

# 1 **A phylogenetic approach to comparative genomics**

2 Anna E. Dewar<sup>1,2,†</sup>, Laurence J. Belcher<sup>1</sup> & Stuart A. West<sup>1</sup>

3 <sup>1</sup>Department of Biology, University of Oxford, Oxford, OX1 3SZ, United Kingdom

4 <sup>2</sup>St John's College, Oxford, OX1 3JP, United Kingdom

5 These authors contributed equally: Anna E. Dewar, Laurence J. Belcher

6 †Corresponding author e-mail: [anna.dewar@biology.ox.ac.uk](mailto:anna.dewar@biology.ox.ac.uk)

7

## 8 **Abstract**

9 Insights from comparative genomics, whereby the genomes of different species are compared,  
10 have the potential to address broad and fundamental questions at the intersection of genetics  
11 and evolution. However, species, genomes and genes cannot be considered as independent data  
12 points within statistical tests. Closely related species tend to be similar because they share  
13 genes by common descent, which must be accounted for in analyses. This problem of non-  
14 independence can be exacerbated when examining genomes or genes. The application of  
15 phylogeny-based methods to comparative genomics can address this problem. These methods  
16 must be considered an essential part of the comparative genomics toolkit, to prevent the  
17 accumulation of studies that produce incorrect conclusions. Here, we review how controlling  
18 for phylogeny can change the conclusions of comparative genomics studies. We address  
19 common questions about how to apply these methods and illustrate how they can be used to  
20 test causal hypotheses. The combination of rapidly expanding genomic datasets and  
21 phylogenetic comparative methods is set to revolutionize our understanding of biology.

22

## 23 **[H1] Introduction**

24 Since the beginning of the century, the number of sequenced genomes has increased to many  
25 hundreds of thousands, providing a full genetic record for thousands of species across the tree  
26 of life<sup>1-5</sup>. Comparing the genomes of different species can provide insight into which genes  
27 control particular phenotypic traits and how those traits evolved. For example, comparative  
28 genomics has identified bacterial genes that harm crops<sup>6,7</sup>, tracked antibiotic resistance through  
29 hospitals<sup>8,9</sup>, and unravelled how echolocation evolved independently in dolphins and bats<sup>10-12</sup>.

30 Comparative genomics datasets often contain biases that, unless fully accounted for,  
31 can lead to incorrect conclusions. Statistical tests assume that data points are 'independent'  
32 from one another. This means that each data point represents an individual replicate drawn  
33 from an underlying distribution, and the value of any one data point does not depend on or

34 influence the value of any other data point. Even a slight lack of independence can lead to bias  
35 in statistical analyses and incorrect conclusions<sup>13–15</sup>.

36 A problem for biologists wanting to conduct any form of across-species comparative  
37 analyses is that species do not represent independent data points<sup>16–20</sup> (Fig. 1a). Closely related  
38 species tend to be similar because they share traits by common descent rather than through  
39 independent evolution (Fig. 1b). To give an extreme example, passerine birds all have wings  
40 because they inherited the genes to produce them from a common ancestor, not because each  
41 species independently evolved to fly. In genomic datasets, data points can correspond to  
42 genomes or even individual genes, meaning the problem of species not being independent can  
43 be exacerbated by additional levels of non-independence, including genes within genomes,  
44 such as those carried on the same replicon, and genomes of the same or closely related  
45 species<sup>16,19,21–26</sup>. The problem of non-independence is complicated even further by how  
46 unrepresentative most genomic datasets are of natural communities, owing to biases in which  
47 species have been sequenced more frequently (Fig. 2)<sup>27–29</sup>.

48 The fact that species cannot be considered independent data points was first recognized  
49 in the 1970s by evolutionary biologists studying the evolution of life history traits such as  
50 sexual dimorphism and home range in primates<sup>16,18,30–32</sup>. This led to the development of  
51 numerous methods that take phylogeny into account and control for the non-independence of  
52 species<sup>16,18,19,33</sup>. These phylogeny-based comparative methods have revolutionized behavioural  
53 and evolutionary ecology research, including research into sexual selection, resource  
54 competition and social evolution<sup>34</sup>. However, the application of phylogeny-based comparative  
55 methods to genomic data has not been consistent, and has been limited in some areas such as  
56 microbial genomics. The potential problem of non-independence remains to be fully  
57 recognized across the entire field of comparative genomics. In contrast, many of those  
58 developing and using phylogeny-based methods focus on phenotypic, rather than genomic,  
59 traits. Consequently, some sections of these fields have developed separately, with little  
60 discourse, although the problems and solutions are very similar. It is crucial that these existing  
61 phylogenetic methods are now consistently applied to genomic data to prevent the  
62 accumulation of studies with spurious conclusions.

63 Here, we review how phylogeny-based comparative methods can be used to analyse  
64 genomic data to resolve two key problems: species are not independent, and using genes and  
65 genomes as data points can multiply the problem of non-independence. We provide illustrative  
66 case studies in which controlling for non-independence led to changes in the biological  
67 conclusions drawn from genomic data (Box 1). We address common questions that arise

68 regarding when and how to control for phylogeny, and show how phylogeny-based methods  
69 can be used to test causal hypotheses. We focus on comparative genomic analyses across  
70 species, although comparative genomics can also be conducted within species; the questions  
71 and methods used are generally analogous, and problems encountered are similar to those we  
72 discuss here<sup>35–40</sup>. Of note, phylogenetic methods can be used to also address other questions,  
73 such as how features of the genome influence diversification<sup>41–43</sup>.

74

## 75 **[H1] Control for phylogenetic relationships**

76 The problem of species' non-independence can be solved by taking the evolutionary history of  
77 species into account. A phylogeny is a tree that shows the evolutionary relationships between  
78 species, and many methods allow us to incorporate phylogenies into statistical models and  
79 tests. While each approach is different, they all ultimately account and control for phylogenetic  
80 non-independence. Here, we summarize the most used methods (Table 1).

81

## 82 **[H2] Phylogenetic regressions and mixed models**

83 One class of methods to control for phylogeny uses phylogenetic regressions<sup>19,44,45</sup>.  
84 These are like standard regression approaches but with an additional term to account for a  
85 trait's correlation with species' evolutionary history ('phylogenetic signal'). Similar to other  
86 linear models, phylogenetic regressions can be generalized to allow for non-Gaussian response  
87 variables. These phylogenetic regressions can also be expanded to phylogenetic mixed models,  
88 to analyse the influence of multiple variables. As with standard mixed models, phylogenetic  
89 mixed models allow us to include all explanatory variables as either: a fixed effect, meaning  
90 one would like to test the effect on the response variable; or a random effect, meaning one  
91 would like to control for any effect on the response but not explicitly test it in the analysis. A  
92 phylogenetic mixed model includes the phylogeny as a random effect, often as a matrix  
93 corresponding to the structure of the phylogeny. This approach can examine whether any of  
94 the potential explanatory variables has a significant effect on the response variable, while  
95 controlling for the phylogenetic signal. In terms of controlling for phylogeny, using a  
96 phylogenetic regression or including the phylogeny as a random effect within a mixed model  
97 are largely analogous. Which method is most appropriate will depend on factors such as the  
98 type and distribution of the data and the number of variables of interest, and it is often useful  
99 to analyse the data multiple ways to confirm that results are robust.

100           There are many methods to run phylogenetic regressions and mixed models. Classic  
101 phylogenetic regressions can be run in R packages such as ‘phylolm’ and ‘ape’. More recently,  
102 Bayesian approaches have become popular, using R packages such as ‘MCMCglmm’ and  
103 ‘brms’, and the computer package ‘BayesTraits’ (Table 1)<sup>42,46–52</sup>. These different approaches  
104 all share the same underlying rationale: to account for phylogeny within statistical analyses.  
105 Which approach to use depends on the research question and the nature of the data<sup>42</sup>. In  
106 addition, as with any regression, the impact of any deviation from homogeneous variance  
107 should be considered<sup>53</sup>.

108

## 109 **[H2] Rapidly evolving traits**

110           The problem of species non-independence still arises when examining evolutionary  
111 labile traits that evolve very fast or are the product of very recent events. For example,  
112 mutation–selection balance, polymorphism at highly variable alleles, or the recent acquisition  
113 of genes by horizontal gene transfer. A total lack of phylogenetic signal is at one end of a  
114 spectrum and not a reason not to avoid phylogenetic comparative methods<sup>44,54–56</sup>. Even if a  
115 response variable shows low (or no) phylogenetic signal, an explanatory variable, known or  
116 unknown, could still be correlated with phylogeny. Phylogenetic comparative analyses provide  
117 a method to control for spurious results that could arise due to even unknown variables. It is  
118 rare that one can be sure that all biologically important variables are being considered, or that  
119 the unknown unknown variables are completely independent with respect to phylogeny.  
120 Bergeron *et al.* analysed how the fast-evolving trait of germline mutation rate varied across  
121 vertebrate species, while controlling for phylogenetic non-independence<sup>57</sup>. The key point here  
122 is that controlling for phylogeny is not simply controlling for phylogenetic signal or inertia  
123 (constraint) — it is about the philosophy of hypothesis testing and whether data points are  
124 independent.

125           Closely related species share many features of their ecology and life history.  
126 Consequently, there is an almost unlimited number of unknown ‘third’ variables which might  
127 influence the trait of interest<sup>20</sup>. For example, closely related species could show similar  
128 mutation–selection balance or genetic variation at certain loci because of selection imposed by  
129 a similar feature of their ecology or some other trait<sup>58–60</sup>. It is unrealistic to assume that every  
130 possible confounding factor can be identified and controlled for, but phylogenetic methods  
131 offer a way to control for unknowns. Similar problems, and a need for phylogenetic methods,  
132 arise when using data from other omic methods, such as transcriptomics, metagenomics or  
133 metabolomics<sup>61–64</sup>.

