

Uneven Missing Data Skew Phylogenomic Relationships within the Lories and Lorikeets

Brian Tilston Smith^{1,*}, William M Mauck III^{1,2}, Brett W Benz³, and Michael J Andersen⁴

¹Department of Ornithology, American Museum of Natural History, New York, New York

²New York Genome Center, New York, New York

³Museum of Zoology and Department of Ecology and Evolutionary Biology, University of Michigan

⁴Department of Biology and Museum of Southwestern Biology, University of New Mexico

*Corresponding author: E-mail: brianstilstonsmith@gmail.com.

Accepted: 26 May 2020

Abstract

The resolution of the Tree of Life has accelerated with advances in DNA sequencing technology. To achieve dense taxon sampling, it is often necessary to obtain DNA from historical museum specimens to supplement modern genetic samples. However, DNA from historical material is generally degraded, which presents various challenges. In this study, we evaluated how the coverage at variant sites and missing data among historical and modern samples impacts phylogenomic inference. We explored these patterns in the brush-tongued parrots (lories and lorikeets) of Australasia by sampling ultraconserved elements in 105 taxa. Trees estimated with low coverage characters had several clades where relationships appeared to be influenced by whether the sample came from historical or modern specimens, which were not observed when more stringent filtering was applied. To assess if the topologies were affected by missing data, we performed an outlier analysis of sites and loci, and a data reduction approach where we excluded sites based on data completeness. Depending on the outlier test, 0.15% of total sites or 38% of loci were driving the topological differences among trees, and at these sites, historical samples had 10.9× more missing data than modern ones. In contrast, 70% data completeness was necessary to avoid spurious relationships. Predictive modeling found that outlier analysis scores were correlated with parsimony informative sites in the clades whose topologies changed the most by filtering. After accounting for biased loci and understanding the stability of relationships, we inferred a more robust phylogenetic hypothesis for lories and lorikeets.

Key words: bird, phylogeny, parrot, likelihood, museum specimen, museum DNA.

Introduction

Historical and ancient DNA from museum specimens is widely employed for incorporating rare and extinct taxa into phylogenetic studies (e.g., Thomas et al. 1989; Mitchell et al. 2014; Fortes et al. 2016). The inclusion of these samples has helped discover and delimit species (Helgen et al. 2013; Pajmans et al. 2017), resolve phylogenetic relationships (Mitchell et al. 2016), and clarify biogeographic history (Kehlmaier et al. 2017; Yao et al. 2017). DNA sequences obtained from dry and alcohol-preserved museum specimens have been collected using a range of techniques, including Sanger sequencing (Sorenson et al. 1999), restriction site-associated DNA sequencing (Tin et al. 2015; Ewart et al. 2019), and sequence capture of reduced (McCormack et al. 2016; Linck et al. 2017; Ruane and

Austin 2017) or whole genomes (Enk et al. 2014; Hung et al. 2014). However, DNA sequences collected from these museum specimens are subject to errors associated with contamination (Malmström et al. 2005), DNA degradation (Briggs et al. 2007; Sawyer et al. 2012), and low coverage in read depth (Tin et al. 2015), which all present challenges in distinguishing evolutionary signal from noise.

Sequence capture of ultraconserved elements (UCEs) is a popular approach for collecting orthologous genomic markers in phylogenomic studies (Faircloth et al. 2015; Chakrabarty et al. 2017; Esselstyn et al. 2017) and is increasingly used for historical specimens (Hosner et al. 2016; McCormack et al. 2016; Ruane and Austin 2017). A common finding is that the loci recovered are typically shorter in older

samples (Hosner et al. 2016; McCormack et al. 2016; Ruane and Austin 2017). Shorter loci are potentially problematic because the sequence capture approach targets the invariable UCE core, limiting the portion of the flanking region that contains polymorphic sites. Another factor that may cause differences among historical and modern samples is that phylogenomic pipelines that do not involve variant calling typically employ read-specific filtering, where the average read depth across all positions along a locus is used to determine whether the locus is excluded (e.g., Faircloth 2016). Under this type of scenario, low coverage characters may pass typical filters, exacerbating differences among historical and modern samples. Although some studies only use DNA sequences collected from historical or ancient samples (e.g., Hung et al. 2014), most phylogenetic approaches involving noncontemporaneous samples combine with those from modern samples. For those that do use DNA from both sample types, additional challenges in downstream analyses may arise due to an asymmetry in the phylogenetic signal caused by nonrandom missing data (e.g., Hosner et al. 2016).

The impact of missing data on phylogenetic inference remains contentious (Lemmon et al. 2009; Wiens and Morrill 2011; Simmons 2012, 2014; Hovmöller et al. 2013; Jiang et al. 2014; Streicher et al. 2016). Missing data have been shown to bias phylogenetic relationships, particularly when the missing characters are nonrandomly distributed (e.g., Lemmon et al. 2009; Simmons 2012, 2014). However, findings also suggest that even when some taxa have a large proportion of characters with no data, phylogenetic signal is retained if enough characters are present (Philippe et al. 2004; Roue et al. 2013; Shavit-Grievink et al. 2013; Molloy and Warnow 2018). Bias may manifest as inflated support values and erroneous branch lengths, or as inconsistencies between optimality criteria or phylogenomic approaches (i.e., concatenation vs. the multispecies coalescent). The increased availability of phylogenomic data has provided a more nuanced look at missing data's effect on phylogenetic inference (Philippe et al. 2004; Huang and Knowles 2016; Streicher et al. 2016; Xi et al. 2016). One means of dealing with missing data in phylogenomic data sets is to filter loci based on the proportion of either missing characters or missing species in the data set (Hosner et al. 2016). However, this approach may not directly target problematic regions of an alignment, and phylogenetically informative signal may be discarded unnecessarily. A more direct approach would entail identifying which specific sites or genes are influenced by missing data.

Analyses of outlier sites or loci in phylogenomic data indicate that a few genes can have a large impact on a topology (Arcila et al. 2017; Brown and Thomson 2017; Shen et al. 2017; Walker et al. 2018). These conflicting genealogies can be due to biological processes (e.g., incomplete lineage sorting, introgression, and horizontal gene transfer) or to spurious phylogenetic signal caused by poor alignments, paralogy, and/or sequencing error. Putative outlier loci have been

identified using topology tests (Arcila et al. 2017; Esselstyn et al. 2017), Bayes factors (Brown and Thomson 2017), and site/locus-wise log-likelihood differences among alternative topologies (Shen et al. 2017; Walker et al. 2018). Support for particular phylogenetic hypotheses may be driven by a small subset of loci (Brown and Thomson 2017; Walker et al. 2018), and the targeted removal of outlier loci can reconcile differences among topologies. Outlier analyses provide a framework for assessing how differences between historical and modern DNA sequences impact phylogenetic inference. In this study, we performed site and locus likelihood outlier analyses to evaluate whether sequence coverage and missing data impact phylogenetic relationships in our focal group, the Loriini.

Lories and lorikeets, commonly known as the brush-tongued parrots, are a speciose clade (Tribe: Loriini) of colorful birds that are widely distributed across the Australasian region (Forshaw and Cooper 1989). The characterization of distributional ranges, phenotypic variation, and systematics of the clade was the product of expansive biological inventories that peaked during the early 1900s (Mivart 1896; Forshaw and Cooper 1989). The geographical extent of this work encompasses thousands of large and small islands spread across many countries in the Australasian region. Given these immense logistical constraints, modern collecting expeditions that aim to produce voucher specimens with genetic samples for continued systematic work (e.g., Kratter et al. 2006; Andersen et al. 2017) have been much more focused in scope relative to the pioneering work of the 20th century that produced extensive series of specimens across species' entire ranges (e.g., Mayr 1933, 1938, 1942; Amadon 1943). Thus, the lack of modern genetic samples means that phylogenetic relationships in many groups, like the Loriini, remain unresolved. To get around this constraint, phylogenomic studies have sourced DNA from historical specimens to fill modern sampling gaps (Moyle et al. 2016; Andersen et al. 2018).

Prior phylogenetic work on the Loriini showed evidence for at least three paraphyletic genera (*Trichoglossus*, *Psittuteutes*, and *Chamosyna*) and highlighted the need for increased taxon and genomic sampling to fully resolve relationships among taxa (Schweizer et al. 2015). To this end, we collected UCEs from 105 described taxa in the Loriini, including species and subspecies. Our sampling design used DNA isolated from fresh tissues (hereafter modern) and historical specimens, including some over 100 years old (hereafter historical; [supplementary fig. S1](#) and [table S1](#), [Supplementary Material](#) online). We anticipated challenges with processing, recovering, and analyzing UCEs from historical specimens and expected that biases in the DNA sequence data might yield misleading relationships. To evaluate biased phylogenetic signal and explore options for maximizing the amount of data recovered, we produced alignments using different site coverage thresholds that produced alignments with varying levels of missing data. We then estimated phylogenies with particular sites and loci

removed, and with varying percentages of data completeness to evaluate topological stability. To target specific sites or loci that may be influencing relationships, we used site-wise and locus-wise likelihoods to identify which portions of the alignment drive topological differences among trees estimated with and without low coverage characters, and with and without missing data. From these analyses, we produced a series of trees with different subsets of putative outliers removed and quantified the change in topology using a tree distance metric and support values and summarized the information content of each locus. Next, we assessed whether the likelihood scores from the outlier analyses could be predicted by locus-specific alignment statistics. Finally, we took a more general approach evaluating how data completeness impacted the estimated topology by producing a series of alignments with varying levels of missing data. The alternative data reduction approaches we employed allowed us to compare the utility of precise versus general filtering of missing data on phylogenetic inference. After rigorously assessing potential biases in the data, we propose a phylogenetic hypothesis for lories and lorikeets.

Materials and Methods

We sampled all 12 genera, 58/59 species, and 102/112 named taxa (species and subspecies; Clements et al. 2019) within the Loriini, and three additional subspecies (*Glossopsitta concinna concinna*, *Glossopsitta concinna didimus*, and *Trichoglossus haematodus caeruleiceps*) recognized by Gill and Donsker (2019) and Forshaw (2010), respectively. In total, we sampled 105 taxa within Loriini. *Chamosyna diadema* is the only species not included in our study, which is extinct and known from a single female specimen (Forshaw and Cooper 1989). Two additional taxa (*Eos histrio talautensis* and *Eos squamata riciniata*) produced few loci with high missing data in those loci and were excluded from final analyses. We did not obtain samples from the following taxa: *Chamosyna rubronotata kordoana*, *Psitteuteles iris rubripileum*, *Neopsittacus pullicauda socialis*, *Eos histrio challengerii*, *Trichoglossus haematodus brooki*, and *Trichoglossus moluccanus septentrionalis*. We treated *Trichoglossus haematodus rosenbergii* as *Trichoglossus rosenbergii* and *Trichoglossus haematodus intermedius* as *Trichoglossus haematodus haematodus* following Gill and Donsker (2019). We also followed Gill and Donsker (2019) and used *Parvipsitta* for *P. pusilla* and *P. porphyrocephala*. When possible, we sampled more than one individual per species to verify the phylogenetic position of a taxon. For outgroups, we used *Melopsittacus undulatus*, *Psittaculirostris edwardsii*, and *Cyclopsitta diophthalma*, which together with the Loriini form the clade Loriinae (Joseph et al. 2012; Provost et al. 2018). Sampling map and specimen details and locality information are available in [supplementary figure 1](#) and [table S1, Supplementary Material](#) online.

