

# Dissecting Genomic Determinants of Positive Selection with an Evolution-Guided Regression Model

Yi-Fei Huang \*<sup>1,2</sup>

<sup>1</sup>Department of Biology, Pennsylvania State University, University Park, PA, USA

<sup>2</sup>Huck Institutes of the Life Sciences, Pennsylvania State University, University Park, PA, USA

\*Corresponding author: E-mail: yuh371@psu.edu.

Associate editor: Jeffrey Townsend

## Abstract

In evolutionary genomics, it is fundamentally important to understand how characteristics of genomic sequences, such as gene expression level, determine the rate of adaptive evolution. While numerous statistical methods, such as the McDonald–Kreitman (MK) test, are available to examine the association between genomic features and the rate of adaptation, we currently lack a statistical approach to disentangle the independent effect of a genomic feature from the effects of other correlated genomic features. To address this problem, I present a novel statistical model, the MK regression, which augments the MK test with a generalized linear model. Analogous to the classical multiple regression model, the MK regression can analyze multiple genomic features simultaneously to infer the independent effect of a genomic feature, holding constant all other genomic features. Using the MK regression, I identify numerous genomic features driving positive selection in chimpanzees. These features include well-known ones, such as local mutation rate, residue exposure level, tissue specificity, and immune genes, as well as new features not previously reported, such as gene expression level and metabolic genes. In particular, I show that highly expressed genes may have a higher adaptation rate than their weakly expressed counterparts, even though a higher expression level may impose stronger negative selection. Also, I show that metabolic genes may have a higher adaptation rate than their nonmetabolic counterparts, possibly due to recent changes in diet in primate evolution. Overall, the MK regression is a powerful approach to elucidate the genomic basis of adaptation.

**Key words:** adaptive evolution, positive selection, McDonald–Kreitman test, statistical inference, causal inference.

## Introduction

Understanding the genetic basis of positive selection is a fundamental problem in evolutionary biology. Numerous statistical approaches have been developed to detect loci under positive selection. A popular framework is codon substitution models that seek to infer positively selected genes solely from interspecies sequence divergence (Goldman and Yang 1994; Muse and Gaut 1994; Yang et al. 2000). By contrasting the rate of nonsynonymous substitutions ( $dN$ ) against the rate of synonymous substitutions ( $dS$ ), codon substitution models can identify positively selected genes with a  $dN/dS$  ratio greater than 1. However, because negative (purifying) selection can dramatically reduce  $dN/dS$  ratios, codon substitution models may be underpowered to detect genes that experienced both positive selection and strong negative selection (Hughes 2007).

Unlike codon substitution models, the McDonald–Kreitman (MK) test utilizes both interspecies divergence and intraspecies polymorphism to elucidate positive selection in a species of interest (McDonald and Kreitman 1991; Fay et al. 2001; Smith and Eyre-Walker 2002). By contrasting the levels of divergence and polymorphism at functional sites

and putatively neutral sites, the MK test seeks to identify positively selected genes that show an excess of interspecies divergence at functional sites. Because highly deleterious mutations can neither segregate nor reach fixation in a population, the MK test is intrinsically robust to the presence of strong negative selection. On the other hand, weak negative selection may lead to biased results in the MK test because mutations under weak selection can segregate in a population but not reach fixation. To address this problem, several recent studies have extended the MK test to account for the effects of weak negative selection on intraspecies polymorphism (Eyre-Walker and Keightley 2009; Messer and Petrov 2013; Galtier 2016; Haller and Messer 2017; Uricchio et al. 2019), ensuring that the inference of positive selection is not biased by the presence of weak selection. Thus, the MK test and its extensions are powerful methods to disentangle positive selection from ubiquitous negative selection.

Because MK-based methods use relatively sparse divergence and polymorphism data from closely related species, they may be underpowered to pinpoint individual genes under positive selection. To boost statistical power, MK-based methods often are applied to a collection of genes or

nucleotide sites with similar genomic features. Using this pooling strategy, previous studies have identified numerous genomic features associated with positive selection in *Drosophila* and primates. The features associated with positive selection in *Drosophila* include local mutation rate (Campos et al. 2014; Castellano et al. 2016; Rousselle et al. 2020), local recombination rate (Marais and Charlesworth 2003; Campos et al. 2014; Castellano et al. 2016), gene expression specificity (Fraïsse et al. 2019), residue exposure to solvent (Moutinho, Trancoso et al. 2019), X linkage (Avila et al. 2015; Campos et al. 2018), and sex-biased expression (Pröschel et al. 2006; Avila et al. 2015; Campos et al. 2018). The features associated with positive selection in primates include protein disorder (Afanasyeva et al. 2018), virus–host interaction (Enard et al. 2016; Uricchio et al. 2019), protein–protein interaction (PPI) degree (Luisi et al. 2015), and X linkage (Hvilsom et al. 2012).

While existing MK-based methods can identify genomic features associated with the signatures of positive selection, they may not be able to distinguish genomic features *independently affecting* the rate of adaptive evolution from spurious features *without independent effects* on adaptation (Moutinho, Trancoso et al. 2019; Fraïsse et al. 2019). For instance, MK-based methods often are applied to one genomic feature at a time. If a genomic feature with an independent effect on the rate of adaptive evolution is strongly correlated with a second feature without an independent effect, MK-based methods may report a spurious association between the second feature and adaptive evolution.

Before the current study, two simple heuristic methods have been previously used to estimate the independent effect of a genomic feature on the rate of adaptation by controlling for other potentially correlated features. If we are interested in estimating the independent effect of a gene-level feature, such as tissue specificity, we may estimate the rate of adaptation at the gene level and then fit a standard linear regression model, in which we treat the feature of interest and correlated genomic features as covariates and treat the gene-level rate of adaptation as a response variable (Luisi et al. 2015; Castellano et al. 2016; Moutinho, Trancoso et al. 2019; Fraïsse et al. 2019). The regression coefficient associated with the feature of interest can be interpreted as its independent effect on positive selection, holding constant all other genomic features. Although this strategy is powerful and elegant, it cannot be applied to species with low levels of polymorphism, such as primates, due to the challenge of estimating the rate of adaptation at the gene level. Alternatively, we may first stratify genes into a “treatment” group and a “control” group based on the genomic feature of interest. Then, we may use statistical matching algorithms to match each gene from the “treatment” group with a gene of similar characteristics from the “control” group. A significant difference in the rate of adaptation between the two groups of matched genes indicates that the feature of interest has an independent effect on positive selection. Although this method has been successfully used in previous studies (Enard et al. 2016; Campos et al. 2018; Castellano et al. 2019), it is difficult to match genes when there are a large number of genomic

features to control for. Therefore, we currently lack a general and powerful statistical framework to estimate the independent effects of genomic features on positive selection by adjusting for a large number of correlated genomic features.

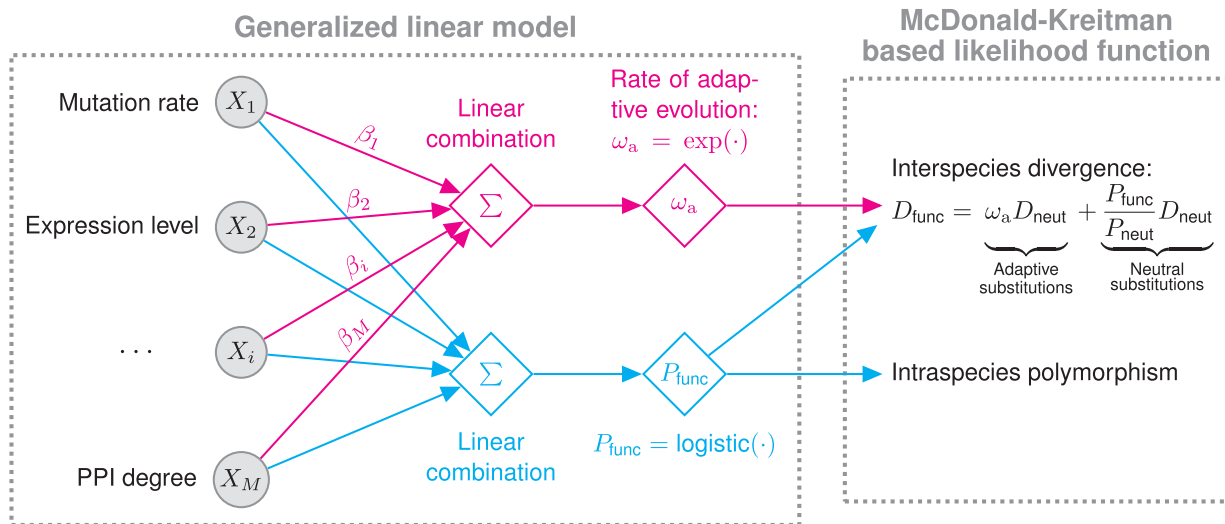
In the current study, I present a novel statistical method, the MK regression, to estimate the independent effects of genomic features on the rate of adaptive evolution. The MK regression is a hybrid of the MK test and the generalized linear regression. Unlike standard linear regression models and statistical matching algorithms, the MK regression can control for a large number of correlated genomic features and is applicable to species with a low level of polymorphism. To the best of my knowledge, the MK regression is the first evolutionary model tailored to characterize the independent effects of genomic features on adaptive evolution. Using synthetic data, I show that the MK regression can unbiasedly estimate the independent effect of a genomic feature even when it is strongly correlated with another genomic feature. Applying the MK regression to polymorphism and divergence data in the chimpanzee lineage, I corroborate previous findings that local mutation rate, residue exposure level, tissue specificity, and immune system genes are key determinants of positive selection in protein-coding genes. In addition, I show that highly expressed genes and metabolic genes may have a higher rate of adaptive evolution than other genes after controlling for several correlated genomic features, which has not been widely reported in previous studies. Taken together, the MK regression is a valuable addition to evolutionary biologists’ arsenal for investigating the genetic basis of adaptation.

## Results

### The MK Regression is a Generalized Linear Model Tailored to Estimate the Effects of Genomic Features on Positive Selection

The key idea behind the MK regression is to model the site-wise rate of adaptive evolution as a linear combination of local genomic features (fig. 1). I use  $\omega_a$ , the relative rate of adaptive substitutions at a functional nucleotide site with respect to the average substitution rate at neutral nucleotide sites, as a measure of the rate of adaptation (Booker et al. 2017; Moutinho, Bataillon et al. 2019). Unlike previous MK-based models that treat  $\omega_a$  as a gene-level measure, I treat  $\omega_a$  as a measure of adaptive evolution at an individual nucleotide site and assume that it can be predicted from local genomic features. To integrate the effects of multiple features on the rate of adaptive evolution, I assume that  $\omega_a$ , in a site-wise manner, is a linear combination of local genomic features, such as local mutation rate, local recombination rate, and gene expression level. For each genomic feature, the MK regression seeks to estimate a regression coefficient indicating its independent effect on the rate of adaptation, holding constant all other genomic features.

