

# UCLA

## UCLA Previously Published Works

### Title

Greater strength of selection and higher proportion of beneficial amino acid changing mutations in humans compared with mice and *Drosophila melanogaster*

### Permalink

<https://escholarship.org/uc/item/43q8s3pr>

### Journal

Genome Research, 31(1)

### ISSN

1088-9051

### Authors

Zhen, Ying

Huber, Christian D

Davies, Robert W

et al.

### Publication Date

2021

### DOI

10.1101/gr.256636.119

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial License, available at <https://creativecommons.org/licenses/by-nc/4.0/>

Peer reviewed

# Greater strength of selection and higher proportion of beneficial amino acid changing mutations in humans compared with mice and *Drosophila melanogaster*

Ying Zhen,<sup>1,2,3</sup> Christian D. Huber,<sup>1,4</sup> Robert W. Davies,<sup>5,6</sup> and Kirk E. Lohmueller<sup>1,7</sup>

<sup>1</sup>Department of Ecology and Evolutionary Biology, University of California, Los Angeles, California 90095, USA; <sup>2</sup>Zhejiang Provincial Laboratory of Life Sciences and Biomedicine, Key Laboratory of Structural Biology of Zhejiang Province, School of Life Sciences, Westlake University, Hangzhou, Zhejiang, 310024, China; <sup>3</sup>Institute of Biology, Westlake Institute for Advanced Study, Hangzhou, Zhejiang, 310024, China; <sup>4</sup>School of Biological Sciences, The University of Adelaide, Adelaide, South Australia 5005, Australia; <sup>5</sup>Program in Genetics and Genome Biology and The Centre for Applied Genomics, The Hospital for Sick Children, Toronto, Ontario, M5G 0A4, Canada; <sup>6</sup>Department of Statistics, University of Oxford, Oxford, OX1 3LB, United Kingdom; <sup>7</sup>Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, California 90095, USA

Quantifying and comparing the amount of adaptive evolution among different species is key to understanding how evolution works. Previous studies have shown differences in adaptive evolution across species; however, their specific causes remain elusive. Here, we use improved modeling of weakly deleterious mutations and the demographic history of the outgroup species and ancestral population and estimate that at least 20% of nonsynonymous substitutions between humans and an outgroup species were fixed by positive selection. This estimate is much higher than previous estimates, which did not correct for the sizes of the outgroup species and ancestral population. Next, we jointly estimate the proportion and selection coefficient ( $p^+$  and  $s^+$ , respectively) of newly arising beneficial nonsynonymous mutations in humans, mice, and *Drosophila melanogaster* by examining patterns of polymorphism and divergence. We develop a novel composite likelihood framework to test whether these parameters differ across species. Overall, we reject a model with the same  $p^+$  and  $s^+$  of beneficial mutations across species and estimate that humans have a higher  $p^+s^+$  compared with that of *D. melanogaster* and mice. We show that this result cannot be caused by biased gene conversion or hypermutable CpG sites. We discuss possible biological explanations that could generate the observed differences in the amount of adaptive evolution across species.

[Supplemental material is available for this article.]

Since the inception of molecular population genetics, there has been tremendous interest in quantifying the amount of adaptive evolution in different organisms. The neutral theory of molecular evolution postulated that beneficial mutations are rare and that many of the substitutions between species are neutral (Kimura 1983). One early challenge to this theory originated from a comparison of polymorphisms and divergence at synonymous and nonsynonymous sites in *Drosophila* (Fay et al. 2002; Smith and Eyre-Walker 2002). Under neutral models, the ratio of nonsynonymous to synonymous changes should remain equal when comparing polymorphisms (i.e., differences within species) and divergence (i.e., differences between species). In contrast to this prediction, a genome-wide excess of nonsynonymous divergence between species was observed, a pattern indicative of an abundance of positive selection in *Drosophila*. Smith and Eyre-Walker (2002) proposed a statistic,  $\alpha$ , which is the proportion of nonsynonymous substitutions between species that can be attributed to positive selection. Their approach has found that at least 40% of nonsynonymous substitutions have been fixed by positive selection in *Drosophila* (Smith and Eyre-Walker 2002).

Since the publication of the original study,  $\alpha$  has been estimated from different species across the tree of life (Fay 2011).

Estimates of  $\alpha$  vary tremendously across species, tending to be higher in insects (Andolfatto 2005; Eyre-Walker and Keightley 2009) but much lower in primates and plants (Boyko et al. 2008; Eyre-Walker and Keightley 2009; Gossmann et al. 2010). In these latter species, formal tests have been unable to reject the hypothesis that  $\alpha$  is zero (i.e., no positive selection) (Boyko et al. 2008; Foxe et al. 2008; Eyre-Walker and Keightley 2009). It is not clear why  $\alpha$  varies across species. One possibility is that  $\alpha$  is higher for species with larger population sizes, which could occur if adaptation is mutation limited. Here, species with larger population sizes would have a higher rate of beneficial mutations. The fixation probability of a given beneficial mutation also would be higher in species with larger population size, but this effect is likely to only be important for very weakly beneficial mutations. Evidence indicates that, in some cases,  $\alpha$  is correlated with population size. For example, Phifer-Rixey et al. (2012) found that estimates of  $\alpha$  were higher for species of mice that have larger population sizes compared with species with smaller population sizes. Further, there is a positive correlation between  $\alpha$  and population size when comparing different species of sunflowers (Strasburg et al. 2011) and from phylogenetically diverse taxa

**Corresponding authors:** klohmue@ucla.edu, zhenying@westlake.edu.cn

Article published online before print. Article, supplemental material, and publication date are at <http://www.genome.org/cgi/doi/10.1101/gr.256636.119>.

© 2021 Zhen et al. This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <http://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

(Gossmann et al. 2012). More recently, Galtier (2016) found a positive correlation between  $\alpha$  and effective population size for 44 animal species. Additional evidence that there is more positive selection in larger populations stems from analyses of selection at linked sites. Corbett-Detig et al. (2015) found increased evidence for selection at linked sites in species with larger population sizes, although the mechanism driving this pattern is not entirely clear. Further, Nam et al. (2017) have suggested that across primates, species with larger population sizes experienced more selective sweeps.

Although evidence suggests that adaptation could be mutation limited and that this could be driving the variation in  $\alpha$  across species, other factors can influence the  $\alpha$  statistic (Messer and Petrov 2013; Rousselle et al. 2018). As the denominator of  $\alpha$  is the total number of observed differences between species, it is sensitive to the fixation of weakly deleterious mutations. For two populations with the same number of beneficial substitutions, the one with a higher number of substitutions owing to weakly deleterious mutations will have a lower  $\alpha$ . Indeed, because the number of fixed weakly deleterious mutations is inversely related to population size (Ohta 1973, 1992), this effect could drive the correlation between  $\alpha$  and population size. In support of this prediction, Galtier (2016) found that the rate of adaptive divergence relative to neutral divergence (omega-a) (Gossmann et al. 2012) showed no correlation with population size. Similar arguments have been made by Phifer-Rixey et al. (2012). Further, for humans, it has been suggested that  $\alpha$  has been underestimated owing to the presence of weakly beneficial mutations segregating as polymorphisms (Galtier 2016; Uricchio et al. 2019). Indeed, methods that account for weakly beneficial mutations segregating as polymorphisms infer slightly higher values of  $\alpha=0.24$  between humans and chimpanzee and  $\alpha=0.135$  in the human lineage (Galtier 2016; Uricchio et al. 2019). In addition, population sizes of the outgroup and ancestral population determine the rate of fixation of weakly deleterious mutations in the outgroup lineage and in the ancestral population and could potentially influence the estimate of  $\alpha$  (McDonald and Kreitman 1991; Eyre-Walker 2002; Rousselle et al. 2018).

Other studies quantified positive selection by focusing on the proportion of beneficial mutations ( $p^+$ ) and their selection coefficients ( $s^+$ ). Boyko et al. (2008) found that by assuming a fraction (0%–1.86%) of new mutations is positively selected, they could better match the frequency spectrum of polymorphisms and the counts of human–chimpanzee divergence. Models with weaker se-

lection coefficients for beneficial mutations tended to have a higher proportion of positively selected mutations than models with stronger selection (Boyko et al. 2008). Several studies also have estimated  $p^+$  and  $s^+$  in *Drosophila* species and mice (Sella et al. 2009; Schneider et al. 2011; Elyashiv et al. 2016; Keightley et al. 2016; Campos et al. 2017; Booker and Keightley 2018). However, there has not been a systematic comparison across species.

Here we compare the amount of adaptive evolution in primates, rodents, and *Drosophila*. We use two complementary approaches that quantify different aspects of the adaptive process. First, we use improved modeling of weakly deleterious mutations and demographic models, particularly correcting for the sizes of outgroup species and ancestral population to infer  $\alpha$ . Second, we jointly estimate  $p^+$  and  $s^+$  of newly arising beneficial mutations by examining patterns of polymorphism and divergence. We develop a composite likelihood framework to test whether these parameters differ across taxa. This approach enables a more direct comparison of beneficial mutations across species and is less confounded by the fixation of weakly deleterious mutations.

## Results

### Estimates of $\alpha$ for multiple species using the MK method

We first estimated  $\alpha$  from coding regions of primates, rodents, and *Drosophila*. We analyzed published genomic data sets to obtain counts of synonymous and nonsynonymous polymorphisms ( $P_S$  and  $P_N$ ) and divergent sites between species ( $D_S$  and  $D_N$ ). For computation of  $\alpha$  in humans, we used chimpanzee and macaque as outgroup species. For mice and *D. melanogaster*, we used rat and *Drosophila simulans* as the outgroup species, respectively.