134

## 135 **[H1] Control for genomes and genes in a phylogenetic context**

136 Analysing genes or genomes as if they were independent data points can be equivalent to  
137 having more technical (pseudo) replicates in an experiment, but then treating each as if they  
138 were a truly independent sample (biological replicate)<sup>22</sup>. For example, if examining whether a  
139 gene's GC content is correlated with its length, we could analyse all 41 million genes across  
140 all 1,800 eukaryote genomes in the RefSeq database. However, multiple genomes from the  
141 same species are not independent – in the extreme, they may even have arisen via duplication.  
142 In addition, the number of genes per genome varies from 465 to 155,000 genes, which means  
143 that considering each gene as an independent data point would bias the data towards species  
144 with larger genomes.

145 This problem is made worse because genomic databases are biased towards taxonomic  
146 groups that have been sequenced more frequently<sup>27–29</sup> (Fig. 2). For example, of 28,504  
147 complete bacterial genomes in both the RefSeq and Genome Taxonomy Databases, 26%  
148 (7,376) are from 0.1% of species (Supplementary Table 1). By contrast, 73% of bacterial  
149 species have only one representative genome in the databases. Across all genomes in the newly  
150 curated AllTheBacteria database, the 10 most sequenced species make up 77% of the 1.9  
151 million genomes<sup>5</sup>. More broadly, there is also a bias towards bacterial genomes relative to  
152 eukaryotes<sup>3</sup>. If all genomes were used as independent data points, this would increase the bias  
153 towards species that have been sequenced the most.

154 As with non-independence due to the phylogenetic history of species, non-  
155 independence of genes and genomes should be controlled for in statistical analyses to avoid  
156 incorrect conclusions. One possibility is to retain all the data in a mixed model — individual  
157 genomes could be included as data points, with species' phylogeny and also species' sample  
158 size as random effects<sup>22,23,49,65–67</sup>. Where possible and required for the question being asked,  
159 non-independence within the genome due to duplication should also be controlled for, such as  
160 by grouping genes into homologous groups.

161 Another simple solution is to calculate average values for all genes in a genome and/or  
162 all genomes in a species before controlling for species phylogeny<sup>22,26,65,68</sup>. This would control  
163 for genome and/or species-level phylogenetic non-independence without the additional  
164 confounding effect of sample size differences. Data could be weighted according to the number  
165 of genomes per species, because the average for species with more samples would have a lower  
166 error variance<sup>69</sup>. The 'average value' for a species could alternatively be some other statistic

167 calculated from the genomes of that species, such as a correlation coefficient or an effect  
168 size<sup>65,67</sup>. For our GC content and gene length example given above, one could calculate a  
169 correlation coefficient between GC content and gene length for each genome or species.

170 In terms of sampling bias, phylogenetic methods will control for when taxa are  
171 overrepresented, so that better sampled taxa will not have a disproportionately large influence  
172 on analyses. It is, however, only possible to control for biases in representation of taxa currently  
173 present in a dataset. If taxa are missing, then their underrepresentation cannot of course be  
174 accounted for. Current and future efforts to sequence more taxa will be incredibly useful to  
175 tackle broad across-species questions<sup>3,4,70–74</sup>.

176

## 177 **[H1] Testing causal hypotheses**

178 Comparative genomic analyses usually test for correlations between traits. For example, recent  
179 genomic studies have examined correlations between factors such as pathogenicity and genome  
180 size, GC content and growth temperature, or the presence of certain genes and ecology<sup>37,67,75–</sup>  
181 <sup>83</sup>. However, these results are open to multiple causal explanations. For example, a correlation  
182 between pathogenicity and smaller genome size could arise because pathogenicity favours the  
183 evolution of a reduced genome size or because a smaller genome size facilitates the evolution  
184 of pathogenicity.

185 More recently developed methods allow analyses to go further and test competing  
186 causal hypotheses<sup>42,84</sup>. One approach is to use the characteristics present in extant species,  
187 together with a phylogeny, to reconstruct the characteristics of their ancestors ('ancestral  
188 states'). These ancestral states can be used to examine whether the evolution of certain  
189 characteristics is correlated. For example, do two characteristics tend to vary together along  
190 the branches of the phylogeny? If yes, this could be because they are both independently  
191 correlated with phylogenetic history or another third trait, or at least one characteristic directly  
192 influences the evolution of the other. To distinguish which scenario is more likely, we could  
193 test whether an ancestral change in one trait consistently leads to the same change in the other  
194 trait, which would suggest the co-evolution is due to a causal link<sup>42,43,51,85–87</sup> (Fig. 3).

195 What if there are more than two traits we are interested in? For example, there could be  
196 multiple traits that are correlated with one another. In this case, causal inference methods such  
197 as phylogenetic path analysis could be used to test causal hypotheses<sup>84,88,89</sup>. These methods are  
198 based on the fact that if correlation is not due to random chance, correlation must instead be  
199 due to some underlying causal relationship(s), and we can use information from correlations to

200 examine which relationship(s) is most likely<sup>84,88</sup>. A general point here is that correlational  
201 results from observational data do not prove causation, but correlations can be used within  
202 comparative analyses to test causal hypotheses and examine evidence for causality<sup>42,90</sup>.

203         Recent studies have provided examples of the insights that can be made with this  
204 approach. For example, species of snapping shrimp that are more social tend to have larger  
205 genomes with more transposable elements; a phylogenetic path analysis and ancestral state  
206 reconstruction supported the hypothesis that the acquisition of more transposable elements  
207 leads to both larger genome size and higher rates of sociality<sup>68</sup>. As another example, bacterial  
208 species with more fluid pangenomes have both more variable lifestyles and larger effective  
209 population sizes; a phylogenetic path analysis supported the hypothesis that this was due to  
210 lifestyle traits together influencing gene gain and loss, and that effective population size did  
211 not have a casual influence on pangenome fluidity<sup>24</sup>. There are numerous other possible  
212 applications for this approach.

213

## 214 **[H1] Common questions**

215 Several questions commonly arise about when and how to control for phylogeny in across-  
216 species comparative studies.

217

### 218 ***[H2] Do we always need to control for phylogeny?***

219 “Phylogenies are fundamental to comparative biology; there is no doing it without taking them  
220 into account.”<sup>18</sup> There is no simple and robust alternative to controlling for phylogeny, and  
221 consequently, the best effort possible must be made to take phylogeny into account, to be able  
222 to detect real effects and avoid false positives. There are almost always approaches which can  
223 reduce, even if not completely remove, the risk of incorrect conclusions.

224         If there is uncertainty in the phylogeny, one solution is to include multiple trees with  
225 alternative structures in the analysis, to estimate the effect of phylogenetic uncertainty on any  
226 result (that is, to determine for how many alternative structures the results are robust)<sup>91</sup>.  
227 Another solution is to use datasets with a narrower range of genomes, so that phylogenetic  
228 relationships can be reconstructed. For example, Frigols *et al.* focused on a particular group of  
229 viruses, phages that infect staphylococcal bacteria, to reconstruct and control for phylogenetic  
230 relationships<sup>92</sup>. Different methods can also be used depending upon the information  
231 available<sup>38,66,93,94</sup>. Murray *et al.* examined how genome size correlated with pathogenicity in  
232 bacteria by comparing ‘phylogenetically independent’ pairs of pathogenic and non-pathogenic

233 species within the same genera<sup>37</sup>. Similarly different methods can be used to investigate or  
234 account for factors such as different modes of evolution or changing evolutionary rates<sup>95-97</sup>.

235 New methods are required to examine phylogenetic relationships in cases where  
236 generating a meaningful phylogeny is hard or impossible, such as fast-evolving mobile genetic  
237 elements and viruses. For such cases, current approaches include using the phylogeny of the  
238 host genome, controlling for similarity using conserved regions such as plasmid relaxases, or  
239 similarity scores from network analyses<sup>22,23,65,77,98-100</sup>. There are pros and cons of these  
240 approaches. The phylogeny of the host genome can be useful for broad, across-species phyla  
241 datasets, where host range often acts as a meaningful barrier, but not between more closely  
242 related species. The extent to which similarity scores correspond to evolutionary history  
243 remains an open question. Results based on such methods will be tentative, and it can be useful  
244 to try different methods to test the robustness of conclusions. The usefulness of controlling for  
245 phylogenetic non-independence will also be dependent on the accuracy of the phylogeny,  
246 emphasising the advantage of new data and methods which allow more reliable phylogenies<sup>101</sup>.

247 The problem of non-independence is not solved by alternative methods such as  
248 grouping data by a higher taxonomic rank (for example, genera or family), or by sampling the  
249 same number of genomes from taxonomic groups to minimize oversampling  
250 ('downsampling')<sup>102,103</sup>. Grouping by an arbitrary taxonomic rank does not remove the problem  
251 of non-independence between the ranks<sup>19,104</sup>. Uniform sampling of genomes reduces bias in  
252 the dataset towards more sequenced species, but has no effect on the phylogenetic non-  
253 independence between those species. A related issue arises with methods that use a phylogeny  
254 to extract data from genomes such as orthologous or co-occurring genes<sup>93,105-107</sup>. Although a  
255 phylogeny has been used during the collection of such data, it is still important to control for  
256 any phylogenetic non-independence when analysing those data across genes, genomes and  
257 species<sup>93</sup>.

258 What about methods used in the fields of phylogenomics and evolutionary genomics?  
259 Researchers in these fields are developing tools to extract data from genomes such as  
260 orthologous or co-occurring genes, and to do this explicitly use a phylogeny<sup>93,95,105,106,108-110</sup>.  
261 However, just because a phylogeny has been used during the collection of any data, it is still  
262 important to control for any phylogenetic non-independence when analysing that data across  
263 genes, genomes and species<sup>93,111,112</sup>.

