AF304779	AF304709	AF304744	AF303209	AF304813	AF304847
<i>Diopsis longicornis</i> (Macquart)	M'Be-Bouoko, Ivory Coast	AF304780	AF304710	AF304745	AF303210	AF304814	AF304848
<i>Eurydiopsis argentifera</i> (Bigot)	Ulu Gombak, Malaysia	AF304764	AF304694	AF304729	AF303194	AF304799	AF304833
<i>Sphyracephala beccarii</i> (Rondani)	Rivulets, South Africa	AF304772	AF304702	AF304737	AF303202	AF304806	AF304840
<i>Sphyracephala bipunctipennis</i> (Senior-White)	Gombak River, Malaysia	AF304770	AF304700	AF304735	AF303200	AF324431 ^d	AF304838
<i>Sphyracephala brevicornis</i> (Say)	College Park, Maryland	AF304765	AF304695	AF304730	AF303195	AF304800	AF304834
<i>Sphyracephala munroi</i> (Curran)	Pietermaritzburg, South Africa	AF304762	AF304692	AF304727	AF303192	AF304797	AF304831
<i>Teleopsis breviscriptum</i> (Rondani)	Ulu Langat, Malaysia	AF304763	AF304693	AF304728	AF303193	AF304798	AF324435 ^e
<i>Teleopsis quadriguttata</i> (Walker)	Ulu Gombak, Malaysia	AF304768	AF304698	AF304733	AF303198	AF304802	AF304836
<i>Teleopsis rubicunda</i> (Wulp)	Ulu Gombak, Malaysia	AF304769	AF304699	AF304734	AF303199	AF304803	AF304837
<i>Teleolabrus entabeneis</i> (Feijen) ^b	Pietermaritzburg, South Africa	AF304766	AF304696	AF304731	AF303196	AF324433 ^f	AF304835
<i>Teleolabrus milleri</i> (Feijen) ^b	Pietermaritzburg, South Africa	AF304767	AF304697	AF304732	AF303197		
<i>Trichodiopsis minuta</i> (Séguy)	Lope, Gabon	AF304759	AF304689	AF304724	AF303189	AF304794	AF304828

^aVoucher specimens are archived at AMNH molecular lab # nos. 942-1 through 942-35. Identifications were made by Richard Baker, Sabine Hilger, and Marion Kotrba.

^bOutgroup taxa.

^cand AF324430.

^dand AF324432.

^eand AF324434.

^fand AF324436.

TABLE 2. Primers used to sequence the various gene regions. All positional information refers to nucleotide location relative to *Drosophila melanogaster* sequence. The direction of each primer is provided with respect to the sense (S) or antisense (A) strands.

Gene region	Primer name or sequence	Reference or position	
COII	A3772 ^a	Brower, 1994	
	S3291 ^a		
	A3661		
12S	12Sai ^a	Simon et al., 1994	
	12Sbi ^a		
16S	5'-AATTTATTGCACTAATCTGCC-3' ^a	12727-12747 (S)	
	5'-GCTGGAATGAATGGTTGGACG-3' ^a	13270-13290 (A)	
	5'-TATAATTTTGGGTGTAGCCG-3'	12923-12942 (S)	
	5'-TAATCCAACATCGAGGTCGC-3'	12946-12965 (A)	
EF-1 α ^b	M44-1 ^a	Cho et al., 1995	
	M64-1		
	rcM53-2		
	rcM4 ^a		
	5'-TATYGCTTTRTGGAAATTCG-3' ^c		2284-2303 (S)
	5'-CTTGCTTTCACHTTGGGTG-3' ^c		2477-2495 (S)
	5'-GGTGTDTTGAACCAGGTTG-3' ^c		2869-2889 (S)
	5'-CTTCGTGATGCATTTCAACGG-3' ^c		2934-2954 (A)
	5'-GAAATGCGNCARGARTGYAA-3' ^a		1099-1118 (S)
	5'-ACYTRCARCACCARTGRAA-3' ^a		1756-1775 (A)
<i>wingless</i>	5'-GTTAGAACWTGYTGGATGCG-3' ^c	1147-1166 (S)	
	5'-AYAGTGATCACGNAATTCGG-3' ^c	1451-1471 (S)	
	5'-GAATTNCGTGATACACTRTTCG-3' ^c	1448-1469 (A)	
	5'-ATTYTTTTCRCAAAARCTTGG-3' ^c	1597-1617 (A)	
	5'-CGYTCNACNACAATRACCTC-3' ^c	1723-1742 (A)	
	5'-TGYGCTATGTNCARCARGAYGA-3' ^a	11404-11426 (S)	
<i>white</i>	5'-ACYTGNACRTAAAARTCNGCNGG-3' ^a	11975-11997 (A)	

^aDenotes primers used for PCR amplification.

^bThe Cho et al. (1995) primers used in the present study did not contain M13 primer sites.

^cDiopsid specific.

in ALIGN. For the 12S region, culling reduced the number of characters from ~305 (depending on the alignment) to 268 and for 16S from ~500 to 396.

Before an alignment of the protein-coding sequences, intron in both the *wingless* and *white* gene fragments was removed. Placement of the *wingless* intron for all the diopsids corresponds to the intron between the fourth and fifth exon of *Drosophila melanogaster* (Rijsewijk et al., 1987), which is ~100–180 bp long. Placement of the *white* intron in diopsids corresponds to the intron between exons 3 and 4a in *Bactrocera tryoni* (Bennett and Frommer, 1997), an intron that is ~50–90 bp long. Because the extreme sequence divergence in these introns provided little alignment stability, they were therefore excluded from all subsequent analyses. Alignment of the protein-coding gene fragments through their amino acid sequences was performed with ClustalX (Gibson et al., 1994). The alignments of COII, EF-1 α , and *white* were trivial because no indels were required; *wingless*, however,

contains a hypervariable region corresponding to a large (~85 amino acids) insertion found in Diptera (Rijsewijk et al., 1987). The amino acids were aligned by using gap:change costs of 2:1, 5:1, 10:1, 20:1, and 30:1 and the Blossum protein cost matrix (Gibson et al., 1994). These cost parameters produced only two different alignments, which varied slightly in the placement of three small indels. To discriminate between these two alignments, we placed the hypothesized gaps into the corresponding regions of the original DNA sequences and performed a parsimony search to determine which alignment required the fewest steps. That alignment was then used in the final matrix (this matrix, as well as the unaligned ribosomal fragments, is available at <http://research.amnh.org/molecular/sequence.html>)

Phylogenetic analysis.—Phylogenetic trees were generated by using parsimony in PAUP* versions 4.0d64 and 4.0b1 (Swofford, 1998). Unless otherwise specified, all characters were weighted equally and all

tree statistics were calculated with uninformative characters excluded. In most cases, heuristic searches using parsimony were performed with 100 random sequence-addition repetitions and TBR branch swapping. Bootstrap (Felsenstein, 1985) and Bremer support (Bremer, 1998, 1994) values were used to assess branch stability. Bootstrap values were calculated in PAUP by using 1,000 replicates and excluding uninformative characters. Bremer support values were calculated by constructing constraint commands specific to each node on the most-parsimonious tree and then searching for the most-parsimonious tree that did not contain that node. Data decisiveness (DD) values (Goloboff, 1991) were also calculated for the various separate and combined analyses, providing an index that serves as a measure of the information content of a given data set. This index is calculated as $DD = (\bar{L} - S) / (\bar{L} - M)$ where \bar{L} is the average length of a given data matrix on all possible trees, S is the actual length of the most-parsimonious tree (or trees), and M is the minimum possible number of steps for that matrix. In this case, the average length of each matrix was determined by using 100,000 random trees generated in PAUP.

Incongruence among the various gene regions was assessed by calculating incongruence length differences (ILD; Mickevich and Farris, 1981) and testing for their statistical significance by using a permutation procedure (Farris et al., 1994, 1995). This test was implemented by using the ARNIE program in the Random Cladistics Software Package (Siddall, 1995). All tests initially performed 999 permutations and used the Hennig86 search commands "mh" and "bb*" (Farris, 1988). Because numerous permutation tests were conducted, statistical significance was assessed by using Bonferroni-corrected critical values (Rice, 1989). Given the low alpha values required in these tests, the use of 999 permutations was insufficient to determine their significance accurately. Therefore, for comparisons that were initially significantly different at the 0.05 level, an additional 9,000 permutations were conducted for the Bonferroni-corrected test. ILD tests were performed on several different combinations of data from the entire matrix. In some cases, ILD tests were performed with certain nodes constrained to be present; for these runs, the partition homogeneity test in PAUP* was used.

