

RESEARCH ARTICLE

# Mutation Rate Variation is a Primary Determinant of the Distribution of Allele Frequencies in Humans

Arbel Harpak<sup>1</sup>\*, Anand Bhaskar<sup>2,3</sup>, Jonathan K. Pritchard<sup>1,2,3</sup>

**1** Department of Biology, Stanford University, Stanford, California, United States of America, **2** Department of Genetics, Stanford University, Stanford, California, United States of America, **3** Howard Hughes Medical Institute, Stanford University, Stanford, California, United States of America

\* These authors contributed equally to this work.

\* [arbelh@stanford.edu](mailto:arbelh@stanford.edu)



 OPEN ACCESS

**Citation:** Harpak A, Bhaskar A, Pritchard JK (2016) Mutation Rate Variation is a Primary Determinant of the Distribution of Allele Frequencies in Humans. *PLoS Genet* 12(12): e1006489. doi:10.1371/journal.pgen.1006489

**Editor:** Adam Eyre-Walker, University of Sussex, UNITED KINGDOM

**Received:** April 12, 2016

**Accepted:** November 16, 2016

**Published:** December 15, 2016

**Copyright:** © 2016 Harpak et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All human polymorphism data are available from the ExAC database ([http://ftp.broadinstitute.org/pub/ExAC\\_release/release0.2/](http://ftp.broadinstitute.org/pub/ExAC_release/release0.2/)). All primate divergence files are available from the UCSC Genome browser database (<https://genome.ucsc.edu/>). All human mutation rate measurements from Kong et al. 2012 are available from <http://www.nature.com/nature/journal/v488/n7412/extref/nature11396-s2.xls>. All chromHMM annotations are available from the Encode database (<https://genome.ucsc.edu/ENCODE/>).

## Abstract

The site frequency spectrum (SFS) has long been used to study demographic history and natural selection. Here, we extend this summary by examining the SFS conditional on the alleles found at the same site in other species. We refer to this extension as the “phylogenetically-conditioned SFS” or cSFS. Using recent large-sample data from the Exome Aggregation Consortium (ExAC), combined with primate genome sequences, we find that human variants that occurred independently in closely related primate lineages are at higher frequencies in humans than variants with parallel substitutions in more distant primates. We show that this effect is largely due to sites with elevated mutation rates causing significant departures from the widely-used infinite sites mutation model. Our analysis also suggests substantial variation in mutation rates even among mutations involving the same nucleotide changes. In summary, we show that variable mutation rates are key determinants of the SFS in humans.

## Author Summary

The site frequency spectrum (SFS, i.e., the distribution of allele frequencies) is a summary of natural variation, used to study demographic history and natural selection. Here, we extended the SFS by conditioning on phylogenetic divergence patterns. We refer to this extension as the “phylogenetically-conditioned SFS” or cSFS. Exploring the determinants of the cSFS revealed two main findings. First, we found that mutations that have been independently fixed in another species are more likely to be benign for contemporary humans if the fixation occurred in a closely related species. The background on which a mutation occurs is therefore an important feature to consider when predicting the fitness consequences of a mutation. Second, we found that the SFS is substantially affected by repeat mutations within the human population. The extent of repeat mutations implies that some sites must have particularly high mutation rates, beyond the known variation across the different possible nucleotide changes. Our observations contradict the “infinite sites” mutation model, which is commonly used in population genetic analyses, and imply

**Funding:** This work was supported by grants R01MH084703 and U01HG007036 from the Howard Hughes Medical Institute (HHMI) to JKP. AB was supported in part by a postdoctoral fellowship from the Stanford Center for Computational, Evolutionary and Human Genomics (CEHG). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing Interests:** The authors have declared that no competing interests exist.

that future SFS-based analyses of human populations should account for mutation rate variation.

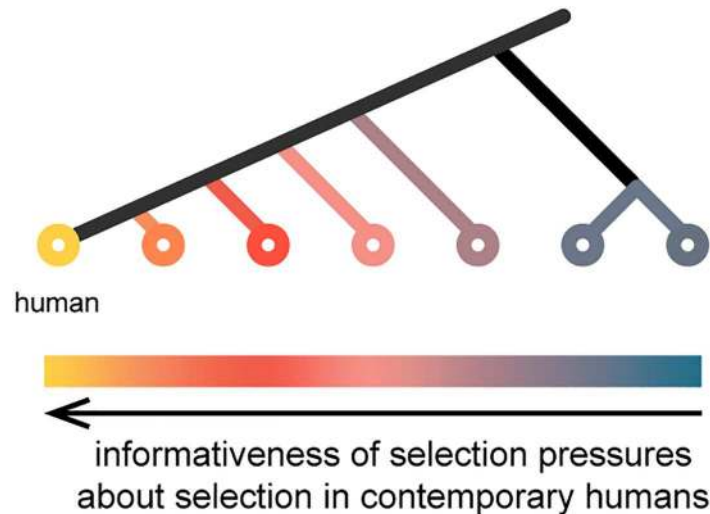
## Introduction

The distribution of allele frequencies across segregating sites, commonly referred to as the Site Frequency Spectrum (SFS), is a central focus of population genetics research as it can reflect a wide range of evolutionary processes, including demographic history as well as positive and purifying selection [1–8]. Until recently, the SFS was usually measured in samples of tens or hundreds of people, but advances in sequencing technology have enabled the collection of sequence data at much larger scales [9–14]. Notably, the Exome Aggregation Consortium (ExAC) recently released high quality, exome-wide allele counts for over 60,000 people [12].

Large sample sizes are valuable because they make it possible to detect many more segregating sites, and to estimate the frequencies of rare variants. For example, the recent dramatic expansion of human populations leaves little signal in the SFS in small samples [15], but is readily detected in large samples, where there is a huge excess of low frequency variants compared to model-predictions without growth [13,14,16,17]. Similarly, large samples enable the detection of deleterious variants that are held at very low frequencies by purifying selection [18–22].

In this paper, we extend the SFS by considering the SFS *conditional* on the observed alleles at a given site in other species (specifically, other primates in our analysis). Our original motivation was that this could allow us to measure the effects of sequence context on the selective constraint of missense variants. In general, sites with strong levels of average constraint across mammals tend to be less polymorphic within humans [16,23,24] but, to the best of our knowledge, there has not been extensive consideration of the joint distribution of the substitutions across other lineages and the human SFS. In particular, we hypothesized that if an identical substitution has occurred independently in a closely related species—e.g., in a great ape—then this is strong evidence that the same variant is unlikely to be deleterious in humans. However, an identical substitution in a more distantly related species may be much less informative, as substitutions at other positions within the same gene may change the set of preferred alleles due to epistatic interactions [25–31] (Fig 1). For example, it has been shown that, in a handful of cases, likely disease-causing variants in humans are actually wildtype alleles in mouse, presumably rendered harmless by parallel substitutions at interacting positions [26].

As we show below, the human SFS varies greatly depending on patterns of substitutions in other species. In part, this does appear to be due to differences in fitness effects; however a more important factor is mutation rate variation across sites. Under the widely-used infinite sites model, the SFS is independent of mutation rate; but in the ExAC dataset we observe a clear breakdown of this model. Mutation rates are known to vary across sites due to a variety of different mechanisms, leading to differences between CpGs, transitions and transversions, as well as additional effects that correlate with broader sequence context, replication timing, transcription, recombination rate and chromatin environment [32–39]. We show here that mutation rates are much more variable than generally appreciated, and that rates at some sites are high enough to generate substantial deviations from infinite sites predictions. The main ExAC paper [12] also recently reported that the SFS varies substantially across mutation types, and also noted that this implies departures from the infinite sites model, especially for CpG transitions.



**Fig 1. Hypothesis about information in parallel mutations.** If an identical substitution occurred independently in a closely related species, then the variant is unlikely to be deleterious in humans. An identical substitution in a more distantly related species may be less informative because sequence divergence at interacting sites may change the set of preferred alleles, and hence the selective constraint at the site.

doi:10.1371/journal.pgen.1006489.g001

In summary, our results suggest more variation in mutation rates across sites than is generally appreciated, and further that the infinite sites model provides a poor fit for population genetic analyses in large modern data sets.

## Results

To investigate the properties of the human cSFS, we combined exome sequence data from 60,706 humans from ExAC version 0.2 [12,40] and orthologous reference alleles for 6 nonhuman primate species from the UCSC genome browser [41]. After applying several filters (see [Materials and methods](#) for details) we were left with 6,002,065 single nucleotide polymorphisms (SNPs) for which we had orthologous data in at least one nonhuman species.

We examined how the human SFS changes as we condition on various divergence patterns observed in primates. There are many possible ways to condition on variation across the non-human primates. We focus here on sites that are variable in human only (denoted *human-private*), as well as sites where exactly one other species carries the human minor allele (and all others match the human major allele); see [S1 Fig](#) for an alternative conditioning based on the most closely related species carrying the human minor allele. Throughout, we assume that the observation of the human minor allele as the reference allele in another primate implies that the mutation arose independently and fixed in that primate. This assumption may be violated for a small fraction of SNPs when comparing human to our closest relatives (notably, chimpanzee and gorilla [42]), but the overall patterns that we report here are maintained when we consider more distant species for which shared ancestral polymorphism is unlikely (see [Materials and Methods](#) and [S1 text](#) for further discussion). The SFS presented here, unless otherwise stated, are constructed using minor allele frequencies.

Henceforth, we will use the term “substituted species” to refer to the single species in which the human minor allele is observed, and the corresponding *species cSFS* to refer to the human SFS conditional on a substituted-species divergence pattern. For example, “substituted-orangutan” refers to human variants for which the human minor allele is observed in orangutan,

and the human major allele is observed in all other primates; “orangutan cSFS” will refer to the human SFS at these sites (Fig 2A). There were 5,286,937 human-private sites in the data set, and the number of substituted-species sites ranged from 22,209 (substituted-chimpanzee) to 66,254 (substituted-gibbon).

Fig 2B shows a comparison of the human-private cSFS and the orangutan cSFS for nonsynonymous and synonymous sites, respectively. Within each cSFS class, the nonsynonymous spectrum has more rare variants than the synonymous spectrum, as expected given that nonsynonymous variants are more likely to have deleterious effects. Secondly, if we compare the human-private versus orangutan cSFS at nonsynonymous sites, we see more rare variants in the human-private set. Again, this matches expectations, as the presence of a parallel substitution in orangutan implies that a substitution at this position is tolerated.

However, we were surprised to see that substituted-orangutan synonymous sites also segregate at much higher frequencies than both synonymous and nonsynonymous human-private sites. Taken at face value, this would seem to imply that a large fraction of synonymous sites are functionally constrained. While it is known that some synonymous sites play roles in functions such as splicing [43,44], it is generally believed that most synonymous variants in mammals are effectively neutral. We were thus curious to understand whether this result is primarily driven by a surprising degree of constraint at synonymous sites, or by some other factors.