264 A possible exception arises for certain types of purely descriptive questions. For  
265 example, does the GC content of sequenced species vary between two orders of insects, such  
266 as Hymenoptera and Coleoptera<sup>113</sup>? While there is no evolutionary power to explain this

267 variation ( $n=2$  groups), one can still ask whether they differ. However, even with this kind of  
268 question, it could still be useful to determine an evolutionary average, which accounted for  
269 sampling across the phylogeny<sup>113</sup>. Furthermore, controlling for phylogeny would be required  
270 for other questions, such as does GC content correlate with chromosome size, or how can we  
271 explain variation in GC content<sup>113</sup>? Wording is key, to make clear what kind of question is  
272 being asked.

273

## 274 ***[H2] Does controlling for phylogeny reduce statistical power?***

275 Statistical power refers to how well a dataset can detect real effects (that is, avoiding Type I  
276 errors or ‘false negatives’). Controlling for confounding variables such as phylogeny does not  
277 reduce statistical power, because such variables should be accounted for before making any  
278 estimates of power<sup>18,114,115</sup>. Instead, controlling for phylogeny can reduce the chance of  
279 detecting a significant result that is not real (Type II error or ‘false positive’)<sup>18,42</sup>. For example,  
280 when a result that is driven by a phylogenetically correlated confounding factor, or a bias in  
281 genome sampling (Fig. 1). The general point here is that prioritizing the ability to produce  
282 significant results ahead of controlling for potential confounding variables can lead to incorrect  
283 conclusions not actually supported by the data. Statistical power is based upon the ability to  
284 detect *real* effects, not just significant effects. Phylogenetic analyses could also make it more  
285 likely to detect a real effect, by controlling for a confounding variable.

286

## 287 ***[H2] Is collecting more data a solution?***

288 More data is always preferable, but not all data are equal. Collecting many more data points  
289 from the same phylogenetic cluster could increase any bias due to phylogenetic non-  
290 independence (Fig. 1)<sup>16,18,20</sup>. If phylogeny is not controlled for, then increasing the number of  
291 data points in this way could make obtaining a spurious result more likely.

292         There can be trade-offs between data quantity and quality. It is important to sample  
293 across taxonomic breadth, especially where there have been changes (transitions) in key  
294 variables. Other possible factors to consider are the number of genomes per species, the quality  
295 of the genomes (for example, compare the large-scale database GenBank with the smaller, but  
296 higher-quality database RefSeq), the methods or tools used to assign gene function, or the  
297 assignment of species to a certain environment or lifestyle<sup>65,82,105,106,116,117</sup>. Applying higher  
298 thresholds towards any of these factors can lead to smaller but higher-quality datasets. The  
299 relative costs and benefits of different methods can depend upon the question being asked. For

300 example, population genetic analyses have found it easier to detect signatures of selection in  
301 smaller but higher-quality data sets<sup>117,118</sup>.

302

303 **[H2] Does statistical significance reflect biological importance?**

304 Statistical significance does not necessarily reflect biological importance. Usually, the level of  
305 significance (the size of a *P*-value) will depend on the size of the effect and the number of data  
306 points, tested against a null hypothesis of *exactly* zero difference, or *exactly* zero correlation<sup>119–</sup>  
307 <sup>121</sup>. Hence, in the extreme, very large datasets will almost always produce significant *P*-values,  
308 even if the estimated effect size is relatively small<sup>122</sup> (Fig. 4). Two groups are unlikely to ever  
309 have *identical* means, and a correlation of a finite number of data points is unlikely to ever be  
310 *exactly* zero. This can be a real problem for comparative genomic analyses, where datasets can  
311 be massive, made up of thousands of individual genes, replicons or genomes.

312         Consequently, it is crucial to consider the size of an effect, not just its significance<sup>122</sup>.  
313 Is an effect large enough to be biologically meaningful? One approach would be to examine  
314 the percentage of variance explained ( $R^2$ ). For example, do the explanatory variables explain  
315 at least 5–10% of variance in the response variable?<sup>123</sup> The maximum possible variance that  
316 could be explained will depend on the extent of measurement error<sup>69</sup>. To consider a specific  
317 case, in Fig. 4, even for the very small *P*-value of  $P=1 \times 10^{-4}$ , datasets of ~350 data points or  
318 more will produce significant results when the  $R^2$  value of the effect is less than 5%. At  
319  $n=10,000$ , a relationship explaining only 0.027% of the variance would be significant to  
320  $P<0.05$ . A result can be statistically significant even when it is unlikely to be biologically  
321 important.

322         It can also be helpful to compare the size of an effect to those typically found in a  
323 research field. What is the average % explained? What % is explained by a successful study?  
324 Considering the fields of ecology and evolutionary biology, the average % of variance  
325 explained in an analysis is 3.6%<sup>123</sup>. Particularly successful areas, such as the evolution of sex  
326 ratios and cooperation, have produced comparative studies that can explain 20–40% of the  
327 variation in data across species<sup>124,125</sup>.

328

329 **[H2] What can we conclude if there is little or no variation in an explanatory variable?**

330 The usual purpose of a comparative study is to explain variation that has been observed in  
331 nature. If there is a characteristic or trait that exhibits variation, but the cause of this variation  
332 has yet to be explained, one can use comparative analyses to examine what might explain that  
333 variation.

334 Imagine a scenario where a response variable (Y) showed considerable variation across  
335 species, but one of the possible explanatory variables (X) did not vary appreciably across  
336 species (Supplementary Fig. 2a). Does this mean that the possible role of the explanatory  
337 variable X cannot be tested? On the one hand, if we were interested in understanding the  
338 consequences of variation in X, then this dataset would not allow us to test this, because X does  
339 not vary. On the other hand, if we were interested in what explains variation in Y, it is clear  
340 that variation in Y is not explained by variation in X. This doesn't mean that X has no influence  
341 on Y, but rather that variation in X is not important for explaining the variation in Y that has  
342 been observed in nature. A lack of variation in a potential explanatory variable (X), while the  
343 response variable (Y) does vary, is an important result, not a failing of the dataset<sup>24</sup>.

344 In this case, the next step would be to look for other explanatory variables that could  
345 explain the variation in Y (Z etc.) (Supplementary Fig. 2b). Returning to the issue of data  
346 quality versus quantity, this also emphasises the benefit of trying to capture as much of the  
347 observed variation in nature as possible (for both response and explanatory variables).

348

## 349 **[H1] Conclusions**

350 Comparative analyses in the genomic age have the potential to answer broad and fundamental  
351 questions at the intersection of genetics and evolution. The explosion of genomic data during  
352 the first quarter of the 21<sup>st</sup> century calls for a rapid adoption of appropriate approaches and  
353 methods, or else risks an accumulation of spurious results and conclusions.

354 Methods are still being developed and debated for certain types of question, and there  
355 are cases where the lack of a good phylogeny will pose a problem<sup>112</sup>. In cases where  
356 complications arise it can be extremely useful to analyse the data with multiple methods to test  
357 the robustness of conclusions<sup>24,126</sup>. Nonetheless, there are many situations where relatively  
358 standard and accepted methods can be used. And testing different ways of controlling for  
359 phylogeny is better than not controlling for phylogeny (Box 1). It is also worth looking back  
360 historically — while there was initially resistance and misunderstanding about the application  
361 of phylogenetic methods to study adaptation at the organismal level, there is now no doubt that  
362 they have revolutionized the field<sup>20,34,42,127,128</sup>.

363 Finally, although comparative genomics can provide powerful insights, the success of  
364 comparative genomics relies on experiments and observations, both to generate hypotheses and  
365 data. Insights from comparative analyses can reveal broad patterns which can then explicitly  
366 be tested experimentally. If comparative analyses can increase the chance of choosing the

367 correct target for laboratory work, then this could provide an efficiency benefit by preventing  
368 wasted experimental work. Vice versa, experimental insights can generate hypotheses which  
369 can then be examined in comparative analyses across species to look for generality. Both  
370 approaches have different pros and cons, and the greatest insights can often be made by their  
371 combination<sup>34</sup>.

372

## 373 **References**

- 374 1. Binnewies, T. T. *et al.* Ten years of bacterial genome sequencing: comparative-  
375 genomics-based discoveries. *Funct Integr Genomics* **6**, 165–185 (2006).
- 376 2. Land, M. *et al.* Insights from 20 years of bacterial genome sequencing. *Funct Integr*  
377 *Genomics* **15**, 141–161 (2015).
- 378 3. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,  
379 taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733-745 (2016).
- 380 4. The Darwin Tree of Life Project Consortium. Sequence locally, think globally: The  
381 Darwin Tree of Life Project. *Proceedings of the National Academy of Sciences* **119**,  
382 e2115642118 (2022).
- 383 5. Hunt, M., Lima, L., Shen, W., Lees, J. & Iqbal, Z. AllTheBacteria - all bacterial genomes  
384 assembled, available and searchable. 2024.03.08.584059 Preprint at  
385 <https://doi.org/10.1101/2024.03.08.584059> (2024).
- 386 6. David, S. *et al.* Epidemic of carbapenem-resistant *Klebsiella pneumoniae* in Europe is  
387 driven by nosocomial spread. *Nat Microbiol* **4**, 1919–1929 (2019).
- 388 7. León-Sampedro, R. *et al.* Pervasive transmission of a carbapenem resistance plasmid in  
389 the gut microbiota of hospitalized patients. *Nat Microbiol* **6**, 606–616 (2021).
- 390 8. Xin, X.-F., Kvitko, B. & He, S. Y. *Pseudomonas syringae*: what it takes to be a pathogen.  
391 *Nat Rev Microbiol* **16**, 316–328 (2018).
- 392 9. Sarkar, S. F., Gordon, J. S., Martin, G. B. & Guttman, D. S. Comparative Genomics of  
393 Host-Specific Virulence in *Pseudomonas syringae*. *Genetics* **174**, 1041–1056 (2006).

