1   **Article Type:** Discoveries

2

3   **Whole-genome phylogenomics of the tinamous (Aves: Tinamidae): comparing gene tree**
4   **estimation error between BUSCOs and UCEs illuminates rapid divergence with**
5   **introgression**

6

7   **Authors:**

8   Lukas J. Musher[1], Therese A. Catanach[1], Thomas Valqui[2], Robb T. Brumfield[3], Alexandre
9   Aleixo[4], Kevin P. Johnson[5], Jason D. Weckstein[1,6]

10

11  **Author affiliations:**

12  [1] The Academy of Natural Sciences of Drexel University, Department of Ornithology
13  Philadelphia, PA, 19103, USA.
14  [2] Facultad de Ciencias Forestales, Universidad Nacional Agraria La Molina, Lima, Peru and
15  CORBIDI-Centro de Ornitología y Biodiversidad, Lima, Peru.
16  [3] Department of Biological Sciences and Museum of Natural Science, Louisiana State
17  University, Baton Rouge, LA, 70803, USA.
18  [4] Instituto Tecnológico Vale - ITV, Belém, Brazil.
19  [5] Illinois Natural History Survey, Prairie Research Institute, University of Illinois Urbana-
20  Champaign, Champaign, IL, USA.
21  [6] Department of Biodiversity, Earth, and Environmental Sciences, Drexel University,
22  Philadelphia, PA, 19103, USA.

23  **Abstract:**
24  Incomplete lineage sorting (ILS) and introgression increase genealogical discordance across
25  the genome, which complicates phylogenetic inference. In such cases, identifying orthologs that
26  result in gene trees with low estimation error is crucial because phylogenomic methods rely on
27  accurate gene histories. We sequenced whole genomes for the tinamous (Aves: Tinamidae) to
28  dissect the sources of gene and species-tree discordance and reconstruct their
29  interrelationships. We compared results based on four ortholog sets: (1) coding genes
30  (BUSCOs), (2) ultraconserved elements (UCEs) with short flanking regions, (3) UCEs with
31  intermediate flanks, and (4) UCEs with long flanks. We hypothesized that orthologs with more
32  phylogenetically informative sites would result in more accurate species trees because the
33  resulting gene trees contain lower error. Consistent with our hypothesis, we found that long
34  UCEs had the most informative sites and lowest rates of error. However, despite having many
35  informative sites, BUSCO gene trees contained high error compared to long UCEs. Unlike
36  UCEs, BUSCO gene sequences showed a positive association between the proportion of
37  parsimony informative sites and gene tree error. Thus, BUSCO and UCE datasets have
38  different underlying properties of molecular evolution, and these differences should be
39  considered when selecting loci for phylogenomic analysis. Still, species trees from different
40  datasets were mostly congruent. Only one clade, with a history of ILS and introgression,
41  exhibited substantial species-tree discordance across the different data sets. Overall, we
42  present the most complete phylogeny for tinamous to date, identify a new species, and provide
43  a case study for species-level phylogenomic analysis using whole-genomes.
44

45  **Introduction**
46
47        Although the growth of high-throughput sequencing approaches over the past decade
48  has greatly improved our understanding of evolutionary relationships, reconstructing the tree of
49  life remains an enduring challenge. Analyses that utilize alternative datasets, methodological
50  frameworks, or substitution models to answer the same phylogenetic questions often yield
51  conflicting results, which is surprising given that phylogenomic datasets often contain hundreds
52  of thousands or even millions of phylogenetically informative characters (Dunn et al. 2008;
53  Jarvis et al. 2014; Prum et al. 2015; Reddy et al. 2017; Simion et al. 2017; Franz et al. 2019;
54  Schultz et al. 2023). In many cases, alternative phylogenomic topologies resulting from different
55  methods or data are equally well-supported. Thus, two key questions for molecular biologists
56  are, (1) why do phylogenomic datasets yield discordant topologies? and (2) how should one
57  choose among conflicting well-supported topologies?
58        One way in which discordant phylogenies can result from the same or similar datasets is
59  from the use of different methods for reconstructing species trees. One frequently employed
60  method is to concatenate all alignments into a single supermatrix and treat the resulting
61  phylogeny as a genome-wide average (henceforth, concatenation). Another common approach
62  is to estimate a gene tree for each molecular marker independently and summarize their
63  histories to estimate the species tree based on expectations from a multi-species coalescent
64  (MSC) model (Zhang and Mirarab 2022). Unlike the concatenation approach, MSC methods are
65  designed to account for expected variation in gene histories due to incomplete lineage sorting
66  (ILS). ILS causes some gene trees to differ from the species tree due to retained ancestral

67    variation during speciation, resulting in alleles coalescing prior to speciation events in ways that
68    result in incongruence between gene and species trees (Maddison 1997).
69         Each of these two methods, concatenation and MSC, has theoretical advantages, but
70    can result in erroneous species trees under certain conditions. For example, there is an
71    expectation that concatenation may be the best method when ILS is weak but is likely to be
72    statistically inconsistent when ILS is strong (Mendes and Hahn 2018; Bryant and Hahn 2020).
73    This is because when ILS is strong, gene genealogies are expected to differ more from the
74    species tree (i.e., there is increased heterogeneity) compared to situations when ILS is weak,
75    thus reducing the probability that concatenation will find the true tree. ILS is more likely to occur
76    when successive divergence times are short and effective population sizes are large, such as
77    during rapid or adaptive radiations (Maddison 1997; Mclean et al. 2019; Lescroart et al. 2023;
78    Tan et al. 2023). Thus, discordance among concatenation and MSC methods may be in part
79    driven by the presence of rapid diversification that leads to ILS and high gene tree heterogeneity
80    that biases concatenation results.
81         When applying MSC methods, identifying orthologous markers that lead to accurate,
82    unbiased gene tree estimation is crucial to properly infer phylogeny (Meiklejohn et al. 2016;
83    Springer and Gatesy 2016; Tea et al. 2021). MSC methods that utilize gene trees as input (gene
84    tree summarization methods) may fail if gene tree heterogeneity is driven by estimation error
85    rather than coalescence events. In other words, the MSC model is intended to incorporate gene
86    tree heterogeneity due to ILS, but erroneous gene trees resulting from low quality, biased, or
87    inconsistent sequence data will increase species tree estimation error (Xi et al. 2015). Thus,
88    dissecting biological and artifactual sources of gene tree discordance to more accurately infer
89    species trees is a major challenge in the age of phylogenomics, where more data does not
90    automatically translate to improved inferences (Meiklejohn et al. 2016; Blom et al. 2017; Smith
91    et al. 2023). Because MSC methods assume that all gene tree discordance is driven by ILS,
92    eliminating erroneous gene trees becomes important (Roch and Warnow 2015; Springer and
93    Gatesy 2016). This point has led to a variety of methods for filtering genomic datasets for
94    phylogenetic efficacy (Doyle et al. 2015; Kuang et al. 2018; B.T. Smith et al. 2018; S.A. Smith et
95    al. 2018; Zhao et al. 2023) but has generally lacked consensus on the types of data that contain
96    low error to begin with.
97         Given the complexities associated with phylogenetic inference, identifying datasets of
98    gene orthologs containing minimal bias is a principal goal of systematics research and has
99    implications across diverse fields in molecular biology. One widely employed data type for
100   phylogenomics includes coding sequences obtained via anchored hybrid enrichment,
101   transcriptomics, or low coverage whole-genome sequencing (Lemmon et al. 2012; Prum et al.
102   2015; Allen et al. 2017; Waterhouse et al. 2017; Johnson 2019; Zhang et al. 2019; Burbrink et
103   al. 2020; Boyd et al. 2022; Van Damme et al. 2022). Another data type includes ultraconserved
104   elements (UCEs), which target conserved genomic regions to identify orthology, but also yield
105   sequences of their more variable flanking regions for phylogenetic reconstruction across a
106   range of timescales (Faircloth et al. 2012; McCormack et al. 2012; Faircloth et al. 2013; Musher
107   and Cracraft 2018; Gueuning et al. 2020; Catanach et al. 2021; Ostrow et al. 2023). Both data
108   types are widely applied in phylogenomics but can sometimes result in different phylogenetic
109   topologies. For instance, coding and non-coding markers have resulted in conflicting topologies
110   for the backbone of Neoaves (McCormack et al. 2013; Jarvis et al. 2014; Prum et al. 2015), and

111 some authors have suggested that this discordance can be attributed to data biases associated
112 with coding genes (e.g., elevated GC content and model misspecification) that also bias
113 downstream species tree inferences (Reddy et al. 2017). Nevertheless, there have been few
114 studies comparing the efficacy of these marker types to one another. Thus, the ability of coding
115 versus primarily non-coding markers to resolve phylogenomic trees remains an open question.
116         One method for evaluating ortholog quality and gene tree error is through quantification
117 of alignment information content, such as the number of parsimony informative sites
118 (henceforth, PIS), or by comparing levels of gene tree heterogeneity among data types (Fong
119 and Fujita 2011; Harris et al. 2014; Burbrink et al. 2020; Leite et al. 2021; Smith et al. 2023). For
120 example, alignments with few PIS often result in erroneous or poorly resolved gene trees
121 because of a lack of phylogenetic signal in the data (Xi et al. 2015; Meiklejohn et al. 2016). In
122 such cases, we expect datasets with fewer PIS to result in increased gene tree heterogeneity
123 that is driven by gene tree estimation error. For example, alignments of coding genes, or UCEs
124 that include long flanking regions, should contain many PIS, and therefore should result in lower
125 rates of gene tree error than UCEs with short flanking regions that contain relatively few variable
126 sites. Some studies have found that loci with more PIS result in more "clock-like" gene trees,
127 suggesting that more informative loci result in gene trees with branch lengths that represent real
128 biological signal rather than data-driven artifacts (Musher et al. 2019). Other studies have
129 shown a negative correlation between the number of PIS and Robinson-Foulds distance (a
130 measure of phylogenetic dissimilarity; henceforth, RF distance) between gene and species trees
131 (Burbrink et al. 2020). This relationship implies that more informative alignments result in
132 reduced gene tree estimation error. Alternatively, one might also expect that phylogenetic
133 markers that are too variable can result in erroneous gene trees, if phylogenetic signal becomes
134 lost due to sequence saturation (Felsenstein 1978; Brinkmann et al. 2005). Thus, documenting
135 the relationship between alignment information content and RF distance, as well as the variation
136 in overall RF distances among different datasets, should inform future phylogenomic study
137 design and illuminate data-type efficacy.
138         Whole-genome sequencing offers a way to test the robustness of a phylogenetic
139 topology given multiple types of phylogenetic markers, including both coding and non-coding
140 sequences, from the same underlying data. For example, whole genomes are often used to
141 harvest large numbers of single-copy orthologous protein-coding markers called BUSCO
142 (Benchmarking Universal Single Copy Orthologs) genes (Waterhouse et al. 2017; Alaei Kakhki
143 et al. 2023). Given their conservation of function, BUSCO genes can theoretically tackle both
144 relatively deep and shallow phylogenetic problems (Timilsena et al. 2022; Van Damme et al.
145 2022; Alaei Kakhki et al. 2023). Similarly, there are multiple pipelines that can easily be used to
146 obtain UCEs and other conserved elements for phylogenomics from whole-genome datasets
147 (Faircloth 2016; Edwards et al. 2017). Thus, sequencing whole genomes enables a robust test
148 of the efficacy of different data types for phylogenomics.
149         In this study, we apply whole-genome sequencing to reconstruct a species-level
150 phylogeny of the tinamous (Aves: Tinamidae). There has been limited work on tinamou
151 phylogenomics to date. Most past phylogenetic studies utilized either morphological data with
152 relatively few genetic markers (Bertelli et al. 2002; Bertelli and Porzecanski 2004; Valqui 2008;
153 Bertelli 2017; Almeida et al. 2022), or phylogenomic data with large-scale molecular sampling
154 but minimal taxonomic sampling (Cloutier et al. 2019). Thus, the species-level phylogenetic

