

# Gene Tree Discordance Does Not Explain Away the Temporal Decline of Convergence in Mammalian Protein Sequence Evolution

Zhengting Zou<sup>1</sup> and Jianzhi Zhang<sup>\*,2</sup>

<sup>1</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI

<sup>2</sup>Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI

\*Corresponding author: E-mail: jianzhi@umich.edu.

Associate editor: Jeffrey Thorne

## Abstract

Several authors reported lower frequencies of protein sequence convergence between more distantly related evolutionary lineages and attributed this trend to epistasis, which renders the acceptable amino acids at a site more different and convergence less likely in more divergent lineages. A recent primate study, however, suggested that this trend is at least partially and potentially entirely an artifact of gene tree discordance (GTD). Here, we demonstrate in a genome-wide data set from 17 mammals that the temporal trend remains (1) upon the control of the GTD level, (2) in genes whose genealogies are concordant with the species tree, and (3) for convergent changes, which are extremely unlikely to be caused by GTD. Similar results are observed in a comparable data set of 12 fruit flies in some but not all of these tests. We conclude that, at least in some cases, the temporal decline of convergence is genuine, reflecting an impact of epistasis on protein evolution.

**Key words:** convergent evolution, epistasis, incomplete lineage sorting, parallel evolution.

## Diminishing Convergence over Time in Protein Evolution

Protein sequence convergence refers to independent amino acid substitutions at the same site in two or more evolutionary lineages that result in the same end state. It may signal molecular adaptation and therefore has long been of interest (Stewart et al. 1987; Doolittle 1994; Zhang and Kumar 1997; Castoe et al. 2009; Christin et al. 2010; Storz 2016). Recent genomic analyses found that protein sequence convergence is widespread (Bazykin et al. 2007; Rokas and Carroll 2008; Parker et al. 2013; Foote et al. 2015; Zou and Zhang 2015a), but the vast majority appears to have occurred by chance rather than by adaptive selection (Foote et al. 2015; Thomas and Hahn 2015; Zou and Zhang 2015a, 2015b). Interestingly, a number of authors reported that the frequency of sequence convergence between two lineages decreases as their genetic distance increases (Rogozin et al. 2008; Povolotskaya and Kondrashov 2010; Naumenko et al. 2012; Goldstein et al. 2015; Zou and Zhang 2015a). It was proposed that epistasis, or interaction among amino acid residues within or between proteins, causes the selective constraint at the same site to vary among species depending on the genetic background. Consequently, the probability of convergence between two lineages declines as they become more divergent from each other (Goldstein et al. 2015; Zou and Zhang 2015a). Indeed, the same amino acid sites were found to be subject to different selective constraints in different lineages, and a computer simulation confirmed that epistasis can produce diminishing convergence over time (Zou and Zhang 2015a). Nevertheless,

Mendes and colleagues recently proposed an alternative explanation of the temporal decline of convergence (Mendes et al. 2016). Specifically, incomplete lineage sorting (ILS) and introgression can cause a gene genealogy to differ from the underlying species tree, a phenomenon called gene tree discordance (GTD). GTD creates artificial signals of convergence when sequence changes are inferred using the species tree. Because the probability of ILS and introgression declines with species divergence, the amount of artificial convergence created by GTD is expected to reduce as the two lineages compared become more distant from each other. Indeed, Mendes et al. found several lines of evidence supporting their hypothesis, including the disappearance of the temporal decline of convergence in a 12-primate data set of 5,264 genes when the influence of GTD is excluded (Mendes et al. 2016). It is clear that GTD cannot explain the temporal decline of convergence observed in mitochondrial genes (Goldstein et al. 2015) due to the unique features of mitochondrial inheritance (Mendes et al. 2016). What is unclear, however, is whether GTD is fully responsible for the temporal declines of convergence in other nuclear gene data sets, because the relative contributions of GTD and epistasis to the temporal trend likely depend on the level of species divergence, which varies among data sets. It is important to clarify the above question, because if the temporal trend is always fully explainable by GTD in nuclear genes, there would be no genuine diminishing convergence over time for these genes and the role of epistasis in protein evolution might be substantially smaller than is currently thought. We therefore reanalyzed the two nuclear gene data sets (17 mammals and 12 fruit flies,

respectively) where we previously discovered the temporal decline of convergence (Zou and Zhang 2015a).

## Convergence Measures

Hereinafter, independent amino acid substitutions at the same site that have the same ancestral state and the same end state are referred to as parallel changes while those with different ancestral states are referred to as convergent changes (Zhang and Kumar 1997). These two categories are collectively referred to as convergence. The distinction between parallel and convergent changes is important, because GTD is expected to create artificial parallel changes but not artificial convergent changes (Mendes et al. 2016). The reason for the latter notion is that for ILS to create artificial convergent changes at a site, the site must be polymorphic with at least three distinct, high-frequency alleles, which is extremely unlikely. Similarly, for introgression to create artificial convergent changes at a site, the site must experience at least two different amino acid substitutions within a time that is sufficiently short to allow introgression, which is improbable.