- 394 10. Li, Y., Liu, Z., Shi, P. & Zhang, J. The hearing gene Prestin unites echolocating bats and  
395 whales. *Current Biology* **20**, R55–R56 (2010).
- 396 11. Liu, Y. *et al.* Convergent sequence evolution between echolocating bats and dolphins.  
397 *Current Biology* **20**, R53–R54 (2010).
- 398 12. Yuan, Y. *et al.* Comparative genomics provides insights into the aquatic adaptations of  
399 mammals. *Proceedings of the National Academy of Sciences* **118**, e2106080118 (2021).
- 400 13. Kruskal, W. Miracles and Statistics: The Casual Assumption of Independence. *null* **83**,  
401 929–940 (1988).
- 402 14. Ives, A. R. & Zhu, J. Statistics for correlated data: phylogenies, space, and time. *Ecol*  
403 *Appl* **16**, 20–32 (2006).
- 404 15. Whitney, K. D. & Jr, T. G. Did Genetic Drift Drive Increases in Genome Complexity?  
405 *PLOS Genetics* **6**, e1001080 (2010).
- 406 16. Harvey, P. H. & Pagel, M. D. *The Comparative Method in Evolutionary Biology*. (Oxford  
407 University Press, Oxford, New York, 1991).
- 408 17. Harvey, P. H. & Purvis, A. Comparative methods for explaining adaptations. *Nature* **351**,  
409 619–624 (1991).
- 410 18. Felsenstein, J. Phylogenies and the Comparative Method. *The American Naturalist* **125**,  
411 1–15 (1985).
- 412 19. Grafen, A. The phylogenetic regression. *Philos Trans R Soc Lond B Biol Sci* **326**, 119–  
413 157 (1989).
- 414 20. Ridley, M. Why not to use species in comparative tests. *Journal of Theoretical Biology*  
415 **136**, 361–364 (1989).
- 416 21. Hardison, R. C. Comparative Genomics. *PLOS Biology* **1**, e58 (2003).
- 417 22. Dewar, A. E. *et al.* Plasmids do not consistently stabilize cooperation across bacteria but  
418 may promote broad pathogen host-range. *Nat Ecol Evol* **5**, 1624–1636 (2021).

- 419 23. Dewar, A. E., Belcher, L. J., Scott, T. W. & West, S. A. Genes for cooperation are not  
420 more likely to be carried by plasmids. *Proceedings of the Royal Society B: Biological*  
421 *Sciences* **291**, 20232549 (2024).
- 422 24. Dewar, A. E., Hao, C., Belcher, L. J., Ghoul, M. & West, S. A. Bacterial lifestyle shapes  
423 pangenomes. *Proceedings of the National Academy of Sciences* **121**, e2320170121  
424 (2024).
- 425 25. Bouvier, J. W. & Kelly, S. Response to Tcherkez and Farquhar: Rubisco adaptation is  
426 more limited by phylogenetic constraint than by catalytic trade-off. *Journal of Plant*  
427 *Physiology* **287**, 154021 (2023).
- 428 26. Bouvier, J. W. *et al.* Rubisco Adaptation Is More Limited by Phylogenetic Constraint  
429 Than by Catalytic Trade-off. *Molecular Biology and Evolution* **38**, 2880–2896 (2021).
- 430 27. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable snapshot  
431 of archived DNA sequences. *PLOS Biology* **19**, e3001421 (2021).
- 432 28. Feng, S. *et al.* Dense sampling of bird diversity increases power of comparative  
433 genomics. *Nature* **587**, 252–257 (2020).
- 434 29. Upham, N. S. & Landis, M. J. Genomics expands the mammalverse. *Science* **380**, 358–  
435 359 (2023).
- 436 30. Clutton-Brock, T. H. & Harvey, P. H. Primate ecology and social organization. *Journal*  
437 *of Zoology* **183**, 1–39 (1977).
- 438 31. Clutton-Brock, T. H. & Harvey, P. H. Comparison and Adaptation. *Proceedings of the*  
439 *Royal Society of London. Series B, Biological Sciences* **205**, 547–565 (1979).
- 440 32. Ridley, M. *The Explanation of Organic Diversity: The Comparative Method and*  
441 *Adaptations for Mating.* (Clarendon Press, 1983).
- 442 33. Pagel, M. Inferring the historical patterns of biological evolution. *Nature* **401**, 877–884  
443 (1999).

- 444 34. Davies, N. B., Krebs, J. R. & West, S. A. *An Introduction to Behavioural Ecology*.  
445 (Wiley-Blackwell, 2013).
- 446 35. Loos, R. J. F. 15 years of genome-wide association studies and no signs of slowing down.  
447 *Nat Commun* **11**, 5900 (2020).
- 448 36. Tam, V. *et al.* Benefits and limitations of genome-wide association studies. *Nat Rev*  
449 *Genet* **20**, 467–484 (2019).
- 450 37. Murray, G. G. R. *et al.* Genome Reduction Is Associated with Bacterial Pathogenicity  
451 across Different Scales of Temporal and Ecological Divergence. *Molecular Biology and*  
452 *Evolution* **38**, 1570–1579 (2021).
- 453 38. Beavan, A., Domingo-Sananes, M. R. & McInerney, J. O. Contingency, repeatability, and  
454 predictability in the evolution of a prokaryotic pangenome. *Proceedings of the National*  
455 *Academy of Sciences* **121**, e2304934120 (2024).
- 456 39. Godfroid, M. *et al.* Evo-Scope: Fully automated assessment of correlated evolution on  
457 phylogenetic trees. *Methods in Ecology and Evolution* **15**, 282–289 (2024).
- 458 40. Martinez, J., Klasson, L., Welch, J. J. & Jiggins, F. M. Life and Death of Selfish Genes:  
459 Comparative Genomics Reveals the Dynamic Evolution of Cytoplasmic Incompatibility.  
460 *Molecular Biology and Evolution* **38**, 2–15 (2021).
- 461 41. Nee, S. Birth-Death Models in Macroevolution. *Annual Review of Ecology, Evolution,*  
462 *and Systematics* **37**, 1–17 (2006).
- 463 42. Cornwallis, C. K. & Griffin, A. S. A Guided Tour of Phylogenetic Comparative Methods  
464 for Studying Trait Evolution. (2024) doi:10.1146/annurev-ecolsys-102221-050754.
- 465 43. Revell, L. J. & Harmon, L. J. *Phylogenetic Comparative Methods in R*. (Princeton  
466 University Press, 2022).
- 467 44. Revell, L. J. Phylogenetic signal and linear regression on species data. *Methods in*  
468 *Ecology and Evolution* **1**, 319–329 (2010).

- 469 45. Ives, A. R. & Garland, T. Phylogenetic logistic regression for binary dependent variables.  
470 *Syst Biol* **59**, 9–26 (2010).
- 471 46. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and  
472 evolutionary analyses in R. *Bioinformatics* **35**, 526–528 (2019).
- 473 47. Tung Ho, L. si & Ané, C. A Linear-Time Algorithm for Gaussian and Non-Gaussian  
474 Trait Evolution Models. *Systematic Biology* **63**, 397–408 (2014).
- 475 48. Orme, D. *et al.* CAPER: comparative analyses of phylogenetics and evolution in R.  
476 *Methods in Ecology and Evolution* **3**, 145–151 (2013).
- 477 49. Hadfield, J. D. MCMC Methods for Multi-Response Generalized Linear Mixed Models:  
478 The MCMCglmm R Package. *Journal of Statistical Software* **33**, 1–22 (2010).
- 479 50. Pagel, M. & Meade, A. Bayesian Analysis of Correlated Evolution of Discrete Characters  
480 by Reversible-Jump Markov Chain Monte Carlo. *The American Naturalist* **167**, 808–825  
481 (2006).
- 482 51. Pagel, M., Meade, A. & Barker, D. Bayesian estimation of ancestral character states on  
483 phylogenies. *Systematic biology* **53**, 673–684 (2004).
- 484 52. Bürkner, P.-C. brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal*  
485 *of Statistical Software* **80**, 1–28 (2017).
- 486 53. Mundry, R. Chapter 6 ‘Statistical Issues and Assumptions of Phylogenetic Generalized  
487 Least Squares’ in *Modern Phylogenetic Comparative Methods and Their Application in*  
488 *Evolutionary Biology: Concepts and Practice*. (Springer, 2014).
- 489 54. Losos, J. B. Seeing the Forest for the Trees: The Limitations of Phylogenies in  
490 Comparative Biology. *The American Naturalist* **177**, 709–727 (2011).
- 491 55. Blomberg, S. P., Garland, T. & Ives, A. R. Testing for Phylogenetic Signal in  
492 Comparative Data: Behavioral Traits Are More Labile. *Evolution* **57**, 717–745 (2003).

- 493 56. Revell, L. J., Harmon, L. J. & Collar, D. C. Phylogenetic Signal, Evolutionary Process,  
494 and Rate. *Systematic Biology* **57**, 591–601 (2008).
- 495 57. Bergeron, L. A. *et al.* Evolution of the germline mutation rate across vertebrates. *Nature*  
496 **615**, 285–291 (2023).
- 497 58. Imrit, M. A., Dogantzis, K. A., Harpur, B. A. & Zayed, A. Eusociality influences the  
498 strength of negative selection on insect genomes. *Proceedings of the Royal Society B:*  
499 *Biological Sciences* **287**, 20201512 (2020).
- 500 59. Rubin, B. E. R. Social insect colony size is correlated with rates of molecular evolution.  
501 *Insect. Soc.* **69**, 147–157 (2022).
- 502 60. Ruis, C. *et al.* Mutational spectra are associated with bacterial niche. *Nat Commun* **14**,  
503 7091 (2023).
- 504 61. Wyatt, C. D. R. *et al.* Social complexity, life-history and lineage influence the molecular  
505 basis of castes in vespid wasps. *Nat Commun* **14**, 1046 (2023).
- 506 62. Raulo, A. *et al.* Social and environmental transmission spread different sets of gut  
507 microbes in wild mice. *Nat Ecol Evol* **8**, 972–985 (2024).
- 508 63. Ghoul, M., Andersen, S. B. & West, S. A. Sociomics: Using Omic Approaches to  
509 Understand Social Evolution. *Trends in Genetics* **33**, 408–419 (2017).
- 510 64. Downing, T. & Angelopoulos, N. A primer on correlation-based dimension reduction  
511 methods for multi-omics analysis. *Journal of The Royal Society Interface* **20**, 20230344  
512 (2023).
- 513 65. Hao, C., Dewar, A. E., West, S. A. & Ghoul, M. Gene transferability and sociality do not  
514 correlate with gene connectivity. *Proceedings of the Royal Society B: Biological Sciences*  
515 **289**, 20221819 (2022).
- 516 66. McNally, L., Viana, M. & Brown, S. P. Cooperative secretions facilitate host range  
517 expansion in bacteria. *Nat Commun* **5**, 4594 (2014).