An extension of the McDonald–Kreitman (MK) test was used to estimate  $\alpha$  (Table 1; Supplemental Table S1; Smith and Eyre-Walker 2002). To examine the effect of slightly deleterious mutations on  $\alpha$ , we filtered the data with several minor allele frequency (MAF) cutoffs (Supplemental Table S2; Fay et al. 2001; Charlesworth and Eyre-Walker 2008; Messer and Petrov 2013). For example, after removing low-frequency polymorphisms with MAF <20%, the estimated  $\alpha$  is close to zero for the human–chimpanzee comparison (Table 1), consistent with previous estimates (Boyko et al. 2008; Eyre-Walker and Keightley 2009). However, for the human–macaque comparison,  $\alpha$  is  $-0.22$  (Table 1), suggesting that the choice of outgroup species could greatly influence  $\alpha$ .

**Table 1.** Estimates of  $\alpha$  and omega-a using different methods

Species	Outgroup	Method of inference								
		MK-method		Model-based				DFE $\alpha$		
		All	MAF > 20%	Simple model		Complex model		$\alpha$	omega-a	
				$\alpha$	omega-a	$\alpha$	omega-a	$\alpha$	omega-a	
Full data										
Human	Chimpanzee	-0.41	-0.01	0.11	0.03	0.25	0.07	0.24	0.07	
Human	Macaque	-0.70	-0.22	-0.08	-0.02	0.26	0.06	0.02	0.01	
Human lineage	—	—	—	0.06	—	0.16	—	—	—	
<i>D. melanogaster</i>	<i>D. simulans</i>	-0.13	0.49	0.53	0.08	0.60	0.09	0.71	0.14	
Mice	Rat	0.25	0.40	0.45	0.11	0.41	0.10	0.51	0.12	
SSWW only										
Human	Chimpanzee	-0.37	0.08	0.13	0.04	0.26	0.08	—	—	
Mice	Rat	-0.14	0.08	0.20	0.03	0.10	0.02	—	—	

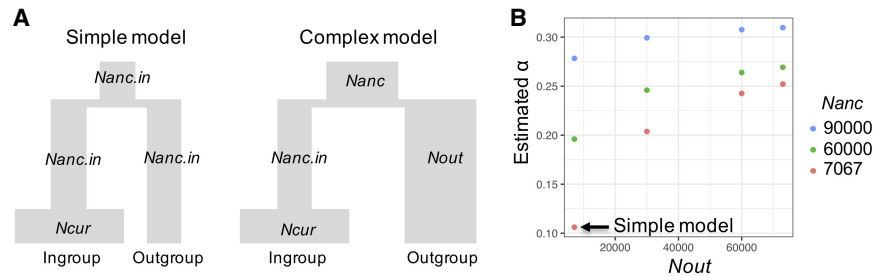
In contrast, for the *D. melanogaster*–*D. simulans* and mouse–rat comparisons, the estimated values of  $\alpha$  are 49% and 40%, respectively, with the MAF filter at 20% (Table 1). Both of these estimates are comparable to previous studies (Andolfatto 2005; Eyre-Walker and Keightley 2009; Phifer-Rixey et al. 2012). These results suggest that the proportion of substitutions fixed by positive selection varies drastically across species.

### Model-based inference of $\alpha$

Model-based approaches estimate  $\alpha$  by contrasting the observed number of nonsynonymous differences between species ( $D_{NO}$ ) with the number expected under a demographic model and a distribution of fitness effects (DFE) including only neutral and deleterious mutations ( $D_{NE}$ ) (Supplemental Table S3; Loewe et al. 2006; Boyko et al. 2008; Eyre-Walker and Keightley 2009; Haddrill et al. 2010). The excess of observed  $D_{NO}$  compared with the predicted  $D_{NE}$  is attributed to fixations driven by positive selection. These methods assume that the population sizes of outgroup species and the ancestral population of the ingroup and outgroup are the same as the ancestral size of the ingroup population (*Nanc.in*). We refer to this demographic model as the Simple model (Fig. 1A).

By using the Simple model, we estimated that  $\alpha = 11\%$  for humans and chimpanzees. This estimate is comparable to the inference by Boyko et al. (2008). When we use macaque as the outgroup species, we estimated that  $\alpha$  is negative (Table 1). Although it is possible that this difference could reflect distinct evolutionary events experienced by different outgroups, it could also be an artifact of the modeling assumptions. Previous estimates of the human–chimpanzee ancestral population size (*Nanc*) range from 12,000 to 125,000 using different data sets and methods (Chen and Li 2001; Hobolth et al. 2007; Burgess and Yang 2008; Prado-Martinez et al. 2013; Schrago 2014). The estimated chimpanzee population size (*Nout*) is around 30,000 (Supplemental Table S4; Fischer et al. 2004; Hvilson et al. 2012; Prado-Martinez et al. 2013). However, the ancestral size inferred for the human lineage using human polymorphism data is 7067, which is much smaller than all previous estimates of *Nanc* and *Nout*. Using values of *Nanc* and *Nout* that are too small likely biases estimates of  $\alpha$  because more of the nonsynonymous substitutions are incorrectly attributed to the fixation of weakly deleterious mutations, causing  $\alpha$  to be underestimated.

To more accurately model the larger *Nanc* and *Nout*, we use the Complex model, which allows *Nout* and *Nanc* to differ from *Nanc.in* (Fig. 1A). For humans, the larger *Nanc* is very ancient; thus, it does not affect the polymorphism pattern within humans (confirmed by coalescent simulations) (see Supplemental Text). We calculated the number of substitutions fixed in the ingroup, outgroup, and ancestral populations (Methods). We first explored how *Nout* and *Nanc* affect the inferred  $\alpha$ . For the human–chimpanzee comparison, the inferred  $\alpha$  changes with different values of *Nout* and *Nanc*. Larger values result in larger estimates of  $\alpha$  (Fig. 1B). For example, when *Nout* = 30,000 (Fischer et al. 2004; Hvilson et al. 2012; Prado-Martinez et al. 2013) and *Nanc* = 60,000 (Chen and Li 2001; Hobolth et al. 2007; Prado-Martinez et al. 2013), supported by many previous studies (Supplemental



**Figure 1.** The ancestral and outgroup population sizes greatly influence  $\alpha$ . (A) Schematic demographic models illustrate the Simple and Complex models with associated parameters. In the Simple model, the size of the ancestral population (*Nanc*) and the size of the outgroup (*Nout*) are assumed to be the same as the ancestral size of the ingroup (*Nanc.in*). This assumption is relaxed in the Complex model. (B) Effect of *Nanc* and *Nout* on estimates of  $\alpha$  for humans using the chimpanzee as an outgroup. Colors denote different values of *Nanc*. Arrow points to the estimate of  $\alpha$  from the Simple model, where *Nanc* = *Nout* = *Nanc.in* = 7067.

Table S4),  $\alpha$  is approximately 24.6%. When using other values of *Nout* and *Nanc* within the ranges of previous estimates (Supplemental Table S4), the corresponding estimates of  $\alpha$  may differ, but remain above 20% (Fig. 1B). We next revisited the human–macaque comparison under the Complex model (Supplemental Fig. S1), using *Nout* = 73,000 (Hernandez et al. 2007; Xue et al. 2016), *Nanc* = 48,000 (McVicker et al. 2009), and changing the ingroup population size to 60,000 at human–chimpanzee divergence time to match what we modeled in human–chimpanzee (Supplemental Table S5). Here we infer that  $\sim 26.0\%$  of human–macaque nonsynonymous substitutions were fixed by positive selection, which is comparable to what we found for the human–chimpanzee analysis. These estimates of  $\alpha$  for primates using the more realistic Complex demographic model are much higher than previous estimates, implying that there is a greater contribution of positive selection to nonsynonymous divergence than previously appreciated.

Similarly, we estimated  $\alpha$  for *Drosophila* and rodents using the Complex demographic model. For *Drosophila*, it had been inferred that *D. simulans* have slightly larger  $N_e$  than *D. melanogaster* (Andolfatto et al. 2011), so we set *Nout* to be 1.5 $\times$  the current population size of *D. melanogaster* (*Ncur*) at the species' split (Supplemental Table S5). For mice, a previous study estimated that the outgroup rat species has an effective population size about fivefold lower than that of wild house mice (Ness et al. 2012). Thus, we set the *Nout* to be 0.2 $\times$  *Ncur* of mice (Supplemental Table S5). Because there is limited knowledge of the population sizes of the ancestor of *D. melanogaster* and *D. simulans* and of the ancestor of mice and rat, for these two comparisons, we assume *Nanc* = *Nout*. By using these Complex models for *D. melanogaster*–*D. simulans* and mouse–rat, we estimated their  $\alpha$  to be 60% and 41%, compared with 53% and 45% using the Simple model, respectively (Table 1; for 95% confidence intervals, see Supplemental Table S1). The differences between these estimates reflect the importance of accurately modeling the population size of outgroup species for calculations of  $\alpha$ .

We estimated  $\alpha$  for substitutions that occurred exclusively on the human lineage, using the human–macaque alignment to polarize sequence differences between human and chimpanzee. We estimate that  $\alpha = 6.2\%$  using the Simple model and 16.0% using the Complex model (Table 1; Supplemental Text).

