1              Research article

## 2 Special care is needed in applying phylogenetic comparative methods to

## 3 gene trees with speciation and duplication nodes

4

5 Tina Begum[1,2], Marc Robinson-Rechavi[1,2]*

6

7

8 [1]Department of Ecology and Evolution, University of Lausanne, Lausanne, Switzerland

9 [2]SIB Swiss Institute of Bioinformatics, Lausanne, Switzerland

10

11 *Corresponding author

12 E-mail: marc.robinson-rechavi@unil.ch

13

## 14  **Abstract**

15  How gene function evolves is a central question of evolutionary biology. It can be

16  investigated by comparing functional genomics results between species and between genes.

17  Most comparative studies of functional genomics have used pairwise comparisons. Yet it

18  has been shown that this can provide biased results, since genes, like species, are

19  phylogenetically related. Phylogenetic comparative methods should allow to correct for

20  this, but they depend on strong assumptions, including unbiased tree estimates relative to

21  the hypothesis being tested. Such methods have recently been used to test the "ortholog

22  conjecture", the hypothesis that functional evolution is faster in paralogs than in orthologs.

23  Whereas pairwise comparisons of tissue specificity ($\tau$) provided support for the ortholog

24  conjecture, phylogenetic independent contrasts did not. Our reanalysis on the same gene

25  trees identified problems with the time calibration of duplication nodes. We find that the

26  gene trees used suffer from important biases, due to the inclusion of trees with no

27  duplication nodes, to the relative age of speciations and duplications, to systematic

28  differences in branch lengths, and to non-Brownian motion of tissue-specificity on many

29  trees. We find that incorrect implementation of phylogenetic method in empirical gene

30  trees with duplications can be problematic. Controlling for biases allows to successfully

31  use phylogenetic methods to study the evolution of gene function, and provides some

32  support for the ortholog conjecture using three different phylogenetic approaches.

33  Keywords: ortholog; paralog; gene expression; phylogenetic comparative methods;

34  Brownian; Ornstein-Uhlenbeck

35

2

## Introduction

36

37 The "ortholog conjecture", a standard model of phylogenomics, has become a topic of

38 debate in recent years (Koonin 2005; Studer and Robinson-Rechavi 2009; Nehrt et al. 2011;

39 Altenhoff et al. 2012; Chen and Zhang 2012; Gabaldón and Koonin 2013; Rogozin et al.

40 2014; Kryuchkova-Mostacci and Robinson-Rechavi 2016; Dunn et al. 2018; Stamboulian

41 et al. 2020). The ortholog conjecture is routinely used by both experimental and

42 computational biologists in predicting or understanding gene function. According to this

43 model, orthologs (i.e. homologous genes which diverged by a speciation event) retain

44 equivalent or very similar functions, whereas paralogs (i.e. homologous genes which

45 diverged by a duplication event) share less similar functions (Studer and Robinson-Rechavi

46 2009). This is linked to the hypothesis that paralogs evolve more rapidly. This hypothesis

47 was challenged by results suggesting that paralogs would be functionally more similar than

48 orthologs (Nehrt et al. 2011). Such findings not only raised questions on the evolutionary

49 role of gene duplication but also questioned the reliability of using orthologs to annotate

50 unknown gene functions in different species (Sonnhammer et al. 2014). Several studies

51 (Altenhoff et al. 2012; Chen and Zhang 2012; Rogozin et al. 2014; Kryuchkova-Mostacci

52 and Robinson-Rechavi 2016) later found support for the ortholog conjecture, mostly based

53 on comparisons of gene expression data.

54 While all previous studies of the ortholog conjecture had used pairwise comparisons of

55 orthologs and paralogs, a recent article suggested that this was flawed, and that

56 phylogenetic comparative methods should be used (Dunn et al. 2018). Phylogenetic

57 structure can violate the fundamental assumption of independent observations in statistics,

58 and thus ignoring it can lead to mistakes (Felsenstein 1985). A solution is to use phylogeny-

59    based methods. Phylogenetic Independent Contrast (PIC) (Felsenstein 1985), and

60    Phylogenetic Generalized Least-Square (PGLS) (Martins and Hansen 1997; Grafen 1989;

61    Rohlf 2001) are the most commonly used phylogenetic comparative methods. They were

62    developed under a purely neutral model of evolution, i.e. Brownian motion (BM). Such

63    Brownian process have been extended using a maximum likelihood approach, to allow for

64    different rates of evolution on different branches of a phylogeny (O'Meara et al. 2006;

65    Thomas et al. 2006), and to include stabilizing selection in which the trait is shifted towards

66    a single fitness optimum, or multiple different adaptive optima (i.e. "Ornstein-Uhlenbeck"

67    or OU process) (Hansen 1997; Butler and King 2004; Beaulieu et al. 2012). These

68    phylogenetic data modeling with different modes of trait evolution (e.g. BM, OU) require

69    a priori knowledge of different states on the tree. Other approaches implemented a Markov

70    chain Monte Carlo (MCMC) sampling in a Bayesian framework to accurately estimate the

71    number, location, and magnitude of shifts in evolutionary rates, or in optimal trait values

72    without a priori assignment of states (Eastman et al. 2011; Pennell et al. 2014; Uyeda and

73    Harmon 2014; Catalan et al. 2019). Bayesian approaches are time consuming, while OU

74    modeling with phylogenetic lasso algorithm allows a faster detection of directional

75    selection due to a shift in optimal trait value (Khabbazian et al. 2016). Moreover, OU has

76    been used to model gene expression evolution (Rohlfs and Nielsen 2015; Chen et al. 2019).

77    Among all the phylogenetic methods, PIC is widely adopted for its relative simplicity, and

78    its applicability to a wide range of statistical procedures (Cooper et al. 2016a; Dunn et al.

79    2018). The performance of PIC relies on three basic assumptions: a correct tree topology;

80    accurate branch lengths; and trait evolution following Brownian motion (where trait

81    variance accrues as a linear function of time) (Felsenstein 1985; Garland 1992; Garland et

4

82    al. 1992; Díaz-Uriarte and Garland 1998; Freckleton and Harvey 2006; Cooper et al.

83    2016a). If any of these assumptions is incorrect, this can lead to incorrect interpretation of

84    results without control for biases (Diaz-Uriarte and Garland 1996; Díaz-Uriarte and

85    Garland 1998). While previous applications of PIC used multivariate traits on pure

86    speciation trees to explore the relationship between them, Dunn et al. (2018) took an

87    innovative approach in applying PIC to compare the divergence rates of a univariate trait

88    between two different node events ("speciation" and "duplication"), to test the ortholog

89    conjecture. They performed extensive analyses in support of their results. However, such

90    an application might be problematic since the time of occurrence of gene duplication, one

91    of the two types of events compared, is unknowable by external information (e.g. no fossil

92    evidence). Therefore, further study is required to understand why Dunn et al. (2018)

93    obtained results which are inconsistent with previous studies. It is possible that all the

94    conclusions drawn by previous studies on gene duplication are incorrect due to overlooking

95    phylogenetic tree structure. If so, it should be well supported.

96     We re-examined the data of Dunn et al., after reproducing their results using the resources

97    and scripts provided by the authors (Dunn et al. 2018). We have uncovered problems with

98    the use of PIC on biased calibrated gene trees, violation of the underlying assumptions, and

99    the inclusion of pure speciation gene trees. We used PIC on gene trees after fixing

100   calibration bias for old duplication nodes. With proper controls, the phylogenetic method

101   supports the ortholog conjecture. To verify this result, we also applied data modeling

102   approaches using a maximum likelihood framework, and using a reversible-jump Bayesian

103   MCMC algorithm. Support for the ortholog conjecture still holds with proper controls.

104

5

## Results

**Issues with straightforward application of Phylogenetic Independent Contrasts (PICs)**

Dunn et al. (2018) have made a relevant argument that the test should be done in a phylogenetic framework, since closely related species or genes tend to share more similar traits. They applied PIC method to a processed dataset of 8520 time-calibrated trees (details in the Materials and Methods, Table 1) by assuming that the computed node contrasts (PICs) are always phylogenetically independent, and reported evidence in contradiction with the ortholog conjecture for tissue-specificity $\tau$ (median: $\text{PIC}_{\text{speciation}}$ = 0.0072, $\text{PIC}_{\text{duplication}}$ = 0.0051, one-sided Wilcoxon test $P$ = 1). Yet the same data supported the ortholog conjecture when analyzed by pairwise comparisons, both in Kryuchkova-Mostacci and Robinson-Rechavi (2016), and in the re-analysis by Dunn et al. (2018). To understand the incongruence between results of PIC and of pairwise comparison approaches, they performed simulations of $\tau$ on their trees under the OC (ortholog conjecture), and under a null of uniform Brownian motion. PICs and pairwise comparisons have different expectations under the null ($\sigma^2_{\text{duplication}} = \sigma^2_{\text{speciation}}$) and under the ortholog conjecture ($\sigma^2_{\text{duplication}} > \sigma^2_{\text{speciation}}$) (Supplementary fig. S1). The simulation results of Dunn et al. (2018) indicated that the pairwise comparisons of events could not distinguish the two scenarios (null and OC), unlike the PIC method. As the result on their empirical data resembled their null simulation result, they questioned both the use of pairwise comparisons, and the support for the ortholog conjecture from tissue specificity data.

126   To understand their results, we first reproduced and reanalyzed the data of Dunn et al.

127   (2018) by focusing on the phylogenetic approach. Dunn et al. reported a non-significant

128   result ($P = 1$) for the PIC under the null simulation as well as for the empirical data, using

129   a Wilcoxon one-tailed rank test to check if the contrasts of duplication events are higher

130   than the contrasts of speciation events. Surprisingly, our reanalysis with a Wilcoxon two-

131   tailed rank test on the same data shows that the PIC rejects the null hypothesis on the null

132   simulations (Fig. 1A), with significant support for higher contrasts after speciation than

133   duplication. This means that the PIC method supports a trend opposite to the trend expected

134   under the ortholog conjecture in a null simulation. This was robust to repeating the

135   simulations with different random seed number (Supplementary fig. S2). This indicates

136   that neither of the approaches, PIC or pairwise, worked properly for these calibrated trees,

137   since both the approaches reject the null hypothesis when simulations are performed under

138   the null. Similarly, when we used a Wilcoxon two-tailed rank test instead of a one-tailed

139   test on the empirical data, the non-significant result ($P = 1$) (Dunn et al. 2018) was also

140   significant ($P < 2.2e^{-16}$) in the same unexpected direction as the null simulation results.