We extracted total genomic DNA from muscle tissue using QIAamp DNeasy extraction kits (Qiagen, Valencia, CA). For historical samples, we used a modified DNeasy extraction protocol that used QIAquick PCR filter columns that size selected for smaller fragments of DNA. The modified protocol also included washing the sample with H₂O and EtOH prior to extracting as well as extra time for digestion. DNA extraction from historical samples was done in a dedicated lab for working with degraded samples to reduce contamination risk. We quantified DNA extracts using a Qubit 2.0 Fluorometer (Thermo Fisher Scientific). Library preparation of UCEs and enrichment, and Illumina sequencing were performed by RAPID Genomics (Gainesville, FL). The Tetrapod UCE 5K probe set was used to enrich 5,060 UCE loci (Faircloth et al. 2012). Variant bases increase with distance from the UCE core and these variant sites are phylogenetically informative (Faircloth et al. 2012). Even at shallow phylogenetic scales (i.e., within species), the majority of loci have been shown to be polymorphic with an average of two to three variant sites per locus (Smith et al. 2014). The wet-lab component of this study was carried out over 3 years and the number of individuals multiplexed per lane ranged from 48 to 384. Sequencing was done on an Illumina HiSeq 2500 PE 125 or HiSeq 3000 PE 150. Fastq files are available on the Sequence Read Archive (SRA Bioproject ID: 498485).

We used a modified data-processing pipeline that incorporated PHYLUCE (Faircloth 2016), a software package developed for analyzing UCE data, and seqcap_pop (Smith et al. 2014; Harvey et al. 2016). We used FastQ Screen to map raw reads to bacterial genomes and filter contaminant DNA (Wingett and Andrews 2018). Low-quality bases and adapter sequences were trimmed from multiplexed fastq files using Illumiprocessor v1 (Faircloth 2013; Bolger et al. 2014). Next, reads were assembled into contigs with Trinity v2.0.6 (Grabherr et al. 2011) and contigs were mapped to UCE probes. We chose the sample that produced the largest number of UCEs as the reference for subsequent mapping for all individuals. We generated an index of the reference sequence and independently mapped reads from each sample to the same reference sequence using BWA v0.7.13-r1126 (Li and Durbin 2009). SAM files produced from the BWA mapping were converted to BAM files, which were sorted with SAMtools (Li et al. 2009), and cleaned with Picard v1.106 (<http://broadinstitute.github.io/picard>). Then, we used the mpileup function in SAMtools to call variant sites and produce a VCF file (-C 30; -Q 20), vcfutils to convert from VCF to fastq (excluding sites with quality scores <20), and seqtk (github.com/lh3/seqtk) to convert fastq to fasta. From this last step, we produced two sets of DNA sequences that were analyzed independently. One data set retained all variant sites irrespective of coverage (hereafter Low Coverage data set), and a second set that excluded variant sites with <6× coverage using bcftools (hereafter Filtered data set). These collective steps produced single fasta files containing all UCE loci for

each individual sample. The following steps were independently performed for both data sets (Low Coverage and Filtered). Loci with >30% missing characters were removed from each individual before alignment. In PHYLUCe, we concatenated fasta files of each sample, aligned sequences in MAFFT (Katoh and Standley 2013), and retained loci where 75% of the samples were present in a locus for the final concatenated alignment. Both the conserved UCE core and variable flanking region of each locus were retained for all analyses.

Low Coverage and Filtered Outlier Analyses

We estimated trees with 171 tips, which included multiple individuals per taxon. This data set was used to check if samples from the same taxon grouped together as a means of identifying problematic samples. To focus on phylogenetic relationships among named taxa, all subsequent analyses are based on a reduced data set that contained one sample per taxon with 105 ingroup taxa and three outgroup samples. We estimated phylogenomic trees for both the Low Coverage and Filtered concatenated alignments containing only unique taxa in IQ-TREE (Nguyen et al. 2015) using ModelFinder (Kalyaanamoorthy et al. 2017) to select the best-fit substitution model for each locus partition (Chernomor et al. 2016). To assess support, we estimated 1,000 rapid bootstraps (BSs). The trees from the two different alignments did not produce the same phylogenetic relationships, so we performed an outlier site/locus analysis to identify which sites were causing the topologies to be different. We performed a two-topology, site-specific log-likelihood test that estimated the site-likelihoods based on the locus partition of the Low Coverage alignment using topologies estimated from the Low Coverage (T_1) and Filtered alignments (T_2) in RAxML (Stamatakis 2014). We then estimated the change in site-wise log-likelihoods (hereafter Δ s-lk = T_1 site log-likelihood – T_2 site log-likelihood). We binned putative outlier sites into bins representing Δ s-lk: >20, >10, >2, <–2, <–10, and <–20. We produced new concatenated alignments that corresponded to each Δ s-lk threshold bin, for which outlier sites were converted to ambiguous characters (N) in all individuals. This approach allowed us to estimate trees with different levels of outlier sites removed from the alignment. Next, we converted the DNA sequence of each locus alignment into only parsimony informative sites using FASconCAT-G (Kück and Longo 2014) and summarized the amount of parsimony informative sites and missing data at these sites for modern and historical samples. To visualize how different the trees were, we measured the distance among 100 BS trees using Robinson–Foulds distances (Robinson and Foulds 1981) with the multiRF function in phytools (Revell 2012) and used multi-dimensional scaling to plot the distances in two-dimensional space. All trees were processed and visualized using phytools and ape (Paradis et al. 2004) in R (R Core Team 2019). We

classified samples into two categories: 1) samples that were collected within the last 30 years and came from frozen or ethanol preserved tissue (hereafter Modern) and 2) samples that came from dry museum skins with ages ranging from the late 1800s through the 1960s (hereafter Historical). To visualize the distribution of each sample type on the tree, we colored tips blue (historical) or red (modern).

Subclade Outlier Analyses

We performed a complementary outlier analysis assessing subclades, but in this set of analyses, we compared trees estimated from alignments with and without missing data. By performing this analysis on subclades, we were able to examine how missing data impacted different portions of the tree. This approach could not be applied to alignments containing all clades because at least one individual had missing data at every site in the alignment. The six clades including a single outgroup sample were based on preliminary phylogenetic analysis and were 1) *Eos*, *Trichoglossus*, *Glossopsitta concinna*, and *Psittuteles iris* ($n = 22$), 2) *Parvipsitta* and *Psittuteles* ($n = 7$), 3) *Neopsittacus* ($n = 7$), 4) *Chalcopsitta* and *Pseudeos* ($n = 14$), 5) *Lorius* ($n = 18$), and 6) *Chamosyna*, *Vini*, and *Phigys* ($n = 32$). To produce a concatenated alignment for a clade, we followed the same steps listed above. To retain more characters in the larger clades we did not include redundant taxa or samples. We further reduced the sample size from 58 to 22 in the diverse clade containing *Eos*, *Trichoglossus*, *Ps. iris*, and *Glossopsitta* because the amount of missing data in this clade was high. We estimated subclade trees in IQ-TREE following the same procedures described above to produce alternative topologies (T_1 and T_2) estimated from alignments with (T_1) and without missing data (T_2). Each tree was rooted with a single outgroup (*Oreopsittacus arfaki* for *Chamosyna*, *Vini*, and *Phigys*; *Psittuteles goldiei* for all other clades) using phyx (Brown et al. 2017). We performed the same site-specific log-likelihood procedure described above for the two alternative topologies (T_1 and T_2), except that site-likelihoods were converted to locus-wise log-likelihoods to assess the impact of missing data across an entire locus by summing the site-likelihoods for each locus using the scripts in Walker et al. (2018). We then estimated the Δ locus-wise log-likelihood (hereafter Δ l-lk).

To explore how these putatively biased loci impacted phylogenetic inference, we grouped loci into bins representing Δ l-lk scores of >2, >10, >20, <–2, and <–10 (there were no genes where Δ l-lk < –20 in the Filtered data set). We followed the same procedure for producing concatenated alignments corresponding to each likelihood threshold bin, but for this step, we excluded the outlier loci from the global concatenated alignment. We then estimated phylogenies from each alignment to assess how sensitive phylogenetic relationships were to excluding loci in each of these approaches. If missing data bias phylogenetic relationships then the exclusion of loci with positive Δ l-lk should alter

relationships driven by missing data. The removal of loci with negative Δ I-lk will enhance relationships driven by biases in missing data.

Trees with Varying Levels of Complete Data

To determine which percentage of data completeness was necessary to produce a topology similar to the Filtered tree, we generated a series of alignments from the Low Coverage data set with increasing levels of complete data. Using trimal (Capella-Gutiérrez et al. 2009), we converted all missing characters to gaps (-) to conform to the software requirements and trimmed alignments by setting the percentage of individuals required to have an unambiguous site to retain the position in the alignment. In increments of 10%, we removed all sites where 0–100% of the sites had no missing data. This approach produced a range of alignments keeping all sites (0%) through no missing data (100%). We estimated phylogenies for each of the 11 data sets in IQ-TREE using the same approach previously described, except that we estimated the best-fit model across the entire alignment because the locus partitions were not retained after filtering.

Manipulating Modern Samples to Mimic Historical Samples

To provide a complementary approach for assessing whether missing data versus data quality were biasing phylogeny, we converted a percentage of characters in five modern samples (*Trichoglossus rubritorquis* KU22839, *Trichoglossus chlorolepidotus* DOT2422, *Trichoglossus ornatus* DOT7930, *Phigys solitarius* KU22543, and *Chamosyna placentis pallidior* DOT20055) to missing data. If the position of these samples was sensitive to the addition of missing data and the samples clustered with historical samples then missing data are the more likely culprit of the bias. We explored different percentages of missing data (50–99.9%) and converting random characters versus only parsimony informative sites. We found that the position of these samples only changed when we converted 99.9% of parsimony informative sites to missing data, and we do not present results from the other thresholds. The reason why such a high percentage was necessary is because the positions of these taxa were supported by a small number of parsimony informative sites, and the only way to influence the sites driving their relationships was to use a high threshold. We converted sites independently ten times and estimated phylogenies in IQ-TREE from the alignments using the same approach previously described.

Summarizing Phylogenetic Signal in Modern and Historical Samples

To explicitly compare the information content in DNA sequence from historical and modern samples, we calculated alignment statistics for each locus and for alignments of only

parsimony informative sites partitioned into sample types using AMAS (Borowiec 2016). To determine if older samples had more missing data, we also regressed the amount of missing data in each sample versus the age of the sample. Because we sequenced samples over multiple years, batch effects or biases attributable to differences among sequencing runs could also bias our results (Leigh et al. 2018). To provide a qualitative assessment of batch effects, we provide plots of trees where tips have been colored according to one of three plates that they were sequenced on. If there were substantial batch effects in the data then phylogenetic relationships could be, in part, due to whether samples were sequenced together.

We built a neural network in the R package caret v. 6.0.79 (Kuhn 2008) to test whether the Δ I-lk of each locus partition could be predicted by the alignment statistics. The alignment statistics (alignment length, the number of undetermined characters, the number of parsimony informative sites, the number of variable sites, and GC content) were specified as the input neurons, and the output neuron was the Δ log-likelihood. The input data were scaled to the minimum and maximum for each statistic, and the percentage of training/test data was set to 75%/25%, respectively. We produced 100 training/test data sets, independently ran each analysis, and reported mean R^2 , root-mean-square-error, and variable importance. We performed this analysis on the Low Coverage alignment that included all taxa and independently on the six subclades using the data from both filtering schemes.