155  relationships of the tinamous remain uncertain. Here, we build on past studies by sampling
156  thousands of orthologous markers across 45 out of 46 described species and a total of 65
157  named taxa (monotypic species plus subspecies). Specifically, we compare levels of PIS and
158  gene tree estimation error in two types of orthologous markers: complete single-copy protein-
159  coding genes (BUSCO's) and UCE's. In doing so, our objectives are to (1) examine the effect of
160  alignment information content on gene tree estimation error, (2) use this information to identify
161  high quality (low error) orthologs, (3) identify the drivers of phylogenetic discordance among
162  data types and methods, and (4) use these inferences to robustly reconstruct the evolutionary
163  history of the tinamous.
164
165  **Results**
166
167  *Assembly metrics, completeness, and ortholog metrics*
168
169  We successfully assembled 61 of 62 newly sequenced genomes and extracted BUSCO
170  and UCE targets from these assemblies plus 10 publicly available tinamou and 2 publicly
171  available outgroup whole genome assemblies (Table S1). Genome wide sequence coverage
172  varied within and among samples but was generally high; the mean of average coverage across
173  samples was 40.41x (standard deviation = 13.55x). Most genomes were also relatively
174  complete, containing an average of 89.69% (standard deviation = 10.88%) of 8,338 complete
175  single-copy BUSCO genes.
176  From these whole genomes, we extracted and aligned two types of orthologous
177  markers, BUSCO genes and UCEs. Specifically, we compiled four different ortholog sets, which
178  after filtering for complete occupancy (all 72 samples represented in each alignment) included
179  (1) 2,507 BUSCOs (6,368,028 bp in concatenated alignment), (2) 2,887 UCEs with 100 bp
180  flanking regions (969,279 bp in concatenated alignment), (3) 2,879 UCEs with 300 bp flanking
181  regions (2,227,921 bp in the concatenated alignment), and (4) 2,803 UCEs with 1000 bp flanks
182  (6,902,014 bp in the concatenated alignment). Henceforth, we refer to these datasets as
183  BUSCO, UCE100Flank, UCE300Flank, and UCE1000Flank datasets, respectively.
184
185  *Phylogenomics of Tinamous*
186
187  Each of the four ortholog datasets was analyzed twice: first we reconstructed the
188  phylogenetic history of tinamous using all loci concatenated into a single alignment (henceforth,
189  concatenated phylogeny) and then using a gene tree summarization approach implemented in
190  ASTRAL (Mirarab et al. 2014; Zhang et al. 2018) (henceforth, MSC phylogeny). The
191  reconstructed concatenated and MSC phylogenies for tinamous were well-supported and
192  broadly congruent across data types (Figures 1, 2, S1–S6). The phylogenies recovered from all
193  four datasets were entirely congruent except for relationships within a single clade (henceforth,
194  Clade A) of *Crypturellus* containing 13 named taxa across nine recognized species (Figure 3).
195  One genome, downloaded from NCBI (GCA 013389825) was labeled as *C. undulatus*, but did
196  not cluster with other members of that species, instead clustering with either *C. strigulosus*
197  (MSC trees) or *C. erythropus* (concatenated trees). We could not confirm the species identity of
198  this sample because the voucher was unavailable, and metadata indicated it was missing its

199  head. Moreover, another downloaded genome, GCA 013398335, was listed as *Nothoprocta*
200  *ornata* on NCBI but after including newly sequenced material, this *N. ornata* clustered with *N.*
201  *pentlandii* samples from Peru. After examining the voucher specimen (CORBIDI 168605), the
202  sample was indeed confirmed as *N. pentlandii oustaleti* and not *N. ornata* as indicated in NCBI.
203        Our results were broadly consistent with recognized taxonomic classifications. For
204  example, we recovered monophyletic subfamilies Nothurinae and Tinaminae and most genera
205  and species were also monophyletic. However, we recovered two genera as non-monophyletic:
206  the monotypic genus *Taoniscus* was embedded within *Nothura* and *Rhynchotus* was embedded
207  within *Nothoprocta*, rendering *Nothura* and *Nothoprocta* paraphyletic. We also recovered
208  *Nothoprocta pentlandii* as polyphyletic; one specimen from Bolivia (voucher LSUMZ 123403)
209  was sister to *N. perdicaria*, whereas the remaining samples from Peru formed a clade that was
210  sister to all remaining *Nothoprocta* (consistent with mitochondrial results in Valqui 2008). All
211  other species that were represented by multiple samples were recovered as monophyletic.
212
213  *Assessment of ortholog information content and gene tree estimation error*
214
215        We recovered significant differences in the number and proportion of PIS per alignment
216  as well as levels of gene tree estimation error among the four datasets. First, smilograms
217  revealed key differences in patterns of within-locus variability, wherein BUSCOs showed a
218  bimodal distribution of variable sites (likely reflecting increased variability at third codon
219  positions) but UCE datasets showed an expected pattern of gradually increasing variability
220  moving away from conserved UCE cores (Figure 4). As expected, the UCE100Flank dataset
221  had the fewest PIS, averaging just 45 PIS per alignment (standard deviation = 28.84), the
222  UCE300Flank dataset averaged an intermediate number of PIS with a mean of 207 per
223  alignment (standard deviation = 88.37), and the UCE1000Flank dataset contained the most PIS
224  with a mean of 1,036 per alignment (standard deviation = 250.57). The BUSCO alignments also
225  contained a high number of PIS, though with very high variance, indicating a range of
226  informative and uninformative loci (mean = 923, standard deviation = 690.08 PIS per locus)
227  (Table 1; Figure 5). Kruskall-Wallis tests confirmed that differences in both the number ($X^2$ =
228  578, $df$ = 3, $P$ < 0.00001) and proportion ($X^2$ = 6408.9, $df$ = 3, $P$ < 0.00001) of PIS significantly
229  differed among datasets. Wilcoxon rank sum tests also indicated that differences between
230  pairwise comparisons of these datasets were all significantly different ($P$ < 0.00001 for all).
231        To measure gene tree estimation error, we examined gene tree heterogeneity across
232  datasets using Robinson-Foulds (RF) distances between gene and species trees. Although
233  some heterogeneity is expected due to the stochastic nature of the coalescent process, we
234  assumed that increases in the mean and variance of RF distances between datasets were
235  attributable to increased gene tree estimation error. The UCE100Flank dataset had the highest
236  RF distances between gene and species trees, whereas the UCE1000Flank dataset had the
237  lowest mean and variance in RF distances (Figure 5; Table 1). The BUSCOs and UCE300Flank
238  datasets had intermediate RF distances, with similar means and variances (Table 1).
239  Interestingly, RF distances from the BUSCO dataset were highly variable, with estimated gene
240  trees from some loci showing relatively low (akin to the UCE1000Flank dataset) and some
241  showing exceptionally high (akin to the UCE100Flank dataset) RF distances relative to both
242  MSC and concatenated species trees. A Kruskall-Wallis test confirmed that RF-distances ($X^2$ =

243  2740.1, $df$ = 3, $P$ < 0.00001) differed significantly among datasets, and a Wilcoxon rank sum
244  test also revealed that differences between pairwise comparisons of these datasets were
245  significantly different ($P$ < 0.00001) except that RF-distances did not significantly differ between
246  the UCE300Flank and BUSCO datasets ($P$ = 0.25).
247      To quantify the relationship between alignment information content and gene tree
248  estimation error and test the hypothesis that increased information content led to reduced gene
249  tree estimation error, we modeled RF distances as a function of the number of parsimony
250  informative sites (PIS) per locus using generalized linear models. These models were overall
251  consistent with our hypothesized pattern of a negative association between both the number
252  and proportion of PIS and RF distance, with lower RF distances in gene trees built from more
253  informative alignments (Figure 6; Table S2). However, there was notable variation in the slope
254  and tightness of fit of these models among datasets. For example, the UCE300Flank dataset,
255  with intermediate numbers of PIS overall, had a tight negative association between RF distance
256  and both the number of PIS and percentage of PIS per locus (Table S2). The same negative
257  associations were revealed for UCE100Flank and UCE1000Flank datasets, but these
258  relationships were much noisier. These differences were better visualized when the three UCE
259  datasets were combined, revealing a strong and tightly fitting negative logarithmic relationship
260  between the number of PIS per alignment and RF distance. Although this relationship was
261  similar for the BUSCO dataset using the number of PIS as the predictor variable, we found the
262  opposite pattern when using the percentage of PIS. Thus, based on multiple assessments,
263  BUSCOs and UCEs behave differently with regards to gene tree to species tree discordance
264  and patterns of nucleotide site variation.
265
266  *Tests of incomplete lineage sorting and introgression*
267
268      Finally, because we found conflicting phylogenetic relationships between datasets for a
269  clade of species in the genus *Crypturellus* (Figure 3), we assessed the extent to which this
270  conflict could be explained by ILS due to rapid divergences or introgression between non-sister
271  taxa. To do so, we looked at the relative frequencies of alternative quartets for five short internal
272  branches within Clade A (Figure 7A). The MSC model assumes that all gene tree heterogeneity
273  arises via ILS (i.e., no introgression, no gene tree error, no recombination, etc.), and it makes
274  explicit predictions about the relative frequencies of alternative quartet topologies (unrooted four
275  taxon statements) for all internal branches in a phylogeny. Specifically, the MSC model predicts
276  one majority topology that is consistent with the species tree at a relative frequency >⅓ and two
277  minority topologies of equal relative frequencies <⅓ (Pamilo and Nei 1988). As ILS increases,
278  the relative frequencies of all three alternative quartets approach the ⅓ threshold. Deviations
279  from these expectations (i.e., minority topologies are not equivalent in frequency) suggests that
280  processes such as introgression may be influencing gene tree heterogeneity. All five branches
281  that we examined showed evidence of either ILS or deviations from the MSC model. For
282  branches 1 and 2 (two of the three shortest branches within Clade A) we recovered frequencies
283  for three alternative quartets that closely approached the ⅓ threshold, implying ILS has
284  promoted phylogenetic conflict at these branches. For branches 3–5, we found possible
285  deviations from the MSC model, though this could be an artifact of sample size.