141   Statistical non-independence among species trait values because of their phylogenetic

142   relatedness can be measured by phylogenetic signal (Pagel 1999; Freckleton et al. 2002;

143   Blomberg et al. 2003; Münkemüller et al. 2012; Molina-Venegas and Rodríguez 2017).

144   Use of the PIC is mainly important for the data sets with strong phylogenetic signal, where

145   it allows to recover phylogenetically independence. Dunn et al. (2018) used Blomberg's

146   K. Its value ranges from 0 to ∞ for each tree, where a value of 0 indicates no phylogenetic

147   signal for the trait studied, and a value close to 1 or higher indicates strong phylogenetic

148   signal (Pagel 1999; Freckleton et al. 2002; Blomberg et al. 2003; Münkemüller et al. 2012;

7

149     Molina-Venegas and Rodríguez 2017). With a cutoff of K > 0.551, Dunn et al. (2018)

150     obtained only 2082 trees (Table 1), 24.4% of the total, with strong phylogenetic signal. The

151     phylogenetic method still rejects the null hypothesis under null simulations for those 2082

152     trees using a Wilcoxon two-tailed rank test (Fig. 1B), showing that the problem is not

153     simply due to low phylogenetic signal. Using a cut-off of $P < 0.05$ together with K > 0.551

154     leads to 1135 statistically significant trees with strong phylogenetic signals, for which we

155     obtained a similar result (Supplementary fig. S3). This means that the bias is not limited to

156     the selection of tree sets or to the number of speciation or duplication events used for the

157     analyses. Since the trend was similar for these 1135 trees we continued analyses with the

158     2082 trees of Dunn et al. (2018) for consistency.

159     The accuracy and performance of the PIC method largely depend on proper branch length

160     calibration in absolute time (e.g. in Million Years – My) (Garland 1992; Díaz-Uriarte and

161     Garland 1998; Cooper et al. 2016a). We thus investigated possible biases created during

162     calibration of gene trees. Due to non-availability of external references for duplication time

163     points (e.g. no fossils), Dunn et al. (2018) used only 7 speciation time points to calibrate

164     gene trees. The ages of other node events are estimated using the penalized likelihood

165     method (Sanderson 2002) by the chronos() function of the "ape" R package (Paradis et al.

166     2004), and varies for the same duplication clade labels even within the same gene trees.

167     The oldest speciation age for their calibrated trees was 296 My (Table 1), corresponding to

168     the use of chicken as the outgroup. Surprisingly, the calibrated node age of the oldest

169     duplication event was 11799977 My (Table 1, Supplementary table S1), that is, 2600 times

170     older than the Earth. This is indicative of issues with calibration. The tree pruning to species

171     with $\tau$ data (details in Materials and Methods) lead to trees for which all nodes older than

172    296 My are duplication or NA events, even if there were older speciation events present

173    before pruning (Supplementary fig. S4A). If the root node of a pruned tree is a speciation,

174    the duplication ages are constrained by speciation ages. Otherwise, there are no constraints

175    for the duplication events older than the oldest speciation events (Supplementary fig. S4,

176    Supplementary Table S1), which can introduce a calibration bias. This unreliable branch

177    length estimation for the old duplication nodes eventually led to much larger expected

178    variances for gene duplication events than for speciation events (Supplementary figs. S5A

179    and S5B).

180    PIC of a node is a ratio of changes in trait values ($\tau$ here) for descendant nodes to their

181    expected variance, i.e. the lengths of the two branches that connect the node to its two

182    descendants. This means that similar changes in $\tau$ for two nodes can produce different PIC

183    values, with the lower contrast for the node with higher expected variance (i.e., calibrated

184    branch length). In the null simulations only the $\tau$ values are simulated, while the branch

185    lengths (hence the expected variances) are taken from the empirical data, and thus share its

186    biases. This explains why contrasts are lower for duplications than for speciations under

187    null simulations as well as with empirical data. Such calibration bias in branch lengths

188    violates the second assumption of PIC applicability, and inflates type I error rates (Diaz-

189    Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998).

190

191    **Randomization tests to assess the performance of phylogenetic method**

192    We used randomization tests to assess bias in different analyses of the empirical dataset.

193    Our expectation is that the trend of the empirical result should differ from the randomized

9

194    ones. In a first randomization test, we permuted the $\tau$ values across the tips of each tree

195    without altering the node events of the trees. By such randomization, the real phylogenetic

196    relationships between trait values are removed for each tree. When we compared the node

197    contrasts of the speciation and duplication events computed based on these 8420

198    randomized $\tau$ trees (Fig. 2A), we found the same pattern as reported for the empirical gene

199    trees by Dunn et al. (2018), contrary to expectation. It confirms that results are driven by

200    their large differences in branch lengths (i.e. in expected variances) (Fig. 2B), as on

201    simulated null data. Any effect of trait divergence rates of speciation and duplication events

202    is always masked by this branch length difference of node events. This violates the basic

203    assumption of applicability of the PIC method to Brownian trait evolution. To remove the

204    problem of difference in expected variances of the two events, we performed a second

205    randomization test: we kept the original $\tau$ value for tips but randomly shuffled the events

206    (duplication, speciation, or NA) of internal nodes of the 8420 empirical gene trees to

207    maintain the original proportions of speciation and duplication events. The resulting trend

208    (Fig. 2C) still resembled the empirical gene trees data. This appears due to the fact that the

209    majority of the nodes are speciations (Fig. 2D, Table 1) with node ages $\leq$ 296 My. Most

210    of the trees with many duplication events on the other hand have ancient duplication events

211    for which the evolutionary rates of duplication are often masked by the effect of longer

212    branch lengths. Opposite to our expectation, the calibrated trees with no or few duplications

213    have higher overall nodes contrast (apparent fast evolution) than trees with many

214    duplications (apparent slow evolution). This might be due to greater difficulty in detecting

215    paralogs for fast evolving genes. Therefore, reshuffling of the events may not change the

216    observed pattern of higher speciation contrasts than duplication contrasts.

217    Out of 8520 calibrated trees, 2990 were pure speciation trees with no duplication events.

218    For these 2990 trees, random shuffling of events had no impact. To avoid this bias, we

219    removed those 2990 speciation trees as well as trees with negative branch lengths, and

220    randomized the trait or the internal node events 100 times on the remaining 5479 trees.

221    However, we still always obtained significantly higher contrasts of speciation than of

222    duplication (Supplementary figs. S6A and S6B). The randomization tests pattern is the

223    same when we used 2082 trees with strong phylogenetic signals (Supplementary figs. S6C

224    and S6D).

225    All these analyses indicate that the results reported by Dunn et al. (2018) are biased by the

226    calibrated phylogeny structures, and that this bias is not easy to correct. We propose three

227    approaches to correct for this bias and recover a proper phylogenetic signal of trait

228    evolution.

229

230    **Approach-1: PIC with  diagnostic tests**

231    Diagnostic tests (details in the Materials and Methods) for each tree are essential to ensure

232    phylogenetic independence of node contrasts, especially since there is evidence of bias in

233    the calibrated trees. This can be verified by the lack of correlation between the absolute

234    value of PICs of $\tau$ and their standard deviations, node height, node age, or node depth

235    (Garland 1992; Garland et al. 1992; Diaz-Uriarte and Garland 1996; Díaz-Uriarte and

236    Garland 1998; Freckleton 2000; Freckleton and Harvey 2006; Cooper et al. 2016a). A

237    statistically significant negative or positive correlation in any of the diagnostic tests

238    confirms that the PICs for that tree are non-independent (Garland 1992; Garland et al. 1992;

11

239    Diaz-Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998; Freckleton 2000;

240    Freckleton and Harvey 2006; Cooper et al. 2016a); in practice, we used $P < 0.05$ for

241    significance.

242    We performed such diagnostic tests on 4288 trees, for which calibration biases are fixed

243    for old duplication nodes (see Materials and Methods, Table 1). Among them only 2088

244    (48.7%), which includes 15321 speciation and 6213 duplication nodes, passed all 4

245    diagnostics tests for $\tau$ evolution. We performed our PIC analyses separately for 3948 young

246    ($\leq$ 296 My, the oldest speciation in the trees) and 2265 old ($>$ 296 My) duplication events.

247    Analyses on young duplicates after diagnostic tests provided support for the ortholog

248    conjecture (Fig. 3), but old duplicates did not. Randomization tests showed patterns distinct

249    from real data only for the young duplicates (Supplementary figs. S7A and S7B), indicating

250    a biological pattern rather than a data bias. Thus PIC on the trees after diagnostic plot tests

251    supports the ortholog conjecture for young duplicates, whereas the inference remains

252    biased for older duplicates.

253

254    **Approach-2: PIC with branch length transformation**

255    Most phylogenetic methods are developed for the Brownian model of trait evolution,

256    including the PIC method (Felsenstein 1985; Cornwell and Nakagawa 2017). Deviations

257    from pure BM violate the fundamental assumptions of PIC applicability and can affect its

258    performance for testing hypotheses about correlated evolution (Garland 1992; Garland et

259    al. 1992; Diaz-Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998). Using model-

260    fitting (see Materials and Methods), we found that 75.6% gene trees (Supplementary fig.

12

261    S8) supported the Ornstein-Uhlenbeck (OU) model. Remedial measures such as branch

262    length transformations along with diagnostic tests, can substantially recover the

263    performance of the PIC methods when character evolution is not BM or when contrasts are

264    non-independent of the phylogeny (Garland et al. 1992; Diaz-Uriarte and Garland 1996;

265    Díaz-Uriarte and Garland 1998).

266    We applied branch length transformation (details in the Materials and Methods) on all 4288

267    trees, along with diagnostic tests for consistency. We found substantial support for the

268    ortholog conjecture for the 4190 trees (97.7%) which pass diagnostic tests after branch

269    length transformation (Fig. 4A). Due to the lack of absolute age for these transformed trees,

270    we did not distinguish young and old duplicates. Applying such branch length

271    transformation then diagnostic tests to the gene trees of Dunn et al. we also found support

272    for the ortholog conjecture in 98.8% (8417 out of 8520) (Supplementary fig. S9A), as well

273    as for 99.9% (2080 out of 2082) of their trees with strong phylogenetic signal

274    (Supplementary fig. S10A). Randomization tests on all these sets of trees following branch

275    length transformations clearly showed distinct patterns compared to the empirical data

276    (Figs. 4B and 4C, Supplementary figs. S9B, S9C, S10B, S10C), indicating that results are

277    not due to inference bias once the data is properly transformed.