- 518 67. Simonet, C. & McNally, L. Kin selection explains the evolution of cooperation in the gut  
519 microbiota. *PNAS* **118**, (2021).
- 520 68. Chak, S. T. C., Harris, S. E., Hultgren, K. M., Jeffery, N. W. & Rubenstein, D. R.  
521 Eusociality in snapping shrimps is associated with larger genomes and an accumulation  
522 of transposable elements. *Proceedings of the National Academy of Sciences* **118**,  
523 e2025051118 (2021).
- 524 69. Grafen, A., Hails, R., Grafen, A. & Hails, R. *Modern Statistics for the Life Sciences*.  
525 (Oxford University Press, Oxford, New York, 2002).
- 526 70. Zoonomia Consortium. A comparative genomics multitool for scientific discovery and  
527 conservation. *Nature* **587**, 240–245 (2020).
- 528 71. Lloyd-Price, J. *et al.* Strains, functions and dynamics in the expanded Human  
529 Microbiome Project. *Nature* **550**, 61–66 (2017).
- 530 72. Duarte, C. M. *et al.* Sequencing effort dictates gene discovery in marine microbial  
531 metagenomes. *Environ Microbiol* **22**, 4589–4603 (2020).
- 532 73. Arikawa, K. & Hosokawa, M. Uncultured prokaryotic genomes in the spotlight: An  
533 examination of publicly available data from metagenomics and single-cell genomics.  
534 *Comput Struct Biotechnol J* **21**, 4508–4518 (2023).
- 535 74. Rhie, A. *et al.* Towards complete and error-free genome assemblies of all vertebrate  
536 species. *Nature* **592**, 737–746 (2021).
- 537 75. Garcia-Garcera, M. & Rocha, E. P. C. Community diversity and habitat structure shape  
538 the repertoire of extracellular proteins in bacteria. *Nature Communications* **11**, 758  
539 (2020).
- 540 76. Shaw, L. P. *et al.* Niche and local geography shape the pangenome of wastewater- and  
541 livestock-associated Enterobacteriaceae. *Science Advances* **7**, eabe3868 (2021).

- 542 77. Shaw, L. P., Rocha, E. P. C. & MacLean, R. C. Restriction-modification systems have  
543 shaped the evolution and distribution of plasmids across bacteria. *Nucleic Acids Research*  
544 **51**, 6806–6818 (2023).
- 545 78. Hall, R. J. *et al.* Gene-gene relationships in an Escherichia coli accessory genome are  
546 linked to function and mobility. *Microbial Genomics* **7**, 000650 (2021).
- 547 79. Whelan, F. J., Hall, R. J. & McInerney, J. O. Evidence for selection in the abundant  
548 accessory gene content of a prokaryote pangenome. *Molecular Biology and Evolution*  
549 (2021) doi:10.1093/molbev/msab139.
- 550 80. Hu, E.-Z., Lan, X.-R., Liu, Z.-L., Gao, J. & Niu, D.-K. A positive correlation between  
551 GC content and growth temperature in prokaryotes. *BMC Genomics* **23**, 110 (2022).
- 552 81. Haudiquet, M. *et al.* Capsules and their traits shape phage susceptibility and plasmid  
553 conjugation efficiency. *Nat Commun* **15**, 2032 (2024).
- 554 82. Rendueles, O., Sousa, J. A. M. de, Bernheim, A., Touchon, M. & Rocha, E. P. C. Genetic  
555 exchanges are more frequent in bacteria encoding capsules. *PLOS Genetics* **14**, e1007862  
556 (2018).
- 557 83. Rendueles, O., Garcia-Garcerà, M., Néron, B., Touchon, M. & Rocha, E. P. C.  
558 Abundance and co-occurrence of extracellular capsules increase environmental breadth:  
559 Implications for the emergence of pathogens. *PLoS Pathog* **13**, e1006525 (2017).
- 560 84. Garamszegi, L. Z. *Chapter 8 in Modern Phylogenetic Comparative Methods and Their*  
561 *Application in Evolutionary Biology: Concepts and Practice.* (Springer, 2014).
- 562 85. Pagel, M. Detecting Correlated Evolution on Phylogenies: A General Method for the  
563 Comparative Analysis of Discrete Characters. *Proceedings: Biological Sciences* **255**, 37–  
564 45 (1994).
- 565 86. Boyko, J. D. & Beaulieu, J. M. Generalized hidden Markov models for phylogenetic  
566 comparative datasets. *Methods in Ecology and Evolution* **12**, 468–478 (2021).

- 567 87. Harmon, L. J. *Phylogenetic Comparative Methods: Learning from Trees*. (CreateSpace  
568 Independent Publishing Platform, 2018).
- 569 88. van der Bijl, W. phylopath: Easy phylogenetic path analysis in R. *PeerJ* **6**, e4718 (2018).
- 570 89. Hardenberg, A. von & Gonzalez-Voyer, A. Disentangling Evolutionary Cause-Effect  
571 Relationships with Phylogenetic Confirmatory Path Analysis. *Evolution* **67**, 378–387  
572 (2013).
- 573 90. Cornwell, W. & Nakagawa, S. Phylogenetic comparative methods. *Curr. Biol.* **27**, R333–  
574 R336 (2017).
- 575 91. Cornwallis, C. K. *et al.* Cooperation facilitates the colonization of harsh environments.  
576 *Nat Ecol Evol* **1**, 1–10 (2017).
- 577 92. Frígols, B. *et al.* Virus Satellites Drive Viral Evolution and Ecology. *PLOS Genetics* **11**,  
578 e1005609 (2015).
- 579 93. Gavriilidou, A. *et al.* Goldfinder: Unraveling Networks of Gene Co-occurrence and  
580 Avoidance in Bacterial Pangenomes. 2024.04.29.591652 Preprint at  
581 <https://doi.org/10.1101/2024.04.29.591652> (2024).
- 582 94. Leeks, A., Young, P. G., Turner, P. E., Wild, G. & West, S. A. Cheating leads to the  
583 evolution of multipartite viruses. *PLOS Biology* **21**, e3002092 (2023).
- 584 95. Hu, Z., Sackton, T. B., Edwards, S. V. & Liu, J. S. Bayesian Detection of Convergent  
585 Rate Changes of Conserved Noncoding Elements on Phylogenetic Trees. *Molecular*  
586 *Biology and Evolution* **36**, 1086–1100 (2019).
- 587 96. Eastman, J. M., Alfaro, M. E., Joyce, P., Hipp, A. L. & Harmon, L. J. A novel  
588 comparative method for identifying shifts in the rate of character evolution on trees.  
589 *Evolution* **65**, 3578–3589 (2011).
- 590 97. Garamszegi, L. Z. *Modern Phylogenetic Comparative Methods and Their Application in*  
591 *Evolutionary Biology: Concepts and Practice*. (Springer, 2014).

- 592 98. Coluzzi, C., Garcillán-Barcia, M. P., de la Cruz, F. & Rocha, E. P. C. Evolution of  
593 Plasmid Mobility: Origin and Fate of Conjugative and Nonconjugative Plasmids.  
594 *Molecular Biology and Evolution* **39**, msac115 (2022).
- 595 99. Acman, M., van Dorp, L., Santini, J. M. & Balloux, F. Large-scale network analysis  
596 captures biological features of bacterial plasmids. *Nat Commun* **11**, 2452 (2020).
- 597 100. Matlock, W. *et al.* Genomic network analysis of environmental and livestock F-type  
598 plasmid populations. *ISME J* **15**, 2322–2335 (2021).
- 599 101. Adams, R. *et al.* A tale of too many trees: a conundrum for phylogenetic regression.  
600 2024.02.16.580530 Preprint at <https://doi.org/10.1101/2024.02.16.580530> (2024).
- 601 102. Nogueira, T., Touchon, M. & Rocha, E. P. C. Rapid Evolution of the Sequences and  
602 Gene Repertoires of Secreted Proteins in Bacteria. *PLoS One* **7**, e49403 (2012).
- 603 103. Maddamsetti, R. *et al.* Duplicated antibiotic resistance genes reveal ongoing selection  
604 and horizontal gene transfer in bacteria. *Nat Commun* **15**, 1449 (2024).
- 605 104. Pagel, M. D. & Harvey, P. H. Recent Developments in the Analysis of Comparative  
606 Data. *The Quarterly Review of Biology* **63**, 413–440 (1988).
- 607 105. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for  
608 comparative genomics. *Genome Biology* **20**, 238 (2019).
- 609 106. Whelan, F. J., Rusilowicz, M. & McInerney, J. O. Coinfinder: detecting significant  
610 associations and dissociations in pangenomes. *Microb Genom* **6**, e000338 (2020).
- 611 107. Bernheim, A., Bikard, D., Touchon, M. & Rocha, E. P. C. Atypical organizations and  
612 epistatic interactions of CRISPRs and cas clusters in genomes and their mobile genetic  
613 elements. *Nucleic Acids Research* **48**, 748–760 (2020).
- 614 108. Pond, S. L. K., Frost, S. D. W. & Muse, S. V. HyPhy: hypothesis testing using  
615 phylogenies. *Bioinformatics* **21**, 676–679 (2005).