Partitioned Bremer support (PBS) was used to evaluate the contribution of a given partition to the overall support of a combined analysis (Baker and DeSalle, 1997). This index divides the Bremer support at each node of the most-parsimonious tree from a combined analysis among the various partitions used to construct that tree. A partition can have either a positive or negative score at any given node, and the sum of the PBS scores of all the partitions for a given node will always equal the original Bremer support for that node. For the diopsid data, PBS scores were calculated as described in Baker and DeSalle (1997) and Baker et al. (1998). Twenty random sequence-addition repetitions were performed for each constrained search.

RESULTS

Combined and Separate Analyses of Gene Regions

Trees.—A combined analysis of all the data weighted equally resulted in a single most-parsimonious cladogram (Fig. 1). Statistics for this tree are presented in Table 3. In general, relationships on the tree are very strongly supported. Of the 32 nodes, 21 are characterized by bootstrap values ≥ 95 and by Bremer support values ≥ 10 . Three of the five major genera sampled are monophyletic. *Teleopsis* is paraphyletic and imbedded within *Cyrtodiopsis*. A monophyletic *Teleopsis* requires an additional seven steps on the tree, and a monophyletic *Cyrtodiopsis* requires an additional 12 steps. For both of these genera to be monophyletic, 21 extra steps are required. This tree also addresses the status of the monotypic genera *Chaetodiopsis* and *Trichodiopsis*. Originally described by Seguy (1955), these genera were considered dubious by Shillito (1971), who designated as synonymous with *Diasemopsis*. Feijen (1984) considered them both to be valid genera, but the phylogeny presented here firmly embeds them within the genus *Diasemopsis*.

The most important systematic result concerns the placement of *Sphyracephala*. Flies in this genus have eye stalks substantially smaller (about one-third to one-half of their body length) than those of other genera so their position within the family is critical for polarizing eye stalk evolution. On the basis of their eye span and a few other plesiomorphic morphological features, this genus has been presumed to be basal within the family

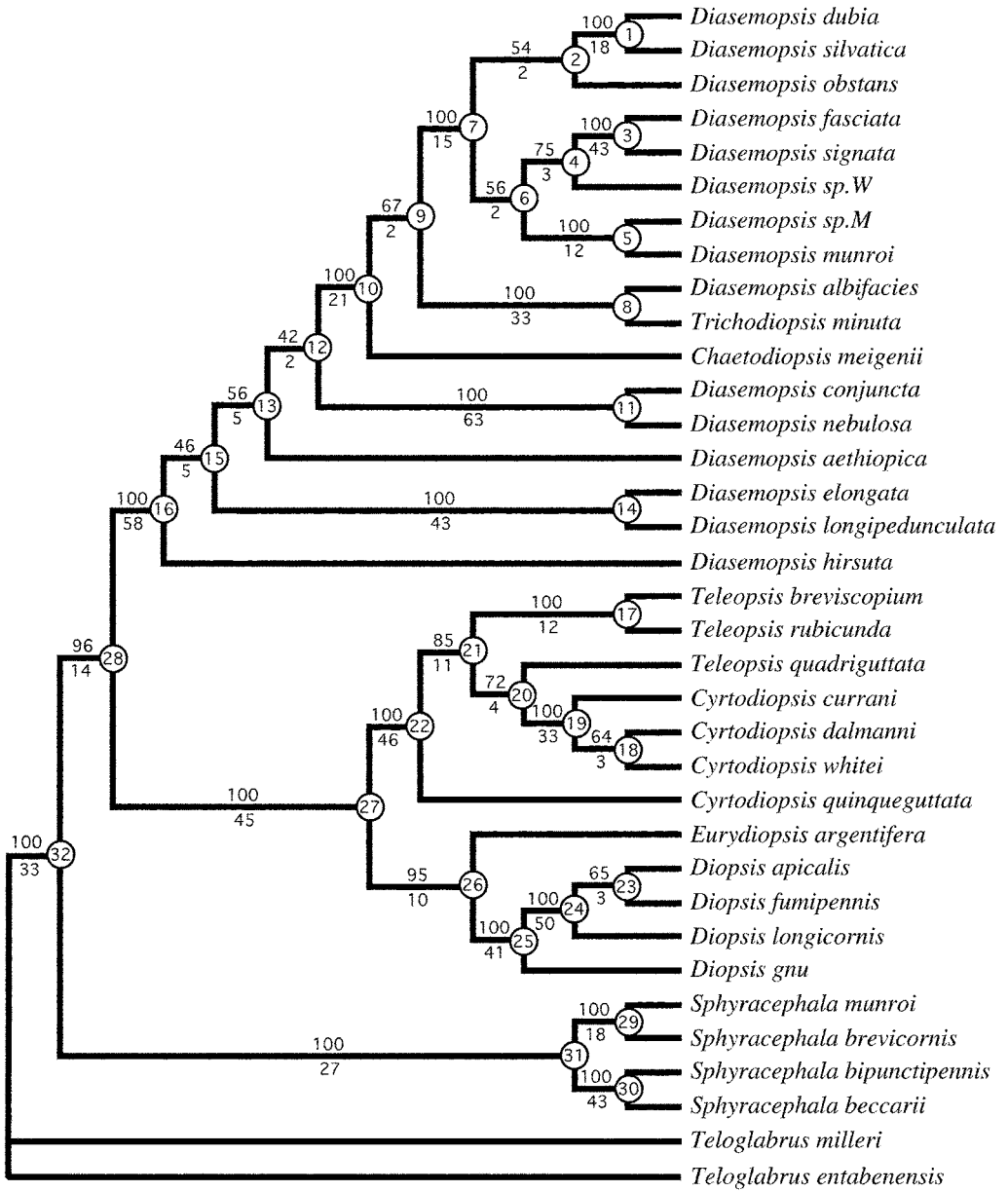


FIGURE 1. The single most-parsimonious cladogram from a combined analysis of all the data equally weighted. Tree statistics are presented in Table 3. Circled numbers at each node are reference numbers used throughout the text. Bootstrap values are provided above each node and Bremer support values below each node.

(Hennig, 1965; Shillitto, 1971; Feijen, 1989). The molecular data both corroborate the monophyly of the genus (bootstrap value of 100 and Bremer support value of 27) and support its basal placement on the tree (96 and 14, respectively).

Analysis of the various gene partitions making up the combined matrix reveals various degrees of disagreement. Figure 2

presents the strict consensus tree for each gene analyzed separately, for the mitochondrial genes alone, and for the nuclear genes alone. Results for these analyses vary in terms of topology, degree of resolution, and amount of homoplasy (Table 3). None of the eight trees is identical to any other or to the combined analysis tree. COII is most striking in its topological differences with the

TABLE 3. Tree statistics for various separate and combined analyses. chars. = characters, PI = phylogenetically informative, MP trees = number of most parsimonious trees, L = length, CI = consistency index, RI = retention index, DD = data decisiveness. Resolved nodes = the number of nodes resolved in a strict consensus of that partition's most-parsimonious trees, Congruent nodes = the number of these resolved nodes that also appear in the combined analysis topology (Fig. 1). PBS values for each data partition were summed across all the nodes on the combined analysis tree and standardized by the minimum possible number of steps for each partition.

Data partition	Total chars.	PI chars.	MP trees	L	CI	RI	DD	Resolved nodes	Congruent nodes	Summed PBS	PBS/min. steps
COII	436	164	8	827	.304	.472	.377	24	14	45.2	0.180
12S	268	56	756	212	.392	.588	.535	8	8	30.7	0.369
16S	396	79	72	304	.391	.668	.605	26	20	48.1	0.405
Mitochondrial	1100	299	12	1384	.327	.525	.443	23	20	124	0.274
EF-1 α	1031	224	24	758	.406	.712	.656	21	19	137	0.445
<i>wingless</i>	619	257	16	786	.533	.812	.776	27	25	294	0.702
<i>white</i>	486	186	6	758	.426	.719	.664	28	25	165	0.511
Nuclear ^a	2136	662	6	2272	.449	.755	.702	29	29	596	0.565
Combined	3236	966	1	3724	.404	.684	.625	32	—	720	0.479

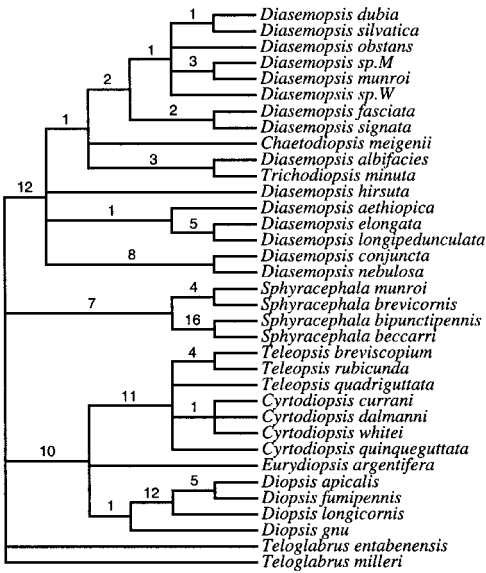
^aExcluding *Telolabrus milleri*.

other genes. Despite relatively good resolution (only eight most-parsimonious trees), the COII gene is characterized by the most homoplasy (CI = 0.304, RI = 0.472), the least decisiveness (DD = 0.377), and the largest number of nodes (10) in conflict with the combined analysis topology. 12S alone simply has very little structure, showing 756 equally most-parsimonious trees and only eight resolved nodes, although all of these nodes are consistent with the combined analysis tree. In contrast, the three nuclear-protein coding genes exhibit less homoplasy, more resolution, and greater topological similarities with each other and with the total evidence tree than do the mitochondrial genes (Table 3).