To compare our estimates of  $\alpha$  to those from another model-based method that assumes *Nanc* = *Nout* = *Nanc.in*, we used *DFE-alpha*, to infer  $\alpha$  for our three species pairs. Using this method,  $\alpha$  is

estimated to be 24% and 2% for humans using chimpanzee and macaque as outgroup, respectively; 71% for *D. melanogaster*–*D. simulans*; and 51% for mouse–rat (Table 1). These estimates are all higher compared with estimates from Simple demographic models. However, the estimates of  $\alpha$  for primates differ significantly depending on whether the macaque or chimpanzee is used as the outgroup for humans.

In addition, we calculated omega-a for both the Simple model and the Complex model for all taxa (Table 1). We find similar patterns of omega-a as for  $\alpha$ . First, using the Complex models, estimates of omega-a for primates are similar (i.e., 0.06 or 0.07) regardless of outgroup. Second, the estimate of omega-a for primates is higher using the Complex model than the Simple model. Across taxa, primates have the lowest estimate of omega-a, compared with *Drosophila* (0.09) and rodents (0.10) using the full data set.

### Testing whether $p^+$ and $s^+$ differ across species

We next estimated the proportion and selection coefficient of new beneficial mutations. For each species, we estimate the proportion of new mutations that are beneficial ( $p^+$ ) and their selection coefficient ( $s^+$ ) jointly using a grid search approach. We then test whether these two parameters differ across species.

The number of nonsynonymous differences between a pair of species ( $D_N$ ) is assumed to be Poisson-distributed (Sawyer and Hartl 1992), with rate parameter equal to

$$E[D_N] = 2N\mu[G(s)u(s)(1 - p^+) + u(s^+)p^+], \quad (1)$$

where  $G(s)$  is the DFE of deleterious and neutral mutations from Huber et al. (2017),  $u(s)$  is the fixation probability (Kimura 1962) of deleterious and neutral mutations,  $u(s^+)$  is the fixation probability of beneficial mutations, and  $p^+$  is the proportion of new mutations that are beneficial. We then use a Poisson log-likelihood (LL) function for  $D_N$  in each species and a series of likelihood ratio tests (LRTs) to determine whether  $p^+$  and  $s^+$  differ across primates, rodents, and *Drosophila* (see Methods).

By using this framework, we find that the full model H1, in which each taxon is allowed to have its own  $p^+$  and  $s^+$ , fits  $D_N$  significantly better than the constrained null model, where  $p^+$  and  $s^+$  are constrained to be the same across all three taxa (LRT statistic  $\Lambda = 124,974$ ,  $df = 4$ ,  $P < 10^{-16}$ ) (Fig. 2A–D; Supplemental Fig. S2;

Supplemental Table S6). We also compared  $p^+$  and  $s^+$  between pairs of taxa (i.e., primate vs. rodent; primate vs. *Drosophila*; rodent vs. *Drosophila*). In all pairwise tests, the model in which each species has its own  $p^+$  and  $s^+$  fits the observed  $D_N$  significantly better than a model in which  $p^+$  and  $s^+$  are constrained to be the same in the tested taxa (Supplemental Table S6). Many combinations of  $p^+$  and  $s^+$  values show similar LLs (Fig. 2A–C), suggesting the two parameters are not separable. Models with a larger  $s^+$  have a lower proportion of positively selected mutations ( $p^+$ ) than models with a smaller  $s^+$ .

We next investigated whether when allowing  $p^+$  to differ across primates, rodents, and *Drosophila*, a model with the same  $s^+$  could fit all species. This is shown in conditional likelihood plots, where when assuming the same  $s^+$  for all species, humans would need a higher proportion of beneficial mutations compared with mice and *D. melanogaster* to match the observed  $D_N$  (Fig. 3A). When we constrain  $p^+$  to be the same for all species,  $s^+$  is larger in humans compared with that in mice and *D. melanogaster* (Fig. 3B).

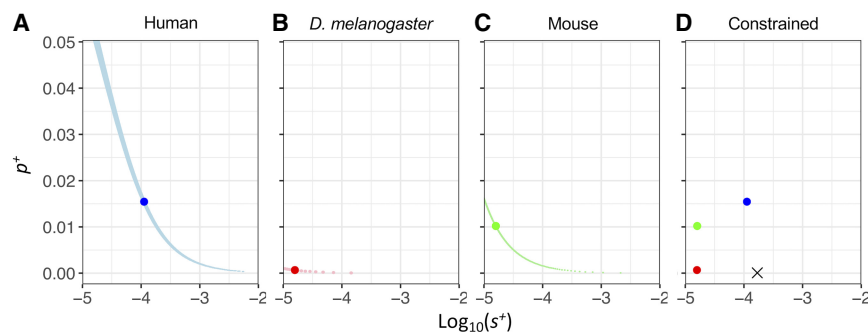
Motivated by the observation that  $p^+$  and  $s^+$  cannot be reliably estimated separately, we estimated a composite parameter  $p^+s^+$ , the product of  $p^+$  and  $s^+$ . We find clear evidence that, regardless of the model of ancestral demography or outgroup, humans have a significantly higher  $p^+s^+$  than do *D. melanogaster* and mice (Fig. 4). The LL curves are quite peaked and have little overlap across species. Thus, approximate 95% confidence intervals on  $p^+s^+$  do not overlap across species, implying more positive selection in humans than in *D. melanogaster* and mice.

### Testing whether $\gamma^+$ and $p^+$ differ across species

Humans, *D. melanogaster*, and mice have drastically different population sizes, which can influence the efficacy of selection within each species. Thus, we next examined whether the selection coefficient scaled by current population size ( $\gamma^+ = 2Ns^+$ ) and  $p^+$  differ across primates, rodents, and *Drosophila*.

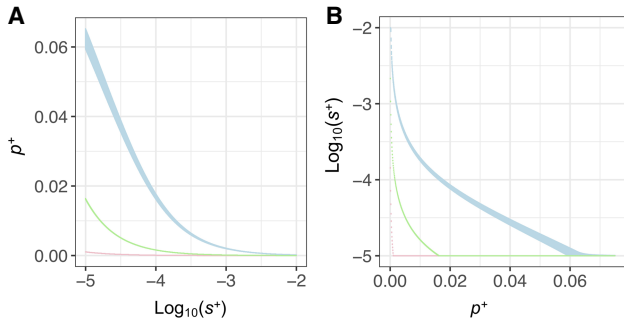
We find that the full model (H1) in which each taxon has its own  $\gamma^+$  and  $p^+$  fits the observed  $D_N$  significantly better than the constrained model (H0) in which  $\gamma^+$  and  $p^+$  are constrained to be the same across all three taxa (LRT statistic  $\Lambda = 3541$ ,  $df = 4$ ,  $P < 10^{-16}$ ) (Fig. 5A–D; Supplemental Fig. S2; Supplemental Table S6). We also compared  $\gamma^+$  and  $p^+$  between pairs of taxa. For the human–*D. melanogaster* and human–mice pairs, the models in which

each species has its own  $\gamma^+$  and  $p^+$  fit the observed  $D_N$  significantly better than a model in which  $\gamma^+$  and  $p^+$  are constrained to be the same in the tested two taxa (Supplemental Table S6), regardless of outgroup. However, we cannot reject the hypothesis that mice and *D. melanogaster* have same  $\gamma^+$  and  $p^+$  (Supplemental Table S6). Similar to what was seen above for  $p^+$  and  $s^+$ , a ridge of parameter values have similar LLs (Fig. 5A–C). Models with larger  $\gamma^+$  tended to have a lower proportion of positively selected mutations than models with smaller  $\gamma^+$ , suggesting the two parameters are not separable. The relative ordering of the MLEs of  $p^+$  across species considering  $\gamma^+$  is not necessarily the same as that for the  $s^+$  results because the population sizes differ across species.



**Figure 2.** Log-likelihood (LL) surfaces for  $p^+$  and  $s^+$  for different species: (A) human; (B) *D. melanogaster*; and (C) mouse. (D) The constrained model (H0), in which  $p^+$  and  $s^+$  are constrained to be the same across all three taxa. LLs are calculated using a grid search method of  $\log_{10}(s^+)$  in the range of  $-5$  to  $-2$  and  $p^+$  in the range of  $0\%$ – $7.5\%$ . Blue denotes human; red, *D. melanogaster*; and green, mouse. The large points represent the MLE for each species; the black cross in panel D represents the MLE of the constrained model; and the lighter colors show grid points within three LL units of each MLE. The Complex model is used for each species, and we use the chimpanzee as the outgroup for humans.





**Figure 3.** Conditional LL surfaces: (A) maximizing  $p^+$  given particular values of  $s^+$  and (B) maximizing  $s^+$  given particular values of  $p^+$ . Only grid points within three LL units of the MLEs for each parameter for each species are shown. Light blue denotes human; pink, *D. melanogaster*; and light green, mouse.

We next investigated whether it is possible that either  $\gamma^+$  or  $p^+$  is the same across taxa while the other parameter varies. With the same  $\gamma^+$  for all species, humans would need a lower proportion of beneficial mutations compared with that of mice and *D. melanogaster* to fit the observed  $D_N$  (Supplemental Fig. S3A). When we constrain  $p^+$  to be the same for all taxa, humans have a smaller  $\gamma^+$  for beneficial mutations (Supplemental Fig. S3B). Furthermore, because it is challenging to reliably estimate  $p^+$  and  $\gamma^+$  separately, we also estimated a composite parameter  $p^+\gamma^+$ , the product of  $p^+$  and  $\gamma^+$ . We find that  $p^+\gamma^+$  is not as distinct across different taxa as is  $p^+s^+$ . In the Complex models, humans have a smaller  $p^+\gamma^+$  than do mice. However, the LL curves of  $p^+\gamma^+$  overlap between *D. melanogaster* and mice and between *D. melanogaster* and humans, suggesting they are not significantly different (Supplemental Fig. S4).