7

286  To test whether the deviations from the MSC model observed in Clade A could be
287  attributed to introgression, we applied a phylogenomic network analysis, which revealed multiple
288  nodes involved in reticulation that may explain gene tree discordance in *Crypturellus* Clade A
289  (Figure 7B; electronic supplementary material). Specifically, we found two reticulate nodes. In
290  the first, *C. cinnamomeus* was recovered as reticulate between *C. transfasciatus* (inheritance
291  probability = 0.44) and *C. erythropus* (inheritance probability = 0.56). In the second, *C.
292  undulatus* was reticulate between the basal node of the clade (inheritance probability = 0.76)
293  and *C. erythropus* (inheritance probability = 0.24). Thus, some gene tree heterogeneity within
294  the low-error UCE1000Flank dataset may be attributable to historical introgression between
295  non-sister taxa.
296
297  **Discussion**
298
299  Here we used whole-genome sequencing to explore the biological and artifactual
300  sources of phylogenomic conflict and reconstruct the species-level history of a relatively old
301  Neotropical avian group of broad interest in evolutionary biology (Bertelli and Porzecanski 2004;
302  Prum et al. 2015; Altimiras et al. 2017; Sackton et al. 2019; Li et al. 2023), the tinamous.
303  Although we found that alignment information content and gene tree error varied considerably
304  within and among datasets, our results based on different analytical approaches (concatenated
305  and MSC methods) and datasets (UCEs and BUSCOs) were largely congruent. The topology of
306  only a single clade (Clade A) varied among methods and datasets (Figures 1, 2, and 3),
307  indicating that species tree estimation is, in general, robust to gene tree estimation error when
308  biological sources of gene tree heterogeneity are limited. Nevertheless, the concatenated tree
309  shows that Clade A contained multiple successive short internodal branch lengths (Figure 1),
310  which suggests ILS may be a key driver of gene tree heterogeneity in this clade, and thus an
311  important factor driving phylogenetic conflict. Indeed, quartet analysis revealed multiple internal
312  branches with relative quartet frequencies approaching the ⅓ threshold predicted by high ILS
313  (Figures 7A and S7). Cases such as this, where ILS is strong, necessitate the use of datasets
314  without significant gene tree estimation error to accurately infer the species tree because
315  concatenation is expected to fail when ILS is strong (Roch and Warnow 2015; Xi et al. 2015;
316  Springer and Gatesy 2016). To complicate matters, we found evidence of historical
317  introgression in this clade that may have further elevated levels of gene tree discordance
318  (Figure 7B). For this reason, we suggest the MSC tree based on the UCE1000Flank dataset,
319  which had the lowest rates of gene tree estimation error (lowest median and variance in RF
320  distances between gene and species trees) is likely the most reliable (Xi et al. 2015).
321  Nevertheless, in addition to ILS, rapid divergence events can be difficult to reconstruct because
322  the short timeframe also means there is limited phylogenetic signal for those divergence events
323  within individual gene trees (Leaché et al. 2015; Springer and Gatesy 2016; Mclean et al. 2019).
324  Thus, given the lack of agreement across datasets and methods, we would argue that more
325  analysis is needed to confirm the relationships of this clade that diversified rapidly with
326  introgression between non-sister taxa.
327
328  *The relationship between alignment information content and gene tree error*
329

330   We examined the effects of data type (coding versus primarily non-coding datasets) and
331   information content (number and proportion of PIS) on gene tree estimation error and found that
332   coding sequences and genes with relatively few PIS contained high rates of gene tree
333   estimation error. Although UCEs conformed to the expected negative association between PIS
334   and gene tree estimation error as measured using RF distance (Figure 6), BUSCOs did not.
335   Instead, we recovered a positive association between the proportion of PIS per alignment and
336   RF distance, possibly because coding genes with more variable sites are evolving faster and
337   thus are more prone to multiple substitutions (i.e., saturation effects), especially at third codon
338   positions which are less constrained. This is also consistent with the bimodal distribution of
339   variable sites in this dataset (Figure 4). Overall, these results show that data type and alignment
340   information content have the potential to compromise phylogenomic inference from gene tree
341   summarization methods that rely on the MSC model, even for datasets that contain relatively
342   large amounts of data (Roch and Warnow 2015; Springer and Gatesy 2016).
343   We measured rates of gene tree estimation error using RF distances between gene
344   trees and the inferred species tree, under the assumption that increases in RF distance were
345   indicative of higher rates of phylogenetic error. Although gene histories are not expected to be
346   congruent with the species tree in many cases (Tajima 1983; Pamilo and Nei 1988; Maddison
347   1997), the assumption that RF distances are good indicators of error rates is valid because
348   different datasets contained different means and variances in RF distances among gene trees.
349   This was evident when the three UCE datasets were combined into a single analysis, which
350   showed asymptotic convergence on relatively low mean and variance RF distances after
351   reaching about 500 PIS (Figure 6). The asymptotic shape of this relationship implies that as PIS
352   are added to an alignment, the resulting gene trees converge on a level of heterogeneity that is
353   representative of or approaching real biological signal (i.e., ILS and/or introgression) rather than
354   methodological or data-driven artifacts. Moreover, if all gene tree heterogeneity was due to
355   biological processes such as ILS, one would expect the mean and variance in RF distances to
356   be constant among datasets because these are derived from the same underlying samples.
357   Although some level of gene tree heterogeneity is expected in *all* phylogenomic datasets, even
358   if one assumes the impossible, where empirical gene tree estimation error is absent (Maddison
359   1997; Gatesy and Springer 2014; Edwards et al. 2016), it is reasonable to assume that
360   differences in the mean and variance in RF distances among datasets are good proxies for
361   differences in gene tree estimation error.
362
363   *Comparisons of estimation error between data types and the benefits of whole genome*
364   *phylogenomics*
365
366   We scrutinized orthologous data harvested from whole genome assemblies to identify
367   the characteristics of phylogenomic datasets that might lead to lower bias. We found that short
368   alignments with relatively few PIS had increased error (measured as difference between gene
369   tree and species tree), a finding consistent with other studies (Meiklejohn et al. 2016; Burbrink
370   et al. 2020). Likewise, we found that long UCEs tend to have low rates of gene tree estimation
371   error and are efficacious markers to use in phylogenomic studies. We also found that the gene
372   trees from the protein-coding dataset derived from BUSCO harvesting and splicing of exons
373   from the tinamou genomes were very noisy. Despite containing large numbers and proportions

374  of PIS (Figure 5) available across many genes, we found high variance in RF distances for
375  these genes and thus conclude that such datasets might be more prone to error than datasets
376  of long UCEs, which had lower gene tree estimation error and behaved more predictably
377  (Figures 4 and 5). As a comparison, the UCE300Flank dataset also had high variance in RF
378  distances (i.e., it contained a mix of low and high error trees), but this variance was well
379  explained by information content. Thus, coding and UCE datasets have different underlying
380  properties of molecular evolution, and these differences should be considered when selecting
381  loci for study and in data analysis. Despite these differences, BUSCO and UCE300Flank
382  datasets converged on the same phylogenetic topology for Clade A, perhaps related to their
383  similar levels of gene tree heterogeneity (Wilcoxon rank sum tests accepted the null hypothesis
384  that RF distances for both datasets did not differ).
385          Although we derived datasets from whole genome sequences, the UCE100Flank and
386  UCE300Flank datasets are likely to most closely match datasets obtained from typical target
387  capture approaches (Smith et al. 2014; Musher and Cracraft 2018; Tea et al. 2021) and are
388  therefore indicative of rates of error in datasets derived from those protocols. For example,
389  many studies utilize sequences from degraded DNA such as historical museum specimens,
390  which may result in many relatively short UCE alignments. Still, even datasets sequenced from
391  fresh tissues typically only capture about 100-300 bp of flanking region around the conserved
392  UCE core. Thus, we offer multiple recommendations for future phylogenomic studies. First, if
393  available, WGS data are highly preferred to sequence capture datasets because they allow for
394  analysis of a much wider array of data types that can be compared as done herein (Jarvis et al.
395  2014; Zhang et al. 2019). Such data also allow studies to achieve longer flanking regions
396  around conserved UCE cores that significantly reduce bias during gene tree estimation (Figures
397  5 and 6), and enable analyses of many other processes of interest, such as evolution of gene
398  families. Moreover, we did not include intron data in the current study, but these data would also
399  be available from WGS and may behave more like UCE data rather than like the nearby exons
400  themselves.
401          Second, if a very large genome size is cost prohibitive for WGS, then target capture
402  datasets such as UCEs will likely contain high variance in RF distances (i.e., a mix of "good"
403  and "bad" loci), but this variance is likely to be well explained by the number of PIS in each
404  alignment (Figure 6). We therefore suggest that these datasets should be scrutinized and
405  filtered, especially if a rapid radiation is involved (Doyle et al. 2015; S.A. Smith et al. 2018;
406  Mclean et al. 2019; Smith et al. 2023). As the number and proportion of PIS were strongly
407  predictive of gene tree estimation error, we suggest that simply removing alignments with too
408  few PIS may be useful in many cases, as has been suggested elsewhere (Harris et al. 2014;
409  Hosner et al. 2016; Blom et al. 2017; Leite et al. 2021). Given these findings, common
410  phylogenomic approaches that sequence loci with intermediate or low levels of information
411  content such as sequence capture of UCEs may be insufficient to resolve the relationships of
412  taxa that evolved rapidly, such as adaptive radiations, where rates of ILS and short internodal
413  distances caused by rapid diversification are expected to be extreme (Gatesy and Springer
414  2014; Mclean et al. 2019). In these situations, the number of variable sites may also be too few
415  to resolve gene trees involving short branches.
416          Finally, although our BUSCO dataset was noisy and included a range of high- and low-
417  error gene trees, we do not suggest that datasets of coding genes are inherently poor quality.

418 Indeed, protein-coding datasets such as BUSCOs, transcriptomes, and others are likely to
419 remain useful for a range of phylogenomic problems. For example, given their conservation of
420 function, datasets of protein-coding genes will be important for reconstructing phylogenetic
421 relationships at deeper timescales using either nucleotide or amino acid sequences (Dunn et al.
422 2008; Allen et al. 2017; Zhang et al. 2019; Boyd et al. 2022). Still, given that our BUSCO
423 dataset was noisy compared with our UCE datasets, we suggest exercising caution when
424 applying protein-coding genes to difficult phylogenomic problems (Reddy et al. 2017). Careful
425 fitting of substitution models, data filtering, and perhaps even manual inspection of alignments
426 may reduce bias in these and other datasets of protein-coding sequence (Anisimova and Kosiol
427 2009; Springer and Gatesy 2016).

