

# Quantifying the relationship between genetic diversity and population size suggests natural selection cannot explain Lewontin's Paradox

Vince Buffalo\*

Institute for Ecology and Evolution, University of Oregon, Eugene, United States

**Abstract** Neutral theory predicts that genetic diversity increases with population size, yet observed levels of diversity across metazoans vary only two orders of magnitude while population sizes vary over several. This unexpectedly narrow range of diversity is known as Lewontin's Paradox of Variation (1974). While some have suggested selection constrains diversity, tests of this hypothesis seem to fall short. Here, I revisit Lewontin's Paradox to assess whether current models of linked selection are capable of reducing diversity to this extent. To quantify the discrepancy between pairwise diversity and census population sizes across species, I combine previously-published estimates of pairwise diversity from 172 metazoan taxa with newly derived estimates of census sizes. Using phylogenetic comparative methods, I show this relationship is significant accounting for phylogeny, but with high phylogenetic signal and evidence that some lineages experience shifts in the evolutionary rate of diversity deep in the past. Additionally, I find a negative relationship between recombination map length and census size, suggesting abundant species have less recombination and experience greater reductions in diversity due to linked selection. However, I show that even assuming strong and abundant selection, models of linked selection are unlikely to explain the observed relationship between diversity and census sizes across species.

\*For correspondence:  
vsbuffalo@gmail.com

**Competing interests:** The author declares that no competing interests exist.

**Funding:** See page 19

**Received:** 12 February 2021

**Accepted:** 16 August 2021

**Published:** 19 August 2021

**Reviewing editor:** Guy Sella, Columbia University, United States

© Copyright Buffalo. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

## Introduction

A longstanding mystery in evolutionary genetics is that the observed levels of genetic variation across sexual species span an unexpectedly narrow range. Under neutral theory, the average number of nucleotide differences between sequences (pairwise diversity,  $\pi$ ) is determined by the balance of new mutations and their loss by genetic drift (*Kimura and Crow, 1964; Malécot, 1948; Wright, 1931*). In particular, expected pairwise diversity at neutral sites in a panmictic population of  $N_c$  diploids is  $\pi \approx 4N_c\mu$ , where  $\mu$  is the per basepair per generation mutation rate. Given that metazoan germline mutation rates only differ 10-fold ( $10^{-8}$ – $10^{-9}$ , *Kondrashov and Kondrashov, 2010; Lynch, 2010*), and census sizes vary over several orders of magnitude, under neutral theory one would expect that pairwise diversity also vary over several orders of magnitude. However, early allozyme surveys revealed that diversity levels across a wide range of species varied just an order of magnitude (*Lewontin, 1974*, p. 208); this is known as Lewontin's 'Paradox of Variation'. With modern sequencing-based estimates of  $\pi$  across taxa ranging over only three orders of magnitude (0.01–10%, *Leffler et al., 2012*), Lewontin's paradox remains unresolved through the genomics era.

Early on, explanations for Lewontin's Paradox have been framed in terms of the neutralist–selectionist controversy (*Lewontin, 1974; Kimura, 1984; Gillespie, 1991; Gillespie, 2001*). The neutralist view is that beneficial alleles are sufficiently rare and deleterious alleles are removed sufficiently quickly, that levels of genetic diversity are shaped predominantly by genetic drift and mutation

(Kimura, 1984). Specifically, non-selective processes decouple the effective population size implied by observed levels of diversity  $\hat{\pi}$ ,  $\tilde{N}_e = \hat{\pi}/4\mu$ , from the census size,  $N_c$ . By contrast, the selectionist view is that direct selection and the indirect effects of selection on linked neutral diversity suppress diversity levels across taxa, specifically because the impact of linked selection is greater in large populations. Undoubtedly, these opposing views represent a false dichotomy, as population genomic studies have uncovered evidence for the substantial impact of both demographic history (e.g. Zhao et al., 2013; Palkopoulou et al., 2015) and linked selection on genome-wide diversity (e.g. Elyashiv et al., 2016; Begun and Aquadro, 1992; Aguade et al., 1989; McVicker et al., 2009).

## Possible resolutions of Lewontin's Paradox

A resolution of Lewontin's Paradox would involve a mechanistic description and quantification of the evolutionary processes that prevent diversity from scaling with census sizes across species. This would necessarily connect to the broader literature on the empirical relationship between diversity and population size (Frankham, 1996; Nei and Graur, 1984; Soulé, 1976; Leroy et al., 2021), and the ecological and life history correlates of genetic diversity (Nevo, 1978; Powell, 1975; Nevo et al., 1984). Three categories of processes stand out as potentially capable of decoupling census sizes from diversity: non-equilibrium demography, variance and skew in reproductive success, and selective processes.

It has long been appreciated that effective population sizes are typically less than census population sizes, tracing back to early debates between R.A. Fisher and Sewall Wright (Fisher and Ford, 1947; Wright, 1948). Possible causes of this divergence between effective and census population sizes include demographic history (e.g. population bottlenecks), extinction and recolonization dynamics, or the breeding structure of populations (e.g. the variance in reproductive success and population substructure). Early explanations for Lewontin's Paradox suggested bottlenecks during the last glacial maximum severely reduced population sizes (Kimura, 1984; Ohta and Kimura, 1973; Nei and Graur, 1984), and emphasized that large populations recover to equilibrium diversity levels more slowly (Nei and Graur, 1984, Kimura, 1984 p. 203–204). Another explanation is that cosmopolitan species repeatedly endure extinction and recolonization events, which reduces effective population size (Maruyama and Kimura, 1980; Slatkin, 1977).

While intermittent demographic events like bottlenecks and recent expansions have long-term impacts on diversity (since mutation-drift equilibrium is reached on the order of size of the population), characteristics of the breeding structure such as high variance ( $V_w$ ) or skew in reproductive success continuously suppress diversity below the levels predicted by the census size (Wright, 1938). For example, in many marine animals, females are highly fecund, and dispersing larvae face extremely low survivorship, leading to high variance in reproductive success (Waples et al., 2018; Waples et al., 2013; Hedgecock and Pudovkin, 2011; Hauser and Carvalho, 2008). Such "sweepstakes" reproductive systems can lead to remarkably small ratios of effective to census population size (e.g.  $N_e/N_c$  can range from  $10^{-6}$ – $10^{-2}$ ), since  $N_e/N_c \approx 1/V_w$  (Hedgecock, 1994; Wright, 1938; Nunney, 1993), and require multiple-merger coalescent processes to describe their genealogies (Eldon and Wakeley, 2006). Overall, these reproductive systems diminish the diversity in some species, but seem unlikely to explain Lewontin's Paradox broadly across metazoans.

Alternatively, selective processes, and in particular the indirect effects of selection on linked neutral variation, could potentially explain the observed narrow range of diversity. The earliest mathematical model of hitchhiking was proffered as an explanation of Lewontin's Paradox (Smith and Haigh, 1974). Since, linked selection has been shown to impact diversity levels in a variety of species, as evidenced by the correlation between recombination and diversity (Aguade et al., 1989; Begun and Aquadro, 1992; Cutter and Payseur, 2003; Stephan and Langley, 1998; Cai et al., 2009). Theoretic work to explain this pattern has considered the impact of a steady influx of beneficial mutations (recurrent hitchhiking; Stephan et al., 1992; Stephan, 1995), and purifying selection against deleterious mutations (background selection, BGS; Charlesworth et al., 1993; Nordborg et al., 1996; Hudson and Kaplan, 1994). Indeed, empirical work indicates background selection diminishes diversity around genic regions in a variety of species (McVicker et al., 2009; Hernandez et al., 2011; Charlesworth, 1996), and now efforts have shifted towards teasing apart the effects of positive and negative selection on genomic diversity (Elyashiv et al., 2016).

A class of models that are of particular interest in the context of Lewontin's Paradox are recurrent hitchhiking models that decouple diversity from the census population size. These models predict diversity levels when strongly selected beneficial mutations regularly enter and sweep through the population, trapping lineages and forcing them to coalesce (Kaplan et al., 1989; Gillespie, 2000). In general, decoupling occurs under these hitchhiking models when the rate of coalescence due to selection is much greater than the rate of neutral coalescence (e.g. Coop and Ralph, 2012, Equation 22). In contrast, under other linked selection models, the resulting effective population size is proportional to population size; these models cannot decouple diversity, all else equal. For example, models of background selection and polygenic fitness variation predict diversity is proportional to population size, mediated by the total recombination map length and the deleterious mutation rate or fitness variation (Charlesworth et al., 1993; Nicolaisen and Desai, 2012; Nordborg et al., 1996; Robertson, 1961; Santiago and Caballero, 1995).

## Recent approaches towards resolving Lewontin's Paradox

Recently, Corbett-Detig et al., 2015 used population genomic data to estimate the reduction in diversity due to background selection and hitchhiking across 40 species, and showed that the impact of selection increases with two proxies of census population size, species range and the inverse of body size. Based on this evidence, they argued that selection could explain Lewontin's Paradox; however, in a re-analysis, Coop, 2016 demonstrated that the observed magnitude of these reductions is insufficient to explain the orders-of-magnitude shortfall between observed and expected levels of diversity across species. Other recent work has found that life history characteristics related to parental investment, such as propagule size, are good predictors diversity in animals (Romiguier et al., 2014; Chen et al., 2017). Nevertheless, while these diversity correlates are important clues, they do not propose a mechanism by which these traits act to constrain diversity within a few orders of magnitude.

Here, I revisit Lewontin's Paradox by integrating several data sets in order to compare the observed relationship between diversity and census size with the predicted relationship under different selection models. Prior surveys of genetic diversity either lacked census population size estimates, used allozyme-based measures of heterozygosity, or included fewer species. To address these shortcomings, I first estimate census sizes by combining predictions of population density based on body size with ranges estimated from geographic occurrence data. Using these estimates, I quantify the relationship between census size and previously-published genomic diversity estimates across 172 metazoan taxa within nine phyla, thus characterizing the relationship between  $\pi$  and  $N_c$  that underlies Lewontin's Paradox.

Past work looking at the relationship between  $\pi$  and  $N_c$  has been unable to fully account for phylogenetic non-independence across taxa (Felsenstein, 1985). To address this, I use phylogenetic comparative methods (PCMs) with a synthetic time-calibrated phylogeny to account for shared phylogenetic history. Moreover, it is disputed whether considering phylogenetic non-independence is necessary in population genetics, given that coalescent times within species are much less than divergence times (Whitney and Garland, 2010; Lynch, 2011). Using PCMs, I address this by estimating the degree of phylogenetic signal in the diversity census size relationship, and investigating how these traits evolve along the phylogeny.

Finally, I explore whether the predicted reductions of diversity under background selection and recurrent hitchhiking are sufficiently strong to resolve Lewontin's Paradox. I do so using selection parameters from *Drosophila melanogaster*, a species known to be strongly affected by linked selection. Given the effects of linked selection are mediated by recombination map length, I also investigate how map lengths vary with census population size using data from a previously-published survey (Stapley et al., 2017). I find map lengths are typically shorter in large-census-size species, increasing the effects of linked selection in these species, which could further decouple diversity from census size. Still, I find the combined impact of these modes of linked selection fall short in explaining Lewontin's Paradox, and discuss future avenues through which the Paradox of Variation could be fully resolved.

## Results

### Estimates of census population size

An impediment in resolving Lewontin's Paradox is characterizing the relationship between diversity and census population sizes. This is difficult because census population sizes are unavailable for many taxa, especially for extremely abundant, cosmopolitan species that define the upper limit of ranges. Previous work has surveyed the literature for census size estimates (Nei and Graur, 1984; Soulé, 1976; Frankham, 1996), or used range, body size, or qualitative categories as proxies for census size (Corbett-Detig et al., 2015; Leffler et al., 2012). To quantify the relationship between genomic estimates of diversity and census population sizes, I first approximate census population sizes for 172 metazoan taxa (Figure 1). I estimate population densities based on an empirical linear relationship between body sizes and density that holds across metazoans (see Figure 1—figure supplement 1; Damuth, 1981; Damuth, 1987). Then, from geographic occurrence data, I estimate range sizes. Finally, I estimate population size as the product of these predicted densities and range estimates (see Materials and methods: Macroecological Estimates of Population Size). Note that the relationship between population density and body size is driven by energy budgets, and thus reflects macroecological equilibria (Damuth, 1987). Consequently, population sizes are underestimated for taxa like humans and their domesticated species, and overestimated for species with anthropogenically reduced densities or fragmented ranges. For example, the population size of *Lynx lynx* is likely around 50,000 (IUCN, 2020) which is around two orders of magnitude smaller than my estimate. Additionally, the range size estimates do not consider whether an area has unsuitable habitat, and thus may be overestimated for species with particular niches or patchy habitats. While my approach produces approximate and sometimes crude estimates, it has the advantage that it can be efficiently calculated for numerous taxa, which is sufficient to estimate the magnitude of Lewontin's Paradox (see Population Size Validation for more on validation based on biomass and other approaches).

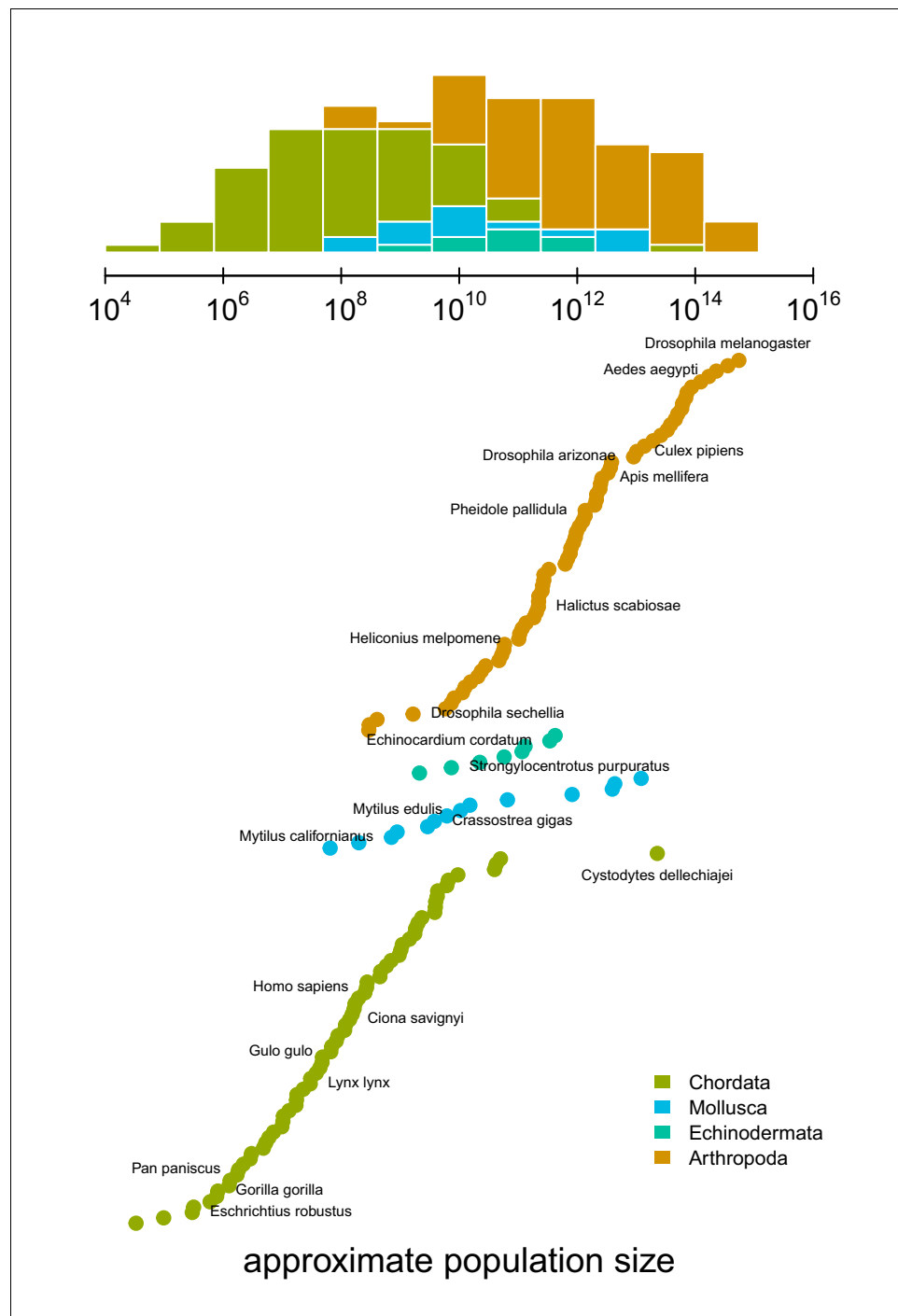
### Characterizing the Diversity–Census-size Relationship

To determine which ecological or evolutionary processes could decouple diversity from census population size, we first need to quantify this relationship across a wide variety of taxa. Previous work has found there is a significant relationship between heterozygosity and the logarithm of population size or range size, but these studies relied on heterozygosity measured from allozyme data (Soulé, 1976; Frankham, 1996; Nei and Graur, 1984). I confirm these findings using pairwise diversity estimates from genomic sequence data and the estimated census sizes (Figure 2). The pairwise diversity estimates are from three sources: Leffler et al., 2012, Corbett-Detig et al., 2015, and Romiguier et al., 2014, and are predominantly from either synonymous or non-coding DNA (see Methods and Materials: 4.1 Diversity and Map Length Data). Overall, an ordinary least squares (OLS) relationship on a log-log scale fits the data well (Figure 2, gray dashed line). The OLS slope estimate is significant and implies a 13% percent increase in differences per basepair for every order of magnitude census size grows (95% confidence interval [12%, 14%], adjusted  $R^2 = 0.26$ ; see also the OLS fit per-phyla, Figure 2—figure supplement 2).

