1    **A real data-based simulation procedure to select an imputation strategy for mixed-type**

2    **trait data**

3

4    Jacqueline A. May[1]*, Zeny Feng[2], Sarah J. Adamowicz[1]

5

6

7

8    *¹ Department of Integrative Biology & Biodiversity Institute of Ontario, University of Guelph,*

9    *Guelph, Ontario, Canada.*

10   *² Department of Mathematics & Statistics, University of Guelph, Guelph, Ontario, Canada.*

11

12

13   * Corresponding author

14   Email: mayj@uoguelph.ca

# Abstract

15

16      Missing observations in trait datasets pose an obstacle for analyses in myriad biological

17   disciplines. Imputation offers an alternative to removing cases with missing values from datasets.

18   Imputation techniques that incorporate phylogenetic information into their estimations have

19   demonstrated improved accuracy over standard techniques. However, previous studies of

20   phylogenetic imputation tools are largely limited to simulations of numerical trait data, with

21   categorical data not evaluated. It also remains to be explored whether the type of genetic data

22   used affects imputation accuracy. We conducted a real data-based simulation study to compare

23   the performance of imputation methods using a mixed-type trait dataset (lizards and

24   amphisbaenians; order: Squamata). Selected methods included mean/mode imputation, $k$-nearest

25   neighbour, random forests, and multivariate imputation by chained equations (MICE). Known

26   values were removed from a complete-case dataset to simulate different missingness scenarios:

27   missing completely at random (MCAR), missing at random (MAR), and missing not at random

28   (MNAR). Each method (with and without phylogenetic information derived from mitochondrial

29   and nuclear gene trees) was used to impute the removed values. The performances of the

30   methods were evaluated for each trait and in each missingness scenario. A random forest method

31   supplemented with a nuclear-derived phylogeny performed best overall, and this method was

32   used to impute missing values in the original squamate dataset. Data with imputed values better

33   reflected the characteristics and distributions of the original data compared to the complete-case

34   data. However, phylogeny did not always improve performance for every trait and in every

35   missingness scenario, and caution should be taken when imputing trait data, particularly in cases

36   of extreme bias. Ultimately, these results support the use of a real data-based simulation

37   procedure to select a suitable imputation strategy for a given mixed-type trait dataset. Moreover,

38   they highlight the potential biases that complete-case usage may introduce into analyses.

## Author summary

40    The issue of missing data is problematic in trait datasets as observations for rare or threatened

41    species are often missing disproportionately. When only complete cases are used in an analysis,

42    derived results may be biased. Imputation is an alternative to complete-case analysis and entails

43    filling in the missing values using known observations. It has been demonstrated that including

44    phylogenetic information in the imputation process improves accuracy of predicted values.

45    However, most previous evaluations of imputation methods for trait datasets are limited to

46    numerical, simulated data, with categorical traits not considered. Using a reptile dataset

47    comprised of both numerical and categorical trait data, we employed a real data-based simulation

48    strategy to select an optimal imputation method for the dataset. We evaluated the performance of

49    four different imputation methods across different missingness scenarios (e.g. missing

50    completely at random, values missing disproportionately for smaller species. Results indicate

51    that imputed data better reflected the original dataset characteristics compared to complete-case

52    data; however, the optimal imputation strategy for a given scenario was contingent on

53    missingness scenario and trait type. As imputation performance varies depending on the

54    properties of a given dataset, a real data-based simulation strategy can be used to provide

55    guidance on best imputation practices.

56

# Introduction

Trait data are used in a wide variety of biological disciplines, including evolutionary biology, community ecology, and biodiversity conservation. For instance, trait data pertaining to the life history of a species, such as longevity, metabolic rate, and generation time, are integral in studies of biological aging (1,2). Environmental trait data, such as latitude, temperature, and habitat type, may be used to identify those species most at risk of extinction (3,4). However, an extensive proportion of these trait data are often missing. Missingness may stem from a taxonomic bias: data are available in copious amounts for well-researched or charismatic species and are lacking for endangered species or those that inhabit remote environments (e.g. deep sea) (5–7). Mammal and bird taxa tend to be well sampled, and data for a large and diverse array of traits are available for many groups (8,9). However, regional and phylogenetic biases are common in trait data for groups such as reptiles and amphibians, and observations are largely limited to body size and habitat traits (9). Species traits are often tied to evolutionary history, a concept referred to as phylogenetic signal (10). Closely related species can share the characteristics that render them elusive or difficult to study (e.g. small body size), resulting in sparse or unreliable data for entire taxonomic clades (5,6,8). Certain types of trait data may also be easier to quantify (e.g. morphometric data) as opposed to traits that require arduous or invasive data collection techniques (e.g. age or reproductive data) (11; see Fig 1 for a visualization of missingness in reptiles). When trait datasets are used in studies, these biases can lead researchers to make erroneous conclusions about the data. Consequently, the development of approaches for handling missing data is an important area of research that spans across multiple biological disciplines.

79    **Fig 1. Visualization of missingness.** Visualization of missingness (proportion of present vs.

80    missing observations) in Squamata trait data obtained from the primary literature. Superscripts

81    indicate the original sources of the trait data: 1) amniote life history database (12,13), 2)

82    vertebrate home range sizes dataset (14,15), 3) traits of lizards of the world (16,17) and 4)

83    AnAge (18,19). See S1 File for further detail on trait sources.

84        The use of complete-case datasets can result in a large proportion of information being

85    discarded (7,20). If data are "missing completely at random" (MCAR), the removal of cases

86    leads to a reduction in the size of the dataset, and in turn, a reduction in statistical power (7,21).

87    Trait data, however, are often "missing at random" (MAR): observations that are missing for a

88    particular trait are related to known values for some other traits. Simply removing incomplete

89    cases when data are MAR can result in biased estimations of model parameters (7,11,22). In

90    more extreme cases, trait data may be "missing not at random" (MNAR): the reason data are

91    missing is related to the unobserved data themselves. In such scenarios, the reason for

92    missingness may be unclear to the researcher and thus difficult to verify empirically (23).

93        Imputing missing observations is a common alternative to the complete-case analysis.

94    Imputation techniques use known observations to estimate the missing and unobserved values of

95    a variable (or variables) of interest. Single imputation techniques such as hot deck imputation or

96    $k$-nearest neighbour (KNN; 16) offer an efficient means for estimating missing values; however,

97    these methods provide only a single estimate of the missing value. Random forest methods such

98    as missForest (25) are also growing in popularity as they make no prior assumptions about the

99    distributions of variables. Multiple imputation techniques have been developed that perform

100    single imputation several times and are therefore capable of providing a measure of uncertainty

101    of the imputed values (7,26). An example of a multiple imputation method is multivariate

102  imputation by chained equations (MICE; 19), which offers numerous models for imputing data

103  of different types. Incorporating phylogenetic information into the imputation process has also

104  been shown to increase the accuracy of imputed values (11,28). This increase in accuracy is a

105  result of the phylogenetic signal that is often inherent in trait data. A commonly used method for

106  incorporating phylogenetic information into the imputation process is the use of phylogenetic

107  eigenvectors. More specifically, methods such as phylogenetic eigenvector regression (PVR)

108  (29) and phylogenetic eigenvector mapping (PEM) (30) employ a principal coordinates analysis

109  (PCoA) to derive eigenvectors from a phylogenetic tree. PEM expands on the PVR method by

110  applying an additional branch length transformation based on the Ornstein-Uhlenbeck

111  evolutionary model (30,31). Phylogenetic eigenvectors may then be used as additional predictor

112  variables in the imputation process (see 11,24,25).

113        As missing data are a major concern in trait datasets, we are motivated to consider

114  imputing these missing values. The correlative nature of trait data makes them suitable

115  candidates for imputation, particularly when phylogenetic signal is also present (34). In an

116  evaluation of imputation methods using mammalian trait data, Penone *et al.* (11) found that

117  supplementing the imputation process with phylogenetic information improved the accuracy of

118  KNN, missForest, and MICE for several life history traits. Kim *et al.* (24) similarly found that

119  adding phylogenetic information to MICE improved accuracy rates of estimated functional

120  diversity metrics. However, when imputing bird demographic traits with moderate phylogenetic

121  signal (Pagel's $\lambda < 0.8$), Johnson *et al.* (27) found that use of phylogenetic information improved

122  error rates by a margin of less than 1%. Moreover, they suggest that the use of auxiliary traits

123  (traits that are present in the dataset but not the target of imputation) were often sufficient for

124  accurate imputations. In sum, these findings indicate that improvements conferred by

125    phylogenetic imputation methods are context-dependent, contingent upon the presence of

126    phylogenetic signal and relationships among traits in the dataset.