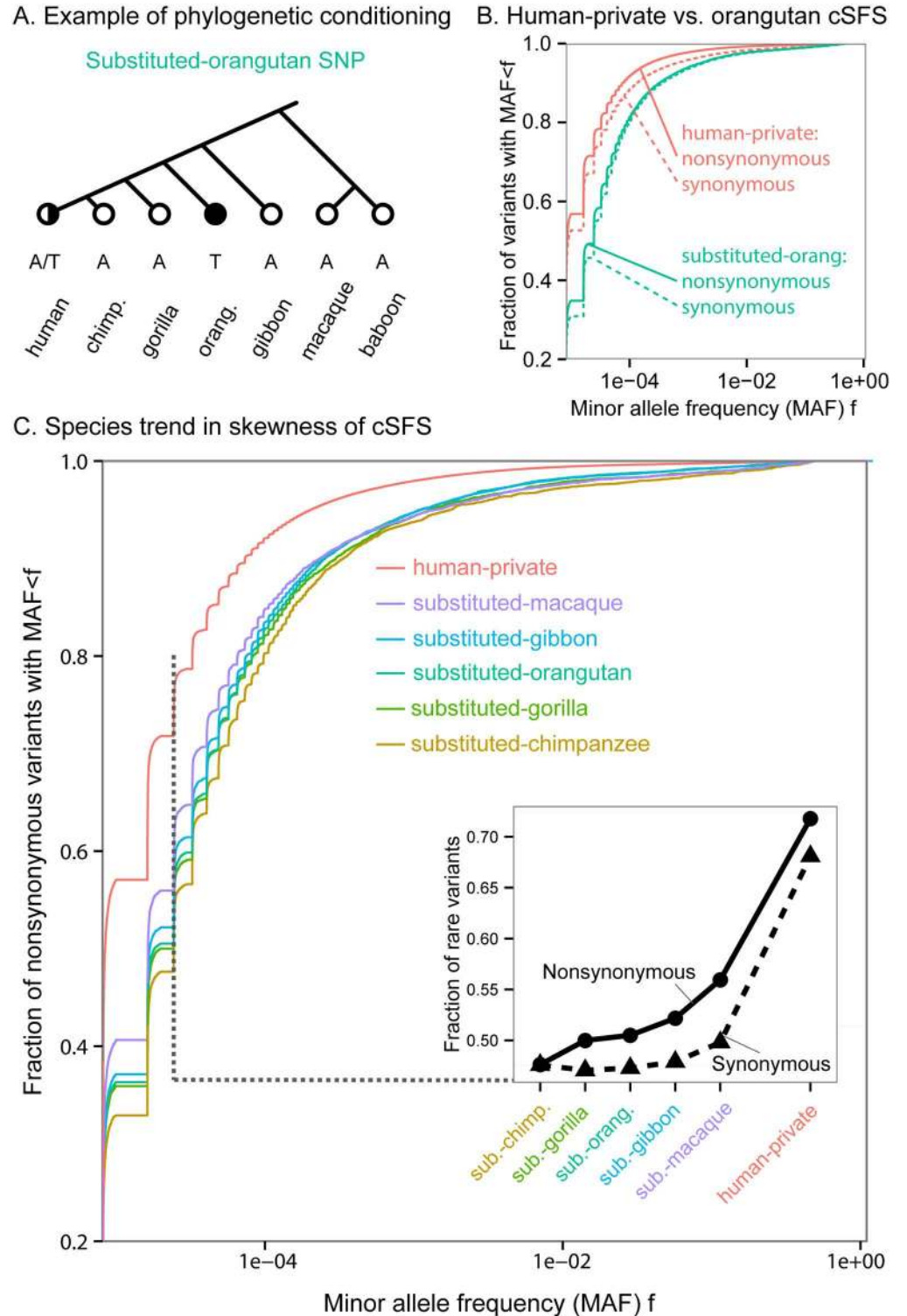
Looking more broadly across the primates, we observed a clear trend of cSFS across substituted species (Fig 2C): the more closely related the substituted species, the greater the skew towards high frequency variants. This trend is most easily noticeable in the fraction of rare variants (defined here, arbitrarily, as singletons and doubletons; Fig 2C, inset). In the following sections, we try to understand the factors driving these observations.

## Effect of mutation rate variation on the human SFS

In this section we consider whether mutation rate variation may contribute to the observed trend across cSFS. Under the standard infinite sites assumption, the SFS is independent of mutation rate. However, we conjectured that in the very large sample size of ExAC, the infinite sites assumption may no longer be a good model for the data [12].

To examine this, we stratified the human SFS by mononucleotide mutation types (as well as the dinucleotide mutation type CpG->TpG), for which there are well-characterized differences in mutation rates. For this analysis we focused on intronic sites, to reduce potential effects of selective constraint. We found that the different mutation types have significantly distinct spectra. The fraction of rare variants among CpG->TpG mutations (36%) was roughly half that of non-CpG transitions (71%, see Fig 3A). Similarly, non-CpG transitions have higher mutation rates than transversions and indeed, the SFS for transitions is also skewed towards higher frequencies than transversions (Fig 3B). Overall, the fraction of rare variants in the subsample of Europeans was significantly negatively correlated with germline mutation rates estimated from the deCode project dataset [45] (weighted linear regression  $p = 4.9 \times 10^{-6}$ , see Fig 4A and S1 text).

If multiple hits are prevalent within the ExAC sample, then some of them should occur in different subpopulations. Higher mutation rates should then lead to excess sharing of low frequency variants among subpopulations. To verify the occurrence of recurrent mutations, we examined the sharing between the European and East Asian ExAC subsamples. Indeed, at low frequencies, non-CpG transitions exhibited a higher sharing rate than transversions, and CpG transitions exhibited much higher sharing rate than non-CpG sites (For example, for sites with



**Fig 2. The human SFS conditioned on primate substitution patterns.** (A) An example of the phylogenetic conditioning that defines what we denote as “substituted-orangutan” sites. (B) The cumulative distribution functions (CDF) of orangutan cSFS (i.e. the SFS of substituted-orangutan sites), and the SFS of phylogenetically conserved sites. The cSFS are more skewed towards common variants than the SFS of conserved sites. These skews are much more pronounced than in the comparison of synonymous and nonsynonymous sites. (C) The more closely related the substituted-species, the higher the skew of the cSFS



towards common variants (only nonsynonymous mutations shown). The inset shows the rare variants slice of the CDF for each species, for both synonymous and nonsynonymous variants.

doi:10.1371/journal.pgen.1006489.g002

a minor allele count of 10, we get a t-test  $p < 2.2 \times 10^{-16}$  for both comparisons; see [Fig 3C](#), and a similar analysis performed in [Fig 2d](#) in the main ExAC paper [12]).

As an additional test of whether mutation rate affects the fraction of rare variants, we turned to sites in transcribed regions. It is known that in such regions, A->G and A->T mutations occur at higher rates on the template (non-coding) strand than on the non-template (coding) strand, due to the effects of transcription-coupled repair or other transcription-associated mutational asymmetries [46–48]. Indeed, as predicted from these rate asymmetries, we observed a 1% difference between the template and the coding strands in the fraction of rare variants in introns (t-test  $p < 2.2 \times 10^{-16}$  for A->G,  $p = 6.0 \times 10^{-7}$  for A->T). C->T mutations also exhibit a small but significant difference (t-test  $p = 3.0 \times 10^{-4}$ ) between the strands, even though, to our knowledge, no previous work has observed a rate asymmetry for C->T mutations ([Fig 3D](#)).

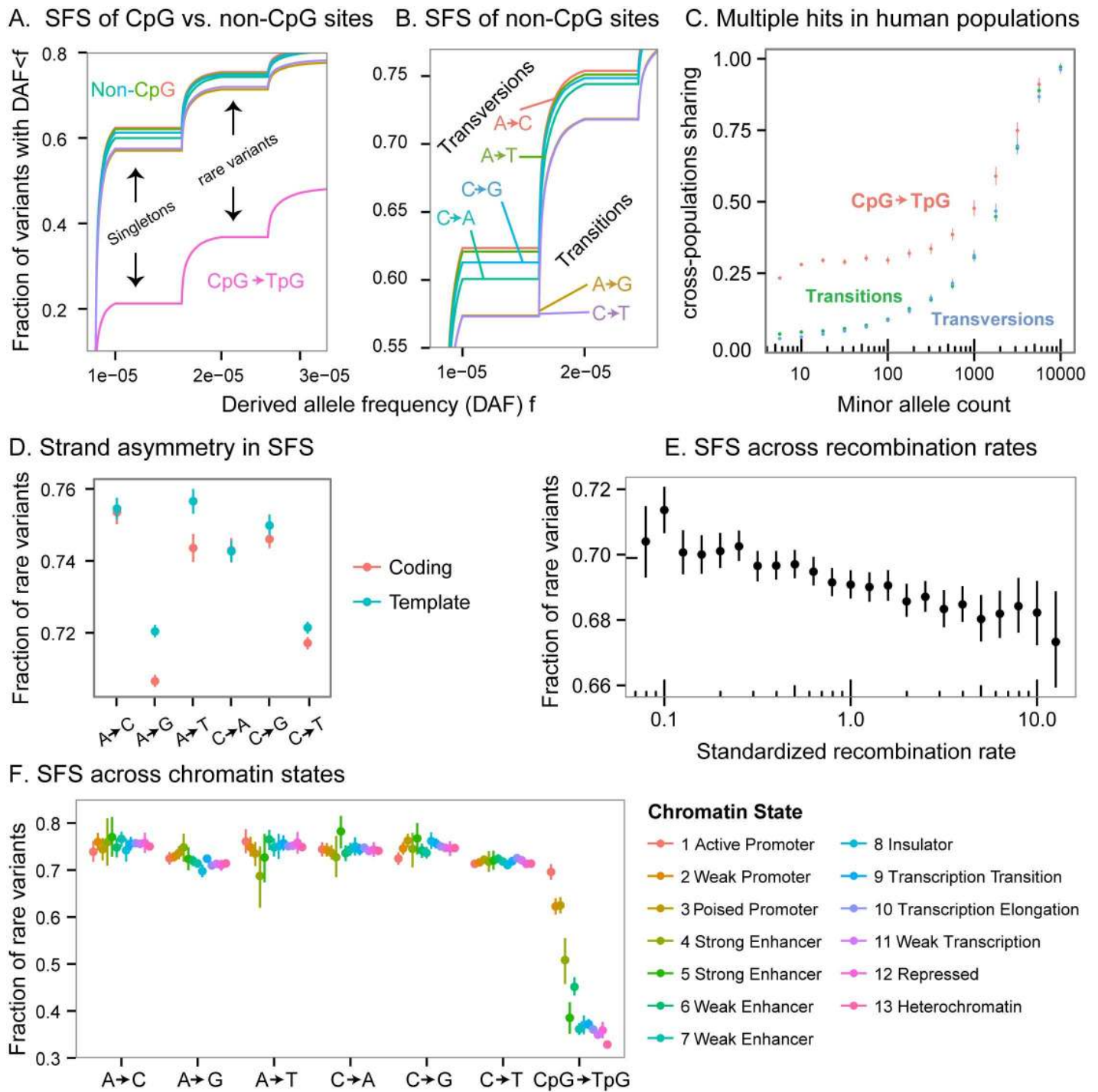
Similarly, we hypothesized that the SFS at CpG sites might also depend on chromatin environment. Specifically, CpG sites experience high mutation rates only when they are methylated [49–52]. We thus examined the effect of chromatin states in H1 human embryonic stem cell lines, inferred by ChromHMM [53] (as a proxy for germline chromatin states) on the SFS across different mutation types. Methylation levels are expected to be low in active regions including promoters and enhancers and high in repressed regions such as heterochromatin. Indeed, we find highly significant differences in the SFS at CpGs (see [S1 text](#) for details), consistent with this expectation: i.e., fewer rare variants in heterochromatin, where methylation levels are high. In contrast, the other mononucleotide mutations showed only modest variation across chromatin states ([Fig 3F](#)).