Notably, this relationship has few outliers and is relatively homoscedastic. This is in part because of the log-log scale, in contrast to previous work (Nei and Graur, 1984; Soulé, 1976); see Figure 2—figure supplement 1 for a version on a log-linear scale. However, it is noteworthy that few taxa have diversity estimates below  $10^{-3.5}$  differences per basepair. Those that do, lynx (*Lynx lynx*), wolverine (*Gulo gulo*), and Massasauga rattlesnake (*Sistrurus catenatus*) face habitat loss and declining population sizes. These three species are all in the IUCN Red List, but are listed as least concern (though their presence in the Red List indicates they are of conservation interest). In Appendix D, Appendix D Diversity and IUCN Red List Status, I explore the relationships between IUCN Red List status, diversity, and population size.

### Phylogenetic non-independence and the population size diversity relationship

One limitation of using ordinary least squares is that shared phylogenetic history can create correlation structure in the residuals, which violates an assumption of the regression model and can lead to bias (Felsenstein, 1985; Revell, 2010). To address this shortcoming, I fit the diversity–census-size



**Figure 1.** The distribution of approximate census population sizes estimated by this study. Some phyla containing few species were excluded for clarity.

The online version of this article includes the following source data and figure supplement(s) for figure 1:

**Source data 1.** The population size estimates for 172 metazoan taxa.

**Figure supplement 1.** The relationship between body mass and population density found by *Damuth, 1987*, which is used to predict population densities.

**Figure supplement 2.** The fraction of total species per class on earth included in this study's sample, per class.

**Figure supplement 3.** Comparison of this paper's range estimates procedure against the IUCN Red List's range estimates.

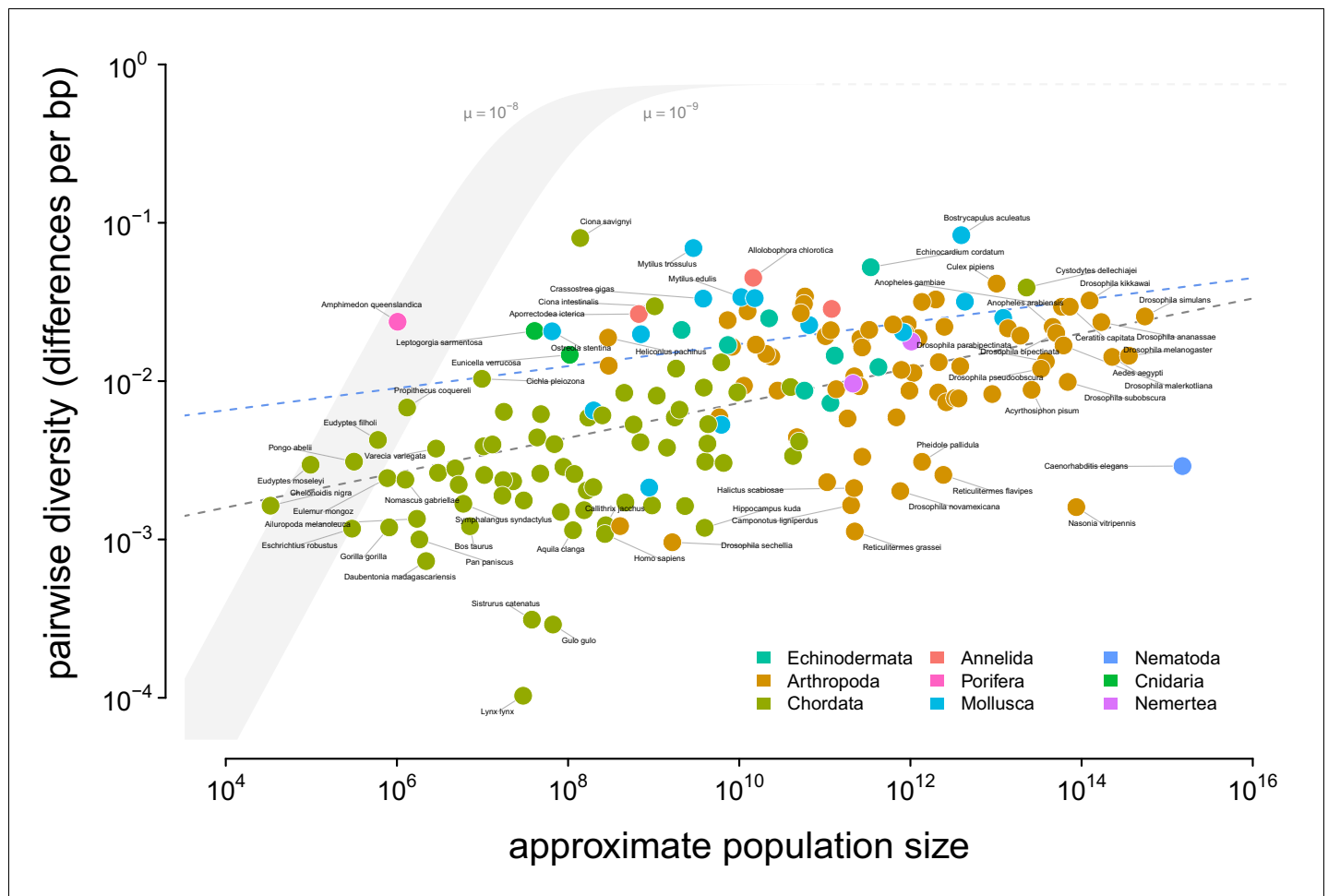
**Figure supplement 4.** Validation of this paper's range estimates against the categorical labels of *Leffler et al., 2012*.

Figure 1 continued on next page

Figure 1 continued

**Figure supplement 5.** The relationship between body length (meters) and body mass (grams) in the Romiguier et al., 2014 data set.

relationship using a phylogenetic mixed-effects model, investigated whether there is a signal of phylogenetic non-independence, estimated the continuous trait values on the phylogeny, and explored how diversity and population size evolve. Prior population genetic comparative studies have lacked time-calibrated phylogenies and assumed unit branch lengths (Whitney and Garland, 2010), a



**Figure 2.** A visualization of Lewontin’s Paradox of Variation. Pairwise diversity (data from Leffler et al., 2012, Corbett-Detig et al., 2015, and Romiguier et al., 2014), which varies over three orders of magnitude, shows a weak relationship with approximate population size, which varies over 12 orders of magnitude. The shaded curve shows the range of expected neutral diversity if  $N_e$  were to equal  $N_c$  under the four-alleles model,  $\log_{10}(\pi) = \log_{10}(\theta) - \log_{10}(1 + 4\theta/3)$  where  $\theta = 4N_c\mu$ , for two mutation rates,  $\mu = 10^{-8}$  and  $\mu = 10^{-9}$ , and the light gray dashed line represents the maximum pairwise diversity under the four alleles model. The dark gray dashed line is the OLS regression fit, and the blue dashed line is the regression fit using a phylogenetic mixed-effects model. Points are colored by phylum. The species *Equus ferus przewalskii* ( $N_e \approx 10^3$  and  $\pi = 3.6 \times 10^{-3}$ ) was an outlier and excluded from this figure for visual clarity.

The online version of this article includes the following source data and figure supplement(s) for figure 2:

**Source data 1.** The diversity and population size dataset for 172 metazoan taxa.

**Figure supplement 1.** A linear-log version of Figure 2.

**Figure supplement 2.** A version of Figure 2 with OLS estimates per phylum.

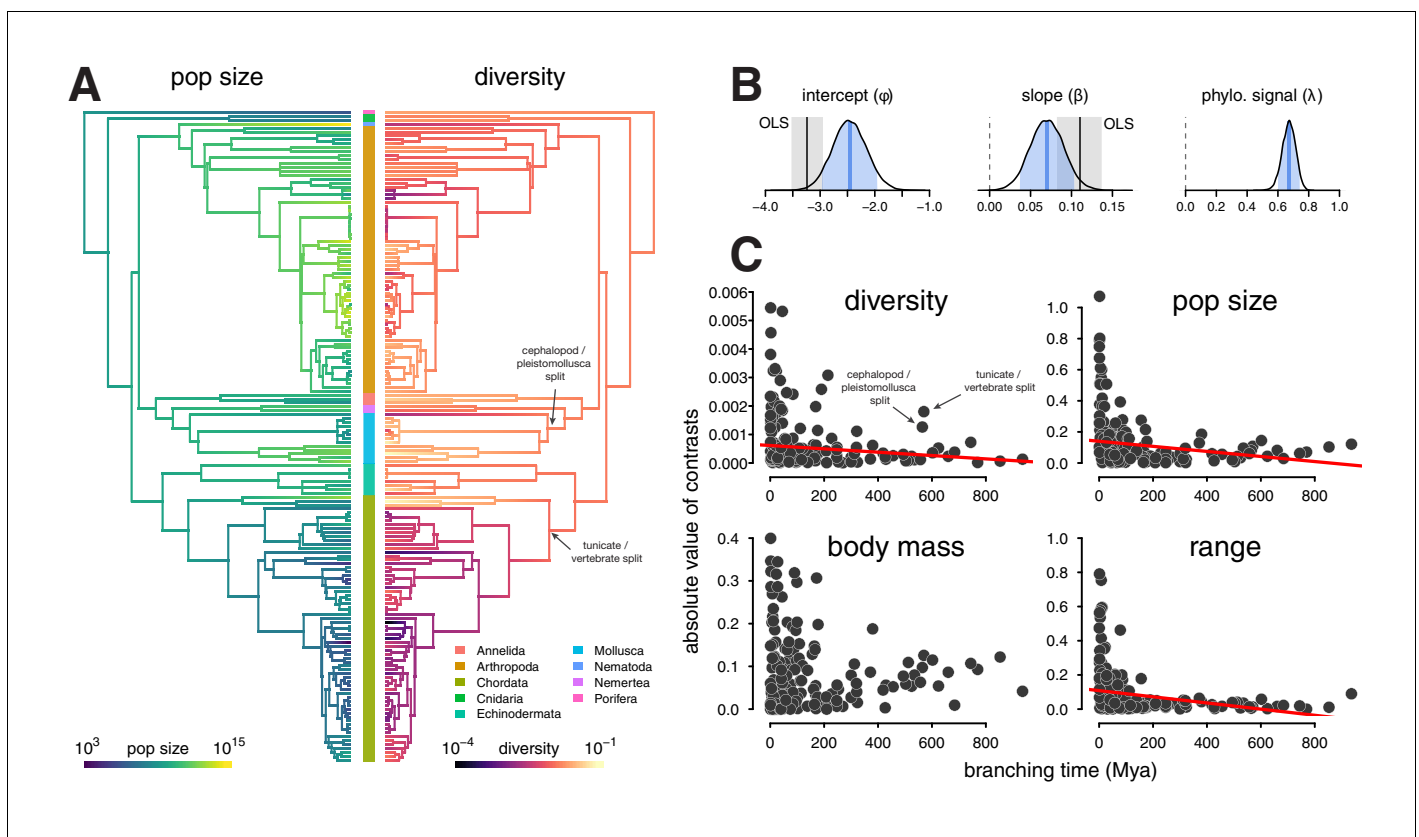
**Figure supplement 3.** The posterior distributions and fitted relationship between diversity and both body mass and range size.

**Figure supplement 4.** Pairwise diversity grouped by the range categories from Leffler et al., 2012, with point size indicating the predicted population density.



shortcoming that has drawn criticism (Lynch, 2011). I use a synthetic time-calibrated phylogeny created from the DateLife project (O'Meara et al., 2020) to account for shared phylogenetic history (see Materials and methods: Phylogenetic Comparative Methods).

Using a phylogenetic mixed-effects model (Lynch, 1991; Hadfield and Nakagawa, 2010; de Villemereuil and Nakagawa, 2014) implemented in Stan (Carpenter et al., 2017; Stan Development Team, 2020), I estimated the linear relationship between diversity and population size (on a log-log scale) accounting for phylogeny, for the 166 taxa without missing data and present in the synthetic chronogram. This type of model is needed because closely-related species may differ from the average trend between  $N_e$  and  $\pi$  in similar ways due to shared phylogenetic history, similar life history traits, etc., and thus do not represent independent observations as is assumed by the standard regression model. This is a form of phylogenetic pseudoreplication, and can be accounted for with a phylogenetic mixed-effects model. The phylogenetic mixed-effects model does not assume that there is phylogenetic structure in either  $N_e$  or  $\pi$  (which itself is not a violation of the standard regression model, Revell, 2010 and Uyeda et al., 2018), but rather accounts



**Figure 3.** Phylogenetic comparative models of diversity and population size. (A) The ancestral continuous trait estimates for the population size and diversity (differences per bp, log scaled) across the phylogeny of 166 taxa. The phyla of the tips are indicated by the color bar in the center. (B) The posterior distributions of the intercept, slope, and phylogenetic signal ( $\lambda$ , de Villemereuil and Nakagawa, 2014) of the phylogenetic mixed-effects model of diversity and population size (log scaled). Also shown are the 90% credible interval (light blue shading), posterior mean (blue line), OLS estimate (gray solid line), and bootstrap OLS confidence intervals (light gray shading). (C) The node-height tests of diversity, population size, and the two components of the population size estimates, body mass, and range (all traits on log scale before contrast was calculated). Each point shows the standardized phylogenetic independent contrast and branching time for a pair of lineages. Red lines are robust regression estimates (and are only shown for statistically significant relationships at the  $\alpha = 0.05$  level). Note that some outlier pairs with very high phylogenetic independent contrasts were excluded (in all cases, these outliers were in the genus *Drosophila*).

The online version of this article includes the following figure supplement(s) for figure 3:

**Figure supplement 1.** The posterior distributions for the parameters of the phylogenetic mixed-effects model of diversity and population size (this is analogous to Figure 3B) fit separately on chordates ( $n = 68$ ), molluscs ( $n = 13$ ), and arthropods ( $n = 68$ ).

**Figure supplement 2.** The ancestral continuous trait estimates for diversity and population size with species labels.

**Figure supplement 3.** The ancestral continuous trait estimates for recombination map length and diversity and population size with species labels.

for phylogenetic correlation structure in the residuals if any is present. Importantly, phylogenetic mixed-effects models simultaneously estimate the degree of phylogenetic structure in the residuals while fitting the relationship between  $N_c$  and  $\pi$ . If the residuals are distributed independently, the estimated relationship would be similar to that found by ordinary least squares, and the estimated phylogenetic signal would be zero. Overall, this approach is conservative, making no assumptions about the source of the phylogenetic signal while accounting for violations of the regression model due to dependence among the residuals if present (see [Revell, 2010](#) for a discussion of this).

As with the linear regression, I find this relationship is positive and significant (95% credible interval 0.03, 0.11), though somewhat attenuated compared to the OLS estimates ([Figure 3B](#)). Since the population size estimates are based on range and body mass, they are essentially a composite trait; fitting phylogenetic mixed-effects models separately on body mass and range indicates these have significant positive and negative effects, respectively ([Figure 2—figure supplement 3](#); see also [Figure 2—figure supplement 4](#) for the relationship between diversity and the range categories of [Leffler et al., 2012](#)).

Since the phylogenetic mixed-effects model simultaneously estimates the variance of the phylogenetic effect ( $\sigma_p^2$ ) and the residual variance ( $\sigma_r^2$ ), these can be used to estimate the phylogenetic signal,  $\lambda = \sigma_p^2 / (\sigma_p^2 + \sigma_r^2)$  ([Lynch, 1991](#); [de Villemereuil and Nakagawa, 2014](#); see [Freckleton et al., 2002](#) for a comparison to Pagel's  $\lambda$ ). When residuals are free of correlations due to shared phylogenetic history, then  $\lambda = 0$  and all the variance could be explained by evolution or noise on the tips. In the relationship between population size and diversity, the posterior mean of  $\lambda = 0.67$  (90% credible interval [0.58, 0.75]) indicates a majority of the variance perhaps might be due to shared phylogenetic history ([Figure 3B](#)).

This high degree of phylogenetic signal substantiates Gillespie's concern ([Gillespie, 1991](#)) that the  $\pi$ - $N_c$  relationship may be driven by chordate-arthropod differences. A visual inspection of the estimated ancestral continuous values for diversity and population size on the phylogeny indicates the high phylogenetic signal seems to be driven in part by chordates having low diversity and small population sizes compared to non-chordates ([Figure 3A](#)). This problem resembles Felsenstein's worst-case scenario ([Felsenstein, 1985](#); [Uyeda et al., 2018](#)), where a singular event on a lineage separating two clades generates a spurious association between two traits.

To investigate whether clade-level differences dominated the relationship between diversity and population size, I fit phylogenetic mixed-effects models to phyla-level subsets of the data for clades with sufficient sample sizes (see Methods: 4.4 Phylogenetic Comparative Methods). This analysis shows a significant positive relationship between diversity and population size in arthropods, and positive weak relationships in molluscs and chordates ([Figure 3—figure supplement 1](#)). Each of the 90% credible intervals for slope overlap, suggesting the relationship between  $\pi$  and  $N_c$  is similar across these clades.