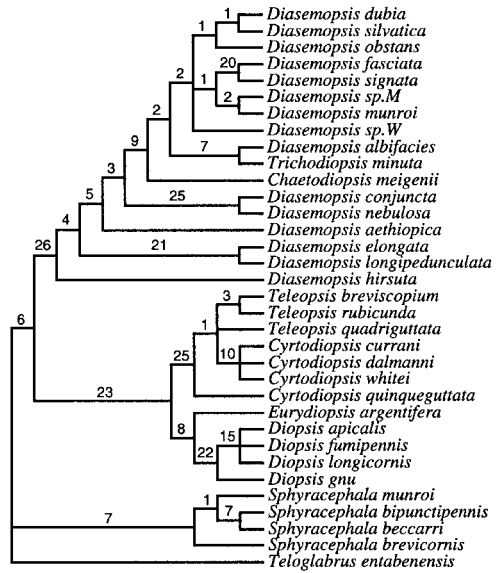
Incongruence.—The disagreement among data partitions can be examined more accurately by using incongruence measures, that account for the amount of nodal support present in each data set. An ILD test of all six gene partitions combined indicates significant heterogeneity among the data sets (ILD = 79 extra steps, $P < 0.0005$). Table 4 provides ILD values for all pairwise gene comparisons and for each gene examined against all the rest of the genes combined. Similar to the results presented for an analysis of Hawaiian *Drosophila* (Baker and DeSalle, 1997), the ILD matrix is asymmetrical with respect to the significance values of the various scores. In this case, neither 12S or 16S is significantly different from COII, but they are significantly different from each other. This pattern necessitates that any separation made among data partitions will require either the

combination of genes that are significantly different from one another or the separation of genes that are not significantly different from each other. Despite this asymmetry, the ILD matrix clearly indicates that the COII and 12S genes are significantly more problematic than the others. Both genes are significantly different from each of the nuclear genes (although the EF-1 α vs. COII comparison is not significant when Bonferroni-corrected values are used) and from the rest of the data combined. In contrast, pairwise comparisons among the other four genes show no significant heterogeneity, and none of these genes differs significantly from the rest of the data combined. Because of the influence of 12S and COII genes, the ILD test between the mitochondrial data and nuclear data also indicates significant incongruence (ILD = 24, $P < 0.001$).

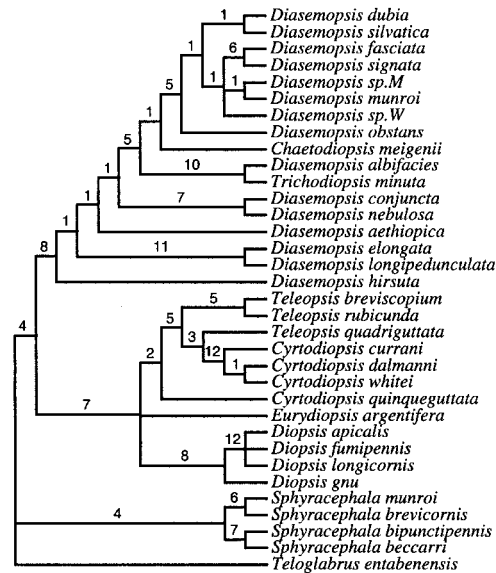
Partitioned Bremer support.—Given the clear differences in signal among the gene partitions, it is important to assess how these data sets interact in a combined analysis. PBS provides one means for examining this by partitioning the Bremer support at each node among the various genes. Just as all the Bremer support on a tree can be added to obtain a "total support" score (Källersjö et al., 1992), PBS scores can be summed across the entire tree to evaluate the contribution of a given partition to the overall support of the tree. Summed PBS scores are presented in Table 3. To control for differences in size between each data partition, we divided the PBS values by the minimum possible number of steps for each partition. The



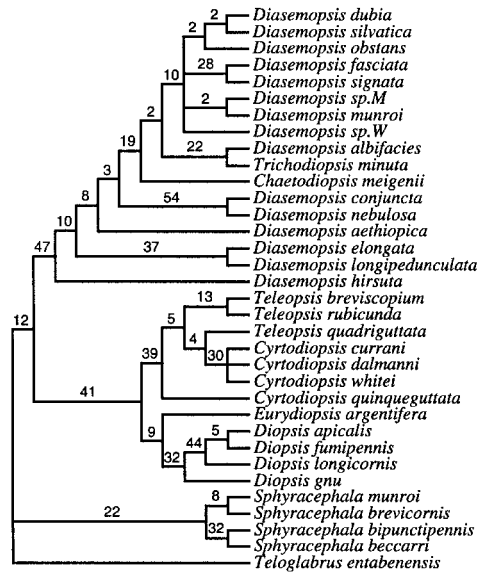
EF-1 alpha



wingless



white



nuclear

FIGURE 2. (Continued).

the mitochondrial data relative to the nuclear data. The standardized values for the mitochondrial data are about half that of the nuclear data; in the most extreme comparisons, *wingless* contributes more than threefold, and *white* and EF-1 α more than twofold, the amount of support provided by COII.

Phylogenetic Effects of Faster- and Slower-Evolving Characters

Nuclear genes.—Besides the dichotomy between mitochondrial and nuclear genes, we focused on the behavior of the more rapidly evolving nucleotides within each category.

TABLE 4. Incongruence length differences for each pairwise gene combination. The upper half of the matrix provides the number of extra steps that result from combining partitions, and the lower half provides the extra steps standardized by the length of the most-parsimonious tree or trees for each pairwise combined analysis. The Total column gives ILD values for each gene compared with the rest of the data combined.

Genes	COII	12S	16S	EF-1 α	<i>wingless</i>	<i>white</i>	Total
COII	—	24	19	20*	24**	24**	32**
12S	0.023	—	21**	23**	22**	18**	20**
16S	0.017	0.039	—	9	11	9	11
EF-1 α	0.012	0.023	0.008	—	8	8	14
<i>wingless</i>	0.015	0.022	0.01	0.005	—	4	12
<i>white</i>	0.015	0.019	0.008	0.005	0.003	—	2

*Statistically significant ($P < 0.01$); **Statistically significant when using a Bonferroni-corrected alpha value.

As expected, for the nuclear genes, third position sites for each codon are clearly changing at a faster rate than are the first or second position sites. The uncorrected pairwise divergence for third positions (range, from 0.008 to 0.45) is much greater than for first and second positions (0 to 0.07). This divergence is most extreme in the *white* gene, which exhibits some third position pairwise differences >60%. A plot of sequence divergence for the nuclear third positions relative to first and second positions (Fig. 3) suggests saturation in the third position data. The asymptotic point on this graph roughly corresponds to intergeneric comparisons. Given the rapid evolution of these characters, they have been generally assumed to provide little useful phylogenetic information, particularly for more distantly related taxa (Swofford et al., 1996). Therefore, their phylogenetic utility deserves investigation.

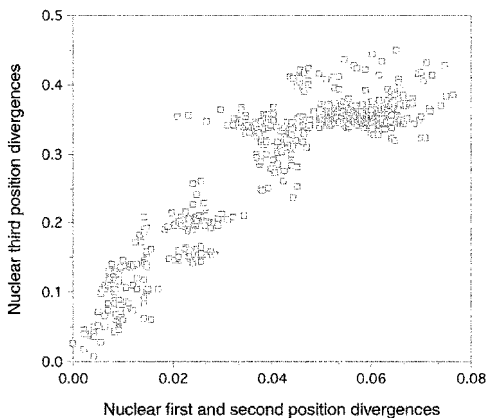


FIGURE 3. Uncorrected pairwise sequence divergence plot of third position nucleotide sites relative to first and second nucleotide sites for the three nuclear genes combined.