Specifically, the MK regression consists of two components: a generalized linear model and an MK-based likelihood function (fig. 1). First, I assume that  $\omega_a$ , in a site-wise manner, is a linear combination of local genomic features followed by an exponential transformation,



**Fig. 1.** Schematic of the MK regression. The MK regression consists of two components: a generalized linear model and a McDonald–Kreitman-based likelihood function. First, I assume that, in a site-wise manner, the rate of adaptive evolution ( $\omega_a$ ) at a functional site is a linear combination of local genomic features followed by an exponential transformation, in which regression coefficient  $\beta_i$  indicates the effect of the  $i$ th feature on adaptive evolution. Similarly, I assume that the probability of observing a SNP ( $P_{\text{func}}$ ) at the same functional site is another linear combination of the same set of genomic features, followed by a logistic transformation. Second, in the McDonald–Kreitman-based likelihood function, I combine  $\omega_a$  and  $P_{\text{func}}$  at every functional site with two neutral parameters,  $D_{\text{neut}}$  and  $P_{\text{neut}}$ , to calculate the probability of observed divergence and polymorphism data given model parameters.  $D_{\text{neut}}$  and  $P_{\text{neut}}$  denote the expected number of substitutions and the probability of observing a SNP at a neutral site, respectively.  $D_{\text{func}}$  denotes the expected number of substitutions at a functional site.

$$\omega_a = \exp(\beta_0 + \beta_1 X_1 + \cdots + \beta_i X_i + \cdots + \beta_M X_M).$$

In this equation,  $X_i$  is the  $i$ th feature at a functional nucleotide site;  $\beta_0$  is an intercept indicating the baseline rate of adaptive evolution when all genomic features are equal to 0;  $\beta_i$  is a regression coefficient indicating the  $i$ th feature's effect on the rate of adaptation;  $M$  is the total number of genomic features;  $\exp$  is an exponential inverse link function which ensures  $\omega_a$  is positive. If  $\beta_i$  is statistically different from 0, I consider that feature  $i$  may have an independent effect on adaptation after adjusting for the other features. Similarly, to accommodate the effects of local genomic features on polymorphism data, I assume that the probability of observing intraspecies polymorphism at a functional nucleotide site,  $P_{\text{func}}$ , is another linear combination of local genomic features followed by a logistic transformation (fig. 1).

Second, in the component of MK-based likelihood function (fig. 1), I combine  $\omega_a$  and  $P_{\text{func}}$  at every functional site with two neutral parameters to calculate the probability of observed polymorphism and divergence data at both functional and neutral sites, which allows for a maximum likelihood estimation of model parameters. Finally, I use the Wald test to examine whether the estimated regression coefficient,  $\hat{\beta}_i$ , is significantly different from 0 for each feature  $i$ . It is worth noting that, unlike the standard linear regression, the MK regression does not assume that response variables, that is, polymorphism and divergence data, follow a normal distribution. Instead, the likelihood function of the MK regression uses the Jukes–Cantor substitution model (Jukes and Cantor 1969) and the Bernoulli distribution to describe the generation of divergence and

polymorphism in a site-wise manner. Thus, the MK regression can naturally describe the evolution of functional and neutral sites.

### Joint Analysis of Multiple Features Distinguishes Independent Effects from Spurious Associations

I conducted two simulation experiments to assess the MK regression's validity and its power to infer the independent effects of genomic features on the rate of adaptive evolution. The simulation experiments consisted of two steps. First, I randomly sampled genomic features from a bivariate normal distribution at each functional site. Second, I generated synthetic polymorphism and divergence data at both functional and neutral sites based on the MK regression model.

In the first simulation experiment, I assumed that there were two genomic features of interest. The first genomic feature had an independent effect on the rate of adaptive evolution, and its regression coefficient,  $\beta_1$ , was equal to 1. On the other hand, the second feature had no independent effect on selection. Thus, its regression coefficient,  $\beta_2$ , was equal to 0 by definition. The other parameters required for the simulation experiment were chosen to ensure that genome-wide levels of polymorphism and divergence are comparable between synthetic data and empirical data from chimpanzees (see details in the Materials and Methods section). To systematically assess the MK regression's performance with respect to various degrees of correlation between genomic features, I generated four sets of synthetic data with different correlation coefficients between features (0.0, 0.2, 0.4, and 0.6). In each synthetic data set, I generated 10 independent replicates each of which consisted of 10 Mb functional sites and 10 Mb neutral sites.

I applied two different versions of the MK regression to the synthetic data. The first version was the simple MK regression that analyzed one genomic feature at a time, which was designed to mimic previous MK-based methods. The second one was the multiple MK regression that analyzed two features simultaneously. As shown in [figs. 2A and 2B](#), both the simple MK regression and the multiple MK regression produced unbiased estimates of regression coefficients when there was no correlation between features. However, the simple MK regression frequently estimated that  $\hat{\beta}_2$  was positive when the two features were correlated with each other, whereas the true value of  $\beta_2$  was equal to 0. On the other hand, the multiple MK regression always produced unbiased estimates of regression coefficients regardless of the degree of correlation between features.

In the second simulation experiment, I evaluated the extent to which the correlation between two causal features complicates the estimation of their independent effects. I set the regression coefficients of the two features to  $\beta_1 = 1$  and  $\beta_2 = -0.2$ , respectively. Then, I followed the same procedure described in the first simulation experiment to generate synthetic data. As shown in [figs. 2C and 2D](#), the multiple MK regression accurately estimated regression coefficients without any noticeable bias, whereas the simple MK regression produced biased results when the two features were correlated with each other. Importantly, when the correlation was strong, the simple MK regression estimated that  $\hat{\beta}_2$  was positive while the true value of  $\beta_2$  was equal to  $-0.2$ .

Furthermore, using the same synthetic data, I evaluated the performance of a previous MK-based method ([Smith and Eyre-Walker 2002](#); [Fraïsse et al. 2019](#)), which can only analyze one feature at a time. For each genomic feature, I stratified functional sites into two equal-sized groups. The first group included the top half of functional sites with higher feature value, whereas the second group included the bottom half of functional sites with lower feature value. I estimated  $\omega_a$  for each group separately. Then, I calculated  $\Delta\omega_a$ , that is, the difference in  $\omega_a$  between the two groups of functional sites. If the previous MK-based method can unbiasedly estimate the effects of genomic features, the sign of  $\Delta\omega_a$  should match the sign of the true regression coefficient. However, the previous MK-based method frequently produced wrong estimates of the sign of  $\Delta\omega_a$  when genomic features were strongly correlated with each other ([supplementary fig. 1, Supplementary Material online](#)). In summary, it is critical to jointly analyze multiple genomic features for an unbiased estimation of their independent effects on adaptive evolution.

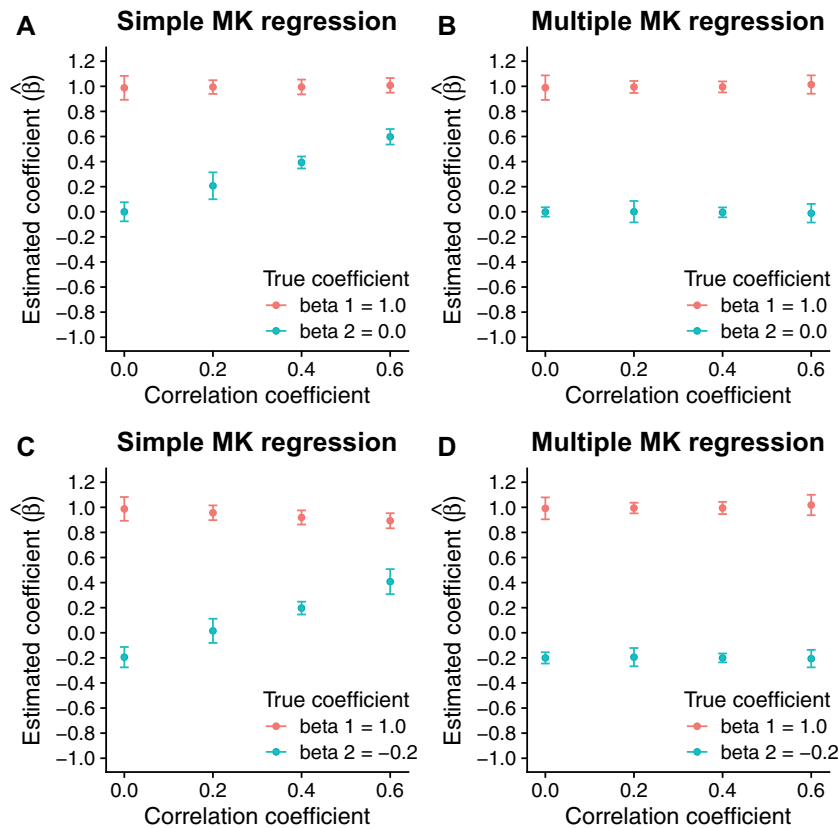
### The Multiple MK Regression Elucidates Genomic Determinants of Positive Selection in Chimpanzees

I investigated positive selection in chimpanzee autosomal genes using the MK regression. Because gene annotations were of high quality in the human genome, I converted 4-fold degenerate (4D) and 0-fold degenerate (0D) sites annotated in dbNSFP ([Liu et al. 2013, 2016](#)) from the human genome to the chimpanzee genome. Because all point mutations at 4D sites are synonymous, I assume that they are putatively neutral. On the other hand, because all point

mutations at 0D sites are nonsynonymous, I assume that they are potentially functional. I obtained a genome-wide map of single nucleotide polymorphisms (SNPs) in 18 central chimpanzee (*Pan troglodytes troglodytes*) individuals ([de Manuel et al. 2016](#)), and inferred ancestral alleles using a reconstructed chimpanzee ancestral genome ([Herrero et al. 2016](#); [Yates et al. 2020](#)). Because the SNP data set consisted of samples from both females and males, the number of sampled sequences was different between autosomes and sex chromosomes. Thus, I retained only autosomal genes for downstream analysis. To mitigate the impact of weak negative selection on the inference of positive selection, I filtered out SNPs with a derived allele frequency lower than 50% for downstream analysis. In addition, I reconstructed fixed substitutions at 4D and 0D sites in the chimpanzee lineage by comparing the reconstructed ancestral genome with the chimpanzee reference genome. I estimated that the proportion of adaptive nonsynonymous substitutions ( $\alpha$ ) was equal to 15.3% in chimpanzee autosomal genes, which is similar to the estimate in a previous study ([Tataru et al. 2017](#)).

I collected six genomic features in chimpanzee autosomal genes, including local mutation rate, local recombination rate, residue exposure level, gene expression level, tissue specificity, and the number of unique protein–protein interaction partners per gene (PPI degree). Specifically, I obtained a chimpanzee-based map of local recombination rates from a previous study ([Auton et al. 2012](#)) and constructed a map of local mutation rates using putatively neutral substitutions in the chimpanzee lineage. Because functional genomic data were more complete and of higher quality in humans than in chimpanzees, I obtained tissue-based gene expression data from the Human Protein Atlas ([Uhlen et al. 2015](#)) and utilized the expression level averaged across all tissues and a summary statistic, tau ([Yanai et al. 2005](#)), as measures of gene expression level and tissue specificity, respectively. I also obtained predicted levels of residue exposure to solvent and experimentally determined PPI degrees in the human genome from previous studies ([Wong et al. 2011](#); [Luck et al. 2020](#)). I converted human-based annotations of gene expression level, tissue specificity, residue exposure level, and PPI degree to the chimpanzee genome (panTro4) using liftOver ([Haeussler et al. 2019](#)).

I first employed the simple MK regression to analyze the effect of one genomic feature at a time, with no attempt to distinguish independent effects from spurious associations. Because the MK regression used a logarithmic link function for  $\omega_a$ , I explored if a logarithmic transformation of genomic features can improve model fitting. I found that the logarithmic transformation improved the fitting of the simple MK regression for all features but tissue specificity ([supplementary table 1, Supplementary Material online](#)). Therefore, I applied the logarithmic transformation to all features except tissue specificity throughout this study. In the simple MK regression, the regression coefficients of local mutation rate, residue exposure level, and PPI degree were significantly higher than 0, whereas the regression coefficient of gene expression level was significantly lower than 0 ([fig. 3A and supplementary table 2, Supplementary Material online](#)). On the other hand, local



**Fig. 2.** Simulation results. (A) Estimates of regression coefficients in the simple MK regression. The true coefficients are  $\beta_1 = 1$  and  $\beta_2 = 0$ . (B) Estimates of regression coefficients in the multiple MK regression. The true coefficients are  $\beta_1 = 1$  and  $\beta_2 = 0$ . (C) Estimates of regression coefficients in the simple MK regression. The true coefficients are  $\beta_1 = 1$  and  $\beta_2 = -0.2$ . (D) Estimates of regression coefficients in the multiple MK regression. The true coefficients are  $\beta_1 = 1$  and  $\beta_2 = -0.2$ . In each plot, dots and error bars indicate the means and the 2-fold standard deviations of estimated coefficients across 10 independent replicates, respectively.