428    However, our results do not suggest that BUSCOs, or protein-coding genes in general,
429 should be filtered by information content or RF distances. Information content was a weak and
430 inconsistent predictor of gene tree estimation error for BUSCOs. Alignments with a fewer
431 number of PIS or a higher proportion of PIS were both associated with increased gene tree
432 error. Likewise, some amount of gene tree heterogeneity is expected due to the coalescent
433 process. Thus, assuming that all high RF gene trees are erroneous and removing them from a
434 dataset could mislead MSC analyses. Thus, more research is needed to understand how best
435 to filter BUSCO datasets, and perhaps protein-coding datasets in general, for phylogenomics
436 using gene tree summarization methods.

437

438 *Phylogeny and taxonomy of the Tinamous*

439

440    Our phylogenies based on multiple datasets with varying amounts of gene tree
441 estimation error and information content all resulted in highly similar and well-supported
442 relationships among tinamou species (Figures 1–3, S1–S6). Moreover, concatenation and MSC
443 approaches largely agreed on the overarching relationships at both shallow and deeper
444 timescales. However, for one clade, *Crypturellus* Clade A, we recovered phylogenetic conflict
445 among datasets and species tree approaches. Topologies based on analysis of concatenated
446 data tended to be less supported and less consistent than those estimated from gene trees in
447 ASTRAL, and the UCE100Flank trees showed the lowest support overall. However, MSC trees
448 based on longer UCEs and BUSCOs mostly agreed on the topology of Clade A, only varying in
449 the placement of *C. atrocapillus* and a misidentified downloaded genome (GCA 013389825).
450 The concatenated UCE1000Flank species tree was also most similar to these MSC results,
451 though the placement of *C. kerriae* differed. Therefore, for this subclade, the concatenation
452 results were statistically inconsistent, but the MSC results converged on a nearly identical
453 topology across datasets (other than UCE100Flank, which contained high error) (Bryant and
454 Hahn 2020).

455    Importantly, our phylogenetic results imply that multiple taxonomic changes are
456 necessary to appropriately classify taxa within Tinamidae. First, *Nothoprocta cinerascens* was
457 recovered as sister to the genus *Rhynchotus*. *Rhynchotus* (Spix 1825) has nomenclatural
458 priority over *Nothoprocta* (Sclater and Salvin 1873)*,* and therefore we suggest moving *N.*
459 *cinerascens* to the genus *Rhynchotus (Bertelli and Porzecanski 2004; Valqui 2008).* Second,
460 *Taoniscus nanus* was recovered as embedded within the genus *Nothura*, and thus, *Taoniscus*
461 (Gloger 1842) is a junior synonym of *Nothura* (Wagler 1827)*.* We therefore suggest transferring

462  *T. nanus* to *Nothura.* Finally, we found evidence for what is likely a new species of *Nothoprocta*
463  (see also Valqui 2008); specimens ascribed to *N. pentlandii* were polyphyletic, with a specimen
464  of the nominate subspecies from Bolivia recovered as sister to *N. perdicaria*, and multiple
465  specimens from Peru recovered as a clade sister to all other *Nothoprocta.* Thus, there appear to
466  be at least two species-level taxa within *N. pentlandii*.

468  **Materials and methods**

470  *Sampling*

472  A total of 46 species of tinamous classified into nine genera are currently recognized
473  (Clements et al. 2023). Of these, 22 are monotypic and the remaining species contain two or
474  more subspecies. Altogether there are 175 named taxa (monotypic species plus subspecies)
475  recognized in Tinamidae. We sampled 50 fresh tissues, 12 toe pads from historical museum
476  study skins, and downloaded 10 whole-genome assemblies from the NCBI Assembly Archive,
477  spanning a total of 66 named taxa (subspecies plus monotypic species) and 45 recognized
478  tinamou species. Our sampling included all recognized species except one; a sample initially
479  identified as *Crypturellus boucardi* turned out to be *C. soui meserythrus* (Voucher LSUMZ Birds
480  180663, Tissue LSUMNS B-53413).
481  Tinamous belong to the infraclass, Palaeognathae, which includes many extant flightless
482  ratites such as ostriches (Struthionidae) and rheas (Rheidae) along with extinct forms such as
483  Moas (Dinorninithiformes). Recent work has indicated that moas are the sister group to
484  tinamous and that rheas belong to a clade that is sister to tinamous plus moas (Cloutier et al.
485  2019). Thus, as outgroup taxa for tree rooting, we also downloaded whole genome assemblies
486  for *Rhea pennata* (Greater Rhea) and *Anomalopteryx didiformis* (Little Bush Moa) from the
487  NCBI Assembly Archive. Table S1 lists the 74 samples (72 tinamous plus two outgroups)
488  included in this study.

490  *DNA extraction, library preparation, and whole genome sequencing*

492  For fresh tissues, we extracted genomic DNA (gDNA) using the MagAttract High
493  Molecular Weight kit from Qiagen (Valencia, California). Toe pads extraction of historical study
494  skins was carried out in a dedicated historical DNA extraction laboratory at the Academy of
495  Natural Sciences of Drexel University to reduce contamination risk. Toe pad samples were first
496  washed in a brief bath of EtOH to help remove superficial contaminants, and then soaked in
497  $H_2O$ for three hours to hydrate the desiccated flesh for DNA lysis. Samples were then digested
498  using 180 μL ATL and 30 μL proteinase K for each sample and incubated at 56º C overnight.
499  DNA isolation was then performed using the QiaQuick spin columns and extraction kit from
500  Qiagen (Valencia, CA).
501  Shotgun sequencing libraries were prepared for each extract using the Hyper library
502  construction kit (Kapa Biosystems). These libraries were sequenced using 150 bp paired-end
503  reads on an S4 lane of an Illumina NovaSeq 6000. These libraries were pooled and tagged with
504  unique dual-end adaptors, and pooling consisted of between 13 and 18 samples per lane aimed

505    to achieve at least 30X coverage of the nuclear genome. Adapters were trimmed and
506    demultiplexed using bcl2fastq v.2.20. We deposited raw reads in the NCBI SRA (Table S1).
507
508
509    *Whole genome assembly*
510
511       We cleaned and then mapped raw reads for each sample (Table S1) to a closely-related
512    scaffold-level assembly from NCBI. Decisions about which downloaded genome to use for read
513    mapping were made based on a previous tinamou phylogeny (Almeida et al. 2022). Specifically,
514    we used fastp (Chen et al. 2018) to clean fastq files, trim adapters and remove low quality
515    reads. Next, we used Burrows-Wheeler Aligner (BWA) version 0.7.17 (Li and Durbin 2009) to
516    map the cleaned reads to reference genomes. Once reads were mapped, we used samtools
517    version 1.6 to convert the resulting sam-files into sorted bam-files (Li et al. 2009) and Analysis
518    of Next Generation Sequencing Data (ANGSD) version 1.2.13 (Korneliussen et al. 2014) to
519    convert bam-files into fasta format. We used '-doFasta 2' to ensure that the consensus
520    nucleotide was written for each polymorphic site.
521       To assess genome quality, completeness, and redundancy for each assembly, we used
522    Benchmarking Universal Single Copy Orthologs (BUSCO) version 5.3.0 (Simão et al. 2015).
523    BUSCO searches genome assemblies and identifies genes present in single copy using a
524    database of known single-copy orthologs from a clade-specific database of genes. We used the
525    'aves_odb10' lineage dataset, which utilizes 8,338 genes in the chicken genome. We used the '-
526    - augustus' flag to obtain nucleotide sequences for genes in addition to amino acid sequences.
527    This setting uses augustus version 3.2 (Hoff et al. 2019) to annotate each assembly, and
528    outputs full nucleotide gene sequences for all complete single-copy orthologs in fasta format.
529    This step was necessary to obtain our coding gene dataset for phylogenomic analysis. We also
530    used samtools to estimate mean and standard deviation of sequence coverage for each
531    genome.
532
533    *Ortholog identification and alignment*
534
535       We identified and extracted two types of orthologous markers from the WGS scaffolds.
536    First, we utilized nucleotide sequences for the 8,338 single copy orthologs obtained from the
537    BUSCO 'aves_odb10' dataset (henceforth, BUSCO dataset). We used custom scripts to append
538    orthologous genes into the same text file and then used MACSE (Multiple Alignment of Coding
539    SEquences) version 2.06 to align nucleotide sequences (Ranwez et al. 2018). MACSE is a
540    codon-aware alignment algorithm that aligns nucleotide sequences with respect to their protein
541    translation, but allows nucleotide sequences to contain exceptions to this underlying codon
542    structure. This allowed us to accurately align each of the 8,338 gene sequences for all 74
543    samples.
544       Second, we harvested ultraconserved elements (UCEs) from each genome assembly.
545    UCEs are highly conserved, typically single-copy, sequences in the genome that are flanked by
546    neutral sites that increase in variability moving away from the conserved core (McCormack et al.
547    2012). These may overlap with coding or non-coding sequences, but the general pattern of
548    conserved UCE core with variable flanking region should result in very different distribution of