- 616 109. Kowalczyk, A. *et al.* RERconverge: an R package for associating evolutionary rates  
617 with convergent traits. *Bioinformatics* **35**, 4815–4817 (2019).
- 618 110. De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE: a computational tool  
619 for the study of gene family evolution. *Bioinformatics* **22**, 1269–1271 (2006).
- 620 111. Read, A. F. & Nee, S. Inference from binary comparative data. *Journal of Theoretical*  
621 *Biology* **173**, 99–108 (1995).
- 622 112. Maddison, W. P. & FitzJohn, R. G. The Unsolved Challenge to Phylogenetic  
623 Correlation Tests for Categorical Characters. *Systematic Biology* **64**, 127–136 (2015).
- 624 113. Kyriacou, R. G., Mulhair, P. O. & Holland, P. W. H. GC Content Across Insect  
625 Genomes: Phylogenetic Patterns, Causes and Consequences. *J Mol Evol* **92**, 138–152  
626 (2024).
- 627 114. Boettiger, C., Coop, G. & Ralph, P. Is your phylogeny informative? Measuring the  
628 power of comparative methods. *Evolution* **66**, 2240–2251 (2012).
- 629 115. Uyeda, J. C., Zenil-Ferguson, R. & Pennell, M. W. Rethinking phylogenetic  
630 comparative methods. *Systematic Biology* **67**, 1091–1109 (2018).
- 631 116. Gupta, A., Kapil, R., Dhakan, D. B. & Sharma, V. K. MP3: A Software Tool for the  
632 Prediction of Pathogenic Proteins in Genomic and Metagenomic Data. *PLOS ONE* **9**,  
633 e93907 (2014).
- 634 117. Belcher, L. J. *et al.* SOCfinder: a genomic tool for identifying social genes in bacteria.  
635 *Microbial Genomics* **9**, 001171 (2023).
- 636 118. Belcher, L. J., Dewar, A. E., Ghoul, M. & West, S. A. Kin selection for cooperation in  
637 natural bacterial populations. *Proceedings of the National Academy of Sciences* **119**,  
638 e2119070119 (2022).
- 639 119. Szucs, D. & Ioannidis, J. P. A. When Null Hypothesis Significance Testing Is  
640 Unsuitable for Research: A Reassessment. *Front Hum Neurosci* **11**, 390 (2017).

- 641 120. Cohen, J. The earth is round ( $p < .05$ ). *American Psychologist* **49**, 997–1003 (1994).
- 642 121. Tukey, J. W. The Philosophy of Multiple Comparisons. *Statistical Science* **6**, 100–116  
643 (1991).
- 644 122. Sullivan, G. M. & Feinn, R. Using Effect Size—or Why the P Value Is Not Enough. *J*  
645 *Grad Med Educ* **4**, 279–282 (2012).
- 646 123. Jennions, M. D. & Møller, A. P. A survey of the statistical power of research in  
647 behavioral ecology and animal behavior. *Behavioral Ecology* **14**, 438–445 (2003).
- 648 124. West, S. *Sex Allocation*. (Princeton University Press, 2009).
- 649 125. West, S. A., Shuker, D. M. & Sheldon, B. C. Sex-Ratio Adjustment When Relatives  
650 Interact: A Test of Constraints on Adaptation. *Evolution* **59**, 1211–1228 (2005).
- 651 126. Cornwallis, C. K. *et al.* Symbioses shape feeding niches and diversification across  
652 insects. *Nat Ecol Evol* **7**, 1022–1044 (2023).
- 653 127. Harvey, P. H., Read, A. F. & Nee, S. Why Ecologists Need to be Phylogenetically  
654 Challenged. *The Journal of Ecology* **83**, 535 (1995).
- 655 128. Harvey, P. H., Read, A. F. & Nee, S. Further Remarks on the Role of Phylogeny in  
656 Comparative Ecology. *Journal of Ecology* **83**, 733–734 (1995).
- 657 129. Garland, T. & Ives, A. R. Using the Past to Predict the Present: Confidence Intervals  
658 for Regression Equations in Phylogenetic Comparative Methods. *Am Nat* **155**, 346–364  
659 (2000).
- 660 130. Nakagawa, S., Johnson, P. C. D. & Schielzeth, H. The coefficient of determination  $R^2$   
661 and intra-class correlation coefficient from generalized linear mixed-effects models  
662 revisited and expanded. *J. R. Soc. Interface* **11** (2017).
- 663 131. Ives, A. R.  $R^2$  for Correlated Data: Phylogenetic Models, LMMs, and  
664 GLMMs. *Systematic Biology* **68**, 234–251 (2019).

- 665 132. Beaulieu, J. M., O'Meara, B. C. & Donoghue, M. J. Identifying Hidden Rate Changes  
666 in the Evolution of a Binary Morphological Character: The Evolution of Plant Habit in  
667 Campanulid Angiosperms. *Systematic Biology* **62**, 725–737 (2013).
- 668 133. Bell-Roberts, L. *et al.* Larger colony sizes favoured the evolution of more worker  
669 castes in ants. *Nat Ecol Evol* 1–13 (2024) doi:10.1038/s41559-024-02512-7.
- 670 134. Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity  
671 through a phylogenetically consistent, rank normalized and complete genome-based  
672 taxonomy. *Nucleic Acids Research* **50**, D785–D794 (2022).
- 673 135. Downing, T. & Rahm, A. Bacterial plasmid-associated and chromosomal proteins  
674 have fundamentally different properties in protein interaction networks. *Sci Rep* **12**,  
675 19203 (2022).
- 676 136. West, S. A., Griffin, A. S., Gardner, A. & Diggle, S. P. Social evolution theory for  
677 microorganisms. *Nat Rev Microbiol* **4**, 597–607 (2006).
- 678 137. Smith, J. The social evolution of bacterial pathogenesis. *Proc. R. Soc. Lond. B* **268**,  
679 61–69 (2001).
- 680 138. Mc Ginty, S. É., Lehmann, L., Brown, S. P. & Rankin, D. J. The interplay between  
681 relatedness and horizontal gene transfer drives the evolution of plasmid-carried public  
682 goods. *Proc. R. Soc. B* **280**, 20130400 (2013).
- 683 139. Scott, T. W., West, S. A., Dewar, A. E. & Wild, G. Is cooperation favored by  
684 horizontal gene transfer? *Evolution Letters* **7**, 113–120 (2023).
- 685 140. Savir, Y., Noor, E., Milo, R. & Tlusty, T. Cross-species analysis traces adaptation of  
686 Rubisco toward optimality in a low-dimensional landscape. *Proceedings of the National*  
687 *Academy of Sciences* **107**, 3475–3480 (2010).
- 688 141. Flamholz, A. I. *et al.* Revisiting Trade-offs between Rubisco Kinetic Parameters.  
689 *Biochemistry* **58**, 3365–3376 (2019).

- 690 142. Tcherkez, G. G. B., Farquhar, G. D. & Andrews, T. J. Despite slow catalysis and  
691 confused substrate specificity, all ribulose biphosphate carboxylases may be nearly  
692 perfectly optimized. *Proceedings of the National Academy of Sciences* **103**, 7246–7251  
693 (2006).
- 694 143. Koonin, E. V. Orthologs, Paralogs, and Evolutionary Genomics. *Annual Review of*  
695 *Genetics* **39**, 309–338 (2005).
- 696 144. Kryuchkova-Mostacci, N. & Robinson-Rechavi, M. Tissue-Specificity of Gene  
697 Expression Diverges Slowly between Orthologs, and Rapidly between Paralogs. *PLOS*  
698 *Computational Biology* **12**, e1005274 (2016).
- 699 145. Dunn, C. W., Zapata, F., Munro, C., Siebert, S. & Hejnol, A. Pairwise comparisons  
700 across species are problematic when analyzing functional genomic data. *Proceedings of*  
701 *the National Academy of Sciences* **115**, E409–E417 (2018).

702  
703

## 704 **Acknowledgements**

705 The authors thank A. Grafen, E. Rocha, J. Bouvier, P. Holland and S. Shimeld for useful  
706 comments on the manuscript; A. Griffin, J. Turner, L. Bell-Roberts, M. Brindle, M. Liu, R.  
707 Bonifacii, S. Kershenbaum and Z. Katz for helpful discussion; and three anonymous reviewers  
708 for their feedback. The authors thank the European Research Council (834164) and St John’s  
709 College, Oxford for funding.