### Effects of biased gene conversion and hypermutable CpG sites

Biased gene conversion (BGC) is the preferred transmission of G/C alleles (S indicates strong alleles) at the expense of A/T alleles (W indicates weak alleles). This process is common in mammals (Duret and Galtier 2009; Lachance and Tishkoff 2014; Bolívar et al. 2016) but is not as evident in *D. melanogaster* (Robinson et al. 2014; Jackson et al. 2017). Further, methylated CpG sites in mammals have a higher mutation rate to T alleles owing to deamination of the C nucleotide (Duncan and Miller 1980; Sved and Bird 1990). Both factors could possibly bias comparisons of selection between mammals and insects.

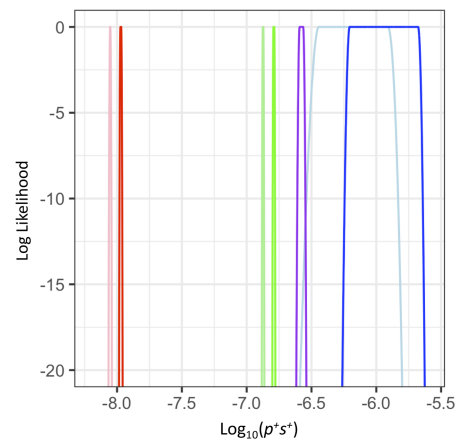
To test whether BGC and hypermutable CpG sites drive the observed pattern of positive selection across primates, rodents, and *Drosophila*, we filtered primate and rodent data to keep only strong-to-strong or weak-to-weak mutations (herein called SSWW mutations), which are not affected by BGC and are not CpG changes (Supplemental Text; Supplemental Fig. S5). We then reinfered  $\alpha$ . For primates,  $\alpha$  is comparable to that without filtering (Table 1). For rodents, however,  $\alpha$  is substantially lower than before filtering, suggesting that BGC and CpG mutational processes may account for some of the nonsynonymous differences between the mouse and rat. We also calculated omega-a using SSWW mutations. When using the Complex models, omega-a for primates is slightly higher than that before filtering, whereas for rodents, it is much lower than before filtering (Table 1).

After removing the effect of BGC and CpG sites in primates and rodents, we again quantify the strength and proportion of new beneficial mutations across all three taxa. A model in which

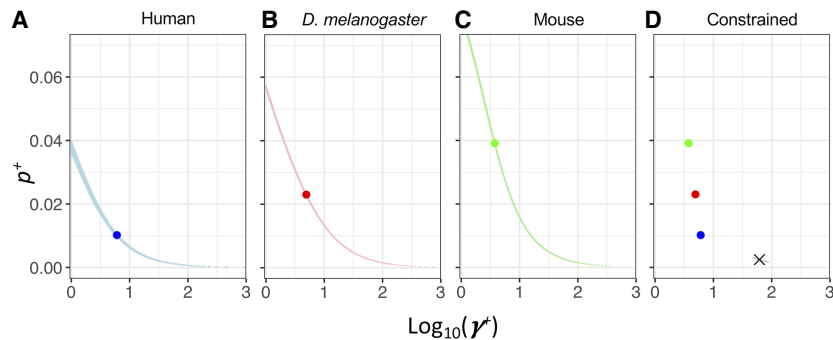
each taxon has its own  $p^+$  and  $s^+$  fits the observed  $D_N$  significantly better than a model in which  $p^+$  and  $s^+$  are constrained to be the same across all three taxa (LRT statistic  $\Lambda = 825$ ,  $df = 4$ ,  $P < 10^{-16}$ ) (Supplemental Fig. S6A–D; Supplemental Table S6). Models comparing each pair of taxa (e.g., primates vs. rodents; primates vs. *Drosophila*; rodents vs. *Drosophila*) suggest that each taxon has its own unique  $p^+$  and  $s^+$ , regardless of demography (Supplemental Table S6). Allowing  $p^+$  to differ across taxa, a model with the same  $s^+$  across all taxa could fit the data and vice versa. When we constrain  $s^+$  to be the same for all species, humans still have the highest proportion of new beneficial mutations (Supplemental Fig. S6E). Constraining  $p^+$  to be the same for all species, humans have the largest selection coefficient (Supplemental Fig. S6F).

Similarly, the model in which each taxon has its own  $p^+$  and  $\gamma^+$  fits the observed  $D_N$  significantly better than a model in which  $p^+$  and  $\gamma^+$  are constrained to be the same across all three taxa (LRT  $\Lambda = 4496$ ,  $df = 4$ ,  $P < 10^{-16}$ ) (Supplemental Table S6; Supplemental Fig. S7). Models comparing each pair of taxa suggest that each taxon has its own unique  $p^+$  and  $\gamma^+$ , regardless of demography (Supplemental Table S6). Allowing  $p^+$  to differ across species, a model with the same  $\gamma^+$  across all taxa could fit the data and vice versa. When we constrain  $\gamma^+$  to be the same for all species, *D. melanogaster* has the highest proportion of beneficial new mutations, and mice have the lowest proportion. When we constrain  $p^+$  to be the same for all species, *D. melanogaster* has the largest  $\gamma^+$ , and mice have the smallest  $\gamma^+$  (Supplemental Fig. S7).

Last, we estimated two composite parameters,  $p^+s^+$  and  $p^+\gamma^+$ , for SSWW mutations. Humans again have a significantly higher  $p^+s^+$  than do *D. melanogaster* and mice (Supplemental Fig. S8), suggesting more positive selection in humans than in *D. melanogaster* and mice. Similar to the full data set,  $p^+\gamma^+$  is not as distinct across different species as  $p^+s^+$ . We find that under the Complex model, humans and *D. melanogaster* have a significantly (i.e., the LL curves



**Figure 4.** The composite parameter  $p^+s^+$ , capturing the proportion of beneficial mutations and the strength of selection, differs across species. LL curves for  $p^+s^+$  in the three species. Red denotes the inference for *D. melanogaster*; green, the inference for mouse; blue, the inference for human using the chimpanzee as the outgroup; and purple, the inference for human using the macaque as the outgroup. Lighter colors denote the Simple model. Darker colors denote the Complex model, which better models the ancestral demography and population size of the outgroup. Note that regardless of which demographic model is used, the LL curves from the different species do not overlap within the top 500 LL units, suggesting  $p^+s^+$  is significantly different among taxa.



**Figure 5.** LL surfaces for  $p^+$  and  $\gamma^+$  for different species: (A) human; (B) *D. melanogaster*; and (C) mouse. (D) The constrained model, H0, in which  $p^+$  and  $\gamma^+$  are constrained to be the same across all three taxa. LLs are calculated using a grid search method of  $\log_{10}(\gamma^+)$  in the range of 0–3 and  $p^+$  in the range of 0%–7.5%. Blue denotes human; red, *D. melanogaster*; and green, mouse. The large points represent the MLE for each species; the black cross in panel D represents the MLE of the constrained model; and the lighter colors show grid points within three LL units of each MLE. The Complex model is used for each species, and we use chimpanzee as the outgroup for humans.

do not overlap within two units of the MLEs) larger  $p^+\gamma^+$  than mice, and the LL curves of  $p^+\gamma^+$  for humans and *D. melanogaster* overlap near the MLEs (Supplemental Fig. S8).

## Discussion

We have quantified the amount of adaptive evolution in multiple species with varying degrees of complexity and population size using two different approaches. First, we examined the proportion of nonsynonymous substitutions fixed by positive selection and found that this proportion is higher in *Drosophila* than in rodents and primates, consistent with prior work (Andolfatto 2005; Boyko et al. 2008; Eyre-Walker and Keightley 2009; Fay 2011). However, after correcting for the population sizes of the outgroup and ancestral population, we infer that  $\alpha$  is much higher in primates than what was inferred in previous studies. Second, we show that the proportion of and strength of positive selection on new beneficial mutations differs across species. The species with the smaller population size and greater complexity (i.e., humans) has stronger and/or more abundant new beneficial mutations than the other two species with much larger population sizes (i.e., mice and *D. melanogaster*). These results are robust to the choice of outgroup, BGC, and hypermutable CpG sites.

One major advantage of our method to infer  $\alpha$  over previous similar approaches from Boyko et al. (2008) and *DFE-alpha* (Eyre-Walker and Keightley 2009) is that we allow the outgroup and ancestral population sizes to differ from the ancestral population size of the ingroup. The population size of the outgroup matters because it affects the fixation probability of weakly deleterious alleles in the outgroup lineage (Ohta 1973, 1992; Kimura 1983). The size of the ancestral population matters because it determines the amount of variation in the ancestral population that can then contribute to divergence between species. Larger values of  $N_{anc}$  mean that the ratio of nonsynonymous to synonymous polymorphism in the ancestral population will be smaller because proportionally fewer deleterious mutations will be segregating. All else being equal, the smaller ratio of nonsynonymous to synonymous polymorphism in the ancestral population will lead to a lower ratio of nonsynonymous to synonymous divergence. As such, the number of nonsynonymous substitutions attributed to weakly deleterious mutations is highly affected by the population sizes of the

outgroup and ancestral population, affecting estimates of  $\alpha$ . By using more realistic population sizes, the  $\alpha$  estimates we obtained for primates are similar when using chimpanzee or macaque as outgroup species. This is strong evidence that our method is more accurate, as all the other methods give drastically different estimates of  $\alpha$  using these two different outgroups. Eyre-Walker and Keightley (2009) also suggested that  $\alpha$  in humans could be as high as 0.31 if the effective population size of humans and macaques was much higher than 10,000 until very recently, foreshadowing our current estimates. The increased accuracy of the Complex model applies to the estimation of other parameters of adaptive evolution, including but not limited to  $\omega$ .