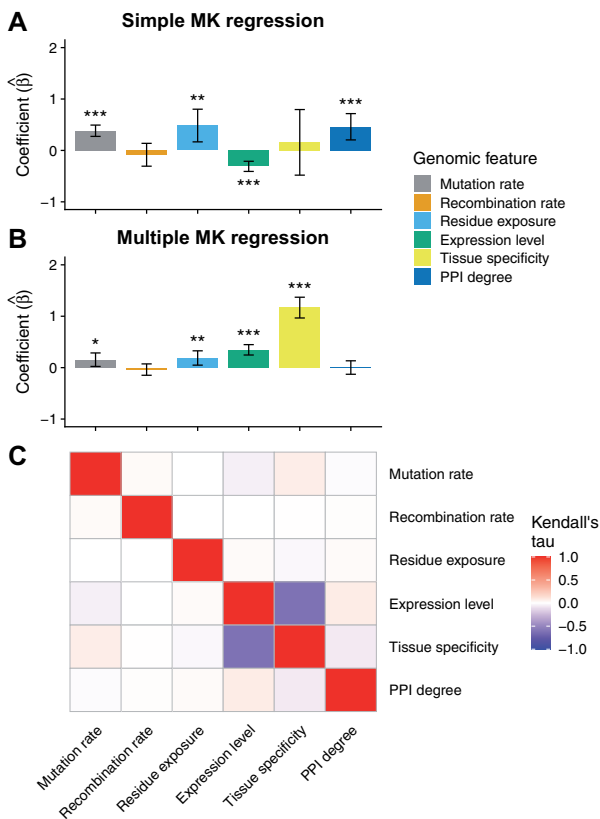
recombination rate and tissue specificity were not significantly associated with the rate of adaptive evolution in the simple MK regression (fig. 3A and supplementary table 2, Supplementary Material online).

As discussed in the simulation experiments, the simple MK regression may produce biased estimates of regression coefficients if genomic features are correlated with each other. To test if this was the case in the chimpanzee data, I used the multiple MK regression to analyze the effects of the six genomic features simultaneously. Surprisingly, while local mutation rate, local recombination rate, and residue exposure level showed similar effects in the multiple MK regression, the regression coefficients of the other features were different between the multiple MK regression and the simple MK regression (fig. 3B and supplementary table 3, Supplementary Material online). Specifically, the regression coefficient of PPI degree was not significant in the multiple MK regression ( $\hat{\beta} = 0.002$ ;  $P$ -value = 0.975), whereas the same coefficient was significant in the simple MK regression ( $\hat{\beta} = 0.460$ ;  $P$ -value =  $3.178 \times 10^{-4}$ ). The coefficient of gene expression level was significantly higher than 0 in the multiple MK regression ( $\hat{\beta} = 0.347$ ;  $P$ -value =  $6.674 \times 10^{-12}$ ), whereas the same coefficient was negative in the simple MK regression ( $\hat{\beta} = -0.310$ ;  $P$ -value =  $3.534 \times 10^{-10}$ ). Also, the regression coefficient of

tissue specificity was significantly higher than 0 in the multiple MK regression ( $\hat{\beta} = 1.168$ ;  $P$ -value =  $6.394 \times 10^{-31}$ ) but not in the simple MK regression ( $\hat{\beta} = 0.157$ ;  $P$ -value = 0.622).

To examine whether these results were robust to different metrics of tissue specificity, I utilized the negative value of Hg (Kryuchkova-Mostacci and Robinson-Rechavi 2017) as an alternative metric of tissue specificity. Similar to tau, a higher value of negative Hg indicates a higher level of tissue specificity. I observed qualitatively similar regression coefficients when I replaced tau with negative Hg in the multiple MK regression (supplementary fig. 2, Supplementary Material online), although the regression coefficient of local mutation rate was not statistically significant when negative Hg was used. Thus, the estimated effects of genomic features may be robust to different metrics of tissue specificity.

To investigate whether correlations between genomic features could explain the differences in estimated coefficients between the multiple and the simple MK regression, I calculated the Kendall rank correlation coefficient for all pairs of genomic features (fig. 3C and supplementary table 4, Supplementary Material online). I found that local mutation rate, local recombination rate, and residue exposure level were weakly correlated with other features, which may explain why the regression coefficients of these features were



**Fig. 3.** Effects of genomic features on the rate of adaptive evolution. (A) Estimated coefficients of genomic features in the simple MK regression. (B) Estimated coefficients of genomic features in the multiple MK regression. In each bar plot, error bars indicate 95% confidence intervals while one, two, and three asterisks indicate  $0.01 \leq P\text{-value} < 0.05$ ,  $0.001 \leq P\text{-value} < 0.01$ , and  $P\text{-value} < 0.001$ , respectively. (C) Correlations between genomic features.

consistent between the multiple and the simple MK regression. In contrast, gene expression level and tissue specificity showed a strong negative correlation, which may cause spurious associations in the simple MK regression. I also found that PPI degree was correlated with gene expression level and tissue specificity, although the correlations were to a lesser extent compared with the correlation between gene expression level and tissue specificity. Thus, the observed association of PPI degree with positive selection in the simple MK regression could be due to its correlation with gene expression level and/or tissue specificity.

### Statistical Matching Analysis Confirms Genomic Determinants Identified by the Multiple MK Regression

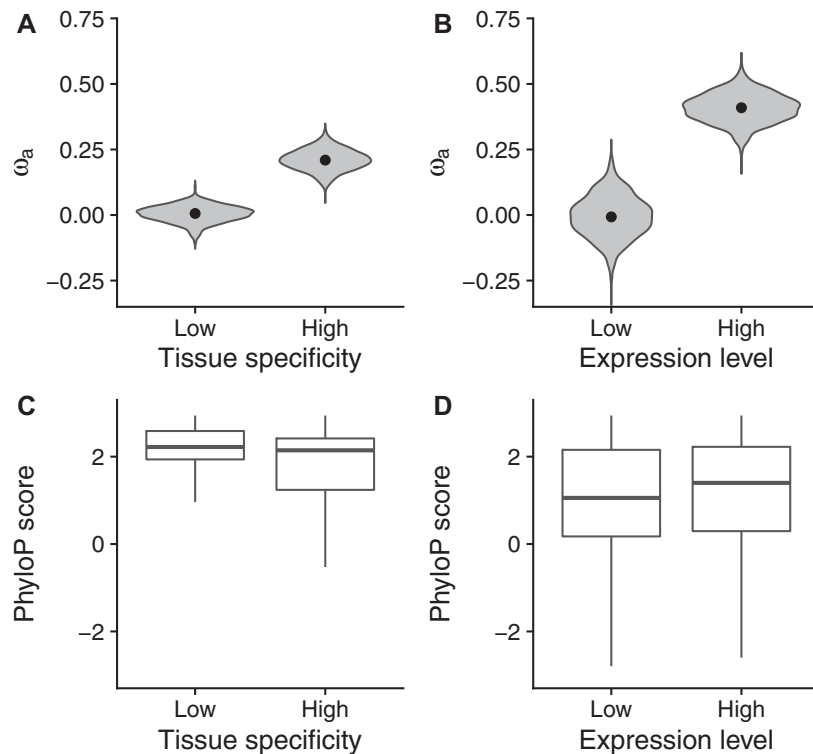
I used statistical matching algorithms to corroborate the results of the multiple MK regression. First, I verified whether the association between PPI degree and the rate of adaptive evolution could be explained away by controlling for gene expression level and tissue specificity. I stratified protein-coding genes into two groups with different PPI degrees, in which 1,556 genes with at least 10 protein interaction partners were assigned to the high PPI-degree group whereas

8,471 genes with no more than one interaction partner were assigned to the low PPI-degree group. Without controlling for gene expression level and tissue specificity, the high PPI-degree group had a higher  $\omega_a$  than the low PPI-degree group (supplementary fig. 3A, Supplementary Material online), but the difference in  $\omega_a$  was not significant ( $P\text{-value} = 0.146$ ; two-tailed permutation test), possibly due to a reduction of sample size in the stratified analysis. Then, using the default propensity score matching algorithm in MatchIt (Ho et al. 2011), I matched each gene from the high PPI-degree group with a gene of similar expression level and tissue specificity from the low PPI-degree group. In the matched data,  $\omega_a$  was not different between the high-PPI and low-PPI groups (supplementary fig. 3B, Supplementary Material online). I observed similar results using two alternative cutoffs, 5 and 20, for the high PPI-degree group (supplementary fig. 4, Supplementary Material online). Thus, PPI degree is unlikely to be a genomic determinant of positive selection in chimpanzees.

I also verified the effect of tissue specificity on the rate of adaptation after adjusting for gene expression level. Due to the strong negative correlation between expression level and tissue specificity (supplementary fig. 5, Supplementary Material online), MatchIt returned few matched genes when I attempted to control for expression level. Therefore, I implemented a different matching approach. By closely examining the relationship between expression level and tissue specificity, I found that the variation in tissue specificity was high among highly expressed genes (supplementary fig. 5, Supplementary Material online). Thus, I stratified 737 highly expressed genes (gene expression level  $> 30$ ) into two equal-sized groups based on the ranking of their tissue specificity. As shown in fig. 4A, highly expressed genes with high tissue specificity had a significantly higher  $\omega_a$  than their counterparts with low tissue specificity ( $P\text{-value} = 0.001$ ; two-tailed permutation test). Therefore, my gene matching analysis confirms the positive effect of tissue specificity on the rate of adaptation in the multiple MK regression.

As an alternative analysis to infer the independent effect of tissue specificity after controlling for expression level, I divided genes into 10 equal-sized groups (deciles) based on their expression levels. In each decile, I further divided genes into two equal-sized subgroups based on the ranking of their tissue specificity within the decile. As expected, the subgroup with high tissue specificity had a significantly higher  $\omega_a$  than the subgroup with low tissue specificity in the decile with the highest expression level (supplementary table 5, Supplementary Material online), which confirms the positive effect of tissue specificity on adaptation rate in highly expressed genes. The difference in  $\omega_a$  was not significant in other deciles, possibly due to the low variation in tissue specificity in lowly expressed genes (supplementary fig. 5, Supplementary Material online).

I observed that the variation in expression level was high among genes with high tissue specificity (supplementary fig. 5, Supplementary Material online). To verify the positive effect of gene expression level after adjusting for tissue specificity, I stratified 993 tissue-specific genes ( $\tau > 0.85$ ) into two



**Fig. 4.** Statistical matching analysis. (A) Estimates of  $\omega_a$  in 369 highly expressed genes with low tissue specificity and 368 highly expressed genes with high tissue specificity. (B) Estimates of  $\omega_a$  in 498 tissue-specific genes with a low expression level and 495 tissue-specific genes with a high expression level. In each violin plot, dots indicate point estimates of  $\omega_a$  while violins depict the distributions of  $\omega_a$  from a gene-based bootstrapping analysis with 1,000 resamplings. (C) Distributions of phyloP scores in 369 highly expressed genes with low tissue specificity and 368 highly expressed genes with high tissue specificity. (D) Distributions of phyloP scores in 498 tissue-specific genes with a low expression level and 495 tissue-specific genes with a high expression level. In each box plot, the bottom, the top, and the middle horizontal bar of the box indicate the first quartile, the third quartile, and the median of phyloP scores, respectively. The whiskers indicate 1.5-fold interquartile ranges.

approximately equal-sized groups based on the ranking of their expression levels. The first group consisted of the top 495 tissue-specific genes with higher expression level, whereas the second group consisted of the bottom 498 tissue-specific genes with lower expression level. The mean expression levels were equal to 5.905 and 0.368 in the two gene groups, which corresponds to a 16-fold difference in mean expression level. As shown in [fig. 4B](#), tissue-specific genes with a high expression level had a significantly higher  $\omega_a$  than their lowly expressed counterparts ( $P$ -value = 0.002; two-tailed permutation test), which confirms the positive effect of gene expression level on the rate of adaptive evolution in the multiple MK regression.