13

549  variation within a locus than the BUSCO genes, which are purely coding and are most likely to
550  vary at third codon positions due to the redundant nature of the genetic code. To identify and
551  extract UCEs from the genome we followed the pipeline outlined in Phyluce version 1.7.1
552  (Faircloth 2016) for harvesting and loci from whole genomes. We specifically harvested UCE
553  loci using the Tetrapods-UCE-5kv1 dataset, which includes 5,060 UCEs.
554       In this pipeline, Phyluce scripts are used to align probe sequences (corresponding to
555  conserved UCE cores) to the genome and then extract those sequences from the genome with
556  a user-specified length of flanking region. Because we were interested in understanding
557  whether and how information content covaried with gene tree estimation error, we harvested
558  UCEs three times, varying the length of the flanking region included with the UCE core each
559  time. We harvested UCEs from each genome that included 100 bp (similar to sequence lengths
560  obtained with target capture from degraded DNA datasets such as historical museum skins),
561  300 bp (similar to sequence lengths obtained with target capture from standard fresh tissue;
562  e.g., Musher and Cracraft, 2018, Tea et al. 2021), and 1,000 bp (longer than typical of target
563  capture datasets) of flanking region. We then used MAFFT (Katoh and Standley 2013) to align
564  orthologous UCE loci for downstream phylogenomic analysis.
565       Because missing data can be an additional source of phylogenetic noise, we evaluated
566  each of the four datasets (three UCE and one BUSCO) using only complete (alignments for
567  which all samples are included) alignments. This helped to eliminate gene tree discordance and
568  other sources of error deriving from variation in missing data content.
569
570  *Phylogenomic analyses*
571
572       For each of the four datasets, we employed both a concatenation approach using
573  RAxML version 8.2.12 (Stamatakis 2014) and an MSC approach using ASTRAL version 5.7.7
574  (Mirarab et al. 2014; Zhang et al. 2018). First, we concatenated all alignments using the
575  'phyluce_align_concatenate_alignments.py' script available in the Phyluce software (Faircloth
576  2016) for each of the four ortholog datasets. We then used RAxML to reconstruct the phylogeny
577  under a GTR + CAT substitution model which approximates the GTR + GAMMA model and is
578  expected to perform well for large datasets (Stamatakis 2006; Abadi et al. 2019) We then
579  examined the robustness of our phylogenies using the autoMRE option to perform
580  bootstrapping but halt replicates when they converged. Next, for our MSC approach, we
581  estimated gene trees (phylogenies derived from presumed independently sorting loci) for each
582  alignment within each of the four datasets using RAxML. Then, we implemented a quartet-
583  based MSC approach in ASTRAL. ASTRAL approximates the MSC model by estimating the
584  proportions of all possible four-taxon statements (quartets) among all estimated gene trees and
585  then summarizing across the genome. This approach assumes that all gene tree heterogeneity
586  (discordance among estimated gene histories) is due to ILS without significant ancestral
587  introgression or gene tree estimation error. We also pruned our resulting trees to a single clade
588  of interest and compared topologies for this clade using Robinson-Foulds (RF) distances using
589  the 'RF.dist' function in phangorn version 2.11.1 (Schliep 2011).
590
591  *Assessing ortholog quality and gene tree error*
592

593    To evaluate general differences among dataset variability, we first estimated the average
594    variability among sites for each alignment within each dataset. To do so, we quantified the
595    average variability for each nucleotide site relative to the center of the alignment using the
596    'phyluce_align_get_smilogram_from_alignments' script available in the Phyluce package
597    (Faircloth 2016). Then, to assess the relationships between alignment information content, data-
598    type, and gene tree estimation error, we wrote custom scripts in R version 4.3.1 (R Core Team
599    2023) to estimate the number and proportion (informative sites/alignment length) of parsimony
600    informative sites at each locus, as well as the RF distances (Robinson and Foulds 1981)
601    between each estimated gene tree and the inferred species tree. Parsimony informative sites
602    were defined as variable sites in the alignment where each variant nucleotide is represented by
603    at least two samples. To estimate RF distances, we compared the gene trees from each dataset
604    inferred using RAxML to both the MSC and concatenated trees based on the UCE1000Flank
605    dataset. To test for differences in these measures among datasets, we used Kruskall-Wallace
606    and pairwise Wilcoxon rank sum tests with a Benjamini-Hochberg p-correction implemented in
607    base R (R Core Team 2023). We then ran generalized linear models using information content
608    as the independent and RF distance as the dependent variable for each of the four datasets.
609    Because we found no differences in RF distance when using MSC and concatenated species
610    trees, for these models, we only used RF distances relative to the MSC species tree resulting
611    from our ASTRAL analysis. We compared AIC values for linear and logarithmic fits for each
612    regression and chose the model with the lowest AIC for each dataset. These regression
613    analyses were done once for each dataset, and once with all three UCE datasets combined.
614    Although the latter analysis contains the same UCE locus multiple times, albeit with additional
615    flanking regions, we believe this analysis helps illuminate how the information content of loci
616    from across the spectrum (i.e., very uninformative to very informative) influences rates of gene
617    tree error.
618
619    *Tests of incomplete lineage sorting and introgression*
620
621    To examine the effects of ILS and introgression on the phylogenetic discordance of
622    Clade A, we first looked at the alternative quartet topologies for five short internal branches and
623    then used a phylogenomic network approach to test for reticulation (phylogenomic non-
624    bifurcation). For the quartet method, we used the custom scripts for the R package
625    'MSCquartets' to identify the relative quartet frequencies for five nodes of interest (Rhodes et al.
626    2021; Allman et al. 2023). Quartets are unrooted four taxon statements for which only three
627    alternative sets of relationships (topologies) exist. Because ILS is expected at relatively short
628    internal branches, we qualitatively evaluated the relative frequencies of alternative topologies at
629    these branches to assess whether ILS might be involved in increasing genealogical
630    heterogeneity, and thus phylogenetic disagreement of Clade A.
631    To test for putative introgression, we applied a Bayesian approach implemented in
632    PhyloNet 3.8.0 (the 'MCMC_GT' algorithm) (Wen et al. 2016). This method employs reverse-
633    jump Markov Chain Monte Carlo (rjMCMC) to sample the posterior distribution of phylogenetic
634    networks under a multi-species network coalescent (MSNC) model using only rooted gene trees
635    as input data. The MSNC is similar to the MSC, but relaxes the assumption of no introgression
636    by modeling genome evolution as a network rather than bifurcating phylogeny. In doing so, it

15

637   accounts for both incomplete lineage sorting and reticulation (multiple ancestral lineages that
638   contribute alleles to a single daughter lineage). Such reticulation is typically attributed to
639   historical introgression. We used the gene trees from the UCE1000Flank dataset as these
640   contained reduced rates of gene tree estimation error. We ran the analysis three times, varying
641   the maximum number of reticulate nodes in each run to assess variability in the results.
642   Specifically, we ran MCMC_GT allowing for a maximum of between zero (m=0) and ten (m=10)
643   reticulate nodes. We assigned all samples to the species level, and ran the analysis using
644   species as tips in the network rather than samples (i.e., some species were represented by
645   multiple taxa). We ran the rjMCMC for 5 x$10^7$ generations with a burn-in of 5x$10^6$ generations,
646   using the pseudo-likelihood calculation to reduce computation time. Because likelihood scores
647   tended to increase with each successive increase in m-value, and the rjMCMC typically found
648   the maximum allowed in each run, we chose the optimal m-value using a breakpoint analysis.
649   Specifically, using the R package 'segmented' (Muggeo et al. 2014), we choose the network
650   with the optimal m-value by fitting a segmented linear model to our likelihoods for each m-
651   value(Muggeo et al. 2014). This allowed us to identify breakpoints in the slope of the regression,
652   where increases in m-value resulted in diminishing gains in likelihood. The m-value at the
653   breakpoint in the segmented model was chosen as the optimal m-value.
654
655   **Acknowledgements:**

667
668   **Data Availability:**
669         All raw reads are available on the NCBI Sequence Read Archive (project and sample
670   accession numbers can be found in Table S1). All data and code needed to replicate the
671   analyses found herein can be found at LJM's github page:
672   https://github.com/lukemusher/Tinamou-phylogenomics.git
673
674   **Literature Cited**

675   Abadi S, Azouri D, Pupko T, Mayrose I. 2019. Model selection may not be a mandatory step for
676         phylogeny reconstruction. *Nat. Commun.* 10:934.

677   Alaei Kakhki N, Schweizer M, Lutgen D, Bowie RCK, Shirihai H, Suh A, Schielzeth H, Burri R.
678         2023. A Phylogenomic Assessment of Processes Underpinning Convergent Evolution in
679         Open-Habitat Chats. *Mol. Biol. Evol.* [Internet] 40. Available from:

680     http://dx.doi.org/10.1093/molbev/msac278

681     Allen JM, Boyd B, Nguyen N-P, Vachaspati P, Warnow T, Huang DI, Grady PGS, Bell KC,
682         Cronk QCB, Mugisha L, et al. 2017. Phylogenomics from Whole Genome Sequences Using
683         aTRAM. *Syst. Biol.* 66:786–798.

684     Allman ES, Baños H, Mitchell JD, Rhodes JA. 2023. MSCquartets: analyzing gene tree quartets
685         under the multi-species coalescent. R package version 1.3.1. Available from:
686         https://CRAN.R-project.org/package=MSCquartets

687     Almeida FC, Porzecanski AL, Cracraft JL, Bertelli S. 2022. The evolution of tinamous
688         (Palaeognathae: Tinamidae) in light of molecular and combined analyses. *Zool. J. Linn.*
689         *Soc.* 195:106–124.

690     Altimiras J, Lindgren I, Giraldo-Deck LM, Matthei A, Garitano-Zavala Á. 2017. Aerobic
691         performance in tinamous is limited by their small heart. A novel hypothesis in the evolution
692         of avian flight. *Sci. Rep.* 7:15964.

693     Anisimova M, Kosiol C. 2009. Investigating protein-coding sequence evolution with probabilistic
694         codon substitution models. *Mol. Biol. Evol.* 26:255–271.

695     Bertelli S. 2017. Advances on tinamou phylogeny: an assembled cladistic study of the volant
696         palaeognathous birds. *Cladistics* 33:351–374.

697     Bertelli S, Giannini NP, Goloboff PA. 2002. A phylogeny of the tinamous (aves:
698         palaeognathiformes) based on integumentary characters. *Syst. Biol.* 51:959–979.

699     Bertelli S, Porzecanski AL. 2004. Tinamou (tinamidae) systematics: A preliminary combined
700         analysis of morphology and molecules. *Ornitol. Neotrop.* 15 (Suppl):293–299.

701     Blom MPK, Bragg JG, Potter S, Moritz C. 2017. Accounting for Uncertainty in Gene Tree
702         Estimation: Summary-Coalescent Species Tree Inference in a Challenging Radiation of
703         Australian Lizards. *Syst. Biol.* 66:352–366.

704     Boyd BM, Nguyen N-P, Allen JM, Waterhouse RM, Vo KB, Sweet AD, Clayton DH, Bush SE,
705         Shapiro MD, Johnson KP. 2022. Long-distance dispersal of pigeons and doves generated
706         new ecological opportunities for host-switching and adaptive radiation by their parasites.
707         *Proc. Biol. Sci.* 289:20220042.

708     Brinkmann H, van der Giezen M, Zhou Y, Poncelin de Raucourt G, Philippe H. 2005. An
709         empirical assessment of long-branch attraction artefacts in deep eukaryotic phylogenomics.
710         *Syst. Biol.* 54:743–757.

711     Bryant D, Hahn MW. 2020. The Concatenation Question. Pages 3.4: 1–3.4: 23 in Phylogenetics
712         in the Genomic Era (C. Scornavacca, F. Delsuc, and N. Galtier, eds.).

713     Burbrink FT, Grazziotin FG, Pyron RA, Cundall D, Donnellan S, Irish F, Keogh JS, Kraus F,
714         Murphy RW, Noonan B, et al. 2020. Interrogating Genomic-Scale Data for Squamata
715         (Lizards, Snakes, and Amphisbaenians) Shows no Support for Key Traditional
716         Morphological Relationships. *Syst. Biol.* 69:502–520.