127         Trait data exist in several forms, ranging from the discrete categories of foraging

128    behaviour to the countable number of eggs in a nest. Available trait datasets are often comprised

129    of mixed types that contain categorical, count, and numerical data. Many contemporary

130    imputation methods are able to estimate both categorical and numerical values. However, most

131    previous studies have only evaluated their performances using simulated trait data, and the few

132    studies that have utilized real data are limited to numerical traits. Additionally, phylogenetic

133    information is usually included in the form of a multigene tree; it remains to be explored whether

134    the type of genetic data used to construct the phylogeny affects imputation accuracy.

135    Phylogenetic resolution varies among gene trees (36,37), and certain genes may be more or less

136    suited for imputation in a given taxon and taxonomic rank. To determine the best-suited

137    imputation method for a given mixed-type dataset, we propose a method-selection strategy that

138    employs real data-based simulations. Results from the real-data simulations will address: 1)

139    whether there is an optimal imputation strategy for a specific data type (continuous, count,

140    categorical) and missingness scenario (MCAR, MAR, and MNAR); 2) which imputation method

141    performs the best for a given dataset containing mixed data types; 3) whether phylogenetic

142    information improves the imputation performance; and 4) which type of phylogenetic

143    information is influential (mitochondrial, nuclear). The strategy proposed here may be

144    considered for future trait-based analyses to reduce biases that may occur if researchers analyze

145    only complete cases, bolster sample size and improve statistical power, and mitigate error rates

146    when imputing missing values. In turn, this will facilitate the pursuit of new research directions,

147    particularly in those fields impeded by sparsely available trait data.

# Results

## Performance comparison without phylogeny

In general, when missing data were generated under MCAR, error rate increased with missingness proportion; this trend was observed for all trait and method combinations (Fig 2). Under the same simulation setting, *k*-nearest neighbour (*KNN;* 24,38), random forests (*RF*; "missForest" R package 25,39) and multivariate imputation by chained equations (*MICE;* 27) outperformed mode and mean imputation for the majority of traits. However, there were exceptions to this pattern. For the categorical trait activity time, mode imputation resulted in a lower error rate than *RF* and *KNN* at 30-40% missingness (Fig 2a). Additionally, for smallest clutch, the mean imputation method outperformed *KNN* (10-40%) and *RF* (10%, 30-40%) (Fig 2d). *MICE* resulted in lower error rates than *KNN* and *RF* for most traits across all missingness proportions. However, *KNN* resulted in the lowest error rate for activity time at 10%, and *RF* resulted in the lowest error rate across all missingness proportions settings for the insular endemic trait (Fig 2b) and at 10% missingness for largest clutch (Fig 2g). In both MAR and MNAR scenarios without phylogenetic information added, *MICE* generally outperformed both *RF* and *KNN* (see Fig 3).

**Fig 2. MCAR performance without phylogeny.** Performance of the methods mean imputation, *KNN*, missForest (*RF*), and multivariate imputation by chained equations (*MICE*) across different proportions of missingness when data were MCAR. *MICE_LR* and *MICE_PMM* signify the use of logistic regression and predictive mean matching for imputing categorical and numerical traits, respectively. Error rate was measured as PFC for the categorical traits a) activity time and b) insular endemic and as MSE for the numerical traits c) largest clutch, d) smallest clutch, e) female snout-vent length (SVL), f) maximum SVL, and g) latitude. In both cases, error rates closer to 0 are indicative of better performance.

172    **Fig 3. Imputation performance across all missingness scenarios.** Comparison of error rates

173    for the methods mode imputation, *KNN*, *RF*, and *MICE* for different missingness scenarios with

174    and without the addition of phylogenetic information. Phylogenetic information was added in the

175    form of trees built from sequence data of mitochondrial cytochrome *c* oxidase subunit I (COI)

176    and nuclear oocyte maturation factor (c-mos) and recombination activating gene 1 (RAG1).

177    Performance was quantified using PFC for the categorical traits a) activity time and b) insular

178    endemic and using MSE for the numerical traits c) largest clutch, d) smallest clutch, e) female

179    SVL, f) maximum SVL, and g) latitude. MCAR = missing completely at random; MAR =

180    missing at random; MNAR = missing not at random.

## Phylogenetic imputation performance

182        All traits exhibited significant phylogenetic signal in all gene trees (S1 Fig; see S1 File

183    for more details on phylogenetic signal measures). However, improvements to imputation

184    performance through the addition of phylogeny were contingent on method, data type, and

185    missingness scenario (Fig 3). For instance, when considering the categorical trait activity time,

186    supplementing phylogenetic information from any of the three genes generally improved

187    performance for each method and in each missingness scenario (Fig 3a). On the contrary, in the

188    case of the binary trait insular endemic, adding phylogenetic information to *MICE* at low

189    missingness levels (10%) resulted in an increased error rate (Fig 3b). For most traits, MAR

190    results reflected those in the MCAR scenarios; however, deviations from the general pattern

191    occurred in some MNAR cases. For example, in the MNAR scenario for insular endemic,

192    phylogeny was only beneficial when nuclear information was added to *KNN*.

193        For the traits largest clutch, smallest clutch, and latitude, *KNN* and *RF* performances were

194    improved by the addition of any type of phylogenetic information in the MCAR and MAR

9

195    scenarios; this was particularly evident in the case of nuclear oocyte maturation factor (c-mos)

196    (Fig 3c-d, g). However, phylogeny did not improve *MICE* performance consistently for these

197    traits. In the MAR scenarios, phylogenetic information improved *MICE* performance for smallest

198    clutch and latitude; conversely, for largest clutch, any type of phylogenetic information increased

199    error rate for *MICE*. In the MNAR scenarios, the addition of any type of phylogenetic

200    information increased error rate for *MICE* imputation for all of these traits drastically in several

201    situations (e.g. more than doubling the error rate for largest clutch and latitude). The traits female

202    snout-vent length (SVL) and maximum SVL displayed somewhat dissimilar patterns from the

203    other traits (Figs 4e-f) as phylogenetic information tended to decrease imputation performance

204    for most methods and in most scenarios.

205        The relationship between phylogenetic signal and error ratio varied depending on data

206    type. For categorical traits, higher error ratio, indicative of better performance due to phylogeny,

207    was associated with higher phylogenetic signal strength (Fig 4a). This same pattern was not

208    observed for numerical traits (Fig 4b). Moreover, in MNAR scenarios for numerical traits, many

209    error ratio values fell below 1 at higher levels of phylogenetic signal, indicative of a reduction in

210    performance due to phylogeny. Generally, the improvement in imputation performance resulting

211    from phylogeny was most apparent for *KNN* and *RF*, as these methods account for the majority

212    of error ratio values greater than 1; error ratio values for *MICE*, however, often fell below 1,

213    particularly in the case of numerical traits.

214    **Fig 4. Association between error ratio and phylogenetic signal.** Association between error

215    ratio (error rate without phylogeny/error rate with phylogeny) and phylogenetic signal for the c-

216    mos gene (Fritz and Purvis' *D* (40) for categorical traits and Pagel's λ (41) for numerical traits)

217    at different proportions of missingness. Error ratio values above 1 (indicated by the gray line)

218    signify an improvement in performance when phylogeny is added. In the case of *D,* lower values

219    are indicative of higher levels of phylogenetic conservation for the trait; conversely, higher

220    values of λ suggest stronger phylogenetic signal. Results are not shown for MAR in a) as only

221    one trait (activity time) was simulated for this scenario. To improve visualization, values were

222    jittered (random noise introduced to data) using the package "ggplot2" (42). Additionally, results

223    are only shown for the c-mos gene as results for cytochrome *c* oxidase subunit I (COI) and

224    recombination activating gene 1 (RAG1) follow similar patterns.

## Imputation of original dataset using best strategy

226    Although results varied considerably, particularly in MNAR scenarios, the method that

227    resulted in the lowest error rates overall was *RF* with c-mos. Consequently, this method was

228    chosen to impute the original dataset. Out of the total species in the original dataset ($n = 6657$),

229    those with available c-mos sequence records were included in the imputed subset ($n = 921$). The

230    proportion of missingness varied for each trait in this subset as 0.16 for activity time, 0 for

231    insular endemic, 0.21 for largest clutch, 0.21 for smallest clutch, 0.23 for female SVL, 0 for

232    maximum SVL, and 0 for latitude. As insular endemic, maximum SVL, and latitude had

233    complete observations in this subset, these traits were not imputed.