Results

Data Characteristics

We sequenced 176 unique samples, including 16 that were resequenced to improve the amount of data recovered. We dropped five individuals that had aberrant relationships and long branches in the tree, patterns that were presumably driven by limited data. The final data set comprised 171 individuals (168 in the ingroup; three outgroups) where 54% and 46% were from historical and modern samples, respectively. Of the 58 species sampled, 27 had intraspecific sampling that included historical and modern samples. Historical samples on average had more reads (mean = 5.5 million; SD = 4.6 million) than modern samples (mean = 3.5 million; SD = 2.3 million), but a higher percentage of the reads in modern samples mapped to the reference (modern: mean = 87.1%; SD = 13.3%; historical: mean = 52.0%; SD = 21.7%). In modern samples, a greater number of positions were masked for having coverage $<6\times$ (modern: mean = 443,180; SD = 202,764; historical: mean = 359,958; SD = 152,304). The mean per-site coverage across individuals was similar between the two sample types (modern: mean = 67.9; SD = 25.2; historical: mean = 72.6; SD = 28.2). Additional read and locus statistics are available in [supplementary tables S2 and S3](#),

Supplementary Material online. We produced a Low Coverage data set that included all variant sites irrespective of coverage, and Filtered data set that excluded variant sites with $<6\times$ coverage. In the Low Coverage and Filtered data sets, loci had a mean length of 498 bp (range: 140–1,708 bp) and 482 bp (range: 105–1,413 bp), respectively. The mean and range number of taxa per gene was 164 for the Low Coverage (128–171) and 152 for the Filtered (5–171) data sets. After retaining loci where 75% of the individuals were present in any one locus, the Low Coverage data set had 4,208 loci, 2,105,994 bp, and 47,338 parsimony informative sites, whereas the Filtered concatenated alignment had 3,765 loci, 1,917,997 bp, and 39,404 parsimony informative sites. Additional supplementary data are available at <https://doi.org/10.5061/dryad.n5tb2rbps>.

Overall, the Low Coverage data set had more parsimony informative sites than the Filtered data set (fig. 1). A comparison of sample types shows that the range in the number of parsimony informative sites among loci was lower in the modern samples in contrast to the historical samples (fig. 1A and B). In the Filtered data set (fig. 1B), there was greater variability in the number of samples per locus in historical samples than the Low Coverage data set (fig. 1A). For each alignment type, the modern samples contained a greater number of parsimony informative sites and less missing data than the historical samples (1.7 \times and 2.1 \times more parsimony informative sites in the Low Coverage and Filtered data sets, respectively; fig. 1C and D). In the Filtered data set, the number of parsimony informative sites dropped, and the range of the number of individuals in each locus alignment increased. Plotting non-parsimony informative sites and missing data at those positions showed a similar pattern where there was more missing data in historical samples (supplementary fig. S2, Supplementary Material online). The percentage of missing data in samples in the Filtered data set decreased with specimen age (adjusted $R^2 = 0.31$; $n = 144$; P value < 0.0001 ; supplementary fig. S2, Supplementary Material online).

Resolving Phylogenomic Relationships among Lorikeets

The backbone phylogeny we inferred for the Loriini generally had high support and the placement of genera was stable (supplementary figs. S3 and S4, Supplementary Material online). Summarizing higher-level relationships, *Oreopsittacus* was sister to all other ingroup taxa, then *Charmosyna* was sister to the clade containing *Neopsittacus*, *Lorius*, *Pseudeos*, *Chalcopsitta*, *Psittuteles*, *Glossopsitta*, *Eos*, *Trichoglossus*, and *Parvipsitta*. The placements of *Neopsittacus*, *Lorius*, *Pseudeos*, and *Chalcopsitta* were well supported in the tree, and each of these genera was monophyletic. *Trichoglossus*, *Charmosyna*, and *Psittuteles* were not monophyletic. *Psittuteles* was found in three separate places in the tree: *Psittuteles versicolor* was sister to the recently erected genus *Parvipsitta*; *Ps. iris* was nested within a clade of *Trichoglossus*

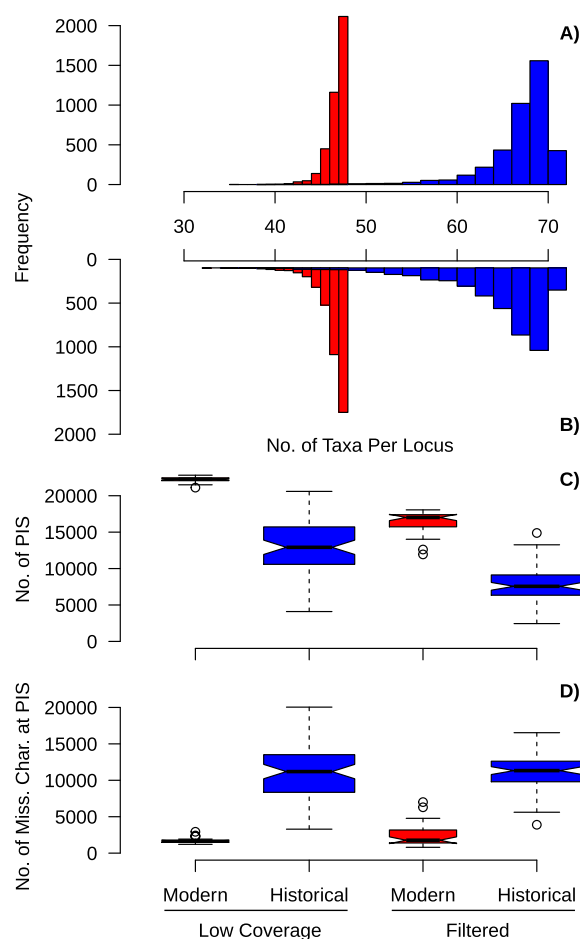


Fig. 1—Modern samples have more parsimony informative sites (PIS), less missing data at PIS, and less variation in number of samples among loci. Shown are histograms of the number of samples per locus in the Low Coverage (A) and Filtered (B) alignments. (C, D) Boxplots showing the number of parsimony informative sites (C) and number of missing characters at parsimony informative sites (D) in the ingroup samples. The data are partitioned into the modern versus historical samples, and Low Coverage versus Filtered alignments. In all plots, modern samples are shown in red and historical samples in blue.

taxa that are from Indonesia; and *Psittuteles goldei* was sister to the clade containing *Glossopsitta*, *Eos*, *Trichoglossus*, and *Ps. iris*. *Vini* and *Phigys* are strongly supported as nested within *Charmosyna*. Relationships within *Charmosyna* (including *Vini* and *Phigys*) and *Chalcopsitta* were generally stable across filtering schemes, as were relationships of the less diverse clades (*Oreopsittacus*, *Neopsittacus*, and *Parvipsitta*). Within the remaining clades, there were several notable differences in topological relationships among the Low Coverage and Filtered trees.

The Filtered tree has four clades containing *Trichoglossus*, *Eos*, *Ps. iris*, and *G. concinna* with varying levels of support (fig. 2A). *Glossopsitta concinna* was sister to a clade containing a monophyletic *Eos*, *Trichoglossus*, and *Ps. iris*. Within this

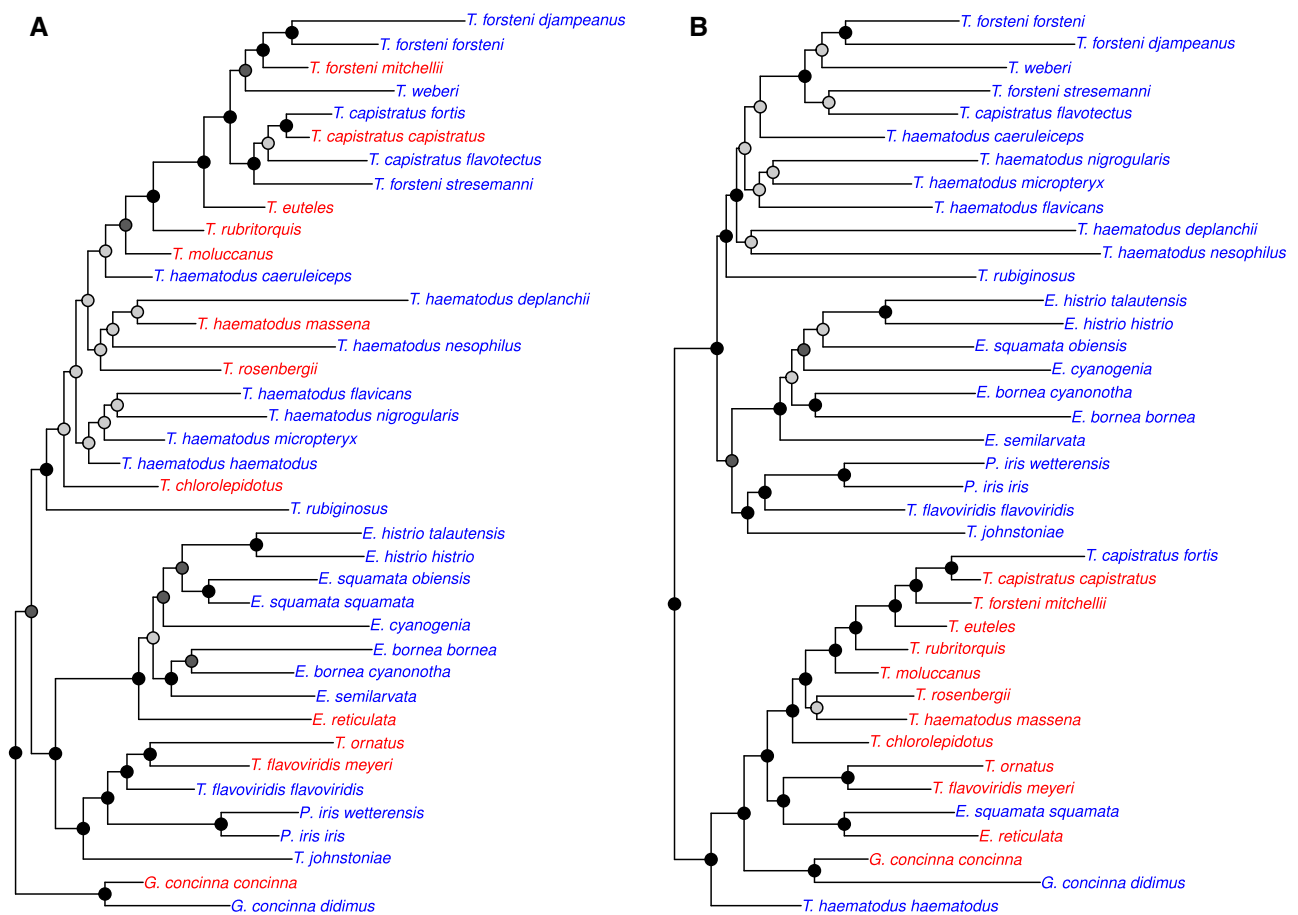


FIG. 2—Alternative topologies for the subclade that differs the most among filtering schemes. Shown is the subclade containing *Trichoglossus*/*Eos*/*Psitteuteles iris*/*Glossopsitta* from trees estimated without (A: Filtered Tree) and with low coverage characters (B: Low Coverage Tree). In the Low Coverage tree are clades composed of mostly historical versus modern samples. Bootstrap nodes are colored on a gradient from 100% (black) to <70% (gray). Taxon names are colored according to whether their DNA came from modern tissues (red) or historical specimens (blue).