Mendes et al. used the ratio between the observed numbers of convergences and divergences ( $C/D$ ) as a measure of convergence level between two lineages. Note that divergences can be separated into two types: those starting from the same ancestral states as in parallel changes and those starting from different ancestral states as in convergent changes. It is known that, when  $C$  is the total number of parallel and convergent changes and  $D$  is the total number of the two types of divergence events,  $C/D$  decreases with the divergence time between the two lineages compared even when neither epistasis nor GTD exists, because the probability of convergent changes relative to that of parallel changes rises as the genetic distance between the two lineages increases (Goldstein et al. 2015). Mendes et al. suggested that  $C/D$  no longer declines with the divergence time when only parallel or convergent (but not both) changes are considered in  $C$  and only the corresponding type of divergence events is considered in  $D$ ; these two  $C/D$  ratios are, respectively, referred to as  $(C/D)_s$  and  $(C/D)_d$ , where the subscript “ $s$ ” stands for the same ancestral states and “ $d$ ” stands for different ancestral states. Our computer simulation in the absence of epistasis and GTD confirmed that  $C/D$ , but not  $(C/D)_s$  or  $(C/D)_d$ , decreases with time (supplementary fig. S1, Supplementary Material online).

In addition to  $(C/D)_s$  and  $(C/D)_d$ , we used the ratio ( $R$ ) between the observed and expected numbers of convergences to measure the convergence level, because  $R$  has a clear biological meaning and, in the absence of epistasis and GTD, is not expected to correlate with the genetic distance between lineages, as was previously demonstrated by simulation (Zou and Zhang 2015a). An amino acid substitution model is needed in computing  $R$ , and we used two models employed in the original study: JTT- $f_{\text{site}}$  and JTT- $f_{\text{gene}}$  (Zou and Zhang 2015a). Both models assume the JTT substitution matrix (Jones et al. 1992) except that the former uses the observed amino acid frequencies at a site as its equilibrium amino acid frequencies whereas the latter uses the observed amino acid frequencies of an entire protein as the equilibrium

frequencies at each site of the protein. To examine the robustness of our results, we used two different distances between evolutionary lineages: (1) the total length of branches linking the descendant nodes of the two branches compared (Zou and Zhang 2015a) and (2) the total length of branches linking the ancestral nodes of the two branches compared (Mendes et al. 2016). They are referred to as  $d_1$  and  $d_2$ , respectively. Using  $d_1$  and  $d_2$  yielded qualitatively similar results in most cases (table 1). We therefore describe only the results with  $d_1$  in the main text unless otherwise mentioned.

## Does GTD Fully Explain the Temporal Decline of Convergence in Mammals: Test I

A straightforward statistical test of the null hypothesis that GTD fully explains the temporal decline of convergence is to conduct a partial correlation between genetic distance and convergence level after controlling the GTD level. The partial correlation should be zero under the null hypothesis. We first inferred the maximum likelihood gene tree for each protein. For each independent branch pair in the species tree, we sampled four species whose tree includes the two focal branches and their respective sister branches (supplementary fig. S2, Supplementary Material online), and estimated the GTD level for the focal branch pair by the proportion of genes whose gene trees differ from the species tree of these four species (see Materials and Methods).

We started with the mammalian data, composed of 2,759 protein sequence alignments of 14 placentals, two marsupials, and one monotreme (Zou and Zhang 2015a). The GTD level can be evaluated for 208 branch pairs (see Materials and Methods). For these branch pairs, we found  $R$  to be negatively correlated with  $d_1$  even after the control of the GTD level ( $r = -0.51$ ,  $P = 0.03$  under JTT- $f_{\text{site}}$ ;  $r = -0.64$ ,  $P = 0.001$  under JTT- $f_{\text{gene}}$ ; table 1), suggesting that GTD does not explain away the diminishing convergence over time.

Because GTD could result in apparent parallel changes, we further tested the null hypothesis by correlating  $(C/D)_s$  with  $d_1$  after the control of the GTD level. This partial correlation is significantly negative ( $r = -0.66$ ,  $P = 0.0003$ ; table 1), consistent with the result based on  $R$ .

If GTD is the primary cause of the temporal decline of protein sequence convergence as Mendes et al. proposed,  $(C/D)_s$  for synonymous sites is also expected to decline with  $d_1$  (Mendes et al. 2016). But this trend is not statistically significant (table 1), suggesting at most a minor influence of GTD on convergence level in our data. Note that the significant negative correlation between  $(C/D)_s$  for synonymous sites and  $d_1$  after the control of GTD (table 1) is due to the unexpected negative correlation between  $(C/D)_s$  and GTD (e.g.,  $r = -0.38$ ,  $P = 0.009$  upon the control of  $d_1$ ), which does not conform to Mendes et al.’s hypothesis.