Future studies of adaptive substitutions should carefully consider ancestral and outgroup population sizes and use statistical methods that model them realistically.

Here we have quantified adaptive evolution from two different perspectives. First, we estimate the proportion of adaptive nonsynonymous substitutions between species,  $\alpha$ . This statistic measures the endpoint where a number of factors such as demography, genetic drift, and natural selection all come into play. Second, we estimate  $p^+$  and  $s^+$  of newly arising beneficial mutations. This approach aims to understand the properties of new beneficial mutations, the beginning point where beneficial mutations appear and enter the population. Thus, it is possible that these two approaches can yield qualitative opposite results. Indeed, we have shown that *D. melanogaster* has the largest values of  $\alpha$ , but the smallest  $p^+s^+$ .

We previously found that 15% of nonsynonymous mutations in humans are weakly beneficial (Huber et al. 2017). Because Huber et al. (2017) only analyzed polymorphism data, they did not consider strongly beneficial mutations that became substitutions between species. Our present study leverages genome sequences of several related species pairs to estimate the proportion of strongly beneficial ( $s^+ > 10^{-5}$ ) mutations along with the strength of selection acting on them. Weakly beneficial mutations are already accounted for in our model using the DFE from Huber et al. (2017) as they likely segregate as polymorphisms. We find that humans have a higher  $p^+s^+$  than *D. melanogaster*, mirroring the qualitative trend seen for weakly beneficial mutations in Huber et al. (2017). In the context of human evolution, it has been shown that segregating weakly beneficial mutations could result in an underestimate of  $\alpha$  (Galtier 2016; Uricchio et al. 2019). Thus,  $\alpha$  in humans could be even higher if adaptive polymorphisms are taken into account. However, our approach of modeling the complex demography in the outgroup and ancestral populations should correct for some of the bias in inferring  $\alpha$  owing to weakly beneficial mutations segregating as polymorphisms. Modeling the complex demography reduces the number of fixations from nearly neutral mutations that are treated in our model as nonadaptive, including those from segregating weakly beneficial mutations.

Our estimate of  $\alpha$  for human lineage is 16.0% using the Complex demographic model, which is comparable to the estimate of human-lineage  $\alpha$  by Uricchio et al. (2019), despite the use of different analytical approaches.  $\alpha$  is expected to be lower

on the human lineage compared with the chimpanzee or macaque lineage owing to the higher proportion of weakly deleterious amino acid substitutions on the human lineage to their smaller population size (Ohta 1973, 1992). For *D. melanogaster*, by using the Complex model, we estimate  $p^+\gamma^+$  to be  $\sim 0.1$ , which is within a factor of two of a previous estimate of Keightley et al. (2016). For mice, using the Complex model, we estimate  $p^+\gamma^+$  to be  $\sim 0.15$ , which is somewhat greater than the value of  $\sim 0.05$  inferred by Booker and Keightley (2018) and 0.05–0.06 by Campos et al. (2017). The differences in estimates could be because of the use of different methods and models, such as our detailed modeling of the ancestral and outgroup population sizes.

We make several key modeling assumptions. First, we assume that the DFE is the same between the ingroup and outgroup for each pair of species (e.g., humans and chimpanzee have the same DFE). This assumption could be relaxed by using only polymorphism data to infer the beneficial DFE (Tataru et al. 2017), although not including divergence to an outgroup species reduces power to infer positive selection (Booker 2020). Second, our inference of  $p^+$  and  $s^+$ , as well as related approaches (Schneider et al. 2011; Campos et al. 2017; Tataru et al. 2017; Booker and Keightley 2018), makes the assumption of selection starting from a single mutation, because we use the fixation probability of a mutation introduced as a single copy. Selection on standing variation (i.e., one type of soft sweep) might bias the parameter estimates, although the meaning of  $p^+$  and  $s^+$  are not as clear under models with soft sweeps. However,  $\alpha$  should be accurately estimated if the beneficial mutations reach fixation. Furthermore, as our approach only focuses on strongly beneficial mutations that contribute to divergence, it is not sensitive to types of selection in which mutations do not reach fixation. These include another type of soft sweeps (Pritchard et al. 2010), in which multiple independent beneficial mutations in the same gene are selected simultaneously but no mutation becomes fixed, as well as polygenic adaptation, in which beneficial alleles only slightly increase in frequency, without reaching fixation. Another limitation comes from the model itself. It has been suggested that models in which  $s^+$  does not change over time may not model the complexity of adaptive walks, in which the first beneficial mutation may have a greater effect on fitness than subsequent beneficial mutations (Gillespie 2004; Lourenço et al. 2013). Thus, our estimates of  $p^+$  and  $s^+$  should be interpreted as the average values over time and over genetic backgrounds, rather than literal values that have stayed constant over time. In fact, our findings show that  $p^+$  and  $s^+$  have changed over deep evolutionary time, pointing to their dynamic nature. Furthermore, as shown by the ridges on the LL surfaces (Figs. 2A–C, 5A–C), our method cannot separately infer  $p^+$  and  $s^+$  (as well as  $p^+$  and  $\gamma^+$ ). Other types of data and analyses could potentially separately estimate  $p^+$  and  $s^+$  or  $\gamma^+$ . Andolfatto (2007) and Campos et al. (2017) suggested that large values of  $s^+$  and  $\gamma^+$  with smaller values of  $p^+$  are required to generate a negative relationship between synonymous diversity and nonsynonymous divergence. Such parameter values could still be consistent with our estimates of composite parameters  $p^+s^+$  and  $p^+\gamma^+$ . In addition, there likely is a distribution of  $s^+$ , which could include a relatively smaller proportion of strongly beneficial mutations that generate the sweep effect on synonymous diversity.

If adaptation is mutation limited, more beneficial mutations are expected in organisms with larger population sizes. This view was not supported by conceptual arguments by Gillespie (2004) who suggested that the rate of environment change will matter more than the population size in determining the rate of adaptive

evolution. The simulation study by Lourenço et al. (2013) that considered a changing DFE over time in the context of a Fisher's geometric model found that the population size was only weakly related to  $\alpha$ . Instead, the rate at which the environment changed was an important predictor of the amount of adaptive evolution, as environmental shifts moved the population from the fitness optimum, creating the opportunity for new beneficial mutations. Further, Rousselle et al. (2020) recently found a weak negative correlation between the amount of adaptive evolution and the amount of genetic diversity in modern populations, which supports the models of Lourenço et al. (2013) and Gillespie (2004). However, not all studies are in agreement on the role of the environment as Connallon and Clark (2015) found that environmental heterogeneity reduces the fraction of beneficial mutations by inflating the standardized mutation size in a Fisher's geometric model. Lourenço et al. (2013) also found that organismal complexity, here defined as the number of phenotypes under selection, was a key predictor of the amount of adaptive evolution within species. Through a "cost of complexity," more complex organisms have a harder time adapting to new environmental conditions because of the additional constraints imposed by the increased number of traits under selection. As such, adaptive walks require more beneficial mutations (Orr 1998).

Our results presented here are in broad agreement with the conceptual model of Lourenço et al. (2013). Specifically, we do not find that species with larger population sizes (i.e., *D. melanogaster*) have more beneficial mutations. Instead, we find that  $p^+s^+$  is higher in humans than in *D. melanogaster* or mice. Second, although it is hard to precisely define organismal complexity, previous work has found more protein–protein interactions in humans than in *Drosophila* (Valentine et al. 1994; Stumpf et al. 2008), suggesting that humans may be more complex than *Drosophila*. If this is the case, then our findings of a higher  $p^+s^+$  in humans than *D. melanogaster* and mice support the arguments from Lourenço et al. (2013) that adaptive walks after an environmental shift are less efficient and require more steps (i.e., beneficial mutations) in more complex organisms, leading to a higher  $p^+s^+$  in complex organisms. Additionally, differences in the degree of environmental change across species could also contribute to the disparate inferences of  $p^+s^+$  across taxa. Although it is hard to say which species have experienced more environmental shifts, species with longer generation times, like primates, may experience more environmental change per generation than species with shorter generation times, like *Drosophila* or rodents (Romiguier et al. 2014; Rousselle et al. 2020). Another possibility is that the selection coefficient per year is similar across species (Charlesworth 1994; Chevin 2011). This could occur if the selection coefficient is related to interactions with the environment on a per-year timescale. Under such a model, the per-generation selection coefficient, as we measure here, would be larger in species with longer generation times. Our finding of a higher  $p^+s^+$  in primates than in *Drosophila* is consistent with these predictions, although it is not possible from our analyses to conclusively favor any particular explanation.