Additionally, I have explored the rate of trait change through time using node-height tests ([Freckleton and Harvey, 2006](#)). Node-height tests regress the absolute values of the standardized contrasts between lineages against the branching time (since present) of these lineages. Under Brownian Motion (BM), standardized contrasts are estimates of the rate of character evolution ([Felsenstein, 1985](#)); if a trait evolves under constant rate BM, this relationship should be flat. For both diversity and population size, node-height tests indicate a significant increase in the rate of evolution towards the present (robust regression p-values 0.023 and 0.00018 respectively; [Figure 3C](#)). Considering the constituents of the population size estimate, range and body mass, separately, the rate of evolution of range but not body mass shows a significant increase (p-value  $1.03 \times 10^{-7}$ ) towards the present.

Interestingly, the diversity node-height test reveals two rate shifts at deeper splits ([Figure 3C](#), top left) around 570 Mya. These nodes represent the branches between tunicates and vertebrates in chordates, and cephalopods and pleistomollusca (bivalves and gastropods) in molluscs. While the cephalopod-pleistomollusca split outlier may be an artifact of having a single cephalopod (*Sepia officinalis*) in the phylogeny, the tunicate-vertebrate split outlier is driven by the low diversity of vertebrates and the previously-documented exceptionally high diversity of tunicates (sea squirts; [Nydam and Harrison, 2010](#); [Small et al., 2007](#)). This deep node representing a rate shift in diversity could reflect a change in either effective population size or mutation rate, and there is some evidence of both in this genus *Ciona* ([Small et al., 2007](#); [Tsagkogeorga et al., 2012](#)). Neither of these



deep rate shifts in diversity is mirrored in the population size node-height test (**Figure 3C**, top right). Rather, it appears a trait impacting diversity but not census size (e.g. mutation rate or offspring distributions) has experienced a shift on the lineage separating tunicates and vertebrates. At nearly 600 Mya, these deep nodes illustrate that expected effective population sizes (and thus coalescence times) can share phylogenetic history, due to phylogenetic inertia in some combination of population size, reproductive system, and mutation rates.

Finally, an important caveat is the increase in rate towards the tips could be caused by measurement noise, or possibly uncertainty or bias in the divergence time estimates deep in the tree. Inspecting the lineage pairs that lead to this increase in rate towards the tips indicates these represent plausible rate shifts, e.g. between cosmopolitan and endemic sister species like *Drosophila simulans* and *Drosophila sechellia*; however, ruling out measurement noise entirely as an explanation would involve modeling the uncertainty of diversity and population size estimates.

### Assessing the impact of linked selection on diversity across taxa

The above analyses reemphasize the drastic shortfall of diversity levels as compared to census sizes. Linked selection has been proposed as the mechanism that acts to reduce diversity levels from what we would expect given census sizes (**Smith and Haigh, 1974; Gillespie, 2000; Corbett-Detig et al., 2015**). Here, I test this hypothesis by estimating the scale of diversity reductions expected under background selection and recurrent hitchhiking, and comparing these to the observed relationship between  $\pi$  and  $N_c$ .

I quantify the effect of linked selection on diversity as the ratio of observed diversity ( $\pi$ ) to the estimated diversity in the absence of linked selection ( $\pi_0$ ),  $R = \pi/\pi_0$ . Here,  $\pi_0$  would reflect only demographic history and non-heritable variation in reproductive success. There are two difficulties in evaluating whether linked selection could resolve Lewontin's Paradox. The first difficulty is that  $\pi_0$  is unobserved. Previous work has estimated  $\pi_0$  using methods that exploit the spatial heterogeneity in recombination and functional density across the genome to fit linked selection models that incorporate both hitchhiking and background selection (**Elyashiv et al., 2016; Corbett-Detig et al., 2015**). The second difficulty is understanding how  $R$  varies across taxa, since we lack estimates of critical model parameters for most species. Still, I can address a key question: if diversity levels were determined by census sizes ( $\pi_0 = 4N_c\mu$ ), would the combined effects of background selection and recurrent hitchhiking be sufficient to reduce diversity to observed levels? Furthermore, does the relationship between census size and predicted diversity under linked selection across species,  $\pi_{BGS+HH} = R\pi_0$ , match the observed relationship in **Figure 2**?

Since we lack estimates of selection parameters across species, I parameterize the hitchhiking and BGS models using estimates from *Drosophila melanogaster*, a species known to be strongly affected by linked selection (**Sella et al., 2009**). Under a generalized model of hitchhiking and background selection (**Elyashiv et al., 2016; Coop and Ralph, 2012**) and assuming  $N_e = N_c$ , the expected diversity is

$$\pi_{BGS+HH} \approx \frac{\theta}{1/B(U,L) + 2N_c S(\gamma,J,L)} \quad (1)$$

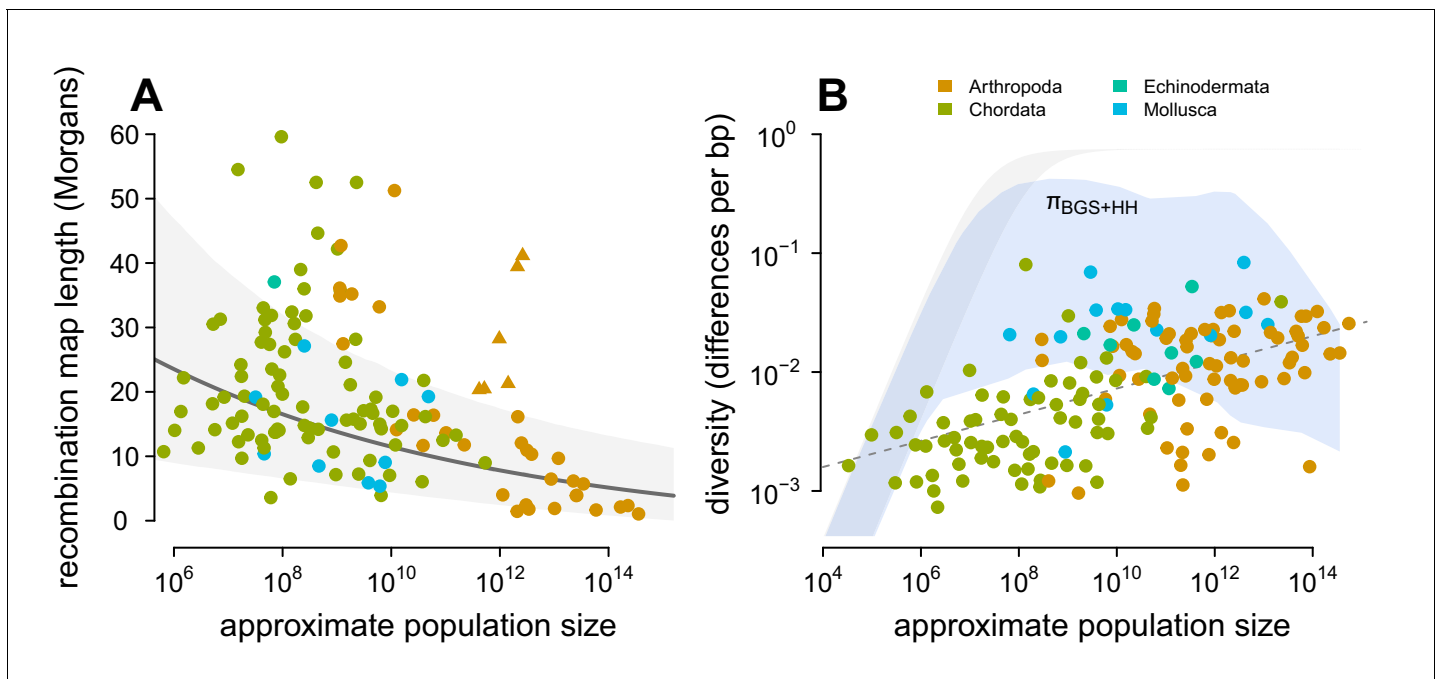
where  $\theta = 4N_c\mu$ ,  $B(U,L)$  is the effect of background selection, and  $S(\gamma,J,L)$  is the rate of coalescence caused by sweeps (**Elyashiv et al., 2016**, Equation 1, **Coop and Ralph, 2012**, Equation 20). Under background selection models with recombination, the reduction is  $B(U,L) = \exp(-U/L)$  where  $U$  is the per diploid genome per generation deleterious mutation rate, and  $L$  is the recombination map length in Morgans (**Hudson and Kaplan, 1994; Nordborg et al., 1996**). This BGS model is similar to models of effective population size under polygenic fitness variation, and can account for other modes of linked selection (**Robertson, 1961; Santiago and Caballero, 1995; Santiago and Caballero, 1998**, see Appendix 2, Background Selection and Polygenic Fitness Models). The coalescence rate due to sweeps is  $S(\gamma,J,L) = \frac{\gamma}{L}J$ , where  $\gamma$  is the number of adaptive substitutions per generation, and  $J$  is the probability a lineage is trapped by sweeps as they occur across the genome ( $J_{2,2}$  in Equation 15 of **Coop and Ralph, 2012**).

Parameterizing the model this way, I then set the key parameters that determine the impact of recurrent hitchhiking and background selection ( $\gamma$ ,  $J$ , and  $U$ ) to strong selection values estimated for *Drosophila melanogaster* by **Elyashiv et al., 2016**. My estimate of the adaptive substitutions per

generation ( $\gamma_{D_{mel}} \approx 2.3 \times 10^{-3}$ ) based Elyashiv et al. implies a rate of sweeps per basepair of  $\nu_{BP, D_{mel}} \approx 2.34 \times 10^{-11}$ , which is close to other estimates from *D. melanogaster* (see **Figure 4—figure supplement 5A**). The rate of deleterious mutations per diploid genome, per generation is parameterized using the estimate from Elyashiv et al.,  $U_{D_{mel}} = 1.6$ , which is slightly greater than previous estimates based on Bateman-Mukai approaches (Mukai, 1985; Mukai, 1988; Charlesworth, 1987). Finally, the probability that a lineage is trapped in a sweep,  $J_{D_{mel}} \approx 4.5 \times 10^{-4}$ , is calculated from the estimated genome-wide average coalescence rate due to sweeps from Elyashiv et al. (see **Figure 4—figure supplement 5B** and Materials and methods: Predicted Reductions in Diversity for more details on parameter estimates). Using these parameters, I then explore how the predicted range of diversity levels varies across species with recombination map length ( $L$ ) and census population size ( $N_c$ ).

Previous work has found that the impact of linked selection increases with  $N_c$  (Corbett-Detig et al., 2015; see **Figure 4—figure supplement 4A**), and it is often thought that this is driven by higher rates of adaptive substitutions in larger populations (Ohta, 1992), despite equivocal evidence (Galtier, 2016). However, there is another mechanism by which species with larger population sizes might experience a greater impact of linked selection: recombination map length,  $L$ , is known to correlate with body mass (Burt and Bell, 1987) and thus varies inversely with population size. As this is a critical parameter that determines the genome-wide impact of both hitchhiking and background selection, I examine the relationship between recombination map length ( $L$ ) and census population size ( $N_c$ ) across taxa, using available estimates of map lengths across species (Stapley et al., 2017; Corbett-Detig et al., 2015). I find a significant non-linear relationship using phylogenetic mixed-effects models (**Figure 4A**; see Methods and materials: 4.4 Phylogenetic Comparative Methods). There is also a correlation between map length and genome size (**Figure 4—figure supplement 2**) and genome size and population size (**Figure 4—figure supplement 1**). These findings are consistent with the hypothesis that non-adaptive processes increase genome size in small- $N_e$  species (Lynch and Conery, 2003) which in turn could increase map lengths, as well as the hypothesis that map lengths are adaptively longer to more efficiently select against deleterious alleles in smaller populations (Roze, 2021). Overall, the negative relationship between map length and census size indicates linked selection is expected to be stronger in species with short map lengths, which are high- $N_c$  species.

Then, I predict the expected diversity ( $\pi_{BGS+HH}$ ) under background selection and hitchhiking, assuming  $N_e = N_c$  and that all species had the rate of sweeps and strength of BGS as *D. melanogaster*. Since neutral mutation rates  $\mu$  are unknown and vary across species, I calculate the range of predicted  $\pi_{BGS+HH}$  estimates for  $\mu = 10^{-9}$ – $10^{-8}$  (using the four-alleles model, Tajima, 1996), and compare this to the observed relationship between  $\pi$  and  $N_c$  in **Figure 4B**. Under these parameters and assumptions, linked selection begins to appreciably constrain diversity for  $N_c \geq 10^7$ , since  $S(\gamma_{D_{mel}}, J_{D_{mel}}, L) \approx 10^{-8}$ – $10^{-7}$  and linked selection dominates drift when  $S(\gamma, J, L) > 1/2N$ . Overall, this reveals two problems for the hypothesis that linked selection could solve Lewontin's Paradox. First, low to mid- $N_c$  species (census sizes between  $10^4$ – $10^7$ ) have sufficiently long map lengths that their diversity levels are only moderately reduced by linked selection, leading to a wide gap between predicted and observed diversity levels. For this not to be the case, the rate of adaptive mutations or the deleterious mutation rate would need to be orders of magnitude higher for species within this range than in *Drosophila melanogaster*, which is incompatible with the rate of adaptive protein substitutions across species (Galtier, 2016) and overall mutation rates (Lynch, 2010). Furthermore, linked selection has been quantified in humans, which fall in this census size range, and has been found to be relatively weak (McVicker et al., 2009; Hernandez et al., 2011; Hellmann et al., 2008; Cai et al., 2009; Boyko et al., 2008). Second, while hitchhiking and BGS can reduce predicted diversity levels for high- $N_c$  species ( $N_c > 10^{12}$ ) to observed levels, this would imply available estimates of  $\pi_0$  are underestimated by several orders of magnitude in *Drosophila* (**Figure 4—figure supplement 4B**). The high reductions in  $\pi$  predicted here (compared to those of Elyashiv et al., 2016) are a result of using  $N_c$ , rather than  $N_e = \pi_0/4\mu$  in the denominator of **Equation (1)**, which leads to a very high rate of sweeps in the population. I do not consider selective interference, though the saturation of adaptive substitutions per Morgan would only act to limit the reduction in diversity (Weissman and Barton, 2012), and thus these results are conservative.



**Figure 4.** Predicting the impact of linked selection on diversity. (A) The observed relationship between recombination map length ( $L$ ) and census size ( $N_c$ ) across 136 species with complete data and known phylogeny. Triangle points indicate six social taxa excluded from the model fitting since these have adaptively higher recombination map lengths (Wilfert et al., 2007). The dark gray line is the estimated relationship under a phylogenetic mixed-effects model, and the gray interval is the 95% posterior average. (B) Points indicate the observed  $\pi$ - $N_c$  relationship across taxa shown in Figure 2, and the blue ribbon is the range of predicted diversity were  $N_e = N_c$  for  $\mu = 10^{-8}$ - $10^{-9}$ , and after accounting for the expected reduction in diversity due to background selection and recurrent hitchhiking under *Drosophila melanogaster* parameters. In both plots, point color indicates phylum.

The online version of this article includes the following source data and figure supplement(s) for figure 4:

**Source data 1.** The map length, population size, and linked selection estimates for 136 metazoan taxa.

**Figure supplement 1.** The relationship between genome size and approximate census population size.

**Figure supplement 2.** The relationship between genome size and recombination map length.

**Figure supplement 3.** The observed  $\pi$ - $N_c$  relationship (points) across species compared to the predicted diversity (ribbons) under different modes of linked selection and parameters, for a range of mutation rates,  $10^{-9}$ - $10^{-8}$ .

**Figure supplement 4.** The relationship between  $N_c$  and diversity in the Corbett-Detig et al., 2015 data, and the relationship between estimated reduction in diversity and census size, for three different approaches.

**Figure supplement 5.** Comparison of the *Drosophila* sweep parameters used in this study with parameters from other studies.

Finally, the poor fit between observed and predicted levels of diversity across species is not remedied by stronger selection parameters. In Figure 4—figure supplement 3B, I increase both selection parameters  $U$  and  $\gamma$  ten-fold each, and find the same qualitative pattern: on a log-log scale the relationship between  $N_c$  and  $\pi$  is linear, while the predicted diversity under linked selection is non-linear with  $N_c$ . Under this ten-fold higher selection regime, there is more overlap between observed and predicted levels of diversity, but diversity is severely under-predicted for high- $N_c$  species. Additionally, this would imply that selection in low-to-mid- $N_c$  species is ten-fold higher than estimated in *Drosophila melanogaster*, which is implausible. Overall, this suggests that present models of linked selection, even with very strong selection across species, are qualitatively incapable of matching the observed relationship between  $N_c$  and  $\pi$  and thus cannot explain Lewontin's Paradox.

## Discussion

Nearly fifty years after Lewontin's description of the Paradox of Variation, how evolutionary, life history, and ecological processes interact to constrain diversity across taxa to a narrow range remains a mystery. I revisit Lewontin's Paradox by first characterizing the relationship between genomic estimates of pairwise diversity and approximate census population size across 172 metazoan species. Previous surveys have used allozyme-based estimates, fewer taxa, or proxies of population size. My

estimates of census population sizes are rough approximates, since they use body size to predict density. An improved estimate might account for vagility (as *Soulé, 1976* did), though this is harder to do systematically across many taxa. Future work might also use other ecological information, such as total biomass, or species distribution modeling to improve census size estimates (*Bar-On et al., 2018; Mora et al., 2011*). Still, it seems more accurate estimates would be unlikely to change the qualitative findings here, which resemble those of early surveys (*Nei and Graur, 1984; Soulé, 1976*).