Together, test I demonstrates that, in the mammalian data, GTD is not the primary cause of the temporal decline in protein convergence. In the original study (Zou and Zhang 2015a), we rejected the hypothesis that potential genome-wide changes in amino acid frequencies cause the observed temporal decline of convergence. Hence, we no longer

**Table 1.** Pearson's Correlations between Genetic Distance and Various Convergence Levels.

	Mammals				Fruit flies			
	Amino acid			Synonymous	Amino acid			Synonymous
	R (JTT- $f_{\text{site}}$ )	R (JTT- $f_{\text{gene}}$ )	C/D	C/D	R (JTT- $f_{\text{site}}$ )	R (JTT- $f_{\text{gene}}$ )	C/D	C/D
<b>Mantel test (all genes)</b>								
$d_1$	-0.73** <sup>a</sup>	-0.79****	-0.74**** <sup>b</sup>	-0.19 <sup>b</sup>	-0.49*	-0.63**	-0.44* <sup>b</sup>	0.12 <sup>b</sup>
$d_2$	-0.88****	-0.74****	-0.57**** <sup>b</sup>	-0.0052 <sup>b</sup>	-0.75****	-0.72****	-0.41** <sup>b</sup>	0.19 <sup>b</sup>
<b>Partial Mantel test (controlling GTD)</b>								
$d_1$	-0.51*	-0.64**	-0.66**** <sup>b</sup>	-0.38* <sup>b</sup>	0.55	0.18	-0.015 <sup>b</sup>	0.049 <sup>b</sup>
$d_2$	-0.72****	-0.48****	-0.40* <sup>b</sup>	-0.21 <sup>b</sup>	-0.015	0.23	0.20 <sup>b</sup>	0.18 <sup>b</sup>
<b>Mantel test (concordant genes)</b>								
$d_1$	-0.53****	-0.60****	-0.54**** <sup>b</sup>	-0.015 <sup>b</sup>	-0.42*	-0.56*	-0.42* <sup>b</sup>	0.094 <sup>b</sup>
$d_2$	-0.68****	-0.55****	-0.38** <sup>b</sup>	0.12 <sup>b</sup>	-0.71****	-0.69****	-0.43** <sup>b</sup>	0.20 <sup>b</sup>
<b>Mantel test (convergent changes)</b>								
$d_1$	-0.45*	-0.51**	-0.31* <sup>c</sup>	0.27 <sup>c</sup>	-0.32	-0.44*	-0.21 <sup>c</sup>	0.46 <sup>c</sup>
$d_2$	-0.52****	-0.47****	-0.33** <sup>c</sup>	0.17 <sup>c</sup>	-0.47**	-0.50****	-0.31* <sup>c</sup>	0.43 <sup>c</sup>

<sup>a</sup>Significance is shown only when  $r < 0$ .

<sup>b</sup>(C/D)<sub>s</sub>.

<sup>c</sup>(C/D)<sub>d</sub>.

\* $P < 0.05$ ;

\*\* $P < 0.01$ ;

\*\*\* $P < 0.001$ ;

\*\*\*\* $P \leq 0.0001$ .

$d_1$ , total length of branches linking the descendant nodes of the two branches compared;  $d_2$ , total length of branches linking the ancestral nodes of the two branches compared.

consider this possibility here. A potential source of error in our analysis arises from ancestral sequence inference. Analyzing 2,759 protein sequence alignments generated by an evolutionary simulation with realistic parameters for the species tree, branch lengths, site-specific evolutionary rates, and JTT- $f_{\text{gene}}$  model with gene-specific amino acid frequencies, we found no significant correlation between (C/D)<sub>s</sub> and genetic distance, confirming that ancestral sequence inference and other steps in our analysis do not create artificial diminishing convergence over time.