As an alternative analysis to infer the independent effect of gene expression level after controlling for tissue specificity, I divided genes into 10 deciles based on their tissue specificity. In each decile, I further divided genes into two equal-sized subgroups based on the ranking of their expression levels within the decile. As expected, I observed that the subgroup with a high expression level had a significantly higher  $\omega_a$  than its counterpart in the decile with the highest tissue specificity ([supplementary table 6, Supplementary Material online](#)). To a lesser extent,  $\omega_a$  was slightly lower in the subgroup with a high expression level than the subgroup with a low expression level in the 7th decile ([supplementary table 6, Supplementary Material online](#)). The difference in  $\omega_a$  was not statistically

significant in other deciles, possibly due to the low variation in expression level among genes with low tissue specificity ([supplementary fig. 5, Supplementary Material online](#)). On average, highly expressed genes showed a higher rate of adaptive evolution than their lowly expressed counterparts.

Also, I used phyloP scores ([Pollard et al. 2010; Hubisz et al. 2011](#)) to examine the effects of gene expression level and tissue specificity on the rate of protein evolution. After controlling for gene expression level, phyloP scores increased with decreasing tissue specificity ([fig. 4C](#)), which is in line with the observation that housekeeping genes tend to evolve at a lower substitution rate than tissue-specific genes ([Zhang and Li 2004; Zhu et al. 2008](#)). On the other hand, after controlling for tissue specificity, phyloP scores increased with increasing expression level ([fig. 4D](#)), which is in line with stronger purifying selection on highly expressed genes ([Zhang and Yang 2015](#)). Taken together, it seems that highly expressed genes may be subject to more frequent positive selection than their lowly expressed counterparts, although a higher expression level may impose stronger purifying selection and reduce the overall rate of protein evolution.

#### The Rate of Adaptive Evolution May Also Increase with Gene Expression Level in *Drosophila*

In the previous section, I showed that the rate of adaptive evolution may increase with increasing gene expression level

in chimpanzees. Recently, Fraïsse et al. (2019) examined the same problem in *Drosophila melanogaster*. Fraïsse et al. (2019) first estimated  $\omega_a$  for each gene separately. Then, they regressed the gene-level  $\omega_a$  on expression level and other potentially correlated genomic features, such as tissue specificity, using the standard linear regression. The coefficients of the standard linear regression were interpreted as the independent effects of genomic features after controlling for other correlated features. Unlike the current study, Fraïsse et al. (2019) observed that the rate of adaptation might decrease with increasing expression level in *D. melanogaster*.

To reconcile the discrepancy between the current study and Fraïsse et al. (2019), I reanalyzed the data from Fraïsse et al. (2019). Following Fraïsse et al. (2019), I used the standard linear regression to regress the gene-level estimate of  $\omega_a$  on gene expression level and tissue specificity, and observed that the coefficient of gene expression was negative ( $-0.041384$ ;  $P$ -value =  $8.93 \times 10^{-9}$ ). However, the standard linear regression may suffer from two critical problems in this data set. First, the residuals of the standard linear regression did not follow a normal distribution, as suggested by a quantile–quantile plot (supplementary fig. 6A, Supplementary Material online). Second, the majority of genes had less than 4 nonsynonymous polymorphisms in this data set (supplementary fig. 6B, Supplementary Material online), so the gene-level estimate of  $\omega_a$  may be highly inaccurate. Thus, I argue that the standard linear regression may not be an appropriate statistical method for analyzing this data set.

On the other hand, gene expression level had a positive regression coefficient in the multiple MK regression (supplementary fig. 7A, Supplementary Material online). In an orthogonal statistical matching analysis, I stratified 681 *Drosophila* tissue-specific genes ( $\tau > 0.85$ ) into two gene groups based on the ranking of their expression levels. The first group consisted of the top 340 tissue-specific genes with higher expression level, whereas the second group consisted of the bottom 341 tissue-specific genes with lower expression level. Again, tissue-specific genes with higher expression level had a higher  $\omega_a$  than tissue-specific genes with lower expression level (supplementary fig. 7B, Supplementary Material online) despite that the difference in  $\omega_a$  was marginally significant ( $P$ -value = 0.067; two-tailed permutation test). Taken together, the rate of adaptation may also increase with increasing gene expression level in *D. melanogaster*.

### Metabolic and Immune Genes May Be under Frequent Positive Selection in Chimpanzees

To explore whether the positive effect of gene expression level on the rate of adaptive evolution can be explained by the functions of highly expressed genes, I examined the enrichment of Reactome pathways (Jassal et al. 2020) and tissue types (Uhlen et al. 2015) in the aforementioned 495 chimpanzee tissue-specific genes with a high expression level, using the 498 chimpanzee tissue-specific genes with a low expression level as a background set. Surprisingly, genes associated with the metabolism of proteins and lipids, and genes with enriched expression in the intestine, liver, and pancreas, showed a strong enrichment in the 495 tissue-specific genes

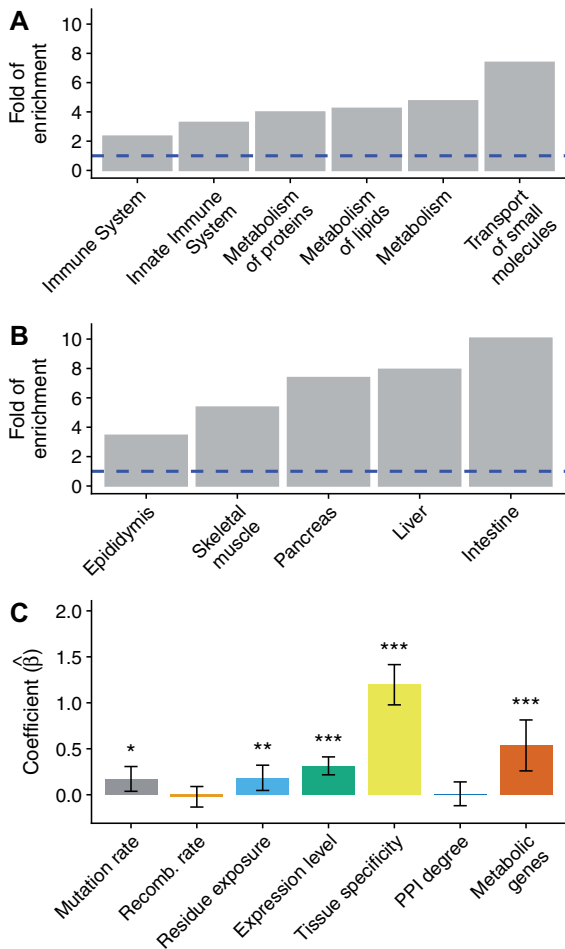
with a high expression level (fig. 5A and 5B; false-discovery rate  $< 0.01$ ). Based on these results, I hypothesized that metabolic genes may be subject to more frequent positive selection than nonmetabolic genes in chimpanzees.

To test this hypothesis, I constructed a genomic feature indicating whether each OD site was located in one of the 2,220 metabolic genes from MSigDB (Subramanian et al. 2005; Liberzon et al. 2011). Then, I used the multiple MK regression to simultaneously estimate the effects of the new feature and the six original features. As shown in fig. 5C and supplementary table 7, Supplementary Material online, the regression coefficient of metabolic genes was significantly higher than 0 in the multiple MK regression ( $\hat{\beta} = 0.537$ ;  $P$ -value =  $1.075 \times 10^{-4}$ ), suggesting that metabolic genes might have a higher rate of adaptation than their nonmetabolic counterparts. Interestingly, the effect of metabolic genes was not significant in the simple MK regression ( $\hat{\beta} = 0.276$ ;  $P$ -value = 0.516). Therefore, controlling for potentially correlated features is critical for revealing elevated positive selection in metabolic genes. Metabolic genes may partially explain the positive effect of gene expression level on adaptive evolution because the regression coefficient of gene expression level reduced moderately from 0.347 to 0.314 after adding metabolic genes as a new feature in the multiple MK regression (supplementary tables 3 and 7, Supplementary Material online).

I observed similar results using propensity score matching. Specifically, I observed that residue exposure level, gene expression level, and tissue specificity showed different distributions between metabolic and nonmetabolic genes (supplementary fig. 8, Supplementary Material online). Before controlling for these correlated genomic features,  $\omega_a$  was similar between metabolic and nonmetabolic genes (supplementary fig. 8, Supplementary Material online). On the other hand,  $\omega_a$  was more than two times higher in metabolic genes than in nonmetabolic genes (0.0550 vs. 0.0234) after controlling for residue exposure level, gene expression level, and tissue specificity (supplementary fig. 8, Supplementary Material online) despite that the difference was not statistically significant due to reduced sample size ( $P$ -value = 0.390; two-tailed permutation test).

To test whether the effect of metabolic genes on adaptation rate can be explained away by their biased expression in digestive organs, I obtained 2,274 genes with biased or enriched expression in digestive organs, including intestine, liver, pancreas, salivary gland, and stomach, from the Human Protein Atlas (Uhlen et al. 2015). As expected, metabolic genes were more likely to have a biased expression in digestive organs than nonmetabolic genes (odds ratio = 2.177;  $P$ -value  $< 2.2 \times 10^{-16}$ ; Fisher's exact test). Then, I constructed a genomic feature indicating whether each OD site was located in the 2,274 genes with biased expression. After adding this genomic feature to the multiple MK regression analysis (supplementary fig. 9, Supplementary Material online), I observed that the regression coefficient of metabolic genes was still significantly higher than 0 ( $\hat{\beta} = 0.475$ ;  $P$ -value =  $8.780 \times 10^{-4}$ ), whereas the regression coefficient of digestive system-biased genes was only marginally significant ( $P$ -value = 0.028). Thus, the effect of metabolic genes on the





**Fig. 5.** Positive selection in metabolic genes. (A) Enrichment of Reactome pathways in 495 tissue-specific genes with a high expression level. (B) Enrichment of tissue types in 495 tissue-specific genes with a high expression level. In each enrichment test, 498 tissue-specific genes with a low expression level are used as a background gene set. A dashed blue line indicates that the fold of enrichment is equal to 1 (no enrichment). (C) Estimates of regression coefficients in the multiple MK regression. This analysis includes a new binary feature indicating whether each OD site is located in a metabolic gene. Error bars indicate 95% confidence intervals while one, two, and three asterisks indicate  $0.01 \leq P\text{-value} < 0.05$ ,  $0.001 \leq P\text{-value} < 0.01$ , and  $P\text{-value} < 0.001$ , respectively.

adaptation rate could not be explained by their biased expression in digestive organs.

In line with frequent positive selection on the immune system (Schlenke and Begun 2003; Nielsen et al. 2005; Kosiol et al. 2008; Barreiro and Quintana-Murci 2010), I observed that genes associated with the immune system had a 2- to 4-fold enrichment in the 495 tissue-specific genes with a high expression level (fig. 5A; false-discovery rate  $< 0.01$ ). To formally test if immune system genes have a higher rate of adaptation than nonimmune genes in chimpanzees, I constructed a genomic feature indicating whether each OD site was located in one of the 3,400 immune system genes from MSigDB (Subramanian et al. 2005; Liberzon et al. 2011). After adding this new feature to the multiple MK regression

analysis, I observed that the regression coefficient of immune genes was significantly higher than 0 (supplementary fig. 10, Supplementary Material online and supplementary table 8, Supplementary Material online). Thus, immune genes may have a higher rate of adaptive evolution than their nonimmune counterparts. Immune genes may also partially explain the positive effect of gene expression level on the rate of adaptive evolution, since the regression coefficient of gene expression level decreased from 0.314 to 0.238 after adding immune genes as a new feature (supplementary tables 7 and 8, Supplementary Material online).