The third position data are generally more homoplasious than the first and second positions, but the topologies between the two partitions are very similar. Both analyses produced three most-parsimonious trees (Table 5), and a comparison of the strict consensus for each partition reveals topological conflict at only five nodes. ILD tests indicate no heterogeneity between the partitions (ILD = 8, $P > 0.05$). Furthermore, the inclusion of third position sites reduces the incongruence among the three different genes. When first and second positions are analyzed alone, the ILD score among the three genes is 21 extra steps, but this drops to 15 when third positions are added to the analysis. Permutation tests indicate that,

TABLE 5. Tree statistics, incongruence values, and support measures for the first/second position sites and the third position sites of the three nuclear genes. MP trees = number of most-parsimonious trees, L = length, CI = consistency index, RI = retention index, DD = data decisiveness. CATree nodes, the number of nodes based on a strict consensus of all the MP trees that are consistent with the combined analysis topology (Fig. 1). ILD, the number of extra steps that result from combining the three nuclear genes by using only the characters from the designated codon positions. PBS values are presented both as summed across the entire combined analysis topology and as standardized by the minimum possible number of steps for each partition.

	Codon position	
	1st and 2nd	3rd
MPtree	3	3
L	370	1,896
CA tree nodes	23	28
CI	0.608	0.431
RI	0.866	0.725
DD	0.837	0.676
ILD	21	14
Summed PBS	155	441
PBS/min. steps	0.689	0.540

compared across genes, the first and second position sites alone exhibit significant conflict (ILD = 21, $P < 0.001$), but there is no heterogeneity for the thirds alone (ILD = 14, $P > 0.6$) or for all the codon positions combined (ILD = 15, $P > 0.3$). PBS scores (Table 5) also indicate that third position characters have substantial phylogenetic signal. The summed PBS values for the third position data account for well over half (441 of 720) of all the support for the combined analysis tree. Relative to their abundance, the informative first and second positions provide slightly more support than third positions but as a class of data, third positions contain substantially more information.

In addition, the distribution of the support provided by third position sites is not limited to the tips of the tree. The divergence plot suggests that thirds become saturated across intergeneric comparisons. If these plots accurately reflect the utility of the third position data, then these characters should be disproportionately informative for more closely related taxa. To evaluate this possibility, we divided all the nodes on the combined analysis tree (Fig. 1) into three categories: (1) between sister species nodes; (2) within genera nodes but not including sister species nodes (*Teleopsis* and *Cyrtodiopsis* were treated as a single genus for this analysis); and (3) generic and intergeneric nodes. Because sequence data for *wingless* and *white* were not available for *T. milleri*, the ingroup-outgroup node (no. 32 in Fig. 1) was excluded from this analysis. Of the 31 nodes on the combined analysis tree, 11 fell into the first category, 13 into the second, and 7 into the third. PBS scores for the third position sites were then summed across all of the nodes in each of the three categories (Table 6). Overall, the third positions pro-

vided more information per node for the basal intergeneric nodes than for the other two node categories. In addition, although relative to the rest of the data, they provide the most support for the relationships in node category 2 (91% of all the support), they still provide the majority of support for the basal relationships (56% of all the support).

Mitochondrial genes.—Given the relatively poor performance of the mitochondrial data in general and the COII and 12S genes specifically, we examined the extent to which this result depended on the presence of rapidly evolving characters. For a combined COII and 12S data set, we used a series of weighting schemes traditionally used for reducing the effect of faster-evolving characters and then assessed the results relative to those for the combined analysis topology. The plots in Figures 4a and 4b show that the divergences between more closely related taxa for COII and 12S are characterized by relatively more transitions and third position changes. This pattern suggests that these characters are evolving at a faster rate than transversions or first and second position sites, a result consistent with several other studies (Brown et al., 1982; DeSalle et al., 1987; Helm-Bychowski and Cracraft, 1993; Bloomer and Crowe, 1998). The plots of third position sites and transitions relative to total divergences (Fig. 4c and 4d) also suggest these data are more saturated. Therefore, the different weighting schemes used the following options: (1) all characters weighted equally, (2) transitions for 12s and for COII third positions eliminated, (3) all transitions eliminated, (4) COII third positions eliminated, or (5) all transitions and COII third positions eliminated.

Table 7 shows the length of the most-parsimonious trees obtained by using the above weighting schemes and also the number of extra steps required to fit the weighted data to the combined analysis topology (Fig. 1). Despite some increases in consistency, Table 7 indicates that using the various weighting schemes adds little improvement in congruence and, moreover, when standardized by the length of the weighted analyses, the equally weighted data exhibit the least incongruence with the combined analysis topology. Likewise, the equally weighted data have the greatest topological similarity (18 shared nodes) to the combined analysis tree.

TABLE 6. Distribution of support on combined analysis tree provided by nuclear third position sites. Node categories: (1) most derived nodes, (2) intermediate nodes, and (3) most basal nodes (see text for specific definitions). PBS values are presented both as summed across all the nodes in a given category and as standardized by the number of nodes in that category. The percentage of total support (resulting from all the data) across these node categories that is provided by the third position data is also shown.

Node category	No. of nodes	Summed PBS	PBS/node	% support
1	11	165.7	15.1	57
2	13	141.3	10.9	91
3	7	134	19.1	56

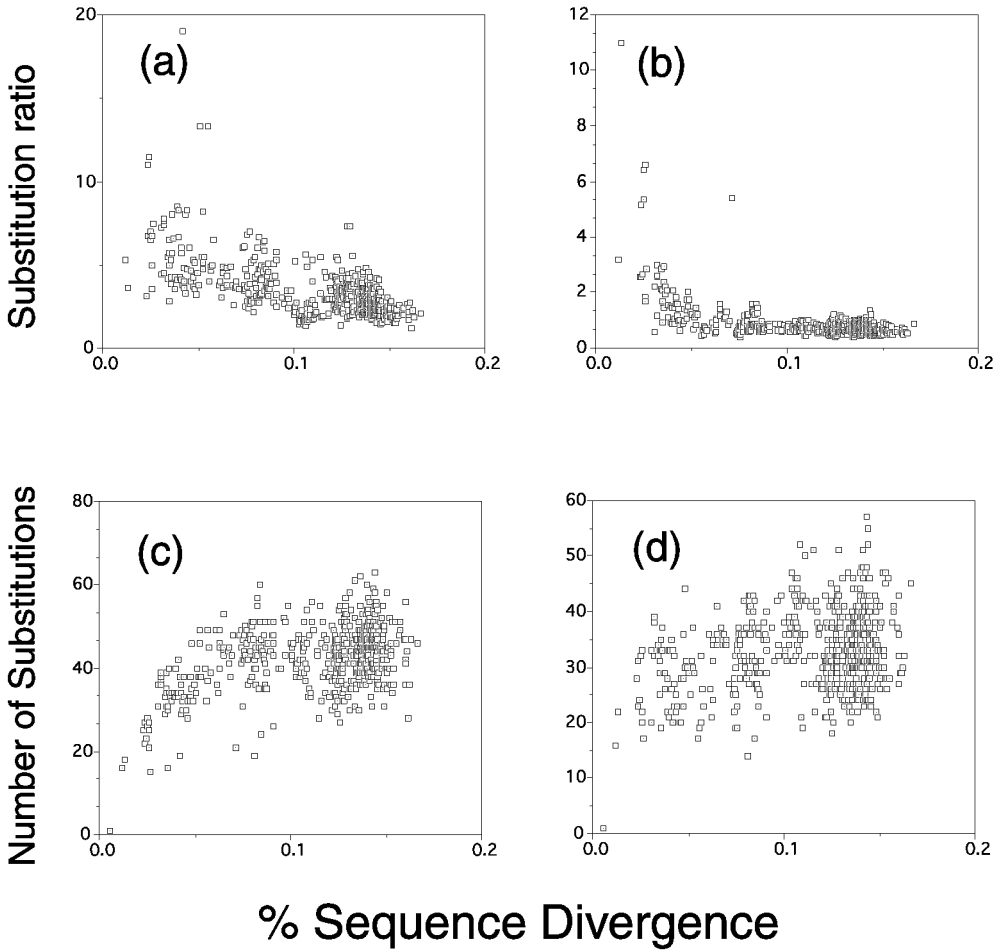


FIGURE 4. Relative substitution rates of transitions and transversions for the COII and 12S genes and of third positions and first and second positions for the COII gene. (a) Third position/first and second position ratio for the COII gene only. (b) Transition/transversion ratio for the COII and 12S genes combined. (c) Third position substitutions for the COII gene. (d) Transition substitutions for the COII and 12S genes combined. All substitution patterns are plotted relative to the total sequence divergence of the six genes combined.

TABLE 7. The effect of various weighting schemes on the congruence of a combined COII and 12S data set to the combined analysis topology (Fig. 1). 3rd, COII third codon position sites; ti, transitions; L, length; CI, consistency index. CA tree nodes, the number of nodes obtained on a strict consensus of all the most-parsimonious trees that also appear in the combined analysis topology (Fig. 1). CA tree ILD, the increase in length that results from constraining the weighted data to the combined analysis topology.