234    Distributions and categorical frequencies of the complete-case, original, and imputed data

235    can be observed in Fig 5. For the trait activity time, when compared to the original data,

236    discrepancies in the categorical frequencies were more apparent in the complete-case data than in

237    the imputed data (Fig 5a). The complete-case data displayed a greater overrepresentation of the

238    rarest category (cathemeral: 11% vs. 8.9%) and underrepresentation of the most common

239    category (diurnal: 57.9% vs. 64.4%). Conversely, the imputed data displayed a greater

240    representation of observations in the most common category compared to the original data

241    (diurnal: 67.5% vs. 64.4%). For all numerical traits, the imputed data distributions followed the

242    distributions of the original more closely than did the complete-case distributions (Figs 6b-d;

243    Table 1). Perhaps most apparent are the discrepancies in the maximum values in the complete-

244    case data compared to those in the original and imputed data (e.g. for largest clutch, 68 vs. 88;

245    for smallest clutch, 8 vs. 30). Although the discrepancies in the complete-case data were greater,

246    both complete-case and imputed data displayed reduced variance relative to the original data for

247    the traits largest clutch, smallest clutch, and female SVL.

248    **Fig 5. Comparison of quantitative characteristics across datasets.** Comparison of a)

249    categorical frequencies for the trait activity time and distributions for the traits b) largest clutch,

250    c) smallest clutch, and d) female SVL of the complete-case, original, and imputed data. The

251    natural logarithm (ln) of the numerical data were taken to improve visualization.

# Discussion

253    In agreement with previous evaluations of imputation methods using trait data (11,34,43),

254    there was no "optimal" method for imputing values in all scenarios. In the absence of phylogeny,

255    the best overall method for imputing mixed-type trait data was *MICE*. This trend was apparent

256    even in cases of MNAR, as *MICE* resulted in the lowest error rates for five out of seven traits in

257    these scenarios when phylogeny was not included. *MICE* demonstrated strong performances in

258    previous evaluations of imputation techniques in mammalian (11) and plant (43) trait datasets.

259    Furthermore, the robustness of predictive mean matching is appealing for the non-linear

260    relationships and non-normal distributions commonly observed in numerical trait data (44,45).

261    This may explain the superior performance of *MICE* in the case of smallest clutch, a count trait

262    with a right-skewed distribution (many species with smallest clutch size = 1).

**Table 1. Summary statistics for the complete-case, original, and imputed datasets.**

| | Largest clutch (# eggs/neonates) | | | Smallest clutch (# eggs/neonates) | | | Female SVL (mm) | | | Maximum SVL (mm) | | Latitude (°) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *CC* | *O* | *I* | *CC* | *O* | *I* | *CC* | *O* | *I* | *CC* | *O* | *CC* | *O* |
| **N** | 141 | 731 | 921 | 141 | 731 | 921 | 137 | 705 | 921 | 152 | 921 | 152 | 921 |
| **Min** | 1 | 1 | 1 | 1 | 1 | 1 | 18.7 | 18.7 | 18.7 | 21.7 | 21.7 | -40.36 | -47.89 |
| **Max** | 68 | 88 | 88 | 8 | 30 | 30 | 499.5 | 534.3 | 534.3 | 1170 | 1170 | 56.6 | 56.6 |
| **Range** | 67 | 87 | 87 | 7 | 29 | 29 | 480.8 | 515.6 | 515.6 | 1148.3 | 1148.3 | 96.96 | 104.49 |
| **Median** | 2 | 3 | 3 | 1 | 2 | 2 | 60.1 | 62.7 | 65.2 | 77 | 80 | -11.36 | -9.48 |
| **Mean** | 5.79 | 6.08 | 6.06 | 1.63 | 2.06 | 2.14 | 75.82 | 83.4 | 84.36 | 103.44 | 110.35 | 1.65 | -3.8 |
| **SE (mean)** | 0.67 | 0.33 | 0.27 | 0.09 | 0.07 | 0.06 | 4.76 | 2.43 | 2.07 | 8.58 | 3.30 | 2.03 | 0.75 |
| **0.95 CI (mean)** | 1.33 | 0.65 | 0.53 | 0.18 | 0.14 | 0.12 | 9.41 | 4.76 | 4.06 | 16.96 | 6.48 | 4.01 | 1.48 |
| **Variance** | 63.85 | 79.46 | 67.41 | 1.22 | 3.90 | 3.46 | 3103.49 | 4151.01 | 3950.16 | 11197.05 | 10048.39 | 627.38 | 521.53 |
| **Standard deviation** | 7.99 | 8.91 | 8.21 | 1.10 | 1.98 | 1.86 | 55.71 | 64.43 | 62.85 | 105.82 | 100.24 | 25.05 | 22.84 |

264     Summary statistics of the complete-case (*CC*), original (*O*), and imputed (*I*) datasets for the numerical traits largest clutch, smallest clutch, female

265     snout-vent length (SVL), and latitude. As the proportion of missingness was 0 for the traits maximum SVL and latitude in the original data subset,

266     these traits were not imputed. Original trait data obtained from Meiri (16).

267   Predictive mean matching has also been shown to perform well on smaller sample sizes (45), as

268   seen in the current study ($n = 152$). Its use in trait imputation is therefore an appealing option

269   when phylogenetic information is scarce.

270        As reported in previous studies (11), imputation error rates tended to increase with

271   missingness proportion and varied amongst different traits. Adding phylogenetic information,

272   however, did not always improve imputation performance; on the contrary, in some instances its

273   inclusion led to increased error rates. The effect of phylogeny therefore appears to be situational

274   and linked to the method used, the underlying mechanism of the missingness in the data, and

275   quantitative attributes and evolutionary history of the target trait. The performances of *KNN* and

276   *RF* were often improved when any type of phylogenetic information was provided, even in some

277   cases of MNAR. This pattern was more prominent at higher missingness proportions, as

278   phylogeny can offset the loss of the trait data. Conversely, phylogeny often increased the error

279   rate for *MICE*. This increase in error rate was also found in Johnson *et al.* (34) when

280   phylogenetic information was added to *MICE*, particularly in MNAR scenarios (e.g. larger

281   values more likely to be missing). The authors suggest this may stem from an issue relating to

282   the large number of eigenvectors used in the imputation process (e.g. more than 20 eigenvectors

283   were included in biased missingness scenarios). Penone *et al.* (11) restricted their maximum

284   number of eigenvectors to 10 and suggest that the use of too many eigenvectors can mask the

285   information provided by other traits in the imputation process. Indeed, in the current study,

286   *MICE* performed well when the number of predictors were low, as in the case of trait-only

287   imputation. As phylogenetic resolution varies between nuclear and mitochondrial gene trees, the

288   number of eigenvectors used for imputation varied in accordance. In this study, the 65%

289   variation method was used to determine the number of eigenvectors to be included; however, it is

14

290   possible that the use of too many eigenvectors (i.e. more than 40; 62), with less information

291   provided by each eigenvector, would introduce more noise or lead to overfitting by the

292   regression-based models. Thus, analyses using phylogenetic eigenvectors for imputation may

293   consider the use of tree-based methods such as *RF* (or recursive partitioning; see Kim *et al.* (32))

294   that are more robust to high-dimensional data. Future studies may also consider exploring

295   whether the optimal number of phylogenetic eigenvectors to use for imputation changes under

296   varying degrees of missingness bias.

297         *RF* with phylogeny demonstrated the strongest performance overall as it resulted in the

298   lowest error rates across all missingness scenarios. This result supports previous evaluations of

299   the effectiveness of *RF* for mixed-type data (25). For both *KNN* and *RF*, adding phylogenetic

300   information reduced imputation error rate for traits of all types (categorical, count, continuous).

301   Nuclear-derived phylogenetic information (i.e. c-mos or RAG1) generally conferred a greater

302   improvement in imputation performance relative to mitochondrial COI. Due to their faster rates

303   of nucleotide substitution, mitochondrial genes are less adept at resolving deeper phylogenetic

304   relationships relative to nuclear genes (47). Consequently, the relationships resolved by nuclear

305   gene trees may more closely follow the evolutionary trajectory of the traits used in this study.