## Discussion

In this work, I have introduced the MK regression, the first evolutionary model for jointly estimating the effects of multiple, potentially correlated genomic features on the rate of adaptive substitutions. Based on similar ideas, my colleagues and I have previously developed statistical approaches to infer negative selection on genetic variants (Huang et al. 2017; Huang and Siepel 2019; Huang 2020) and the evolutionary turnover of cis-regulatory elements (Dukler et al. 2020). Thus, unifying generalized linear models and evolutionary models may be a powerful strategy to address a variety of statistical problems in evolutionary biology.

As shown in the simulation experiments (fig. 2), when two genomic features are correlated with each other, even at a moderate level, the simple MK regression and a previous MK-based method (Smith and Eyre-Walker 2002; Fraïsse et al. 2019) cannot accurately estimate the independent effects of genomic features because they cannot control for correlated genomic features. On the other hand, the multiple MK regression can unbiasedly estimate the independent effects of genomic features if all relevant genomic features are included in the same analysis. Because we are almost always interested in the independent effects of genomic features, the multiple MK regression may be superior to the simple MK regression and other MK-based methods that can only analyze one feature at a time.

Regression coefficients in the multiple MK regression might be interpreted as the direct causal effects of genomic features on the rate of adaptation. However, similar to other linear regression models (Pearl et al. 2016), the causal interpretation of the multiple MK regression relies on two implicit assumptions. First, genomic features of interest and all correlated genomic features should be included in the same MK regression analysis. Second, there should be no reverse causality, that is, the rate of adaptive evolution should not cause changes in genomic features. As discussed in the literature of causal inference (Pearl et al. 2016), these assumptions cannot be verified using observational data alone and, thus, have to be justified by domain knowledge on a case-by-case basis.

It is worth noting that I do not attempt to infer the total causal effects of genomic features on adaptation rate in the current study. According to the theory of causal inference (Pearl et al. 2016), the total effect of a genomic feature of interest includes its direct effect on the rate of adaptation as well as its indirect effect through mediators, that is, others

genomic features that reside on a directed path from the genomic feature of interest to the rate of adaptation in an assumed causal graph. Thus, inferring the total effect of the genomic feature of interest will require me to impose very strong assumptions about the causal relationship between genomic features (see examples in Rosenbaum et al. 2020; Laubach et al. 2021). On the other hand, to infer the direct effect of the genomic feature of interest, I can simply control for all potentially correlated features without specifying their causal relationship with the genomic feature of interest (Laubach et al. 2021). In other words, the MK regression effectively assumes a simplified causal graph where I do not specify the directions of causality between genomic features (supplementary fig. 11, Supplementary Material online; see also Chapter 4.3 in Shipley 2016).

It is also worth noting that I have ignored potential interactions between genomic features in the current study. It is possible to add interaction terms to the multiple MK regression. However, because of the large number of potential interaction terms and the sparseness of polymorphisms in the chimpanzee genome, I may lack statistical power to detect interactions between features in chimpanzees.

Using the multiple MK regression, I have identified numerous genomic features with independent effects on adaptive evolution in the chimpanzee lineage (fig. 3B). First, in line with previous studies (Campos et al. 2014; Castellano et al. 2016; Rousselle et al. 2020), I have shown that the rate of adaptation increases with increasing mutation rate. Because mutations are the ultimate source of genetic variation, a higher mutation rate will increase the genetic variation for positive selection to act on.

Second, previous studies have shown that local recombination rate is positively correlated with the rate of adaptation in *Drosophila* (Marais and Charlesworth 2003; Campos et al. 2014; Castellano et al. 2016), probably due to a reduced effect of Hill-Robertson interference in recombination hotspots. However, I have not observed the same pattern in chimpanzees, which could be explained by a reduced impact of linked selection in species with a small census population size, such as chimpanzees (Corbett-Detig et al. 2015). Alternatively, my analysis may have limited power to detect a weak association between recombination rate and positive selection due to the lower level of polymorphism and the smaller proportion of adaptive substitutions in chimpanzees compared with *Drosophila* (Castellano et al. 2016).

Third, in agreement with a previous study (Moutinho, Trancoso et al. 2019), I have shown that the rate of adaptive evolution increases with the increasing level of residue exposure to solvent. On the other hand, it is well-known that the site-wise rate of protein evolution is positively correlated with the level of residue exposure, possibly due to relaxed negative selection on exposed residues (Goldman et al. 1998; Franzosa and Xia 2009; Liberles et al. 2012; Echave et al. 2016). Taken together, exposed residues on protein surfaces may be subject to both weaker negative selection and more frequent positive selection than their buried counterparts. From a biophysical perspective, missense mutations on protein surfaces are less likely to disrupt protein stabilities than mutations in

hydrophobic cores (Bloom et al. 2005; Bloom, Labthavikul et al. 2006; Bloom, Drummond et al. 2006; Franzosa and Xia 2009). Thus, missense mutations on protein surfaces may be under more frequent positive selection because they are less likely to perturb protein folding. Alternatively, protein surfaces may have a higher rate of adaptive evolution because they may play an important role in host–pathogen interactions (Moutinho, Trancoso et al. 2019).

Fourth, I have shown that the tissue specificity of a gene has a positive effect on the rate of adaptation after controlling for correlated genomic features, such as gene expression level. Thus, tissue-specific genes are more likely to be under positive selection than housekeeping genes. Because nonsynonymous mutations in housekeeping genes have a higher chance to disrupt multiple phenotypes, my findings support that the pleiotropic effect is a key determinant of adaptive evolution (Fraïsse et al. 2019).

Fifth, I have shown that the rate of adaptive evolution increases with increasing gene expression levels in both chimpanzees (figs. 3B and 4B) and *Drosophila* (supplementary fig. 7, Supplementary Material online) after controlling for correlated genomic features, such as tissue specificity. Recently, Fraïsse et al. (2019) reported an opposite trend in *Drosophila* by regressing a gene-level estimate of  $\omega_a$  on gene expression level and potentially correlated features. However, my reanalysis of their data suggests that the standard linear regression used in Fraïsse et al. (2019) may not be an appropriate method for inferring the effects of genomic features in *Drosophila*. Unlike the standard linear regression, the MK regression does not rely on inaccurate estimates of  $\omega_a$  at the gene level, and does not assume that the response variables follow a normal distribution. Thus, the MK regression may be more broadly applicable than the standard linear regression in inferring the effects of genomic features on adaptation.

Sixth, numerous studies have shown that immune genes may have a higher rate of adaptive evolution than other genes in various species (Schlenke and Begun 2003; Nielsen et al. 2005; Kosiol et al. 2008; Barreiro and Quintana-Murci 2010). Using the MK regression, I have observed the same trend in chimpanzees (supplementary fig. 10, Supplementary Material online and supplementary table 8, Supplementary Material online). Thus, immune genes may be subject to constant adaptation in chimpanzees to fight against ever-evolving pathogens and parasites.

Last but not least, I have shown that highly expressed genes are more likely to be associated with metabolic pathways and digestive organs than their lowly expressed counterparts (fig. 5A and 5B), which implies that frequent positive selection in metabolic genes may partially explain the positive effect of gene expression level on the rate of adaptation. In agreement with this hypothesis, I have shown that metabolic genes may have a higher rate of adaptation than their nonmetabolic counterparts after controlling for potentially correlated genomic features (fig. 5C). Similarly, a recent study has reported that metabolic pathways may be subject to more frequent positive selection than nonmetabolic pathways in multiple inner branches of the primate phylogeny (Daub et al. 2017).

Taken together, genes in metabolic pathways may be subject to frequent positive selection in multiple primate species, possibly due to recent changes in diet in primate evolution (Daub et al. 2017; Haygood et al. 2007; Blekhman et al. 2008, 2014). In future studies, it is also interesting to examine whether other genomic features related to metabolism, such as the replication timing of genes (Chen et al. 2010) in digestive organs, have independent effects on the rate of adaptive evolution using the MK regression.

Frequent positive selection in metabolic genes has not been widely reported in primates, except in the current study and in Daub et al. (2017). I speculate that the discrepancy could be explained by the unique design of the MK regression and the method in Daub et al. (2017). First, previous studies focused on identifying individual genes with significant signals of positive selection. If metabolic pathways are under polygenic selection, the signal of selection in a single gene may be too weak to reach genome-wide significance (Csilléry et al. 2018; Barghi et al. 2020). In contrast, the MK regression and the method in Daub et al. (2017) have pooled data across a large number of metabolic genes, which may significantly increase the statistical power to detect diffused signals of polygenic selection (Barghi et al. 2020). Second, if metabolic genes are under lineage-specific adaptation in primates, positive selection may only be detected by statistical approaches tailored for a single branch of the primate phylogeny, such as the MK regression and the branch-site codon substitution model (Daub et al. 2017). Third, many previous methods may not be able to control for the effects of correlated genomic features. In the current study, I have shown that the effect of metabolic genes is manifested in the multiple MK regression but not in the simple MK regression, which highlights the importance of controlling for correlated genomic features. In future studies, it is tempting to test these hypotheses for a comprehensive understanding of adaptive evolution in metabolic genes.

Similar to a previous study (Luisi et al. 2015), I have found that the rate of adaptation increases with increasing PPI degree in the simple MK regression (fig. 3A). However, the same pattern cannot be replicated in the multiple MK regression (fig. 3B). Also, after controlling for gene expression level and tissue specificity using propensity score matching, I have found no difference in the rate of adaptation between genes with high PPI degree and genes with low PPI degree (supplementary fig. 3B, Supplementary Material online). Thus, PPI degree is unlikely to be a key determinant of positive selection in chimpanzees.

Because functional genomic data are scarce in the chimpanzee genome, I have mapped multiple genomic features from the human genome to the chimpanzee genome to examine adaptive evolution in the chimpanzee lineage. It is worth noting that the current study does not require or assume that these genomic features are perfectly correlated between humans and chimpanzees. Unlike studies that aim to identify individual loci under positive selection, I focus on examining genome-wide relationships between genomic features and adaptation rate. As long as genomic features are well correlated between humans and chimpanzees at the

genome-wide scale, my results should be robust to species differences in genomic features in a small set of genes. Because of the short divergence time between humans and chimpanzees, I expect that human features are reasonable proxies of corresponding chimpanzee features at the genome-wide scale. Nevertheless, it is of interest to revisit the results reported in the current study when chimpanzee-specific annotations become available in the future.

While the MK regression is a powerful framework to estimate the effects of multiple genomic features simultaneously, it has a few limitations that are worth of future exploration. Notably, the MK regression is based on the classical MK test and, thus, inherits its limitations (McDonald and Kreitman 1991; Smith and Eyre-Walker 2002). First, the MK regression does not explicitly model the effects of weak negative selection on polymorphism data. To mitigate this problem, I have used a simple strategy to filter out low-frequency SNPs. This strategy may not be optimal, because a large proportion of SNPs cannot be used in the MK regression despite the fact that they are potentially informative of positive selection (Messer and Petrov 2013). Second, the MK regression assumes that positively selected mutations have a negligible contribution to polymorphisms. This assumption is valid when weak positive selection is rare. However, there is evidence for frequent weak positive selection in the human genome, and ignoring the effects of weak positive selection may lead to an underestimation of the adaptation rate in humans (Uricchio et al. 2019). Thus, the MK regression may not be suitable for examining positive selection in the human genome. Third, the MK regression does not explicitly model the impact of demography on polymorphisms. In future studies, it is tempting to extend the MK regression by explicitly modeling the effects of weak negative selection, weak positive selection, and demography on the site-frequency spectrum of polymorphisms (Eyre-Walker and Keightley 2009; Messer and Petrov 2013; Galtier 2016; Haller and Messer 2017; Uricchio et al. 2019).