One limitation of this study is that diversity estimates are collated from a variety of sources rather than estimated with a single bioinformatic pipeline. This leads to technical noise across diversity estimates; perhaps the relationship between  $\pi$  and  $N_c$  found here could be tighter with a standardized bioinformatic pipeline. In addition, there might be systematic bioinformatic sources of bias: for example high-diversity sequences may fail to align to the reference genome and end up unaccounted for, leading to a downward bias. Alternatively, a high-diversity sequences might map to the reference genome, but adjacent mis-matching SNPs might be mistaken for a short insertion or deletion. While these issues might affect estimates in high-diversity species, it is unlikely to change the qualitative relationship between  $\pi$ - $N_c$ .

### Macroevolution and Across-Taxa population genomics

Lewontin's Paradox arises from a comparison of diversity across species, yet it has been disputed whether such comparisons require phylogenetic comparative methods. Extending previous work that has accounted for phylogeny in particular clades (*Leffler et al., 2012*), or using taxonomical-level averages (*Romiguier et al., 2014*), I show that the positive relationship between diversity and census size is significant using a mixed-effects model with a time-calibrated phylogeny. Additionally, I find a high degree of phylogenetic signal, evidence of deep shifts in the rate of evolution of genetic diversity, and that arthropods and chordates form clusters. Overall, this suggests that previous concerns about phylogenetic non-independence in comparative population genetic studies were warranted (*Gillespie, 1991; Whitney and Garland, 2010*). Notably, *Lynch, 2011* has argued that PCMs for pairwise diversity are unnecessary, since mutation rate evolution is fast and thus free of phylogenetic inertia, sampling variance should exceed the variance due to phylogenetic shared history, and coalescence times are much shorter than divergence times. Since my findings suggest PCMs are necessary in some cases, it is worthwhile to address these points.

First, Lynch has correctly pointed out that while coalescence times are much less than divergence times and should be free of phylogenetic shared history, the factors that determine coalescence times (e.g. mutation rates and effective population size) may not be (*Lynch, 2011*). In other words, coalescence times are free from phylogenetic shared history *were we to condition* on these causal factors that could be affected by shared phylogenetic history. My estimates of phylogenetic signal in the residuals, by contrast, are not conditioned on these factors. Importantly, even "correcting for" phylogeny implicitly favors certain causal interpretations over others (*Westoby et al., 1995; Uyeda et al., 2018*). Future work could try to untangle what causal factors determine coalescence times across species, as well as how these factors evolve across macroevolutionary timescales. Second, it is a misconception that a fast rate of trait evolution necessarily reduces phylogenetic signal (*Revell et al., 2008*), and that if either or both variables in a regression are free of phylogenetic signal, PCMs are unnecessary (*Revell, 2010; Uyeda et al., 2018*). The evidence of high phylogenetic signal found in this study suggests PCMs are necessary when fitting the relationship between  $N_c$  and  $\pi$  in order to account for correlated residuals among closely-related species, and to avoid spurious results from phylogenetic pseudoreplication.

Finally, beyond just accounting for phylogenetic non-independence, macroevolution and phylogenetic comparative methods are a promising way to approach across-species population genomic questions. For example, one could imagine that diversification processes could contribute to Lewontin's Paradox. If large- $N_c$  species were to have a rate of speciation that is greater than the rate at which mutation and drift reach equilibrium (which is indeed slower for large  $N_c$  species), this could act to decouple diversity from census population size. That is to say, even if the rate of random demographic bottlenecks were constant across taxa, lineage-specific diversification processes could lead certain clades to be systematically further from demographic equilibrium, and thus have lower diversity than expected for their census population size.

## How could selection still explain Lewontin's Paradox?

Even assuming selection parameters estimated from *Drosophila melanogaster*, where the effects of linked selection are thought to be especially strong, the predicted patterns of diversity under linked selection poorly fit observed patterns of diversity across species. My results support the analysis by [Coop, 2016](#) showing that levels of  $\pi_0$  estimated by [Corbett-Detig et al., 2015](#) are not decoupled from genome-wide average  $\pi$ , as would occur if linked selection were to explain Lewontin's Paradox. Additionally, my analysis goes a step further, showing that current linked selection models under a wide range of selection parameters are incapable of explaining the observed relationship between census size and diversity. This is in part because mid- $N_c$  species have sufficiently long recombination map lengths to diminish the effects of even strong selection. Overall, while this suggests hitchhiking and background selection seem unlikely to explain patterns of diversity across taxa, there are three major potential limitations of my approach that need further evaluation.

First, I approximate the reduction in diversity using homogeneous background selection and recurrent hitchhiking models ([Kaplan et al., 1989](#); [Hudson and Kaplan, 1995](#); [Coop and Ralph, 2012](#)), when in reality, there is genome-wide heterogeneity in functional density, recombination rates, and the adaptive substitutions across species. Each of these factors mediate how strongly linked selection impacts diversity across the genome. Despite these model simplifications, the predicted reduction in diversity in *Drosophila melanogaster* is 85% (when using  $N_e$ , not  $N_c$ ), which is reasonably close to the estimated 77% from the more realistic model of [Elyashiv et al.](#) that accounts for the actual position of substitutions, annotation features, and recombination rate heterogeneity (though it should be noted that these both use the same parameter estimates). Furthermore, even though my model fails to capture the heterogeneity of functionality density and recombination rate in real genomes, it is still conservative, likely overestimating the effects of linked selection to see if it could be capable of decoupling diversity from census size and explain Lewontin's Paradox. This is in part because the strong selection parameter estimates from *Drosophila melanogaster* used, but also because I assume that the effective population size is equal to the census size. Even then, this decoupling only occurs in very high-census-size species, and implies that the diversity in the absence of linked selection,  $\pi_0$ , is currently underestimated by several orders of magnitude. Moreover, the study of [Corbett-Detig et al., 2015](#) did consider recombination rate and functional density heterogeneity in estimating the reduction due to linked selection across species, yet their predicted reductions are orders of magnitude weaker than those considered here by assuming that  $N_e = N_c$  ([Figure 4—figure supplement 4B](#)). Overall, given the effects estimated under more realistic inference models are still orders of magnitude weaker than those used in this study, current models of linked selection seem fundamentally unable to fit the diversity-census-size relationship.

Second, my model here only considers hard sweeps, and ignores the contribution of soft sweeps (e.g. from standing variation or recurrent mutations; [Hermisson and Pennings, 2005](#); [Pennings and Hermisson, 2006](#)), partial sweeps (e.g. those that do not reach fixation), and the interaction of sweeps and spatial processes. While future work exploring these alternative types of sweeps is needed, the predicted reductions in diversity found here under the simplified sweep model are likely relatively robust to these other modes of sweeps for a few reasons. First, the shape of the diversity-recombination curve is equivalent under models of partial sweeps and hard sweeps, though these imply different rates of sweeps ([Coop and Ralph, 2012](#)). Second, in the limit where most fitness variation is due to weak soft sweeps from standing variation scattered across the genome (i.e. due to polygenic fitness variation), levels of diversity are well approximated by quantitative genetic linked selection models ([Robertson, 1961](#); [Santiago and Caballero, 1995](#)). The reduction in diversity under these models is nearly identical to that under background selection models, in part because deleterious alleles at mutation-selection balance constitute a considerable component of fitness variation (see Appendix Section B; [Charlesworth and Hughes, 2000](#); [Charlesworth, 2015](#)). Third, the parameters from [Elyashiv et al., 2016](#) could reflect a mixture of types of sweeps ([Elyashiv et al., 2016](#) p. 14 and p. 19 of their Supplementary Online Materials). Finally, I also disregarded the interaction of sweeps and spatial processes. For populations spread over wide ranges, limited dispersal slows the spread of sweeps, allowing for new beneficial alleles to arise, spread, and compete against other segregating beneficial variants ([Ralph and Coop, 2015](#); [Ralph and Coop, 2010](#)). Through limited dispersal should act to "soften sweeps" and not impact my findings for the reasons described



above, future work could investigate how these processes impact diversity in ways not captured by hard sweep models.

Third, other selective processes, such as fluctuating selection or hard selective events (i.e. selection resulting in a reduction in the population size), could reduce diversity in ways not captured by the background selection and hitchhiking models. Since frequency-independent fluctuating selection reduces diversity under most conditions (Novak and Barton, 2017), this could lead seasonality and other sources of temporal heterogeneity to reduce diversity in large- $N_c$  species with short generation times more than longer-lived species with smaller population sizes. Future work could consider the impact of fluctuating selection on diversity under simple models (Barton, 2000) if estimates of key parameters governing the rate of such fluctuations were known across taxa. Additionally, another mode of selection that could severely reduce diversity across taxa, yet remains unaccounted for in this study, is periodic hard selective events. These selective events could occur regularly in a species' history yet be indistinguishable from demographic bottlenecks with just population genomic data.

### Spatial and demographic processes

One limitation of this study is the inability to quantify the impact of spatial and demographic population genetic processes on the relationship between diversity and census population sizes across taxa. The genomic diversity estimates collated in this study unfortunately lack details about the sampling process and spatial data, which can have a profound impact on population genomic summary statistics (Battey et al., 2020). These issues could systematically bias species-wide diversity estimates; for example, if diversity estimates from a cosmopolitan species were primarily from a single region or subpopulation, diversity would be an underestimate relative to the entire population. However, biased spatial sampling alone seems incapable of explaining the  $\pi$ - $N_c$  divergence in high- $N_c$  taxa. In the extreme scenario in which only one subpopulation was sampled,  $F_{ST}$  would need to be close to one for population subdivision alone to sufficiently reduce the total population heterozygosity to explain the orders-of-magnitude shortfall between predicted and observed diversity levels. This can be seen by rearranging the expression for  $F_{ST}$  as  $H_S = (1 - F_{ST})H_T$ , where  $H_S$  and  $H_T$  are the subpopulation and total population heterozygosities; if  $H_T = 4N_c\mu$ , then only  $F_{ST} \approx 1$  can reduce  $H_S$  several orders of magnitude. Yet, across-taxa surveys indicate that  $F_{ST}$  is almost never this high within species (Roux et al., 2016). Future work could quantify the extent to which more realistic spatial processes contribute to Lewontin's Paradox. For example, high- $N_c$  taxa usually experience range expansions, with repeated founder effects and local extinction/recolonization dynamics that depress diversity (Slatkin, 1977). In particular, with the appropriate data, one could estimate the empirical relationship between dispersal distance, range size, and coalescent effective population size across taxa.

In this study, I have focused entirely on assessing the role of linked selection, rather than demography, in reducing diversity across taxa. In contrast to demographic models, models of linked selection have comparatively fewer parameters and more readily permit rough estimates of diversity reductions across taxa. Given that I find that models of linked selection are incapable of explaining the observed relationship between  $N_c$  and  $\pi$ , this supports the hypothesis the diversity across species are shaped primarily by past demographic fluctuations. Still, a full resolution of Lewontin's Paradox would require understanding how the demographic processes across taxa with incredibly heterogeneous ecologies and life histories transform  $N_c$  into  $N_e$ . With population genomic data becoming available for more species, this could involve systematically inferring the demographic histories of tens of species and looking for correlations in the frequency and size of bottlenecks with  $N_c$  across species.

### Measures of effective population size, Timescales, and Lewontin's Paradox

Lewontin's Paradox describes the extent to which the effective population sizes implied by diversity,  $\tilde{N}_e$ , diverge from census population sizes. However, there are a variety other effective population size estimators calculable from different data and summary statistics (Wang et al., 2016; Caballero, 1994; Galtier and Rousselle, 2020). These include estimators based on the site frequency spectrum, observed decay in linkage disequilibrium, or temporal estimators that use the variance in allele frequency change through time. These various estimators capture different summaries of



effective population size on shorter timescales than coalescent-based estimators (see **Wang, 2005** for a review), and thus could be used to tease apart processes that impact the  $N_e$ - $N_c$  relationship in the more recent past.

Temporal  $N_e$  estimators already play an important role in understanding another summary of the  $N_e$ - $N_c$  relationship: the ratio  $N_e/N_c$ , which is an important quantity in conservation genetics (**Frankham, 1995; Mace and Lande, 1991**) and in understanding evolution in highly fecund marine species. Surveys of the short-term  $N_e/N_c$  relationship across taxa indicate mean  $N_e/N_c$  is on order of  $\approx 0.1$  (**Frankham, 1995; Palstra and Ruzzante, 2008; Palstra and Fraser, 2012**), though the uncertainty in these estimates is high, and some species with sweepstakes reproduction systems like Pacific Oyster (*Crassostrea gigas*) can have  $N_e/N_c \approx 10^{-6}$  (**Hedgcock, 1994**). Estimates of the  $N_e/N_c$  ratio may be an important, yet under appreciated piece of solving Lewontin's Paradox. For example, if  $N_e$  is estimated from the allele frequency change across a single generation (i.e. **Waples, 1989**),  $N_e/N_c$  constrains estimates of the variance in reproductive success (**Wright, 1938; Nunney, 1993; Nunney, 1996**). This implies that apart from species with sweepstakes reproductive systems, the variance in reproductive success each generation (whether heritable or non-heritable) is likely insufficient to significantly contribute to constraining  $\tilde{N}_e$  for most taxa. Still, further work is needed to characterize (1) how  $N_e/N_c$  varies with  $N_c$  across taxa (though see **Palstra and Fraser, 2012, Figure 2**), and (2) the variance of  $N_e/N_c$  over longer time spans (i.e. how periodic sweepstakes reproductive events act to constrain  $N_e$ ). Overall, characterizing how  $N_e/N_c$  varies across taxa and correlates with ecology and life history traits could provide clues into the mechanisms that leads propagule size and survivorship curves to be predictive of diversity levels across taxa (**Romiguier et al., 2014; Hallatschek, 2018; Barry et al., 2020**).

Finally, short-term temporal  $N_e$  estimators may play an important role in resolving Lewontin's Paradox. These estimators, along with short-term estimates of the impact of linked selection (**Buffalo and Coop, 2019; Buffalo and Coop, 2020**), can inform us how much diversity is depressed by selection on shorter timescales, free from the rare strong selective events or severe bottlenecks that impact pairwise diversity. It could be that in any one generation, selection contributes more to the variance of allele frequency changes than drift, yet across-taxa patterns in diversity are better explained processes acting sporadically on longer timescales, such as colonization, founder effects, and bottlenecks. Thus, the pairwise diversity may not give us the best picture of the generation to generation evolutionary processes acting in a population to change allele frequencies. Furthermore, certain observed adaptations occur at a pace that is inexplicable given small effective population sizes implied by diversity, and are only possible if short-term effective population sizes are orders of magnitude larger (**Karasov et al., 2010; Barton, 2010**).

## Conclusions

In *Building a Science of Population Biology* (**Lewontin et al., 2004**), Lewontin laments the difficulty of uniting population genetics and population ecology into a cohesive discipline of population biology. Lewontin's Paradox of Variation remains a major unsolved problem at the nexus of these two different disciplines: we fail to understand the processes that connect a central parameter of population ecology, census size, to a central parameter of population genetics, effective population size across species. Given that selection seems to fall short in resolving Lewontin's Paradox, a full resolution will require a mechanistic understanding the ecological, life history, and macroevolutionary processes that connect  $N_c$  to  $N_e$  across taxa. While I have focused exclusively on metazoan taxa since their population densities are more readily approximated from body mass, a full resolution must also include plant species (with the added difficulties of variation in selfing rates, different dispersal strategies, pollination, etc.).

Looking at Lewontin's Paradox through an macroecological and macroevolutionary lens begets interesting questions outside of the traditional realm of population genetics. Here, I have found that diversity and  $N_c$  have a consistent relationship without many outliers, despite the wildly disparate ecologies, life histories, and evolutionary histories of the taxa included. Furthermore, taxa with very large census sizes have surprisingly low diversity. Is this explained by macroevolutionary processes, such as different rates of speciation for large- $N_c$  taxa? Or, are the levels of diversity we observe today an artifact of our timing relative to the last glacial maximum, or the last major extinction? Did large- $N_c$  prehistoric animal populations living in other geological eras have higher levels of diversity

than our present taxa? Or, does ecological competition occur on shorter timescales such that strong population size contractions transpire and depress diversity, even if a species is undisturbed by climatic shifts or mass extinctions? Overall, patterns of diversity across taxa are determined by many overlaid evolutionary and ecological processes occurring on vastly different timescales. Lewontin's Paradox of Variation may persist unresolved for some time because the explanation requires synthesis and model building at the intersection of all these disciplines.

## Materials and methods

### Diversity and map length data

The data used in this study are collated from a variety of previously published surveys. Of the 172 taxa with diversity estimates, 14 are from *Corbett-Detig et al., 2015*, 96 are from *Leffler et al., 2012*, and 62 are from *Romiguier et al., 2014*. The Corbett-Detig et al. data is estimated from four-fold degenerate sites, the Romiguier et al. data is synonymous sites, and the Leffler et al. data is estimated predominantly from silent, intronic, and non-coding sites. All types of diversity estimates from *Leffler et al., 2012* were included to maximize the taxa in the study, since the variability of diversity across functional categories is much less than the diversity across taxa. Multiple diversity estimates per taxa were averaged. The total recombination map length data were from both (*Stapley et al., 2017*; 127 taxa), and (*Corbett-Detig et al., 2015*; 9 taxa). Both studies used sex-averaged recombination maps estimated with cross-based approaches; in some cases errors in the original data were found, documented, and corrected. These studies also included genome size estimates used to create *Figure 4—figure supplement 2* and *Figure 4—figure supplement 1*.