tree, *Eos* was monophyletic and sister (BS = 100%) to a clade containing *Trichoglossus* taxa that occur in Indonesia and the Philippines (*T. ornatus*, *Trichoglossus flavoviridis*, and *Trichoglossus johnstoniae*) and *Ps. iris*. The *Eos*, *Trichoglossus*, and *Ps. iris* clade was sister (BS = 87%) to a clade containing the remaining *Trichoglossus*, which was supported by a BS value of 92%. This *Trichoglossus* clade had several short internodes and poorly supported relationships, particularly among *Trichoglossus haematodus* subspecies, which primarily came from historical samples. *Trichoglossus euteles*, *Trichoglossus forsteni*, *Trichoglossus capistratus*, *Trichoglossus weberi*, and *Trichoglossus rubritorquis* are nested within *T. haematodus*. *Trichoglossus forsteni stresemanni* was more closely related to *Trichoglossus capistratus* than to other *T. forsteni* taxa. In contrast, the Low Coverage tree has two well-supported (BS \geq 95%) clades composed of *Trichoglossus*, *Eos*, *Ps. iris*, and *G. concinna* (fig. 2B). One clade consists of entirely historical samples ($N=23$), whereas the other was primarily modern samples (13/16). Within each of these clades, tips have similar relationships among taxa as

seen in the Filtered tree. *Trichoglossus* taxa that occur in Indonesia or the Philippines and *Ps. iris* are sister to *Eos* and the remaining *Trichoglossus* form a clade with the exception of one historical sample (*T. haematodus haematodus*). Support values are higher in the clade composed of mostly modern samples. In both trees (Low Coverage and Filtered), *Charmosyna* was composed of four clades, *Charmosyna wilhelminae* was sister to all other taxa in the clade, *Charmosyna rubronotata* and *Charmosyna placentis* are sister and form a clade, *Charmosyna multistriata* was sister to *Charmosyna josefinae* and *Charmosyna papou*, and the remaining *Charmosyna* taxa (*Charmosyna margarethae*, *Charmosyna rubrigularis*, *Charmosyna meeki*, *Charmosyna palmarum*, and *Charmosyna amabilis*) and *Ph. solitarius* and *Vini* form a clade. The position of *Charmosyna pulchella* and *Charmosyna toxopei* will be discussed below.

Similar clustering patterns based on sample type (historical vs. modern) are observed in *Lorius*, *Vini*, and *Charmosyna* in the Low Coverage trees (supplementary fig. S3, Supplementary Material online). The two subspecies of

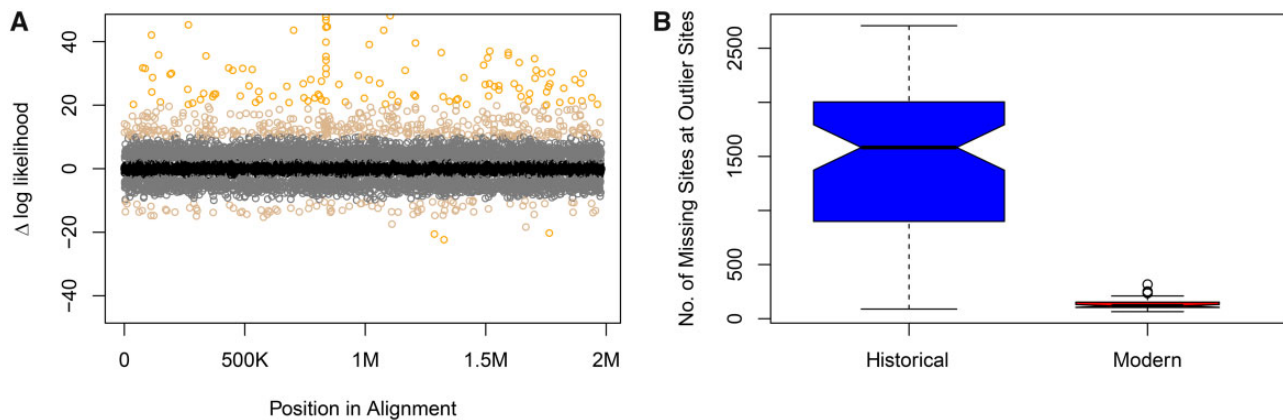


Fig. 3—Outlier sites have high missing data in historical samples. (A) Outlier site plot showing Δ sites-wise log-likelihoods (Δ s-lk) for topologies estimated with and without low coverage sites. The y axis is the Δ s-lk score and the x axis represents individual sites in the concatenated alignment, where K and M represent thousand and million, respectively. Points are colored according to the magnitude of the Δ site-wise log-likelihood scores according to a gradient reflecting the different likelihood thresholds (>2 , >10 , >20 , <-2 , <-10 , and <-20). (B) Boxplot of historical (blue) and modern (red) samples showing the amount of missing data in the 3,084 outlier sites (Δ s-lk > 2) identified in plot A.

Lorius lory that come from modern samples are sister taxa. The one modern sample of *C. placentis* was sister to *Charmosyna rubronotata*. *Charmosyna palmarum* (a modern sample) was strongly supported as sister to *Phigys* and *Vini*. The three *Vini* from historical samples group together. The two *C. papou* subspecies that come from historical samples are sister to a clade containing the remaining *C. papou* subspecies. None of these relationships are observed in the Filtered tree.

A qualitative assessment of batch effects, by coloring each tip in the tree according to sequencing run, did not detect biases whereby samples would have clustered together based on sequencing plate (supplementary fig. S15, Supplementary Material online). In the Low Coverage tree (supplementary fig. S6A, Supplementary Material online), biases in clustering were more apparent when tips are colored according to whether the sample came from a modern or historical source, which was not observed in the Filtered tree (supplementary fig. S6B, Supplementary Material online).

Outlier Sites and Loci

The outlier analyses assessing the change in site-likelihoods scores between the Low Coverage versus the Filtered topology identified 3,084 (3,084 sites: Δ s-lk > 2 ; 473 sites: Δ s-lk > 10 ; 112 sites: Δ s-lk > 20) and 1,925 (1,925 sites: Δ s-lk < -2 ; 89 sites: Δ s-lk < -10 ; and three sites Δ s-lk, -20) outlier sites in the alignment (1,980,082 bp) with positive and negative Δ s-lk values, respectively (fig. 3A). Higher and more positive Δ s-lk are sites that better support the topology estimated from the Filtered alignment, and lower and more negative values favor the tree estimated from the Low Coverage alignment. The 1,925 outlier sites with negative Δ s-lk were found on 1,381 loci, and the 3,084 outlier sites with positive values were on 1,878 loci.

The 3,084 sites with Δ s-lk > 2 , which favored the topology of the Filtered tree, exhibited a disproportionate number of missing sites in the historical versus modern samples (fig. 3B). We plotted Δ s-lk scores versus the best-fit nucleotide substitution models from IQ-TREE to assess whether there was a relationship between particular models and the extent of the score but we observed that high and low Δ s-lk were found across a wide array of models (supplementary fig. S7, Supplementary Material online).

Overall, the subclade outlier analyses for the Low Coverage alignment identified more outlier loci (fig. 2 and supplementary fig. S8, Supplementary Material online). There were 61/2 (Low Coverage/Filtered), 396/47, and 1,608/431 loci in the three bins (Δ l-lk of >20 , >10 , and >2), respectively (fig. 4). There were 121/11 (Low Coverage/Filtered) and 1,309/255 loci in the three bins (Δ l-lk of <-10 and <-2), respectively. The maximum and minimum Δ l-lk were much higher in the Low Coverage (in *Trichoglossus/Eos/Psitteuteles/Glossopsitta*: Δ l-lk = -33.089 to 394.392) versus the Filtered data set (in *Trichoglossus/Eos/Psitteuteles/Glossopsitta*: Δ l-lk = -15.742 to 33.106). There were 1,164 loci identified by both the Low Coverage versus Filtered tree and subclade outlier analyses. In the Low Coverage and Filtered analyses, the outlier sites were found on 444 loci uniquely identified, and 740 loci identified by the subclade clade analyses.

We found that by converting parsimony sites to missing data, modern samples could cluster with historical samples. The extent of the shift of the sample and the position of the manipulated sample in the tree varied across the trees (supplementary fig. S9A–H, Supplementary Material online). For example, *T. ornatus*, which was strongly supported as sister to *Trichoglossus flavoviridis* in the Low Coverage Tree (supplementary figs. S3 and S9A, Supplementary Material online), was nested within the clade containing only historical samples

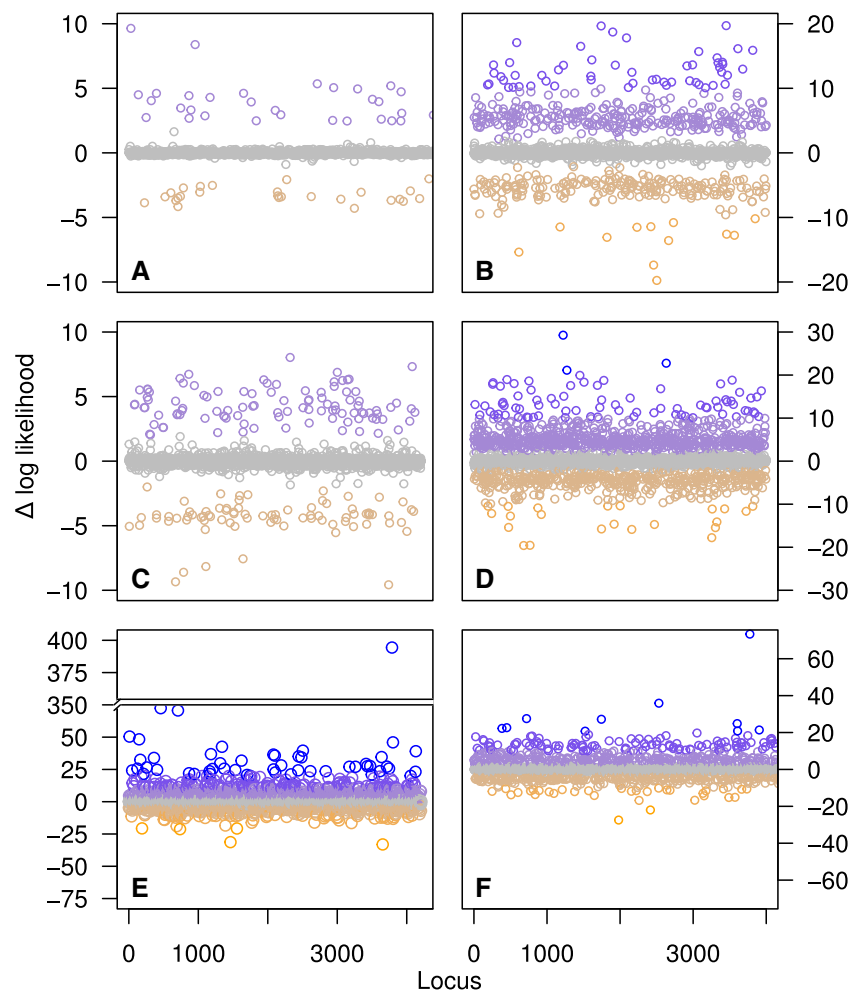


FIG. 4—Likelihood plots showing Δ locus-wise log-likelihood (Δ I-lk) for topologies estimated with and without missing data for the Low Coverage data set. The y axis is the Δ I-lk and the x axis represents individual loci across the full alignment. Shown are the results for six subclades assessed within Loriini using the Low Coverage data set: (A) *Parvipsitta* and *Psitteuteles*, (B) *Chalcopsitta* and *Pseudeos*, (C) *Neopsittacus*, (D) *Charmosyna*, *Vini*, and *Phigys*, (E) *Eos*, *Trichoglossus*, *Glossopsitta concinna*, and *Psitteuteles iris*, and (F) *Lorius*. Points are colored according to the magnitude of the Δ I-lk scores according to a gradient ranging from >20 (blue) through <-10 (orange).

in the *Trichoglossus/Eos/Psitteuteles* clade in some of the trees with manipulated sequences (supplementary fig. S9E and I, Supplementary Material online). *Charmosyna placentis pallidior* was strongly supported as sister to *Charmosyna rubronotata rubronotata*, and when most of its parsimony informative sites are converted to missing data it is nested within in its correct position in *C. placentis* (supplementary fig. S9D, E, and I, Supplementary Material online). In some trees, the position of taxa (e.g., *Trichoglossus chlorolepidotus*) did not change at all (supplementary fig. S9B–F, Supplementary Material online), and in others, the same taxon placed well outside their clade (supplementary fig. S9I and J, Supplementary Material online).