Weighting scheme	L	CI	CA nodes	CA tree ILD	ILD/L
All equal	1063	0.314	18	28	0.026
No 3rd ti, 12S ti	629	0.273	16	28	0.045
No ti	526	0.257	13	27	0.051
No 3rd	364	0.393	15	23	0.063
No 3rd, ti	163	0.393	9	23	0.141

DISCUSSION

Combined Analysis and Incongruence

The importance of using multiple, diverse sources of phylogenetic information has become increasingly clear as numerous studies using several data partitions have consistently demonstrated the limited ability of single data partitions to accurately reconstruct phylogeny (Cao et al., 1994; Olmstead and Sweere, 1994; Cummings et al., 1995; Baker and DeSalle, 1997). The combined and separate analyses of the six different gene regions presented here reinforce this point and reveal the limitations of even a relatively large data set. With 966 phylogenetically informative

characters, the diopsid matrix is certainly larger than most studies and the fact that the resulting tree is completely resolved with relatively strong support at many nodes is encouraging. However, several nodes are still characterized by weak support, suggesting that certain relationships may be inherently difficult to resolve.

As the amount of data in systematic studies increases, issues relating to their analysis become more complex and controversial. Specifically, the issue of data combinability has received the most attention, and both theoretical (Bull et al., 1993; Eernisse and Kluge, 1993; Miyamoto and Fitch, 1995; Nixon and Carpenter, 1996; DeSalle and Brower, 1997) and empirical (Baker and DeSalle, 1997; Cunningham, 1997a, 1997b; Miller et al., 1997; Davis et al., 1998) aspects of this debate have now been extensively treated in the literature. Substantial conflict among data partitions has become a prevalent phenomenon (Baker and DeSalle, 1997; Cunningham, 1997a, 1997b; Baker et al., 1998; Baum et al., 1998; O'Grady et al., 1998; Remsen and DeSalle, 1998), and it has emerged again in this study. Some have argued (Bull et al., 1993; Huelsenbeck et al., 1994) that significantly conflicting data should not be combined in a single analysis because that might increase the chances of obtaining misleading results; they contend being conservative in such situations, is more appropriate.

In this study, however, we found that separating data seems more likely to obscure relationships than either to clarify or to be appropriately conservative about them. Analyzing separately a data partition including both COII and 12S showed either conflict or ambiguity at 15 nodes when compared with the topology derived from the rest of the data (Table 8). This ambiguity does not, however, accurately reflect the disagreement between the two partitions. Table 8 shows the decrease in the incongruence length difference (from a value of 28). Nodes 13, 15, and 17 stand out as being particularly critical to this incongruence, and when these three nodes are constrained together, a permutation test between the partitions no longer indicates heterogeneity (ILD = 12, $p > 0.05$). This local incongruence has been documented in other studies (Mason-Gamer and Kellogg, 1996; Poe, 1996; Baker and DeSalle, 1997) and has important implications for the combinability of data. First, by separating the two partitions

TABLE 8. The decrease in ILD scores between the COII/12S data partition and the nuclear genes/16S data partition when each of the nodes on the combined analysis tree is enforced as a constraint on the tree searches of each partition (node numbers correspond to those in Fig. 2). The ILD is 28 extra steps when no constraints are enforced. A plus sign indicates that the node occurs in the consensus trees of both partitions. These are the only nodes in agreement between the two trees.

Node no.	ILD score decrease
1	0+
2	1
3	0+
4	0
5	0+
6	0
7	0+
8	0+
9	1
10	0+
11	0+
12	3
13	8
14	0+
15	10
16	0+
17	5
18	2
19	1
20	2
21	1
22	0+
23	3
24	0+
25	0+
26	0+
27	3
28	1
29	0+
30	0+
31	0+
32	0+

in this case, we are left with an ambiguous and conservative statement concerning 15 nodes on the tree when, in fact, only 3 of these nodes are actually problematic. This ambiguity is particularly misleading for five of the nodes, which are very strongly supported (bootstraps between 85 and 100, Bremer supports between 11 and 45) on the combined analysis tree. Second, besides being overly conservative about some conflicting relationships, separating data prevents unique resolution of previously unsupported relationships that are not the source of the incongruence between partitions. The relationships among eight *Diasemopsis* species (nodes 1–7 in Fig. 1) are only partly resolved in both separate analyses (COII/12S and 16S/nuclear) but become fully resolved when the data are combined. A similar

relationship between local incongruence and the effects of separating and combining data has been demonstrated in a study of Hawaiian *Drosophila* relationships (Baker and DeSalle, 1997).

Phylogenetic Utility of Different Classes of Molecular Data

Nuclear versus mitochondrial genes.—By virtually any measure used, the utility of the nuclear genes is substantially superior to that of the mitochondrial genes in this study. The nuclear genes exhibit more resolution, less homoplasy, and greater support than do the mitochondrial genes (Table 3). In addition to their consistently weaker signal, mitochondrial genes exhibit more variation in signal among the different partitions. The incongruence analyses demonstrate that nearly all of the conflict among partitions in the combined matrix results from mitochondrial genes. This is especially interesting considering that mitochondrial genes are linked and therefore necessarily share the same history. Moore (1995) suggested that disagreement among nuclear genes should be more common than among mitochondrial genes because of the increased probability of lineage sorting for these loci, but that prediction is clearly contradicted by the genes collected for this study. Overall, the differential utility of protein-coding nuclear genes is an important result, given the number of systematic studies based exclusively on mitochondrial data. Whether this discrepancy is limited to these genes and this group of flies or represents a more general pattern requires additional studies using several genes of each type.

Faster-evolving versus slower-evolving characters.—Saturation plots and divergence estimates have become a prominent feature of molecular systematics (Mindell et al., 1996; Murphy and Collier, 1997; Bloomer and Crowe, 1998; Danforth and Ji, 1998; Martin and Bermingham, 1998), and many of the characters examined in this study show evidence of rapid evolution and saturation. Given this type of evidence, downweighting or eliminating the influence of such characters in the final analysis is not unusual. Whether this strategy consistently improves phylogenetic estimation, however, is increasingly being questioned (Olmstead

et al., 1998; Yang, 1998; Källersjö et al., 1999; Wenzel and Siddall, 1999).

For the nuclear genes in this study, the majority of intergeneric divergences of the *wingless* and *white* third position sites are >50%, a value that exceeds many of the recommended divergence cutoffs for phylogenetic usefulness (Hillis and Dixon, 1991; Friedlander et al., 1994; Graybeal, 1994). The incongruence and support analyses, however, indicate these data still contain substantial phylogenetic signal. Including the third position sites reduces the conflict among the three nuclear genes and contributes more support to the combined analysis tree than all the rest of the data combined. Moreover, this support is present at all levels of the tree. Similar conclusions can be drawn from our analysis of the faster-evolving characters in the COII and 12S regions. The divergence plots suggest substantial saturation in these genes, particularly for the transitions and third position sites. All of the weighting strategies examined for this partition, however, reduce the concordance of the COII/12S data to the rest of the information in terms of both topological similarities (identical nodes) and ILD measures. Given the overall strength of the relationships on the combined analysis tree and the fact that, ultimately, the best measure of accuracy for a phylogenetic hypothesis involves topological similarity and support, the information from PBS and ILD measures should take precedent over less direct indices of data quality, such as divergences and homoplasy.

Overall, the evidence presented in this paper suggests that saturation curves have limited value, especially as a criterion for downweighting or eliminating data. The consistent reliability of divergence estimates as a means for assessing phylogenetic utility is contradicted by the improved phylogenetic estimation obtained by including the more rapidly evolving characters. Recent theoretical analysis (Yang, 1998) has also questioned the utility of saturation plots. He concluded that "pairwise sequence divergence is not a good indicator of the information content in the data, as the accuracy depends on not only the amount of evolution, but also on how many branches the tree has and how the substitutions are distributed among the branches in the tree." Results from the present study, as well as those from Olmstead et al. (1998),

Bjorkland (1999), and Källersjö et al. (1999), provide an important empirical demonstration of this contention.

A primary reason why saturation curves and divergence estimates can be misleading is that they fail to distinguish between variation among sites and variation within a site (Olmstead et al., 1998; Bjorkland, 1999). Rate differences between partitions may be caused by variations in the proportions of informative sites; for instance, a class of data may exhibit high divergences even if it evolves at a relatively slow rate as long as most of its sites are free to vary. In this way, single-measure divergence estimates may markedly overestimate the difference in substitution rates among classes of data. For the nuclear genes in this study, the comparison of the percent divergences for different codon sites suggests substantial rate differences. Nearly all (98.5%) of the pairwise species comparisons show third position divergence estimates at least fivefold those of the corresponding first and second position divergence estimates. Comparisons of the average substitution rate for each informative site, however, indicate that third positions are changing at ~ 1.6 times the rate of first and second positions. Therefore, although more third position sites were variable, the average rate of change for each character was similar between codon positions.