278    **Approach-3: Phylogenetic data modeling**

279    State dependent model-fitting allows to compare the evolutionary rates ($\sigma^2$), and the

280    changes in adaptive optimum value (θ) associated to specific states (speciation or

281    duplication) for each tree (Beaulieu et al. 2012; Clavel et al. 2015). Under the ortholog

282    conjecture, our expectation is that there should be more shifts in optimum value of $\tau$

283     between paralogs than orthologs. Moreover, the evolutionary rates after duplication should

284     be higher than after speciation ($\sigma^2_{\text{duplication}} > \sigma^2_{\text{speciation}}$). Of course, trends on empirical data

285     should differ from randomized ones. When we modeled the evolution of $\tau$ (see Materials

286     and Methods), 32 out of 4288 trees failed to fit any model due to invariance in $\tau$. Among

287     the others, 308 supported BM1, 704 BMM, 2874 OU1, and 370 OUM, as the best fit

288     models (Supplementary fig. S8). We performed our analyses separately for young and old

289     duplicates.

290     On the 8.6% multi optima trees (OUM) the optimum value are significantly higher for both

291     young and old duplications ($\theta_{\text{dup}} > \theta_{\text{spe}}$) (Supplementary Table S2). Thus paralogs regime

292     shift towards higher tissue-specificity. These results are not observed on randomized trees,

293     supporting a biological pattern in the data (Supplementary Table S2).

294     We also applied a Bayesian method (Udeya and Harmon 2014) on them to quantify the

295     number of adaptive optimum shifts, as suggested for small trees (Cooper et al. 2016b).

296     Unlike the other approach, such detection of evolutionary shifts in a phylogeny does not

297     need a priori knowledge of different states on the tree. Using a strict posterior probability

298     threshold of $\geq 0.7$ with this method, we find that most optimum shifts per branch for $\tau$

299     follow duplications (median after speciation: 0%, after duplication: 12.5%, paired two-

300     sided Wilcoxon rank-sum test $P < 2.2e^{-16}$). An OU model can often be incorrectly favored

301     over a BM model in a maximum likelihood framework when applied to trees with < 200

302     tips (Cooper et al. 2016b). Our gene trees have a median of only 15 tips. We thus applied

303     a conservative Bayesian approach on all of the 3244 trees for which OU was the preferred

304     model (OU1 + OUM). Even with such a strict posterior probability threshold of $\geq 0.7$, 1101

305     trees (33.9%) still supported the OUM model, including 901 trees identified as OU1 by

14

306    maximum likelihood. We detected the same trend of optimum shifts per branch (median

307    after speciation: 2.3%, after duplication: 10%, paired Wilcoxon rank-sum test $P < 2.2\text{e}^{-16}$).

308    These results are largely consistent for both young and old duplicates (Table 2;

309    Supplementary Table S3). However, the rates of optimum shifts are faster only for young

310    duplicates (Table 2; Supplementary Table S3).

311    Analyses on the trees where $\sigma^2$ varies between events (BMM) also supports the ortholog

312    conjecture for young duplicates (Table 3). Randomized data showed distinct patterns

313    from empirical data.  However, again there was neither support for the ortholog

314    conjecture nor signal relative to randomization for the old duplicates.

## Discussion

We agree with Dunn et al. (2018) that evolutionary comparisons should be done considering a phylogenetic framework when possible. However, this does not imply that phylogenetic methods can be applied easily to phylogenomics. To get a clear picture, we limited our study to the same gene trees used by Dunn et al. (2018). Our reanalysis identified problems generated by the time calibration of old duplication nodes of pruned trees, the inclusion of pure speciation gene trees, and violations of the Brownian model. The strongest bias was for duplication nodes preceding the oldest speciation nodes. This, in turn, introduced several biases in the analyses, and influenced results.

When we identified and controlled for such biases, PIC results changed to support the ortholog conjecture, consistent with our previous pairwise analysis (Kryuchkova-Mostacci and Robinson-Rechavi 2016) on the same $\tau$ data. Our fundamental point is that the conclusions drawn by Dunn et al, but also by anyone else who will have followed the same approach of applying PIC to gene trees, are not reliable unless extreme care is taken. This is because gene trees with orthologs and paralogs have more complex evolutionary histories, and different sampling biases, than species trees for which these methods were developed.

To date, a few studies have applied phylogenetic comparative methods to understand the effect of gene duplication on functional evolution (Oakley et al. 2005; Oakley et al. 2006; Eng et al. 2009; Rohlfs and Nielsen 2015; Dunn et al. 2018; Fukushima and Pollock 2020). None before Dunn et al. applied PIC method to compare speciation and duplication events on the same trees using a single continuous trait. Such application requires thorough testing

337  of the fundamental assumptions of the method on such time calibrated trees (Garland 1992;

338  Garland et al. 1992; Diaz-Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998;

339  Freckleton 2000; Freckleton and Harvey 2006; Cooper et al. 2016a). Hence, we explored

340  whether the application of a phylogenetic method might inflate errors (e.g. rejection of the

341  null hypothesis in null condition) if applied without assumption testing. Indeed, it is the

342  case (Figs. 1A and 1B). Along with the calibration bias for old duplication nodes, the

343  relative ages of the speciation and duplication events strongly differ in these trees due to

344  the choice of species. Using such trees without control for biases may bring about lack of

345  statistical power to detect the signal of ortholog conjecture, and even bias towards an

346  opposite pseudo-signal.

347  Time calibration of ancient duplication events is one of the major issues we uncovered.

348  The approach of Dunn et al. considered pruned trees with available trait ($\tau$ here) data for

349  time-calibration using speciation time points (see Materials and Methods). Such pruned

350  trees often have many duplication nodes older than the oldest speciation nodes. Sequence

351  based evolutionary rate (e.g., dN/dS) analyses in different species have found higher

352  sequence evolutionary rate following gene duplication (Conant and Wagner 2003; Kim

353  and Yi 2006; Scannell and Wolfe 2008; Han et al. 2009; Studer and Robinson-Rechavi

354  2009; Panchin et al. 2010; Pegueroles et al. 2013; Pich and Kondrashov 2014; Holland et

355  al. 2017). Therefore, calibration bias is not surprising for those duplication nodes in the

356  absence of time constraints (Supplementary figs. S4A-S4C, Supplementary Table S1).

357  Instead, we performed time calibration before pruning, so that the oldest speciation time

358  points can provide upper age limits and reduce calibration bias (Supplementary figs. S4D-

359  S4F). This is strongly recommended since the performance of the phylogenetic methods

17

360    rely on accurate branch length information, especially for multi-states univariate trait

361    analysis.

362    Dunn et al. (2018) performed several analyses (e.g. added random noise in the speciation

363    calibration time points, extended terminal branch length, removed old duplication nodes,

364    etc.) to take into account issues with branch lengths, but their simulations and our

365    randomization tests show that they appear not to have been sufficient to correct for this

366    bias (Figs. 2A and 2C). Dunn et al. also provided the hutan::picx() R function to compute

367    PIC for OU trees. In their simulation-based function, they estimated ancestral states by the

368    'GLS_OUS' method using the bias calibrated phylogeny. Therefore, their method does not

369    add anything specific to deal with the OU trees. Since they did not control for phylogenetic

370    independence of the contrasts, and did not consider the relative ages of the speciation and

371    old duplication events, they always obtained lower PIC of duplication events. Due to such

372    phylogenetic internal parameter dependence, their PIC analyses produced similar trends

373    with real or randomized data.

374    Assumptions of proper branch length information and of Brownian motion of trait

375    evolution are related, so that modifications of branch lengths can change the evolutionary

376    model (Diaz-Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998). Contrasting a

377    single rate OU to BM models, Dunn et al. (2018) identified 99.9% gene trees which favored

378    an OU model, more explicitly an OU1 model. This appears to be 67% when we performed

379    multivariate data modeling in a maximum likelihood framework on trees with less or no

380    calibration bias (Supplementary fig. S8). PIC analyses with diagnostic tests provided weak

381    support for the ortholog conjecture for the young duplicates (Figs. 3A-3C), in contrast to

382    previous results of Dunn et al. Small effect size difference in our inference is not surprising

18

383    since PIC is applied on OU trees. Similar patterns of results from empirical and

384    randomization tests for the old duplicates indicate that one should be extremely careful

385    before integrating them into a phylogenetic analysis. Branch length transformation

386    attempts to transform the OU trees to BM trees to meet the underlying assumption of

387    phylogenetic comparative method (Butler and King 2004). Hence, it can address the issue

388    of low power when underlying assumptions of phylogenetic methods are violated (Diaz-

389    Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998). Following this approach along

390    with the diagnostic tests, we obtained substantial support for the ortholog conjecture (Figs.

391    4A-4C, Supplementary figs. S9 and S10).

392    Phylogenetic data modeling also appears to be a powerful tool for such hypothesis testing,

393    where one can estimate the trait evolutionary rates or optima shift rates per event without

394    transforming OU trees to BM trees. More support for the OU trees (Supplementary fig. S8)

395    could be due to the fact that we performed multivariate evolutionary model-fitting mostly

396    on small trees (Cooper et al. 2016b). Among them only 8.6% trees supported the OUM

397    model. Following the recommendation of Cooper et al. (2016), we applied Bayesian

398    approach on small trees to accurately identify multi optima trees. Although previous studies

399    (Uyeda and Harmon 2014; Khabbazian et al. 2016; Uyeda et al. 2017) have suggested a

400    liberal cutoff of $\geq 0.2$ to detect an optimum shift with a Bayesian approach, we used a strict

401    posterior probability cutoff of $\geq 0.7$. We performed our analyses on the 33.9% OUM trees

402    passing such a strict posterior probability threshold. Our results from the PIC analyses with

403    controls was also supported by the maximum likelihood, and Bayesian data modeling

404    approaches. This shows that once proper precautions are taken, the empirical trends do not

405    depend on the number of selected gene trees or of internal node events included.