710

## 711 **Competing interests**

712 The authors declare no competing interests.

713

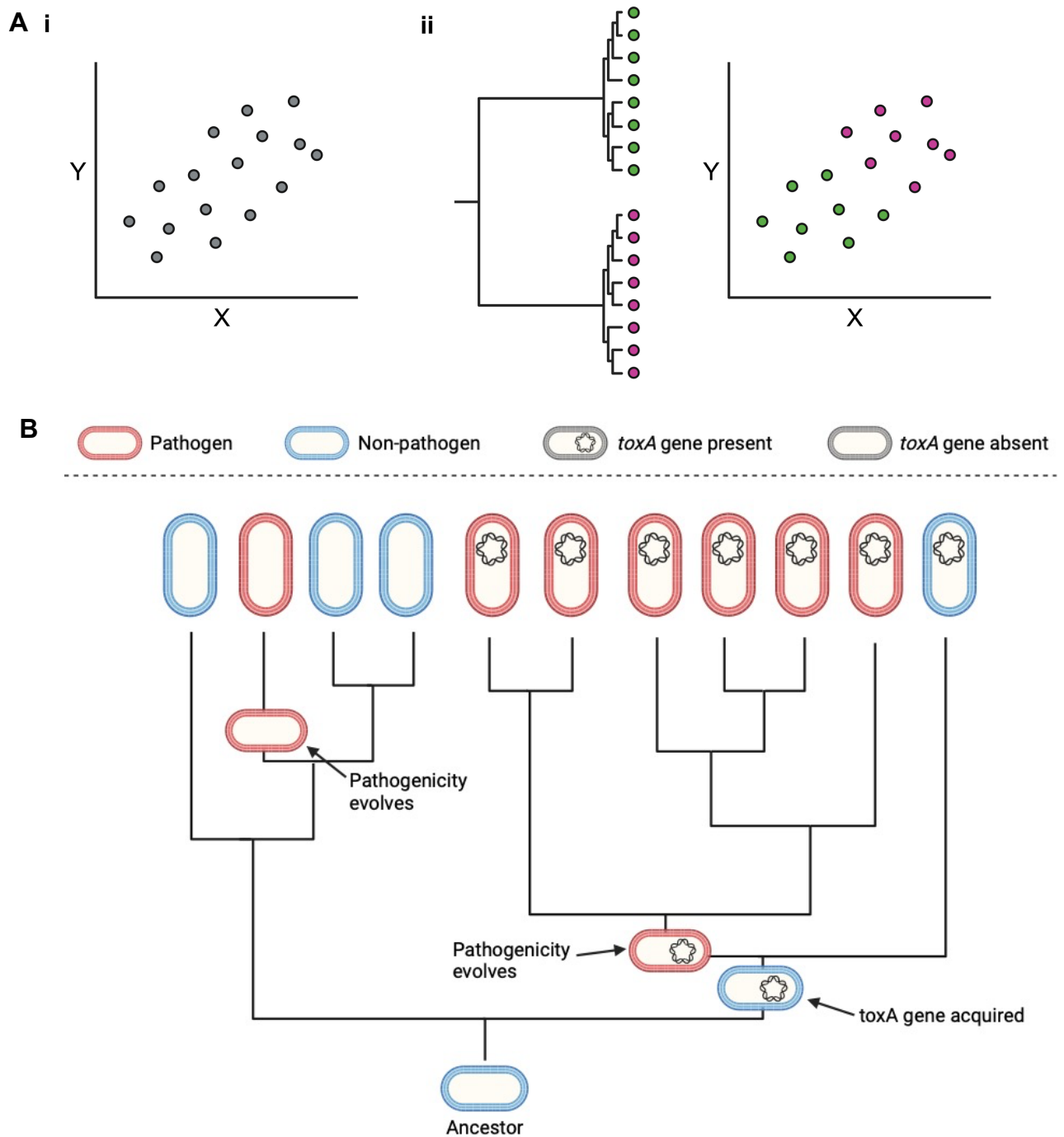
714 **Table 1. Current methods to identify and/or control for phylogenetic non-independence.**

Method	What it does	When to use it	How to use it
Phylogenetic	Adds a trait’s correlation with evolutionary history	Estimating correlation between two variables while controlling for phylogeny.	R packages ‘phylolm’, ‘ape’, ‘caper’.

regression <sup>19,44,45,129</sup> (phylogenetic least squares, PGLS)	as an additional term.	Can be generalised for response variables such as count data. Use phylogenetic logistic regression for binominal or binary response variables. <sup>45</sup> Estimating phylogenetic signal. Continuous or discrete data.	
Phylogenetic mixed model <sup>49,52,86,130,131</sup>	Includes phylogenetic similarity as a random effect within a mixed model.	For multiple explanatory traits of interest. Bayesian and frequentist options. Continuous or discrete data.	R packages 'MCMCglmm', 'brms', 'phyr'. Computer package 'BayesTraits'.
Independent contrasts <sup>18,37</sup>	Tests correlation between traits across pairs of closely related species. Equivalent to phylogenetic regression.	When phylogeny uncertain or dataset includes many closely related pairs of species with different trait values. Largely equivalent to phylogenetic regression. <sup>19</sup> Continuous or discrete data.	R packages 'ape', 'caper', 'castor'.
Ancestral state reconstruction <sup>43,51,85-87,132</sup>	Uses trait values of extant species to reconstruct likely trait values of ancestors.	To examine how many times a trait has evolved, and compare the ancestral states of multiple traits to examine evidence for correlated evolution. Continuous or discrete data.	R packages 'corHMM', 'MCMCglmm' (as described in <sup>91,133</sup> ). Computer package 'BayesTraits'.
Correlated evolution models <sup>33,39,51,85</sup>	Tests if two traits are correlated with one another more than expected due to phylogenetic history.	When both traits are binary, or when continuous traits can be expressed as binary (i.e. 'high' and 'low'). Often includes 'transition rate' analysis, to examine how evolution occurs between states. Discrete data.	R packages 'ape', 'phytools'. Computer package 'BayesTraits', 'RevBayes'.
Phylogenetic path analysis <sup>84,88,89</sup>	Compares support for different causal hypotheses, while taking phylogeny into account.	Use if three or more variables/traits of interest, and goal is to test between different causal hypotheses. Continuous or discrete data.	R package 'phylopath'.

715  
716  
717

### Figures

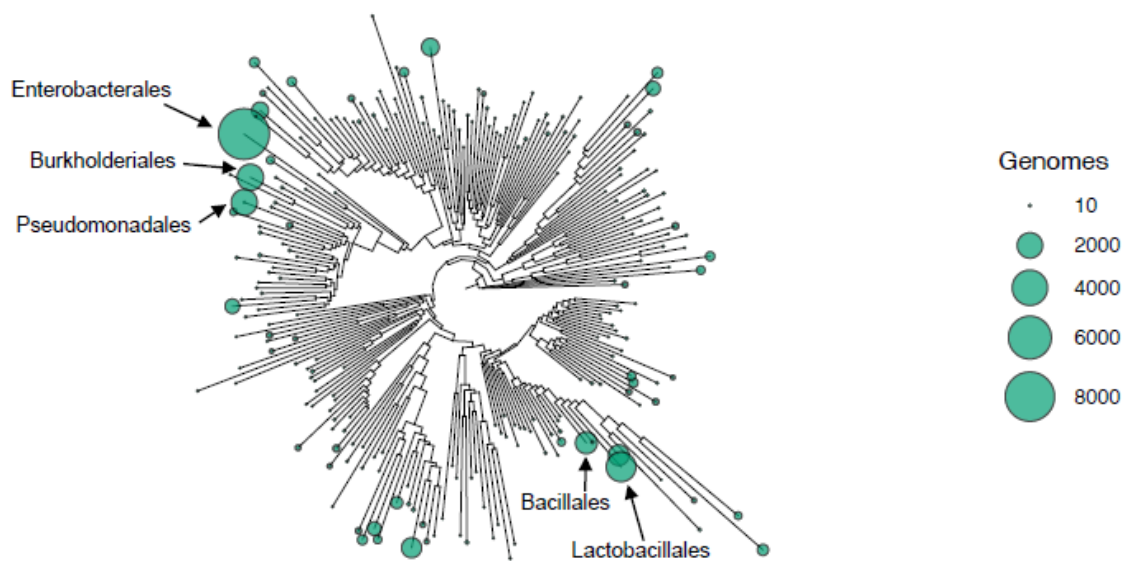


718  
719  
720  
721  
722

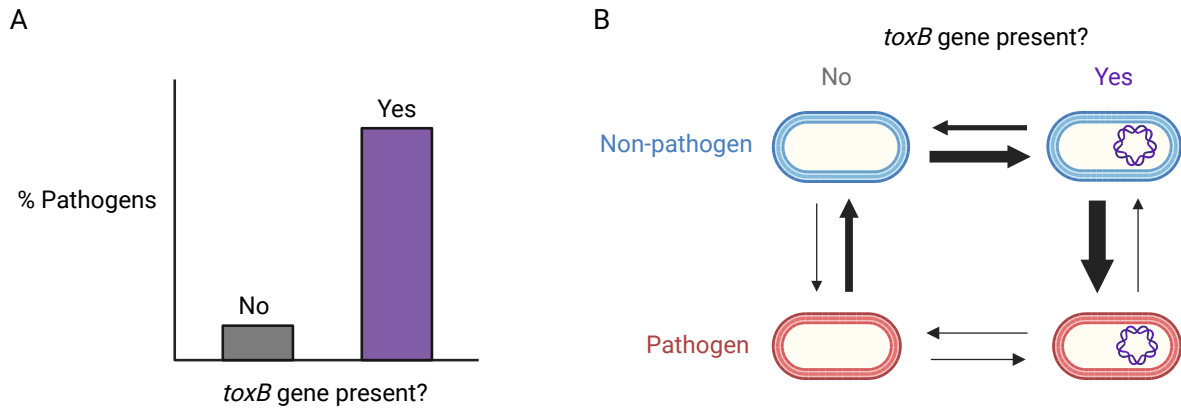
**Figure 1. Species are not independent data points.** A.i. There seems to be a strong positive correlation between X and Y, where each dot is a species ( $n=16$  species). ii. In an extreme scenario those species are from two separate, monophyletic lineages represented by green and pink dots. When mapped onto the original scatterplot, the two lineages form largely separate

723 clusters. Within each of those clusters, there is no relationship between X and Y. The original  
 724 correlation was just an artefact of the mean values of X and Y in the pink lineage being larger  
 725 than the green lineage. Inspired by Figures 5-7 in Felsenstein, 1985<sup>18</sup>. B. A simple hypothetical  
 726 genomic example. Each tip is a bacterial species; red and blue cells correspond to pathogenic  
 727 and non-pathogenic species, respectively, and species which carry *toxA* are indicated by cells  
 728 with the gene present. Using the species at the tips of the tree as independent data points, we  
 729 might spuriously conclude that *toxA* facilitates the evolution of pathogenicity: 85% of  
 730 pathogenic species carry *toxA* (6/7 species), compared to only 25% of non-pathogenic species  
 731 (1/4 species) (Chi-squared=4.05, p<0.05). However, this significant correlation is an artefact  
 732 of shared evolutionary history. Pathogenicity only evolved twice in the phylogeny: once in a  
 733 lineage with *toxA* and once in a lineage without *toxA*. Rather than independent data points, the  
 734 cluster of pathogenic, *toxA* carrying species are analogous to technical pseudoreplicates in an  
 735 empirical study.