## Methods

### Polymorphism and divergence data for humans, mice, and *D. melanogaster*

For humans, we used polymorphism data from 112 individuals from Yoruba in Ibadan, Nigeria (YRI) from The 1000 Genomes



Project (The 1000 Genomes Project Consortium 2012). Published genome alignments of human and chimpanzee (hg19/panTro4), and human and *Macaca mulatta* (hg19/rheMac3) were downloaded from the UCSC Genome Browser (<http://hgdownload.soe.ucsc.edu/goldenPath/hg19/>). As coding regions in hg19 are well annotated, use of more recent genome builds would not affect our conclusions. For *D. melanogaster*, we used the *Drosophila* Population Genomics Project phase 3 data, including 197 African *D. melanogaster* lines from Zambia, Africa (Lack et al. 2015). For divergence, *D. melanogaster* and *D. simulans* genic alignments (Dmel v5/Dsim v2) were extracted from the multispecies alignments from Hu et al. (2013). Only autosomal regions were used in our analysis. Human and *D. melanogaster* polymorphism data were filtered and down-sampled to 100 chromosomes as described by Huber et al. (2017).

For mice, raw data were downloaded for 10 *Mus musculus castaneus* individuals that were collected in the northwest Indian state of Himachal Pradesh (Halligan et al. 2010, 2013). Reads were mapped against mouse genome mm9 using BWA (Li and Durbin 2009) and Stampy (Lunter and Goodson 2011). As coding regions in mm9 are well annotated, use of more recent genome builds would not affect our conclusions. Duplicate reads were marked using Picard, and further preprocessing was performed following GATK best practice guidelines (McKenna et al. 2010). Variants were called using the GATK UnifiedGenotyper and filtered using the GATK VQSR using Affymetrix Mouse Diversity Genotyping Array sites (Yang et al. 2009). We further filtered the data set to only retain sites with a sample size of at least 16 chromosomes and down-sampled all sites to a sample size of 16 chromosomes using the hypergeometric probability distribution. Published genome alignments of mice and rat (mm9/rn5) were downloaded from UCSC (<http://hgdownload.soe.ucsc.edu/goldenPath/mm9/vsRn5/axtNet/>). For each species, polymorphism and divergence data were intersected, and only coding regions shared by both data sets were used in our analysis.

In total, 19.1 Mb of coding sequences for primates, 26.6 Mb of coding sequences for rodents, and 15.8 Mb of coding sequences for *Drosophila* were included. The nonsynonymous and synonymous total sequence lengths ( $L_{NS}$ ,  $L_S$ ) were estimated using multipliers of  $L_{NS} = 2.85 \times L_S$  in *D. melanogaster* and  $L_{NS} = 2.31 \times L_S$  in mammals from Huber et al. (2017). Human variants were annotated using the SeattleSeq Annotation pipeline (<http://snp.gs.washington.edu/SeattleSeqAnnotation138/>). Mice and *Drosophila* variants were annotated using SnpEff v3.6 using the mice NCBIM37.66 annotation database and the *D. melanogaster* BDGP5.75 annotation database, respectively. Sites that are annotated as near-splice, or loss of function, were removed. The ratio of nonsynonymous/synonymous differences between human and chimp sequences in our data set is approximately 0.65, which is consistent with several previous reports from different data sets (Bustamante et al. 2005; Torgerson et al. 2009; Enard et al. 2014).

### Model-based estimates of $\alpha$

Model-based estimates of  $\alpha$  require a demographic model and a DFE for neutral and deleterious mutations to predict the  $D_{NE}$  that is accounted for by such mutations. For primates and *Drosophila*, we used demographic and DFE parameters from Huber et al. (2017; Supplemental Table S3). For rodents, we conducted our own inference of these parameters by summarizing the polymorphism data in mice by the folded SFS (Supplemental Text).

To compute the expected divergence between species, we computed the divergence accumulated in the ingroup population, outgroup population, and ancestral population separately and

summed them. The divergence in the ingroup and outgroup population is a function of divergence time, effective population size, mutation rate, and selection coefficient and was calculated according to equation 13 in Sawyer and Hartl (1992). We computed the expected divergence under a gamma DFE model by Monte Carlo integration using 1 million gamma-distributed selection coefficients. The contribution of the ancestral population to divergence between two species was computed numerically based on the diffusion equation, using *prfreq* (Boyko et al. 2008) and assuming the same gamma DFE. These calculations are implemented in our program *predicDiv* (Supplemental Code). For both the Simple and the Complex models, we first estimate the divergence times (Supplemental Table S5) that fit the observed number of synonymous differences between a pair of species because there is a wide range of divergence times from the literature for each species. Here the number of synonymous differences equals  $2 \times$  divergence time  $\times$  mutation rate. Second, using this divergence time, demography, and DFE inferred from Huber et al. (2017), or as described above for mice, we estimated the expected number nonsynonymous differences ( $D_{NE}$ ) according to equation 13 of Sawyer and Hartl (1992; Boyko et al. 2008). Then,  $\alpha$  and omega-a are calculated as

$$\alpha = \frac{D_{NO} - D_{NE}}{D_{NO}}, \quad (2)$$

$$\omega_a = \frac{D_{NO} - D_{NE}}{D_S} \left( \frac{L_{NS}}{L_S} \right), \quad (3)$$

where  $D_{NO}$  is the observed number of nonsynonymous differences between species, and  $D_S$  is the observed number of synonymous differences between species. The 95% CIs of  $\alpha$  were calculated by parametric bootstrapping through resampling 10,000 draws of  $D_N$  using a Poisson distribution with mean of  $D_{NO}$ .

### DFE-alpha

Data files and the program v2.15 were downloaded from the following link: [http://www.homepages.ed.ac.uk/pkeightl/dfe\\_alpha/download-dfe-alpha.html](http://www.homepages.ed.ac.uk/pkeightl/dfe_alpha/download-dfe-alpha.html). Folded synonymous and nonsynonymous SFSs were used as input in the inferences. The *est\_alpha\_omega* program was used to estimate the proportion of adaptive divergence.

### Composite likelihood approach for testing whether $p^+$ and $s^+$ differ across species

We first used *predicDiv* and *prfreq* to generate a look-up table for the expected number of nonsynonymous differences between species for a range of  $s^+$  for each species. We focused on this range to capture strongly advantageous mutations. Given a demographic model, we assume a substitution rate according to equation 13 in Sawyer and Hartl (1992). The contribution of the ancestral population to the divergence between two species is computed by numerically solving the diffusion equation using *predicDiv* and *prfreq* (Boyko et al. 2008). For each species, we then searched a grid of values of  $\log_{10}(s^+)$  ( $-5$  to  $-2$ ) and  $p^+$  (0%–7.5%). We are interested in this range of strong  $s^+$  because weakly beneficial mutations still segregating as polymorphisms and are accounted for by the DFE fit to the SFS. We use a Poisson LL function to calculate the LL for each combination of  $s^+$  and  $p^+$ . We find the MLE of  $s^+$  and  $p^+$  for each taxon under each demographic model that maximizes the LL and best fit the observed  $D_N$ . This is the full model (H1) in which each taxon (i.e., primates, rodent, and *Drosophila*) is allowed to have its own  $s^+$  and  $p^+$ . In the constrained model (H0), two or three taxa have the same  $s^+$  and  $p^+$ . The LL of the constrained model is the sum of LL for each  $s^+$  and  $p^+$  for each taxon

under comparison. We then found the MLE for the constrained model and the likelihood ratio between H1 and H0.

### Composite likelihood approach for testing whether $p^+$ and $\gamma^+$ differ across species

Similarly, we used *predicDiv* and *prfreq* (Boyko et al. 2008) to generate a look-up table for the expected number of nonsynonymous differences for a range of  $\gamma^+$  for each taxon under each demographic model. For each species, we performed a grid search of  $\log_{10}(\gamma^+)$  (0–3) and  $p^+$  (0%–7.5%). Note that because effective population sizes differ over several orders of magnitude across our three taxa, we are searching across drastically different ranges of  $s^+$  compared with our previous inference described above. We again use a Poisson LL function to calculate the LL for each combination of  $\gamma^+$  and  $p^+$  and use a similar LRT as described above for  $s^+$  and  $p^+$ .

### Composite parameters $p^+s^+$ and $p^+\gamma^+$

We examined two composite parameters, the product of  $p^+$  and  $s^+$  and the product of  $p^+$  and  $\gamma^+$  for each taxon. Multiple combinations of parameters could give the same product  $p^+s^+$  and  $p^+\gamma^+$  and distinct combinations could have different LL. Thus, we found the values of  $p^+$  and  $s^+$  that gave the highest LL for each  $p^+s^+$  value (and similarly each  $p^+$  and  $\gamma^+$  for each  $p^+\gamma^+$ ). Values of  $\log_{10}(p^+s^+)$  and  $p^+\gamma^+$  were rounded to three digits before this comparison.

### Conditional LLs

To examine whether primates, rodents, and *Drosophila* could have the same  $s^+$  or  $p^+$ , we examined the conditional LLs. To make the conditional LL curve for  $p^+$ , for each  $s^+$  value, we find the  $p^+$  that maximizes the LL, as well as the  $p^+$  values that have a LL within three LL units of this maximum LL for this  $s^+$  (i.e.,  $p^+|s^+$ ). To make the conditional LL curve for  $s^+$ , for each  $p^+$  value, we look for the  $s^+$  that maximizes the LL as well as the  $s^+$  values that have a LL within three LL units of this maximum LL for this  $p^+$  (i.e.,  $s^+|p^+$ ). The same approach was used to construct the conditional LL curve for  $p^+$  and  $\gamma^+$ .

### Filtering to only include sites not affected by BGC or CpG hypermutation

We repeated our analyses only including SSWW mutations because these changes were not affected by BGC or increased mutation rates owing to deamination of methylated CpG sites. To do this, we adjusted the mutation rates to include only SSWW sites (Supplemental Text). We also reinferred the demographic parameters and DFE using only SSWW mutations in rodents (Supplemental Text) because the SFS for SSWW mutations differed from that of the full data.