406    Empirical support for the ortholog conjecture has been mixed, with some studies

407    supporting it (Koonin 2005; Studer and Robinson-Rechavi 2009; Altenhoff et al. 2012;

408    Chen and Zhang 2012; Gabaldón and Koonin 2013; Rogozin et al. 2014; Kryuchkova-

409    Mostacci and Robinson-Rechavi 2016; Fukushima and Pollock 2020), and a few failing to

410    do so (Nehrt et al. 2011; Dunn et al. 2018; Stamboulian et al. 2020). Our results provide

411    additional support for the ortholog conjecture using tissue specificity data in a phylogenetic

412    framework after controlling for biases. Due to lack of detailed functional information,

413    many studies are still limited to gene expression data as a proxy of function. Recently,

414    using functional replaceability assay, experimental studies (Kachroo et al. 2015; Laurent

415    et al. 2020) have shown that orthologous genes can be swapped between essential yeast

416    genes and human, although this is rarely the case for all the members of expanded human

417    gene families (Laurent et al. 2020), validating one prediction of the ortholog conjecture.

418

## Materials and Methods

419

### Data reproducibility details

420

421    Our analyses are based on 21124 gene trees obtained from ENSEMBL Compara v.75

422    (Herrero et al. 2016) as used by Dunn et al. (2018). We used the same random seed number

423    as in Dunn et al. (2018) to reproduce the simulation results for reanalysis. All reproduced

424    data of Dunn et al. were stored in the "manuscript_dunn.RData" file

425    (https://doi.org/10.5281/zenodo.4003391). We used the results stored in the 'data' or

426    'phylo' slot of the trees for further analyses. To differentiate our own function from theirs

427    (Dunn et al. 2018), we renamed the original function script of Dunn et al. from

428    "functions.R" to "functions_Dunn.R". We made separate scripts for PIC analyses

429    ("Premanuscript_run_TMRR.R"), and for data modeling analyses ("Model_fitting.R").

430    Some of the analyses were time consuming, so we stored our outputs in

431    "Analyses_TMRR.RData", and in "Model_fitting_TMRR.Rdata" files

432    (https://doi.org/10.5281/zenodo.4003391), to load during analyses. All the details of

433    different functions are provided inside the scripts. We supply all the previously stored data

434    (to reduce computation time during reproduction of result) and function files including our

435    own ("functions_TM_new.R") with this manuscript. All scripts are available on GitHub:

436    https://github.com/tbegum/Testing_the_ortholog_conjecture.

### Fixing time calibration bias of duplication nodes

437

438    We first present the approach that Dunn et al. (2018) used, for clarity. When two speciation

439    nodes had the same label in the gene tree, Dunn et al. edited the more recent one to "NA"

440    rather than "speciation". Indeed the presence of the same clade names at different node

21

441 depths forces all the intervening branches to have length zero when the tree is time

442 calibrated, leading to failure of calibration (Dunn et al. 2018). For trait evolution, they

443 annotated the tips of these modified trees with precomputed tissue specificity data, $\tau$ from

444 8 vertebrate species (human, gorilla, chimpanzee, macaque, mouse, opossum, platypus, and

445 chicken) (from Kryuchkova-Mostacci and Robinson-Rechavi 2016). $\tau$ is a univariate index

446 between 0 and 1 that measures tissue-specificity of gene expression (Yanai et al. 2005): $\tau$

447 close to 1 indicates high tissue specificity, while close to 0 indicates more ubiquitous

448 expression. Here $\tau$ was computed across 6 tissues: brain, cerebellum, heart, kidney, liver,

449 and testis, based on the RNA-seq data of Brawand et al. (2011). Dunn et al. pruned the

450 gene trees to remove tips with missing $\tau$ data, and then time calibrated them using

451 speciation clade ages in the chronos() function with the 'correlated' model from the R

452 package "ape" (Paradis et al. 2004). The modified NA clades were not used for this

453 calibration. They used 7 speciation time points with a maximum age of 296 My. Thus they

454 obtained 8520 calibrated gene trees having at least 4 tips with non-null trait data (Table 1;

455 Supplementary figs. S4A-S4C). Among these trees, 2990 were pure speciation trees, which

456 includes 12919 speciation events, or 19% of all speciation nodes.

457 Relative to Dunn et al., we exchanged the order of pruning and time calibration steps, i.e.,

458 we first time calibrated the 21124 modified (i.e. with NA added) gene trees, followed by

459 pruning to have at least 4 tips with $\tau$ data. This makes use of all 32 available speciations

460 time points, and helps to limit the calibration bias of the old duplication events

461 (Supplementary figs. S4D-S4F). Calibration fails for some trees, and we obtained 7336

462 calibrated gene trees. The maximum node age of old duplication events is 1175.2 My for

463 these trees, as opposed to 11799977 My (older than the universe) for the trees obtained by

464   the original approach (Table 1, Supplementary table S1). Among these 7336 gene trees, we

465   kept 4288 which have at least 1 speciation and 1 duplication events; we removed 39 pure

466   duplication and 3009 pure speciation trees. This 4288 gene tree set is our basis for

467   evaluating phylogenetic methods' capacity to test the ortholog conjecture (Table 1): we

468   compare the evolutionary rates, $\sigma^2$, or PICs of speciation and duplication events of the same

469   genes.

**Model selection for $\tau$ evolution**

471   We followed a state dependent model-fitting approach to identify Brownian motion (BM)

472   or Ornstein-Uhlenbeck (OU) trees. We classified time-calibrated gene duplication nodes

473   as "young" ($\leq$ 296 My, the maximum speciation age) or "old" ($>$ 296 My) before model

474   fitting. We performed stochastic mapping of our gene trees by assigning discrete states

475   ("speciation", "young-duplication", "old-duplication", and "NA") to the branches based on

476   the corresponding ancestral node events using the simmap() function of the phytools R

477   package (Revell 2012). For each mapped tree, we fitted 4 different models of $\tau$ evolution

478   using maximum-likelihood: (i) BM1, a single Brownian motion rate of evolution (i.e.

479   $\sigma^2_{speciation} = \sigma^2_{young-duplication} = \sigma^2_{old-duplication}$), (ii) BMM, a BM with multiple rates of evolution

480   for different events (i.e. different $\sigma^2$ are allowed), (iii) OU1, a single optimum OU model

481   (i.e. $\theta_{speciation} = \theta_{young-duplication} = \theta_{old-duplication}$, $\sigma^2_{speciation} = \sigma^2_{young-duplication} = \sigma^2_{old-duplication}$, $\alpha$

482   $_{speciation} = \alpha_{young-duplication} = \alpha_{old-duplication}$), and (iv) OUM, a multi optimum OU model with

483   identical strength of selection and rate of drift acting on all selective regimes (i.e. like OU1

484   but $\theta_{speciation} \neq \theta_{young-duplication} \neq \theta_{old-duplication}$).

23

485    We used both the mvMORPH (Clavel et al. 2015), and OUwie (Beaulieu et al. 2012) R

486    packages to perform model-fitting. Sometimes the information contained within a tree is

487    insufficient with respect to the complexity of the fitted models. This can lead to poor model

488    choice by returning a log-likelihood that is suboptimal and may provide incorrect

489    estimation of one or more model parameters for that tree (Beaulieu et al. 2012). Hence, we

490    included the diagnostics (diagnostic=T or diagn=T) during model-fitting. The eigen values

491    of the Hessian matrix of the diagnostics indicate whether convergence of the model has

492    been achieved or whether the parameter estimates are reliable (Beaulieu et al. 2012). For

493    the BM1, BMM, OU1, and OUM models, we first fitted the model using mvMORPH for

494    each gene tree. If any of the model failed to converge for the tree or if the eigen values of

495    the Hessian matrix indicated that it was not reliable, we re-fitted that model using OUwie

496    to include it in model comparison. If still it failed, we removed that model for that tree. For

497    model comparisons on each gene tree, we calculated the Akaike weights (ω) for each fitted

498    model by means of the second order Akaike information criteria (AICc), which includes a

499    correction for small sample sizes (Akaike 1974; Burnham and Anderson 2002). The model

500    with highest ω was selected as the best-supported model of $\tau$ evolution for the tree

501    (Burnham and Anderson 2002; Gearty et al. 2018). We estimated model parameters for

502    each tree based on the best fit model.

503    **Bayesian modeling to detect phenotypic optimum shift**

504    Regime shifts, i.e. shifts of optimal $\tau$ values, in OU models were detected by a Bayesian

505    phylogenetic approach of the bayou R package (Uyeda and Harmon 2014). The reversible-

506    jump phylogenetic comparative approach was used to perform MCMC sampling of

507    locations, magnitudes and numbers of shifts in multiple-optima Ornstein–Uhlenbeck

508    models. We ran MCMC chains for 100000 generations, and the first 30% of samples were

509    dropped as burn-in. We used a strict threshold of posterior probability $\geq 0.7$ to detect an

510    adaptive shift at a given branch of the phylogeny. For each event ("speciation" or

511    "duplication"), we used a ratio of the number of optimum shifts to the number of branches

512    for that event to estimate the proportions of shifts in a phylogeny.

**Randomization test of $\tau$ values**

513

514    For each tree, we used $\tau$ data (column name "Tau" in each tree 'data' object) across the

515    tips to carry out our randomization test. To randomize we permuted the actual $\tau$ data

516    without altering internal node events. The pic() function of the "ape" package (Paradis et

517    al. 2004) was used to compute PIC of nodes for each tree using permuted $\tau$ of tips. For

518    each run, we compared the contrasts of speciation and duplication events of the whole set

519    of randomized trees to estimate difference in event contrasts based on Wilcoxon signed

520    rank test. For 100 runs, we repeated the above process 100 times to obtain a distribution

521    plot of 100 independent $P$ values. For our model-fitting approach, we used the same

522    empirical simmap trees with permuted $\tau$ data at the tips. We re-estimated the model

523    parameters of the randomized $\tau$ trees using the best fit model chosen for the corresponding

524    empirical gene trees.