According to our neural network, alignment statistics predicted $\sim 4\%$ of the variation of Δ I-lk scores in the Low Coverage versus Filtered trees (mean and SD; $R^2 = 0.04$ [0.02]; RMSE = 0.02 [0.006] Δ I-lk scores). In the model, GC

content (mean variable importance: GC content = 31.04) and the number of parsimony informative sites (23.07) were more important than the other statistics (alignment length = 18.90; no. of taxa = 18.90; undetermined characters = 11.59; and no. variable sites = 1.30). For the neural networks on the six subclades (supplementary table S4, Supplementary Material online), the models for the *Eos/Trichoglossus/Glossopsitta/Psitteuteles* predicted $\sim 9\%$ of the variation in Δ I-lk scores ($R^2 = 0.088$) and the most important variable in the model was parsimony informative sites (Low Coverage: 52.57; Filtered: 60.72). The *Eos/Trichoglossus/Glossopsitta/Psitteuteles* clade had the most variable topology among filtering schemes. For the remaining subclades, the neural nets performed poorly (supplementary table S4, Supplementary Material online) or had positive R^2 values for clades with limited variation in Δ I-lk scores.

Impacts of Filtering Sites and Loci

The removal of outlier loci with positive Δ I-lk scores broke up some of the same-type clusters, particularly at a threshold value of all loci with $\Delta > 2$ (supplementary fig. S3B–D, Supplementary Material online). However, the removal of this many loci ($n = 1,608$) also reduced the support for other nodes in the tree. Trees estimated with the removal of negative outlier loci retained the apparent sample-type clusters (supplementary fig. S3A, E, and F, Supplementary Material online). Individual taxa whose position varied the most among filtering schemes were *G. concinna* and *Trichoglossus rubiginosus*. In the Filtered data set, which did not exhibit the sample-type clusters, the removal of outlier loci (Δ I-lk > 2 ; supplementary fig. S4D, Supplementary Material online) increased the support for the placement of *G. concinna* as sister to the clade containing *Trichoglossus*, *Eos*, and *Ps. iris*. In contrast, the removal of outlier loci (Δ I-lk < -2) placed *G. concinna* within the clade containing *Trichoglossus*, *Eos*, and *Ps. iris*. This placement received moderate BS support for either being sister to the clade containing *T. haematodus* and allies or the entire clade containing *Trichoglossus*, *Eos*, *Ps. iris*, and *G. concinna*. *Lorius lory* has seven subspecies, which formed a well-supported clade in the Filtered tree (supplementary fig. S9, Supplementary Material online), with the exception of *Lorius lory viridicrissalis*, whose placement was equivocal. The Low Coverage tree has *L. lory viridicrissalis* within the *L. lory* clade with low support (supplementary fig. S3, Supplementary Material online). Filtering of outlier loci changed support values but never unequivocally placed *L. lory viridicrissalis* within *L. lory*. *Charmosyna pulchella* and *C. toxopei* are sister taxa, however, their position within *Charmosyna* varied across trees. Trees estimated with all loci or loci with negative Δ I-lk scores excluded had these taxa as sister (often with high support) to the clade containing the subclades *Charmosyna multistriata*; *C. josefinae* and *C. papou*; and *Charmosyna margarethae*, *Charmosyna rubrigularis*, *Charmosyna meeki*, *C. palmarum*, *Charmosyna amabilis*, *Ph. solitarius*, and *Vini* (supplementary fig. S3E and F, Supplementary Material online). Alternatively, trees where positive Δ I-lk scores > 2 were excluded had these taxa as sister, albeit with lower support (BS = 63%) to a clade containing *Charmosyna*, *Phigys*, and *Vini* (supplementary fig. S3D, Supplementary Material online).

Examining the differences among topologies in multidimensional space showed distances among trees change across filtering schemes (fig. 5). In the trees where outlier sites were excluded (fig. 5A), the Robinson–Foulds distances among the Low Coverage and Filtered trees decreased (supplementary fig. S10, Supplementary Material online). At a filtering threshold of Δ s-lk > 2 (3,084 sites), the distance between the two trees was minimal (fig. 5 and supplementary fig. S10D, Supplementary Material online). Filtering sites with negative Δ s-lk values maintained a topology similar to the

Low Coverage tree, except at a threshold of Δ s-lk < -2 , which was distant from all other trees in multidimensional space. In the subclade analyses, the filtering of loci did not yield similar topologies between the Low Coverage and Filtered trees (fig. 5B). However, the Low Coverage tree where loci were excluded at a threshold of Δ s-lk > 2 was the least distant from the Filtered tree (fig. 5B). The Low Coverage trees with all loci and Δ s-lk < -2 produced similar topologies and were the most distant from the Filtered trees (fig. 5B). Despite some differences in the placement of taxa across the Filtered trees, the Robinson–Foulds distances among trees were comparatively low.

Impact of Data Completeness

The alignment length ranged from 2,105,994 bp (0% or all sites) through 41,504 bp (100% or no sites with missing data [table 1]), and the trees estimated from these alignments are in supplementary figure S11A–K, Supplementary Material online. Across this same range of filtering, there were 30,380 (0%) through 373 (100%) parsimony informative sites (table 1). At 60% completeness, the sample-type clusters started to break-up (supplementary fig. S11G, Supplementary Material online), and at 70% the tree was similar to the Filtered tree (supplementary figs. S4A and S11H, Supplementary Material online). By 90% completeness, some relationships differed from the Filtered tree (supplementary figs. S4A and S11J, Supplementary Material online), and by 100% the tree had lower resolution and support (supplementary fig. S11K, Supplementary Material online). Between the 70% and the 90% data completeness threshold, the alignment was reduced from 1,186,107 to 800,137 bp and 15,404 to 8,632 bp parsimony informative sites.

Discussion

We showed that systematic bias caused by missing informative sites between DNA sequences from modern versus historical specimens can produce aberrant or unstable phylogenetic relationships. To obtain dense taxon sampling in our focal group, the Loriini, we leveraged samples collected over the last 100+ years and assessed how this sampling scheme impacted phylogenetic relationships by producing alignments with low coverage characters included and excluded. These two trees exhibited some striking differences. In the Low Coverage tree, there were numerous cases where historical or modern samples clustered together that were not observed in the Filtered tree (e.g., fig. 2). We employed a targeted and general approach to assess how missing data were influencing these unexpected relationships. The targeted method using a site outlier analysis showed that a small number of sites were driving the topological differences, and at these sites, historical samples had substantially more missing data (fig. 1). Excluding low coverage characters reduced

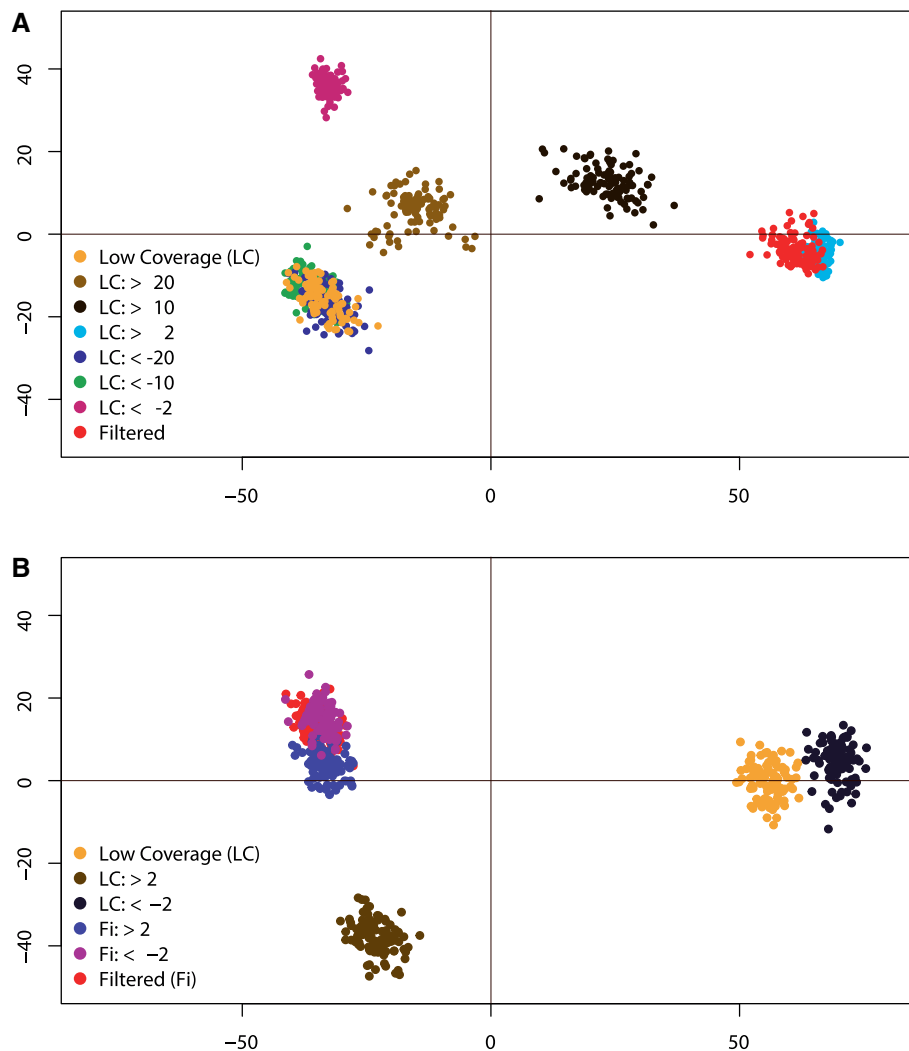


Fig. 5—Multidimensional scaling of Robinson–Foulds distances among 100 bootstrap trees with differing levels of outlier sites or loci excluded. (A) Compares distances among Filtered and Low Coverage trees where outlier sites have been removed at different increments. Outlier sites were excluded in the Low Coverage alignment using Δ site-wise log-likelihood (Δ s-lk) thresholds of >20, >10, >2, <−2, <−10, and <−20. (B) The distances among trees produced from the subclade outlier analyses. Shown is a comparison of the Low Coverage and Filtered trees with topologies estimated with outlier loci excluded using Δ locus-wise log-likelihood (Δ l-lk) thresholds of >2 and <−2.

the discrepancy in missing data between historical and modern samples, and at this level of disparity, the tree did not contain sample-type clusters and was similar in topology to the Filtered tree (fig. 5). A more nuanced look at outlier loci within subclades showed that many loci supported alternative topologies when sites with missing data were excluded, and the position of some branches shifted when these loci were excluded (supplementary fig. S3, Supplementary Material online). Biases in the historical samples could also be observed in modern samples by dropping the majority of their parsimony informative sites, which produced similar sample-type clusters observed in the Low Coverage tree (supplementary fig. S9, Supplementary Material online). The neural network was able to show that the number of parsimony informative sites could

predict likelihood scores for the clade most impacted by missing data (*Eos/Trichoglossus/Glossopsitta/Psitteuteles*; supplementary table S4, Supplementary Material online), but for most clades, which did not have as many outlier loci, the models were poor fits to the data. The more general approach of data reduction using the percentage of data completeness indicated that sites with high data completeness were necessary to avoid spurious relationships, but more stringent conditions of data completeness produced less-resolved trees. After accounting for biased loci and understanding the stability of nodes, we inferred a more robust phylogenetic hypothesis for the Loriini. Taxonomic relationships within the clade can now be revised to reflect natural groupings, but for some groups, additional work is still necessary.