In addition, saturation curves do not account for the number of taxa sampled in a study. In this study, the saturation curves would have shown a similar asymptote if 5 rather than 17 *Diasemopsis* species had been included. Trees that contain more taxa, however, can accommodate more rapid substitution rates because there are more branches among which to distribute these changes (Yang, 1998). In fact, Yang's simulations predicted that for an analysis of 50 taxa, the optimal substitution rate per site is between two and four. The rate of change for the nuclear third position sites (3.76) fits within this limit. Saturation curves also misrepresent the weight of the divergence evidence relative to the nodes on the tree. Relationships between the more divergent taxa provide the majority of points in a saturation curve but relate to only a few nodes on the tree. For instance, comparisons between intergeneric diopsid species represent 71.6% of the data points in the saturation figures, but intergeneric nodes on the tree are only 12.5% of the total.

Arguments concerning the utility of rapidly-evolving and saturated data are part of a larger debate within molecular systematics that concerns the distinction between data quality and data sufficiency. Most attempts to improve the overall quality of phylogenetic information require reducing the overall weight of this information. Theoretical analyses have attempted to identify the critical points in data quality/data sufficiency space where misleading results will occur (Felsenstein, 1978; Bull et al., 1993; Givnish and Sytsma, 1997), but these studies often focus on limited examples, such as four taxon statements or extreme rate differences. Therefore, it is critical to explore the tradeoffs between data quality and data sufficiency with actual empirical systems. This approach has recently been adopted by several authors (Mindell and Thacker, 1996; Russo et al., 1996; Cunningham, 1997a, 1997b; Allard et al., 1999), who are using the entire mitochondrial genome of various vertebrate taxa. These studies represent an important first step but have generally produced ambiguous results. More importantly, these studies fail to capture accurately the real limitations of most systematic studies. As yet, the literature contains few studies in which one could remove all third positions or transitions and still be left with >100 (and up to 400) phylogenetically informative characters to resolve <15 species. In most cases, the effects of data removal or downweighting will be much more acute than in the studies that use the entire mitochondrion. In an attempt to examine these issues with a more realistic data set, we have identified here three different instances in which the utility of including a certain type of data might be viewed as questionable. They are combining the COII and 12S genes with the rest of the data, including the third position sites from the nuclear genes, and including the faster-evolving characters from the COII and 12S genes. In all three cases, the results strongly suggest that phylogenetic estimation is improved by including this information within an equally weighted framework.

ACKNOWLEDGMENTS

We thank John Wenzel, John Gatesy, Patrick O'Grady, Paul Goldstein, and Robert Hanner for helpful comments on the manuscript. We thank two anonymous reviewers, Andy Brower, and Richard Olmstead for suggesting significant improvements to an earlier

version. The collection of field specimens was aided by several individuals and organizations. For their help, we are grateful to Sarah Laird, Terry Sunderland, Rudolf Meier, Lee White, Henri Bourobou-Bourobou, Richard Ruggiero, David Barraclough, Limbe Botanical Gardens, and the Wildlife Conservation Society. Sabine Hilger and Marion Kotrba kindly provided assistance in the identification of many of the specimens. R. H. B. was supported by an AMNH Graduate Fellowship; the research was funded by National Science Foundation awards DEB-9423436 and DEB-9409369.

REFERENCES

- ALLARD, M. W., J. S. FARRIS, AND J. M. CARPENTER. 1999. Congruence among mammalian mitochondrial genes. *Cladistics* 15:75–84.
- BAKER, R. H., AND R. DESALLE. 1997. Multiple sources of character information and the phylogeny of Hawaiian *Drosophilids*. *Syst. Biol.* 46:654–673.
- BAKER, R. H., X. YU, AND R. DESALLE. 1998. Assessing the relative contribution of molecular and morphological characters in simultaneous analysis trees. *Mol. Phylogenet. Evol.* 9:427–436.
- BAKER, R. H., AND G. S. WILKINSON. submitted. Phylogenetic analysis of eye stalk allometry and sexual dimorphism in stalk-eyed flies (diopsidae).
- BAUM, D. A., R. L. SMALL, AND J. F. WENDEL. 1998. Biogeography and floral evolution of Baobabs (*Adansonia*, Bombacaceae) as inferred from multiple data sets. *Syst. Biol.* 47:181–207.
- BENNETT, C. L., AND M. FROMMER. 1997. The white gene of the tephritid fruit fly *Bactrocera tryoni* is characterized by a long untranslated 5' leader and a 12 kb first intron. *Insect Mol. Biol.* 6:343–356.
- BJORKLUND, M. 1999. Are third positions really that bad? A test using vertebrate cytochrome *b*. *Cladistics* 15:191–197.
- BLOOMER, P., AND T. M. CROWE. 1998. Francolin phylogenetics: Molecular, morphobehavioral, and combined evidence. *Mol. Phylogenet. Evol.* 9:236–254.
- BREMER, K. 1988. The limits of amino acid sequence data in angiosperm phylogenetic reconstruction. *Evolution* 42:795–803.
- BREMER, K. 1994. Branch support and tree stability. *Cladistics* 10:295–304.
- BROWER, A. V. Z. 1994. Phylogeny of *Heliconius* butterflies inferred from mitochondrial DNA sequences (Lepidoptera: Nymphalidae). *Mol. Phylogenet. Evol.* 3:159–174.
- BROWER, A. V. Z., AND R. DESALLE. 1994. Practical and theoretical considerations for choice of a DNA sequence region in insect molecular systematics, with a short review of published studies using nuclear gene regions. *Ann. Entomol. Soc. Am.* 87:702–716.
- BROWN, W. M., E. M. PRAGER, A. WANG, AND A. C. WILSON. 1982. Mitochondrial DNA sequences of primates: Tempo and mode of evolution. *J. Mol. Biol.* 18:225–239.
- BULL, J. J., J. P. HUELSENBECK, C. W. CUNNINGHAM, D. L. SWOFFORD, AND P. J. WADDELL. 1993. Partitioning and combining data in phylogenetic analysis. *Syst. Biol.* 42:384–397.
- BURKHARDT, D., AND I. DE LA MOTTE. 1985. Selective pressures, variability, and sexual dimorphism in stalk-eyed flies (Diopsidae). *Naturwissenschaften* 72:204–206.
- BURKHARDT, D., AND I. DE LA MOTTE. 1988. Big 'antlers' are favoured: Female choice in stalk-eyed flies (Diptera, Insecta), field collected harems and laboratory experiments. *J. Comp. Physiol. A* 162:649–652.
- BURKHARDT, D., I. DE LA MOTTE, AND K. LUNAU. 1994. Signalling fitness: Larger males sire more offspring. Studies of the stalk-eyed fly *Cyrtodopsis whitei* (Diopsidae, Diptera). *J. Comp. Physiol. A* 174:61–64.
- CAO, Y., J. ADACHI, A. JANKE, S. PÄÄBO, AND M. HASEGAWA. 1994. Phylogenetic relationships among eutherian orders estimated from inferred sequences of mitochondrial proteins: Instability of a tree based on a single gene. *J. Mol. Evol.* 39:519–527.
- CARPENTER, J. M. 1988. Choosing among multiple equally parsimonious cladograms. *Cladistics* 4:291–296.
- CHIPPINDALE, P. T., AND J. J. WIENS. 1994. Weighting, partitioning, and combining characters in phylogenetic analysis. *Syst. Biol.* 43:278–287.
- CHO, S., A. MITCHELL, J. C. REGIER, C. MITTER, R. W. POOLE, T. P. FRIEDLANDER, AND S. ZHAO. 1995. A highly conserved nuclear gene for low-level phylogenetics: Elongation factor-1 α recovers morphology-based tree for heliothine moths. *Mol. Biol. Evol.* 12:650–656.
- COLLINS, T. M., F. KRAUS, AND G. ESTABROOK. 1994. Compositional effects and weighting of nucleotide sequences for phylogenetic analysis. *Syst. Biol.* 43:449–459.
- CUMMINGS, M. P., S. P. OTTO, AND J. WAKELEY. 1995. Sampling properties of DNA sequence data in phylogenetic analysis. *Mol. Biol. Evol.* 12:814–822.
- CUNNINGHAM, C. W. 1997a. Can three incongruence tests predict when data should be combined? *Mol. Biol. Evol.* 14:733–740.
- CUNNINGHAM, C. W. 1997b. Is congruence between data partitions a reliable predictor of phylogenetic accuracy? Empirically testing an interactive procedure for choosing among phylogenetic methods. *Syst. Biol.* 46:464–478.
- DANFORTH, B. N., AND S. JI. 1998. Elongation factor-1 α occurs as two copies in bees: Implications for phylogenetic analysis of EF-1 α sequences in insects. *Mol. Biol. Evol.* 15:225–235.
- DAVIS, J. I., M. P. SIMMONS, D. W. STEVENSON, AND J. F. WENDEL. 1998. Data decisiveness, data quality, and incongruence in phylogenetic analysis: An example from the monocotyledons using mitochondrial *atp A* sequences. *Syst. Biol.* 47:282–310.
- DESALLE, R., AND A. V. Z. BROWER. 1997. Process partitions, congruence, and the independence of characters: Inferring relationships among closely related Hawaiian *Drosophila* from multiple gene regions. *Syst. Biol.* 46:751–764.
- DESALLE, R., T. FREEDMAN, E. M. PRAGER, AND A. C. WILSON. 1987. Tempo and mode of sequence evolution in mitochondrial DNA of Hawaiian *Drosophila*. *J. Mol. Evol.* 26:157–164.
- EERNISSE, D. J., AND A. G. KLUGE. 1993. Taxonomic congruence versus total evidence, and amniote phylogeny inferred from fossils, molecules, and morphology. *Mol. Biol. Evol.* 10:1170–1195.
- FANG, Q. Q., S. CHO, J. C. REGIER, C. MITTER, M. MATTHEWS, R. W. POOLE, T. P. FRIEDLANDER, AND S. ZHOU. 1997. A new nuclear gene for insect phylogenetics: Dopa decarboxylase is informative of relationships within Heliothinae (Lepidoptera: Noctuidae). *Syst. Biol.* 46:269–283.