**Randomization test of node events**

525

526    Some of the speciation nodes had daughters with same clade names in the gene trees we

527    used for our study. Dunn et al. changed such node events to "NA" to avoid problems during

528    time calibration of the trees. Such annotated node event information ("Speciation",

529    "Duplication", "NA") for each tree was available as "Event" in the tree 'data' slot. To

25

530    randomize, we permuted the internal node events (added as column name "event_new" in

531    the 'data' slot) by maintaining the actual proportion of events for each tree. Then, we used

532    the PIC of actual $\tau$ at tips to estimate contrasts difference between newly assigned

533    speciation and duplication node events by Wilcoxon rank tests. For 100 independent runs,

534    we repeated the same procedure to obtain 100 independent $P$ values. Since the internal

535    node events were changed after such randomization, we reclassified gene duplication nodes

536    as "young" or "old" on the event modified trees, and repainted the trees. We re-estimated

537    the model parameters for the discrete states of the randomized events trees using the best

538    fit model chosen for the corresponding empirical gene trees.

**Checking for contrasts standardization by diagnostic tests**

540    We used several additional diagnostic tests on those trees to identify adequate independent

541    nodes contrast standardization before drawing any inference by PIC method, as

542    recommended in several studies (Garland 1992; Diaz-Uriarte and Garland 1996; Díaz-

543    Uriarte and Garland 1998; Freckleton and Harvey 2006; Cooper et al. 2016a). The most

544    usual method for contrasts standardization is to check a correlation between the absolute

545    values of PICs and their expected standard deviations (i.e. square root of sum of branch

546    lengths) (Garland et al. 1992; Díaz-Uriarte and Garland 1998; Cooper et al. 2016a). Under

547    Brownian motion, there should be no correlation. This test and the correlation between the

548    absolute values of PICs and the logarithm of their node age are model diagnostic plot tests

549    in the caper ("Comparative Analyses of Phylogenetics and Evolution in R") package

550    (Purvis and Rambaut 1995; Cooper et al. 2016a; Orme 2018; R Core Team 2018). We used

551    both of them by using the "crunch" algorithm of the caper package, which implements the

552    methods originally provided in CAIC (Purvis and Rambaut 1995; Cooper et al. 2016a;

26

553    Orme 2018; R Core Team 2018). Correlation of node heights with absolute values of

554    contrasts or PICs has also been reported to be a reliable indicator of deviation from the

555    Brownian model (Freckleton and Harvey 2006). Hence, we computed node height for each

556    node in a tree using the ape package (Paradis et al. 2004). We also used the correlations of

557    node height and node depth to the absolute value of nodes contrasts to rule out significant

558    trend in any of the 4 tests. We used $P < 0.05$ to assess a significant correlation for the

559    diagnostic tests. A significant trend (positive or negative) indicates phylogenetic

560    dependence for that tree (Garland 1992; Garland et al. 1992; Díaz-Uriarte and Garland

561    1998; Freckleton and Harvey 2006; Cooper et al. 2016a), and we removed those trees from

562    our analysis. Contrast calculation on negative branch lengths is not desirable, so we

563    removed trees with negative branch lengths before applying the crunch() function. To

564    assure that nodes contrast standardization is independent of the phylogeny, we considered

565    sets of trees passing all 4 diagnostic tests for further analyses.

566    **Branch length transformation**

567    Transformation of branch lengths has been proposed to restore the performance of PIC

568    method when the true evolutionary model is not BM or is unknown, or when branch lengths

569    are in error (Garland et al. 1992; Diaz-Uriarte and Garland 1996; Díaz-Uriarte and Garland

570    1998). In such cases, branch lengths are transformed by raising a family power of branch

571    length ranging from 0 to 2 in intervals of 0.1, plus the $\log_{10}$ of the branch lengths (Diaz-

572    Uriarte and Garland 1996; Díaz-Uriarte and Garland 1998). For each transformation, the

573    program computes the correlation between the absolute value of the standardized contrasts

574    and their standard deviations until no significant correlation is obtained, to ensure adequate

575    independent contrasts standardization (Diaz-Uriarte and Garland 1996; Díaz-Uriarte and

27

576   Garland 1998). Finally, we excluded trees for which adequate contrasts standardization is

577   not achieved even after raising the branch length power to 2 (Diaz-Uriarte and Garland

578   1996; Díaz-Uriarte and Garland 1998).

579   **Details of other packages used in this study**

580   We used phylosig function() of the phytools package (Revell 2012) to identify trees with

581   phylogenetic signal ($P < 0.05$) using Blomberg's K (Blomberg et al. 2003; Münkemüller

582   et al. 2012; Revell 2012). Analyses and plotting were performed in R version 3.5.1 (R Core

583   Team 2018) using treeio (Guangchuang 2018), ggtree (Guangchuang et al. 2017), stringr

584   (Wickham 2019), digest (Antoine Lucas et al. 2018), dplyr (Wickham et al. 2017),

585   tidyverse (Wickham 2017), ggrepel (Slowikowski 2018), gtools (Warnes et al. 2018),

586   ggplot2 (Wickham 2016), cowplot (Wilke 2019), easyGgplot2 (Kassambara 2014),

587   gridExtra (Auguie 2017), and png (Urbanek 2013) libraries.

588

589   **Acknowledgements**

# References

Akaike H. 1974. New look at statistical-model identification. Automatic Control, IEEE Transactions on 19:716–723.

Altenhoff AM, Studer RA, Robinson-Rechavi M, Dessimoz C. 2012. Resolving the ortholog conjecture: Orthologs tend to be weakly, but significantly, more similar in function than paralogs. PLoS Comput Biol 8:e1002514. Available from: https://www.ncbi.nlm.nih.gov/pubmed/22615551

Antoine Lucas DE with contributions by, Tuszynski J, Bengtsson H, Urbanek S, Frasca M, Lewis B, Stokely M, Muehleisen H, Murdoch D, Hester J, et al. 2018. Digest: Create compact hash digests of r objects. Available from: https://CRAN.R-project.org/package=digest

Auguie B. 2017. GridExtra: Miscellaneous functions for "grid" graphics. Available from: https://CRAN.R-project.org/package=gridExtra

Beaulieu JM, Jhwueng DC, Boettiger C, O'Meara BC. 2012. Modeling stabilizing selection: Expanding the Ornstein-Uhlenbeck model of adaptive evolution. Evolution 66:2369–2383. Available from: https://www.ncbi.nlm.nih.gov/pubmed/22834738

Benjamini Y, Yekutieli D. 2005. Quantitative trait loci analysis using the false discovery rate. Genetics 171:783–790. Available from: https://www.ncbi.nlm.nih.gov/pubmed/15956674

616    Blomberg SP, Garland T, Ives AR. 2003. Testing for phylogenetic signal in comparative

617    data: Behavioral traits are more labile. Evolution 57:717–745. Available from:

618    https://www.ncbi.nlm.nih.gov/pubmed/12778543

619    Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti

620    A, Aximu-Petri A, Kircher M, et al. 2011. The evolution of gene expression levels in

621    mammalian       organs.       Nature       478:343–348.       Available       from:

622    https://www.ncbi.nlm.nih.gov/pubmed/22012392

623    Burnham K, Anderson D. 2002. Model selection and multimodel inference. In: Springer,

624    New York,

625    Butler MA, King AA. 2004. Phylogenetic comparative analysis: a modeling approach for

626    adaptive       evolution.       Am       Nat       164:683-695.       Available       from:

627    https://www.ncbi.nlm.nih.gov/pubmed/29641928

628    Catalán A, Briscoe AD, Höhna, S. 2019. Drift and directional selection are the evolutionary

629    forces driving gene expression divergence in eye and brain tissue of. Genetics 213:581-

630    594. Available from: https://www.ncbi.nlm.nih.gov/pubmed/31467133

631    Chen J, Swofford R, Johnson J, Cummings BB, Rogel N, Lindblad-Toh K, Haerty W,

632    Palma FD, Regev A. 2019. A quantitative framework for characterizing the evolutionary

633    history of mammalian gene expression. Genome Res 29:53–63. Available from:

634    https://www.ncbi.nlm.nih.gov/pubmed/30552105

635  Chen X, Zhang J. 2012. The ortholog conjecture is untestable by the current gene ontology

636  but is supported by rna sequencing data. PLoS Comput Biol 8:e1002784. Available from:

637  https://www.ncbi.nlm.nih.gov/pubmed/23209392

638  Clavel J, Escarguel G, Merceron G. 2015. mvMORPH: An R package for fitting

639  multivariate evolutionary models to morphometric data. Methods Ecol Evol 6:1311–1319.

640  Conant GC, Wagner A. 2003. Asymmetric sequence divergence of duplicate genes.

641  Genome      Res      13:2052–2058.      Available      from:

642  https://www.ncbi.nlm.nih.gov/pubmed/12952876

643  Cooper N, Thomas GH, FitzJohn RG. 2016a. Shedding light on the 'dark side' of

644  phylogenetic comparative methods. Methods Ecol Evol 7:693–699. Available from:

645  https://www.ncbi.nlm.nih.gov/pubmed/27499839

646  Cooper N, Thomas GH, Venditti C, Meade A, Freckleton RP. 2016b. A cautionary note on

647  the use of ornstein uhlenbeck models in macroevolutionary studies. Biol J Linn Soc Lond

648  118:64–77. Available from: https://www.ncbi.nlm.nih.gov/pubmed/27478249

649  Cornwell W, Nakagawa S. 2017. Phylogenetic comparative methods. Curr Biol 27:R333–

650  R336. Available from: https://www.ncbi.nlm.nih.gov/pubmed/28486113

651  Diaz-Uriarte R, Garland T. 1996. Testing hypotheses of correlated evolution using

652  phylogenetically independent contrasts: Sensitivity to deviations from brownian motion.