### Does GTD Fully Explain the Temporal Decline of Convergence in Mammals: Test II

The null hypothesis that GTD fully explains the temporal decline of convergence can be further tested by examining genes whose gene trees are concordant with the species tree, because the temporal pattern of convergence caused by GTD should disappear when only the concordant genes are analyzed. In the mammalian data, only 77 gene trees are concordant with the presumptive species tree. Nonetheless, the negative correlation between  $R$  and  $d_1$  remains significant for these concordant genes ( $r = -0.53$ ,  $P = 5 \times 10^{-4}$  under JTT- $f_{\text{site}}$ ;  $r = -0.60$ ,  $P = 5 \times 10^{-5}$  under JTT- $f_{\text{gene}}$ ; fig. 1a). Similarly, (C/D)<sub>s</sub> decreases with  $d_1$  ( $r = -0.54$ ,  $P = 0.0005$ ; red dots in fig. 1b). Note that removing all genes with discordant gene trees renders the above test conservative, because true convergence, which can also cause GTD, may have been removed too. Although the presumptive mammalian species tree may differ from the true species tree, the fact that we count convergence in all genes having the same gene tree ensures that gene tree variation does not affect our analysis. While recombination within genes may cause a seemingly concordant gene to harbor a discordant segment of DNA,

there is no correlation between  $d_1$  and (C/D)<sub>s</sub> for synonymous sites of concordant genes ( $r = -0.02$ ,  $P = 0.48$ ; gray dots in fig. 1b), suggesting no impact of potential residual GTD in concordant genes.

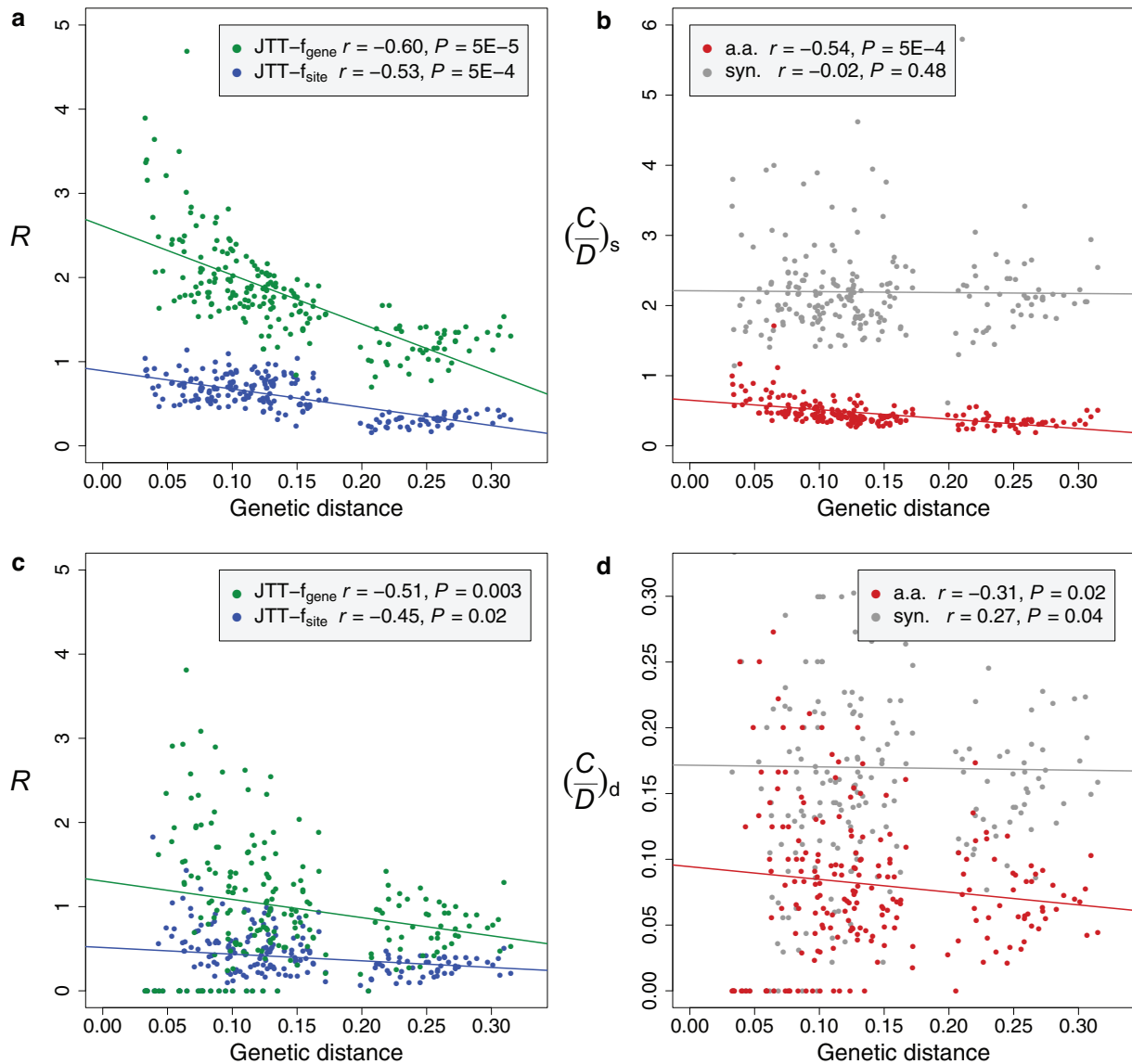
### Does GTD Fully Explain the Temporal Decline of Convergence in Mammals: Test III