717     Catanach TA, Halley MR, Allen JM, Johnson JA, Thorstrom R, Palhano S, Thunder CP,
718         Gallardo JC, Weckstein JD. 2021. Systematics and conservation of an endemic radiation of

Accipiter hawks in the Caribbean islands. *The Auk* 138:1–23.

Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* 34:i884–i890.

Clements JF, Rasmussen PC, Schulenberg TS, Iliff MJ, Fredericks TA, Gerbracht JA, Lepage D, Spencer A, Billerman SM, Sullivan BL, et al. 2023. The eBird/Clements checklist of Birds of the World: v2023. Available from: Downloaded from https://www.birds.cornell.edu/clementschecklist/download/

Cloutier A, Sackton TB, Grayson P, Clamp M, Baker AJ, Edwards SV. 2019. Whole-genome analyses resolve the phylogeny of flightless birds (Palaeognathae) in the presence of an empirical anomaly zone. *Syst. Biol.* 68:937–955.

Doyle VP, Young RE, Naylor GJP, Brown JM. 2015. Can We Identify Genes with Increased Phylogenetic Reliability? *Syst. Biol.* 64:824–837.

Dunn CW, Hejnol A, Matus DQ, Pang K, Browne WE, Smith SA, Seaver E, Rouse GW, Obst M, Edgecombe GD, et al. 2008. Broad phylogenomic sampling improves resolution of the animal tree of life. *Nature* 452:745–749.

Edwards SV, Cloutier A, Baker AJ. 2017. Conserved Nonexonic Elements: A Novel Class of Marker for Phylogenomics. *Syst. Biol.* 66:1028–1044.

Edwards SV, Xi Z, Janke A, Faircloth BC, McCormack JE, Glenn TC, Zhong B, Wu S, Lemmon EM, Lemmon AR, et al. 2016. Implementing and testing the multispecies coalescent model: A valuable paradigm for phylogenomics. *Mol. Phylogenet. Evol.* 94:447–462.

Faircloth BC. 2016. PHYLUCE is a software package for the analysis of conserved genomic loci. *Bioinformatics* 32:786–788.

Faircloth BC, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. *Syst. Biol.* 61:717–726.

Faircloth BC, Sorenson L, Santini F, Alfaro ME. 2013. A Phylogenomic Perspective on the Radiation of Ray-Finned Fishes Based upon Targeted Sequencing of Ultraconserved Elements (UCEs). *PLoS One* 8:e65923.

Felsenstein J. 1978. Cases in which Parsimony or Compatibility Methods will be Positively Misleading. *Syst. Biol.* 27:401–410.

Fong JJ, Fujita MK. 2011. Evaluating phylogenetic informativeness and data-type usage for new protein-coding genes across Vertebrata. *Mol. Phylogenet. Evol.* 61:300–307.

Franz NM, Musher LJ, Brown JW, Yu S, Ludäscher B. 2019. Verbalizing phylogenomic conflict: Representation of node congruence across competing reconstructions of the neoavian explosion. *PLoS Comput. Biol.* 15:e1006493.

Gatesy J, Springer MS. 2014. Phylogenetic analysis at deep timescales: unreliable gene trees, bypassed hidden support, and the coalescence/concatalescence conundrum. *Mol. Phylogenet. Evol.* 80:231–266.

Gueuning M, Frey JE, Praz C. 2020. Ultraconserved yet informative for species delimitation: Ultraconserved elements resolve long-standing systematic enigma in Central European bees. *Mol. Ecol.* 29:4203–4220.

Harris RB, Carling MD, Lovette IJ. 2014. The influence of sampling design on species tree inference: a new relationship for the New World chickadees (Aves: Poecile). *Evolution* 68:501–513.

Hoff KJ, Lomsadze A, Borodovsky M, Stanke M. 2019. Whole-Genome Annotation with BRAKER. *Methods Mol. Biol.* 1962:65–95.

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.* 33:1110–1125.

Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SYW, Faircloth BC, Nabholz B, Howard JT, et al. 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds. *Science* 346:1320–1331.

Johnson KP. 2019. Putting the genome in insect phylogenomics. *Curr Opin Insect Sci* 36:111–117.

Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30:772–780.

Korneliussen TS, Albrechtsen A, Nielsen R. 2014. ANGSD: Analysis of Next Generation Sequencing Data. *BMC Bioinformatics* 15:356.

Kuang T, Tornabene L, Li J, Jiang J, Chakrabarty P, Sparks JS, Naylor GJP, Li C. 2018. Phylogenomic analysis on the exceptionally diverse fish clade Gobioidei (Actinopterygii: Gobiiformes) and data-filtering based on molecular clocklikeness. *Mol. Phylogenet. Evol.* 128:192–202.

Leaché AD, Chavez AS, Jones LN, Grummer JA, Gottscho AD, Linkem CW. 2015. Phylogenomics of phrynosomatid lizards: conflicting signals from sequence capture versus restriction site associated DNA sequencing. *Genome Biol. Evol.* 7:706–719.

Leite RN, Kimball RT, Braun EL, Derryberry EP, Hosner PA, Derryberry GE, Anciães M, McKay JS, Aleixo A, Ribas CC, et al. 2021. Phylogenomics of manakins (Aves: Pipridae) using alternative locus filtering strategies based on informativeness. *Mol. Phylogenet. Evol.* 155:107013.

Lemmon AR, Emme SA, Lemmon EM. 2012. Anchored hybrid enrichment for massively high-throughput phylogenomics. *Syst. Biol.* 61:727–744.

Lescroart J, Bonilla-Sánchez A, Napolitano C, Buitrago-Torres DL, Ramírez-Chaves HE, Pulido-Santacruz P, Murphy WJ, Svardal H, Eizirik E. 2023. Extensive Phylogenomic Discordance and the Complex Evolutionary History of the Neotropical Cat Genus Leopardus. *Mol. Biol. Evol.* [Internet] 40. Available from: http://dx.doi.org/10.1093/molbev/msad255

Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25:1754–1760.

796   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
797       1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map
798       format and SAMtools. *Bioinformatics* 25:2078–2079.

799   Li Q, Chen D, Wang S. 2023. Character displacement of egg colors during tinamou speciation.
800       *Evolution* 77:1874–1881.

801   Maddison WP. 1997. Gene Trees in Species Trees. *Syst. Biol.* 46:523–536.

802   McCormack JE, Faircloth BC, Crawford NG, Gowaty PA, Brumfield RT, Glenn TC. 2012.
803       Ultraconserved elements are novel phylogenomic markers that resolve placental mammal
804       phylogeny when combined with species-tree analysis. *Genome Res.* 22:746–754.

805   McCormack JE, Harvey MG, Faircloth BC, Crawford NG, Glenn TC, Brumfield RT. 2013. A
806       phylogeny of birds based on over 1,500 loci collected by target enrichment and high-
807       throughput sequencing. *PLoS One* 8:e54848.

808   Mclean BS, Bell KC, Allen JM, Helgen KM, Cook JA. 2019. Impacts of Inference Method and
809       Data set Filtering on Phylogenomic Resolution in a Rapid Radiation of Ground Squirrels
810       (Xerinae: Marmotini). *Syst. Biol.* 68:298–316.

811   Meiklejohn KA, Faircloth BC, Glenn TC, Kimball RT, Braun EL. 2016. Analysis of a Rapid
812       Evolutionary Radiation Using Ultraconserved Elements: Evidence for a Bias in Some
813       Multispecies Coalescent Methods. *Syst. Biol.* 65:612–627.

814   Mendes FK, Hahn MW. 2018. Why Concatenation Fails Near the Anomaly Zone. *Syst. Biol.*
815       67:158–169.

816   Mirarab S, Reaz R, Bayzid MS, Zimmermann T, Swenson MS, Warnow T. 2014. ASTRAL:
817       genome-scale coalescent-based species tree estimation. *Bioinformatics* 30:i541–i548.

818   Muggeo VMR, Atkins DC, Gallop RJ, Dimidjian S. 2014. Segmented mixed models with random
819       changepoints: a maximum likelihood approach with application to treatment for depression
820       study. *Stat. Modelling* 14:293–313.

821   Musher LJ, Cracraft J. 2018. Phylogenomics and species delimitation of a complex radiation of
822       Neotropical suboscine birds (Pachyramphus). *Mol. Phylogenet. Evol.* 118:204–221.

823   Musher LJ, Ferreira M, Auerbach AL, Cracraft J. 2019. Why is Amazonia a "source"of
824       biodiversity? Climate-mediated dispersal and synchronous speciation across the Andes in
825       an avian group (Tityrinae). *Proc. Roy. Soc. B* [Internet]. Available from:
826       https://royalsocietypublishing.org/doi/abs/10.1098/rspb.2018.2343

827   Ostrow EN, Catanach TA, Bates JM, Aleixo A, Weckstein JD. 2023. Phylogenomic analysis
828       confirms the relationships among toucans, toucan-barbets, and New World barbets but
829       reveals paraphyly of Selenidera toucanets and evidence for mitonuclear discordance.
830       *Ornithology* 140:ukad022.

831   Pamilo P, Nei M. 1988. Relationships between gene trees and species trees. *Mol. Biol. Evol.*
832       5:568–583.

833   Prum RO, Berv JS, Dornburg A, Field DJ, Townsend JP, Lemmon EM, Lemmon AR. 2015. A
834       comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing.

835    *Nature* 526:569–573.

836    Ranwez V, Douzery EJP, Cambon C, Chantret N, Delsuc F. 2018. MACSE v2: Toolkit for the
837        Alignment of Coding Sequences Accounting for Frameshifts and Stop Codons. *Mol. Biol.*
838        *Evol.* 35:2582–2584.

839    R Core Team. 2023. R: A language and environment for statistical computing. Foundation for
840        Statistical Computing, Vienna, Austria. Available from: https://www.R-project.org/

841    Reddy S, Kimball RT, Pandey A, Hosner PA, Braun MJ, Hackett SJ, Han K-L, Harshman J,
842        Huddleston CJ, Kingston S, et al. 2017. Why Do Phylogenomic Data Sets Yield Conflicting
843        Trees? Data Type Influences the Avian Tree of Life more than Taxon Sampling. *Syst. Biol.*
844        66:857–879.

845    Rhodes JA, Baños H, Mitchell JD, Allman ES. 2021. MSCquartets 1.0: quartet methods for
846        species trees and networks under the multispecies coalescent model in R. *Bioinformatics*
847        37:1766–1768.

848    Robinson DF, Foulds LR. 1981. Comparison of phylogenetic trees. *Math. Biosci.* 53:131–147.

849    Roch S, Warnow T. 2015. On the Robustness to Gene Tree Estimation Error (or lack thereof) of
850        Coalescent-Based Species Tree Methods. *Syst. Biol.* 64:663–676.