### Macroecological estimates of population size

A rough approximation for total population size (census size) is  $N_c = DR$ , where  $D$  is the population density in individuals per km<sup>2</sup> and  $R$  is the range size in km<sup>2</sup>. Since population density estimates are not available for many taxa included in this study, I used the macroecological abundance-body size relationship to predict population density from body size. Since body length measurements are more readily available than body mass, I collated body length data from various sources (see [https://github.com/vsbuffalo/paradox\\_variation](https://github.com/vsbuffalo/paradox_variation); copy archived at [swh:1:rev:8fa6b5834f6536319-b1e5cd9722ca02d317183df](https://www.swh.io/rev/8fa6b5834f6536319-b1e5cd9722ca02d317183df), *Buffalo, 2021*); body lengths were averaged across sexes for sexually dimorphic species, and if only a range of lengths was available, the midpoint was used.

Then, I re-estimated the relationship between body mass and population density using the data in the appendix table of *Damuth, 1987*, which includes 696 taxa with body mass and population density measurements across mammals, fish, reptiles, amphibians, aquatic invertebrates, and terrestrial arthropods. Though the abundance-body size relationship can be noisy at small spatial or phylogenetic scales (Chapter 5, *Gaston and Blackburn, 2008*), across deeply diverged taxa such as those included in this study and *Damuth, 1987*, the relationship is linear and homoscedastic (see *Figure 1—figure supplement 1*). Using Stan (*Stan Development Team, 2020*), I jointly estimated the relationship between body mass from body length using the *Romiguier et al., 2014* taxa, and used this relationship to predict body mass for the taxa in this study. These body masses were then used to predict population density simultaneously, using the *Damuth, 1981* relationship. The code of this routine (`pred_popsizemissing_centered.stan`) is available in the GitHub repository ([https://github.com/vsbuffalo/paradox\\_variation/](https://github.com/vsbuffalo/paradox_variation/)).

To estimate range, I first downloaded occurrence records from Global Biodiversity Information Facility (*Global Biodiversity Information Facility, 2020*) using the `rgbif` R package (*Chamberlain et al., 2014*; *Chamberlain and Boettiger, 2017*). Using the occurrence locations, I inferred whether a species was marine or terrestrial, based on whether the majority of their recorded occurrences overlapped a continent using `rnaturalearth` and the `sf` packages (*South, 2017*; *Pebesma, 2018*). For each taxon, I estimated its range by finding the minimum  $\alpha$ -shape containing these occurrences. The  $\alpha$  parameters were set more permissive for marine species since occurrence data for marine taxa were sparser. Then, I intersected the inferred ranges for terrestrial taxa with continental polygons, so their ranges did not overrun landmasses (and likewise with marine taxa and oceans). I inspected diagnostic plots for each taxa for quality control (all of these plots are available in `paradox_variation` GitHub repository), and in some cases, I manually adjusted the  $\alpha$  parameter or

manually corrected the range based on known range maps (these changes are documented in the code `data/species_ranges.r` and `data/species_range_fixes.r`). The range of *C. elegans* was conservatively approximated as the area of the Western US and Western Europe based on the map in **Frézal and Félix, 2015**. *Drosophila* species ranges are from the *Drosophila* Speciation Patterns website, (**Yukilevich, 2012; Yukilevich, 2017**). To further validate these range estimates, I have compared these to the qualitative range descriptions **Leffler et al., 2012** (**Figure 1—figure supplement 4**) and compared my  $\alpha$ -shape method to a subset of taxa with range estimates from IUCN Red List (**Chamberlain, 2020; IUCN, 2020; Figure 1—figure supplement 3**). Each census population size is then estimated as the product of range and density.

## Population size validation

I validated the approximate census sizes by comparing the implied biomass of these estimates to estimates of the total carbon biomass on earth by phylum (**Bar-On et al., 2018**). For species  $i$  with wet body mass  $m_i$  and census size  $N_i$ , the implied biomass is  $m_i N_i$ . For all species in a phylum  $S$ , this total sample biomass is  $b_S = \sum_{i \in S} m_i N_i$ . I then compare this wet biomass to the carbon biomasses by phylum by **Bar-On et al., 2018**. Across animal species, the ratio of dry to wet body mass, and carbon body mass dry body mass varies little. In their study, Bar-On et al. assume wet body mass has a 70% water content, and 50% of dry body mass is carbon mass, leading to a wet body mass to carbon mass factor of  $1 - 0.7/0.5 = 0.15$ . I use this factor to convert the total wet biomass to carbon biomass per phylum.

First, I compared the relative carbon biomass in this study to the relative carbon biomass on earth per phylum. This shows that this study's sample over represents chordate biomass (by a factor of  $\sim 3$ ), and under represents in arthropod biomass (by a factor of 0.02) relative to the proportion of carbon biomass of these phyla on earth (see column eight of **Table 1**). Second, to check whether the carbon biomass per phylum in the sample was broadly consistent with the total on earth by phylum ( $B_S$  for phylum  $S$ ), I calculated the expected sample biomass if species were sampled randomly from the total species in a phylum, ( $B_S \times n_S / T_S$ , where  $n_S$  is the total number of species in the sample in phylum  $S$ ,  $T_S$  is the total number of species in phylum  $S$  on earth). The fraction of total species on earth included in the sample in this study is depicted in **Figure 1—figure supplement 2**.

Next, I look at the ratio of sample biomass per phylum,  $b_S$ , to this expected biomass per phylum (**Table 1**). The consistency is quite close for this rough approach and the non-random sample of taxa included in this study. The carbon biomass estimates for chordates implied by the census size estimates are  $\sim 24$ -fold higher than expected, but is well within reasonable expectations given that the chordate sample includes many larger-bodied domesticated species (and is a biased sample in other ways). Similarly, the implied arthropod carbon biomass is quite close to what one would expect. Overall, these values indicate that the census size estimates here do not lead to implied biomasses

**Table 1.** How the total carbon biomass estimates by phylum from **Bar-On et al., 2018** compare to the implied biomass estimates from this study.

All biomass estimates are carbon biomass, and the proportions are of total biomass with respect to the study. The proportion of biomass in this study compared to the Bar-On et al. estimates **Bar-On et al., 2018** indicates chordates are overrepresented and arthropods are underrepresented in the present study; the factor that each phylum is overrepresented is given in the eighth column. Total species by phylum estimates are from **Reaka-Kudla et al., 1996; Nicol, 1969; Zhang, 2013; Chapman, 2009**. The ratio column is the ratio of total biomass implied by the  $N_c$  estimates of each species in a phylum to the actual biomass of that phylum.

| phylum     | total species (T)  | Bar-On et al. |               | Present study         |               | num. species (n) | factor overrepresented | prop. total species (f=n/T) | factor (b/fB) |
|------------|--------------------|---------------|---------------|-----------------------|---------------|------------------|------------------------|-----------------------------|---------------|
|            |                    | biomass (B)   | prop. biomass | biomass (b)           | prop. biomass |                  |                        |                             |               |
| Arthropoda | $1.26 \times 10^6$ | 1.20          | 0.4635        | $2.80 \times 10^{-4}$ | 0.0102        | 68               | 0.02                   | $5.41 \times 10^{-5}$       | 4.31          |
| Chordata   | $5.41 \times 10^4$ | 0.87          | 0.3357        | $2.67 \times 10^{-2}$ | 0.9715        | 68               | 2.89                   | $1.26 \times 10^{-3}$       | 24.40         |
| Annelida   | $1.70 \times 10^4$ | 0.20          | 0.0772        | $1.23 \times 10^{-5}$ | 0.0004        | 3                | 0.01                   | $1.76 \times 10^{-4}$       | 0.35          |
| Mollusca   | $9.54 \times 10^4$ | 0.20          | 0.0772        | $4.56 \times 10^{-4}$ | 0.0166        | 13               | 0.21                   | $1.36 \times 10^{-4}$       | 16.70         |
| Cnidaria   | $1.60 \times 10^4$ | 0.10          | 0.0386        | $3.07 \times 10^{-5}$ | 0.0011        | 2                | 0.03                   | $1.25 \times 10^{-4}$       | 2.45          |
| Nematoda   | $2.50 \times 10^4$ | 0.02          | 0.0077        | $4.03 \times 10^{-6}$ | 0.0001        | 1                | 0.02                   | $4.00 \times 10^{-5}$       | 5.03          |

per phylum that are outside the range of plausibility. For other population size consistency checks, see Appendix 3.

## Phylogenetic comparative methods

Of the full dataset of 172 taxa with diversity and population size estimates, a synthetic calibrated phylogeny was created for 166 species that appear in phylogenies in DateLife project (O'Meara et al., 2020; Sanchez-Reyes and O'Meara, 2019). This calibrated synthetic phylogeny was then subset for the analyses based on what species had complete trait data. The diversity-population size relationship assessed by a linear phylogenetic mixed-effects model implemented in Stan (Stan Development Team, 2020), according to the methods described in de Villemereuil and Nakagawa, 2014, (see stan/phylo\_mm\_regression.stan in the GitHub repository). This same Stan model was used to estimate the same relationship between arthropod, chordate, and mollusc subsets of the data, though a reduced model was used for the chordate subset due to identifiability issues leading to poor MCMC convergence (Figure 3—figure supplement 1).

The relationship between recombination map length and the logarithm of population size is non-linear and heteroscedastic, and was fit using a lognormal phylogenetic mixed-effects model on the 130 species with complete data. Since social insects have longer recombination map lengths (Wilfert et al., 2007), social taxa were excluded when fitting this model. All Rhat (Vehtari et al., 2019) values were below 1.01 and the effective number of samples was over 1,000, consistent with good mixing; details about the model are available in the GitHub repository (phylo\_mm\_lognormal.stan). Continuous trait maps (Figure 3A, Figure 3—figure supplement 3, and Figure 3—figure supplement 2) were created using phytools (Revell, 2012). Node-height tests were implemented based on the methods in Geiger (Pennell et al., 2014; Harmon et al., 2008), and use robust regression to fit a linear relationship between phylogenetic independent contrasts and branching times.

## Predicted reductions in diversity

The predicted reductions in diversity due to linked selection are approximated using selection and deleterious mutation parameters from *Drosophila melanogaster*, and the recombination map length estimates from Stapley et al., 2017 and Corbett-Detig et al., 2015. The mathematical details of the simplified sweep model are explained in the Appendix Section A. I use estimates of the number of substitutions,  $m$ , in genic regions between *D. melanogaster* and *D. simulans* from Hu et al., 2013. Following Elyashiv et al., 2016, only substitutions in UTRs and exons are included, since they found no evidence of sweeps in introns. Then, I average over annotation classes to estimate the mean proportion of substitutions that are beneficial,  $\alpha_{Dmel} = 0.42$ , which are consistent with the estimates of Elyashiv et al. and estimates from MacDonald-Kreitman test approaches (see Eyre-Walker, 2006, Table 1). Then, I use divergence time estimates between *D. melanogaster* and *D. simulans* of  $4.2 \times 10^6$  and estimate of ten generations per year (Obbard et al., 2012), calculating there are  $\gamma_{Dmel} = \alpha m / 2T = 2.26 \times 10^{-3}$  substitutions per generation. Given the length of the *Drosophila* autosomes,  $G$ , this implies that the rate of beneficial substitutions per basepair, per generation is  $\nu_{BP,Dmel} = \gamma_{Dmel} / G = 2.34 \times 10^{-11}$ . Finally, I estimate  $J_{Dmel} \approx 4.5 \times 10^{-4}$  from the estimate of genome-wide average rate of sweeps from Elyashiv et al. (Supplementary Table S6) and assuming *Drosophila*  $N_e = 10^6$ . These *Drosophila melanogaster* hitchhiking parameter estimates are close to other previously-published estimates (Figure 4—figure supplement 5). Finally, I use  $U_{Dmel} = 1.6$ , from Elyashiv et al., 2016. With these parameter estimates from *D. melanogaster*, the recombination map lengths across species, and Equation (1), I estimate  $\pi_{BGS+HH}$  (assuming  $N_c = N_e$ ) across all species. This leads to a range of predicted diversity ranges across species corresponding to  $\mu = 10^{-9}$ – $10^{-8}$ ; to visualize these, I take a convex hull of all diversity ranges and smooth this with R's smooth.spline function.

## Acknowledgements

I would like to thank Andy Kern and Peter Ralph for helpful discussions and supporting me during this work, and Graham Coop for inspiration and helpful feedback during socially distanced nature walks at Yolo Basin. I thank Jessica Stapley for kindly providing the recombination map length data, and Yaniv Brandvain, Amy Collins, Doc Edge, Tyler Kent, Chuck Langley, Matt Osmond, Sally Otto,

Molly Przeworski, Jeff Ross-Ibarra, Aaron Stern, Anastasia Teterina, Michael Turelli, Margot Wood, and my Kern-Ralph labmates for helpful discussions. Sarah Friedman, Katherine Corn, and Josef Uyeda provided very useful advice about phylogenetic comparative methods; yet I take full responsibility for any shortcomings of my analysis. Finally, I am indebted to Guy Sella, Matt Pennell, and two other anonymous reviewers for helpful feedback. I would like to also thank UO librarian Dean Walton for helping me track down some rather difficult to find older papers. This work was supported by an NIH Grant (1R01GM117241) awarded to Andrew Kern.

## Additional information

### Funding

| Funder                        | Grant reference number | Author        |
|-------------------------------|------------------------|---------------|
| National Institutes of Health | 1R01GM117241           | Vince Buffalo |

The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

### Author contributions

Vince Buffalo, Conceptualization, Resources, Data curation, Software, Formal analysis, Validation, Investigation, Visualization, Methodology, Writing - original draft, Writing - review and editing

### Author ORCIDs

Vince Buffalo  <https://orcid.org/0000-0003-4510-1609>

### Decision letter and Author response

Decision letter <https://doi.org/10.7554/eLife.67509.sa1>

Author response <https://doi.org/10.7554/eLife.67509.sa2>

## Additional files

### Supplementary files

- Transparent reporting form

### Data availability

All primary datasets collated by this study, including new census size and range estimates, are available on Github at [http://github.com/vsbuffalo/paradox\\_variation](http://github.com/vsbuffalo/paradox_variation) (copy archived at <https://archive.softwareheritage.org/swh:1:rev:8fa6b5834f6536319b1e5cd9722ca02d317183df>). An archived version of this repository is also available at Zenodo.