653  Syst Biol 45:27–47.

31

654  Díaz-Uriarte R, Garland T. 1998. Effects of branch length errors on the performance of

655  phylogenetically independent contrasts. Syst Biol 47:654–672. Available from:

656  https://www.ncbi.nlm.nih.gov/pubmed/12066309

657  Dunn CW, Zapata F, Munro C, Siebert S, Hejnol A. 2018. Pairwise comparisons across

658  species are problematic when analyzing functional genomic data. Proc Natl Acad Sci U S

659  A 115:E409–E417. Available from: https://www.ncbi.nlm.nih.gov/pubmed/29301966

660  Eastman JM, Alfaro ME, Joyce P, Hipp AL, Harmon LJ. 2011. A novel comparative

661  method for identifying shifts in the rate of character evolution on trees. Evolution 65:3578–

662  3589. Available from: https://www.ncbi.nlm.nih.gov/pubmed/22133227

663  Eng KH, Bravo HC, Keleş S. 2009. A phylogenetic mixture model for the evolution of

664  gene      expression.      Mol      Biol      Evol      26:2363–2372.      Available      from:

665  https://www.ncbi.nlm.nih.gov/pubmed/19602540

666  Felsenstein J. 1985. Phylogenies and the comparative method. Am Nat 125:1–15.

667  Available from: https://www.jstor.org/stable/2461605

668  Freckleton R. 2000. Phylogenetic tests of ecological and evolutionary hypotheses:

669  Checking for phylogenetic independence. Func Ecol 14:129–134.

670  Freckleton R, Harvey P, Pagel M. 2002. Phylogenetic analysis and comparative data: A

671  test and review of evidence. Am Nat 160:712–726.

672  Freckleton RP, Harvey PH. 2006. Detecting non-brownian trait evolution in adaptive

673  radiations.          PLoS          Biol          4:e373.          Available          from:

674  https://www.ncbi.nlm.nih.gov/pubmed/17090217

675  Fukushima K, Pollock DD. 2020. Organ-specific propensity drives patterns of gene

676  expression evolution. BioRxiv. doi: https://doi.org/10.1101/409888

677  Gabaldón T, Koonin EV. 2013. Functional and evolutionary implications of gene

678  orthology.      Nat       Rev       Genet       14:360–366.      Available      from:

679  https://www.ncbi.nlm.nih.gov/pubmed/23552219

680  Garland T. 1992. Rate tests for phenotypic evolution using phylogenetically independent

681  contrasts.       Am        Nat        140:509–519.        Available       from:

682  https://www.ncbi.nlm.nih.gov/pubmed/19426053

683  Garland TJ, Harvey P, Ives A. 1992. Procedure for the analysis of comparative data using

684  phylogenetically independent contrasts. Syst Biol 41:18–32.

685  Gearty W, McClain CR, Payne JL. 2018. Energetic tradeoffs control the size distribution

686  of aquatic mammals. Proc Natl Acad Sci U S A 115:4194–4199. Available from:

687  https://www.ncbi.nlm.nih.gov/pubmed/29581289

688  Grafen A. 1989. The phylogenetic regression. Philos Trans R Soc Lond B Biol Sci

689  326:119–157. Available from: https://www.ncbi.nlm.nih.gov/pubmed/2575770

690  Guangchuang Y. 2018. Treeio: Base classes and functions for phylogenetic tree input and

691  output. Available from: https://guangchuangyu.github.io/software/treeio

692  Guangchuang Y, David S, Huachen Z, Yi G, Tommy T-YL. 2017. Ggtree: An R package

693  for visualization and annotation of phylogenetic trees with their covariates and other

694  associated data. Methods Ecol Evol 8:28–36.

695    Han MV, Demuth JP, McGrath CL, Casola C, Hahn MW. 2009. Adaptive evolution of

696    young gene duplicates in mammals. Genome Res 19:859–867. Available from:

697    https://www.ncbi.nlm.nih.gov/pubmed/19411603

698    Hansen TF. 1997. Stabilizing selection and the comparative analysis of adaptation.

699    Evolution                51:1341–1351.              Available              from:

700    https://www.ncbi.nlm.nih.gov/pubmed/28568616

701    Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle

702    SM, Amode R, Brent S, et al. 2016. Ensembl comparative genomics resources. Database.

703    Available from: https://www.ncbi.nlm.nih.gov/pubmed/27141089

704    Hochberg Y, Benjamini Y. 1990. More powerful procedures for multiple significance

705    testing.          Stat          Med          9:811–818.          Available          from:

706    https://www.ncbi.nlm.nih.gov/pubmed/2218183

707    Holland PW, Marlétaz F, Maeso I, Dunwell TL, Paps J. 2017. New genes from old:

708    Asymmetric divergence of gene duplicates and the evolution of development. Philos Trans

709    R        Soc        Lond        B        Biol        Sci        372.        Available        from:

710    https://www.ncbi.nlm.nih.gov/pubmed/27994121

711    Kachroo AH, Laurent JM, Yellman CM, Meyer AG, Wilke CO, Marcotte EM. 2015.

712    Evolution. systematic humanization of yeast genes reveals conserved functions and genetic

713    modularity.          Science          348:921–925.          Available          from:

714    https://www.ncbi.nlm.nih.gov/pubmed/25999509

715    Kassambara A. 2014. EasyGgplot2: Perform and customize easily a plot with ggplot2.

716    Available from: htt://www.sthda.com

717    Khabbazian M, Kriebel R, Rohe K, Ané C. 2016. Fast and accurate detection of

718    evolutionary shifts in ornstein-uhlenbeck models. Methods Ecol Evol 7:811–824.

719    Kim SH, Yi SV. 2006. Correlated asymmetry of sequence and functional divergence

720    between duplicate proteins of saccharomyces cerevisiae. Mol Biol Evol 23:1068–1075.

721    Available from: https://www.ncbi.nlm.nih.gov/pubmed/16510556

722    Koonin EV. 2005. Orthologs, paralogs, and evolutionary genomics. Annu Rev Genet

723    39:309–338. Available from: https://www.ncbi.nlm.nih.gov/pubmed/16285863

724    Kryuchkova-Mostacci N, Robinson-Rechavi M. 2016. Tissue-specificity of gene

725    expression diverges slowly between orthologs, and rapidly between paralogs. PLoS

726    Comput          Biol          12:e1005274.          Available          from:

727    https://www.ncbi.nlm.nih.gov/pubmed/28030541

728    Laurent JM, Garge RK, Teufel AI, Wilke CO, Kachroo AH, Marcotte EM. 2020.

729    Humanization of yeast genes with multiple human orthologs reveals functional divergence

730    between       paralogs.       PLoS       Biol       18:e3000627.       Available       from:

731    https://www.ncbi.nlm.nih.gov/pubmed/32421706

732    Martins E, Hansen T. 1997. Phylogenies and the comparative method: A general approach

733    to incorporating phylogenetic information into the analysis of interspecific data. Am Nat

734    149:646–667.

735    Molina-Venegas R, Rodríguez M. 2017. Revisiting phylogenetic signal; strong or

736    negligible impacts of polytomies and branch length information? BMC Evol Biol 17:53.

737    Available from: https://www.ncbi.nlm.nih.gov/pubmed/28201989

738    Münkemüller T, Lavergne S, Bzeznik B, Dray S, Jombart T, Schiffers K, Thuiller W. 2012.

739    How to measure and test phylogenetic signal. Methods Ecol Evol 3:743–756.

740    Nehrt NL, Clark WT, Radivojac P, Hahn MW. 2011. Testing the ortholog conjecture with

741    comparative functional genomic data from mammals. PLoS Comput Biol 7:e1002073.

742    Available from: https://www.ncbi.nlm.nih.gov/pubmed/21695233

743    Oakley TH, Gu Z, Abouheif E, Patel NH, Li WH. 2005. Comparative methods for the

744    analysis of gene-expression evolution: An example using yeast functional genomic data.

745    Mol        Biol        Evol        22:40–50.        Available        from:

746    https://www.ncbi.nlm.nih.gov/pubmed/15356281

747    Oakley TH, Ostman B, Wilson AC. 2006. Repression and loss of gene expression outpaces

748    activation and gain in recently duplicated fly genes. Proc Natl Acad Sci U S A 103:11637–

749    11641. Available from: https://www.ncbi.nlm.nih.gov/pubmed/16864793

750    O'Meara BC, Ané C, Sanderson MJ, Wainwright PC. 2006. Testing for different rates of

751    continuous trait evolution using likelihood. Evolution 60:922–933. Available from:

752    https://www.ncbi.nlm.nih.gov/pubmed/16817533

753    Orme D. 2018. The caper package: Comparative analysis of phylogenetics and evolution

754    in R. Available from: https://cran.r-project.org/web/packages/caper/vignettes/caper.pdf

755   Pagel M. 1999. Inferring the historical patterns of biological evolution. Nature 401:877–

756   884.

757   Panchin AY, Gelfand MS, Ramensky VE, Artamonova II. 2010. Asymmetric and non-

758   uniform evolution of recently duplicated human genes. Biol Direct 5:54. Available from:

759   https://www.ncbi.nlm.nih.gov/pubmed/20825637

760   Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of phylogenetics and evolution in

761   R        language.        Bioinformatics        20:289–290.        Available        from:

762   https://www.ncbi.nlm.nih.gov/pubmed/14734327

763   Pegueroles C, Laurie S, Albà MM. 2013. Accelerated evolution after gene duplication: A

764   time-dependent process affecting just one copy. Mol Biol Evol 30:1830–1842. Available

765   from: https://www.ncbi.nlm.nih.gov/pubmed/23625888

766   Pennell MW, Eastman JM, Slater GJ, Brown JW, Uyeda JC, FitzJohn RG, Alfaro ME,

767   Harmon LJ. 2014. Geiger v2.0: An expanded suite of methods for fitting

768   macroevolutionary models to phylogenetic trees. Bioinformatics 30:2216–2218. Available

769   from: https://www.ncbi.nlm.nih.gov/pubmed/24728855

770   Pich I Roselló O, Kondrashov FA. 2014. Long-term asymmetrical acceleration of protein

771   evolution after gene duplication. Genome Biol Evol 6:1949–1955. Available from:

772   https://www.ncbi.nlm.nih.gov/pubmed/25070510

773   Purvis A, Rambaut A. 1995. Comparative analysis by independent contrasts (caic): An

774   apple macintosh application for analysing comparative data. Comput Appl Biosci 11:247–

775   251. Available from: https://www.ncbi.nlm.nih.gov/pubmed/7583692

776    R Core Team. 2018. R: A language and environment for statistical computing. Vienna,

777    Austria: R Foundation for Statistical Computing Available from: https://www.R-

778    project.org/

779    Revell LJ. 2012. Phytools: An R package for phylogenetic comparative biology (and other

780    things). Methods Ecol Evol 3:217–223.

781    Rogozin IB, Managadze D, Shabalina SA, Koonin EV. 2014. Gene family level

782    comparative analysis of gene expression in mammals validates the ortholog conjecture.