736  
 737  
 738



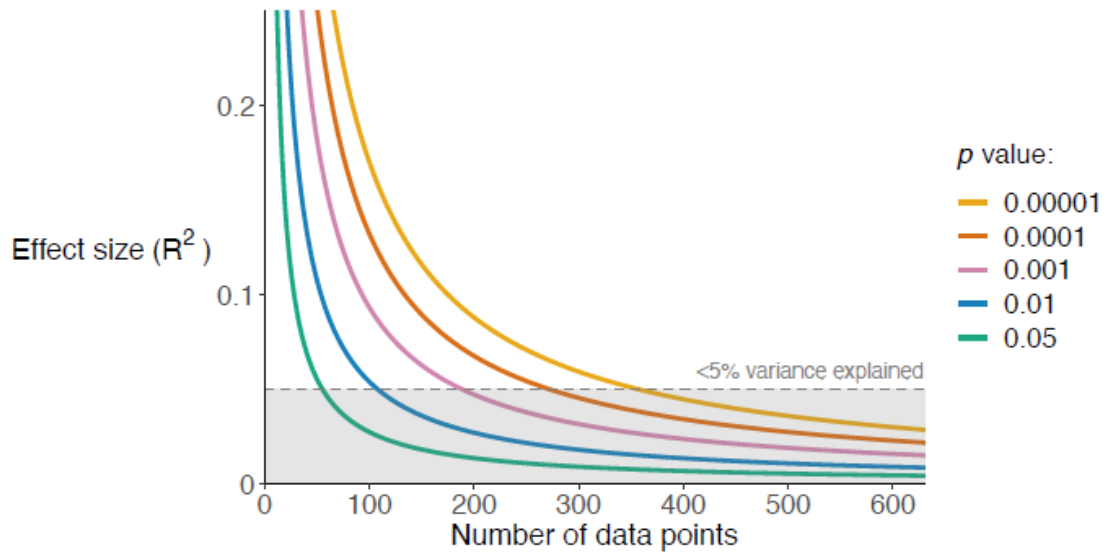
739  
 740 **Figure 2. Phylogenetic bias of genome sequencing.** Order level visualisation of the GTDB  
 741 bacterial phylogeny (v.214)<sup>134</sup>; the size of dots corresponds to the number of complete bacterial  
 742 genomes represented in both the RefSeq and the GTDB database from each taxonomic order.  
 743 Labels correspond to the five orders with the most genome sequences.



744

745 **Figure 3. Testing causal hypotheses.** A. Imagine that bacterial species which carried the gene  
 746 *toxB* were more likely to be pathogens. Can we test if causality is in that direction, with *toxB*  
 747 favouring the transition to pathogenicity? B. Yes, we can use transition rate methods<sup>33,42,51,85</sup>.  
 748 For the two binary traits pathogenicity and *toxB* presence, there are four possible states, each  
 749 represented by a cell. The quantity of evolutionary transitions between these states across the  
 750 phylogeny is indicated by arrows: larger arrows correspond to evolutionary changes which  
 751 occur more frequently. We can see that almost all transitions to pathogenicity (blue → red)  
 752 occur when the species already has the *toxB* gene. Most often, non-pathogens first acquire the  
 753 *toxB* gene, and then evolve pathogenicity, suggesting that *toxB* does help pathogenicity to  
 754 evolve.

755



756

757 **Figure 4. Statistical significance and biological importance.**

758 The effect size ( $R^2$ ; proportion of variance explained) required to produce a statistically  
 759 significant result decreases with the number of data points (sample size). For an unpaired t-  
 760 test, the x-axis is the number of data points in each of the two groups compared within the t-  
 761 test ( $N_1=N_2=N$ ) and the y-axis is the minimum  $R^2$  value which could be significant to one of  
 762 five  $p$ -values for a given  $N$ . The lines show the numerical relationship between the axes for  
 763 each of five  $p$ -values, which are all considered statistically significant. The area corresponding  
 764 to <5% of the variance explained is shaded in grey ( $R^2 < 0.05$ ). Larger datasets can detect  
 765 smaller effects, but very large datasets will assign almost all effects as significant, even when  
 766 they explain far less than 5% of the variance.

767

768 **Box 1. Case studies**

769 Controlling for phylogenetic non-independence has led to different conclusions from the same  
770 or very similar datasets. In this box, we showcase three recent studies. Another example is  
771 provided by two studies which used the same multi-species dataset but found opposite results  
772 when testing whether gene connectivity differed for genes on plasmids compared to on  
773 chromosomes<sup>65,135</sup>.

774

775 **[bH1] Plasmids and cooperation**

776 The growth of many bacteria depends on the secretion of molecules that provide a benefit to  
777 the local group of cells, a form of cooperation (public goods)<sup>136</sup>. It had been hypothesised that  
778 genes for cooperation would be favoured if they could be horizontally transferred on mobile  
779 elements such as plasmids to non-cooperating cells<sup>137,138</sup>. Consistent with this hypothesis, two  
780 comparative studies found that genes for public goods were more likely to be carried on  
781 plasmids<sup>75,102</sup>. One study compared the proportion of genes for secreted proteins on  
782 chromosomes and plasmids from 5,397 genomes, and the other study examined the proportion  
783 of genes for secreted proteins in plasmid compared to chromosomal genes in the pangenomes  
784 of 24 single or multi-species ‘clades’ of bacteria<sup>75,102</sup> (Fig. 4A.i.). These supplementary  
785 analyses assumed that species, genes and genomes represented independent data points.  
786 Inspired by these previous analyses, two recent studies carried out a phylogeny-based analysis  
787 across 51 and 146 species, respectively, which found that genes for the production of public  
788 goods were not more likely to be on plasmids<sup>22,23</sup>. Controlling for phylogeny changed the  
789 results from significant to non-significant, with the number of genomes per species and  
790 phylogeny able to explain 46% and 34% of the variation in the data, respectively. The  
791 significant results in the first two studies seem to be an artefact of these particular analyses  
792 being biased towards the most commonly sampled species<sup>22</sup>. Consequently, even when  
793 analysing similar datasets, the results depend on whether the number and non-independence of  
794 genomes is controlled for. Of note, recent theory has supported the conclusions of the  
795 phylogeny-based analyses, by predicting that cooperation is not appreciably favoured by  
796 horizontal gene transfer<sup>139</sup>.

797

798 **[bH1] Constraints on carbon-fixing enzymes**

799 During photosynthesis, the enzyme rubisco catalyses the fixation of atmospheric carbon  
800 dioxide into glucose. Although rubisco is the most abundant enzyme on earth, it is surprisingly

801 inefficient as a catalyst as it also catalyses a wasteful oxygenation reaction alongside the  
802 desired carboxylation reaction<sup>25,26,140,141</sup>. One hypothesis for its inefficiency is that inherent  
803 catalytic trade-offs force the oxygenase and carboxylase functions of rubisco to be tightly  
804 linked<sup>140–142</sup>. Early quantitative support for this hypothesis came from a study which looked at  
805 rubisco kinetic traits across 27 species and found evidence for a strong trade-off between CO<sub>2</sub>  
806 specificity and carboxylation turnover, suggesting that the more selective an enzyme is for  
807 carbon over oxygen, the slower it is at catalysing the carbon fixation reaction<sup>140</sup>. A later study  
808 examined the same question, with a much larger dataset across 304 diverse species and found  
809 further support for the trade-off hypothesis, with a strong, highly significant negative  
810 correlation between CO<sub>2</sub> specificity and carboxylase turnover<sup>141</sup>.

811         However, a re-analysis of the larger dataset using phylogenetic comparative methods  
812 reduced the strength of the significant negative correlation between CO<sub>2</sub> specificity and  
813 carboxylation turnover from 37.4% to just 2.2% variance explained — much less than  
814 previously thought, and much less than phylogeny, which explained 56.1% of the variation  
815 across all kinetic traits<sup>26</sup>. Additionally, many proposed trade-offs between other kinetic traits  
816 that had support from previous studies disappeared when phylogeny was accounted for<sup>26</sup>. The  
817 authors concluded that the kinetic traits have evolved largely independently of one another,  
818 meaning adaptation of rubisco has only been weakly constrained by catalytic trade-offs<sup>25,26</sup>.  
819 This case study provides a clear example of how more data will not alleviate the problem that  
820 species cannot be considered as independent data points.

821

## 822 **[bH1] Gene expression in orthologs and paralogs**

823 The ‘ortholog conjecture’ states that genes which diverged via a speciation event (orthologs)  
824 should remain more similar compared to genes which diverged via a duplication event  
825 (paralogs)<sup>143</sup>. In support of this, one study found that the tissue specificity of ortholog gene  
826 expression was more similar than for paralogs across several datasets<sup>144</sup>.

827

828 However, another study re-analysed this data, examining gene expression data from six organs  
829 across eight animal species, using phylogenetic comparative methods<sup>145</sup>. This second study  
830 found that, when controlling for evolutionary history and the time since genes diverged, there  
831 was no difference between the similarity of genes which had diverged by speciation compared  
832 to duplication<sup>145</sup>. Instead, differences or similarities in gene expression were better explained  
833 by phylogenetic distance (i.e. how long since the genes had diverged), rather than whether they  
834 had diverged via speciation or duplication<sup>145</sup>.