Table 1

Data Completeness, Alignment Length, and Number of Parsimony Informative Sites at Differing Thresholds of Missing Data Allowance

% of Data		
Completeness	Alignment Length (bp)	PIS
0	2,105,994	30,382
10	1,901,418	30,162
20	1,786,228	28,755
30	1,744,235	27,936
40	1,661,160	25,796
50	1,498,337	21,765
60	1,354,025	18,811
70	1,186,107	15,404
80	1,024,281	12,341
90	800,137	8,632
100	41,504	372

Note.—Shown are the percentage of individuals at each site with nonambiguous characters across the Low Coverage alignment. As the alignment length and number of parsimony informative sites (PIS) decrease, the percentage of data completeness increases and more characters are excluded.

Asymmetric Information Content among Sample Types

We found that alignments with high missing data produced biased phylogenetic relationships. In these trees, subclades consisting of mostly modern samples presumably formed because there was not enough information to place historical samples among the modern samples. Our analyses suggest that an asymmetry in phylogenetic information content among sample types is the primary culprit of the bias because only 3,084 sites (0.15% of total sites) drove the topological differences among trees, and the historical samples had 10.9× more missing data at these sites (fig. 3). By filtering for data completeness, we produced a similar result and inferred the expected phylogeny by including only sites where 70% of the individuals had unambiguous characters. Previous work has shown that ambiguous characters can bias the probability of taxa being sister (Lemmon et al. 2009) and increase the resolution and support of clades (Simmons 2012, 2014). These previous studies did not deal with historical versus modern samples and did not have the magnitude of characters in our data set, but a similar mechanism is likely operating. Although we accounted for among-site rate variation, which has been shown to lead to biases in missing data (Lemmon et al. 2009), we did not evaluate how multispecies coalescent approaches would deal with our data set. We concentrated instead on a concatenated approach because our data met two criteria in which species-tree summary methods perform poorly (Molloy and Warnow 2018); namely, our data comprised 1) many poorly resolved gene trees with high missing data from 2) loci with low information content found in UCEs.

By including low coverage characters, we were able to explore potential biases that can arise between historical and modern samples. Filtering according to a read coverage

threshold at each variant site is common practice in population genomic studies (e.g., Thom et al. 2018), but this approach is less frequently employed in phylogenomic bioinformatic pipelines (e.g., Faircloth 2016). In the Low Coverage tree, we found clusters of historical or modern samples that were not present in the Filtered tree (fig. 2 and [supplementary figs. S3 and S4, Supplementary Material online](#)). Besides an asymmetry in informative sites, these clusters could be caused by sequencing errors present in one sample type, batch effects, or contamination. We address sample type in detail below, but we suspect that biases of batch effects and contamination were minimal. For example, we had limited power to test for batch effects because we did not randomly and evenly sequence samples across runs, therefore, there are portions of the tree where clades are composed almost entirely of samples from the same sequencing run ([supplementary fig. S5, Supplementary Material online](#)). In these cases, we do not interpret these patterns as batch effects because the tips occur in their expected topological position and samples from different sequencing lanes are distributed throughout the tree. We took great care to avoid contamination during wet-lab procedures (Mundy et al. 1997) and we have no strong reason to suggest that contamination is driving the observed pattern, particularly after exploring the impacts of missing data on the topology. The impact of contamination may have been more pronounced on low-quality characters, which were filtered out in all treatments because unreported preliminary trees estimated with these low-quality characters produced trees with long branches. However, more subtle effects of contamination on a small number of characters may not be directly detectable in the approaches we employed. Although we cannot rule out additional artifacts caused by contamination or sequencing error, the topology within each of the most apparent sample-type clusters in the *Trichoglossus/Eos/Psitteuteles* clade exhibited the expected relationships among taxa.

The outlier analysis on subclades also found loci that were impacted by missing data. In the Low Coverage tree, the sample-type clusters were broken up when outlier loci with positive values were excluded but also reduced support values ([supplementary fig. S3, Supplementary Material online](#)). The exclusion of outlier loci with negative values retained the biased relationships. These loci had negative values because the topologies estimated with all sites with missing data removed were more likely given the alignment. The removal of these loci produced alignments that only included loci that were either biased or not impacted by missing data. In the Filtered data set, the number of identified outlier loci was reduced and exclusion of outlier loci was less profound. Nonetheless, the removal of outlier loci in the Filtered data set showed how the placement of *G. conccina*, *Trichoglossus rubiginosus*, and the clade containing *C. pulchella* and *C. toxopei* was sensitive to missing data ([supplementary fig. S4, Supplementary Material online](#)). Interestingly, about 72%

of the loci identified by the subclade outlier analyses were the same loci of the outlier sites identified by the Low Coverage versus Filtered outlier analysis. The information content of the nonoverlapping loci is important because the more targeted site-wise outlier analysis was better at reconciling topological differences among the Low Coverage and Filtered trees than was the subclade approach.

There was a tendency for historical samples to fall outside of their clade or even the ingroup, as evident in previous phylogenomic studies on birds (Hosner et al. 2016; Moyle et al. 2016; Andersen et al. 2019; McCullough et al. 2019). This was the case for seven of our excluded samples, which produced limited data and could not be accurately placed in their genus or higher-level clade. The sample-type cluster within *Trichoglossus/Eos/Psitteuteles* is an extreme example of this pattern, and the pattern is so striking because of the high number of historical samples in this particular clade (fig. 2). In the trees wherein a subset of modern samples we converted most parsimony informative sites to missing data, we observed the same pattern whereby some of the manipulated samples were inferred outside of their expected clade (e.g., [supplementary fig. S9J](#), [Supplementary Material](#) online). Without prior information on whether a taxon is sister or falls outside of a clade of closely related taxa, samples with high missing data in large alignments may not be able to be accurately placed on a phylogeny.

Identifying Biased Samples and Loci

Our neural network models were good predictors of Δ I-lk scores in some tests, but not others. A factor to consider for interpreting our model results is that the range in Δ I-lk scores varied substantially among clades, and typically the models that performed poorly were for the clades with low variation in Δ I-lk scores. In contrast, the *Eos/Trichoglossus/Glossopsitta/Psitteuteles* clade, which had the widest range in Δ I-lk scores, had the best performing models. In both the Filtered and Low Coverage data sets, parsimony informative sites were the most important variable in the models for *Eos/Trichoglossus/Glossopsitta/Psitteuteles*, suggesting that missing information at these sites in historical samples influenced the topological differences. The neural network for Δ I-lk values estimated between the Low Coverage and Filtered trees explained 4% of the variation in outlier scores, and GC content was the most important variable in the model, followed by parsimony informative sites. However, the outlier sites on these loci had high missing data, and when these sites were removed, the estimated phylogeny was similar to relationships in the Filtered tree. Because the magnitude of the Δ I-lk score is going to be partially dictated by how much information there is at a site or across a locus, the outlier analysis is expected to identify sites or loci that have enough information to distinguish alternative trees. Missing data at less informative sites is also known to bias phylogenetic inference (Simmons 2012,

2014), and the outlier analysis we used may not capture the full extent of missing data on our inferred phylogenies.

Preferentially selecting phylogenetically informative loci is expected to produce trees with better support (Gilbert et al. 2018), but our results suggest that this practice can produce less reliable relationships when the data content dramatically varies among samples. Other work has shown that filtering phylogenomic markers by information content had mixed results in terms of resolving discordance among trees estimated with different phylogenetic methods (McClean et al. 2019). Outlier analysis using site and locus likelihood scores (Shen et al. 2017; Walker et al. 2018) provides a rapid means of identifying loci that have a large impact on phylogeny reconstruction, but, as we showed, the resolution of this approach will depend on the trees that are available for comparison (e.g., the a priori expected phylogeny vs. an alternative phylogeny). As mentioned above, a targeted outlier approach will not address all potential biases that missing data can cause, but it can identify sites that are having a strong influence on the phylogeny. Despite the limitations of site-likelihoods, the precision of identifying specific sites/loci may be the more favorable option to filtering data because the alternative of using percentage of data completeness to remove sites resulted in removing positions in the alignment that were important for other portions of the tree. This idiosyncratic behavior of filtering for data completeness to achieve higher topological support for one recalcitrant historical sample occurred in a recent study of honeyeaters. Andersen et al. (2019) increased the filtering stringency toward more complete data sets to improve support for *Gymnomyza aubryana*, however, previously well-supported nodes elsewhere in the tree were negatively impacted due to a reduction in total parsimony informative sites. The optimal percentage of data completeness will vary among data sets and depend on how asymmetric the information content is among sample types. For our data set, there was a narrow window for when data completeness produced a reliable phylogeny (e.g., 70% vs. 90%; [supplementary fig. S11H](#) and [J](#), [Supplementary Material](#) online) because data completeness >70% led to a less-resolved tree.

Taxonomic Implications

Our study builds on previous phylogenetic work on the Loriini by further clarifying relationships and adding 64 previously unsampled taxa (fig. 6). We inferred a backbone phylogeny of relationships among genera that was fairly well resolved with the exception of the clade containing *Trichoglossus*, *Ps. iris*, *Eos*, and *Glossopsitta*, and some nodes in *Charmosyna*. Our analyses corroborated recently proposed taxonomic changes where *Pseudeos cardinalis* was moved into *Pseudeos* from *Chalcopsitta*, and *Parvipsitta* was resurrected to contain *P. pusilla* and *P. porphyrocephala*, which were previously placed in *Glossopsitta* (Schweizer et al.