Because it is very unlikely for GTD to create artificial convergent changes (Mendes et al. 2016), a negative correlation between  $R$  and  $d_1$  for convergent changes would not be explainable by GTD. Indeed, this correlation is significant ( $r = -0.45$ ,  $P = 0.02$  under JTT- $f_{\text{site}}$ ;  $r = -0.51$ ,  $P = 0.003$  under JTT- $f_{\text{gene}}$ ; fig. 1c). The same trend is found between (C/D)<sub>d</sub> and  $d_1$  ( $r = -0.31$ ,  $P = 0.02$ ; red dots in fig. 1d). Similar to using only concordant genes, using only convergent changes renders our test conservative, because all parallel changes are excluded despite that only a fraction of them may be artifacts of GTD. As expected, no negative correlation is observed between  $d_1$  and (C/D)<sub>d</sub> for synonymous sites ( $r = 0.27$ ; gray dots in fig. 1d).

Taken together, the three tests support that the temporal decline of convergence in the mammalian data is not fully attributable to GTD. This finding, in conjunction with the previously published evidence for epistasis (Zou and Zhang 2015a), strongly implicates epistasis in causing diminishing convergence over time in mammals.

### Does GTD Fully Explain the Temporal Decline of Convergence in Fruit Flies?

We next analyzed the fruit fly data, composed of 5,935 protein alignments from 12 *Drosophila* species (Zou and Zhang 2015a). For this data set, GTD level can be evaluated for 84



**Fig. 1.** Correlation between convergence level and genetic distance in mammals. (a) Scatter plot of  $R$  against the genetic distance  $d_1$  for genes having gene trees concordant with the presumptive species tree (“concordant genes”). (b) Scatter plot of  $(C/D)_s$  against  $d_1$  for concordant genes. (c) Scatter plot of  $R$  for convergent changes in all genes against  $d_1$ . (d) Scatter plot of  $(C/D)_d$  for all genes against  $d_1$ . Each dot represents a branch pair and different colors show results under different substitution models or for different types of substitutions, as indicated in inset legends.  $d_1$  is the number of amino acid substitutions per site between the descendant nodes of the two branches considered. The  $r$  values are Pearson’s correlation coefficients. Both  $r$  and  $P$  values are from Mantel tests. Colored lines show linear regressions from data points of the same color. a.a.: amino acid substitutions; syn.: synonymous nucleotide substitutions.