- FARRIS, J. 1988. Hennig86, Program and documentation, version 1.5.
- FARRIS, J. S. 1969. A successive approximations approach to character weighting. *Syst. Zool.* 18:374–385.
- FARRIS, J. S., M. KÄLLERSJÖ, A. G. KLUGE, AND C. BULT. 1994. Testing significance of congruence. *Cladistics* 10:315–320.
- FARRIS, J. S., M. KÄLLERSJÖ, A. G. KLUGE, AND C. BULT. 1995. Constructing a significance test for incongruence. *Syst. Biol.* 44:570–572.
- FEIJEN, H. R. 1983. Systematics and phylogeny of Centronidae, a new Afrotropical family of Diptera (Schizophora). *Zool. Verh. (Leiden)* 202:1–137.
- FEIJEN, H. R. 1984. A further contribution to the genus *Diopsina* Curran, 1928 (Diptera: Diopsidae). *Rev. Zool. Afr.* 98:9–34.
- FEIJEN, H. R. 1989. Diopsidae. Pages 1–122 in *Flies of the Nearctic Region* (G. C. D. Griffiths, ed.). E. Schweizerbartsche Verlagsbuchhandlung, Stuttgart.
- FELSENSTEIN, J. 1978. Cases in which parsimony or compatibility methods will be positively misleading. *Syst. Zool.* 27:401–410.
- FELSENSTEIN, J. 1985. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution* 39:783–791.
- FITCH, W. M., AND J. YE. 1991. Weighted parsimony: Does it work? Pages 147–155 in *Phylogenetic analysis of DNA sequences* (M. M. Miyamoto and J. Cracraft, eds.). Oxford Univ. Press, New York.
- FRIEDLANDER, T. P., K. R. HORST, J. C. REGIER, C. MITTER, R. S. PEIGLER, AND Q. Q. FANG. 1998. Two nuclear genes yield concordant relationships within Attacini (Lepidoptera: Saturniidae). *Mol. Phylogenet. Evol.* 9:131–140.
- FRIEDLANDER, T. P., J. C. REGIER, AND C. MITTER. 1994. Phylogenetic information content of five nuclear gene sequences in animals: Initial assessment of character sets from concordance and divergence studies. *Syst. Biol.* 43:511–525.
- FRIEDLANDER, T. P., J. C. REGIER, C. MITTER, AND D. L. WAGNER. 1996. A nuclear gene for higher level phylogenetics: Phosphoenolpyruvate carboxykinase tracks Mesozoic-age divergences within lepidoptera (Insecta). *Mol. Biol. Evol.* 13:594–604.
- GATESY, J., R. DESALLE, AND W. C. WHEELER. 1994. Alignment-ambiguous nucleotide sites and the exclusion of data. *Mol. Phylogenet. Evol.* 2:152–157.
- GATESY, J., C. HAYASHI, M. CRONIN, AND P. ARCTANDER. 1996. Evidence from milk casein genes that cetaceans are close relatives of hippopotamid artiodactyls. *Mol. Biol. Evol.* 13:954–963.
- GIBSON, T., D. HIGGINS, AND J. THOMPSON. 1994. ClustalX. Program and documentation available at <ftp://ftp.ebi.ac.uk/pub/software/mac/clustalw/clustalx/>.
- GIVNISH, T. J., AND K. J. SYTSMA. 1997. Consistency, characters, and the likelihood of correct phylogenetic inference. *Mol. Phylogenet. Evol.* 7:320–330.
- GOLOBOFF, P. A. 1991. Homoplasy and the choice among cladograms. *Cladistics* 7:215–232.
- GOLOBOFF, P. A. 1993. Estimating character weights during tree search. *Cladistics* 9:83–91.
- GRAYBEAL, A. 1994. Evaluating the phylogenetic utility of genes: A search for genes informative about deep divergences among vertebrates. *Syst. Biol.* 43:174–193.
- GRIMALDI, D., AND G. FENSTER. 1989. Evolution of extreme sexual dimorphisms: Structural and behavioral convergence among broad-headed Drosophilidae (Diptera). *Am. Mus. Nov.* 2939:1–25.
- HELM-BYCHOWSKI, K., AND J. CRACRAFT. 1993. Recovering phylogenetic signal from DNA sequences: Relationships within the corvine assemblage (Class Aves) as inferred from complete sequences of the mitochondrial DNA cytochrome-*b* gene. *Mol. Biol. Evol.* 10:1196–1214.
- HENNIG, W. 1958. Die Familien der Diptera Schizophora und ihre phylogenetischen Verwandtschaftsbeziehungen. *Beitr. Entomol.* 8:505–688.
- HENNIG, W. 1965. Die Acalypterae des Baltischen Bernsteins und ihre Bedeutung für die Erforschung der phylogenetischen Entwicklung dieser Dipteren-Gruppe. *Stuttg. Beitr. Naturkd.* 145:1–215.
- HILLIS, D. M. 1991. Discriminating between phylogenetic signal and random noise in DNA sequences. Pages 278–295 in *Phylogenetic analysis of DNA sequences* (M. M. Miyamoto and J. Cracraft, eds.). Oxford Univ. Press, New York.
- HILLIS, D. M., AND M. T. DIXON. 1991. Ribosomal DNA: Molecular evolution and phylogenetic inference. *Q. Rev. Biol.* 66:411–453.
- HUELSENBECK, J. P., D. L. SWOFFORD, C. W. CUNNINGHAM, J. J. BULL, AND P. J. WADDELL. 1994. Is character weighting a panacea for the problem of data heterogeneity in phylogenetic analysis? *Syst. Biol.* 43:288–291.
- KÄLLERSJÖ, M., V. A. ALBERT, AND J. S. FARRIS. 1999. Homoplasy increases phylogenetic signal. *Cladistics* 8:275–287.
- KÄLLERSJÖ, M., J. S. FARRIS, A. G. KLUGE, AND C. BULT. 1992. Skewness and permutation. *Cladistics* 8:275–287.
- KLUGE, A. G. 1997. Testability and the refutation and corroboration of cladistics hypotheses. *Cladistics* 13:81–96.
- KNIGHT, A., AND D. P. MINDELL. 1993. Substitution bias, weighting of DNA sequence evolution, and the phylogenetic position of Fea's viper. *Syst. Biol.* 42:18–31.
- KNIGHT, A., AND D. P. MINDELL. 1995. Weighting of nucleotide sequences: A reply. *Syst. Biol.* 44:112–116.
- LORCH, P., G. S. WILKINSON, AND P. R. REILLO. 1993. Copulation duration and sperm precedence in the Malaysian stalk-eyed fly, *Cyrtodiopsis whitei* (Diptera: Diopsidae). *Behav. Ecol. Sociobiol.* 32:303–311.
- MARTIN, A. P., AND E. BERMINGHAM. 1998. Systematics and evolution of lower Central American cichlids inferred from analysis of cytochrome *b* gene sequences. *Mol. Phylogenet. Evol.* 9:192–203.
- MASON-GAMER, R. J., AND E. A. KELLOGG. 1996. Testing for phylogenetic conflict among molecular data sets in the tribe Triticeae (Gramineae). *Syst. Biol.* 45:524–545.
- MEIER, R., AND HILGER, S. (2000). On the egg morphology and phylogenetic relationships of Diopsidae (Diptera: Schizophora). *J. Zool. Evol. Research* 38:1–36.
- MICKEVICH, M. F., AND J. S. FARRIS. 1981. The implications of congruence in *Menidia*. *Syst. Zool.* 30:351–370.
- MILINKOVITCH, M. C., R. G. LEDUC, J. ADACHI, F. FARNIR, M. GEORGES, AND M. HASEGAWA. 1996. Effects of character weighting and species sampling on phylogeny reconstruction: A case study based on DNA sequence data in cetaceans. *Genetics* 144:1817–1833.
- MILLER, J. S., A. V. Z. BROWER, AND R. DESALLE. 1997. Phylogeny of the neotropical moth tribe Josiini (Notodontidae: Diopinae): Comparing and