Finally, we found that recombination rate is also negatively correlated with the fraction of rare variants (Pearson correlation  $p < 2.2 \times 10^{-16}$ , and see [Fig 3E](#) and [S5 Fig](#)). This is consistent with the postulated positive correlation between recombination and mutation rates [54,55]. However, linked selection—which is expected to be more pervasive in regions of low recombination—could also contribute to this trend [56–58]. Overall, the SFS variation patterns across chromatin states, recombination rates, and strands, underscores that heterogeneity in mutation rates does exist within mutation types, and that it has a substantial effect on the SFS.

These observations on mutation rate variation led us to conclude that the infinite-sites model provides a poor fit for these large-sample human polymorphism data. We therefore investigated finite-sites mutational models. Below, we describe the fit of various mutational models while using previously-inferred population genetic models of European demography. In particular, we eventually used a modified version of the demography inferred by Nelson et al. [14] (see [Materials and Methods](#) for the other demographic models considered). The assumed demography provides a good fit for the SFS of sites with the lowest mutation rates.

We asked how well different finite-sites models account for the observed relationship between de-novo mutation rates and the SFS. First, we considered the Jukes-Cantor model, which uses a 4 x 4 uniform mutation transition matrix [59]. But we were surprised to find that this finite sites model barely improved the fit to the SFS across the range of estimated mutation rates ([Fig 4A](#)). In our simulations, the probability of obtaining more than one mutation on the genealogy of a segregating site is low enough that the finite sites SFS is similar to the infinite sites SFS, even at the relatively high mutation rate estimated for CpGs.

We hypothesized that we might achieve a better fit if we consider the fact that some sites have higher intrinsic mutation rates than the mean for the particular nucleotide change at that



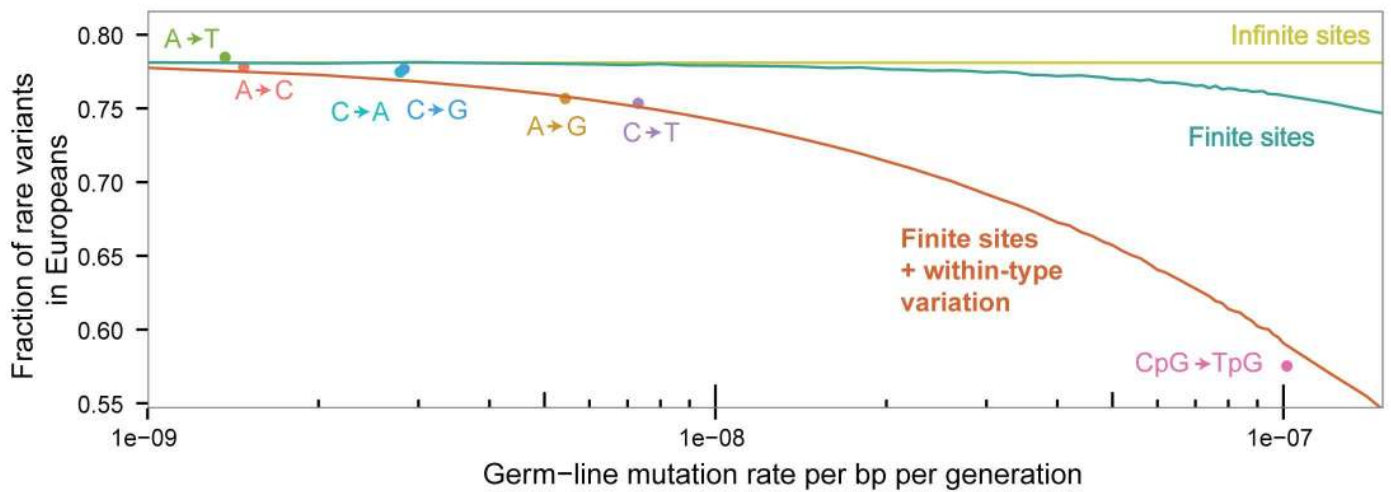
**Fig 3. Mutation rates shape the SFS.** All panels show the SFS of derived alleles constructed from intronic sites. The notation of mutation types refers to mutations on either strand (e.g., A→C indicates an A to C change on either strand). **(A)** SFS stratified by non-CpG mononucleotide mutation types and CpG transitions, represented by different curves. The fraction of rare variants in CpG transitions is nearly half that of other mutations. **(B)** Focusing on non-CpG mutations, transitions have an SFS significantly skewed towards common variants compared with transversions. **(C)** Sharing of polymorphisms between East Asians and Europeans. The excess sharing of CpG polymorphisms at low frequencies is suggestive of multiple occurrences of the mutations. x-axis values are binned on a logarithmic scale. **(D)** Stratification to coding and template strands revealed differences between the two for some mutation types, suggesting transcription-associated mutational mechanisms also affect the SFS. CpG mutations excluded from the analysis in this panel. **(E)** Recombination rates are negatively correlated with the fraction of rare variants; this could be due to a correlation between recombination rates and mutation rates. x-axis values are standardized to the genomewide mean, and are binned on a logarithmic scale. **(F)** SFS across chromatin states. Chromatin states in H1 human embryonic stem cells were inferred by ChromHMM. The chromatin state exhibits substantial association with the fraction of

rare variants in CpG mutations, and modest association in other mononucleotide mutation types. In panels D,E and F: Points show means; lines show 95% confidence interval computed with nonparametric bootstrap.

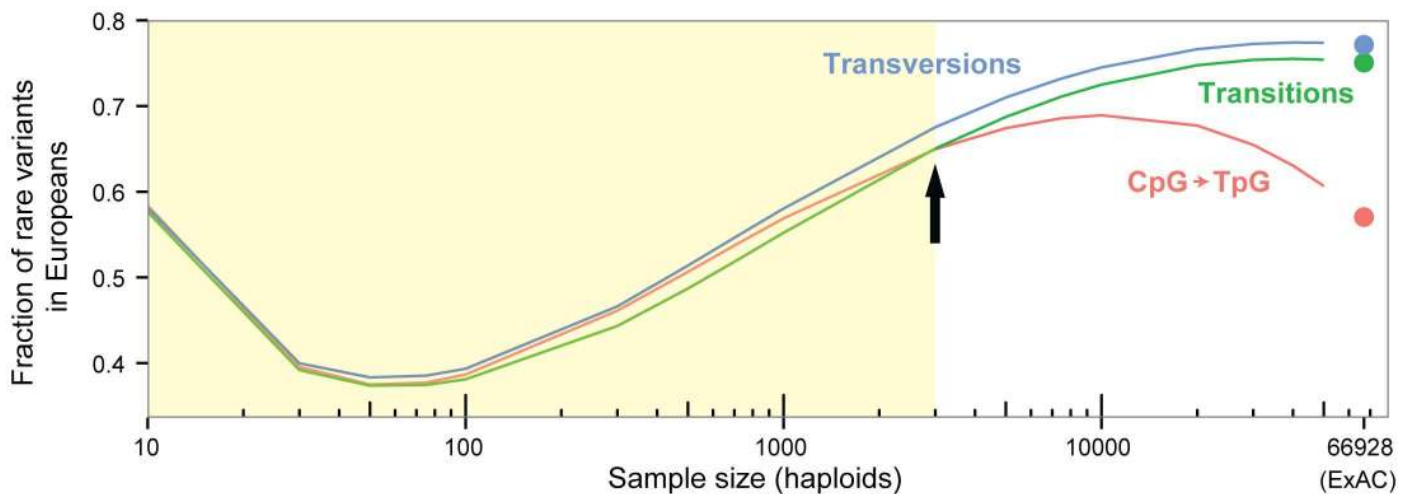
doi:10.1371/journal.pgen.1006489.g003

site; this notion has received increasing support in the recent decade from both evolutionary and family-based studies of human mutation rates [32–35,37,60–62]. We therefore augmented the Jukes-Cantor model by incorporating additional variation in mutation rates across sites belonging to each mutation type (see [Materials and Methods](#)). The augmented Jukes-Cantor model with within-mutation-type variation fitted the data well, including the large difference

**A. Fit of mutational models to observed SFS**



**B. SFS subsampling and the effect of mutation rate**



**Fig 4.** All panels exhibit the unfolded SFS (i.e., constructed using the derived alleles) of intronic sites. **(A)** Fit of mutational models to observed SFS. The x-axis shows previously estimated de-novo germ-line mutation rates [45]. These data illustrate that the fraction of rare variants is strongly negatively correlated with germ-line mutation rates. Lines show expectations under various mutational models: yellow—infinite sites model (SFS independent of mutation rate); teal—Jukes Cantor finite-sites model; red—Jukes-Cantor model with within-mutation-type variation (i.e., variation beyond mutation rate heterogeneity due to the type of mutation in sequence). **(B)** SFS subsampling and the effect of mutation rate. Dots show the fraction of rare variants in the full sample SFS of the European population in ExAC. Lines show the expected fraction of rare variants after subsampling to smaller numbers of individuals. In large samples, the SFS of CpG and non-CpG sites are very different. In smaller samples, these differences shrink. In the shaded region, the trend across mutation types is changed (the inflection point is indicated by an arrow); with these sample sizes, CpG transitions exhibit more rare variation than non-CpG transitions.

doi:10.1371/journal.pgen.1006489.g004



in SFS between CpG and non-CpG sites (Fig 4A). The augmented model suggests that 3% of mutations within a mutation type have a mutation rate of over 5 times the mean rate for that type. This estimate is close to the level of mutation rate variation inferred by Hodgkinson and Eyre-Walker [63].

It is natural to wonder what effect recurrent mutations may have in smaller samples. Small samples have the disadvantages of increased noise and limited temporal resolution of analysis. For example, in demographic inference, larger samples are essential for detecting the signal of recent rapid growth of the human population [17,64,65]. Interestingly, we found that samples much smaller than ExAC may also create an unappreciated bias, as we describe next.

We examined the effect of subsampling the SFS of the European ExAC sample to a smaller number of individuals (see S1 text). SFS differences between non-CpG transitions and transversions remained roughly the same, even with a sample of a few hundred people. Conversely, the difference between CpG and non-CpG sites changed dramatically for smaller samples. For samples smaller than 1500 people, there appears to be more rare variation in CpG than non-CpG transitions (Fig 4B, S4 Fig). This finding exemplifies that if one category of sites has substantially more rare variation in the population than a second category, the sample SFS may actually exhibit more rare variation in the second category. Therefore, a comparison of the amount of rare variation across categories of sites may yield different orderings, depending on the sample size.