851    Sackton TB, Grayson P, Cloutier A, Hu Z, Liu JS, Wheeler NE, Gardner PP, Clarke JA, Baker
852        AJ, Clamp M, et al. 2019. Convergent regulatory evolution and loss of flight in
853        paleognathous birds. *Science* 364:74–78.

854    Schliep KP. 2011. phangorn: phylogenetic analysis in R. *Bioinformatics* 27:592–593.

855    Schultz DT, Haddock SHD, Bredeson JV, Green RE, Simakov O, Rokhsar DS. 2023. Ancient
856        gene linkages support ctenophores as sister to other animals. *Nature* 618:110–117.

857    Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO:
858        assessing genome assembly and annotation completeness with single-copy orthologs.
859        *Bioinformatics* 31:3210–3212.

860    Simion P, Philippe H, Baurain D, Jager M, Richter DJ, Di Franco A, Roure B, Satoh N,
861        Quéinnec É, Ereskovsky A, et al. 2017. A Large and Consistent Phylogenomic Dataset
862        Supports Sponges as the Sister Group to All Other Animals. *Curr. Biol.* 27:958–967.

863    Smith BT, Harvey MG, Faircloth BC, Glenn TC, Brumfield RT. 2014. Target capture and
864        massively parallel sequencing of ultraconserved elements for comparative studies at
865        shallow evolutionary time scales. *Syst. Biol.* 63:83–95.

866    Smith BT, Mauck WM, Benz B, Andersen MJ. 2018. Uneven missing data skews phylogenomic
867        relationships within the lories and lorikeets [Internet].

868    Smith BT, Merwin J, Provost KL, Thom G, Brumfield RT, Ferreira M, Mauck WM, Moyle RG,
869        Wright TF, Joseph L. 2023. Phylogenomic Analysis of the Parrots of the World
870        Distinguishes Artifactual from Biological Sources of Gene Tree Discordance. *Syst. Biol.*
871        72:228–241.

872    Smith SA, Brown JW, Walker JF. 2018. So many genes, so little time: A practical approach to

873     divergence-time estimation in the genomic era. *PLoS One* 13:e0197433.

874     Springer MS, Gatesy J. 2016. The gene tree delusion. *Mol. Phylogenet. Evol.* 94:1–33.

875     Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
876         thousands of taxa and mixed models. *Bioinformatics* 22:2688–2690.

877     Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
878         large phylogenies. *Bioinformatics* 30:1312–1313.

879     Tajima F. 1983. Evolutionary relationship of DNA sequences in finite populations. *Genetics*
880         105:437–460.

881     Tan X, Qi J, Liu Z, Fan P, Liu G, Zhang L, Shen Y, Li J, Roos C, Zhou X, et al. 2023.
882         Phylogenomics Reveals High Levels of Incomplete Lineage Sorting at the Ancestral Nodes
883         of the Macaque Radiation. *Mol. Biol. Evol.* [Internet] 40. Available from:
884         http://dx.doi.org/10.1093/molbev/msad229

885     Tea Y-K, Xu X, DiBattista JD, Lo N, Cowman PF, Ho SYW. 2021. Phylogenomic Analysis of
886         Concatenated Ultraconserved Elements Reveals the Recent Evolutionary Radiation of the
887         Fairy Wrasses (Teleostei: Labridae: Cirrhilabrus). *Syst. Biol.* 71:1–12.

888     Timilsena PR, Wafula EK, Barrett CF, Ayyampalayam S, McNeal JR, Rentsch JD, McKain MR,
889         Heyduk K, Harkess A, Villegente M, et al. 2022. Phylogenomic resolution of order- and
890         family-level monocot relationships using 602 single-copy nuclear genes and 1375 BUSCO
891         genes. *Front. Plant Sci.* 13:876779.

892     Valqui T. 2008. Phylogeogaphy of Nothoprocta Tinamous and the The Phylogeny of the
893         Tinamidae.Van Remsen J, editor. Available from:
894         https://search.proquest.com/openview/8d1be590e409938fa5a5afd6ca0b43b1/1?pq-
895         origsite=gscholar&cbl=18750&diss=y

896     Van Damme K, Cornetti L, Fields PD, Ebert D. 2022. Whole-Genome Phylogenetic
897         Reconstruction as a Powerful Tool to Reveal Homoplasy and Ancient Rapid Radiation in
898         Waterflea Evolution. *Syst. Biol.* 71:777–787.

899     Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, Kriventseva EV,
900         Zdobnov EM. 2017. BUSCO applications from quality assessments to gene prediction and
901         phylogenomics. *Mol. Biol. Evol.* [Internet]. Available from:
902         http://dx.doi.org/10.1093/molbev/msx319

903     Wen D, Yu Y, Nakhleh L. 2016. Bayesian Inference of Reticulate Phylogenies under the
904         Multispecies Network Coalescent. *PLoS Genet.* 12:e1006006.

905     Xi Z, Liu L, Davis CC. 2015. Genes with minimal phylogenetic information are problematic for
906         coalescent analyses when gene tree estimation is biased. *Mol. Phylogenet. Evol.* 92:63–71.

907     Zhang C, Mirarab S. 2022. Weighting by Gene Tree Uncertainty Improves Accuracy of Quartet-
908         based Species Trees. *Mol. Biol. Evol.* [Internet] 39. Available from:
909         http://dx.doi.org/10.1093/molbev/msac215

910     Zhang C, Rabiee M, Sayyari E, Mirarab S. 2018. ASTRAL-III: polynomial time species tree
911         reconstruction from partially resolved gene trees. *BMC Bioinformatics* 19:153.

912   Zhang F, Ding Y, Zhu C-D, Zhou X, Orr MC, Scheu S, Luan Y-X. 2019. Phylogenomics from
913        low-coverage whole-genome sequencing. *Methods Ecol. Evol.* 10:507–517.

914   Zhao M, Kurtis SM, White ND, Moncrieff AE, Leite RN, Brumfield RT, Braun EL, Kimball RT.
915        2023. Exploring Conflicts in Whole Genome Phylogenetics: A Case Study Within Manakins
916        (Aves: Pipridae). *Syst. Biol.* 72:161–178.

917

918  **Tables and Figures:**

919

| Dataset | Mean #PIS | StDev #PIS | Mean %PIS | StDev %PIS | Mean RF | StDev RF |
|---|---|---|---|---|---|---|
| **BUSCOs** | 922.5146 | 690.084 | 0.3717782 | 0.1033898 | 58.5241 | 21.11341 |
| **UCE100Flank** | 45.20298 | 28.84223 | 0.1347156 | 0.0828647 | 106.501 | 20.10551 |
| **UCE300Flank** | 206.8736 | 88.36566 | 0.2648982 | 0.1048726 | 59.1615 | 18.40835 |
| **UCE1000Flank** | 1035.713 | 250.5732 | 0.4178852 | 0.08261385 | 29.8675 | 8.468606 |

920

921  **Table 1: Means and standard deviations of parsimony informative sites (PIS) and**
922  **Robinson-Foulds Distances (RF) for each dataset analyzed for the study.**

**Figure 1: Concatenated Phylogeny of tinamous based on UCEs with 1,000 bp of flanking sequence.** All nodes have bootstrap value of 100% except where noted. Note relatively short internodal branch lengths in Clade A. Tinamou illustrations by TAC.

**Figure 2: Multi-species Coalescent (ASTRAL) Phylogeny of tinamous based on UCEs with 1,000 bp of flanking sequence.** All nodes have posterior probability of 1. Branches labeled B1–B5 denote rapidly diverging branches evaluated using quartet analysis (see Figure 7). Tinamou illustrations by TAC.