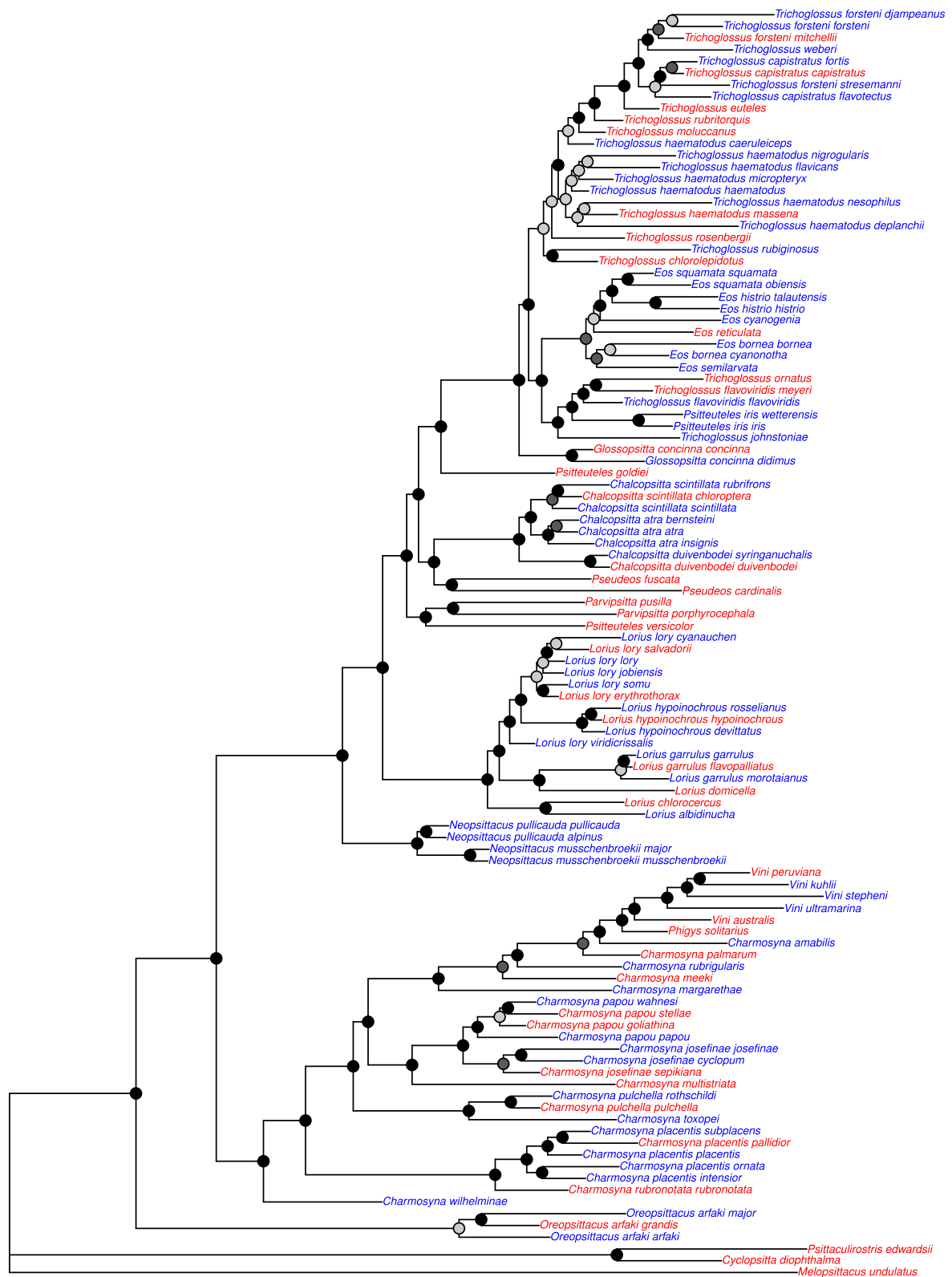


FIG. 6—Maximum likelihood tree containing unique taxa in Loriini. The tree was inferred from a concatenated alignment where loci identified with the locus likelihood analysis with Δ locus-wise log-likelihood (Δ l-lk) values of >10 were excluded. On each node are shown rapid bootstrap values and the taxon names are colored according to whether their DNA came from modern tissues (red) or historical specimens (blue). Bootstrap nodes are colored on a gradient from 100% (black) to $<70\%$ (gray).

2015). In all of our trees, *Pseudeos fuscata* and *Pseudeos cardinalis* were sisters and were in turn sister to *Chalcopsitta*. *Parvipsitta pusilla* and *P. porphyrocephala* were sisters and not closely related to *G. concinna*. However, we found strong support for *P. pusilla* and *P. porphyrocephala* being sister to *Ps. versicolor*, a novel result. *Psitteuteles versicolor* and *Parvipsitta* could be subsumed under a single genus. Irrespective of this taxonomic decision, the polyphyly of *Psitteuteles* will require that *Ps. goldei* and *Ps. iris* be moved into new genera. *Psitteuteles goldei* is sister to the clade containing *Trichoglossus*, *Eos*, *Ps. iris*, and *Glossopsitta*. The taxonomic revision of *Ps. iris* will depend on how *Trichoglossus* is treated because *Ps. iris* is nested within a geographically coherent clade of taxa distributed largely to the west of New Guinea. The clade containing *Charmosyna*, *Phigys*, and *Vini* represents a deep, diverse, and geographically widespread group. The species in these genera are varied in terms of body size and shape, tail length, plumage color, and sexual dimorphism (Forshaw and Cooper 1989; Merwin et al. 2020), and these morphological traits are not found in monophyletic groups in our phylogeny. Species-level relationships among species in *Charmosyna* were well supported and stable with the exception of the placement of *C. toxopei* and *C. pulchella*. Overall, the taxonomic revision of this clade will present challenges regarding when and where to split or lump taxa and how best to circumscribe genera.

Relationships among subspecies within species varied substantially among taxa. Support for relationships among species within *Lorius* were generally stable except the placement of *L. lory viridicrissalis*, which was only nested within *L. lory* in the Low Coverage tree. Our *L. lory viridicrissalis* was a historical sample with a high degree of missing parsimony informative sites and its position as sister to *L. lory* and *L. hypoinochrous* is most likely an artifact. There were also varying levels of support for relationships among the other subspecies in *L. lory*, the most diverse species in the genus. Relationships among subspecies in *C. papou*, *C. josefinae*, and *C. placentis* had high support. Our analyses inferred a paraphyletic *T. haematodus* (with low support) and *T. forsteni*, the latter of which is still included in *T. haematodus* by some taxonomic checklists (Dickinson and Remsen 2013; Clements et al. 2019). This clade had many historical samples, which likely contributed to the clade's low support, but even several of the taxa from modern samples were not placed with high support in the clade. Resolving these challenging relationships within *Trichoglossus* will likely require finer-resolution genetic data and expanded population-level sampling.

Conclusion

Next-generation sequencing has provided systematists with an unprecedented amount of information for inferring phylogenetic relationships (McCormack et al. 2013). However, phylogenomic data sets are being produced faster than the

development of best practices for assembling, processing, and analyzing large data sets for phylogenetic inference, particularly as the use of low-quality museum samples increases. Alignments produced without careful inspection may harbor biased loci that can have a large impact on downstream analyses (Springer and Gatesy 2018). Our findings have general implications for phylogenomic studies where there is an asymmetry in parsimony informative sites among closely related taxa. Although missing data have shown ambiguous impacts on phylogenetic inference (Lemmon et al. 2009; Wiens and Morrill 2011; Simmons 2012, 2014; Hovmöller et al. 2013; Streicher et al. 2016), the combination of a much higher number of informative sites in contemporary phylogenomics with an asymmetry between samples of different quality warrants new investigations on biases that can arise in alignments. The magnitude of biases will likely vary according to clade diversity and age and the number of loci collected. We found that the bias was most extreme in a diverse and rapid radiation where there was likely limited information, even in complete loci, for teasing apart relationships. Shallow systematic and phylogeographic studies are expected to be the most difficult temporal scale for resolving relationships when there are high missing data associated with particular samples. Moving forward, having an understanding of the informational content of a locus, and how that information affects genealogy, will help avoid inferring dubious phylogenomic relationships.

Supplementary Material

Supplementary data are available at *Genome Biology and Evolution* online.

Acknowledgments

We thank the following institutions and people for providing the material used in this study: AMNH (P. Sweet, T. Trombone, G. Rosen, and A. Caragiulo), UWBM (S. Birks, R. Faucett, and J. Klicka), USNM (B. Schmidt, H. James, and G. Graves), ANWC (R. Palmer and L. Joseph), LSUMZ (D. Dittmann, S. Cardiff, R. Brumfield, and F. Sheldon), FMNH (B. Marks, J. Bates, and S. Hackett), and KU (M. Robbins and R. Moyle). We also thank F. Burbrink, K.L. Provost, and J. Merwin for providing code and J. Merwin, K.L. Provost, L.R. Moreira, G. Thom, B. Faircloth, C. Oliveros, M. Harvey, B. Holland, and anonymous reviewers for feedback and suggestions.

Literature Cited

Amadon D. 1943. Birds collected during the Whitney South Sea Expedition. LII, Notes on some non-passerine genera, 3. *Am Mus Novit.* 1237:1–22.

- Andersen MJ, Fatdal L, Mauck WM III, Smith BT. 2017. An ornithological survey of Vanuatu on the islands of Éfaté, Malakula, Gaua, and Vanua Lava. *Check List* 13(6):755–782.
- Andersen MJ, McCullough JM, Mauck WM III, Smith BT, Moyle RG. 2018. A phylogeny of kingfishers reveals an Indomalayan origin and elevated rates of diversification on oceanic islands. *J Biogeogr.* 45(2):269–281.
- Andersen MJ, et al. 2019. Ultraconserved elements resolve genus-level relationships in a major Australasian bird radiation (Aves: Meliphagidae). *Emu Austral Ornithol.* 119(3):218–232.
- Arcila D, et al. 2017. Genome-wide interrogation advances resolution of recalcitrant groups in the tree of life. *Nat Ecol Evol.* 1:0020.
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30(15):2114–2120.
- Borowiec ML. 2016. AMAS: a fast tool for alignment manipulation and computing of summary statistics. *PeerJ* 4:e1660.
- Briggs AW, et al. 2007. Patterns of damage in genomic DNA sequences from a Neandertal. *Proc Natl Acad Sci U S A.* 104(37):14616–14621.
- Brown JM, Thomson RC. 2017. Bayes factors unmask highly variable information content, bias, and extreme influence in phylogenomic analyses. *Syst Biol.* 66(4):517–530.
- Brown JW, Walker JF, Smith SA. 2017. Phyx: phylogenetic tools for unix. *Bioinformatics* 33(12):1886–1888.
- Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25(15):1972–1973.
- Chakrabarty P, et al. 2017. Phylogenomic systematics of ostariophysan fishes: ultraconserved elements support the surprising non-monophyly of characiformes. *Syst Biol.* 66(6):881–895.
- Chernomor O, von Haeseler A, Minh BQ. 2016. Terrace aware data structure for phylogenomic inference from supermatrices. *Syst Biol.* 65(6):997–1008.
- Clements JF, et al. 2019. The eBird/Clements Checklist of Birds of the World: v2019. Available from: <https://www.birds.cornell.edu/clements-checklist/download/>. Accessed January 1, 2020.
- Dickinson EC, Remsen JV Jr, editors. 2013. The Howard and Moore complete checklist of the birds of the world, Volume 1. 4th ed. Eastbourne, United Kingdom: Aves Press.
- Enk JM, et al. 2014. Ancient whole genome enrichment using baits built from modern DNA. *Mol Biol Evol.* 31(5):1292–1294.
- Esselstyn JA, Oliveros CH, Swanson MT, Faircloth BC. 2017. Investigating difficult nodes in the placental mammal tree with expanded taxon sampling and thousands of ultraconserved elements. *Genome Biol Evol.* 9(9):2308–2321.
- Ewart KM, et al. 2019. Museum specimens provide reliable SNP data for population genomic analysis of a widely distributed but threatened cockatoo species. *Mol Ecol Resour.* 19(6):1578–1592.
- Faircloth BC. 2013. illumiprocessor: a trimmomatic wrapper for parallel adapter and quality trimming. Available from: <https://illumiprocessor.readthedocs.io/en/latest/citing.html>. Accessed January 1, 2020.
- Faircloth BC. 2016. PHYLUCE is a software package for the analysis of conserved genomic loci. *Bioinformatics* 32(5):786–788.
- Faircloth BC, Branstetter MG, White ND, Brady SG. 2015. Target enrichment of ultraconserved elements from arthropods provides a genomic perspective on relationships among Hymenoptera. *Mol Ecol Resour.* 15(3):489–501.
- Faircloth BC, et al. 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. *Syst Biol.* 61(5):717–726.
- Forshaw JM. 2010. Parrots of the world. Wiltshire (United Kingdom): D & N Publishing.
- Forshaw JM, Cooper WT. 1989. Parrots of the world. London: Blandford.
- Fortes GG, et al. 2016. Ancient DNA reveals differences in behaviour and sociality between brown bears and extinct cave bears. *Mol Ecol.* 25(19):4907–4918.
- Gilbert PS, Wu J, Simon MW, Sinsheimer JS, Alfaro ME. 2018. Filtering nucleotide sites by phylogenetic signal to noise ratio increases confidence in the Neoaves phylogeny generated from ultraconserved elements. *Mol Phylogenet Evol.* 126:116–128.
- Gill F, Donsker D, editors. 2019. IOC world bird list v9.2. doi:10.14344/IOC.ML.9.2.
- Grabherr MG, et al. 2011. Trinity: reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat Biotechnol.* 29(7):644–652.
- Harvey MG, Smith BT, Glenn TC, Faircloth BC, Brumfield RT. 2016. Sequence capture versus restriction site associated DNA sequencing for shallow systematics. *Syst Biol.* 65(5):910–924.
- Helgen KM, et al. 2013. Taxonomic revision of the olingos (*Bassaricyon*), with description of a new species, the Olinguito. *ZooKeys* 324:1–83.
- Hosner PA, Faircloth BC, Glenn TC, Braun EL, Kimball RT. 2016. Avoiding missing data biases in phylogenomic inference: an empirical study in the landfowl (Aves: Galliformes). *Mol Biol Evol.* 33(4):1110–1125.
- Hovmöller R, Knowles LL, Kubatko LS. 2013. Effects of missing data on species tree estimation under the coalescent. *Mol Phylogenet Evol.* 69(3):1057–1062.
- Huang H, Knowles LL. 2016. Unforeseen consequences of excluding missing data from next-generation sequences: simulation study of RAD sequences. *Syst Biol.* 65(3):357–365.
- Hung CM, et al. 2014. Drastic population fluctuations explain the rapid extinction of the passenger pigeon. *Proc Natl Acad Sci U S A.* 111(29):10636–10641.
- Jiang W, Chen SY, Wang H, Li DZ, Wiens JJ. 2014. Should genes with missing data be excluded from phylogenetic analyses? *Mol Phylogenet Evol.* 80:308–318.
- Joseph L, Toon A, Schirtzinger EE, Wright TF, Schodde R. 2012. A revised nomenclature and classification for family-group taxa of parrots (Psittaciformes). *Zootaxa* 3205(1):26–40.
- Kalyaanamoorthy S, Minh BQ, Wong TK, von Haeseler A, Jermini LS. 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods.* 14(6):587–589.
- Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 30(4):772–780.
- Kehlmaier C, et al. 2017. Tropical ancient DNA reveals relationships of the extinct Bahamian giant tortoise *Chelonoidis alburyorum*. *Proc R Soc B* 284(1846):20162235.
- Kratter AW, Kirchman JJ, Steadman DW. 2006. Upland bird communities on Santo, Vanuatu, Southwest Pacific. *Wilson J Ornithol.* 118(3):295–308.
- Kück P, Longo GC. 2014. FASconCAT-G: extensive functions for multiple sequence alignment preparations concerning phylogenetic studies. *Front Zool.* 11(1):81.
- Kuhn M. 2008. Caret package. *J Stat Softw.* 28:1–26.
- Leigh DM, Lischer HEL, Grossen C, Keller LF. 2018. Batch effects in a multiyear sequencing study: false biological trends due to changes in read lengths. *Mol Ecol Resour.* 18(4):778–788.
- Lemmon AR, Brown JM, Stanger-Hall K, Lemmon EM. 2009. The effect of ambiguous data on phylogenetic estimates obtained by maximum likelihood and Bayesian inference. *Syst Biol.* 58(1):130–145.
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* 25(14):1754–1760.
- Li H, et al. 2009. The Sequence Alignment/Map (SAM) format and SAMtools. *Bioinformatics* 25(16):2078–2079.
- Linck EB, Hanna ZR, Sellas A, Dumbacher JP. 2017. Evaluating hybridization capture with RAD probes as a tool for museum genomics with historical bird specimens. *Ecol Evol.* 7(13):4755–4767.
- Malmström H, Storå J, Dalén L, Holmlund G, Götherström A. 2005. Extensive human DNA contamination in extracts from ancient dog bones and teeth. *Mol Biol Evol.* 22(10):2040–2047.