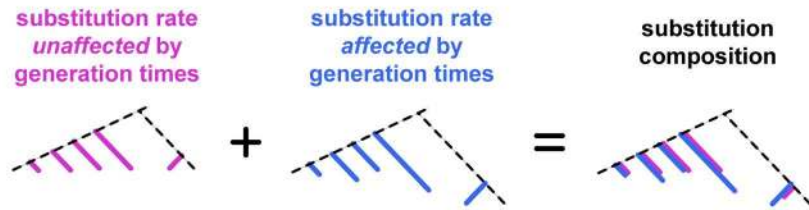
Finally, we returned to the species trend across cSFS that we described earlier (Fig 2C). Given the previous observations on SFS differences between mutation types, we asked whether the trend across substituted-species cSFS we described earlier (Fig 2C) could be explained by differing compositions of the various mutation types. Indeed, most of this trend is due to the fact that CpG transitions make up a higher fraction of sites for more closely related substituted species (Spearman  $\rho = -0.9$ ,  $p = 0.08$ , and see S2 Fig). Since CpG transitions are depleted of rare variants, this results in the cSFS skewness trend. Namely, the fraction of rare variants is strongly negatively correlated with the fraction of CpG transitions across substituted species (Pearson  $r = -0.997$ ,  $p = 9.7 \times 10^{-6}$  for nonsynonymous mutations;  $r = -0.999$ ,  $p = 9.9 \times 10^{-7}$  for synonymous mutations, see Fig 5C).

Why is the fraction of CpG transitions negatively correlated with the relatedness of the substituted-species to humans? Below, we suggest how this could be explained through the mutational mechanism of CpG transitions, which leads to different substitution dynamics on evolutionary timescales than the dynamics at non-CpG sites [66,67].

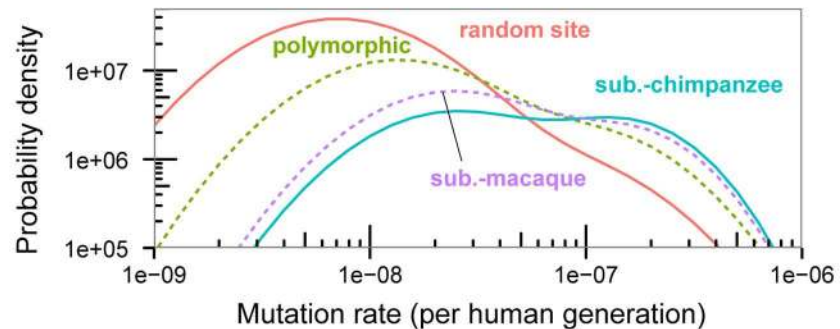
Substituted-species sites likely experienced two independent mutations at the site during primate evolution, and are therefore enriched for hypermutable sites [60,61,63]. A simple model that we develop in S1 text supports this intuition. In this model, we initially assumed a “uniform molecular clock” regime in which substitutions accumulate at the same yearly rate across the primate phylogeny. Under this assumption, differences in the distribution of mutation rates between substituted-species categories should be vanishingly small (S6 Fig).

However, recent work [66,68–71] has demonstrated that while the “uniform clock” assumption is valid for some mutation types—importantly, CpG transitions—the yearly substitution rates of other mutation types depend heavily on life-history traits such as generation time [70,72,73], and thus vary extensively across primates. Notably, Moorjani et al. have pointed out that this difference leads to variable mutational spectra across primates [68]. We therefore augmented our model by including two mutation categories: mutation types that follow a “uniform clock”, and mutation types with rates that depend on generation times (S1 text, and see Fig 5A). The model predicts an enrichment of uniform-clock mutations for substituted-primates with longer generation times. Notably, this translates into a prediction of an enrichment of uniform-clock mutations—like CpG transitions—in substituted species more closely-

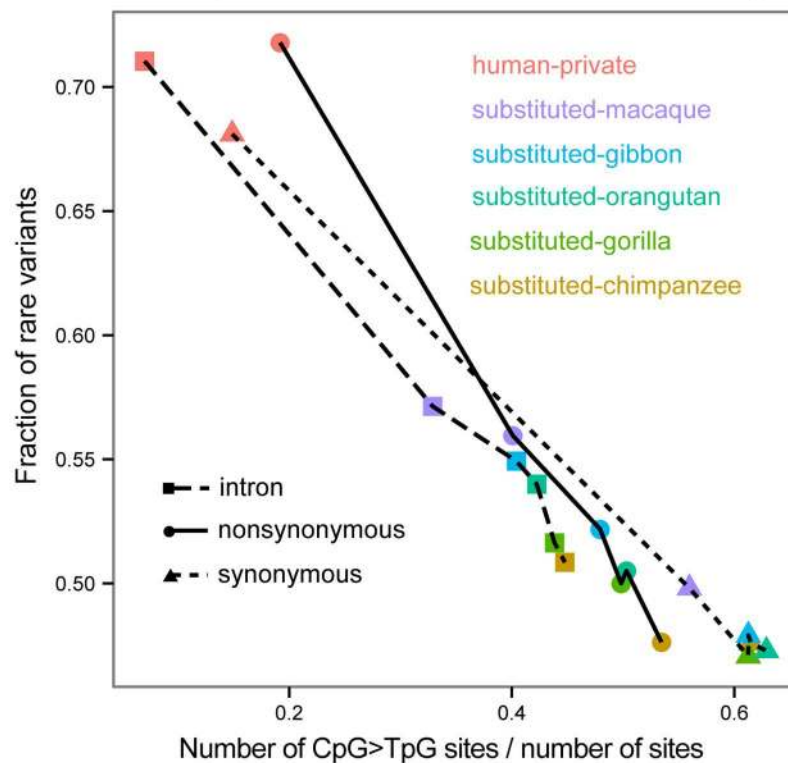
A. Generation time effect and substitution type composition



B. Predicted mutation rate distribution at substituted-species sites



C. Fraction of CpGs and cSFS skewness across sub. species



**Fig 5. (A)** Some mutation types accumulate in a roughly constant yearly rate across different primate lineages. For these mutation types the expected number of substitutions on an evolutionary branch is proportional to the branch length in years (pink). The yearly rates of other substitution types (blue) depends on various life-history traits like generation times (“generation time effect”). As a result, the composition of substitution types in a lineage depends on lineage-specific traits like generation times; this is illustrated by the blue to pink ratio, which differs across lineages. **(B)** Model-based expectations for the distribution of mutation

rates at substituted-species sites. These results were computed using a theoretical model and a set of realistic parameters. At substituted species sites, we expect a distribution skewed towards higher mutation rates compared to random sites, or to random polymorphic sites. In addition, the distribution of mutation rates is skewed towards higher mutation rates for substituted-species with longer generation times; for the primates we considered in this work, this would imply higher mutation rates for more closely-related substituted-species. **(C)** CpG transitions enrichment is a strong predictor of cSFS skewness in real data.

doi:10.1371/journal.pgen.1006489.g005

related to humans (with the exception of orangutan, which is thought to have the longest generation time among the primates considered, although it was only estimated in females [74]). Examining the expected distributions of mutation rates in substituted-species sites, this enrichment leads to a skew towards higher mutation rates for more closely-related substituted-species (Fig 5B).

Overall, this model provides an explanation by which mutational mechanisms underlie the observed correlation between the relatedness of the substituted species and the skew of its cSFS towards common variants.

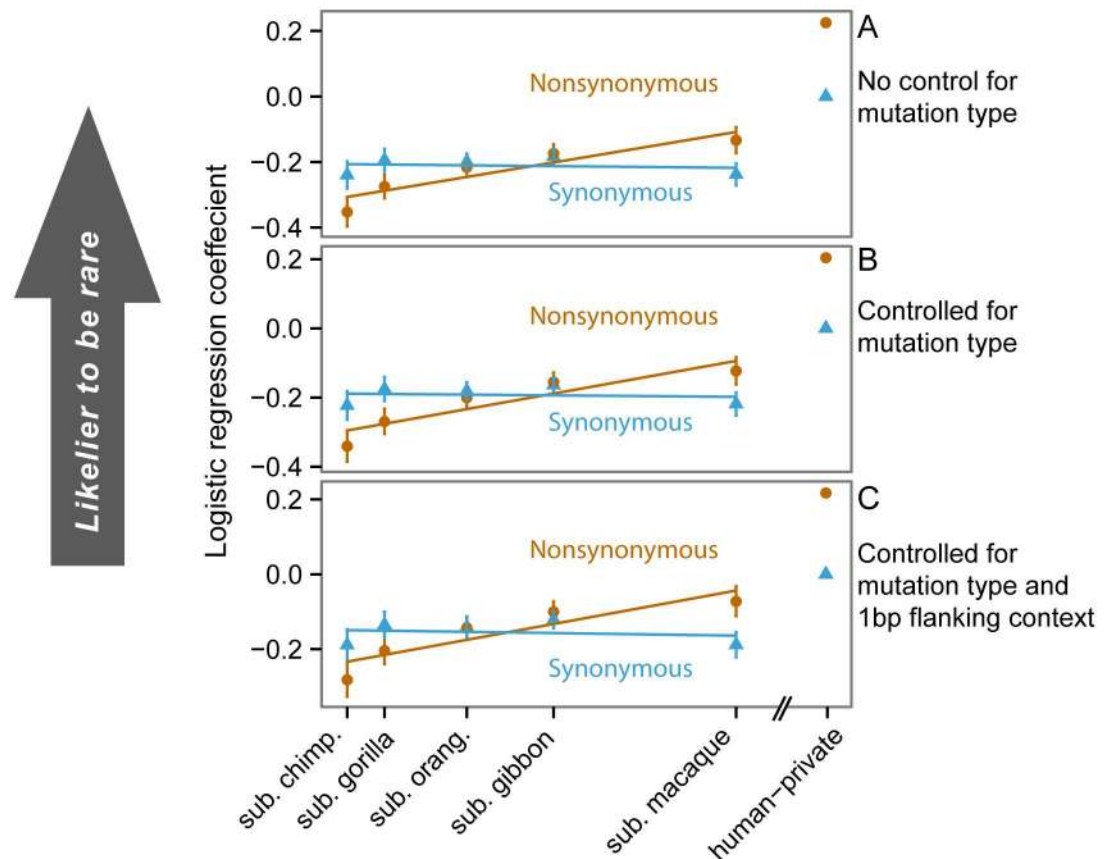
### Residual correlation of substituted-species relatedness and cSFS skewness

Finally, we asked whether additional causes beyond mutation rate variation might also contribute to the species trend across cSFS. To this end, we used a logistic regression model to examine whether the probability of the variant being rare is associated with the relatedness of the substituted species (see [Materials and Methods](#)). A model that included only the relatedness of the substituted species showed a perfectly correlated ordering of the two (Fig 6A, CpG transitions were excluded from this analysis). We then turned to examine whether this correlation persists after controlling for mutational composition differences between substituted-species categories. We controlled for the effect of mononucleotide mutation types on the probability of the variant being rare (Fig 6B). We then further refined the mononucleotide mutation types by using their two flanking nucleotides, and estimated another model with these finer mutation type categories (Fig 6C). The trend persisted even after controlling for mutation type, most noticeably for nonsynonymous sites, which likely involve the strongest purifying selection pressures (Spearman  $\rho = 1$ ,  $p = 0.016$  for the ordering of substituted-species coefficients for both models). We repeated the analysis while including CpG transitions, and found that the perfect correlation persisted (Spearman  $\rho = 1$ ,  $p = 0.016$  for both models; S12 Fig).