306   However, COI often still conferred a reduction in error rate, in some cases more so than the

307   nuclear genes (e.g. smallest clutch); mitochondrial sequences therefore should be used when

308   nuclear data are unavailable and may be more advantageous when studying more closely related

309   species. Strength of phylogenetic signal also appeared to correlate with error ratio (i.e. the

310   magnitude of performance enhancement) for categorical traits. The same pattern was not

311   apparent for numerical traits, however. This may stem from the limited range of phylogenetic

312   signal observed for these other types: all genes displayed significant levels of phylogenetic signal

313     for all traits, many of which verged toward $\lambda = 1$ (higher trait conservation). This may suggest

314     that the boost in performance due to phylogeny is negligible beyond a certain level of

315     phylogenetic signal. However, imputation of a greater number and variety of traits that do not

316     display any evidence of phylogenetic signal would need to be included to test this assertion.

317        The comparison between the distributions and categorical frequencies of the complete-

318     case, original, and imputed trait data support the efficacy of imputation for mixed-type data. A

319     greater than 6-fold increase in sample size when using imputed data ($n = 151$ for complete-case

320     vs. $n = 921$ for imputed data) is striking and illustrates the information loss that can occur when

321     using a complete-case approach. Moreover, complete-case data often do not capture the true

322     variability of the data; instead, they comprise a biased subset and, in turn, the potential for

323     erroneous inferences. Previous studies using clinical data (64) and mammalian trait data (11)

324     found that inferences derived from imputed datasets are less biased when compared to those

325     obtained using complete-case datasets. However, the missing values in these studies were

326     introduced either completely at random (MCAR) or at random (MAR). Although imputation

327     performs well under MCAR and MAR, the mechanism of missingness is often difficult to

328     determine in practice (23,49). Imputation has been shown to perform poorly in scenarios with

329     biased missingness, such as when extreme values or values in the tails of the distribution of the

330     population are disproportionately missing (34). The results from our study provide reason for

331     further discretion in these instances as the most extreme error rates were observed in MNAR

332     scenarios. If data are truly MNAR and the imputation method is not carefully chosen, imputed

333     values and the inferences derived therein may be inaccurate. A recent study completed by Jardim

334     et al. (50) suggests that accurate estimation of phylogenetic signal from imputed datasets is

335     contingent on several variables, including the amount of missing data, missing mechanism, and

336     the evolutionary trajectory of the trait itself. For example, as values closer to the equator were

337     missing in the latitude MNAR scenario, mean imputation outperformed most other imputation

338     methods. Due to the prevalence of allopatric speciation modes in diversification (51,52), closely

339     related species can inhabit different latitudes or distributions; traits with such evolutionary

340     histories may be less suitable for imputation. Therefore, we agree with Johnson *et al.* (34) and

341     Jardim *et al.* (50) that caution should be taken when imputing data and the properties of the

342     dataset of interest be inspected beforehand. Testing imputation methods using a real data-based

343     simulation strategy as we demonstrate here would provide useful insight as to whether

344     imputation is a suitable alternative to complete-case analysis.

345          As is often the case when constructing a complete-case dataset, several traits were

346     excluded from this study. These included many categorical traits that were invariant in the

347     complete-case dataset, such as those containing information about geography or habitat. In turn,

348     the range of phylogenetic signal for traits was also limited. It was therefore not feasible to truly

349     gauge the relationship between error ratio and phylogenetic signal strength in traits as they all

350     exhibited significant levels. The continued collection of high-quality trait data for both known

351     and novel species is necessary to further probe these types of relationships. For instance, in the

352     case of Squamata, snake species are disproportionally undersampled (9) and were thus not

353     included in the current study. An increase in data availability would also facilitate additional

354     research on the use of imputation methods in real datasets. Simulated trait data do not fully

355     capture the nuances of real datasets, and comparative evaluations using real data and different

356     taxonomic groups are needed to test whether imputing values is practical, particularly in cases of

357     severe biases.

358    Missingness in datasets is a pervasive issue in the realm of biological research. It is

359    particularly problematic for those taxonomic groups threatened by extinction, or that are small or

360    reside in understudied areas of the globe. As trait data can take on many forms, methods that can

361    accurately predict missing values for diverse data types are invaluable for the study of these

362    obscure groups. Previous research has focused largely on numerical data, and consideration of

363    imputation performance for categorical traits is imperative in driving this field forward. The

364    results presented here provide support for the use of imputation methods in real mixed-type

365    datasets. Supplementing these methods with phylogenetic information is often beneficial, even if

366    sequence data are available for only one or a limited number of markers. However, researchers

367    should take care to understand the properties of their dataset and consider the ramifications of

368    using imputation. In such situations, a real data-based simulation strategy can provide guidance

369    on best imputation practices for a given biological or ecological dataset. Simulating missingness

370    using real data more accurately reflects the characteristics and the nature of the unobserved

371    values. The imputation method that is robust in these scenarios and across diverse trait types can

372    be used to bolster sample size while simultaneously preserving the original properties of a

373    dataset. Derived inferences may then more accurately represent the biological phenomena under

374    investigation.

## Materials and methods

### Complete-case dataset creation

377    Traits are defined here as characteristics that are typical of a species. These may refer to

378    characteristics relating to the biology of a species or the environment in which it resides. Data for

379    squamates (lizards and amphisbaenians; order: Squamata) were selected for analysis as

380    complete-case observations were available for at least 100 species as well as both categorical and

381    numerical traits. In addition, these species had DNA sequence records publicly available for both

382    mitochondrial and nuclear markers. Squamata represent an incredibly diverse group of

383    vertebrates (~10,000 species; 30), inhabiting disparate environments and boasting a broad range

384    of morphological features. However, trait data for Squamata are undersampled relative to

385    mammal and bird groups, particularly in tropical regions that are home to diverse species at risk

386    (9). As of 2022, 19.6% of squamate species are estimated to be under threat of extinction (54).

387    Imputation may offer additional avenues to identify those traits correlated with risk status in

388    squamates (e.g. 32,33) and in doing so, contribute to biodiversity conservation efforts in

389    vulnerable areas. Trait data were obtained from a dataset published by Meiri (16) (other datasets

390    were also considered, see S1 File). This dataset contains information about the habitat, life

391    history, morphology, behaviour, and conservation threat level of 6,657 squamate species (lizards

392    and amphisbaenians, not including snakes) (34,35). The raw trait data were downloaded into R v.

393    4.0.3 (57).

394        The Barcode of Life Data System (BOLD) (58) was used as the source for mitochondrial

395    sequence data as it contains thousands of published cytochrome $c$ oxidase subunit I (COI) partial

396    gene sequence records (16,676 sequences for over 2000 Squamata species as of July 16[th], 2021).

397    COI sequence data were downloaded into R on March 12[th], 2020 (59). Data were filtered for

398    records that have been identified to the species level, as this information was necessary for trait

399    matching purposes. Additional quality control checks on the sequence data included trimming N

400    and gap content from sequence ends and removing sequences with greater than 1% of internal N

401    and/or gap content across their entire sequence length. Sequences between 650 and 1000 bp were

402    retained to facilitate downstream multiple sequence alignment. As multiple COI sequence

403    records are available for many species, a centroid sequence selection process was employed to

404    find a typical representative sequence for each species (Orton et al., 2019; see S1 File for details

405    on this process). The *AlignTranslation* function from the R package "DECIPHER" v. 2.18.1

406    (61,62) was used to perform a multiple sequence alignment on the centroid sequences.

407    *AlignTranslation* was used as it performs a multiple sequence alignment guided by the translated

408    amino acid sequence, which is more reliable than an alignment based on nucleotide data alone

409    (61). The translated final alignment was visualized using the *ggmsa* function from the R package

410    "ggmsa" v. 0.06 (42) to verify the nucleotides were in the correct reading frame and to check for

411    the presence of stop codons. Nuclear sequence data were obtained from a multigene alignment

412    published in Pyron *et al.* (64,65). This alignment is comprised of sequence data for 12 genes

413    (seven nuclear, five mitochondrial) and 4161 species of Squamata (64). The alignment was

414    partitioned into its constituent gene alignments using RAxML v. 8 (66).