The following dataset was generated:

| Author(s) | Year | Dataset title   | Dataset URL   | Database and Identifier        |
|-----------|------|---|---|--------------------------------|
| Vince B   | 2021 | vsbuffalo/paradox_variation: biorxiv v.1 with minor corrections | <a href="https://doi.org/10.5281/zenodo.4542480">https://doi.org/10.5281/zenodo.4542480</a> | Zenodo, 10.5281/zenodo.4542480 |

The following previously published datasets were used:

| Author(s)  | Year | Dataset title   | Dataset URL   | Database and Identifier                    |
|--|------|---|---|--|
| Stapley J, Feulner PGD, Johnston SE, Santure AW, Smadja CM | 2017 | Supplementary material from "Variation in recombination frequency and distribution across eukaryotes: patterns and processes" | <a href="https://doi.org/10.6084/m9.figshare.c.3904942.v3">https://doi.org/10.6084/m9.figshare.c.3904942.v3</a> | figshare, 10.6084/m9.figshare.c.3904942.v3 |



## References

- Aguade M**, Miyashita N, Langley CH. 1989. Reduced variation in the yellow-achaete-scute region in natural populations of *Drosophila melanogaster*. *Genetics* **122**:607–615. DOI: <https://doi.org/10.1093/genetics/122.3.607>, PMID: 17246506
- Andolfatto P**. 2007. Hitchhiking effects of recurrent beneficial amino acid substitutions in the *Drosophila melanogaster* genome. *Genome Research* **17**:1755–1762. DOI: <https://doi.org/10.1101/gr.6691007>, PMID: 17989248
- Bar-On YM**, Phillips R, Milo R. 2018. The biomass distribution on earth. *PNAS* **115**:6506–6511. DOI: <https://doi.org/10.1073/pnas.1711842115>, PMID: 29784790
- Barry P**, Broquet T, Gagnaire P-A. 2020. Life tables shape genetic diversity in marine fishes. *bioRxiv*. DOI: <https://doi.org/10.1101/2020.12.18.423459>
- Barton NH**. 1995. Linkage and the limits to natural selection. *Genetics* **140**:821–841. DOI: <https://doi.org/10.1093/genetics/140.2.821>, PMID: 7498757
- Barton NH**. 2000. Genetic hitchhiking. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* **355**:1553–1562. DOI: <https://doi.org/10.1098/rstb.2000.0716>, PMID: 11127900
- Barton N**. 2010. Understanding adaptation in large populations. *PLOS Genetics* **6**:e1000987. DOI: <https://doi.org/10.1371/journal.pgen.1000987>, PMID: 20585547
- Batley CJ**, Ralph PL, Kern AD. 2020. Space is the place: effects of continuous spatial structure on analysis of population genetic data. *Genetics* **215**:193–214. DOI: <https://doi.org/10.1534/genetics.120.303143>, PMID: 32209569
- Begun DJ**, Aquadro CF. 1992. Levels of naturally occurring DNA polymorphism correlate with recombination rates in *Drosophila melanogaster*. *Nature* **356**:519–520. DOI: <https://doi.org/10.1038/356519a0>
- Boyko AR**, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR, White TJ, Nielsen R, Clark AG, Bustamante CD. 2008. Assessing the evolutionary impact of amino acid mutations in the human genome. *PLOS Genetics* **4**:e1000083. DOI: <https://doi.org/10.1371/journal.pgen.1000083>, PMID: 18516229
- Buffalo V**. 2021. Code and Data for Why do species get a thin slice of  $\pi$ ? *Software Heritage*. swh:1:rev:8fa6b5834f6536319b1e5cd9722ca02d317183df <https://archive.softwareheritage.org/swh:1:rev:8fa6b5834f6536319b1e5cd9722ca02d317183df>.
- Buffalo V**, Coop G. 2019. The linked selection signature of rapid adaptation in temporal genomic data. *Genetics* **213**:1007–1045. DOI: <https://doi.org/10.1534/genetics.119.302581>, PMID: 31558582
- Buffalo V**, Coop G. 2020. Estimating the genome-wide contribution of selection to temporal allele frequency change. *PNAS* **117**:20672–20680. DOI: <https://doi.org/10.1073/pnas.1919039117>, PMID: 32817464
- Burt A**, Bell G. 1987. Mammalian chiasma frequencies as a test of two theories of recombination. *Nature* **326**:803–805. DOI: <https://doi.org/10.1038/326803a0>, PMID: 3574451
- Caballero A**. 1994. Developments in the prediction of effective population size. *Heredity* **73 (Pt 6)**:657–679. DOI: <https://doi.org/10.1038/hdy.1994.174>, PMID: 7814264
- Caballero A**. 2020. *Quantitative Genetics*. Cambridge University Press.
- Cai JJ**, Macpherson JM, Sella G, Petrov DA. 2009. Pervasive hitchhiking at coding and regulatory sites in humans. *PLOS Genetics* **5**:e1000336. DOI: <https://doi.org/10.1371/journal.pgen.1000336>, PMID: 19148272
- Carpenter B**, Gelman A, Hoffman M, Lee D, Goodrich B, Betancourt M, Brubaker M, Guo J, Li P, Riddell A. 2017. Stan: a probabilistic programming language. *Journal of Statistical Software* **76**:1–32. DOI: <https://doi.org/10.18637/jss.v076.i01>
- Chamberlain S**, Ram K, Barve V, Mcglinn D. 2014. *rgbif: interface to the global biodiversity information facility API*. R package version 0. 7, 7.
- Chamberlain S**. 2020. *rredlist: 'IUCN' red list client*.
- Chamberlain S**, Boettiger C. 2017. R Python, and ruby clients for GBIF species occurrence data. *PeerJ Preprints*. DOI: <https://doi.org/10.7287/peerj.preprints.3304v1>
- Chapman AD**. 2009. *Numbers of Living Species in Australia and the World*. Department of the Environment, Water, Heritage and the Arts Canberra.
- Charlesworth B**. 1987. *Sexual Selection: Testing the Alternatives*. Wiley.
- Charlesworth B**, Morgan MT, Charlesworth D. 1993. The effect of deleterious mutations on neutral molecular variation. *Genetics* **134**:1289–1303. DOI: <https://doi.org/10.1093/genetics/134.4.1289>, PMID: 8375663
- Charlesworth B**. 1996. Background selection and patterns of genetic diversity in *Drosophila Melanogaster*. *Genetical Research* **68**:131–149. DOI: <https://doi.org/10.1017/S0016672300034029>, PMID: 8940902
- Charlesworth B**. 2015. Causes of natural variation in fitness: evidence from studies of *Drosophila* populations. *PNAS* **112**:1662–1669. DOI: <https://doi.org/10.1073/pnas.1423275112>, PMID: 25572964
- Charlesworth B**, Hughes KA. 2000. The maintenance of genetic variation in Life-History traits. In: Singh R. S, Krimbas C (Eds). *Evolutionary Genetics: From Molecules to Morphology*. **1** Cambridge: University Press. p. 369–392.
- Chen J**, Glémin S, Lascoux M. 2017. Genetic diversity and the efficacy of purifying selection across plant and animal species. *Molecular Biology and Evolution* **34**:1417–1428. DOI: <https://doi.org/10.1093/molbev/msx088>, PMID: 28333215
- Coop G**. 2016. Does linked selection explain the narrow range of genetic diversity across species? *bioRxiv*. DOI: <https://doi.org/10.1101/042598>



- Coop G**, Ralph P. 2012. Patterns of neutral diversity under general models of selective sweeps. *Genetics* **192**: 205–224. DOI: <https://doi.org/10.1534/genetics.112.141861>, PMID: 22714413
- Corbett-Detig RB**, Hartl DL, Sackton TB. 2015. Natural selection constrains neutral diversity across a wide range of species. *PLoS Biology* **13**:e1002112. DOI: <https://doi.org/10.1371/journal.pbio.1002112>, PMID: 25859758
- Crow JF**, Kimura M. 1970. *An Introduction to Population Genetics Theory*. New York, Evanston and London: Harper and Row Publishers.
- Cutter AD**, Payseur BA. 2003. Selection at linked sites in the partial selfer *Caenorhabditis elegans*. *Molecular Biology and Evolution* **20**:665–673. DOI: <https://doi.org/10.1093/molbev/msg072>, PMID: 12679551
- Damuth J**. 1981. Population density and body size in mammals. *Nature* **290**:699–700. DOI: <https://doi.org/10.1038/290699a0>
- Damuth J**. 1987. Interspecific allometry of population density in mammals and other animals: the independence of body mass and population energy-use. *Biological Journal of the Linnean Society* **31**:193–246. DOI: <https://doi.org/10.1111/j.1095-8312.1987.tb01990.x>
- de Villemeureuil P**, Nakagawa S. 2014. General quantitative genetic methods for comparative biology. In: Garamszegi L. Z (Ed). *Modern Phylogenetic Comparative Methods and Their Application in Evolutionary Biology: Concepts and Practice*. Berlin: Berlin, Heidelberg. p. 287–303. DOI: <https://doi.org/10.1007/978-3-662-43550-2>
- Eldon B**, Wakeley J. 2006. Coalescent processes when the distribution of offspring number among individuals is highly skewed. *Genetics* **172**:2621–2633. DOI: <https://doi.org/10.1534/genetics.105.052175>, PMID: 16452141
- Elyashiv E**, Sattath S, Hu TT, Strutsovsky A, McVicker G, Andolfatto P, Coop G, Sella G. 2016. A genomic map of the effects of linked selection in *Drosophila*. *PLoS Genetics* **12**:e1006130. DOI: <https://doi.org/10.1371/journal.pgen.1006130>, PMID: 27536991
- Eyre-Walker A**. 2006. The genomic rate of adaptive evolution. *Trends in Ecology & Evolution* **21**:569–575. DOI: <https://doi.org/10.1016/j.tree.2006.06.015>, PMID: 16820244
- FAOSTAT statistics database**. 2021. UN food and agriculture organisation Rome. <http://www.fao.org/faostat/en/> [Accessed May 17, 2021].
- Felsenstein J**. 1985. Phylogenies and the comparative method. *The American Naturalist* **125**:1–15. DOI: <https://doi.org/10.1086/284325>
- Fisher RA**, Ford EB. 1947. The spread of a gene in natural conditions in a colony of the moth *panaxia dominula* L. *Heredity* **1**:143–174. DOI: <https://doi.org/10.1038/hdy.1947.11>
- Frankham R**. 1995. Effective population size/adult population size ratios in wildlife: a review. *Genetical Research* **66**:95–107. DOI: <https://doi.org/10.1017/S0016672300034455>
- Frankham R**. 1996. Relationship of genetic variation to population size in wildlife. *Conservation Biology* **10**:1500–1508. DOI: <https://doi.org/10.1046/j.1523-1739.1996.10061500.x>
- Freckleton RP**, Harvey PH, Pagel M. 2002. Phylogenetic analysis and comparative data: a test and review of evidence. *The American Naturalist* **160**:712–726. DOI: <https://doi.org/10.1086/343873>, PMID: 18707460
- Freckleton RP**, Harvey PH. 2006. Detecting non-Brownian trait evolution in adaptive radiations. *PLoS Biology* **4**: e373. DOI: <https://doi.org/10.1371/journal.pbio.0040373>, PMID: 17090217
- Frézal L**, Félix MA. 2015. *C. elegans* outside the petri dish. *eLife* **4**:e05849. DOI: <https://doi.org/10.7554/eLife.05849>, PMID: 25822066
- Galtier N**. 2016. Adaptive protein evolution in animals and the effective population size hypothesis. *PLoS Genetics* **12**:e1005774. DOI: <https://doi.org/10.1371/journal.pgen.1005774>, PMID: 26752180
- Galtier N**, Rousselle M. 2020. How much does  $\pi$  vary among species? *Genetics* **216**:303622. DOI: <https://doi.org/10.1534/genetics.120.303622>
- Gaston K**, Blackburn T. 2008. *Pattern and Process in Macroecology*. John Wiley & Sons.
- Gillespie JH**. 1991. *The Causes of Molecular Evolution*. Oxford: Oxford University Press Google Scholar.
- Gillespie JH**. 2000. Genetic drift in an infinite population. The pseudohitchhiking model. *Genetics* **155**:909–919. DOI: <https://doi.org/10.1093/genetics/155.2.909>, PMID: 10835409
- Gillespie JH**. 2001. Is the population size of a species relevant to its evolution? *Evolution* **55**:2161–2169. DOI: <https://doi.org/10.1111/j.0014-3820.2001.tb00732.x>
- Global Biodiversity Information Facility**. 2020. (27 August 2020) GBIF Occurrence Download. *GBIF.org*. DOI: <https://doi.org/10.15468/dl.nb3s74>
- Hadfield JD**, Nakagawa S. 2010. General quantitative genetic methods for comparative biology: phylogenies, taxonomies and multi-trait models for continuous and categorical characters. *Journal of Evolutionary Biology* **23**:494–508. DOI: <https://doi.org/10.1111/j.1420-9101.2009.01915.x>, PMID: 20070460
- Hallatschek O**. 2018. Selection-Like biases emerge in population models with recurrent jackpot events. *Genetics* **210**:1053–1073. DOI: <https://doi.org/10.1534/genetics.118.301516>, PMID: 30171032
- Hartman LJ**, Weir JT, Brock CD, Glor RE, Challenger W. 2008. GEIGER: investigating evolutionary radiations. *Bioinformatics* **24**:129–131. DOI: <https://doi.org/10.1093/bioinformatics/btm538>, PMID: 18006550
- Hauser L**, Carvalho GR. 2008. Paradigm shifts in marine fisheries genetics: ugly hypotheses slain by beautiful facts. *Fish and Fisheries* **9**:333–362. DOI: <https://doi.org/10.1111/j.1467-2979.2008.00299.x>
- Hedgecock D**. 1994. Does variance in reproductive success limit effective population sizes of marine organisms. *Genetics and Evolution of Aquatic Organisms* **122**:122–134.
- Hedgecock D**, Pudovkin AI. 2011. Sweepstakes reproductive success in highly fecund marine fish and shellfish: a review and commentary. *Bulletin of Marine Science* **87**:971–1002. DOI: <https://doi.org/10.5343/bms.2010.1051>

- Hellmann I, Mang Y, Gu Z, Li P, de la Vega FM, Clark AG, Nielsen R. 2008. Population genetic analysis of shotgun assemblies of genomic sequences from multiple individuals. *Genome Research* **18**:1020–1029. DOI: <https://doi.org/10.1101/gr.074187.107>, PMID: 18411405
- Hermisson J, Pennings PS. 2005. Soft sweeps: molecular population genetics of adaptation from standing genetic variation. *Genetics* **169**:2335–2352. DOI: <https://doi.org/10.1534/genetics.104.036947>, PMID: 15716498
- Hernandez RD, Kelley JL, Elyashiv E, Melton SC, Auton A, McVean G, Sella G, Przeworski M, 1000 Genomes Project. 2011. Classic selective sweeps were rare in recent human evolution. *Science* **331**:920–924. DOI: <https://doi.org/10.1126/science.1198878>, PMID: 21330547
- Hu TT, Eisen MB, Thornton KR, Andolfatto P. 2013. A second-generation assembly of the *Drosophila simulans* genome provides new insights into patterns of lineage-specific divergence. *Genome Research* **23**:89–98. DOI: <https://doi.org/10.1101/gr.141689.112>, PMID: 22936249
- Hudson RR, Kaplan NL. 1994. Gene trees with background selection. In: Golding B (Ed). *Non-Neutral Evolution: Theories and Molecular Data*. Boston: Springer. p. 140–153. DOI: <https://doi.org/10.1007/978-1-4615-2383-3>
- Hudson RR, Kaplan NL. 1995. Deleterious background selection with recombination. *Genetics* **141**:1605–1617. DOI: <https://doi.org/10.1093/genetics/141.4.1605>, PMID: 8601498
- IUCN. 2020. The IUCN red list of threatened species. <https://www.iucnredlist.org> [Accessed October 31, 2020].
- Jensen JD, Thornton KR, Andolfatto P. 2008. An approximate bayesian estimator suggests strong, recurrent selective sweeps in *Drosophila*. *PLOS Genetics* **4**:e1000198. DOI: <https://doi.org/10.1371/journal.pgen.1000198>, PMID: 18802463
- Kaplan NL, Hudson RR, Langley CH. 1989. The "hitchhiking effect" revisited. *Genetics* **123**:887–899. DOI: <https://doi.org/10.1093/genetics/123.4.887>, PMID: 2612899
- Karasov T, Messer PW, Petrov DA. 2010. Evidence that adaptation in *Drosophila* is not limited by mutation at single sites. *PLOS Genetics* **6**:e1000924. DOI: <https://doi.org/10.1371/journal.pgen.1000924>, PMID: 20585551
- Kim Y, Stephan W. 2000. Joint effects of genetic hitchhiking and background selection on neutral variation. *Genetics* **155**:1415–1427. DOI: <https://doi.org/10.1093/genetics/155.3.1415>, PMID: 10880499
- Kimura M. 1984. *The Neutral Theory of Molecular Evolution*. Cambridge University Press.
- Kimura M, Crow JF. 1964. The number of alleles that can be maintained in a finite population. *Genetics* **49**:725–738. DOI: <https://doi.org/10.1093/genetics/49.4.725>, PMID: 14156929
- Kondrashov FA, Kondrashov AS. 2010. Measurements of spontaneous rates of mutations in the recent past and the near future. *Philosophical Transactions of the Royal Society B: Biological Sciences* **365**:1169–1176. DOI: <https://doi.org/10.1098/rstb.2009.0286>, PMID: 20308091
- Leffler EM, Bullaughey K, Matute DR, Meyer WK, Ségurel L, Venkat A, Andolfatto P, Przeworski M. 2012. Revisiting an old riddle: what determines genetic diversity levels within species? *PLOS Biology* **10**:e1001388. DOI: <https://doi.org/10.1371/journal.pbio.1001388>, PMID: 22984349
- Leroy T, Rousselle M, Tilak MK, Caizergues AE, Scornavacca C, Recuerda M, Fuchs J, Illera JC, De Swardt DH, Blanco G, Thébaud C, Milá B, Nabholz B. 2021. Island songbirds as windows into evolution in small populations. *Current Biology* **31**:1303–1310. DOI: <https://doi.org/10.1016/j.cub.2020.12.040>, PMID: 33476557
- Lewontin RC. 1974. *The Genetic Basis of Evolutionary Change*. New York: Columbia University Press. DOI: <https://doi.org/10.1111/j.1558-5646.1975.tb00851.x>
- Lewontin RC, Singh RS, Uyenoyama MK. 2004. Building a science of population biology. In: Singh RS, Uyenoyama MK (Eds). *The Evolution of Population Biology*. Cambridge University Press. p. 7–20. DOI: <https://doi.org/10.1017/CBO9780511542619.004>
- Li H, Stephan W. 2006. Inferring the demographic history and rate of adaptive substitution in *Drosophila*. *PLOS Genetics* **2**:e166. DOI: <https://doi.org/10.1371/journal.pgen.0020166>, PMID: 17040129
- Lynch M. 1991. Methods for the analysis of comparative data in evolutionary biology. *Evolution* **45**:1065–1080. DOI: <https://doi.org/10.1111/j.1558-5646.1991.tb04375.x>
- Lynch M. 2010. Evolution of the mutation rate. *Trends in Genetics* **26**:345–352. DOI: <https://doi.org/10.1016/j.tig.2010.05.003>, PMID: 20594608
- Lynch M. 2011. Statistical inference on the mechanisms of genome evolution. *PLOS Genetics* **7**:e1001389. DOI: <https://doi.org/10.1371/journal.pgen.1001389>, PMID: 21695228
- Lynch M, Conery JS. 2003. The origins of genome complexity. *Science* **302**:1401–1404. DOI: <https://doi.org/10.1126/science.1089370>, PMID: 14631042
- Mace GM, Lande R. 1991. Assessing extinction threats: toward a reevaluation of IUCN threatened species categories. *Conservation Biology* **5**:148–157. DOI: <https://doi.org/10.1111/j.1523-1739.1991.tb00119.x>
- Macpherson JM, Sella G, Davis JC, Petrov DA. 2007. Genomewide spatial correspondence between nonsynonymous divergence and neutral polymorphism reveals extensive adaptation in *Drosophila*. *Genetics* **177**:2083–2099. DOI: <https://doi.org/10.1534/genetics.107.080226>, PMID: 18073425
- Malécot G. 1948. *Les mathématiques de l'hérédité*. Paris: Masson.
- Maruyama T, Kimura M. 1980. Genetic variability and effective population size when local extinction and recolonization of subpopulations are frequent. *PNAS* **77**:6710–6714. DOI: <https://doi.org/10.1073/pnas.77.11.6710>, PMID: 16592920
- McVicker G, Gordon D, Davis C, Green P. 2009. Widespread genomic signatures of natural selection in hominid evolution. *PLOS Genetics* **5**:e1000471. DOI: <https://doi.org/10.1371/journal.pgen.1000471>, PMID: 19424416
- Mora C, Tittensor DP, Adl S, Simpson AG, Worm B. 2011. How many species are there on earth and in the ocean? *PLOS Biology* **9**:e1001127. DOI: <https://doi.org/10.1371/journal.pbio.1001127>, PMID: 21886479