From a statistical point of view, several aspects of the MK regression may also be improved in future studies. First, it is tempting to relax the strong assumption of a linear relationship between genomic features and the rate of adaptation. For instance, we may replace the generalized linear model in the MK regression by a generalized additive model (Hastie 1990), which can accommodate more complicated relationships between features and selection while maintaining the interpretability of the MK regression. Second, the MK regression is mainly designed to infer the effects of genomic features on the rate of adaptive evolution. Thus, it may not be the best tool for pinpointing individual genes under frequent positive selection. If positive selection at the gene level is of the main interest, and if the numbers of polymorphic and divergent sites are large in a gene, the classical MK test (McDonald and Kreitman 1991; Smith and Eyre-Walker 2002) may be more appropriate for inferring positive selection at the gene level. Third, the MK regression currently can only accommodate the fixed effects of genomic features. It is tempting to extend the MK regression by introducing a gene-level random effect (Huang 2020), which may allow for estimating the rate of

adaptation at the gene level. Fourth, the MK regression currently can only estimate the effects of genomic features on  $\omega_a$ . Based on the extensive literature on the MK test (Booker et al. 2017; Moutinho, Bataillon et al. 2019), it is possible to extend the MK regression to estimate the effects of genomic features on other measures of natural selection, such as the rate of nonadaptive evolution ( $\omega_{na}$ ) and the proportion of adaptive substitutions ( $\alpha$ ). Fifth, similar to other multiple regression models, the MK regression might yield unreliable estimates of regression coefficients if there is strong collinearity between genomic features (Dormann et al. 2013). The current study is unlikely to be affected by strong collinearity because the absolute values of correlation coefficients between genomic features are smaller than 0.7 (supplementary table 4, Supplementary Material online), an established cutoff for strong collinearity (Dormann et al. 2013). Nevertheless, in future studies, it is tempting to explore rigorous statistical techniques, such as the variance inflation factor (Fox and Monette 1992), to detect and handle potentially strong collinearity in the MK regression. Finally, when used as a part of a causal inference pipeline, the MK regression is not an automatic tool and requires domain knowledges of all genomic features being considered in an analysis. I expect that future extensions of the MK regression will enable systematic explorations of the genomic basis of adaptive and nonadaptive evolution in various species, such as humans and *Drosophila*.

## Materials and Methods

### Details of the MK Regression

The MK regression consists of two components: a generalized linear model and an MK-based likelihood function. In the generalized linear model, I assume that  $\omega_a^j$ , the relative rate of adaptive substitutions at functional site  $j$ , is a linear combination of genomic features followed by an exponential transformation,

$$\omega_a^j = \exp(\beta_0 + \beta_1 X_1^j + \cdots + \beta_i X_i^j + \cdots + \beta_M X_M^j), \quad (1)$$

in which  $X_i^j$  is the  $i$ th local genomic feature at site  $j$ ,  $\beta_i$  is a regression coefficient indicating the effect of feature  $i$  on  $\omega_a^j$ ,  $\beta_0$  in an intercept, and  $M$  is the total number of genomic features. Because genomic features may also influence the levels of polymorphism at functional sites, I model the probability of observing a SNP at functional site  $j$ ,  $P_{\text{func}}^j$ , as another linear combination of genomic features followed by a logistic transformation,

$$\begin{aligned} P_{\text{func}}^j &= \text{logistic}(\gamma_0 + \gamma_1 X_1^j + \cdots + \gamma_i X_i^j + \cdots + \gamma_M X_M^j) \\ &= \frac{\exp(\gamma_0 + \gamma_1 X_1^j + \cdots + \gamma_i X_i^j + \cdots + \gamma_M X_M^j)}{1 + \exp(\gamma_0 + \gamma_1 X_1^j + \cdots + \gamma_i X_i^j + \cdots + \gamma_M X_M^j)}, \end{aligned} \quad (2)$$

in which  $\gamma_0$  and  $\gamma_i$  are an intercept and a regression coefficient with respect to  $P_{\text{func}}^j$ , respectively. Similar to the regression coefficients in equation 1,  $\gamma_i$  represents the effect of feature  $i$  on the occurrence of polymorphism at functional site  $j$ . It is worth noting that I effectively assume an infinite-site model

here (Kimura 1969), so no more than one SNP is allowed at a single site. Finally, I introduce two neutral parameters,  $D_{\text{neut}}$  and  $P_{\text{neut}}$ , which represent the expected number of substitutions and the probability of observing a SNP at a neutral site, respectively. These neutral parameters are shared by all neutral sites.

In the MK-based likelihood function, I specify the probability of polymorphism and divergence data at both neutral and functional sites given model parameters ( $\beta_0$ ,  $\beta_i$ ,  $\gamma_0$ ,  $\gamma_i$ ,  $D_{\text{neut}}$ , and  $P_{\text{neut}}$ ). First, I denote  $Y_{\text{neut}}^k$  as a binary response variable indicating the presence/absence of a SNP at neutral site  $k$  and assume that it follows a Bernoulli distribution,

$$\mathbb{P}(Y_{\text{neut}}^k) = \begin{cases} P_{\text{neut}}, & \text{if } Y_{\text{neut}}^k = 1 \\ 1 - P_{\text{neut}}, & \text{otherwise.} \end{cases} \quad (3)$$

Similarly, I denote  $Y_{\text{func}}^j$  as a binary response variable indicating the presence/absence of a SNP at functional site  $j$  and assume that it follows a Bernoulli distribution,

$$\mathbb{P}(Y_{\text{func}}^j) = \begin{cases} P_{\text{func}}^j, & \text{if } Y_{\text{func}}^j = 1 \\ 1 - P_{\text{func}}^j, & \text{otherwise.} \end{cases} \quad (4)$$

Second, I employ the Jukes–Cantor substitution model (Jukes and Cantor 1969) to describe the distribution of interspecies divergence at neutral sites. Denoting  $Z_{\text{neut}}^k$  as a binary response variable indicating if the reference genome and the ancestral genome have different nucleotides at neutral site  $k$ , the Jukes–Cantor model suggests that

$$\mathbb{P}(Z_{\text{neut}}^k) = \begin{cases} \frac{3}{4} - \frac{3}{4} \exp(-\frac{4}{3} D_{\text{neut}}), & \text{if } Z_{\text{neut}}^k = 1 \\ \frac{1}{4} + \frac{3}{4} \exp(-\frac{4}{3} D_{\text{neut}}), & \text{otherwise.} \end{cases} \quad (5)$$

Third, to model the effects of positive selection and neutral factors on interspecies divergence at functional sites, I assume that  $D_{\text{func}}^j$ , that is, the expected number of substitutions at functional site  $j$ , is equal to the sum of the number of adaptive substitutions and the number of neutral substitutions (Bierne and Eyre-Walker 2004; Gossmann et al. 2010),

$$D_{\text{func}}^j = \underbrace{\omega_a^j D_{\text{neut}}}_{\text{substitutions}^{\text{Adaptive}}} + \underbrace{\frac{P_{\text{func}}^j}{P_{\text{neut}}} D_{\text{neut}}}_{\text{substitutions}^{\text{Neutral}}}, \quad (6)$$

in which  $\frac{P_{\text{func}}^j}{P_{\text{neut}}}$  is equal to the relative rate of neutral evolution at functional site  $j$  compared with neutral sites. Given  $D_{\text{func}}^j$  at each functional site, I again employ the Jukes–Cantor substitution model to describe interspecies divergence at functional site  $j$ ,

$$\mathbb{P}(Z_{\text{func}}^j) = \begin{cases} \frac{3}{4} - \frac{3}{4} \exp(-\frac{4}{3} D_{\text{func}}^j), & \text{if } Z_{\text{func}}^j = 1 \\ \frac{1}{4} + \frac{3}{4} \exp(-\frac{4}{3} D_{\text{func}}^j), & \text{otherwise,} \end{cases} \quad (7)$$

in which  $Z_{\text{func}}^j$  indicates if the reference genome and the ancestral genome have different nucleotides at functional site  $j$ .

Finally, I assume that nucleotide sites evolve independently given genomic features and model parameters. Thus, I define the MK-based likelihood function of the whole data set as

$$\prod_{j \in \text{all functional sites}} \mathbb{P}(Y_{\text{func}}^j) \mathbb{P}(Z_{\text{func}}^j) \prod_{k \in \text{all neutral sites}} \mathbb{P}(Y_{\text{neut}}^k) \mathbb{P}(Z_{\text{neut}}^k). \quad (8)$$

It is worth noting that the MK regression currently can only be applied to a nucleotide site where all the potential mutations have similar effects, such as 0D and 4D sites in coding regions.

I estimate model parameters ( $\beta_0$ ,  $\beta_1$ ,  $\beta_2$ ,  $\gamma_0$ ,  $\gamma_1$ ,  $\gamma_2$ ,  $P_{\text{neut}}$ , and  $D_{\text{neut}}$ ) by maximizing the logarithm of the MK-based likelihood function (equation 8). Also, I estimate the standard errors of parameters using the observed Fisher information matrix and compute the  $P$ -values of estimated parameters using the (two-tailed) Wald test.

### Simulation Experiments

In each simulation run, I first sampled genomic features ( $X_i^j$ ) at 10 Mb functional sites from a bivariate normal distribution. I set the means and variances of the bivariate normal distribution to 0 and 1, respectively, and varied its correlation coefficient from 0 to 0.6 to examine the performance of the MK regression with respect to various degrees of correlation between features. Given genomic features, I generated divergence and polymorphism data following the likelihood function of the MK regression. Specifically, I generated  $\omega_a^j$ , the rate of adaptive evolution at each functional site  $j$ , using equation 1 in the MK regression model. Also, I generated  $P_{\text{func}}^j$ , the polymorphic rate at each functional site  $j$ , using equation 2 in the MK regression. Finally, I generated binary response variables of polymorphism and divergence at 10 Mb neutral sites ( $Y_{\text{neut}}^k$  and  $Z_{\text{neut}}^k$ ) and 10 Mb functional sites ( $Y_{\text{func}}^j$  and  $Z_{\text{func}}^j$ ) using equations 3, 4, 5, and 7 in the MK regression.

I performed two simulation experiments each of which consisted of four sets of synthetic data. In the first experiment, I set  $\beta_0 = -2$ ,  $\beta_1 = 1$ ,  $\beta_2 = 0$ ,  $\gamma_0 = -8$ , and sampled  $\gamma_1$  and  $\gamma_2$  from a normal distribution with a mean of 0 and a standard deviation of 0.5 in each simulation run. These parameters were chosen to ensure that genome-wide levels of polymorphism and divergence are comparable between synthetic and empirical data. I generated four sets of synthetic data with different correlation coefficients (0, 0.2, 0.4, and 0.6) between features in the aforementioned bivariate normal distribution. In each data set, I performed 10 independent simulation runs using the method described in the previous paragraph. In the second experiment, I generated synthetic data using the same procedure but replaced  $\beta_2$  with  $-0.2$ .

### Polymorphism and Divergence Data

Throughout this work, I focused on analyzing 0D and 4D sites within previously defined callable regions on autosomal chromosomes in the panTro4 reference genome (de Manuel et al. 2016). I obtained whole-genome sequencing (WGS) based genotypes of 18 central chimpanzee (*Pan troglodytes troglodytes*) individuals from de Manuel et al. (2016). I filtered out

all multiallelic sites and sites with missing genotypes. Then, I filtered out SNPs without a high-confidence ancestral allele (see below) and SNPs with a derived allele frequency below 0.5.

I obtained high-confidence chimpanzee ancestral alleles from Ensembl release 75 (Herrero et al. 2016; Yates et al. 2020). Ensembl reconstructed chimpanzee ancestral alleles using a phylogenetic approach, and defined that an ancestral allele was of high confidence if it was identical to both the human reference allele and the reconstructed allele in the human–chimpanzee–macaque ancestor. To annotate interspecies divergence in the chimpanzee lineage, I compared the chimpanzee ancestral allele with the reference allele in panTro4 at each monoallelic site.