There are a few possible drivers of the residual trend observed in Fig 6B and 6C. First, there may be additional variation in mutation rate that is not accounted for by the 3-mer mutational context [32–39]. However, we note that the estimates in Fig 6 change only slightly between panels, suggesting that the unaccounted variation in mutation rates does not greatly bias the estimates. Furthermore, in S1 text we demonstrate that, in non-CpG sites, mutation rate distributions are expected to be very similar across substituted species categories. Therefore, additional mutation rate variation is not expected to contribute to the correlation between the relatedness of the substituted species and the fraction of rare variants.

Another possibility is that differences in selection pressures across substituted-species categories contribute to the cSFS trend. It is possible that substituted-species sites of more closely related species are on average less deleterious. Furthermore, they may be less deleterious particularly for humans. More-related species have, on average, more similar context on which the mutation occurs. When the sequence context of the substituted species is similar to that of humans, the fixation of the human-minor allele in the substituted species suggests that the mutation is benign for humans. As sequence context diverges, epistatic effects may come into

Depletion of rare variants and relatedness of the substituted species



**Fig 6. Depletion of rare variants is correlated with relatedness to substituted species.** The figure shows logistic regression coefficient estimates with their corresponding standard errors. Substituted-species labels are spaced by their split times from humans. The lines are the least-squares line fitting the coefficients to the split times. **(A)** Estimates from a simple logistic regression to the substituted species. The trend is partly due to mutational composition differences between substituted-species categories. To test whether the trend is driven solely by mutational rate differences, we estimate coefficients in a model including the variation explained by **(B)** mononucleotide mutation type, and **(C)** combinations of focal mononucleotide mutations and upstream and downstream nucleotides. Even after controlling for mutational composition with these models, a significant trend persists for nonsynonymous variants.

doi:10.1371/journal.pgen.1006489.g006

play and change the selective effect of the mutation [28,75,76]. In [S1 text](#), we investigate the effect of sequence context divergence more directly (see [S8 Fig](#)).

### Discussion

Our analysis showed a significant correlation between the probability that a variant is rare in humans and the relatedness of another species in which the same mutation occurred. This trend was largely driven by mutation rate variation, which we have observed to be a primary determinant of the human SFS.

The large effect that mutation rate variation has on the human SFS could have a major impact on any future work involving human polymorphism datasets with large sample sizes. For example, most demographic inference algorithms that use the SFS as a summary statistic [e.g. 6,65,77] rely on the infinite-sites model, which is evidently not a valid assumption for large samples. Adjusting demographic inference schemes to include the effects of recurrent

mutations on the SFS (for examples of recent efforts towards this goal, see [78–82]) has the potential to significantly improve inference accuracy.

We have also seen that the trend across cSFS persisted even after tri-nucleotide mutational composition was taken into account. This remaining correlation is consistent with differences in selection pressures across substituted-species categories.

Substitutions in other lineages have proven to be highly informative for understanding deleterious effects in the contemporary human genome; among numerous features that have been considered, the strongest predictors of the pathogenicity of a mutation are species divergence features [83–86]. Nevertheless, methods used to predict the deleteriousness of a mutation at a site typically rely on a single summary of how variable a site is across the phylogenetic tree. Our analysis suggests that the location of a mutation on the evolutionary tree is informative of how deleterious the same mutation is for humans. It is our hope that the integration of divergence patterns and sequence context into methods that predict the fitness or health effects of human mutations could increase accuracy and predictive power.

## Materials and Methods

### Data

For polymorphism data, we downloaded single nucleotide polymorphism (SNP) data from version 0.2 of the Exome Aggregation Consortium database [40]. This database is a standardized aggregation of several exome sequencing studies amounting to a sample size of over 60,700 individuals and approximately 8 million SNPs. For each SNP we extracted upstream and downstream 30 nucleotides in the coding sequence of the human reference genome hg19 build. For simplicity, we excluded sites that are tri-allelic (6.5% of all SNPs) or quad-allelic (0.2% of all SNPs).

For divergence data, we used the following reference genome builds downloaded from the UCSC genome browser [41]: chimpanzee (panTro4), gorilla (gorGor3), orangutan (ponAbe2), gibbon (nomLeu1), macaque (rheMac3), and baboon (papHam1). We used the UCSC genome browser's liftOver program to align each ExAC SNP along with its 60bp sequence context to the six aforementioned reference genomes. We used the baboon reference genome solely for the ascertainment of all other substituted-species categories (rather than including a substituted-baboon category in the analysis).

For gene annotations, we downloaded the refGene table of the RefSeq Genes track from the UCSC genome browser. For each SNP in our data, we extracted all gene isoforms in which the position was included. We kept all ExAC SNPs that fell in a coding exon, intron or untranslated region. We excluded from the analysis non-autosomal SNPs, SNPs that had multiple annotations corresponding to different transcript models, and SNPs with a sample size of less than 100,000 chromosomes. After applying the filters, we were left with 6,002,065 SNPs.

For recombination rates, we downloaded the sex-averaged recombination rate map constructed by Kong et al. [87], which estimates rates at a resolution of 10kbp bins.

### The probability of ancestrally shared polymorphisms

In order to construct an upper bound on the probability of a human polymorphic site being ancestrally shared with another species, we consider the case of a selectively neutral polymorphism shared with chimpanzees. A polymorphism observed in the human sample at the current time is an ancestral polymorphism at the time of the human-chimpanzee split only if there are at least two lineages ancestral to the human sample at the human-chimpanzee split time. Leffler et al. [42] assume a constant human effective population size of 10,000 people throughout history, and estimate a probability of about  $1.0 \times 10^{-5}$ . In [S1 text](#), we augment Leffler



et al.'s approximation with more complex demographic models for recent human history and derive an upper bound of  $1.4 \times 10^{-5}$  for this probability. Multiplying this probability by the number of exonic sites (3,531,936) in our data, we get an expected number of 49 sites in our data that are ancestrally shared with chimpanzees.

However, our derivations are based on a pre-out-of-Africa effective population size ( $N_e^{\text{pOOA}}$ ) of 10,000 people. Very little is known about human demographic history prior to the out-of-Africa event, and as we show in [S1 text](#), the probability of an ancestral polymorphism rises very quickly with increasing  $N_e^{\text{pOOA}}$ . Estimates of  $N_e^{\text{pOOA}}$  range between 7300 [88] and 12,500 [89] people and are continually revised as estimates of human mutation rate and demographic history are refined. With  $N_e^{\text{pOOA}} = 12,500$ , we get an upper bound of 283 polymorphisms in the dataset that are expected to be shared between human and chimpanzee, which compose at most 1.2% of the substituted-chimpanzee sites. The upper bound for other species, or sites under purifying selection, should be even smaller, and are overall too few to affect our results. We therefore conclude that ancestral polymorphisms are too few to significantly affect our analysis.

## Simulating mutational models and demographic models

To get a theoretical expectation for the fraction of rare variants under different mutational models, we used various software for computing the expected sample SFS of 33,750 diploid individuals, corresponding to the size of the non-Finnish European subsample in the ExAC dataset. For all mutational models, which we describe below, we generated predictions under various demographic models from recent literature: Gazave et al. [90] (model 2 in their work), Tennessen et al. [16] and Nelson et al. [14].

For the infinite-sites model, we computed the expected sample SFS analytically using fastNeutrino [65]. The infinite-sites model corresponds to an upper bound for the fraction of rare variants, but nonetheless predicted a fraction of rare variants much lower than that observed in data (75%-78%) for all non-CpG mutations under the Gazave et al. (59%) and Tennessen et al. (60%) demographies. The Nelson et al. model, which was inferred using a larger sample size of 11,000 people predicted 75% of biallelic polymorphisms would be rare under the infinite-sites model. In order to fit the highest observed fraction of rare variants for non-CpG sites in the ExAC data, we modified the parameters of the most recent epoch of exponential growth in Nelson et al. We estimated these parameters using fastNeutrino [65] on all A->C intronic mutations from ExAC. The inferred parameters were: current effective population size of 4,009,877 diploids, and an exponential growth onset time of 119.47 generations in the past with a growth rate of 5.38% per generation. The more ancient demographic parameters were fixed to the same values as in the model of Nelson et al.

We assume multiple mergers (non-Kingman merger events) have negligible effect on the SFS since the sample size is significantly smaller than the current effective population size. A similar demographic model [91] with a four-fold smaller current effective population size exhibited a relative difference of only about 1.3% and 0.3% in the proportion of singletons and doubletons respectively for a comparable sample size of 50,000 people. Hence, we felt confident in using the Kingman coalescent for drawing genealogies.

For the finite sites model, we first simulated independent coalescent trees using ms [92] and then generated 1kb non-recombining sequences for each coalescent tree using the desired recurrent mutation rate with the 4 x 4 Jukes-Cantor model of mutation [59]. We used the program Seq-Gen [93] to drop recurrent mutations on coalescent trees drawn from ms. We used mutation rates in a uniform logarithmically spaced grid of 40 points ranging from  $10^{-9}$  to  $5.3 \times 10^{-5}$  mutations per basepair per generation per haploid. For each value of the mutation

rate, we simulated enough sequence data so that at least 100,000 biallelic polymorphic sites were available to reliably estimate the expected fraction of rare variants. If we indicate whether a variant is rare by  $Y$ , then for each mutation rate  $\mu$ , the expected fraction of rare variants is

$$E[Y|S = 1; \mu],$$

where  $S$  is an indicator variable indicating whether a site is polymorphic and, specifically, biallelic. Finally, we considered a model with additional, within-mutation-type heterogeneity in mutation rate. Specifically, we considered a model in which sites of a particular mutation type (e.g., C->A sites) have a mean mutation rate  $\mu$  as before, but the mutation rate itself,  $M$ , is no longer fixed (and equal to  $\mu$ ), but rather a random variable with mean  $\mu$ . Let  $f(M|S = 1; \mu)$  be the probability density function of  $M$  in a site with mean mutation rate  $\mu$  conditional on it being biallelic. Then, by the law of total expectation we have:

$$E[Y|S = 1; \mu] = \int E[Y|M, S = 1]f(M|S = 1; \mu)dM.$$

By Bayes' rule,  $f(M|S = 1; \mu)$  is determined by both the within-mutation-type distribution of mutation rates,  $g(M; \mu)$ , and the probability of a site with mutation rate  $M$  being biallelic, as follows:

$$f(M|S = 1; \mu) = \frac{P(S = 1|M)g(M; \mu)}{P(S = 1; \mu)}.$$

Therefore,

$$E[Y|S = 1; \mu] = \int E[Y|S = 1, M] \frac{P(S = 1|M)g(M; \mu)}{\int P(S = 1|M')g(M'; \mu)dM'} dM.$$

For a large range of  $M$ , we have already estimated  $E[Y|M]$  as described above. From the same simulations we have estimated the probability of a site with mutation rate  $M$  being a biallelic polymorphism,  $P(S = 1|M)$ . Lastly, the distribution of mutation rates due to within-mutation-type variance was modeled using a lognormal distribution:

$$\log_{10} M; \mu \sim N\left(\log_{10} \mu - \frac{\sigma^2}{2} \ln(10), \sigma^2\right).$$

The mean parameter in the lognormal distribution above ensures that  $E[M] = \mu$ .  $\sigma$  was arbitrarily chosen to be 0.57 (red line in Fig 4A). Notably, Hodgkinson et al. also fit a lognormal distribution of mutation rates to their dataset of co-occurrence of SNPs in chimpanzees and humans, and estimate a similar value of  $\sigma = 0.83$  for non-CpG mutations [63] and  $\sigma = 0.8$  CpG transitions (personal correspondence).