### Software availability

*predicDiv* is available as Supplemental Code and also at GitHub (<https://github.com/LohmuellerLab/predicDiv>).

### Competing interest statement

The authors declare no competing interests.

### Acknowledgments

We thank Lawrence Uricchio and David Enard for advice and Tanya Phung, Jazlyn Mooney, Clare Marsden, and Jesse Garcia

for helpful comments on our manuscript. This work was supported by a Searle Scholars Program fellowship (20143955) and National Institutes of Health (National Institute of General Medical Sciences) grant R35GM119856 (to K.E.L.) and foundation of Westlake University and National Natural Science Foundation of China (NSFC) 31900315 (to Y.Z.). We acknowledge support from a UCLA Institute for Quantitative and Computational Biosciences (QCB) postdoctoral fellowship to Y.Z. and the QCB Collaboratory community directed by Matteo Pellegrini.

**Author contributions:** K.E.L. conceived of and supervised the study. Y.Z. analyzed positive selection. C.D.H. inferred demography and gamma-DFEs. R.W.D. processed mice raw data to genotypes. Y.Z., C.D.H., R.W.D., and K.E.L. participated in manuscript preparation.

### References

- The 1000 Genomes Project Consortium. 2012. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**: 56–65. doi:10.1038/nature11632
- Andolfatto P. 2005. Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* **437**: 1149–1152. doi:10.1038/nature04107
- Andolfatto P. 2007. Hitchhiking effects of recurrent beneficial amino acid substitutions in the *Drosophila melanogaster* genome. *Genome Res* **17**: 1755–1762. doi:10.1101/gr.6691007
- Andolfatto P, Wong KM, Bachtrog D. 2011. Effective population size and the efficacy of selection on the X chromosomes of two closely related *Drosophila* species. *Genome Biol Evol* **3**: 114–128. doi:10.1093/gbe/evq086
- Bolívar P, Mugal CF, Nater A, Ellegren H. 2016. Recombination rate variation modulates gene sequence evolution mainly via GC-biased gene conversion, not Hill–Robertson interference, in an avian system. *Mol Biol Evol* **35**: 216–227. doi:10.1093/molbev/msv214
- Booker TR. 2020. Inferring parameters of the distribution of fitness effects of new mutations when beneficial mutations are strongly advantageous and rare. *G3 Genes Genomes Genet* **10**: 2317–2326. doi:10.1534/g3.120.401052
- Booker TR, Keightley PD. 2018. Understanding the factors that shape patterns of nucleotide diversity in the house mouse genome. *Mol Biol Evol* **35**: 2971–2988. doi:10.1093/molbev/msy188
- Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR, et al. 2008. Assessing the evolutionary impact of amino acid mutations in the human genome. *PLoS Genet* **4**: e1000083. doi:10.1371/journal.pgen.1000083
- Burgess R, Yang Z. 2008. Estimation of hominoid ancestral population sizes under Bayesian coalescent models incorporating mutation rate variation and sequencing errors. *Mol Biol Evol* **25**: 1979–1994. doi:10.1093/molbev/msn148
- Bustamante CD, Fedel-Alon A, Williamson S, Nielsen R, Hubisz MT, Gnanowski S, Tanenbaum DM, White TJ, Sninsky JJ, Hernandez RD, et al. 2005. Natural selection on protein-coding genes in the human genome. *Nature* **437**: 1153–1157. doi:10.1038/nature04240
- Campos JL, Zhao L, Charlesworth B. 2017. Estimating the parameters of background selection and selective sweeps in *Drosophila* in the presence of gene conversion. *Proc Natl Acad Sci* **114**: E4762–E4771. doi:10.1073/pnas.1619434114
- Charlesworth B. 1994. *Evolution in age-structured populations*. Cambridge University Press, Cambridge.
- Charlesworth J, Eyre-Walker A. 2008. The McDonald–Kreitman test and slightly deleterious mutations. *Mol Biol Evol* **25**: 1007–1015. doi:10.1093/molbev/msn005
- Chen F-C, Li W-H. 2001. Genomic divergences between humans and other hominoids and the effective population size of the common ancestor of humans and chimpanzees. *Am J Hum Genet* **68**: 444–456. doi:10.1086/318206
- Chevin L-M. 2011. On measuring selection in experimental evolution. *Biol Lett* **7**: 210–213. doi:10.1098/rsbl.2010.0580
- Connallon T, Clark AG. 2015. The distribution of fitness effects in an uncertain world. *Evol Int J Org Evol* **69**: 1610–1618. doi:10.1111/evo.12673
- Corbett-Detig RB, Hartl DL, Sackton TB. 2015. Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol* **13**: e1002112. doi:10.1371/journal.pbio.1002112
- Duncan BK, Miller JH. 1980. Mutagenic deamination of cytosine residues in DNA. *Nature* **287**: 560–561. doi:10.1038/287560a0