783    Genome        Biol        Evol        6:754–762.        Available        from:

784    https://www.ncbi.nlm.nih.gov/pubmed/24610837

785    Rohlf FJ. 2001. Comparative methods for the analysis of continuous variables: Geometric

786    interpretations.        Evolution        55:2143–2160.        Available        from:

787    https://www.ncbi.nlm.nih.gov/pubmed/11794776

788    Rohlfs RV, Nielsen R. 2015. Phylogenetic ANOVA: The expression variance and

789    evolution model for quantitative trait evolution. Syst Biol 64:695–708. Available from:

790    https://www.ncbi.nlm.nih.gov/pubmed/26169525

791    Sanderson MJ. 2002. Estimating absolute rates of molecular evolution and divergence

792    times: A penalized likelihood approach. Mol Biol Evol 19:101–109. Available from:

793    https://www.ncbi.nlm.nih.gov/pubmed/11752195

794    Scannell DR, Wolfe KH. 2008. A burst of protein sequence evolution and a prolonged

795    period of asymmetric evolution follow gene duplication in yeast. Genome Res 18:137–147.

796    Available from: https://www.ncbi.nlm.nih.gov/pubmed/18025270

38

797      Slowikowski K. 2018. Ggrepel: Automatically position non-overlapping text labels with

798      'ggplot2'. Available from: https://CRAN.R-project.org/package=ggrepel

799      Sonnhammer EL, Gabaldón T, Sousa da Silva AW, Martin M, Robinson-Rechavi M,

800      Boeckmann B, Thomas PD, Dessimoz C, consortium Q for O. 2014. Big data and other

801      challenges in the quest for orthologs. Bioinformatics 30:2993–2998. Available from:

802      https://www.ncbi.nlm.nih.gov/pubmed/25064571

803      Stamboulian M, Guerrero RF, Hahn MW, Radivojac P. 2020. The ortholog conjecture

804      revisited: The value of orthologs and paralogs in function prediction. Bioinformatics

805      36:i219–i226. Available from: https://www.ncbi.nlm.nih.gov/pubmed/32657391

806      Studer RA, Robinson-Rechavi M. 2009. How confident can we be that orthologs are

807      similar, but paralogs differ? Trends Genet 25:210–216. Available from:

808      https://www.ncbi.nlm.nih.gov/pubmed/19368988

809      Thomas GH, Freckleton RP, Székely T. 2006. Comparative analyses of the influence of

810      developmental mode on phenotypic diversification rates in shorebirds. Proc Biol Sci

811      273:1619–1624. Available from: https://www.ncbi.nlm.nih.gov/pubmed/16769632

812      Urbanek S. 2013. Png: Read and write png images. Available from: https://CRAN.R-

813      project.org/package=png

814      Uyeda JC, Harmon LJ. 2014. A novel bayesian method for inferring and interpreting the

815      dynamics of adaptive landscapes from phylogenetic comparative data. Syst Biol 63:902–

816      918. Available from: https://www.ncbi.nlm.nih.gov/pubmed/25077513

817   Uyeda JC, Pennell MW, Miller ET, Maia R, McClain CR. 2017. The evolution of energetic

818   scaling across the vertebrate tree of life. Am Nat 190:185–199. Available from:

819   https://www.ncbi.nlm.nih.gov/pubmed/28731792

820   Warnes GR, Bolker B, Lumley T. 2018. Gtools: Various R programming tools. Available

821   from: https://CRAN.R-project.org/package=gtools

822   Wickham H. 2016. Ggplot2: Elegant graphics for data analysis. Springer-Verlag New York

823   Available from: https://ggplot2.tidyverse.org

824   Wickham H. 2017. Tidyverse: Easily install and load the 'tidyverse'. Available from:

825   https://CRAN.R-project.org/package=tidyverse

826   Wickham H. 2019. Stringr: Simple, consistent wrappers for common string operations.

827   Available from: https://CRAN.R-project.org/package=stringr

828   Wickham H, Francois R, Henry L, Müller K. 2017. dplyr: A grammar of data manipulation.

829   Available from: https://CRAN.R-project.org/package=dplyr

830   Wilke CO. 2019. Cowplot: Streamlined plot theme and plot annotations for 'ggplot2'.

831   Available from: https://CRAN.R-project.org/package=cowplot

832   Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, Ophir R, Bar-Even A, Horn-

833   Saban S, Safran M, Domany E, et al. 2005. Genome-wide midrange transcription profiles

834   reveal expression level relationships in human tissue specification. Bioinformatics 21:650–

835   659. Available from: <Go to ISI>://WOS:000227241200012

836

# Figure captions

**Figure 1**: **Reanalyses of phylogenetic simulation data of Dunn et al. (2018)**. *P* values are from Wilcoxon two-tailed tests. Values inside boxplots denote median PIC values of the corresponding events. In null simulations, there should be no difference in contrasts between events. In OC (Ortholog Conjecture) simulations, contrasts are expected to be higher for duplication than for speciation. (A) Higher contrasts for speciation than duplication reject the null hypothesis under null simulation scenario for all empirical time calibrated gene trees. (B) Results are similar with a subset of trees with strong phylogenetic signal for $\tau$.

**Figure 2: Analyses on calibrated empirical gene trees of Dunn et al. (2018).** *P* values are from Wilcoxon two-tailed tests. (A) Randomly shuffling the $\tau$ values of the tips for 8520 gene trees does not alter the empirical trend of an opposite trend to the ortholog conjecture. (B) The expected variance is much higher for duplication than speciation events irrespective of the number of tips considered for the study. (C) Using the original $\tau$ data, if we permute the events (Speciation or Duplication or NA) of the nodes, the trend of result remains. (D) The proportions of speciation events is much higher than duplication events for all time-calibrated trees; the dotted line represents the median proportion of both events; a high proportion of trees have no duplication events.

**Figure 3**: **The ortholog conjecture test on $\tau$ for trees passing diagnostic plot tests.** *P* values are from Wilcoxon two-tailed tests. Values inside boxplots denote median PIC values of the corresponding events. Young duplicates: age $\leq$ 296 My, the maximum speciation age; old duplicates: age > 296 My.

859   **Figure 4: The ortholog conjecture test for contrasts standardized branch transformed**

860   **trees.** *P* values are from Wilcoxon two-tailed tests. Values inside boxplots denote median

861   PIC value of the corresponding event. (A) Using 4190 out of 4288 calibrated trees that

862   passed diagnostic tests following branch length transformation. (B) Permuting τ, and (C)

863   permuting internal events on contrasts standardized branch length transformed trees

864   produces distinct patterns compared to the empirical gene trees of (A).

865

866 **Table 1: Information on different tree sets, number of internal node events, and node**

867 **ages used in this reanalysis.**

| Datasets | Number of trees | Number of speciation events | Number of duplication events | Number of NA events | Maximum speciation node age (My) | Maximum duplication node age (My) |
|---|---|---|---|---|---|---|
| Dunn et al.: full set | 8520 | 67911 | 21071 | 26794 | 296 | 11799977 |
| Dunn et al.: trees with strong phylogenetic signals | 2082 | 13118 | 4056 | 5186 | 296 | 1342 |
| This study: after excluding pure speciation trees | 4288 | 38882 | 15274 (8556 young + 6718 old) | 15201 | 296 | 1175 |

868    *Note*- My: Million years; young: age ≤ 296 My; old: age > 296 My.

**Table 2: Summary statistics on 1101 OUM trees passing a posterior probability cutoff of ≥ 0.7 in a Bayesian framework.**

| Duplication Age | Proportions of regime shifts per branch | | Paired two-sided Wilcoxon rank sum test | Regime shift rates (shifts/My) | | Two-sided Wilcoxon rank test |
|---|---|---|---|---|---|---|
| | After speciation | After duplication | | After speciation | After duplication | |
| Young | 3.1% | 4.5% | $3.4e^{-12}$ | 0.013 | 0.031 | $1.7e^{-11}$ |
| Old | 2.6% | 10% | $< 2.2e^{-16}$ | 0.013 | 0.0023 | $< 2.2e^{-16}$ |

*Note*- Above analyses include 13824 speciation, 3027 young and 2814 old duplication events. Values shown in the table indicate median values. The difference in proportions of regime shifts per branch after speciation events for two types of duplications is due to the different sets of trees used. Few trees shared both types of duplicates. Proportions of regime shifts per branch of events is estimated for each tree, and thus paired Wilcoxon test is used to compare the difference. A single gene tree can have multiple optima shift rates for events, and thus two-sided Wilcoxon rank test was used for comparison.

44

880 **Table 3: Summary statistics for Brownian trees.**

| Duplication age | Data | $\sigma^2_{\text{Speciation}}$ | $\sigma^2_{\text{Duplication}}$ | $\sigma^2_{\text{Duplication}} / \sigma^2_{\text{Speciation}}$ | *P*-value |
|---|---|---|---|---|---|
| Young | Empirical ($n_{\text{Speciation}}$= 4642; $n_{\text{Duplication}}$ = 1742) | $9e^{-5}$ | $1.4e^{-4}$ | 1.5 | $5e^{-12}$ |
| | Randomized $\tau$ ($n_{\text{Speciation}}$= 4618; $n_{\text{Duplication}}$ = 1723) | $6.9e^{-4}$ | $2.2e^{-4}$ | 0.32 | $1.4e^{-13}$ |
| | Randomized events ($n_{\text{Speciation}}$= 3215; $n_{\text{Duplication}}$ = 1438) | $1.7e^{-4}$ | $8.5e^{-5}$ | 0.5 | 0.02 |
| Old | Empirical ($n_{\text{Speciation}}$= 5356; $n_{\text{Duplication}}$ = 1295) | $1.7e^{-4}$ | $2e^{-9}$ | $1.2e^{-5}$ | $< 2.2e^{-16}$ |
| | Randomized $\tau$ ($n_{\text{Speciation}}$= 5337; $n_{\text{Duplication}}$ = 1291) | $9.1e^{-4}$ | $2.5e^{-10}$ | $2.7e^{-7}$ | $< 2.2e^{-16}$ |
| | Randomized events ($n_{\text{Speciation}}$= 2788; $n_{\text{Duplication}}$ = 800) | $1.8e^{-4}$ | $2.1e^{-9}$ | $1.2e^{-5}$ | $< 2.2e^{-16}$ |

881 *Note:* Median values of $\sigma^2$ are shown. *P*-value from paired two-sided Wilcoxon test.