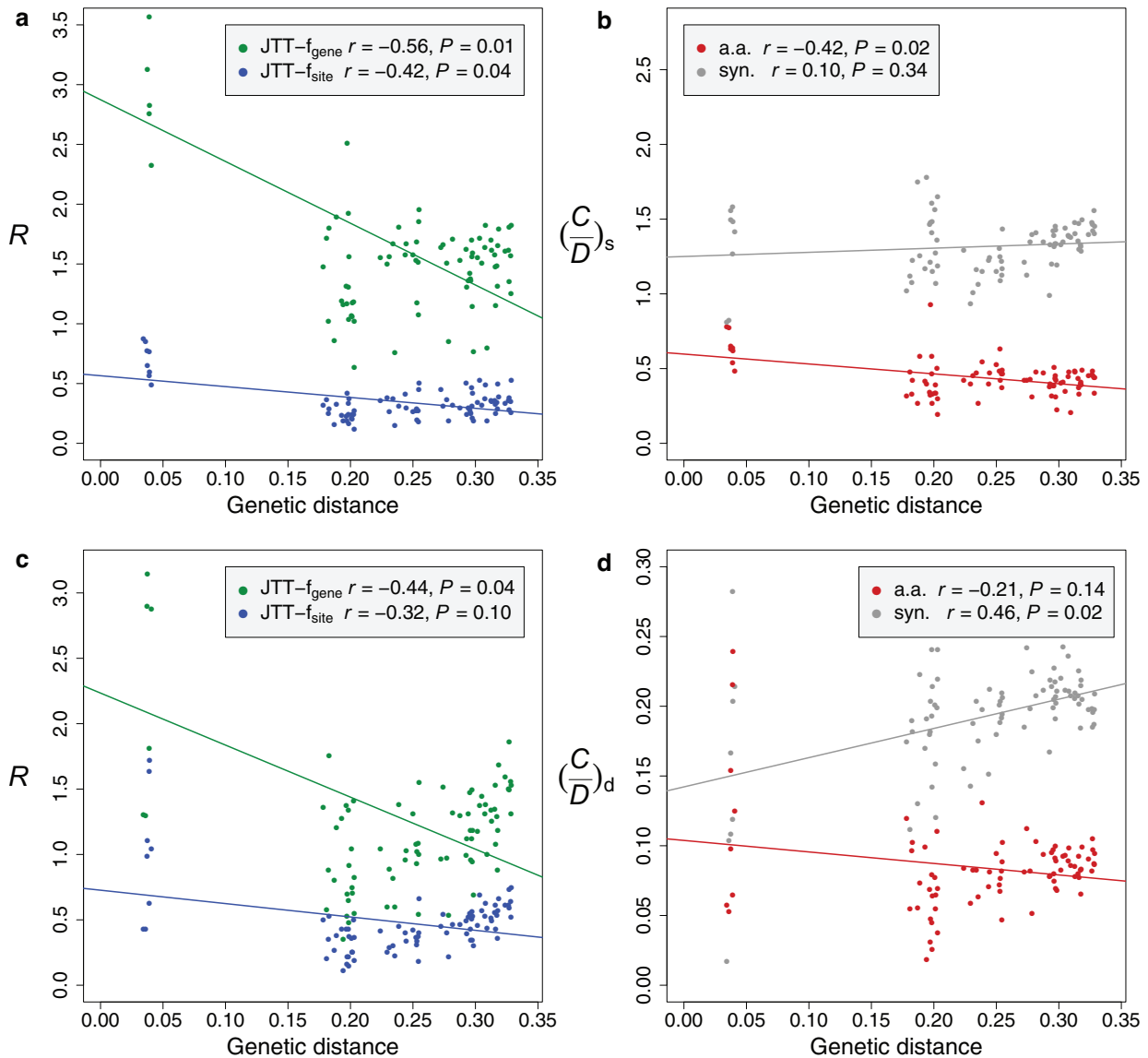
branch pairs. The null hypothesis that GTD fully explains the temporal decline of convergence in fruit flies is refuted by some but not all of the three tests. First, upon control of the GTD level, no significant negative partial correlation was observed between genetic distance and  $R$  or  $(C/D)_s$  (table 1), failing to reject the null hypothesis. Second, for concordant genes, a significant negative correlation was detected between  $d_1$  and  $R$  ( $r = -0.42$ ,  $P = 0.04$  under JTT- $f_{\text{site}}$ ;  $r = -0.56$ ,  $P = 0.01$  under JTT- $f_{\text{gene}}$ ; fig. 2a) or  $(C/D)_s$  ( $r = -0.42$ ,  $P = 0.02$ ; red dots in fig. 2b) but not  $(C/D)_s$  for synonymous sites ( $r = 0.10$ ,  $P = 0.34$ ; gray dots in fig. 2b), refuting the null hypothesis. Finally, when only convergent changes are considered, the null hypothesis is rejected when  $d_2$  is used (table 1), but is rejected in some but not all analyses when  $d_1$  is used

(fig. 2c and d and table 1). These inconsistent results suggest that GTD is more important than epistasis in creating the pattern of diminishing convergence over time in the *Drosophila* data.

### Relative Contributions of GTD and Epistasis to Temporal Declines of Protein Convergence

The difference in the relative contributions of GTD and epistasis to the temporal decline of convergence among the three data sets (primates in Mendes et al. (2016); mammals and flies in this study) is at least in part caused by different frequencies of GTD in the three groups of organisms analyzed. For three species, the probability of GTD due to ILS is  $P = \frac{2}{3}e^{-\frac{T}{2N}}$ , where





**Fig. 2.** Correlation between convergence level and genetic distance in fruit flies. (a) Scatter plot of  $R$  against the genetic distance  $d_1$  for genes having gene trees concordant with the presumptive species tree (“concordant genes”). (b) Scatter plot of  $(C/D)_s$  against  $d_1$  for concordant genes. (c) Scatter plot of  $R$  for convergent changes in all genes against  $d_1$ . (d) Scatter plot of  $(C/D)_d$  for all genes against  $d_1$ . Each dot represents a branch pair and different colors show results under different substitution models or for different types of substitutions, as indicated in inset legends.  $d_1$  is the number of amino acid substitutions per site between the descendant nodes of the two branches considered. The  $r$  values are Pearson’s correlation coefficients. Both  $r$  and  $P$  values are from Mantel tests. Colored lines show linear regressions from data points of the same color. a.a.: amino acid substitutions; syn.: synonymous nucleotide substitutions.