415        Species names from the COI alignment were matched against the species names in the

416    trait dataset. Those species that had available data for at least five traits (both categorical and

417    numerical) and a corresponding COI sequence record were then matched against the species

418    names in the nuclear multigene alignment. The nuclear markers oocyte maturation factor (c-mos)

419    and recombination activating gene 1 (RAG1) had the largest number of available records for the

420    species in the complete-case dataset and were selected for analyses (see S1 Table for sequence

421    identifiers of those records selected). Final checks were performed on the trait data in the

422    complete-case subset. Categorical traits with severe class imbalances and very low variability

423    (e.g. more than 90% of observations in one of the categories and/or the remaining observations

424    sparsely dispersed across other categories), such as reproductive mode, geographic range, and

425    substrate, were excluded from the study. The distributions of numerical trait data were visualized

426    to check for the presence of severe outliers. For each numerical trait, an upper threshold was

427    calculated as follows: quartile 3 + (3 × the interquartile range of the data). Severe outliers are

428    defined here as those values that exceed the upper threshold. If identified, these values were

429    verified in the primary literature to ensure they were real datapoints and not the result of data

430    entry error. The final dataset contains information for the seven most complete traits, including

431    the categorical traits: activity time and insular endemic, the count traits: largest clutch and

432    smallest clutch, and the continuous traits: female snout-vent length (SVL), maximum snout-vent

433    length (SVL), and latitude (geographic centroid for the species; Roll et al. 2017). The final

434    dataset is referred to as the "complete-case dataset" including, 152 species, representing 25

435    Squamata families (S2 Table). To maintain a sufficient sample size, we permitted some missing

436    values (no more than 10% for each trait) present in the so-called "complete-case dataset";

437    otherwise, the sample size will drop to 121 if only species without missing values in their traits

438    are included. For further details on these traits, see S3 Table.

## Phylogenetic information

439

440    The alignments for the COI, c-mos, and RAG1 sequences were used to build maximum

441    likelihood gene trees in RAxML v. 8 (66). The model GTRGAMMAI was specified (option -m),

442    and the alignment was partitioned based on codon position (option -q). The gene trees were then

443    read into R and made ultrametric using the *chronos* function in the R package "ape" v. 5.4.1

444    (68). Phylogenetic eigenvectors were extracted from each gene tree and for each trait using the

445    "MPSEM" package v. 0.3.6 in R (47). To prevent overfitting, the number of eigenvectors that

446    explained greater than or equal to 65% of the phylogenetic structure variance was used (see S1

447    File for further details on this process). Following the method of Penone *et al.* (11), the

448    phylogenetic eigenvectors were appended to the complete-case dataset and treated as predictors

449    in the model to impute the missing value of a given trait.

450    Previous studies have suggested that phylogenetic signal strength in simulated trait data is

451    positively correlated with imputation accuracy (32,70). To assess this association using real data,

452    we measured phylogenetic signal for each trait using Pagel's λ (41) for numerical traits and the *D*

453    metric (40) for categorical traits. Pagel's λ is estimated using maximum likelihood and represents

454    the value that optimally transforms a phylogenetic variance-covariance matrix to fit the observed

455    trait data structure. A λ value of 0 indicates no phylogenetic signal (star-shaped phylogeny),

456    whereas a λ value of 1 suggests that the trait data adhere to a Brownian motion (BM) model of

457    evolution (41). The *D* metric represents whether the number of transitions of a binary trait varies

458    from the expected number under a BM model (40). A *D* value of 0 indicates that the trait data

459    adhere to a BM model, and a *D* value of 1 indicates that there is no phylogenetic signal in the

460    trait data. A *D* value greater than 1 signifies phylogenetic overdispersion. Alternatively, a *D*

461    value less than 0 suggests the trait is phylogenetically conserved (40). These metrics were

462    calculated separately for each trait using each gene tree (S1 File). The *phylosig* function in the R

463    package "phytools" v. 0.7.70 (51) and *phylo.d* function in the R package "caper" v. 1.0.1 (52)

464    were used to measure λ and *D*, respectively.

## Imputation process

466    Four imputation methods were considered: mean/mode imputation, *k*-nearest neighbour

467    (*KNN*) ("VIM" package v. 6.1.0; 16), random forests (*RF*) ("missForest" package v. 1.4; 53,54),

468    and multivariate imputations by chained equations (*MICE*) ("mice" package v. 3.13.0; 19). Mean

469    (for numerical traits) / mode (for categorical traits) imputation, the simplest method, was used as

470    a baseline for comparison. The remaining methods were chosen due to their popularity in trait-

471    based studies (e.g. 27,55) and capacity to impute both continuous and categorical traits. These

472    methods have also been evaluated in previous studies of trait data imputation (11,34,43). *KNN*

473    and *RF* are single imputation methods as they provide a single estimation of the missing value.

474    *MICE* is a multiple imputation method that performs imputation *m* times on the dataset with

475    missing values, resulting in *m* imputed datasets. The *MICE* algorithm utilizes chained equations

476    to estimate missing values and offers several different models for imputing data. In this study,

477    the predictive mean matching model was used to estimate missing continuous data. Predictive

478    mean matching is the default model for continuous data in *MICE* and performed well in previous

479    evaluations using trait data (43,75). Predictive mean matching fills the missing observation with

480    a random value selected from a "donor" pool for the missing observations. This pool is created

481    by fitting a regression model on the observed data and selecting *k* fitted values that are closest to

482    the predicted value for the missing observation (44,45). Logistic regression is a common

483    approach for predicting missing categorical data and is the default method for imputing

484    categorical data in *MICE*. Logistic regression and polytomous logistic regression models were

485    used to impute values for the binary trait insular endemic and the nominal multi-categorical trait

486    activity time, respectively. To obtain a final imputed value for *MICE*, the mean and mode values

487    were taken across the *m* datasets for numerical traits and categorical traits, respectively. See S1

488    File for further details on imputation algorithms.

489        When imputing the missing values of each trait ("target trait") using the observed values

490    of the other traits ("auxiliary traits"), not all of the auxiliary traits are useful for imputing the

491    missing values of the target trait. Association tests between each pair of traits were used to filter

492    out irrelevant auxiliary traits and build a more parsimonious imputation model for the target trait.

493    Regression models were used in the association tests in which the target trait was specified as the

494    response variable and each one of the auxiliary traits was specified as the covariate. Linear

495    regression, Poisson regression, and logistic regression models were used for continuous, count,

496    and categorical target traits, respectively. Only auxiliary traits with a coefficient not significantly

497    equal to zero were retained in the imputation model for a particular target trait. Finally, as

498    methods such as *KNN* are sensitive to the range of the data, numerical traits were natural log-

499    transformed prior to imputation.

## Simulation study

501          To simulate missing data, three different missingness scenarios were considered: 1)

502    missing completely at random ("MCAR"); 2) missing at random ("MAR"); and 3) missing not at

503    random ("MNAR"). Within the MCAR scenario, missing values were randomly introduced into

504    the complete-case dataset at different proportions (0.10, 0.20, 0.30, and 0.40). In cases where

505    traits had values that were already missing (up to 10%), missing values were introduced on top

506    of these (i.e. up to 50% missingness). To reduce stochasticity and maintain a fair comparison of

507    imputation performance across different missing proportions, and not introducing variability

508    relating to species identity, missing data for each increase in proportion (e.g. from 0.10 to 0.20

509    missingness) were added upon the missing values of the previous proportion. To simulate MAR

510    scenarios using real data, logistic regression models were fitted to the original Meiri (16) dataset

511    ($n = 6657$) to identify which auxiliary traits were significantly associated with the missingness

512    for each target trait. In the fitted model, the indicator of whether an observation is missing or not

513    was treated as the response variable and auxiliary traits specified as predictors. The fitted models

514    were then used to introduce missing values into the complete-case datasets (for further details see

515    S1 File). To test how the imputation methods perform in cases of extreme bias, MNAR scenarios

516    were simulated for each trait. Values were removed from the 10th percentile of the tail of data

517    distribution for numerical biological traits, e.g., the 10th percentile of the lower latitudes

518    (between 10° and -10°); and from a single category for categorical traits, e.g., "nocturnal"

519    category for activity time and "yes" category for insular endemic. These values were removed to

520    emulate realistic MNAR scenarios for Squamata (see S1 File for further information).

521    A range of parameters and their values were considered for the different imputation

522    methods (see S1 File for details on this process). The parameters that resulted in the lowest error

523    rate were used in the imputation model. Imputations using only trait data were first performed on

524    the simulated missing dataset. Imputations were again performed using trait data and

525    phylogenetic eigenvectors derived from either COI, RAG1, or c-mos gene trees. This amounted

526    to 78 different combination settings with respect to method and missingness scenario. The entire

527    process was repeated 100 times for each combination of settings, resulting in 7,800 runs of the

528    simulation and imputation pipeline procedure (see Fig 6 for a visualization of the process).