- Duret L, Galtier N. 2009. Biased gene conversion and the evolution of mammalian genomic landscapes. *Annu Rev Genomics Hum Genet* **10**: 285–311. doi:10.1146/annurev-genom-082908-150001
- 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 Genet* **12**: e1006130. doi:10.1371/journal.pgen.1006130
- Enard D, Messer PW, Petrov DA. 2014. Genome-wide signals of positive selection in human evolution. *Genome Res* **24**: 885–895. doi:10.1101/gr.164822.113
- Eyre-Walker A. 2002. Changing effective population size and the McDonald–Kreitman test. *Genetics* **162**: 2017–2024.
- Eyre-Walker A, Keightley PD. 2009. Estimating the rate of adaptive molecular evolution in the presence of slightly deleterious mutations and population size change. *Mol Biol Evol* **26**: 2097–2108. doi:10.1093/molbev/msp119
- Fay JC. 2011. Weighing the evidence for adaptation at the molecular level. *Trends Genet* **27**: 343–349. doi:10.1016/j.tig.2011.06.003
- Fay JC, Wyckoff GJ, Wu CI. 2001. Positive and negative selection on the human genome. *Genetics* **158**: 1227–1234.
- Fay JC, Wyckoff GJ, Wu C-I. 2002. Testing the neutral theory of molecular evolution with genomic data from *Drosophila*. *Nature* **415**: 1024–1026. doi:10.1038/4151024a
- Fischer A, Wiebe V, Pääbo S, Przeworski M. 2004. Evidence for a complex demographic history of chimpanzees. *Mol Biol Evol* **21**: 799–808. doi:10.1093/molbev/msh083
- Foxe JP, Dar V-N, Zheng H, Nordborg M, Gaut BS, Wright SI. 2008. Selection on amino acid substitutions in *Arabidopsis*. *Mol Biol Evol* **25**: 1375–1383. doi:10.1093/molbev/msn079
- Galtier N. 2016. Adaptive protein evolution in animals and the effective population size hypothesis. *PLoS Genet* **12**: e1005774. doi:10.1371/journal.pgen.1005774
- Gillespie JH. 2004. Why  $k = 4N\mu$  is silly. In *The evolution of population biology* (ed. Singh RS, Uyenoyama MK), pp. 178–192. Cambridge University Press, Cambridge.
- Gossmann TI, Song B-H, Windsor AJ, Mitchell-Olds T, Dixon CJ, Kapralov MV, Filatov DA, Eyre-Walker A. 2010. Genome wide analyses reveal little evidence for adaptive evolution in many plant species. *Mol Biol Evol* **27**: 1822–1832. doi:10.1093/molbev/msq079
- Gossmann TI, Keightley PD, Eyre-Walker A. 2012. The effect of variation in the effective population size on the rate of adaptive molecular evolution in eukaryotes. *Genome Biol Evol* **4**: 658–667. doi:10.1093/gbe/evs027
- Haddrill PR, Loewe L, Charlesworth B. 2010. Estimating the parameters of selection on nonsynonymous mutations in *Drosophila pseudoobscura* and *D. miranda*. *Genetics* **185**: 1381–1396. doi:10.1534/genetics.110.117614
- Halligan DL, Oliver F, Eyre-Walker A, Harr B, Keightley PD. 2010. Evidence for pervasive adaptive protein evolution in wild mice. *PLoS Genet* **6**: e1000825. doi:10.1371/journal.pgen.1000825
- Halligan DL, Kousathanas A, Ness RW, Harr B, Eöry L, Keane TM, Adams DJ, Keightley PD. 2013. Contributions of protein-coding and regulatory change to adaptive molecular evolution in murid rodents. *PLoS Genet* **9**: e1003995. doi:10.1371/journal.pgen.1003995
- Hernandez RD, Hubisz MJ, Wheeler DA, Smith DG, Ferguson B, Rogers J, Nazareth L, Indap A, Bourquin T, McPherson J, et al. 2007. Demographic histories and patterns of linkage disequilibrium in Chinese and Indian rhesus macaques. *Science* **316**: 240–243. doi:10.1126/science.1140462
- Hobolth A, Christensen OF, Mailund T, Schierup MH. 2007. Genomic relationships and speciation times of human, chimpanzee, and gorilla inferred from a coalescent hidden Markov model. *PLoS Genet* **3**: e7. doi:10.1371/journal.pgen.0030007
- 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 Res* **23**: 89–98. doi:10.1101/gr.141689.112
- Huber CD, Kim BY, Marsden CD, Lohmueller KE. 2017. Determining the factors driving selective effects of new nonsynonymous mutations. *Proc Natl Acad Sci* **114**: 4465–4470. doi:10.1073/pnas.1619508114
- Hvilsom C, Qian Y, Bataillon T, Li Y, Mailund T, Sallé B, Carlsen F, Li R, Zheng H, Jiang T, et al. 2012. Extensive X-linked adaptive evolution in central chimpanzees. *Proc Natl Acad Sci* **109**: 2054–2059. doi:10.1073/pnas.1106877109
- Jackson BC, Campos JL, Haddrill PR, Charlesworth B, Zeng K. 2017. Variation in the intensity of selection on codon bias over time causes contrasting patterns of base composition evolution in *Drosophila*. *Genome Biol Evol* **9**: 102–123. doi:10.1093/gbe/evw291
- Keightley PD, Campos JL, Booker TR, Charlesworth B. 2016. Inferring the frequency spectrum of derived variants to quantify adaptive molecular evolution in protein-coding genes of *Drosophila melanogaster*. *Genetics* **203**: 975–984. doi:10.1534/genetics.116.188102
- Kimura M. 1962. On the probability of fixation of mutant genes in a population. *Genetics* **47**: 713–719.
- Kimura M. 1983. *The neutral theory of molecular evolution*. Cambridge University Press, Cambridge.
- Lachance J, Tishkoff SA. 2014. Biased gene conversion skews allele frequencies in human populations, increasing the disease burden of recessive alleles. *Am J Hum Genet* **95**: 408–420. doi:10.1016/j.ajhg.2014.09.008
- Lack JB, Cardeno CM, Crepeau MW, Taylor W, Corbett-Detig RB, Stevens KA, Langley CH, Pool JE. 2015. The *Drosophila* genome nexus: a population genomic resource of 623 *Drosophila melanogaster* genomes, including 197 from a single ancestral range population. *Genetics* **199**: 1229–1241. doi:10.1534/genetics.115.174664
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinforma Oxf Engl* **25**: 1754–1760. doi:10.1093/bioinformatics/btp324
- Loewe L, Charlesworth B, Bartolomé C, Nöel V. 2006. Estimating selection on nonsynonymous mutations. *Genetics* **172**: 1079–1092. doi:10.1534/genetics.105.047217
- Lourenço JM, Glémin S, Galtier N. 2013. The rate of molecular adaptation in a changing environment. *Mol Biol Evol* **30**: 1292–1301. doi:10.1093/molbev/mst026
- Lunter G, Goodson M. 2011. Stampy: a statistical algorithm for sensitive and fast mapping of Illumina sequence reads. *Genome Res* **21**: 936–939. doi:10.1101/gr.111120.110
- McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the *Adh* locus in *Drosophila*. *Nature* **351**: 652–654. doi:10.1038/351652a0
- McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K, Altshuler D, Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297–1303. doi:10.1101/gr.107524.110
- McVicker G, Gordon D, Davis C, Green P. 2009. Widespread genomic signatures of natural selection in hominid evolution. *PLoS Genet* **5**: e1000471. doi:10.1371/journal.pgen.1000471
- Messer PW, Petrov DA. 2013. Frequent adaptation and the McDonald–Kreitman test. *Proc Natl Acad Sci* **110**: 8615–8620. doi:10.1073/pnas.1220835110
- Nam K, Munch K, Mailund T, Nater A, Greminger MP, Krützen M, Marquès-Bonet T, Schierup MH. 2017. Evidence that the rate of strong selective sweeps increases with population size in the great apes. *Proc Natl Acad Sci* **114**: 1613–1618. doi:10.1073/pnas.1605660114
- Ness RW, Zhang Y-H, Cong L, Wang Y, Zhang J-X, Keightley PD. 2012. Nuclear gene variation in wild brown rats. *G3 Genes Genomes Genet* **2**: 1661–1664. doi:10.1534/g3.112.004713
- Ohta T. 1973. Slightly deleterious mutant substitutions in evolution. *Nature* **246**: 96–98. doi:10.1038/246096a0
- Ohta T. 1992. The nearly neutral theory of molecular evolution. *Annu Rev Ecol Syst* **23**: 263–286. doi:10.1146/annurev.es.23.110192.001403
- Orr HA. 1998. The population genetics of adaptation: the distribution of factors fixed during adaptive evolution. *Evolution (N Y)* **52**: 935–949. doi:10.1111/j.1558-5646.1998.tb01823.x
- Phifer-Rixey M, Bonhomme F, Boursot P, Churchill GA, Piálek J, Tucker PK, Nachman MW. 2012. Adaptive evolution and effective population size in wild house mice. *Mol Biol Evol* **29**: 2949–2955. doi:10.1093/molbev/mss105
- Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, Veeramah KR, Woerner AE, O'Connor TD, Santpere G, et al. 2013. Great ape genetic diversity and population history. *Nature* **499**: 471–475. doi:10.1038/nature12228
- Pritchard JK, Pickrell JK, Coop G. 2010. The genetics of human adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Curr Biol* **20**: R208–R215. doi:10.1016/j.cub.2009.11.055
- Robinson MC, Stone EA, Singh ND. 2014. Population genomic analysis reveals no evidence for GC-biased gene conversion in *Drosophila melanogaster*. *Mol Biol Evol* **31**: 425–433. doi:10.1093/molbev/mst220
- Romiguier J, Gayral P, Ballenghien M, Bernard A, Cahais V, Chenuil A, Chiari Y, Deraat R, Duret L, Favre N, et al. 2014. Comparative population genomics in animals uncovers the determinants of genetic diversity. *Nature* **515**: 261–263. doi:10.1038/nature13685
- Rousselle M, Mollion M, Nabholz B, Bataillon T, Galtier N. 2018. Overestimation of the adaptive substitution rate in fluctuating populations. *Biol Lett* **14**: 20180055. doi:10.1098/rsbl.2018.0055
- Rousselle M, Simion P, Tilak M-K, Figueu E, Nabholz B, Galtier N. 2020. Is adaptation limited by mutation? A timescale-dependent effect of genetic diversity on the adaptive substitution rate in animals. *PLoS Genet* **16**: e1008668. doi:10.1371/journal.pgen.1008668
- Sawyer SA, Hartl DL. 1992. Population genetics of polymorphism and divergence. *Genetics* **132**: 1161–1176.
- Schneider A, Charlesworth B, Eyre-Walker A, Keightley PD. 2011. A method for inferring the rate of occurrence and fitness effects of advantageous mutations. *Genetics* **189**: 1427–1437. doi:10.1534/genetics.111.131730

- Schrager CG. 2014. The effective population sizes of the anthropoid ancestors of the human–chimpanzee lineage provide insights on the historical biogeography of the great apes. *Mol Biol Evol* **31**: 37–47. doi:10.1093/molbev/mst191
- Sella G, Petrov DA, Przeworski M, Andolfatto P. 2009. Pervasive natural selection in the *Drosophila* genome? *PLoS Genet* **5**: e1000495. doi:10.1371/journal.pgen.1000495
- Smith NGC, Eyre-Walker A. 2002. Adaptive protein evolution in *Drosophila*. *Nature* **415**: 1022–1024. doi:10.1038/4151022a
- Strasburg JL, Kane NC, Raduski AR, Bonin A, Michelmore R, Rieseberg LH. 2011. Effective population size is positively correlated with levels of adaptive divergence among annual sunflowers. *Mol Biol Evol* **28**: 1569–1580. doi:10.1093/molbev/msq270
- Stumpf MPH, Thorne T, de Silva E, Stewart R, An HJ, Lappe M, Wiuf C. 2008. Estimating the size of the human interactome. *Proc Natl Acad Sci* **105**: 6959–6964. doi:10.1073/pnas.0708078105
- Sved J, Bird A. 1990. The expected equilibrium of the CpG dinucleotide in vertebrate genomes under a mutation model. *Proc Natl Acad Sci* **87**: 4692–4696. doi:10.1073/pnas.87.12.4692
- Tataru P, Mollion M, Glémin S, Bataillon T. 2017. Inference of distribution of fitness effects and proportion of adaptive substitutions from polymorphism data. *Genetics* **207**: 1103–1119. doi:10.1534/genetics.117.300323
- Torgerson DG, Boyko AR, Hernandez RD, Indap A, Hu X, White TJ, Sninsky JJ, Cargill M, Adams MD, Bustamante CD, et al. 2009. Evolutionary processes acting on candidate cis-regulatory regions in humans inferred from patterns of polymorphism and divergence. *PLoS Genet* **5**: e1000592. doi:10.1371/journal.pgen.1000592
- Uricchio LH, Petrov DA, Enard D. 2019. Exploiting selection at linked sites to infer the rate and strength of adaptation. *Nat Ecol Evol* **3**: 977–984. doi:10.1038/s41559-019-0890-6
- Valentine JW, Collins AG, Meyer CP. 1994. Morphological complexity increase in metazoans. *Paleobiology* **20**: 131–142. doi:10.1017/S0094837300012641
- Xue C, Raveendran M, Harris RA, Fawcett GL, Liu X, White S, Dahdouli M, Rio Deiros D, Below JE, Salerno W, et al. 2016. The population genomics of rhesus macaques (*Macaca mulatta*) based on whole-genome sequences. *Genome Res* **26**: 1651–1662. doi:10.1101/gr.204255.116
- Yang H, Ding Y, Hutchins LN, Szatkiewicz J, Bell TA, Paigen BJ, Graber JH, de Villena FP-M, Churchill GA. 2009. A customized and versatile high-density genotyping array for the mouse. *Nat Methods* **6**: 663–666. doi:10.1038/nmeth.1359

Received August 31, 2019; accepted in revised form November 10, 2020.