### Logistic model for the probability of a variant to be rare

We tested whether the species trend across cSFS is due solely to the effect of mutation rate variation. We used a logistic regression model to examine whether a residual substituted-species trend remains after controlling for mutation type. Let  $Y$  be a binary-valued random variable indicating whether a variant is rare,  $\vec{\mu}$  be a vector of mutually exclusive indicator (dummy) variables for each mutation type,  $\vec{s}$  be a vector of mutually exclusive indicator variables for the divergence pattern for the variant (substituted in one of the primates or human-private) and  $Z$  be an indicator of whether the variant is nonsynonymous (we only considered coding

variants). We fitted the logistic regression model

$$\text{logit}(P(Y = 1|\vec{\mu}, \vec{s}, Z)) = \beta_0 + \vec{\beta}_\mu \cdot \vec{\mu} + (1 - Z)\vec{\beta}_s^{\text{syn}} \cdot \vec{s} + Z\vec{\beta}_s^{\text{ns}} \cdot \vec{s},$$

where the parameters  $\beta_0, \vec{\beta}_\mu, \vec{\beta}_s^{\text{syn}}$ , and  $\vec{\beta}_s^{\text{ns}}$  were learned from the data. We tested whether the coefficients  $\vec{\beta}_s^{\text{syn}}, \vec{\beta}_s^{\text{ns}}$  exhibit a trend across  $s$ , i.e. whether the probability of the variant being rare is associated with the relatedness of the substituted species. When ignoring the mutation rate effect (i.e. fixing  $\vec{\beta}_\mu \equiv 0$ ), the  $\vec{\beta}_s^{\text{ns}}$  estimates were perfectly anti-correlated with the relatedness of the substituted species to human, consistent with the observation in data (Fig 6A). We then allowed for an effect for the mutation type by estimating  $\vec{\beta}_\mu$  for the different categories of mononucleotide mutation types (Fig 6B). We also estimated a model with a finer resolution of mutational categories, further partitioning the mononucleotide mutation types by their two flanking nucleotides (Fig 6C). For nonsynonymous sites, which likely involve the strongest purifying selection pressures, the trend persisted even after controlling for mutation rate variation (Spearman  $\rho = 1$ ,  $p = 0.016$  for both mononucleotide correction and for the correction including flanking nucleotides context).

## Supporting Information

### S1 text. Supplementary text.

(PDF)

### S1 Fig. Alternative conditioning of the human SFS on primate divergence patterns.

Changes in the human SFS as we condition on divergence patterns in primates. Here, we label the SFS by the primate most closely related to human carrying the human minor allele. (A) An example of the phylogenetic conditioning for sites in which orangutan is the closest primate which carries the human minor allele. (B) The cumulative distribution functions (CDF) of the SFS of sites substituted in orangutan, and the SFS of human-private mutations. The SFS of sites substituted in orangutan have a skew towards common alleles compared with human-private sites. (C) The more closely related the species with the substitution, the higher the skew of the SFS towards common variants (only nonsynonymous mutations shown).

(TIF)

**S2 Fig. Mutational composition of cSFS.** The fractions of three mutational categories associated with different mutability are shown for each substituted-species category.

(TIF)

**S3 Fig. The probability of an ancestral polymorphism as a function of the effective population size prior to the out-of-Africa event.** The y-axis shows an upper bound on the probability for a random polymorphism originating prior to the human-chimpanzee split time. Both models shown assume a constant effective population size prior to the out-of-Africa event (OOA); the pink line gives the upper bound for the Nelson et al. [7] demographic model after the OOA event, whereas the teal “worst-case” line gives the bound assuming no coalescence events occur between the OOA event and the present time.

(TIF)

**S4 Fig. SFS subsampling and the effect of mutation rate.** This figure summarizes the same analysis as in Fig 4B, with a summary of the fraction of variants with minor allele frequency (MAF) below or equal to 1%, instead of the fraction of rare variants. Dots show the fraction of variants with  $\text{MAF} \leq 1\%$  in the full sample SFS of the European population in ExAC. Lines show the expected fraction of variants with  $\text{MAF} \leq 1\%$  after subsampling to a smaller number

of individuals. The trend between mutation types changes as the sample size varies (with an inflection point between CpG and non-CpG transitions marked by the border of the shaded region).

(TIF)

**S5 Fig. Recombination rate and the skewness of the SFS.** Recombination rate is positively correlated with mutation rate, and is likely driving the negative correlation of recombination rate and the fraction of rare variants. The different panels demonstrate that mutations that are subject to biased gene conversion, as well as those that are not, exhibit a negative correlation between recombination rate and the fraction of rare variants. x-axis values are binned on a logarithmic scale, and are standardized to a genomewide mean.

(TIF)

**S6 Fig. Distributions of mutation rates at substituted-species sites, in a constant molecular-clock model.** These results were computed using a simple analytic model and a set of realistic parameters. At substituted-species sites, we expect a distribution skewed towards higher mutation rates compared to random sites, or to random polymorphic sites. However, the distribution of mutation rate changes only slightly across substituted species, as exemplified by the purple (substituted-gibbon) and teal (substituted-chimpanzee) lines.

(TIF)

**S7 Fig. cSFS species trend stratified by mutation type.** This figure summarizes the same analysis as in the inset of [Fig 2C](#), stratified by non-CpG mononucleotide mutation types and CpG. Red stars denote nonsynonymous trends that exhibited significant Spearman correlation at a significance level of 10%.

(TIF)

**S8 Fig. Effects of sequence context at nonsynonymous sites.** Dots show means. Lines and shaded regions show simple logistic model fits to the data and associated confidence bands. The more diverged the sequence context in the substituted-species is from humans, the higher the fraction of rare variants. The two panels exhibit the same trend with different measures of sequence context divergence in a window of 9 residues upstream and 9 downstream of the SNP. In human-private sites, the trend is reversed: the fraction of rare variants decreased with sequence context divergence from gibbon, consistent with higher regional mutation rates and lower constraint implied by higher sequence context divergence.

(TIF)

**S9 Fig. Epistatic effects of sequence context—amino acid context divergence.** Dots show means. Lines and shaded regions show predictions of a simple logistic model fits to the data and their associated confidence bands. The more diverged the sequence context in the substituted-species is from humans, the higher the fraction of rare variants. Sequence context divergence is calculated as the number of amino acid substitutions in a window of 9 residues upstream and 9 downstream of the SNP. In human-private sites, the trend is reversed: the fraction of rare variants decreases with sequence context divergence from the substituted species, consistent with higher regional mutation rates and lower constraint implied by higher sequence context divergence.

(TIF)

**S10 Fig. Epistatic effects of sequence context—distance from nearest substitution.** Dots show means. Lines and shaded regions show predictions of a simple logistic model fits to the data and their associated confidence bands. The more diverged the sequence context in the substituted-species is from humans, the higher the fraction of rare variants. Sequence context

divergence is calculated as the distance to the nearest amino acid substitutions in a window of 9 residues upstream and 9 downstream of the SNP.

(TIF)

**S11 Fig. Branch length parameters for the model of mutation rate distributions at substituted-species sites.**

(TIF)

**S12 Fig. Depletion of rare variants is correlated with relatedness to substituted species including CpG transitions.** Like Fig 6, this figure shows logistic regression coefficient estimates and their corresponding standard errors. Here, however, CpG transitions were included in the analysis. Substituted-species labels are spaced by their split times from humans. The lines are the least-squares line fitting the coefficients to the split times. (A) Estimates from a simple logistic regression to the substituted species. The trend is partly due to mutational composition differences between substituted-species categories. To test whether the trend is driven solely by mutational rate differences, we estimate coefficients in a model including the variation explained by (B) mononucleotide mutation type, and (C) combinations of focal mononucleotide mutations and upstream and downstream nucleotides. Even after controlling for mutational composition with these models, a significant trend persists for nonsynonymous variants.

(TIF)

## Acknowledgments

We thank ExAC and the groups that provided exome variant data for comparison. A full list of contributing groups can be found at <http://exac.broadinstitute.org/about>. We also thank Molly Przeworski, Ziyue Gao, Doc Edge, Xun Lan, David Golan, Kelley Harris, Anil Raj, Eyal Elyashiv, Adam Eyre-Walker and two anonymous reviewers for helpful comments on the manuscript and/or valuable discussions.

## Author Contributions

**Conceptualization:** AH AB JKP.

**Formal analysis:** AH AB.

**Methodology:** AH AB JKP.

**Writing – original draft:** AH.

**Writing – review & editing:** AH AB JKP.

## References

1. Tajima F (1989) Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* 123: 585–595. PMID: [2513255](#)
2. Kimura M (1984) *The neutral theory of molecular evolution*: Cambridge University Press.
3. Jensen JD, Kim Y, DuMont VB, Aquadro CF, Bustamante CD (2005) Distinguishing between selective sweeps and demography using DNA polymorphism data. *Genetics* 170: 1401–1410. doi: [10.1534/genetics.104.038224](#) PMID: [15911584](#)
4. Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, et al. (2005) Genomic scans for selective sweeps using SNP data. *Genome Research* 15: 1566–1575. doi: [10.1101/gr.4252305](#) PMID: [16251466](#)
5. Fu Y-X (1997) Statistical tests of neutrality of mutations against population growth, hitchhiking and background selection. *Genetics* 147: 915–925. PMID: [9335623](#)



6. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD (2009) Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet* 5: e1000695. doi: [10.1371/journal.pgen.1000695](https://doi.org/10.1371/journal.pgen.1000695) PMID: [19851460](https://pubmed.ncbi.nlm.nih.gov/19851460/)
7. Bustamante CD, Wakeley J, Sawyer S, Hartl DL (2001) Directional selection and the site-frequency spectrum. *Genetics* 159: 1779–1788. PMID: [11779814](https://pubmed.ncbi.nlm.nih.gov/11779814/)
8. Nielsen R (2000) Estimation of population parameters and recombination rates from single nucleotide polymorphisms. *Genetics* 154: 931–942. PMID: [10655242](https://pubmed.ncbi.nlm.nih.gov/10655242/)
9. 1000-Genomes-Project-Consortium (2015) A global reference for human genetic variation. *Nature* 526: 68–74. doi: [10.1038/nature15393](https://doi.org/10.1038/nature15393) PMID: [26432245](https://pubmed.ncbi.nlm.nih.gov/26432245/)
10. Fu W, O'Connor TD, Jun G, Kang HM, Abecasis G, et al. (2013) Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature* 493: 216–220. doi: [10.1038/nature11690](https://doi.org/10.1038/nature11690) PMID: [23201682](https://pubmed.ncbi.nlm.nih.gov/23201682/)
11. UK10K-Consortium (2015) The UK10K project identifies rare variants in health and disease. *Nature* 526: 82–90. doi: [10.1038/nature14962](https://doi.org/10.1038/nature14962) PMID: [26367797](https://pubmed.ncbi.nlm.nih.gov/26367797/)
12. Lek M, Karczewski K, Minikel E, Samocha K, Banks E, et al. (2016) Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536: 281–291.
13. Coventry A, Bull-Otterson LM, Liu X, Clark AG, Maxwell TJ, et al. (2010) Deep resequencing reveals excess rare recent variants consistent with explosive population growth. *Nature Communications* 1: 131. doi: [10.1038/ncomms1130](https://doi.org/10.1038/ncomms1130) PMID: [21119644](https://pubmed.ncbi.nlm.nih.gov/21119644/)
14. Nelson MR, Wegmann D, Ehm MG, Kessner D, Jean PS, et al. (2012) An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science* 337: 100–104. doi: [10.1126/science.1217876](https://doi.org/10.1126/science.1217876) PMID: [22604722](https://pubmed.ncbi.nlm.nih.gov/22604722/)
15. Adams AM, Hudson RR (2004) Maximum-likelihood estimation of demographic parameters using the frequency spectrum of unlinked single-nucleotide polymorphisms. *Genetics* 168: 1699–1712. doi: [10.1534/genetics.104.030171](https://doi.org/10.1534/genetics.104.030171) PMID: [15579718](https://pubmed.ncbi.nlm.nih.gov/15579718/)
16. Tennessen JA, Bigham AW, O'Connor TD, Fu W, Kenny EE, et al. (2012) Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science* 337: 64–69. doi: [10.1126/science.1219240](https://doi.org/10.1126/science.1219240) PMID: [22604720](https://pubmed.ncbi.nlm.nih.gov/22604720/)
17. Keinan A, Clark AG (2012) Recent explosive human population growth has resulted in an excess of rare genetic variants. *Science* 336: 740–743. doi: [10.1126/science.1217283](https://doi.org/10.1126/science.1217283) PMID: [22582263](https://pubmed.ncbi.nlm.nih.gov/22582263/)
18. Nielsen R (2005) Molecular signatures of natural selection. *Annu Rev Genet* 39: 197–218. doi: [10.1146/annurev.genet.39.073003.112420](https://doi.org/10.1146/annurev.genet.39.073003.112420) PMID: [16285858](https://pubmed.ncbi.nlm.nih.gov/16285858/)
19. Evans SN, Shvets Y, Slatkin M (2007) Non-equilibrium theory of the allele frequency spectrum. *Theoretical Population Biology* 71: 109–119. doi: [10.1016/j.tpb.2006.06.005](https://doi.org/10.1016/j.tpb.2006.06.005) PMID: [16887160](https://pubmed.ncbi.nlm.nih.gov/16887160/)
20. Ewens WJ (2012) *Mathematical Population Genetics 1: Theoretical Introduction*: Springer Science & Business Media.
21. Cooper GM, Goode DL, Ng SB, Sidow A, Bamshad MJ, et al. (2010) Single-nucleotide evolutionary constraint scores highlight disease-causing mutations. *Nature methods* 7: 250–251. doi: [10.1038/nmeth0410-250](https://doi.org/10.1038/nmeth0410-250) PMID: [20354513](https://pubmed.ncbi.nlm.nih.gov/20354513/)
22. Goode DL, Cooper GM, Schmutz J, Dickson M, Gonzales E, et al. (2010) Evolutionary constraint facilitates interpretation of genetic variation in resequenced human genomes. *Genome Research* 20: 301–310. doi: [10.1101/gr.102210.109](https://doi.org/10.1101/gr.102210.109) PMID: [20067941](https://pubmed.ncbi.nlm.nih.gov/20067941/)
23. Simons YB, Turchin MC, Pritchard JK, Sella G (2014) The deleterious mutation load is insensitive to recent population history. *Nature Genetics* 46: 220. doi: [10.1038/ng.2896](https://doi.org/10.1038/ng.2896) PMID: [24509481](https://pubmed.ncbi.nlm.nih.gov/24509481/)
24. Cooper GM, Stone EA, Asimenos G, Green ED, Batzoglu S, et al. (2005) Distribution and intensity of constraint in mammalian genomic sequence. *Genome Research* 15: 901–913. doi: [10.1101/gr.3577405](https://doi.org/10.1101/gr.3577405) PMID: [15965027](https://pubmed.ncbi.nlm.nih.gov/15965027/)
25. de Visser JAG, Cooper TF, Elena SF (2011) The causes of epistasis. *Proceedings of the Royal Society of London B: Biological Sciences* 278: 3617–3624.
26. Kondrashov AS, Sunyaev S, Kondrashov FA (2002) Dobzhansky—Muller incompatibilities in protein evolution. *Proceedings of the National Academy of Sciences* 99: 14878–14883.
27. Hansen TF (2013) Why epistasis is important for selection and adaptation. *Evolution* 67: 3501–3511. doi: [10.1111/evo.12214](https://doi.org/10.1111/evo.12214) PMID: [24299403](https://pubmed.ncbi.nlm.nih.gov/24299403/)
28. Breen MS, Kemena C, Vlasov PK, Notredame C, Kondrashov FA (2012) Epistasis as the primary factor in molecular evolution. *Nature* 490: 535–538. doi: [10.1038/nature11510](https://doi.org/10.1038/nature11510) PMID: [23064225](https://pubmed.ncbi.nlm.nih.gov/23064225/)
29. Anderson BR, Howell DN, Soldano K, Garrett ME, Katsanis N, et al. (2015) In vivo modeling implicates APOL1 in nephropathy: evidence for dominant negative effects and epistasis under anemic stress. *PLoS Genet* 11: e1005349. doi: [10.1371/journal.pgen.1005349](https://doi.org/10.1371/journal.pgen.1005349) PMID: [26147622](https://pubmed.ncbi.nlm.nih.gov/26147622/)

30. Kulathinal RJ, Bettencourt BR, Hartl DL (2004) Compensated deleterious mutations in insect genomes. *Science* 306: 1553–1554. doi: [10.1126/science.1100522](https://doi.org/10.1126/science.1100522) PMID: [15498973](https://pubmed.ncbi.nlm.nih.gov/15498973/)
31. Jordan DM, Frangakis SG, Golzio C, Cassa CA, Kurtzberg J, et al. (2015) Identification of cis-suppression of human disease mutations by comparative genomics. *Nature Genetics* 524: 225–229.
32. Hodgkinson A, Ladoukakis E, Eyre-Walker A (2009) Cryptic variation in the human mutation rate. *PLoS Biology* 7: e1000027. doi: [10.1371/journal.pbio.1000027](https://doi.org/10.1371/journal.pbio.1000027) PMID: [19192947](https://pubmed.ncbi.nlm.nih.gov/19192947/)
33. Hodgkinson A, Eyre-Walker A (2011) Variation in the mutation rate across mammalian genomes. *Nature Reviews Genetics* 12: 756–766. doi: [10.1038/nrg3098](https://doi.org/10.1038/nrg3098) PMID: [21969038](https://pubmed.ncbi.nlm.nih.gov/21969038/)
34. Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, et al. (2012) Whole-genome sequencing in autism identifies hot spots for de novo germline mutation. *Cell* 151: 1431–1442. doi: [10.1016/j.cell.2012.11.019](https://doi.org/10.1016/j.cell.2012.11.019) PMID: [23260136](https://pubmed.ncbi.nlm.nih.gov/23260136/)
35. Smith T, Ho G, Christodoulou J, Price EA, Onadim Z, et al. (2016) Extensive Variation in the Mutation Rate Between and Within Human Genes Associated with Mendelian Disease. *Human mutation*.
36. Ségurel L, Wyman MJ, Przeworski M (2014) Determinants of mutation rate variation in the human germline. *Annual Review of Genomics and Human Genetics* 15: 47–70. doi: [10.1146/annurev-genom-031714-125740](https://doi.org/10.1146/annurev-genom-031714-125740) PMID: [25000986](https://pubmed.ncbi.nlm.nih.gov/25000986/)
37. Aggarwala V, Voight BF (2016) An expanded sequence context model broadly explains variability in polymorphism levels across the human genome. *Nature Genetics*.
38. Rahbari R, Wuster A, Lindsay SJ, Hardwick RJ, Alexandrov LB, et al. (2015) Timing, rates and spectra of human germline mutation. *Nature Genetics*.
39. Stamatoyannopoulos JA, Adzhubei I, Thurman RE, Kryukov GV, Mirkin SM, et al. (2009) Human mutation rate associated with DNA replication timing. *Nature Genetics* 41: 393–395. doi: [10.1038/ng.363](https://doi.org/10.1038/ng.363) PMID: [19287383](https://pubmed.ncbi.nlm.nih.gov/19287383/)
40. Exome Aggregation Consortium (ExAC), Cambridge, MA. URL: <http://exac.broadinstitute.org> [accessed Jan 2015].
41. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, et al. (2002) The human genome browser at UCSC. *Genome Research* 12: 996–1006. doi: [10.1101/gr.229102](https://doi.org/10.1101/gr.229102) PMID: [12045153](https://pubmed.ncbi.nlm.nih.gov/12045153/)
42. Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, et al. (2013) Multiple instances of ancient balancing selection shared between humans and chimpanzees. *Science* 339: 1578–1582. doi: [10.1126/science.1234070](https://doi.org/10.1126/science.1234070) PMID: [23413192](https://pubmed.ncbi.nlm.nih.gov/23413192/)
43. Chamary J, Parmley JL, Hurst LD (2006) Hearing silence: non-neutral evolution at synonymous sites in mammals. *Nature Reviews Genetics* 7: 98–108. doi: [10.1038/nrg1770](https://doi.org/10.1038/nrg1770) PMID: [16418745](https://pubmed.ncbi.nlm.nih.gov/16418745/)
44. Parmley JL, Chamary J, Hurst LD (2006) Evidence for purifying selection against synonymous mutations in mammalian exonic splicing enhancers. *Molecular Biology and Evolution* 23: 301–309. doi: [10.1093/molbev/msj035](https://doi.org/10.1093/molbev/msj035) PMID: [16221894](https://pubmed.ncbi.nlm.nih.gov/16221894/)
45. Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, et al. (2012) Rate of de novo mutations and the importance of father's age to disease risk. *Nature* 488: 471–475. doi: [10.1038/nature11396](https://doi.org/10.1038/nature11396) PMID: [22914163](https://pubmed.ncbi.nlm.nih.gov/22914163/)
46. Green P, Ewing B, Miller W, Thomas PJ, Green ED (2003) Transcription-associated mutational asymmetry in mammalian evolution. *Nature Genetics* 33: 514–517. doi: [10.1038/ng1103](https://doi.org/10.1038/ng1103) PMID: [12612582](https://pubmed.ncbi.nlm.nih.gov/12612582/)
47. Mugal CF, von Grünberg H-H, Peifer M (2009) Transcription-induced mutational strand bias and its effect on substitution rates in human genes. *Molecular Biology and Evolution* 26: 131–142. doi: [10.1093/molbev/msn245](https://doi.org/10.1093/molbev/msn245) PMID: [18974087](https://pubmed.ncbi.nlm.nih.gov/18974087/)
48. Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, et al. (2015) Genome-wide patterns and properties of de novo mutations in humans. *Nature Genetics* 47: 822–826. doi: [10.1038/ng.3292](https://doi.org/10.1038/ng.3292) PMID: [25985141](https://pubmed.ncbi.nlm.nih.gov/25985141/)
49. Fryxell KJ, Zuckerkandl E (2000) Cytosine deamination plays a primary role in the evolution of mammalian isochores. *Molecular Biology and Evolution* 17: 1371–1383. PMID: [10958853](https://pubmed.ncbi.nlm.nih.gov/10958853/)
50. Bell JT, Pai AA, Pickrell JK, Gaffney DJ, Pique-Regi R, et al. (2011) DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biol* 12: R10. doi: [10.1186/gb-2011-12-1-r10](https://doi.org/10.1186/gb-2011-12-1-r10) PMID: [21251332](https://pubmed.ncbi.nlm.nih.gov/21251332/)
51. Gibbs JR, van der Brug MP, Hernandez DG, Traynor BJ, Nalls MA, et al. (2010) Abundant quantitative trait loci exist for DNA methylation and gene expression in human brain. *PLoS Genet* 6: e1000952. doi: [10.1371/journal.pgen.1000952](https://doi.org/10.1371/journal.pgen.1000952) PMID: [20485568](https://pubmed.ncbi.nlm.nih.gov/20485568/)
52. Bird AP (1980) DNA methylation and the frequency of CpG in animal DNA. *Nucleic Acids Research* 8: 1499–1504. PMID: [6253938](https://pubmed.ncbi.nlm.nih.gov/6253938/)
53. Ernst J, Kellis M (2012) ChromHMM: automating chromatin-state discovery and characterization. *Nature Methods* 9: 215–216. doi: [10.1038/nmeth.1906](https://doi.org/10.1038/nmeth.1906) PMID: [22373907](https://pubmed.ncbi.nlm.nih.gov/22373907/)