529    **Fig 6. Workflow of the pipeline for a particular combination of variables.** 1) 20% of the trait

530    observations are removed missing completely at random (MCAR) from the complete-case

531    dataset; 2) missing values are imputed using *k*-nearest neighbour (*KNN*). Phylogenetic

532    information in the form of a cytochrome *c* oxidase subunit I (COI) gene tree and known trait data

533    are used to estimate the missing trait data; and 3) the imputed values are compared to those in the

534    complete-case dataset. Mean squared error (MSE) or proportion falsely classified (PFC) are

535    calculated for numerical and categorical traits, respectively, and averaged across 100 replicates.

## Evaluation of methods

536

537    To assess imputation accuracy, imputed values were compared against the known values

538    in the complete-case dataset. Mean squared error (MSE) rates and proportion falsely classified

539    (PFC) rates were computed for numerical and categorical traits, respectively. These rates were

540    averaged across the 100 replicates for each combination of methods for each trait. For both

541     metrics, values closer to 0 are indicative of better performance. The packages "ggplot2" v. 3.3.5

542     (42) and "plotly" v. 4.10.0 (76) were used to visualize results in R.

## Real data imputation application and comparison

544     To select the most suitable method for imputing missing values in the original trait

545     dataset, the results of the MAR simulations were first considered as these mimic realistic

546     biological scenarios. In case of more than one method performing equally well, the method that

547     was most robust across different missingness scenarios and that resulted in the lowest average

548     error rate for the majority of traits was selected. To investigate whether imputed values alter the

549     quantitative distributional characteristics of the data, summary statistics for each trait were

550     calculated using the dataset that includes imputed values and compared with the corresponding

551     summary statistics of both the original and complete-case datasets. To investigate whether the

552     phylogenetic information improves the imputation accuracy for a given trait and imputation

553     method, the following error ratio was calculated for each trait and each method:

554
$$Error\ ratio = \frac{Error\ rate\ (MSE\ or\ PFC)\ without\ phylogeny}{Error\ rate\ (MSE\ or\ PFC)\ with\ phylogeny}$$

555     An error ratio value greater than 1 indicates an improvement in imputation performance resulting

556     from the addition of phylogenetic information. To observe the trend of the effect of phylogenetic

557     signal strength on the imputation of different traits, the error ratio values were plotted against the

558     $\lambda$ and $D$ metrics for numerical and categorical traits, respectively.

# Acknowledgements

We would like to thank Dr. Cameron Nugent and Dr. Karl Cottenie for their helpful insights

regarding the design and structure of the simulation pipeline. We also thank Dr. Tyler Elliott for

his valuable comments on the manuscript and code. Finally, we thank many researchers who

have collected trait and sequence data and made them publicly available. This work would not

have been possible without you.

# References

1.  Voituron Y, de Fraipont M, Issartel J, Guillaume O, Clobert J. Extreme lifespan of the human fish (*Proteus anguinus*): a challenge for ageing mechanisms. Biol Lett. 2011;7(1):105–7.

2.  Valcu M, Dale J, Griesser M, Nakagawa S, Kempenaers B. Global gradients of avian longevity support the classic evolutionary theory of ageing. Ecography. 2014 Oct 1;37(10):930–8.

3.  Howard SD, Bickford DP. Amphibians over the edge: silent extinction risk of Data Deficient species. Divers Distrib. 2014 Jul 1;20(7):837–46.

4.  Pacifici M, Visconti P, Butchart SHM, Watson JEM, Cassola FM, Rondinini C. Species' traits influenced their response to recent climate change. Nat Clim Change. 2017 Mar 1;7(3):205–8.

5.  Garamszegi LZ, Møller AP. Nonrandom variation in within-species sample size and missing data in phylogenetic comparative studies. Syst Biol. 2011;60(6):876–80.

6.  González-Suárez M, Lucas PM, Revilla E. Biases in comparative analyses of extinction risk: mind the gap. J Anim Ecol. 2012 Nov 1;81(6):1211–22.

7.  Nakagawa S, Freckleton RP. Missing inaction: the dangers of ignoring missing data. Trends Ecol Evol. 2008;23:592–6.

8.  Titley MA, Snaddon JL, Turner EC. Scientific research on animal biodiversity is systematically biased towards vertebrates and temperate regions. PLOS ONE. 2017 Dec 14;12(12):e0189577.

9.  Etard A, Morrill S, Newbold T. Global gaps in trait data for terrestrial vertebrates. Glob Ecol Biogeogr. 2020;29(12):2143–58.

10. Blomberg SP, Garland TJr, Ives AR. Testing for phylogenetic signal in comparative data: Behavioral traits are more labile. Evolution. 2003;57:717–45.

590   11.   Penone C, Davidson AD, Shoemaker KT, Di Marco M, Rondinini C, Brooks TM, et al.
591         Imputation of missing data in life-history trait datasets: which approach performs the best?
592         Methods Ecol Evol. 2014;5:961–70.

593   12.   Myhrvold NP, Baldridge E, Chan B, Sivam D, Freeman DL, Ernest SKM. An amniote life-
594         history database to perform comparative analyses with birds, mammals, and reptiles.
595         Ecology. 2015;96(11):3109.

596   13.   Nathan P. Myhrvold, Elita Baldridge, Benjamin Chan, Dhileep Sivam, Daniel L. Freeman,
597         S. K. Morgan Ernest. Data from: An amniote life-history database to perform comparative
598         analyses with birds, mammals, and reptiles [Internet]. Wiley. Collection.; 2016. Available
599         from: https://wiley.figshare.com/articles/dataset/Full_Archive/3563457

600   14.   Tamburello N, Côté IM, Dulvy NK. Energy and the Scaling of Animal Space Use. Am Nat.
601         2015 Aug 1;186(2):196–211.

602   15.   Tamburello N, Côté IM, Dulvy NK. Data from: Energy and the Scaling of Animal Space
603         Use. Dryad Dataset. 2015;

604   16.   Meiri S. Traits of lizards of the world: Variation around a successful evolutionary design.
605         Glob Ecol Biogeogr. 2018;27(10):1168–72.

606   17.   Meiri S. Data from: Traits of lizards of the world: Variation around a successful
607         evolutionary design. Dryad Dataset [Internet]. 2019; Available from:
608         https://doi.org/10.5061/dryad.f6t39kj

609   18.   de Magalhães JP, Costa J. A database of vertebrate longevity records and their relation to
610         other life-history traits. J Evol Biol. 2009;22:1770–4.

611   19.   Tacutu R, Thornton D, Johnson E, Budovsky A, Barardo D, Craig T, et al. Human Ageing
612         Genomic Resources: new and updated databases. Nucleic Acids Res. 2018 Jan
613         4;46(D1):D1083–90.

614   20.   Zhang Z. Missing data imputation: focusing on single imputation. Ann Transl Med.
615         2015;4(1):9.

616   21.   Rubin DB. Inference and missing data. Biometrika. 1976;63(3):581–92.

617   22.   Donders ART, van der Heijden GJMG, Stijnen T, Moons KGM. Review: A gentle
618         introduction to imputation of missing values. J Clin Epidemiol. 2006;59:1087–91.

619   23.   van Buuren S. Flexible Imputation of Missing Data. Boca Raton, FL: CRC Press, Taylor &
620         Francis Group; 2012.

621   24.   Kowarik A, Templ M. Imputation with the R Package VIM. J Stat Softw. 2016;74(7):1–16.

622   25.   Stekhoven DJ, Bühlmann P. MissForest—non-parametric missing value imputation for
623         mixed-type data. Bioinformatics. 2012 Jan 1;28(1):112–8.

624    26.   Schafer JL. Multiple imputation: a primer. Stat Methods Med Res. 1999;8:3–15.

625    27.   Van Buuren S, Groothuis-Oudshoorn K. MICE: Multivariate Imputation by Chained
626          Equations in R. J Stat Softw. 2011;45(3):1–67.

627    28.   Swenson NG. Phylogenetic imputation of plant functional trait databases. Ecography.
628          2014;37:105–10.

629    29.   Diniz-Filho JAF, Ramos de Sant'ana CE, Bini LM. An eigenvector method for estimating
630          phylogenetic inertia. Evolution. 1998;52:1247–62.