$T$  is the number of generations between the two relevant speciation events and  $N$  is the effective population size (Hudson 1983; Pamilo and Nei 1988). Let us assume that for mammals ( $M$ ),  $N_M = 10^4$  and generation time  $t_M = 5$  years, and for *Drosophila* fruit flies ( $D$ )  $N_D = 10^6$  and  $t_D = 0.1$  year (Charlesworth 2009). Given the same time interval of  $T'$  million years between relevant speciation events,  $P_D/P_M = e^{\frac{T'}{2N_M t_M}} - e^{\frac{T'}{2N_D t_D}} = e^{5T'}$ . Hence, the probability of GTD is expected to be higher in fruit flies than in mammals given equal speciation frequencies between the two groups. For example, when  $T' = 0.5$  My, the probability of GTD due to ILS is 5.5% in fruit flies but only 0.45% in mammals. Introgression occurs in both mammals and flies

(Ballard 2000; Bachtrog et al. 2006; Mallet et al. 2016), although their rates are unclear. Consequently, the impact of GTD is expected to be higher in fruit flies than in mammals if ILS is an important contributor to GTD. The primate data have relatively short speciation intervals compared with the mammalian data and are thus expected to be influenced more by GTD. The impact of epistasis should depend on sequence divergence; data sets with larger ranges of sequence divergence are expected to be more influenced by epistasis. This factor may render epistasis more influential in the mammalian data than in the primate data. In the mammalian data, depending on the distance ( $d_1$  or  $d_2$ ) and convergence measures ( $R$  under  $JTT-f_{site}$ ,  $R$  under  $JTT-f_{gene}$  or  $(C/D)_s$ ) used, the

partial correlation between  $d$  and convergence after the control for GTD (in the Mantel test) is on average 76% of the corresponding correlation without the control for GTD (table 1), suggesting that the contribution of epistasis is at least as important as GTD.

In conclusion, we showed that, at least for the mammalian data analyzed, GTD cannot fully explain the temporal decline of convergence, which implicates the contribution of epistasis. The different results obtained from different data sets demonstrate that the relative roles of GTD and epistasis in creating diminishing convergence over time depend on speciation intervals and sequence divergences and are thus data-dependent.

## Materials and Methods

In a species tree, let branch  $X$  connect an interior node  $X_0$  and one of its two immediate descendants  $X_1$  and let branch  $Y$  connect an interior node  $Y_0$  and one of its two immediate descendants  $Y_1$  (supplementary fig. S2, Supplementary Material online). The GTD level can be estimated for all branch pairs where  $Y_0$  is not on the path from  $X_0$  to the tree root and  $X_0$  is not on the path from  $Y_0$  to the tree root. Let  $X_2$  be the other immediate descendant of  $X_0$  and let  $Y_2$  be the other immediate descendant of  $Y_0$ . Let exterior nodes  $X_1'$ ,  $X_2'$ ,  $Y_1'$ , and  $Y_2'$  be randomly picked descendants of  $X_1$ ,  $X_2$ ,  $Y_1$ , and  $Y_2$ , respectively. The four exterior nodes have a phylogenetic relationship of  $((X_1', X_2'), (Y_1', Y_2'))$  in the species tree. For a gene, if the topology of  $X_1'$ ,  $X_2'$ ,  $Y_1'$ , and  $Y_2'$  in the gene tree is inconsistent with that in the species tree, this gene is defined as showing GTD for branch pair  $(X, Y)$ . The overall GTD level for the branch pair  $(X, Y)$  is the proportion of genes that show GTD for  $(X, Y)$ . Gene trees were inferred using RAxML v8.2.4 under the JTT- $f_{\text{gene}}$  model with substitution rate variation following a gamma distribution (Stamatakis 2014). The species trees of the mammals and fruit flies analyzed here, respectively, follow those in figures 2a and 3a of Zou and Zhang (2015a). Mantel tests and partial Mantel tests were conducted using the R package “nCF”. In all matrices used for these tests, entries that do not correspond to a branch pair with an available GTD level were set as “NA”. The partial Mantel test used method 1 of permutation, which permutes the entire matrix of  $R$  or  $C/D$  values (Legendre 2000). Protein sequences were acquired from Zou and Zhang (2015a), who obtained them from OrthoMaM v8 (Douzery et al. 2014) and Flybase (in October 2013). The corresponding coding DNA sequences were retrieved from OrthoMaM v9 and Flybase (in September 2016). The protein sequences and nucleotide sequences have consistent lengths after the removal of ambiguous sites as described in Zou and Zhang (2015a), and can be accessed from [http://www.umich.edu/~zhanglab/download/Zou\\_201702/index.htm](http://www.umich.edu/~zhanglab/download/Zou_201702/index.htm) (last accessed March 18, 2017). The 2,759 alignments of mammalian proteins have a median length of 315 amino acids, while the 5,935 alignments of fly proteins have a median length of 289 amino acids.

## Supplementary Material

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

## Acknowledgments

We thank Matt Hahn and three anonymous reviewers for valuable comments. This work was supported in part by U.S. National Institutes of Health research grant R01GM103232 to J.Z.