- Mayr E. 1933. Birds collected during the Whitney South Sea Expedition. 24, Notes on Polynesian flycatchers and a revision of the genus *Clytorhynchus* Elliot. *Am Mus Novit.* 628:1–21.
- Mayr E. 1938. Birds collected during the Whitney South Seas Expedition, XL. *Am Mus Novit.* 522:1–22.
- Mayr E. 1942. Birds collected during the Whitney South Sea Expedition. 48, Notes on the Polynesian species of *Aplonis*. *Am Mus Novit.* 1166:1–8.
- McCormack JE, Hird SM, Zellmer AJ, Carstens BC, Brumfield RT. 2013. Applications of next-generation sequencing to phylogeography and phylogenetics. *Mol Phylogenet Evol.* 66(2):526–538.
- McCormack JE, Tsai WL, Faircloth BC. 2016. Sequence capture of ultra-conserved elements from bird museum specimens. *Mol Ecol Resour.* 16(5):1189–1203.
- McCullough JM, Moyle RG, Smith BT, Andersen MJ. 2019. A Laurasian origin for a pantropical bird radiation is supported by genomic and fossil data (Aves: Coraciiformes). *Proc R Soc B* 286(1910):20190122.
- Mclean BS, Bell KC, Allen JM, Helgen KM, Cook JA. 2019. Impacts of inference method and data set filtering on phylogenomic resolution in a rapid radiation of ground squirrels (Xerinae: Marmotini). *Syst Biol.* 68(2):298–316.
- Merwin JT, Seeholzer GF, Smith BT. 2020. Macroevolutionary bursts and constraints generate a rainbow in a clade of tropical birds. *BMC Evol Biol.* 20:32.
- Mitchell KJ, et al. 2014. Ancient DNA reveals elephant birds and kiwi are sister taxa and clarifies ratite bird evolution. *Science* 344(6186):898–900.
- Mitchell KJ, et al. 2016. Ancient DNA from the extinct South American giant glyptodont *Doedicurus* sp. (Xenarthra: Glyptodontidae) reveals that glyptodonts evolved from Eocene armadillos. *Mol Ecol.* 25(14):3499–3508.
- Mivart SG. 1896. A monograph of the lories, or brush-tongued parrots, composing the family Loridae. London: R.H. Porter.
- Molloy EK, Warnow T. 2018. To include or not to include: the impact of gene filtering on species tree estimation methods. *Syst Biol.* 67(2):285–303.
- Moyle RG, et al. 2016. Tectonic collision and uplift of Wallacea triggered the global songbird radiation. *Nat Commun.* 7:12709.
- Mundy NI, Unitt P, Woodruff DS. 1997. Skin from feet of museum specimens as a non-destructive Source of DNA for avian genotyping. *Auk* 114(1):126–129.
- Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 32(1):268–274.
- Paijmans JL, et al. 2017. Evolutionary history of saber-toothed cats based on ancient mitogenomics. *Curr Biol.* 27(21):3330–3336.
- Paradis E, Claude J, Strimmer K. 2004. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics* 20(2):289–290.
- Philippe H, et al. 2004. Phylogenomics of eukaryotes: impact of missing data on large alignments. *Mol Biol Evol.* 21(9):1740–1752.
- Provost KL, Joseph L, Smith BT. 2018. Resolving a phylogenetic hypothesis for parrots: implications from systematics to conservation. *Emu Austral Ornithol.* 118(1):7–21.
- R Core Team. 2019. R: a language and environment for statistical computing. Vienna (Austria): R Foundation for Statistical Computing. Available from: <http://www.R-project.org/>
- Revell LJ. 2012. phytools: an R package for phylogenetic comparative biology (and other things). *Methods Ecol Evol.* 3(2):217–223.
- Robinson DF, Foulds LR. 1981. Comparison of phylogenetic trees. *Math Biosci.* 53(1–2):131–147.
- Roure B, Baurain D, Philippe H. 2013. Impact of missing data on phylogenies inferred from empirical phylogenomic data sets. *Mol Biol Evol.* 30(1):197–214.
- Ruane S, Austin CC. 2017. Phylogenomics using formalin-fixed and 100+ year-old intractable natural history specimens. *Mol Ecol Resour.* 17(5):1003–1008.
- Sawyer S, Krause J, Guschanski K, Savolainen V, Pääbo S. 2012. Temporal patterns of nucleotide misincorporations and DNA fragmentation in ancient DNA. *PLoS One* 7(3):e34131.
- Schweizer M, Wright TF, Peñalba JV, Schirtzinger EE, Joseph L. 2015. Molecular phylogenetics suggests a New Guinean origin and frequent episodes of founder-event speciation in the nectarivorous lories and lorikeets (Aves: Psittaciformes). *Mol Phylogenet Evol.* 90:34–48.
- Shavit Grievink L, Penny D, Holland BR. 2013. Missing data and influential sites: choice of sites for phylogenetic analysis can be as important as taxon sampling and model choice. *Genome Biol Evol.* 5(4):681–687.
- Shen XX, Hittinger CT, Rokas A. 2017. Contentious relationships in phylogenomic studies can be driven by a handful of genes. *Nat Ecol Evol.* 1(5):126.
- Simmons MP. 2012. Radical instability and spurious branch support by likelihood when applied to matrices with non-random distributions of missing data. *Mol Phylogenet Evol.* 62(1):472–484.
- Simmons MP. 2014. A confounding effect of missing data on character conflict in maximum likelihood and Bayesian MCMC phylogenetic analyses. *Mol Phylogenet Evol.* 80:267–280.
- Smith BT, Harvey MG, Faircloth BC, Glenn TC, Brumfield RT. 2014. Target capture and massively parallel sequencing of ultraconserved elements for comparative studies at shallow evolutionary time scales. *Syst Biol.* 63(1):83–95.
- Sorenson MD, et al. 1999. Relationships of the extinct moa-nalos, flightless Hawaiian waterfowl, based on ancient DNA. *Proc R Soc Lond B.* 266(1434):2187–2193.
- Springer MS, Gates J. 2018. On the importance of homology in the age of phylogenomics. *Syst Biodivers.* 16(3):210–228.
- Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30(9):1312–1313.
- Streicher JW, Schulte JA, Wiens JJ. 2016. How should genes and taxa be sampled for phylogenomic analyses with missing data? An empirical study in iguanian lizards. *Syst Biol.* 65(1):128–145.
- Thom G, et al. 2018. Phenotypic and genetic structure support gene flow generating gene tree discordances in an Amazonian floodplain endemic species. *Syst Biol.* 67(4):700–718.
- Thomas RH, Schaffner W, Wilson AC, Pääbo S. 1989. DNA phylogeny of the extinct marsupial wolf. *Nature* 340(6233):465–467.
- Tin MMY, Rheindt FE, Cros E, Mikheyev AS. 2015. Degenerate adaptor sequences for detecting PCR duplicates in reduced representation sequencing data improve genotype calling accuracy. *Mol Ecol Resour.* 15(2):329–336.
- Walker JF, Brown JW, Smith SA. 2018. Analyzing contentious relationships and outlier genes in phylogenomics. *Syst Biol.* 67(5):916–924.
- Wiens JJ, Morrill MC. 2011. Missing data in phylogenetic analysis: reconciling results from simulations and empirical data. *Syst Biol.* 60(5):719–731.
- Wingett SW, Andrews S. 2018. FastQ Screen: a tool for multi-genome mapping and quality control. *F1000Research* 7:1338.
- Xi Z, Liu L, Davis CC. 2016. The impact of missing data on species tree estimation. *Mol Biol Evol.* 33(3):838–860.
- Yao L, Li H, Martin RD, Moreau CS, Malhi RS. 2017. Tracing the phylogeographic history of Southeast Asian long-tailed macaques through mitogenomes of museum specimens. *Mol Phylogenet Evol.* 116:227–238.

Associate editor: Barbara Holland