631    30.   Guénard G, Legendre P, Peres-Neto P. Phylogenetic eigenvector maps: a framework to
632          model and predict species traits. Methods Ecol Evol. 2013 Dec 1;4(12):1120–31.

633    31.   Guénard G. A phylogenetic modelling tutorial using Phylogenetic Eigenvector Maps
634          (PEM) as implemented in R package MPSEM (0.3-6). 2019.

635    32.   Kim SW, Blomberg SP, Pandolfi JM. Transcending data gaps: a framework to reduce
636          inferential errors in ecological analyses. Ecol Lett. 2018;21(8):1200–10.

637    33.   Fournier A, Penone C, Pennino MG, Courchamp F. Predicting future invaders and future
638          invasions. Proc Natl Acad Sci U S A. 2019/03/29 ed. 2019 Apr 16;116(16):7905–10.

639    34.   Johnson TF, Isaac NJB, Paviolo A, González-Suárez M. Handling missing values in trait
640          data. Glob Ecol Biogeogr. 2021 Jan 1;30(1):51–62.

641    35.   James TD, Salguero-Gómez R, Jones OR, Childs DZ, Beckerman AP. Bridging gaps in
642          demographic analysis with phylogenetic imputation. Conserv Biol. 2021;35(4):1210–21.

643    36.   Keck BP, Near TJ. Assessing phylogenetic resolution among mitochondrial, nuclear, and
644          morphological datasets in Nothonotus darters (Teleostei: Percidae). Mol Phylogenet Evol.
645          2008 Feb 1;46(2):708–20.

646    37.   Blom MPK, Bragg JG, Potter S, Moritz C. Accounting for Uncertainty in Gene Tree
647          Estimation: Summary-Coalescent Species Tree Inference in a Challenging Radiation of
648          Australian Lizards. Syst Biol. 2017 May 1;66(3):352–66.

649    38.   Templ M, Kowarik A, Alfons A, de Cillia G, Prantner B, Rannetbauer W. R package
650          "VIM": Visualization and Imputation of Missing Values [Internet]. 2021. Available from:
651          https://cran.r-project.org/web/packages/VIM/VIM.pdf

652    39.   Stekhoven DJ. missForest: Nonparametric Missing Value Imputation using Random Forest.
653          2013.

654    40.   Fritz SA, Purvis A. Selectivity in Mammalian Extinction Risk and Threat Types: a New
655          Measure of Phylogenetic Signal Strength in Binary Traits. Conserv Biol. 2010 Aug
656          1;24(4):1042–51.

657    41.  Pagel M. Inferring the historical patterns of biological evolution. Nature. 1999 Oct
658         1;401(6756):877–84.

659    42.  Wickham H. ggplot2: Elegant Graphics for Data Analysis. [Internet]. New York: Springer-
660         Verlag; 2016. Available from: https://ggplot2.tidyverse.org

661    43.  Poyatos R, Sus O, Badiella L, Mencuccini M, Martínez-Vilalta J. Gap-filling a spatially
662         explicit plant trait database: comparing imputation methods and different levels of
663         environmental information. Biogeosciences. 2018;15:2601–17.

664    44.  Morris TP, White IR, Royston P. Tuning multiple imputation by predictive mean matching
665         and local residual draws. BMC Med Res Methodol. 2014 Jun 5;14(1):75.

666    45.  Kleinke K. Multiple Imputation by Predictive Mean Matching When Sample Size Is
667                                  Small. Methodology. 2018 Jan 1;14(1):3–15.

668    46.  Diniz-Filho JAF, Bini LM, Rangel TF, Morales-Castilla I, Olalla-Tárraga MÁ, Rodríguez
669         MÁ, et al. On the selection of phylogenetic eigenvectors for ecological analyses.
670         Ecography. 2012 Mar 1;35(3):239–49.

671    47.  Springer MS, DeBry RW, Douady C, Amrine HM, Madsen O, de Jong WW, et al.
672         Mitochondrial Versus Nuclear Gene Sequences in Deep-Level Mammalian Phylogeny
673         Reconstruction. Mol Biol Evol. 2001;18(2):132–43.

674    48.  Madley-Dowd P, Hughes R, Tilling K, Heron J. The proportion of missing data should not
675         be used to guide decisions on multiple imputation. J Clin Epidemiol. 2019;110:63–73.

676    49.  Enders CK. Applied Missing Data Analysis. New York: The Guilford Press; 2010.
677         (Methology in the Social Sciences).

678    50.  Jardim L, Bini LM, Diniz-Filho JAF, Villalobos F. A Cautionary Note on Phylogenetic
679         Signal Estimation from Imputed Databases. Evol Biol. 2021 Jun 1;48(2):246–58.

680    51.  Esquerré D, Brennan IG, Catullo RA, Torres-Pérez F, Keogh JS. How mountains shape
681         biodiversity: The role of the Andes in biogeography, diversification, and reproductive
682         biology in South America's most species-rich lizard radiation (Squamata: Liolaemidae).
683         Evolution. 2019;73(2):214–30.

684    52.  Skeels A, Cardillo M. Reconstructing the Geography of Speciation from Contemporary
685         Biodiversity Data. Am Nat. 2019 Feb 1;193(2):240–55.

686    53.  Uetz P, Aguilar P, Hošek J, editors. The Reptile Database. 2021; Available from:
687         http://www.reptile-database.org

688    54.  Cox N, Young BE, Bowles P, Fernandez M, Marin J, Rapacciuolo G, et al. A global reptile
689         assessment highlights shared conservation needs of tetrapods. Nature [Internet]. 2022 Apr
690         27; Available from: https://doi.org/10.1038/s41586-022-04664-7

691  55.  Böhm M, Williams R, Bramhall HR, McMillan KM, Davidson AD, Garcia A, et al.
692        Correlates of extinction risk in squamate reptiles: the relative importance of biology,
693        geography, threat and range size. Glob Ecol Biogeogr. 2016;25(4):391–405.

694  56.  Munstermann MJ, Heim NA, McCauley DJ, Payne JL, Upham NS, Wang SC, et al. A
695        global ecological signal of extinction risk in terrestrial vertebrates. Conserv Biol [Internet].
696        n/a(n/a). Available from: https://doi.org/10.1111/cobi.13852

697  57.  R Core Team. R: A language and environment for statistical computing. [Internet]. Vienna,
698        Austria: R Foundation for Statistical Computing; 2020. Available from: https://www.R-
699        project.org/

700  58.  Ratnasingham S, Hebert PDN. bold: The Barcode of Life Data System
701        (http://www.barcodinglife.org). Mol Ecol Notes. 2007 May 1;7(3):355–64.

702  59.  Data from: Barcode of Life Data System: DS-IMPMIX2: Squamata cytochrome c oxidase
703        subunit I (COI) dataset. [Internet]. 2020. Available from: dx.doi.org/10.5883/DS-IMPMIX2

704  60.  Orton MG, May JA, Ly W, Lee DJ, Adamowicz SJ. Is molecular evolution faster in the
705        tropics? Heredity. 2019 May 1;122(5):513–24.

706  61.  Wright ES. DECIPHER: harnessing local sequence context to improve protein multiple
707        sequence alignment. BMC Bioinformatics. 2015 Oct 6;16(1):322.

708  62.  Wright ES. RNAconTest: comparing tools for noncoding RNA multiple sequence
709        alignment based on structural consistency. RNA. 2020 May 1;26(5):531–40.

710  63.  Yu G, Zhou L, Huang H. Package "ggmsa". Plot Multiple Sequence Alignment using
711        'ggplot2. 2021.

712  64.  Pyron RA, Burbrink FT, Wiens JJ. A phylogeny and revised classification of Squamata,
713        including 4161 species of lizards and snakes. BMC Evol Biol. 2013 Apr 29;13(1):93.

714  65.  Pyron RA, Burbrink FT, Wiens JJ. Data from: A phylogeny and revised classification of
715        Squamata, including 4161 species of lizards and snakes. Dryad Dataset [Internet]. 2013;
716        Available from: https://doi.org/10.5061/dryad.82h0m

717  66.  Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of
718        large phylogenies. Bioinformatics. 2014;30:1312–3.

719  67.  Roll U, Feldman A, Novosolov M, Allison A, Bauer AM, Bernard R, et al. The global
720        distribution of tetrapods reveals a need for targeted reptile conservation. Nat Ecol Evol.
721        2017 Nov 1;1(11):1677–82.