## References

- Bachtrog D, Thornton K, Clark A, Andolfatto P. 2006. Extensive introgression of mitochondrial DNA relative to nuclear genes in the *Drosophila yakuba* species group. *Evolution* 60:292–302.
- Ballard JW. 2000. When one is not enough: introgression of mitochondrial DNA in *Drosophila*. *Mol Biol Evol.* 17:1126–1130.
- Bazykin GA, Kondrashov FA, Brudno M, Poliakov A, Dubchak I, Kondrashov AS. 2007. Extensive parallelism in protein evolution. *Biol Direct* 2:20.
- Castoe TA, de Koning AP, Kim HM, Gu W, Noonan BP, Naylor G, Jiang ZJ, Parkinson CL, Pollock DD. 2009. Evidence for an ancient adaptive episode of convergent molecular evolution. *Proc Natl Acad Sci U S A.* 106:8986–8991.
- Charlesworth B. 2009. Fundamental concepts in genetics: effective population size and patterns of molecular evolution and variation. *Nat Rev Genet.* 10:195–205.
- Christin PA, Weinreich DM, Besnard G. 2010. Causes and evolutionary significance of genetic convergence. *Trends Genet.* 26:400–405.
- Doolittle RF. 1994. Convergent evolution: the need to be explicit. *Trends Biochem Sci.* 19:15–18.
- Douzery EJ, Scornavacca C, Romiguier J, Belkhir K, Galtier N, Delsuc F, Ranwez V. 2014. OrthoMaM v8: a database of orthologous exons and coding sequences for comparative genomics in mammals. *Mol Biol Evol.* 31:1923–1928.
- Foote AD, Liu Y, Thomas GW, Vinar T, Alföldi J, Deng J, Dugan S, van Elk CE, Hunter ME, Joshi V, et al. 2015. Convergent evolution of the genomes of marine mammals. *Nat Genet.* 47:272–275.
- Goldstein RA, Pollard ST, Shah SD, Pollock DD. 2015. Nonadaptive amino acid convergence rates decrease over time. *Mol Biol Evol.* 32:1373–1381.
- Hudson RR. 1983. Testing the constant-rate neutral allele model with protein-sequence data. *Evolution* 37:203–217.
- Jones DT, Taylor WR, Thornton JM. 1992. The rapid generation of mutation data matrices from protein sequences. *Comput Appl Biosci.* 8:275–282.
- Legendre P. 2000. Comparison of permutation methods for the partial correlation and partial Mantel tests. *J Stat Comput Simul.* 67:37–73.
- Mallet J, Besansky N, Hahn MW. 2016. How reticulated are species?. *BioEssays* 38:140–149.
- Mendes FK, Hahn Y, Hahn MW. 2016. Gene tree discordance can generate patterns of diminishing convergence over time. *Mol Biol Evol.* 33:3299–3307.
- Naumenko SA, Kondrashov AS, Bazykin GA. 2012. Fitness conferred by replaced amino acids declines with time. *Biol Lett.* 8:825–828.
- Pamilo P, Nei M. 1988. Relationships between gene trees and species trees. *Mol Biol Evol.* 5:568–583.
- Parker J, Tsagkogeorga G, Cotton JA, Liu Y, Provero P, Stupka E, Rossiter SJ. 2013. Genome-wide signatures of convergent evolution in echolocating mammals. *Nature* 502:228–231.
- Povolotskaya IS, Kondrashov FA. 2010. Sequence space and the ongoing expansion of the protein universe. *Nature* 465:922–926.
- Rogozin IB, Thomson K, Csuros M, Carmel L, Koonin EV. 2008. Homoplasy in genome-wide analysis of rare amino acid replacements: the molecular-evolutionary basis for Vavilov’s law of homologous series. *Biol Direct* 3:7.
- Rokas A, Carroll SB. 2008. Frequent and widespread parallel evolution of protein sequences. *Mol Biol Evol.* 25:1943–1953.
- Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313.

- Stewart CB, Schilling JW, Wilson AC. 1987. Adaptive evolution in the stomach lysozymes of foregut fermenters. *Nature* 330:401–404.
- Storz JF. 2016. Causes of molecular convergence and parallelism in protein evolution. *Nat Rev Genet.* 17:239–250.
- Thomas GW, Hahn MW. 2015. Determining the null model for detecting adaptive convergence from genomic data: a case study using echolocating mammals. *Mol Biol Evol.* 32:1232–1236.
- Zhang J, Kumar S. 1997. Detection of convergent and parallel evolution at the amino acid sequence level. *Mol Biol Evol.* 14:527–536.
- Zou Z, Zhang J. 2015a. Are convergent and parallel amino acid substitutions in protein evolution more prevalent than neutral expectations?. *Mol Biol Evol.* 32:2085–2096.
- Zou Z, Zhang J. 2015b. No genome-wide protein sequence convergence for echolocation. *Mol Biol Evol.* 32:1237–1241.