- Mukai T. 1985. Experimental verification of the neutral theory. In: Ohta T, Aoki K (Eds). *Population Genetics and Molecular Evolution*. Berlin: Springer-Verlag. p. 125–145.
- Mukai T. 1988. Genotype-environment interaction in relation to the maintenance of genetic variability in populations of *Drosophila melanogaster*. Proceedings of the Second International Conference on Quantitative Genetics.
- Nei M, Graur D. 1984. Extent of protein polymorphism and the neutral mutation theory. *Evolutionary Biology* **17**: 73–118.
- Nevo E. 1978. Genetic variation in natural populations: patterns and theory. *Theoretical Population Biology* **13**: 121–177. DOI: [https://doi.org/10.1016/0040-5809\(78\)90039-4](https://doi.org/10.1016/0040-5809(78)90039-4), PMID: 347627
- Nevo E, Beiles A, Ben-Shlomo R. 1984. The evolutionary significance of genetic diversity: Ecological, demographic and life history correlates. In: Mani GS (Ed). *Evolutionary Dynamics of Genetic Diversity*. Heidelberg: Springer Berlin. p. 13–213. DOI: [https://doi.org/10.1007/978-3-642-51588-0\\_2](https://doi.org/10.1007/978-3-642-51588-0_2)
- Nicol D. 1969. The number of living species of molluscs. *Systematic Zoology* **18**:251–254. DOI: <https://doi.org/10.2307/2412618>
- Nicolaisen LE, Desai MM. 2012. Distortions in genealogies due to purifying selection. *Molecular Biology and Evolution* **29**:3589–3600. DOI: <https://doi.org/10.1093/molbev/mss170>, PMID: 22729750
- Nordborg M, Charlesworth B, Charlesworth D. 1996. The effect of recombination on background selection. *Genetical Research* **67**:159–174. DOI: <https://doi.org/10.1017/S0016672300033619>, PMID: 8801188
- Novak S, Barton NH. 2017. When does Frequency-Independent selection maintain genetic variation? *Genetics* **207**:653–668. DOI: <https://doi.org/10.1534/genetics.117.300129>, PMID: 28798062
- Nunney L. 1993. The influence of mating system and overlapping generations on effective population size. *Evolution* **47**:1329–1341. DOI: <https://doi.org/10.1111/j.1558-5646.1993.tb02158.x>
- Nunney L. 1996. The influence of variation in female fecundity on effective population size. *Biological Journal of the Linnean Society* **59**:411–425. DOI: <https://doi.org/10.1111/j.1095-8312.1996.tb01474.x>
- Nydam ML, Harrison RG. 2010. Polymorphism and divergence within the ascidian genus *Ciona*. *Molecular Phylogenetics and Evolution* **56**:718–726. DOI: <https://doi.org/10.1016/j.ympev.2010.03.042>, PMID: 20403444
- Obbard DJ, Maclennan J, Kim KW, Rambaut A, O’Grady PM, Jiggins FM. 2012. Estimating divergence dates and substitution rates in the *Drosophila* phylogeny. *Molecular Biology and Evolution* **29**:3459–3473. DOI: <https://doi.org/10.1093/molbev/mss150>, PMID: 22683811
- Ohta T. 1992. The nearly neutral theory of molecular evolution. *Annual Review of Ecology and Systematics* **23**: 263–286. DOI: <https://doi.org/10.1146/annurev.es.23.110192.001403>
- Ohta T, Kimura M. 1973. A model of mutation appropriate to estimate the number of electrophoretically detectable alleles in a finite population. *Genetical Research* **22**:201–204. DOI: <https://doi.org/10.1017/S0016672300012994>, PMID: 4777279
- O’Meara B, Sanchez-Reyes LL, Eastman J, Heath T, ril Wright A, Schliep K, Chamberlain S, Midford P, Harmon L, Brown J, Pennell M, Alfaro M. 2020. *Datelife: Go from a List of Taxa or a Tree to a Chronogram using Open Scientific Data*. 0.3.2. <https://github.com/phyloastic/datelife>
- Palkopoulou E, Mallick S, Skoglund P, Enk J, Rohland N, Li H, Omrak A, Vartanyan S, Poinar H, Götherström A, Reich D, Dalén L. 2015. Complete genomes reveal signatures of demographic and genetic declines in the woolly mammoth. *Current Biology* **25**:1395–1400. DOI: <https://doi.org/10.1016/j.cub.2015.04.007>, PMID: 25913407
- Palstra FP, Fraser DJ. 2012. Effective/census population size ratio estimation: a compendium and appraisal. *Ecology and Evolution* **2**:2357–2365. DOI: <https://doi.org/10.1002/ece3.329>, PMID: 23139893
- Palstra FP, Ruzzante DE. 2008. Genetic estimates of contemporary effective population size: what can they tell us about the importance of genetic stochasticity for wild population persistence? *Molecular Ecology* **17**:3428–3447. DOI: <https://doi.org/10.1111/j.1365-294X.2008.03842.x>, PMID: 19160474
- Pebesma E. 2018. Simple features for R: standardized support for spatial vector data. *The R Journal* **10**:439. DOI: <https://doi.org/10.32614/RJ-2018-009>
- Pennell MW, Eastman JM, Slater GJ, Brown JW, Uyeda JC, FitzJohn RG, Alfaro ME, Harmon LJ. 2014. Geiger v2.0: an expanded suite of methods for fitting macroevolutionary models to phylogenetic trees. *Bioinformatics* **30**: 2216–2218. DOI: <https://doi.org/10.1093/bioinformatics/btu181>, PMID: 24728855
- Pennings PS, Hermisson J. 2006. Soft sweeps II—molecular population genetics of adaptation from recurrent mutation or migration. *Molecular Biology and Evolution* **23**:1076–1084. DOI: <https://doi.org/10.1093/molbev/msj117>, PMID: 16520336
- Pershing AJ, Christensen LB, Record NR, Sherwood GD, Stetson PB. 2010. The impact of whaling on the ocean carbon cycle: why bigger was better. *PLOS ONE* **5**:e12444. DOI: <https://doi.org/10.1371/journal.pone.0012444>, PMID: 20865156
- Powell JR. 1975. Protein variation in natural populations of animals. In: Theodosius D, Hecht M, William C. S (Eds). *Evolutionary Biology*. **8** New York: Plenum Press. p. 79–199.
- Ralph P, Coop G. 2010. Parallel adaptation: one or many waves of advance of an advantageous allele? *Genetics* **186**:647–668. DOI: <https://doi.org/10.1534/genetics.110.119594>, PMID: 20660645
- Ralph PL, Coop G. 2015. The role of standing variation in geographic convergent adaptation. *The American Naturalist* **186** Suppl 1:S5–S23. DOI: <https://doi.org/10.1086/682948>, PMID: 26656217
- Reaka-Kudla ML, Wilson DE, Wilson EO. 1996. *Biodiversity II: Understanding and Protecting Our Biological Resources*. Joseph Henry Press.
- Revell LJ, Harmon LJ, Collar DC. 2008. Phylogenetic signal, evolutionary process, and rate. *Systematic Biology* **57**:591–601. DOI: <https://doi.org/10.1080/10635150802302427>, PMID: 18709597

- Revell LJ. 2010. Phylogenetic signal and linear regression on species data: phylogenetic regression. *Methods in Ecology and Evolution* **1**:319–329. DOI: <https://doi.org/10.1111/j.2041-210X.2010.00044.x>
- Revell LJ. 2012. Phytools: an R package for phylogenetic comparative biology (and other things). *Methods in Ecology and Evolution* **3**:217–223. DOI: <https://doi.org/10.1111/j.2041-210X.2011.00169.x>
- Robertson A. 1961. Inbreeding in artificial selection programmes. *Genetical Research* **2**:189–194. DOI: <https://doi.org/10.1017/S001667230000690>
- Robinson TP, Wint GR, Conchedda G, Van Boeckel TP, Ercoli V, Palamara E, Cinardi G, D’Aielli L, Hay SI, Gilbert M. 2014. Mapping the global distribution of livestock. *PLOS ONE* **9**:e96084. DOI: <https://doi.org/10.1371/journal.pone.0096084>, PMID: 24875496
- Romiguier J, Gayral P, Ballenghien M, Bernard A, Cahais V, Chenuil A, Chiari Y, Dernas R, Duret L, Favre N, Loire E, Lourenco JM, Nabholz B, Roux C, Tsagkogeorga G, Weber AA, Weinert LA, Belkhir K, Bierne N, Glémin S, et al. 2014. Comparative population genomics in animals uncovers the determinants of genetic diversity. *Nature* **515**:261–263. DOI: <https://doi.org/10.1038/nature13685>, PMID: 25141177
- Roux C, Fraïsse C, Romiguier J, Anciaux Y, Galtier N, Bierne N. 2016. Shedding light on the grey zone of speciation along a continuum of genomic divergence. *PLOS Biology* **14**:e2000234. DOI: <https://doi.org/10.1371/journal.pbio.2000234>, PMID: 28027292
- Roze D. 2021. A simple expression for the strength of selection on recombination generated by interference among mutations. *PNAS* **118**:e2022805118. DOI: <https://doi.org/10.1073/pnas.2022805118>, PMID: 33941695
- Sanchez-Reyes LL, O’Meara B. 2019. Datalife: leveraging databases and analytical tools to reveal the dated tree of life. *bioRxiv*. DOI: <https://doi.org/10.1101/782094>
- Santiago E, Caballero A. 1995. Effective size of populations under selection. *Genetics* **139**:1013–1030. DOI: <https://doi.org/10.1093/genetics/139.2.1013>, PMID: 7713405
- Santiago E, Caballero A. 1998. Effective size and polymorphism of linked neutral loci in populations under directional selection. *Genetics* **149**:2105–2117. DOI: <https://doi.org/10.1093/genetics/149.4.2105>, PMID: 9691062
- Sella G, Petrov DA, Przeworski M, Andolfatto P. 2009. Pervasive natural selection in the *Drosophila* genome? *PLOS Genetics* **5**:e1000495. DOI: <https://doi.org/10.1371/journal.pgen.1000495>, PMID: 19503600
- Shapiro JA, Huang W, Zhang C, Hubisz MJ, Lu J, Turissini DA, Fang S, Wang HY, Hudson RR, Nielsen R, Chen Z, Wu CI. 2007. Adaptive genic evolution in the *Drosophila* genomes. *PNAS* **104**:2271–2276. DOI: <https://doi.org/10.1073/pnas.0610385104>, PMID: 17284599
- Shirihai H. 2008. *The Complete Guide to Antarctic Wildlife: Birds and Marine Mammals of the Antarctic Continent and the Southern Ocean*. Princeton University Press.
- Slatkin M. 1977. Gene flow and genetic drift in a species subject to frequent local extinctions. *Theoretical Population Biology* **12**:253–262. DOI: [https://doi.org/10.1016/0040-5809\(77\)90045-4](https://doi.org/10.1016/0040-5809(77)90045-4), PMID: 601717
- Small KS, Brudno M, Hill MM, Sidow A. 2007. Extreme genomic variation in a natural population. *PNAS* **104**:5698–5703. DOI: <https://doi.org/10.1073/pnas.0700890104>, PMID: 17372217
- Smith JM, Haigh J. 1974. The hitch-hiking effect of a favourable gene. *Genetical Research* **23**:23–35. DOI: <https://doi.org/10.1017/S0016672300014634>, PMID: 4407212
- Soulé ME. 1976. Allozyme variation, its determinants in space and time. In: Ayala F. J (Ed). *Molecular Evolution*. Sunderland, Massachusetts: Sinauer Associates. p. 60–77.
- South A. 2017. *Rnaturalearth: World Map Data From Natural Earth*. Natural Earth.
- Spielman D, Brook BW, Frankham R. 2004. Most species are not driven to extinction before genetic factors impact them. *PNAS* **101**:15261–15264. DOI: <https://doi.org/10.1073/pnas.0403809101>, PMID: 15477597
- Stan Development Team. 2020. *Stan Modeling Language Users Guide and Reference Manual*. Stan Developer.
- Stapley J, Feulner PGD, Johnston SE, Santure AW, Smadja CM. 2017. Variation in recombination frequency and distribution across eukaryotes: patterns and processes. *Philosophical Transactions of the Royal Society B: Biological Sciences* **372**:20160455. DOI: <https://doi.org/10.1098/rstb.2016.0455>, PMID: 29109219
- Stephan W, Wiehe THE, Lenz MW. 1992. The effect of strongly selected substitutions on neutral polymorphism: analytical results based on diffusion theory. *Theoretical Population Biology* **41**:237–254. DOI: [https://doi.org/10.1016/0040-5809\(92\)90045-U](https://doi.org/10.1016/0040-5809(92)90045-U)
- Stephan W. 1995. An improved method for estimating the rate of fixation of favorable mutations based on DNA polymorphism data. *Molecular Biology and Evolution* **12**:959–962. DOI: <https://doi.org/10.1093/oxfordjournals.molbev.a040274>, PMID: 7476143
- Stephan W, Langley CH. 1998. DNA polymorphism in lycopersicon and crossing-over per physical length. *Genetics* **150**:1585–1593. DOI: <https://doi.org/10.1093/genetics/150.4.1585>, PMID: 9832534
- Tajima F. 1996. The amount of DNA polymorphism maintained in a finite population when the neutral mutation rate varies among sites. *Genetics* **143**:1457–1465. DOI: <https://doi.org/10.1093/genetics/143.3.1457>, PMID: 8807315
- Tsagkogeorga G, Cahais V, Galtier N. 2012. The population genomics of a fast evolver: high levels of diversity, functional constraint, and molecular adaptation in the tunicate *Ciona intestinalis*. *Genome Biology and Evolution* **4**:852–861. DOI: <https://doi.org/10.1093/gbe/evs054>, PMID: 22745226
- Uyeda JC, Zenil-Ferguson R, Pennell MW. 2018. Rethinking phylogenetic comparative methods. *Systematic Biology* **67**:1091–1109. DOI: <https://doi.org/10.1093/sysbio/syy031>, PMID: 29701838
- Vehtari A, Gelman A, Simpson D, Carpenter B, Bürkner P-C. 2019. Rank-normalization, folding, and localization: an improved for assessing convergence of MCMC. *arXiv*. <https://arxiv.org/abs/1903.08008>.
- Wang J. 2005. Estimation of effective population sizes from data on genetic markers. *Philosophical Transactions of the Royal Society B: Biological Sciences* **360**:1395–1409. DOI: <https://doi.org/10.1098/rstb.2005.1682>