**Figure 3: Phylogenetic topologies for a difficult to resolve clade (Clade A) of *Crypturellus* tinamous.** The phylogenetic placement of multiple taxa was discordant among datasets and species tree approaches. These include (1) *C. erythropus* (black short-dashed branch), (2) *C. transfasciatus* (black long-dashed branch), (3) *C. atrocapillus* (pink solid branch), and (4) a misidentified genome downloaded from NCBI that was labeled as *C. undulatus*, but shows phylogenetic affinities with *C. strigulosus* and *C. erythropus* (GCA 013389825; pink long-hashed branch). The top row of trees shows topologies for each dataset estimated with RAxML after concatenation, whereas the bottom row of trees shows topologies estimated with ASTRAL. Numbers on the nodes correspond to posterior probabilities (MSC trees only) or bootstrap percentage values (concatenation trees only). Nodes with PP = 1.0 or bootstrap = 100% not shown. Scale bars under each topology indicate substitutions per site for concatenated trees or coalescent units for MSC trees. Bottom right heatmap shows Robinson-Foulds distances between all pairs of trees.
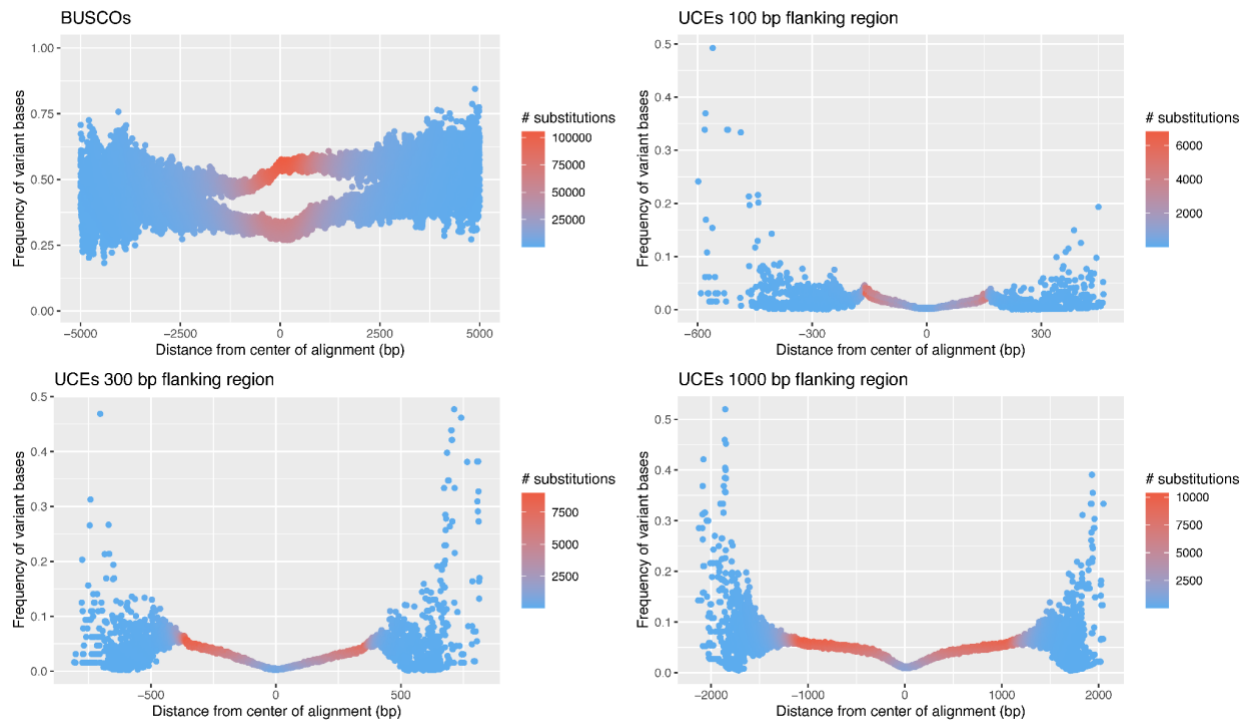
**Figure 4: Smilograms for each dataset showing that variation within coding genes is distributed differently than within UCEs.** Each point represents a base pair position, defined by its distance from the center, across all alignments in each dataset. UCEs show increasing frequency of variant sites with increasing distance from the UCE core. BUSCOs do not differ by distance from core, but instead show a bimodal distribution of frequency of variant bases likely associated with differences in variability between first and second versus third codon positions. Points are colored by the total number of substitutions at a given site across all samples and alignments. The outer regions of each plot show relatively few total substitutions because relatively few alignments in each dataset are very long.
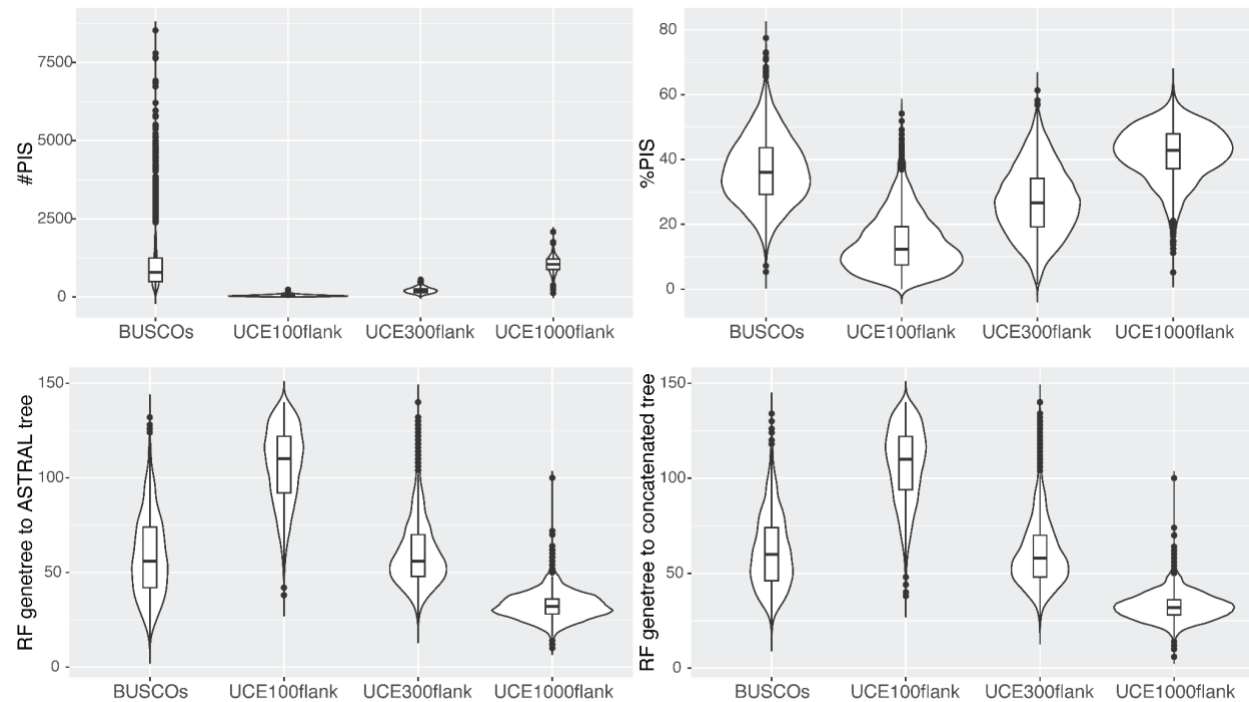
**Figure 5: Information content and gene tree heterogeneity for each dataset.** The top row of violin plots shows information content for each dataset using either the number of parsimony informative sites (#PIS) or the percentage of parsimony informative sites (%PIS) for each dataset. The bottom row of violin plots shows the gene tree heterogeneity for each dataset measured by the Robinson-Foulds distance between each gene tree and the inferred species tree (assuming the UCE1000Flank dataset). The bottom left plot assumes the MSC tree as the species tree, whereas the bottom right plot assumes the concatenated tree as the species tree. Because there is variation in mean and variance of RF distance among datasets, higher RF distances are indicative of more erroneous gene trees.
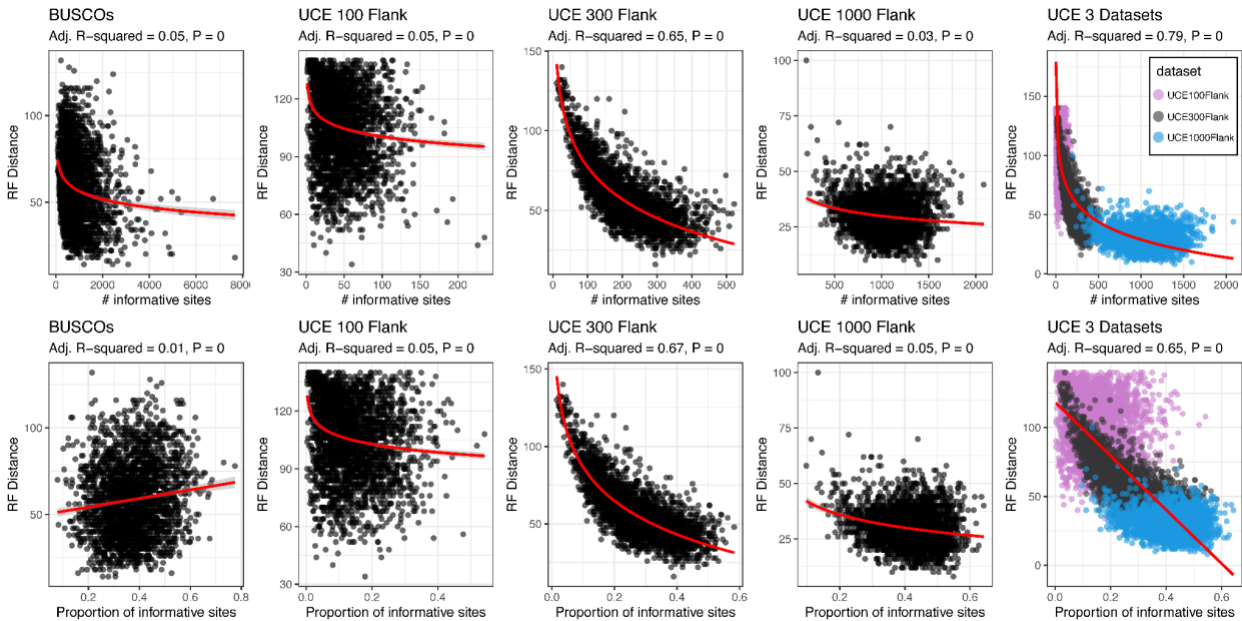
**Figure 6: Covariation between information content and gene tree estimation error.** Linear regression models exploring the relationship between information content and gene tree heterogeneity. The negative correlation in most instances suggests that alignments with less information content result in increased gene tree estimation error.
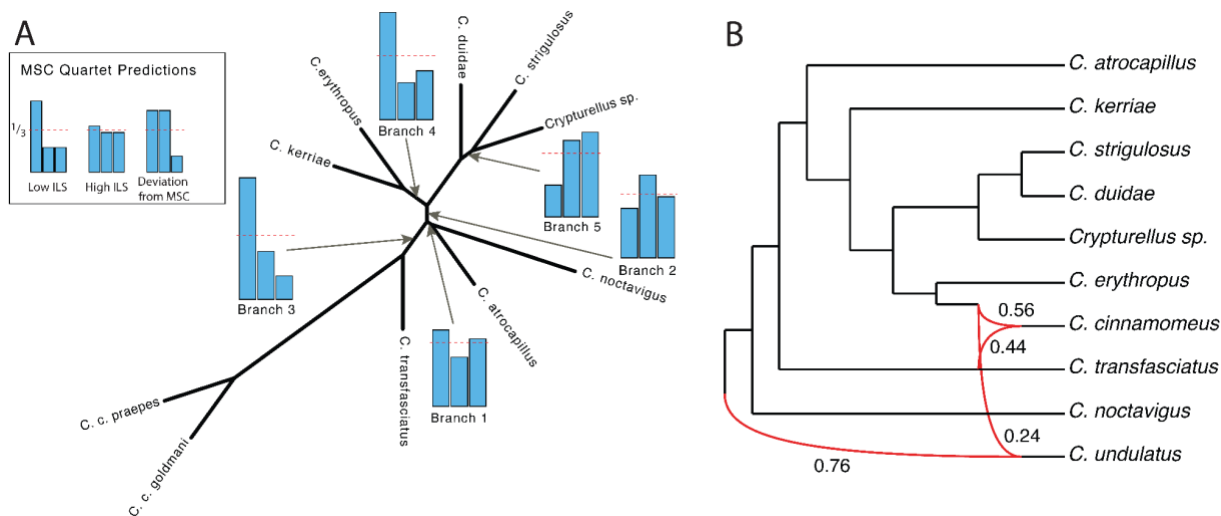
973
974 **Figure 7: Incomplete lineage sorting and genomic reticulation in the tinamou phylogeny.**
975 (A) Relative quartet frequencies for five short internal branches in Clade A based on the
976 UCE1000Flank dataset. Bar graphs depict the relative frequencies of each of three alternative
977 unrooted quartet topologies (for any unrooted four taxon statement there exist only three
978 alternative sets of relationships). The stippled lines indicate the ⅓ threshold for the frequency of
979 gene trees given a multispecies coalescent model, which predicts a single majority quartet
980 topology consistent with the true species tree with a frequency >⅓, and two minority topologies
981 with equivalent relative frequencies < ⅓. As quartet frequencies approach the ⅓ threshold, they
982 indicate stronger ILS. Major deviations from the expectations of the MSC model indicate
983 violations of the model may be present, such as gene tree estimation bias or introgression. (B)
984 Phylogenetic network results for the network with the optimal m-value (m=2) from PhyloNet
985 showing non-bifurcating relationships among a difficult to resolve clade of *Crypturellus* tinamous
986 (Clade A). Red branches demarcate reticulation, multiple ancestors contributing to a single
987 daughter branch. Inheritance probability, or the proportion of the discordant gene trees that can
988 be attributed to introgression, are marked on each reticulate branch. Branches are labeled 1–5
989 and correspond to labeled branches in Figure 2.