722  68.  Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and evolutionary
723        analyses in R. Bioinformatics. 2019;35:526–8.

724  69.  Guénard G, Legendre P. Modeling Phylogenetic Signals using Eigenvector Maps. 2019.

725  70.  Molina-Venegas R, Moreno-Saiz JC, Castro Parga I, Davies TJ, Peres-Neto PR, Rodríguez
726       MÁ. Assessing among-lineage variability in phylogenetic imputation of functional trait
727       datasets. Ecography. 2018 Oct 1;41(10):1740–9.

728  71.  Revell L. phytools: An R package for phylogenetic comparative biology (and other things).
729       Methods Ecol Evol. 2012;3:217–23.

730  72.  Orme D, Freckleton RP, Thomas G, Petzoldt T, Fritz S, Isaac N, et al. Package "caper":
731       Comparative Analyses of Phylogenetics and Evolution in R. 2018.

732  73.  Stekhoven DJ, Buehlmann P. MissForest - non-parametric missing value imputation for
733       mixed-type data. Bioinformatics. 2012;28(1):112–8.

734  74.  Richards C, Cooke RSC, Bates AE. Biological traits of seabirds predict extinction risk and
735       vulnerability to anthropogenic threats. Glob Ecol Biogeogr. 2021 May 1;30(5):973–86.

736  75.  Taugourdeau S, Villerd J, Plantureux S, Huguenin-Elie O, Amiaud B. Filling the gap in
737       functional trait databases: use of ecological hypotheses to replace missing data. Ecol Evol.
738       2014;4(7):944–58.

739  76.  Sievert C. Interactive Web-Based Data Visualization with R, plotly, and shiny. [Internet].
740       Florida: Chapman and Hall/CRC; 2020. Available from: https://plotly-r.com

741

# Supporting information

743

744  **S1 File. Supplementary Information.**

745  **S1 Fig. Phylogenetic signal measurements.** Measures of phylogenetic signal for a) categorical

746  and b) numerical traits in gene trees constructed for mitochondrial COI and nuclear c-mos and

747  RAG1. Asterisks indicate significance at the 0.05 level, according to results from hypothesis

748  tests comparing the results to a null model (no phylogenetic signal). Fritz and Purvis' *D* metric

749  (40) and Pagel's λ (41) were used to measure phylogenetic signal for categorical and numerical

750  traits, respectively. In the case of *D,* lower values are indicative of higher levels of phylogenetic

751  conservation for the trait; conversely, higher values of λ suggest stronger phylogenetic signal. As

752  the *D* metric only measures the phylogenetic signal of binary traits, the three-level categorical

753  trait AT was broken down into the binary traits "AT: Diurnal" and "AT: Nocturnal".

754    **S1 Table. Sequence identifiers.**

755    **S2 Table. Taxonomic composition of complete-case trait dataset (n = 152).**

756    **S3 Table. Descriptions and additional details for traits in the complete-case dataset.**

Fig 1

**Method**
- *Mode*
- *KNN*
- *RF*
- *MICE_LR*

a) Activity time

b) Insular endemic

Error Rate (PFC)

Proportion of missingness

c) Largest clutch    d) Smallest clutch

**Method**
- *Mean*
- *KNN*
- *RF*
- *MICE_PMM*

Error Rate (MSE)

Proportion of missingness

e) Female SVL

f) Maximum SVL

**Method**
- *Mean*
- *KNN*
- *RF*
- *MICE_PMM*

Error Rate (MSE)

Proportion of missingness

g) Latitude

| | MCAR 10% | MCAR 40% | MAR | MNAR |

**MCAR 10%**     **MCAR 40%**     **MAR**     **MNAR**

c) Largest clutch

Legend
- Mean
- KNN
- RF
- MICE

d) Smallest clutch

Error Rate (MSE)

Method

**MCAR 10%**  **MCAR 40%**  **MAR**  **MNAR**

e) Female SVL

Legend
- Mean
- KNN
- RF
- MICE

f) Maximum SVL

NA

Error Rate (MSE)

Method

g) Latitude

a) Categorical

0.10 MCAR

0.40 MCAR

MAR

MNAR

NA

**Method**
- KNN (blue triangle)
- MICE_LR (pink square)
- RF (green circle)

Better performance ↑

Error ratio (PFC without phylogeny/ PFC with phylogeny)

Phylogenetic signal (Fritz and Purvis' *D* metric)

← Increasing phylogenetic signal

b) Numerical

**Method**
- KNN (blue triangle)
- MICE_PMM (pink square)
- RF (green circle)

Better performance

Error ratio (MSE without phylogeny/ MSE with phylogeny)

0.10 MCAR

0.40 MCAR

MAR

MNAR

Phylogenetic signal (Pagel's λ)

Increasing phylogenetic signal

**a) Activity time**

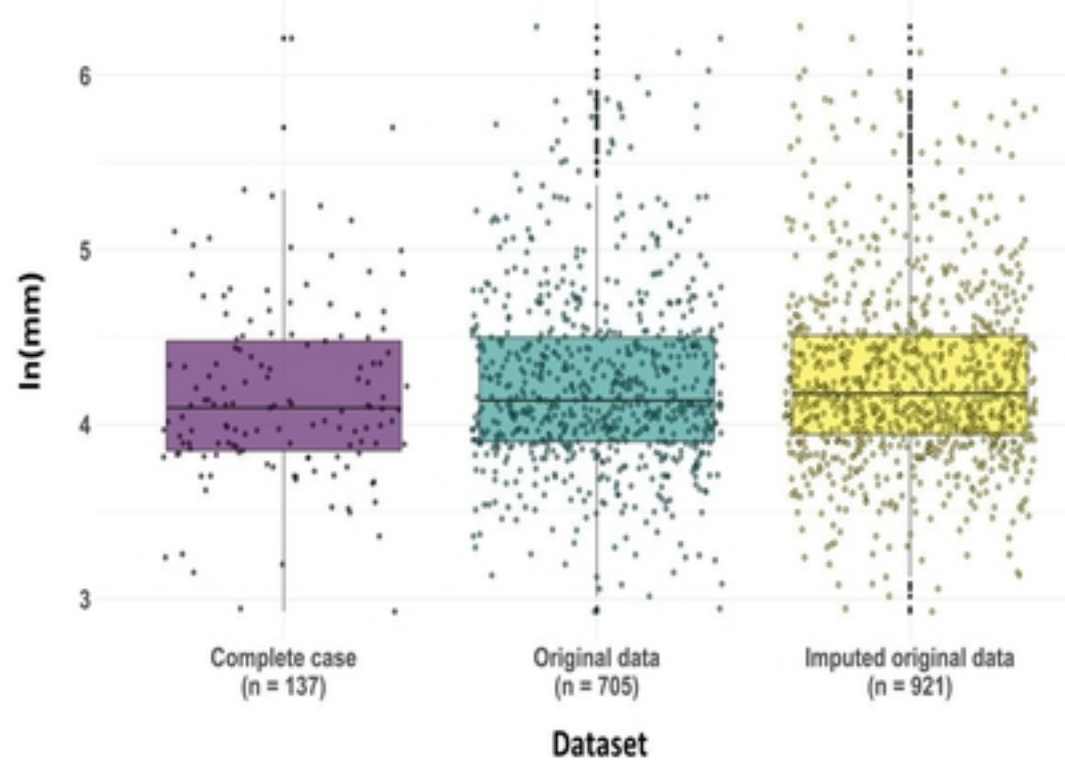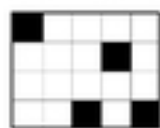**b) Largest clutch**

**c) Smallest clutch**

**d) Female SVL**

Fig 5

**1) Simulate missingness**
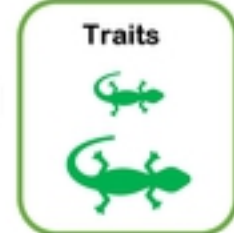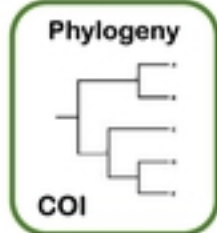
Complete-case dataset

→

Missing dataset

**2) Imputation**

Impute | Missing dataset | using | *KNN*

Phylogeny / COI → ← Traits

Imputed dataset

**3) Precision evaluations**

Imputed dataset | Vs. | Complete-case dataset

Mean squared error (MSE) | **or** | Proportion falsely classified (PFC)

**Result is average MSE or PFC across 100 replicates for _each trait_**

Fig 6