- Wang J, Santiago E, Caballero A. 2016. Prediction and estimation of effective population size. *Heredity* **117**:193–206. DOI: <https://doi.org/10.1038/hdy.2016.43>, PMID: 27353047
- Waples RS. 1989. A generalized approach for estimating effective population size from temporal changes in allele frequency. *Genetics* **121**:379–391. DOI: <https://doi.org/10.1093/genetics/121.2.379>, PMID: 2731727
- Waples RS, Luikart G, Faulkner JR, Tallmon DA. 2013. Simple life-history traits explain key effective population size ratios across diverse taxa. *Proceedings of the Royal Society B: Biological Sciences* **280**:20131339. DOI: <https://doi.org/10.1098/rspb.2013.1339>, PMID: 23926150
- Waples RS, Grewe PM, Bravington MW, Hillary R, Feutry P. 2018. Robust estimates of a high  $N_e/N$  ratio in a top marine predator, southern bluefin tuna. *Science Advances* **4**:eaar7759. DOI: <https://doi.org/10.1126/sciadv.aar7759>, PMID: 30035218
- Weissman DB, Barton NH. 2012. Limits to the rate of adaptive substitution in sexual populations. *PLOS Genetics* **8**:e1002740. DOI: <https://doi.org/10.1371/journal.pgen.1002740>, PMID: 22685419
- Westoby M, Leishman MR, Lord JM. 1995. On Misinterpreting the 'Phylogenetic Correction'. *The Journal of Ecology* **83**:531–534. DOI: <https://doi.org/10.2307/2261605>
- Whitney KD, Garland T. 2010. Did genetic drift drive increases in genome complexity? *PLOS Genetics* **6**:e1001080. DOI: <https://doi.org/10.1371/journal.pgen.1001080>, PMID: 20865118
- Wilfert L, Gadau J, Schmid-Hempel P. 2007. Variation in genomic recombination rates among animal taxa and the case of social insects. *Heredity* **98**:189–197. DOI: <https://doi.org/10.1038/sj.hdy.6800950>, PMID: 17389895
- Wright S. 1931. Evolution in Mendelian Populations. *Genetics* **16**:97–159. DOI: <https://doi.org/10.1093/genetics/16.2.97>, PMID: 17246615
- Wright S. 1938. Size of population and breeding structure in relation to evolution. *Science* **87**:430–431.
- Wright S. 1948. On the roles of directed and random changes in gene frequency in the genetics of populations. *Evolution* **2**:279–294. DOI: <https://doi.org/10.1111/j.1558-5646.1948.tb02746.x>, PMID: 18104586
- Yukilevich R. 2012. Asymmetrical patterns of speciation uniquely support reinforcement in *Drosophila*. *Evolution* **66**:1430–1446. DOI: <https://doi.org/10.1111/j.1558-5646.2011.01534.x>, PMID: 22519782
- Yukilevich R. 2017. *Drosophila* speciation patterns. <http://www.Drosophila-speciation-patterns.com> [Accessed May 27, 2020].
- Zhang Z-Q. 2013. Animal biodiversity: An update of classification and diversity in 2013. In: Zhang Z. Q (Ed). *Animal Biodiversity: An Outline of Higher-Level Classification and Survey of Taxonomic Richness (Addenda 2013)*. **3703** Zootaxa. p. 5–11. DOI: <https://doi.org/10.11646/zootaxa.3703.1.3>
- Zhao S, Zheng P, Dong S, Zhan X, Wu Q, Guo X, Hu Y, He W, Zhang S, Fan W, Zhu L, Li D, Zhang X, Chen Q, Zhang H, Zhang Z, Jin X, Zhang J, Yang H, Wang J, et al. 2013. Whole-genome sequencing of giant pandas provides insights into demographic history and local adaptation. *Nature Genetics* **45**:67–71. DOI: <https://doi.org/10.1038/ng.2494>, PMID: 23242367

## Appendix 1

### Simplified sweep effects model

I use a simplified model of the effects of recurrent hitchhiking and background selection (BGS) occurring uniformly along a genome. Expected diversity is given by

$$E(\pi) = \frac{\theta}{\theta + 1/B + 2NS} \quad (2)$$

$$\approx \frac{\theta}{1/B + 2NS} \quad (3)$$

(Equation 1 *Elyashiv et al., 2016*, Equation 4 of *Kim and Stephan, 2000*, and Equation 20 of *Coop and Ralph, 2012*). The BGS component is given by *Hudson and Kaplan, 1995*,

$$B(U, L) = N_e \exp\left(-\frac{U}{L}\right) \quad (4)$$

and the hitchhiking component is

$$S = \frac{\nu_{BP}}{r_{BP}} J \quad (5)$$

(*Coop and Ralph, 2012*, Equation 20) where  $\nu_{BP}$  and  $r_{BP}$  are the substitutions and recombination per basepair respectively,  $J$  is the probability that two lineages coalesce down to one, given sweeps occur uniformly along the genome. Under this homogeneous sweep model,  $J$  is

$$J = \int_0^L q_f(r)^2 dr \quad (6)$$

where  $q_f(r)$  is the approximate probability that a lineage is trapped by a sweep to frequency  $f$  when it is  $r$  recombination fraction away from this sweep (*Coop and Ralph, 2012*; Equation 15).

Since I use *Drosophila melanogaster* parameter estimates from *Elyashiv et al., 2016*, I now reconcile their model's  $S$  term with the simple model above. They estimate  $S$  in *Drosophila melanogaster* using a composite likelihood model that considers hitchhiking and background selection simultaneously, using substitutions and stratifying by annotation. For a neutral position at site  $x$ , the coalescence rate due to sweeps is given by Elyashiv et al.'s Equation 3,

$$S(x) = \frac{1}{T} \sum_{i_S} \alpha(i_S) \sum_{y \in a(i_S)} \int \exp(-r(x, y) \tau(s, N)) g(s|i_S) ds \quad (7)$$

where  $T$  is the length of the lineage (in generations) on which substitutions accrue,  $i_S = 1, \dots, I_S$  is the annotation class (e.g. exons, introns, UTRs),  $\alpha(i_S)$  is the fraction of substitutions in annotation class  $i_S$  that are beneficial,  $a(i_S)$  is the set of all substitutions in annotation class  $i_S$ ,  $\tau(s, N)$  is the fixation time of a site with additive effect  $s$ , and  $g(s|i_S)$  is the distribution of selection coefficients for annotation class  $i_S$ .

Note, that we can recover the model of *Coop and Ralph, 2012* from this expression. Suppose there is only one annotation class, and  $\alpha$  fraction of substitutions are beneficial, and one selection coefficient  $\bar{s}$ , (i.e.  $g(s) = \delta_0(s - \bar{s})$ ), then

$$S(x) = \frac{\alpha}{T} \sum_{y \in a} \exp(-r(x, y) \tau(\bar{s}, N)). \quad (8)$$

Let the number of substitutions be  $m := |a|$ , and imagine their positions are uniformly distributed on a segment of length  $G$  basepairs with the focal site is the middle at position  $x = 0$ . Then, each substitution  $y$  is a random distance  $l_y \sim U(-G/2, G/2)$  away from the focal site. Assuming the recombination rate is a constant  $r_{BP}$  per basepair, and approximating the sum with an integral, we have,



$$S = \frac{\alpha}{T} \sum_{i=1}^m E_i(\exp(-r_{BP} \ell_i \tau(\bar{s}, N))) \quad (9)$$

$$= \frac{\alpha}{TG} \sum_{i=1}^m \int_0^G \exp(-r_{BP} \ell \tau(\bar{s}, N)) d\ell \quad (10)$$

$$= \frac{\alpha m}{TG} \int_0^G \exp(-r_{BP} \ell \tau(\bar{s}, N)) d\ell \quad (11)$$

Using  $u$ -substitution with  $r = \ell r_{BP}$  this simplifies to

$$S = \frac{\alpha m}{TGr_{BP}} \int_0^L \exp(-r \tau(\bar{s}, N)) dr \quad (12)$$

where  $L = Gr_{BP}$ .

To simplify this notation, note that the rate of adaptive substitutions per basepair per generation is  $\nu_{BP} = \alpha m/GT$ , so

$$S = \frac{\nu_{BP}}{r_{BP}} \int_0^L \exp(-r \tau(\bar{s}, N)) dr \quad (13)$$

This is analogous to the second term of **Coop and Ralph, 2012**, Equation 17, with  $k = i = 2$  and  $x = 1$  (e.g. conditioning on a sweep to fixation). Note that there appears to be a factor of two error in **Elyashiv et al., 2016** compared to **Coop and Ralph, 2012**; here I include the factor of two. Then,

$$S = \frac{\nu_{BP}}{r_{BP}} \underbrace{\int_0^L \exp(-2r \tau(\bar{s}, N)) dr}_J \quad (14)$$

where the integral is equal to  $J$  ( $J_{2,2}$  of Equation 15 in **Coop and Ralph, 2012**) since a simple model of  $q_f(r) = f \exp(-2r \tau(s, N))$  and if we condition on fixation,  $f = 1$ . This expression is useful to generalize across species, since we know  $N$  and  $L$ . Additionally, we have estimates of  $\alpha$  and  $m/T$  in *Drosophila* and other species. In Elyashiv et al, they consider the number of substitutions per generation in genic regions only; it should be noted that the number of coding basepairs varies little across species. For convenience, I define  $\gamma = \alpha m/T$  as the number of adaptive substitutions per generation per entire genome, such that  $S(\gamma, J, L) = \frac{\gamma}{L} J$  used in the main text. Using the estimates of  $m \approx 4.5 \times 10^5$ ,  $\alpha \approx 0.42$ , and  $T \approx 8.4 \times 10^7$  from the Supplementary Material of Elyashiv et al., I arrive at  $\gamma \approx 0.00226$  adaptive substitutions per generation, per genome. For a  $\approx 100$  megabase genome, this translates to a  $\nu_{BP} \approx 2.34 \times 10^{-11}$ , which is close to previous estimates (**Figure 4—figure supplement 5A**). For  $J$ , I use an empirical estimate calculated from the genome-wide average of the rate of coalescent events due to sweeps, from Supplementary Table S6 of Elyashiv et al. ( $r_s = 2NS \approx 0.92$ ; see **Figure 4—figure supplement 5B**). This implies  $J \approx 4.46 \times 10^{-4}$ . Alternatively, I have tried using the estimated distribution of selection coefficients from Elyashiv et al., but this led to a weaker estimate of  $J$ , since the adaptive substitutions considered tend to cluster around genic regions.

## Appendix 2

### Background selection and polygenic fitness models

Throughout the main text, I use recurrent hitchhiking and background selection models to estimate the reduction in diversity due to linked selection. Another class of linked selection models, which I refer to as quantitative genetic linked selection models (QGLS; *Robertson, 1961; Santiago and Caballero, 1995*), can also depress genome-wide diversity. Furthermore, these models may depress diversity at neutral sites unlinked to the regions containing fitness variation. While I did not explicitly incorporate these models into my estimates of the diversity reductions, their effect is implicit in background selection models because they are analytically nearly identical. Here, I briefly sketch out the connection between BGS and QGLS models.

Under the *Santiago and Caballero, 1998* model, the effective population size is  $N_e^{SC98} = N \exp(-C^2/(1-Z)L)$ , where  $C^2$  is the standardized heritable fitness variation,  $1-Z$  is the decay of genetic variance through time, and  $L$  is the recombination map length. This model can accommodate a variety of modes of selection such as selection on an infinitesimal trait (*Santiago and Caballero, 1995*, p. 1016), and the flux of either weakly advantageous or deleterious alleles (*Santiago and Caballero, 1998*, p. 2109). If the source of fitness variation is entirely the input of new deleterious mutations with heterozygous effect  $sh$  at rate  $U$  per diploid genome per generation, then under mutation-selection balance, the equilibrium relative variance in reproductive success  $C^2 = Ush$  (*Crow and Kimura, 1970; Caballero, 2020*, p. 167), and  $Z = 1 - sh - 1/2N_c$  (*Santiago and Caballero, 1998*). Thus, if  $1/2N_c \ll sh \ll 1$ , then  $C^2/(1-Z) \approx U$  and  $N_e^{SC98} \approx N \exp(-U/L)$ , which is the BGS model used in the main text and is a result of many background selection models with similar assumptions (*Hudson and Kaplan, 1994*, Equation 15; *Hudson and Kaplan, 1995*, Equation 9; *Nordborg et al., 1996*, Equation 4; *Barton, 1995*, Equation 22b). Intuitively, the similarity of these models reflects the fact that a substantial proportion of heritable fitness variation is caused by the continual flux of deleterious alleles across the genome under mutation-selection balance (*Charlesworth, 2015; Charlesworth and Hughes, 2000*).

## Appendix 3

### Additional population size validation

In addition to the biomass-based validation described in the main text, I also conducted a few other consistency checks. First, note that the body-mass-based estimates of density for *Drosophila* are similar to previously used estimates in surveys of census size and diversity. **Nei and Graur, 1984** suggested a maximum of 5 *Drosophila* per m<sup>2</sup>, including regions of the range that are not inhabitable. Across *Drosophila*, the body mass based estimates suggest 10<sup>6.7</sup>–10<sup>7.6</sup> individuals per km<sup>2</sup>, or 4.5 – 36.3 individuals per m<sup>2</sup>, which are consistent with this previous estimate. Nei and Graur's estimates of *Drosophila pseudoobscura*'s census size are four orders of magnitude smaller than mine, but their approach uses a speculated ratio of population sizes of different *Drosophila* species rather than range sizes (**Nei and Graur, 1984**, p. 81).

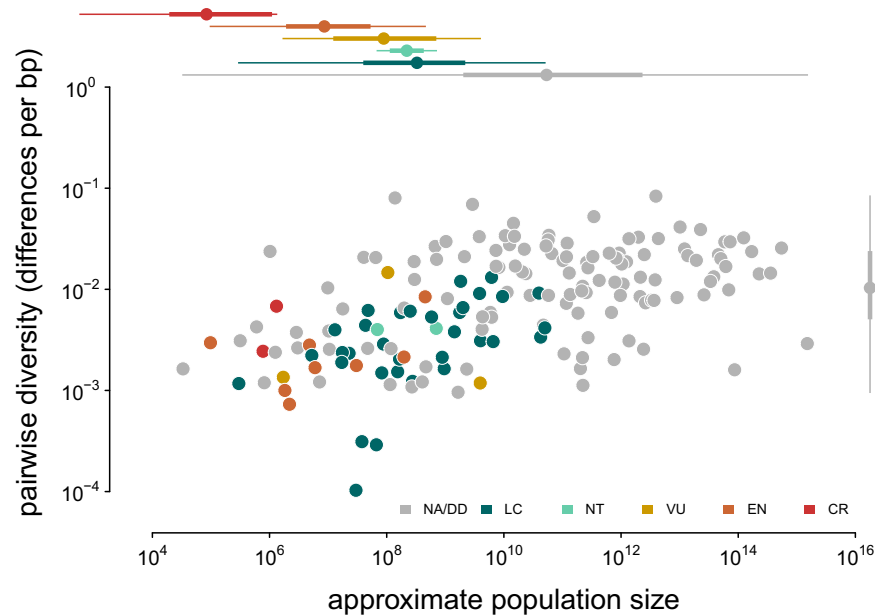
As another consistency check, I looked at the rank order of mammals by biomass. Whale species have the first and third highest biomass with 11.4 and 3.9 megatons of carbon biomass (for *Balaenoptera bonaerensis* and *Eschrichtius robustus*, respectively). While this seems high, a recent study shows that across whale species, pre-whaling carbon biomass was at the tens of megatons level (**Pershing et al., 2010, Table 1 and Figure 1**). Given that my census size estimates represent populations at a macroecological equilibrium, they would not reflect reduced density due to whaling or other anthropogenic causes. Humans had the second largest biomass, followed by wolf species (*Canis lupus* and *C. latrans*); as with whales, the population sizes for wolf species represent pre-anthropogenic densities and are overestimates compared to current population sizes, as expected.

Finally, there are other estimates of approximate population sizes for some species that I compared my estimates to. The United Nation's FAOSTAT database estimates the total number of horses (*Equus caballus*) on earth as ~60 million; the estimate in this study is close to 40 million. For other domesticated species like chicken (*Gallus gallus*), estimates range from 25 million to 19.6 billion (**FAOSTAT statistics database, 2021; Robinson et al., 2014**); the present study's estimate lies in the middle at ~175 million. Again, this is a known limitation of this method, as the range is estimated from occurrence data and does not consider species' niches. This present study's estimate of the number of king penguins (*Aptenodytes patagonicus*) is about 3 million; the population size was recently estimated as 2.23 million pairs (**Shirihai, 2008**).

## Appendix 4

### Diversity and IUCN Red List Status

I also investigated the relationship between species' IUCN Red List categories (an ordinal scale of how threatened a species is) and both diversity and population size, finding that species categorized as more threatened have both smaller population sizes and reduced diversity, compared to non-threatened species (**Appendix 4—figure 1**) consistent with past work (*Spielman et al., 2004*). A linear model of diversity regressed on population size has lower AIC when the IUCN Red List categories are included, and the estimates of the effect of IUCN status are all negative on diversity, though not all are significant in part because some categories have three or fewer species (**Appendix 4—table 1**).



**Appendix 4—figure 1.** A version of **Figure 2** with points colored by their IUCN Red List conservation status. Margin boxplots show the diversity and population size ranges (thin lines) and interquartile ranges (thick lines) for each category. NA/DD indicates no IUCN Red List entry, or Red List status Data Deficient; LC is Least Concern, NT is Near Threatened, VU is Vulnerable, EN is Endangered, and CR is Critically Endangered.

**Appendix 4—table 1.** The regression estimates of full IUCN Red List population size model for diversity,  $\log_{10}(\pi) = \beta_0 + \beta_{LC}LC + \beta_{NT}NT + \beta_{VU}VU + \beta_{EN}EN + \beta_{CR}CR + \beta_{N_c} \log_{10}(N_c)$ ;  $df = 165$ . Using AIC to compare this full model to a reduced model of  $\log_{10}(\pi) = \beta_0 + \beta_{N_c} \log_{10}(N_c)$ ,  $AIC_{full} = 204.9$ ,  $AIC_{reduced} = 216.4$ .

|               | Mean  | 2.5 % | 97.5 % |
|---------------|-------|-------|--------|
| $\beta_0$     | -2.80 | -3.20 | -2.50  |
| $\beta_{LC}$  | -0.39 | -0.57 | -0.21  |
| $\beta_{NT}$  | -0.22 | -0.83 | 0.39   |
| $\beta_{VU}$  | -0.34 | -0.84 | 0.16   |
| $\beta_{EN}$  | -0.40 | -0.73 | -0.07  |
| $\beta_{CR}$  | -0.03 | -0.65 | 0.59   |
| $\beta_{N_c}$ | 0.08  | 0.05  | 0.11   |