### Estimation of $\alpha$ and $\omega_a$

Based on the aforementioned polymorphism and divergence data in chimpanzee autosomal genes, I computed  $d_n$ ,  $p_n$ ,  $d_s$ , and  $p_s$ , which are the proportions of 0D sites with divergence, 0D sites with polymorphism, 4D sites with divergence, and 4D sites with polymorphism, respectively. Each of these proportions was calculated as the ratio of the number of sites with divergence/polymorphism to the total number of sites. I estimated the proportion of adaptive substitutions as (Charlesworth 1994; Smith and Eyre-Walker 2002)

$$\alpha = 1 - \frac{d_s p_n}{d_n p_s}. \quad (9)$$

Similarly, given the polymorphism and divergence data in a gene group of interest, I estimated the relative rate of adaptive substitutions with respect to neutral evolution as (Smith and Eyre-Walker 2002; Booker et al. 2017; Fraïsse et al. 2019)

$$\begin{aligned} \omega_a &= \frac{d_n \cdot \alpha}{d_s} \\ &= \frac{d_n}{d_s} - \frac{p_n}{p_s}. \end{aligned} \quad (10)$$

I also used equation 10 to estimate  $\omega_a$  for synthetic data, where  $d_n$ ,  $p_n$ ,  $d_s$ , and  $p_s$  were interpreted as the proportions of simulated functional sites with divergence, simulated functional sites with polymorphism, simulated neutral sites with divergence, and simulated neutral sites with polymorphism, respectively.

### Genomic Features

Because gene annotations and functional genomic data were more complete and of higher quality in humans than in chimpanzees, I obtained human-based annotations of 0D sites, 4D sites, residue exposure level, gene expression level, tissue specificity, and protein–protein interactions. Then, I converted these annotations from the human reference genome to the panTro4 assembly using liftOver (Haeussler et al. 2019). Because of the very short divergence time between humans and chimpanzees, human-based genomic features should serve as an accurate proxy for the corresponding genomic features in chimpanzees. Specifically, I obtained the coordinates of 0D and 4D sites in the human genome from

dbNSFP version 4.0 (Liu et al. 2013, 2016) and converted the coordinates from the original hg19 assembly to the panTro4 assembly using liftOver (Haeussler et al. 2019). I obtained predicted probabilities of residue exposure (PredRSAE) in the hg19 assembly from SNVBox (Wong et al. 2011). I obtained human-based consensus RNA expression levels across 62 tissues from the Human Protein Atlas version 19.3 (Uhlen et al. 2015). For each human protein-coding gene, I computed its (mean) expression level,

$$\text{expression level} = \frac{\sum_1^K R_k}{K}, \quad (11)$$

in which  $R_k$  is the gene's consensus RNA expression level in tissue  $k$  and  $K$  is the total number of tissues. I computed the tissue specificity ( $\tau$ ) of each human gene,

$$\tau = \frac{\sum_1^K 1 - R_k/\max(R_k)}{K - 1}, \quad (12)$$

in which  $\max(R_k)$  is the gene's maximum expression level across all tissues (Yanai et al. 2005). Also, I computed an alternative metric of tissue specificity (the negative value of Hg; Kryuchkova-Mostacci and Robinson-Rechavi 2017),

$$-\text{Hg} = \sum_1^K p_k \cdot \log_2(p_k), \quad (13)$$

where  $p_k = \frac{R_k}{\sum_1^K R_k}$  is the normalized expression level of a gene

in tissue  $k$ . I obtained 2,274 genes with biased or enriched expression in digestive organs, including intestine, liver, pancreas, salivary gland, and stomach, from the Human Protein Atlas (Uhlen et al. 2015). I obtained human protein-protein interaction data from HuRI (Luck et al. 2020) and computed each gene's PPI degree, that is, the total number of unique interaction partners. Also, I obtained 2,220 human metabolic genes involved in one or more curated metabolic pathways and 3,400 human genes involved in one or more immune system pathways from MSigDB release 7.1 (Subramanian et al. 2005; Liberzon et al. 2011). Finally, I converted human-based annotations of residue exposure level, gene expression level, tissue specificity, PPI degree, metabolic genes, and immune system genes from the hg19 assembly to the panTro4 assembly using liftOver (Haeussler et al. 2019).

I used chimpanzee-based data to build a map of local recombination rates and a map of local mutation rates. Specifically, I obtained a fine-scale chimpanzee genetic map from panMap (Auton et al. 2012) and converted the data from panTro2 to panTro4 using liftOver. Then, I constructed a map of local recombination rates by averaging the recombination rates from the chimpanzee genetic map with a 1 Mb nonoverlapping sliding window. Also, I utilized interspecies divergence in the chimpanzee lineage to construct a map of local mutation rates in the panTro4 assembly. To do so, I converted putatively neutral regions defined in Huang

(2020) from hg19 to panTro4 using liftOver. Then, I computed the density of chimpanzee-specific substitutions in putatively neutral regions using a 100 Kb nonoverlapping sliding window, which was used as a proxy of local mutation rates.

### Estimating the Effects of Genomic Features on the Rate of Adaptation

I fit the MK regression to one genomic feature at a time, which I named as the simple MK regression. To evaluate if a logarithmic transformation can improve model fitting, I carried out two analyses for each feature. In the first analysis, I standardized the feature by subtracting its mean and dividing by its standard deviation, and then fit the simple MK regression to the standardized feature. In the second analysis, I calculated the logarithm of each feature, standardized the output, and then fit the simple MK regression to the transformed data. Because the logarithm of PPI degree is undefined if the PPI degree is equal to 0, I added a pseudocount of 1 to the PPI degree before the logarithmic transformation. I computed the log likelihood of the simple MK regression in each analysis, and used the transformation with a higher log likelihood for each feature throughout this work. I also fit the MK regression to all the features simultaneously, which I named as the multiple MK regression.

### Statistical Matching Analysis

I used statistical matching algorithms to estimate the effects of PPI degree, gene expression level, and tissue specificity after adjusting for other correlated genomic features. To estimate the effect of PPI degree, I stratified protein-coding genes into two groups based on PPI degree. The first group consisted of 1,556 genes with 10 or more interaction partners (PPI degree  $\geq 10$ ) while the second group consisted of 8,471 genes with 1 or less interaction partners (PPI degree  $\leq 1$ ). Here I chose a 10-fold difference in PPI degree between the two gene groups, because a larger difference in PPI degree led to significantly smaller gene groups, whereas a smaller difference in PPI degree may attenuate the potential difference in  $\omega_a$  between the two gene groups. I calculated  $\omega_a$  for the two groups of genes separately using equation 10, and calculated the  $P$ -value of the difference in  $\omega_a$  between the two gene groups using a two-tailed permutation test with 1,000 resamplings. Also, I used MatchIt to match each gene from the high-PPI group with a gene of similar log expression level and tissue specificity from the low-PPI group (Ho et al. 2011). Then, I repeated the calculation of  $\omega_a$  and  $P$ -value for the two groups of matched genes.

To estimate the effect of tissue specificity after controlling for gene expression level, I stratified highly expressed genes (mean expression level  $> 30$ ) into two approximately equal-sized groups based on the ranking of their tissue specificity. The first gene group consisted of 368 highly expressed genes with higher tissue specificity while the second group consisted of 369 highly expressed genes with lower tissue specificity. Then, I calculated  $\omega_a$  for the two groups of genes separately using equation 10, and calculated the  $P$ -value of the difference in  $\omega_a$  using a two-tailed permutation test with 1,000 resamplings.

Finally, I estimated the effect of gene expression level after controlling for tissue specificity. I stratified tissue-specific genes ( $\tau > 0.85$ ) into two approximately equal-sized groups based on the ranking of their expression levels. The first group consisted of 495 tissue-specific genes with higher expression level while the second group consisted of 498 tissue-specific genes with lower expression level. I calculated  $\omega_a$  for the two groups of genes separately using equation 10, and calculated the  $P$ -value of the difference in  $\omega_a$  using a two-tailed permutation test with 1,000 resamplings.

### Reanalysis of Data from Fraïsse et al. (2019)

I estimated the effects of gene expression level and tissue specificity on the rate of adaptive evolution in *D. melanogaster* by reanalyzing data from Fraïsse et al. (2019). I obtained gene-level estimates of  $\omega_a$ , mean expression levels, tissue-by-stage specificity ( $\tau$ ), polymorphism data, and divergence data from <http://doi.org/10.15479/atista:/5757>. In the analysis of standard linear regression, I regressed the gene-level estimate of  $\omega_a$  on the logarithm of mean expression level and the tissue-by-state specificity using the *lm* function in R (R Core Team 2017). In the analysis of multiple MK regression, I used 0D sites and 4D sites annotated by SIFT 4G (Vaser et al. 2016) as functional and putatively neutral sites, respectively, and used the logarithm of mean expression level and the tissue-by-state specificity as input features. In the statistical matching analysis, I stratified 681 tissue-specific genes ( $\tau > 0.85$ ) into a group of 340 genes with higher expression level and a group of 341 genes with lower expression level. I calculated  $\omega_a$  for the two groups of genes separately using equation 10, and calculated the  $P$ -value of the difference in  $\omega_a$  using a two-tailed permutation test with 1,000 resamplings.

### Tissue and Pathway Enrichment

I downloaded annotations of tissue-enriched genes from the Human Protein Atlas version 19.3 (Uhlen et al. 2015). To reduce the burden of multiple testing, I focused on analyzing tissues with at least 50 tissue-enriched genes. I then analyzed the enrichment of each set of tissue-enriched genes in the 495 tissue-specific genes with a high expression level, in which the 498 tissue-specific genes with a low expression level were used as a background gene set. The  $P$ -value of each enrichment test was computed using the Fisher's exact test and then adjusted for multiple testing with false-discovery-rate correction. Similarly, I analyzed the enrichment of Reactome pathways in the 495 tissue-specific genes with a high expression level using PANTHER (Mi et al. 2017), in which I used the 498 tissue-specific genes with a low expression level as a background gene set.

### Supplementary Material

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

### Acknowledgments

The author acknowledges David McCandlish, Xinru Zhang, Jui-Shan Lin, and Zhihan Liu for useful discussions. Research

reported in this publication was supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number R35GM142560 and by startup funds from Pennsylvania State University. The content is solely the responsibility of the author and does not necessarily represent the official views of the National Institutes of Health.

### Data Availability

The MK regression model and companion data are available at <https://github.com/yifei-lab/MK-regression>.