54. Arbeithuber B, Betancourt AJ, Ebner T, Tiemann-Boege I (2015) Crossovers are associated with mutation and biased gene conversion at recombination hotspots. *Proceedings of the National Academy of Sciences* 112: 2109–2114.
55. Yang S, Wang L, Huang J, Zhang X, Yuan Y, et al. (2015) Parent-progeny sequencing indicates higher mutation rates in heterozygotes. *Nature*.
56. Charlesworth D, Charlesworth B, Morgan M (1995) The pattern of neutral molecular variation under the background selection model. *Genetics* 141: 1619–1632. PMID: [8601499](#)
57. Przeworski M, Charlesworth B, Wall JD (1999) Genealogies and weak purifying selection. *Molecular Biology and Evolution* 16: 246–252. PMID: [10084898](#)
58. Good BH, Walczak AM, Neher RA, Desai MM (2014) Genetic diversity in the interference selection limit. *PLoS Genet* 10: e1004222. doi: [10.1371/journal.pgen.1004222](#) PMID: [24675740](#)
59. Jukes TH, Cantor CR (1969) Evolution of protein molecules. *Mammalian Protein Metabolism* 3: 21–132.
60. Johnson PL, Hellmann I (2011) Mutation rate distribution inferred from coincident SNPs and coincident substitutions. *Genome biology and evolution* 3: 842–850. doi: [10.1093/gbe/evr044](#) PMID: [21572094](#)
61. Seplyarskiy VB, Kharchenko P, Kondrashov AS, Bazykin GA (2012) Heterogeneity of the transition/transversion ratio in *Drosophila* and Hominidae genomes. *Genome Biology and Evolution* 29: 1943–1955.
62. Eyre-Walker A, Eyre-Walker YC (2014) How much of the variation in the mutation rate along the human genome can be explained? *G3: Genes|Genomes|Genetics* 4: 1667–1670. doi: [10.1534/g3.114.012849](#) PMID: [24996580](#)
63. Hodgkinson A, Ladoukakis E, Eyre-Walker A (2009) Cryptic variation in the human mutation rate. *PLoS Biol* 7: e1000027. doi: [10.1371/journal.pbio.1000027](#) PMID: [19192947](#)
64. Schraiber JG, Akey JM (2015) Methods and models for unravelling human evolutionary history. *Nature Reviews Genetics*.
65. Bhaskar A, Wang YR, Song YS (2015) Efficient inference of population size histories and locus-specific mutation rates from large-sample genomic variation data. *Genome Research* 25: 268–279. doi: [10.1101/gr.178756.114](#) PMID: [25564017](#)
66. Hwang DG, Green P (2004) Bayesian Markov chain Monte Carlo sequence analysis reveals varying neutral substitution patterns in mammalian evolution. *Proceedings of the National Academy of Sciences of the United States of America* 101: 13994–14001. doi: [10.1073/pnas.0404142101](#) PMID: [15292512](#)
67. Crow JF (2000) The origins, patterns and implications of human spontaneous mutation. *Nature Reviews Genetics* 1: 40–47. doi: [10.1038/35049558](#) PMID: [11262873](#)
68. Moorjani P, Amorim CEG, Arndt PF, Przeworski M (2016) Variation in the molecular clock of primates. *bioRxiv*: 036434.
69. Kim S-H, Elango N, Warden C, Vigoda E, Soojin VY (2006) Heterogeneous genomic molecular clocks in primates. *PLoS Genet* 2: e163. doi: [10.1371/journal.pgen.0020163](#) PMID: [17029560](#)
70. Gao Z, Wyman MJ, Sella G, Przeworski M (2016) Interpreting the dependence of mutation rates on age and time. *PLoS Biol* 14: e1002355. doi: [10.1371/journal.pbio.1002355](#) PMID: [26761240](#)
71. Moorjani P, Gao Z, Przeworski M (2016) Human germline mutation and the erratic molecular clock. *bioRxiv*: 058024.
72. Sayres MAW, Venditti C, Pagel M, Makova KD (2011) Do variations in substitution rates and male mutation bias correlate with life-history traits? A study of 32 mammalian genomes. *Evolution* 65: 2800–2815. doi: [10.1111/j.1558-5646.2011.01337.x](#) PMID: [21967423](#)
73. Amster G, Sella G (2016) Life history effects on the molecular clock of autosomes and sex chromosomes. *Proceedings of the National Academy of Sciences* 113: 1588–1593.
74. Wich S, De Vries H, Ancrenaz M, Perkins L, Shumaker R, et al. (2009) *Orangutans: geographic variation in behavioral ecology and conservation* Oxford University Press, New York: 65–75.
75. Povolotskaya IS, Kondrashov FA (2010) Sequence space and the ongoing expansion of the protein universe. *Nature* 465: 922–926. doi: [10.1038/nature09105](#) PMID: [20485343](#)
76. Callahan B, Neher RA, Bachtrog D, Andolfatto P, Shraiman BI (2011) Correlated evolution of nearby residues in *Drosophilid* proteins. *PLoS Genet*.
77. Li H, Durbin R (2011) Inference of human population history from individual whole-genome sequences. *Nature* 475: 493–496. doi: [10.1038/nature10231](#) PMID: [21753753](#)
78. Cutter AD, WANG GX, Ai H, Peng Y (2012) Influence of finite-sites mutation, population subdivision and sampling schemes on patterns of nucleotide polymorphism for species with molecular hyperdiversity. *Molecular Ecology* 21: 1345–1359. doi: [10.1111/j.1365-294X.2012.05475.x](#) PMID: [22320847](#)

79. Jenkins PA, Song YS (2011) The effect of recurrent mutation on the frequency spectrum of a segregating site and the age of an allele. *Theoretical Population Biology* 80: 158–173. doi: [10.1016/j.tpb.2011.04.001](https://doi.org/10.1016/j.tpb.2011.04.001) PMID: [21550359](https://pubmed.ncbi.nlm.nih.gov/21550359/)
80. Bhaskar A, Kamm JA, Song YS (2012) Approximate sampling formulas for general finite-alleles models of mutation. *Advances in Applied Probability* 44: 408. doi: [10.1239/aap/1339878718](https://doi.org/10.1239/aap/1339878718) PMID: [24634516](https://pubmed.ncbi.nlm.nih.gov/24634516/)
81. Jenkins PA, Mueller JW, Song YS (2014) General triallelic frequency spectrum under demographic models with variable population size. *Genetics* 196: 295–311. doi: [10.1534/genetics.113.158584](https://doi.org/10.1534/genetics.113.158584) PMID: [24214345](https://pubmed.ncbi.nlm.nih.gov/24214345/)
82. Charlesworth B, Jain K (2014) Purifying selection, drift, and reversible mutation with arbitrarily high mutation rates. *Genetics* 198: 1587–1602. doi: [10.1534/genetics.114.167973](https://doi.org/10.1534/genetics.114.167973) PMID: [25230951](https://pubmed.ncbi.nlm.nih.gov/25230951/)
83. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, et al. (2010) A method and server for predicting damaging missense mutations. *Nature Methods* 7: 248–249. doi: [10.1038/nmeth0410-248](https://doi.org/10.1038/nmeth0410-248) PMID: [20354512](https://pubmed.ncbi.nlm.nih.gov/20354512/)
84. Adzhubei I, Jordan DM, Sunyaev SR (2013) Predicting functional effect of human missense mutations using PolyPhen-2. *Current Protocols in Human Genetics*: 7.20. 21–27.20. 41.
85. Ng PC, Henikoff S (2003) SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Research* 31: 3812–3814. PMID: [12824425](https://pubmed.ncbi.nlm.nih.gov/12824425/)
86. Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, et al. (2005) Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Research* 15: 1034–1050. doi: [10.1101/gr.3715005](https://doi.org/10.1101/gr.3715005) PMID: [16024819](https://pubmed.ncbi.nlm.nih.gov/16024819/)
87. Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, et al. (2010) Fine-scale recombination rate differences between sexes, populations and individuals. *Nature* 467: 1099–1103. doi: [10.1038/nature09525](https://doi.org/10.1038/nature09525) PMID: [20981099](https://pubmed.ncbi.nlm.nih.gov/20981099/)
88. Gravel S, Henn BM, Gutenkunst RN, Indap AR, Marth GT, et al. (2011) Demographic history and rare allele sharing among human populations. *Proceedings of the National Academy of Sciences* 108: 11983–11988.
89. Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ, et al. (2005) Calibrating a coalescent simulation of human genome sequence variation. *Genome Research* 15: 1576–1583. doi: [10.1101/gr.3709305](https://doi.org/10.1101/gr.3709305) PMID: [16251467](https://pubmed.ncbi.nlm.nih.gov/16251467/)
90. Gazave E, Ma L, Chang D, Coventry A, Gao F, et al. (2014) Neutral genomic regions refine models of recent rapid human population growth. *Proceedings of the National Academy of Sciences* 111: 757–762.
91. Bhaskar A, Clark AG, Song YS (2014) Distortion of genealogical properties when the sample is very large. *Proceedings of the National Academy of Sciences* 111: 2385–2390.
92. Hudson RR (2002) Generating samples under a Wright—Fisher neutral model of genetic variation. *Bioinformatics* 18: 337–338. PMID: [11847089](https://pubmed.ncbi.nlm.nih.gov/11847089/)
93. Rambaut A, Grass NC (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Computer Applications in the Biosciences* 13: 235–238. PMID: [9183526](https://pubmed.ncbi.nlm.nih.gov/9183526/)