882

## Supporting Information

883

884

885    **Figure S1: Expectations from phylogenetic and pairwise comparison approaches**

886    **under null and ortholog conjecture scenarios.** PIC: Phylogenetic Independent Contrast,

887    OC: Ortholog Conjecture. We present 4 time-calibrated gene trees of Dunn et al. (2018) as

888    illustration. Trees A and B are well calibrated, with the duplication ages are constrained by

889    speciation ages, as shown by the time scales below each phylogeny. Trees C and D

890    represent biased calibrated trees, where old duplication branches are inaccurately calibrated

891    due to lack of age constraints. To evaluate the impacts of gene duplication and speciation

892    events in trait evolution, pairwise comparisons do not rely on the branch lengths of a

893    calibrated phylogeny, but phylogenetic methods do. If time calibration of old duplication

894    nodes has no influence in the inference of phylogenetic approaches, we expect to obtain

895    patterns under a null and OC scenarios as shown in the right part of the figure. This means

896    that the phylogenetic contrasts or pairwise correlations of different events should be drawn

897    from the same distribution under a null model, while the expectation differs under the OC

898    model. We used 2 times higher rates of trait evolution ($\tau$ here) following duplications than

899    speciations (i.e. $\sigma^2_{duplication} = 2 * \sigma^2_{speciation}$) in this example for the OC model.

900    **Figure S2.: Repeating simulations on all calibrated trees with different random seed**

901    **number.** *P* values are from Wilcoxon two-tailed tests. Simulations with different seed

902    number did not change the trend of results as reported in Fig. 1A.

903    **Figure S3: Simulation analyses on 1135 trees with strong phylogenetic signals.** *P*

904    values are from Wilcoxon two-tailed tests. Dunn et al. used a cutoff of K > 0.551 to identify

905    trees with strong phylogenetic signals. However, trees with higher K statistic can have

906    corresponding *P* values which are non-significant. Considering both K statistic and *P* value,

907    we found similar trends as was observed with 2082 trees.

908    **Figure S4: Difference between time calibration approaches of Dunn et al. (2018) and**

909    **of this study.** In this example, we used the phylogeny of ACP1 gene. The top panel (A-C)

910    shows the steps used by Dunn et al. (2018), while the bottom panel (D-F) shows the steps

911    used in this study. Gene trees obtained from Ensembl (Herrero et al. 2016) have branch

912    lengths in substitutions per site. (A) and (D) are the same gene tree, where Dunn et al.

913    (2018) edited few speciation events to 'NA' to pass the time calibration step. (B) The gene

914    trees are pruned to species with available τ. (C) The pruned tree is time calibrated using

915    speciation time points. Pruning before time calibration produces tree with many

916    duplications, and NA nodes older to the oldest speciation nodes as in (B). This leads to

917    using only 7 speciation time points for calibration. Due to unavailable age constraints on

918    the old duplication nodes, the time scale of the phylogeny in (C) reaches 880 million years

919    (My). When we performed time calibration before pruning as in (E), we could use 32

920    speciation nodes for time calibration. This means that we could use many speciation nodes

921    for time calibration, although τ data was unavailable for species at tips due to the choice of

922    species in this study. Hence, the old duplication nodes are constrained by the age of

923    speciation nodes older to them, and thus the maximum age is now of 356 My (F).

924    **Figure S5: Re-analyses of expected variances of calibrated trees considered by Dunn**

925    **et al.** The expected variance plots of (A) all 8520 calibrated trees, and (B) 2082 trees with

926    strong phylogenetic signal. The dotted line represents the mean expected variance of the

927    events. These plots show why duplication nodes preceding ancient speciation nodes can be

928    problematic for PIC.

929    **Figure S6: *P* value distribution plots after 100 independent runs on each set of trees.**

930    Wilcoxon two-tailed test with 95% confidence interval was used to compare the speciation

931    and duplication contrasts after randomization tests. (A) and (B) applied to trees with at least

932    one speciation and one duplication event. (C) and (D) applied to trees with strong

933    phylogenetic signal. (A) and (C) randomization of trait ($\tau$) over the trees. (B) and (D)

934    randomization of internal node events. The inset plots show *P* values adjusted with

935    Benjamini-Hochberg (Benjamini and Yekutieli 2005; Hochberg and Benjamini 1990).

936    Supporting our observations of Figs. 2A and 2C, all the plots confirm that the empirical

937    result of Dunn et al. (2018) is not different from randomized test results.

938    **Figure S7: The ortholog conjecture test after randomizations of contrasts**

939    **standardized trees.** *P* values are from Wilcoxon two-tailed tests. 'PICs': Phylogenetic

940    Independent Contrasts. Values inside boxplots denote median PIC value of the

941    corresponding event. (A-B) Plots after randomizing $\tau$, and after randomizing events using

942    the same trees as in Fig. 3.

943    **Figure S8: Multivariate model fitting result using a maximum likelihood framework.**

944    BM1: Single rate Brownian; BMM: Multi rates Brownian; OU1: Single optimum Ornstein-

945    Uhlenbeck; and OUM: Multi optima Ornstein-Uhlenbeck models.

946    **Figure S9: The ortholog conjecture test for $\tau$ on calibrated trees of Dunn et al.** *P* values

947    are from Wilcoxon two-tailed tests. 'PICs': Phylogenetic Independent Contrasts. Values

948    inside boxplots denote median PIC value of the corresponding event. (A) Using 8417 out

949    of 8520 calibrated trees that passed diagnostic tests following branch length transformation.

48

950  (B) Plot after randomizing $\tau$, and (C) after randomizing events using the same branch

951  transformed trees as in (A).

952  **Figure S10: The ortholog conjecture test for $\tau$ on branch transformed trees with**

953  **strong phylogenetic signals.** *P* value are from Wilcoxon two-tailed tests. 'PICs':

954  Phylogenetic Independent Contrasts. Values inside boxplots denote median PIC value of

955  the corresponding event. (A) using 2080 out of 2082 calibrated trees that passed diagnostic

956  tests following branch length transformation. (B) Plot after randomizing $\tau$, and (C) after

957  randomizing events using the same branch transformed trees as in (A).

958

959    **Table S1: Summary statistics of calibrated old duplication nodes for 8420 trees of**

960    **Dunn et al. (2018).**

| Maximum age group in Million Years (My) | Count |
|---|---|
| 296-500 | 2917 |
| 501-900 | 7548 |
| 901-3000 | 49 |
| 3001-10000 | 16 |
| 10001-11799977 | 9 |

961

962

963 **Table S2: Analyses on multi optima OU trees.**

| Duplication age | Data | $\theta_{Speciation}$ | $\theta_{Duplication}$ | $\theta_{Duplication} / \theta_{Speciation}$ | *P*-value |
|---|---|---|---|---|---|
| Young | Empirical ($n_{Speciation}$= 2690; $n_{Duplication}$ = 842) | 0.41 | 0.74 | 1.8 | $8.6e^{-10}$ |
| | Randomized τ ($n_{Speciation}$= 2690; $n_{Duplication}$ = 842) | 0.53 | 0.55 | 1.03 | 0.97 |
| | Randomized events ($n_{Speciation}$= 1872; $n_{Duplication}$ = 698) | 0.50 | 0.53 | 1.06 | 0.75 |
| Old | Empirical ($n_{Speciation}$= 4152; $n_{Duplication}$ = 847) | 0.42 | 0.92 | 2.19 | $2.4e^{-4}$ |
| | Randomized τ ($n_{Speciation}$= 4152; $n_{Duplication}$ = 847) | 0.54 | $1.5e^{-11}$ | $2.8e^{-11}$ | $1.3e^{-07}$ |
| | Randomized events ($n_{Speciation}$= 2081; $n_{Duplication}$ = 482) | 0.51 | 0.66 | 1.29 | 0.73 |

964 *Note:* Median values of $\sigma^2$ are shown. *P*-value from paired two-sided Wilcoxon test.

965 **Table S3: Summary statistics on OUM trees, passing both the maximum likelihood**

966 **and the Bayesian approaches with a posterior probability cutoff of ≥ 0.7.**

967

| Duplication Age | Proportions of regime shifts per branch | | Paired two-sided Wilcoxon rank sum test | Regime shift rates (shifts/My) | | Two-sided Wilcoxon rank test |
|---|---|---|---|---|---|---|
| | After speciation | After duplication | | After speciation | After duplication | |
| Young | 2.8% | 8.3% | $8.3e^{-4}$ | 0.012 | 0.032 | $6.7e^{-6}$ |
| Old | 0% | 16.7% | $< 2.2e^{-16}$ | 0.012 | 0.0025 | $< 2.2e^{-16}$ |

968 *Note*- Above analyses include 2779 speciation, 486 young, and 548 old duplication events.

969 Values shown in the table indicate median values. The difference in proportions of regime

970 shifts per branch after speciation events for two types of duplications is due to the different

971 sets of trees used. Few trees shared both types of duplicates. Proportions of regime shifts

972 per branch of events is estimated for each tree, and thus paired Wilcoxon test is used to

973 compare the difference. A single gene tree can have one or many optima shift rate(s) for

974 events, and thus two-sided Wilcoxon rank test was used for comparison.
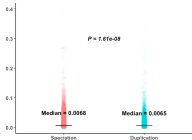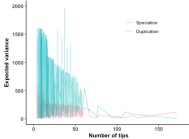
975

976

52

**A. Plots on 8520 trees**

P = 2.42e-04    P < 2.2e-16

**B. Plots on 2082 trees**

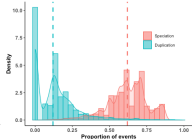P = 7.45e-04    P < 2.2e-16

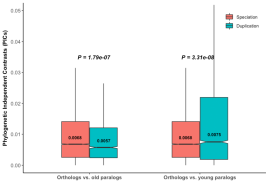**A.** Plot using randomized trait, $\tau$

**C.** Plot using randomized events
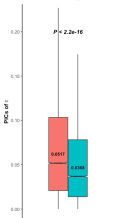
**B.** Variance distribution for events

**D.** Event proportion plot

Plots on contrasts standardized trees

**A. Plot on 4190 trees**   $P < 2.2e-16$

**B. Randomized τ plot**   $P = 2.2e-16$

**C. Randomized events plot**   $P = 1.03e-05$

Speciation
Duplication