- combining evidence from DNA sequences and morphology. *Biol. J. Linn. Soc.* 60:297–316.
- MINDELL, D. P., A. KNIGHT, C. BAER, AND C. J. HUDDLESTON. 1996. Slow rates of molecular evolution in birds and the metabolic rate and body temperature hypotheses. *Mol. Biol. Evol.* 13:422–426.
- MINDELL, D. P., AND C. E. THACKER. 1996. Rates of molecular evolution: Phylogenetic issues and applications. *Annu. Rev. Ecol. Syst.* 1996:279–303.
- MIYAMOTO, M. M., M. W. ALLARD, R. M. ADKINS, L. L. JANECEK, AND R. L. HONEYCUTT. 1994. A congruence test of reliability using linked mitochondrial DNA sequences. *Syst. Biol.* 43:236–249.
- MIYAMOTO, M. M., AND W. M. FITCH. 1995. Testing species phylogenies and phylogenetic methods with congruence. *Syst. Biol.* 44:64–76.
- MOORE, W. S. 1995. Inferring phylogenies from mtDNA variation: Mitochondrial-gene trees versus nuclear-gene trees. *Evolution* 49:718–726.
- MURPHY, W. J., AND G. E. COLLIER. 1997. A molecular phylogeny for Aplocheiloid fishes (Atherinomorpha, Cyprinodontiformes): The role of vicariance and the origins of annualism. *Mol. Biol. Evol.* 18:790–799.
- NAYLOR, G. J. P., AND W. M. BROWN. 1998. Amphioxus mitochondrial DNA, chordate phylogeny, and the limits of inference based on comparisons of sequences. *Syst. Biol.* 47:61–76.
- NIXON, K. C., AND J. M. CARPENTER. 1996. On simultaneous analysis. *Cladistics* 12:221–241.
- O'GRADY, P. M., J. B. CLARK, AND M. G. KIDWELL. 1998. Phylogeny of the *Drosophila saltans* species group based on combined analysis of nuclear and mitochondrial DNA sequences. *Mol. Biol. Evol.* 15:656–664.
- OLMSTEAD, R. G., P. A. REEVES, AND A. C. YEN. 1998. Patterns of sequence evolution and implications for parsimony analysis of chloroplast DNA. Pages 164–187 in *Molecular systematics of plants II: DNA sequencing* (D. E. Soltis, P. S. Soltis, and J. J. Doyle, eds.). Kluwer Press, Boston.
- OLMSTEAD, R. G., AND J. A. SWEERE. 1994. Combining data in phylogenetic systematics: An empirical approach using three molecular data sets in the Solanaceae. *Syst. Biol.* 43:467–481.
- PHILIPPE, H., G. LECOINTRE, H. L. V. LE, AND H. L. GUYANDER. 1996. A critical study of homology in molecular data with the use of a morphologically based cladogram, and its consequences for character weighting. *Mol. Biol. Evol.* 13:1174–1186.
- POE, S. 1996. Data set incongruence and the phylogeny of crocodylians. *Syst. Biol.* 45:393–414.
- PRESGRAVES, D. C., R. H. BAKER, AND G. S. WILKINSON. 1999. Coevolution of sperm and female reproductive tract morphology in stalk-eyed flies. *Proc. R. Soc. London Ser. B* 266:1041–1047.
- RAMOS-ONSINS, S., C. SEGARRA, J. ROZAZ, AND M. AGUADE. 1998. Molecular and chromosomal phylogeny in the *Obscura* group of *Drosophila* inferred from sequences of the *rp49* gene region. *Mol. Phylogenet. Evol.* 9:33–41.
- REMSEN, J., AND R. DESALLE. 1998. Character congruence of multiple data partitions and the origin of the Hawaiian Drosophilidae. *Mol. Phylogenet. Evol.* 9:225–235.
- RICE, W. R. 1989. Analyzing tables of statistical tests. *Evolution* 43:223–225.
- RIJSEWIJK, F., M. SCHUERMANN, E. WAGENAAR, P. PARREN, D. WEIGEL, AND R. NUSSE. 1987. The *Drosophila* homolog of the mouse mammary oncogene *int-1* is identical to the segment polarity gene *wingless*. *Cell* 50:649–657.
- RUSSO, C. A. M., N. TAKEZAKI, AND M. NEI. 1996. Efficiencies of different genes and different tree-building methods in recovering a known vertebrate phylogeny. *Mol. Biol. Evol.* 13:525–536.
- SEGUY, E. 1955. Dipteres Diopsidae africains nouveaux ou peu connus. *Bull. Inst. Franc. Afr. Noire A* 17:1102–1124.
- SHILLITO, J. F. 1950. A note on Speiser's genus *Centriuncus* and a revised definition of Diopsidae (Diptera: Acalypterae). *Proc. R. Entomol. Soc. London B* 19:109–113.
- SHILLITO, J. F. 1971. The genera of Diopsidae (Insecta: Diptera). *Zool. J. Linn. Soc.* 50:287–295.
- SIDDALL, M. E. 1995. Random Cladistics software package. Software available via ftp://zoo.toronto.edu/pub.
- SIMON, C., F. FRATI, A. BECKENBACH, B. CRESPI, H. LIU, AND P. FLOOK. 1994. Evolution, weighting and phylogenetic utility of mitochondrial gene sequences and a compilation of conserved polymerase chain reaction primers. *Ann. Entomol. Soc. Am.* 87:651–701.
- STEYSKAL, G. 1972. A catalogue of species and key to the genera of the family Diopsidae (Diptera: Acalypterae). *Stuttg. Beitr. Naturkd. Ser. A.* 234:1–20.
- SWOFFORD, D. L. 1998. PAUP*: Phylogenetic analysis using parsimony (and other methods), version 4.0. Sinauer Associates, Sunderland, Massachusetts.
- SWOFFORD, D. L., G. L. OLSEN, P. J. WADDELL, AND D. M. HILLIS. 1996. Phylogenetic inference. Pages 407–514 in *Molecular systematics* (D. M. Hillis, C. Morowitz, and B. K. Mable, eds.). Sinauer Associates, Sunderland, Massachusetts.
- VOGLER, A. P., R. DESALLE, T. ASSMANN, C. B. KNISLEY, AND T. D. SCHULTZ. 1993. Molecular population genetics of the endangered tiger beetle *Cicindela dorsalis* (Coleoptera: Cicindelidae). *Ann. Entomol. Soc. Am.* 86:142–152.
- WENZEL, J. W., AND M. E. SIDDALL. 1999. *Noise. Cladistics* 15:51–64.
- WHEELER, W. C., AND D. S. GLADSTEIN. 1994. MALIGN: A multiple sequence alignment program. *J. Hered.* 85:417–418.
- WILKINSON, G. S. 1993. Artificial sexual selection alters allometry in stalk-eyed fly *Cyrtodiopsis dalmanni* (Diptera: Diopsidae). *Genet. Res.* 62:213–222.
- WILKINSON, G. S., AND G. N. DODSON. 1997. Function and evolution of antlers and eye-stalks in flies. Cambridge Univ. Press, Cambridge.
- WILKINSON, G. S., H. KAHLER, AND R. H. BAKER. 1998. Evolution of female mating preference in stalk-eyed flies. *Behav. Ecol.* 9:525–533.
- WILKINSON, G. S., AND P. R. REILLO. 1994. Female preference response to artificial selection on an exaggerated male trait in a stalk-eyed fly. *Proc. R. Soc. London Ser. B* 255:1–6.
- YANG, Z. 1996. Among-site rate variation and its impact on phylogenetic analyses. *Trends Ecol. Evol.* 11:367–372.
- YANG, Z. 1998. On the best evolutionary rate for phylogenetic analysis. *Syst. Biol.* 47:125–133.