### References

- Afanasyeva A, Bockwoldt M, Cooney CR, Heiland I, Gossmann TI. 2018. Human long intrinsically disordered protein regions are frequent targets of positive selection. *Genome Res.* 28(7):975–982.
- Auton A, Fledel-Alon A, Pfeifer S, Venn O, Ségurel L, Street T, Leffler EM, Bowden R, Aneas I, Broxholme J, et al. 2012. A fine-scale chimpanzee genetic map from population sequencing. *Science* 336(6078):193–198.
- Avila V, Campos JL, Charlesworth B. 2015. The effects of sex-biased gene expression and x-linkage on rates of adaptive protein sequence evolution in drosophila. *Biol Lett.* 11(4):20150117.
- Barghi N, Hermisson J, Schlötterer C. 2020. Polygenic adaptation: a unifying framework to understand positive selection. *Nat Rev Genet.* 21(12):769–781.
- Barreiro LB, Quintana-Murci L. 2010. From evolutionary genetics to human immunology: how selection shapes host defence genes. *Nat Rev Genet.* 11(1):17–30.
- Bierne N, Eyre-Walker A. 2004. The genomic rate of adaptive amino acid substitution in drosophila. *Mol Biol Evol.* 21(7):1350–1360.
- Blekhman R, Oshlack A, Chabot AE, Smyth GK, Gilad Y. 2008. Gene regulation in primates evolves under tissue-specific selection pressures. *PLoS Genet.* 4(11):e1000271.
- Blekhman R, Perry GH, Shahbaz S, Fiehn O, Clark AG, Gilad Y. 2014. Comparative metabolomics in primates reveals the effects of diet and gene regulatory variation on metabolic divergence. *Sci Rep.* 4(1):5809.
- Bloom JD, Silberg JJ, Wilke CO, Drummond DA, Adami C, Arnold FH. 2005. Thermodynamic prediction of protein neutrality. *Proc Natl Acad Sci U S A.* 102(3):606–611.
- Bloom JD, Labthavikul ST, Otey CR, Arnold FH. 2006a. Protein stability promotes evolvability. *Proc Natl Acad Sci U S A.* 103(15):5869–5874.
- Bloom JD, Drummond DA, Arnold FH, Wilke CO. 2006b. Structural determinants of the rate of protein evolution in yeast. *Mol Biol Evol.* 23(9):1751–1761.
- Booker TR, Jackson BC, Keightley PD. 2017. Detecting positive selection in the genome. *BMC Biol.* 15(1):98.
- Campos JL, Halligan DL, Haddrill PR, Charlesworth B. 2014. The relation between recombination rate and patterns of molecular evolution and variation in *Drosophila melanogaster*. *Mol Biol Evol.* 31(4):1010–1028.
- Campos JL, Johnston KJA, Charlesworth B. 2018. The effects of sex-biased gene expression and X-linkage on rates of sequence evolution in *Drosophila*. *Mol Biol Evol.* 35(3):655–665.
- Castellano D, Coronado-Zamora M, Campos JL, Barbadilla A, Eyre-Walker A. 2016. Adaptive evolution is substantially impeded by Hill-Robertson interference in *Drosophila*. *Mol Biol Evol.* 33(2):442–455.
- Castellano D, Uricchio LH, Munch K, Enard D. 2019. Viruses rule over adaptation in conserved human proteins. bioRxiv. doi: 10.1101/555060. Available from: <https://www.biorxiv.org/content/10.1101/555060v1>.
- Charlesworth B. 1994. The effect of background selection against deleterious mutations on weakly selected, linked variants. *Genet Res.* 63(3):213–227.

- Chen C-L, Rappailles A, Duquenne L, Huvet M, Guilbaud G, Farinelli L, Audit B, d'Aubenton-Carafa Y, Arneodo A, Hyrien O, et al. 2010. Impact of replication timing on non-CPG and CPG substitution rates in mammalian genomes. *Genome Res.* 20(4):447–457.
- Corbett-Detig RB, Hartl DL, Sackton TB. 2015. Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol.* 13(4):e1002112.
- Csilléry K, Rodríguez-Verdugo A, Rellstab C, Guillaume F. 2018. Detecting the genomic signal of polygenic adaptation and the role of epistasis in evolution. *Mol Ecol.* 27(3):606–612.
- Daub JT, Moretti S, Davydov I, Excoffier L, Robinson-Rechavi M. 2017. Detection of pathways affected by positive selection in primate lineages ancestral to humans. *Mol Biol Evol.* 34(6):1391–1402.
- de Manuel M, Kuhlwillm M, Frandsen P, Sousa VC, Desai T, Prado-Martinez J, Hernandez-Rodríguez J, Dupanloup I, Lao O, Hallast P, et al. 2016. Chimpanzee genomic diversity reveals ancient admixture with bonobos. *Science* 354(6311):477–481.
- Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, Marquéz JRG, Gruber B, Lafourcade B, Leitão PJ, et al. 2013. Collinearity: a review of methods to deal with it and a simulation study evaluating their performance. *Ecography* 36(1):27–46.
- Dukler N, Huang Y-F, Siepel A. 2020. Phylogenetic modeling of regulatory element turnover based on epigenomic data. *Mol Biol Evol.* 37(7):2137–2152.
- Echave J, Spielman SJ, Wilke CO. 2016. Causes of evolutionary rate variation among protein sites. *Nat Rev Genet.* 17(2):109–121.
- Enard D, Cai L, Gwennap C, Petrov DA. 2016. Viruses are a dominant driver of protein adaptation in mammals. *eLife.* 5:e12469.
- 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(9):2097–2108.
- Fay JC, Wyckoff GJ, Wu CI. 2001. Positive and negative selection on the human genome. *Genetics* 158(3):1227–1234.
- Fox J, Monette G. 1992. Generalized collinearity diagnostics. *J Am Stat Assoc.* 87(417):178–183.
- Franzosa EA, Xia Y. 2009. Structural determinants of protein evolution are context-sensitive at the residue level. *Mol Biol Evol.* 26(10):2387–2395.
- Fraïsse C, Puixeu Sala G, Vicoso B. 2019. Pleiotropy modulates the efficacy of selection in drosophila melanogaster. *Mol Biol Evol.* 36(3):500–515.
- Galtier N. 2016. Adaptive protein evolution in animals and the effective population size hypothesis. *PLoS Genet.* 12(1):e1005774.
- Goldman N, Yang Z. 1994. A codon-based model of nucleotide substitution for protein-coding dna sequences. *Mol. Biol. Evol.* 11(5):725–736.
- Goldman N, Thorne JL, Jones DT. 1998. Assessing the impact of secondary structure and solvent accessibility on protein evolution. *Genetics* 149(1):445–458.
- 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(8):1822–1832.
- Haeussler M, Zweig AS, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, Lee CM, Lee BT, Hinrichs AS, Gonzalez JN, et al. 2019. The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res.* 47(D1):D853–D858.
- Haller BC, Messer PW. 2017. asymptoticMK: a web-based tool for the asymptotic mcdonald-kreitman test. *G3 (Bethesda).* 7(5):1569–1575.
- Hastie T. 1990. Generalized additive models. Boca Raton: Routledge.
- Haygood R, Fedrigo O, Hanson B, Yokoyama K-D, Wray GA. 2007. Promoter regions of many neural- and nutrition-related genes have experienced positive selection during human evolution. *Nat Genet.* 39(9):1140–1144.
- Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle SMJ, Amode R, Brent S, et al. 2016. Ensembl comparative genomics resources. *Database (Oxford).* 2016:bav096.
- Ho DE, Imai K, King G, Stuart EA. 2011. MatchIt: nonparametric pre-processing for parametric causal inference. *J Stat Soft.* 42(8):1–28.
- Huang Y-F. 2020. Unified inference of missense variant effects and gene constraints in the human genome. *PLoS Genet.* 16(7):e1008922.
- Huang Y-F, Siepel A. 2019. Estimation of allele-specific fitness effects across human protein-coding sequences and implications for disease. *Genome Res.* 29(8):1310–1321.
- Huang Y-F, Gulko B, Siepel A. 2017. Fast, scalable prediction of deleterious noncoding variants from functional and population genomic data. *Nat Genet.* 49(4):618–624.
- Hubisz MJ, Pollard KS, Siepel A. 2011. PHAST and RPHAST: phylogenetic analysis with space/time models. *Brief Bioinform.* 12(1):41–51.
- Hughes AL. 2007. Looking for darwin in all the wrong places: the misguided quest for positive selection at the nucleotide sequence level. *Heredity (Edinb).* 99(4):364–373.
- 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 U S A.* 109(6):2054–2059.
- Jassal B, Matthews L, Viteri G, Gong C, Lorente P, Fabregat A, Sidiropoulos K, Cook J, Gillespie M, Haw R, et al. 2020. The reactome pathway knowledgebase. *Nucleic Acids Res.* 48(D1):D498–D503.
- Jukes TH, Cantor CR. 1969. Evolution of protein molecules. In: Munro HN, editor. Mammalian protein metabolism. New York and London: Academic Press. p. 21–132.
- Kimura M. 1969. The number of heterozygous nucleotide sites maintained in a finite population due to steady flux of mutations. *Genetics* 61(4):893–903.
- Kosiol C, Vinař T, da Fonseca RR, Hubisz MJ, Bustamante CD, Nielsen R, Siepel A. 2008. Patterns of positive selection in six mammalian genomes. *PLoS Genet.* 4(8):e1000144.
- Kryuchkova-Mostacci N, Robinson-Rechavi M. 2017. A benchmark of gene expression tissue-specificity metrics. *Brief Bioinform.* 18(2):205–214.
- Laubach ZM, Murray EJ, Hoke KL, Safran RJ, Perng W. 2021. A biologist's guide to model selection and causal inference. *Proc Biol Sci.* 288(1943):20202815.
- Liberles DA, Teichmann SA, Bahar I, Bastolla U, Bloom J, Bornberg-Bauer E, Colwell LJ, de Koning APJ, Dokholyan NV, Echave J, et al. 2012. The interface of protein structure, protein biophysics, and molecular evolution. *Protein Sci.* 21(6):769–785.
- Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. 2011. Molecular signatures database (MSigDB) 3.0. *Bioinformatics* 27(12):1739–1740.
- Liu X, Jian X, Eric B. 2013. dbNSFP v2.0: a database of human non-synonymous SNVs and their functional predictions and annotations. *Hum Mutat.* 34(9):E2393–E2402.
- Liu X, Wu C, Li C, Boerwinkle E. 2016. dbNSFP v3.0: a one-stop database of functional predictions and annotations for human nonsynonymous and splice-site snvs. *Hum Mutat.* 37(3):235–241.
- Luck K, Kim D-K, Lambourne L, Spirohn K, Begg BE, Bian W, Brignall R, Cafarelli T, Campos-Laborie FJ, Charleatoux B, et al. 2020. A reference map of the human binary protein interactome. *Nature* 580(7803):402–408.
- Luisi P, Alvarez-Ponce D, Pybus M, Fares MA, Bertranpetit J, Laayouni H. 2015. Recent positive selection has acted on genes encoding proteins with more interactions within the whole human interactome. *Genome Biol Evol.* 7(4):1141–1154.
- Marais G, Charlesworth B. 2003. Genome evolution: recombination speeds up adaptive evolution. *Curr Biol.* 13(2):R68–R70.
- McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the Adh locus in drosophila. *Nature* 351(6328):652–654.
- Messer PW, Petrov DA. 2013. Frequent adaptation and the mcdonald-kreitman test. *Proc Natl Acad Sci U S A.* 110(21):8615–8620.
- Mi H, Huang X, Muruganujan A, Tang H, Mills C, Kang D, Thomas PD. 2017. PANTHER version 11: expanded annotation data from gene ontology and reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* 45(D1):D183–D189.
- Moutinho AF, Trancoso FF, Duthel JY. 2019. The impact of protein architecture on adaptive evolution. *Mol Biol Evol.* 36(9):2013–2028.



- Moutinho AF, Bataillon T, Dutheil JY. 2019. Variation of the adaptive substitution rate between species and within genomes. *Evol Ecol*. 34(3):315–338.
- Muse SV, Gaut BS. 1994. A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome. *Mol Biol Evol*. 11(5):715–724.
- Nielsen R, Bustamante C, Clark AG, Glanowski S, Sackton TB, Hubisz MJ, Fedel-Alon A, Tanenbaum DM, Civello D, White TJ, et al. 2005. A scan for positively selected genes in the genomes of humans and chimpanzees. *PLoS Biol*. 3(6):e170.
- Pearl J, Glymour M, Jewell N. 2016. Causal inference in statistics: a primer. Chichester: Wiley.
- Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. 2010. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res*. 20(1):110–121.
- Pröschel M, Zhang Z, Parsch J. 2006. Widespread adaptive evolution of drosophila genes with sex-biased expression. *Genetics* 174(2):893–900.
- R Core Team 2017. R: a language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing.
- Rosenbaum S, Zeng S, Campos FA, Gesquiere LR, Altmann J, Alberts SC, Li F, Archie EA. 2020. Social bonds do not mediate the relationship between early adversity and adult glucocorticoids in wild baboons. *Proc Natl Acad Sci U S A*. 117(33):20052–20062.
- Rousselle M, Simion P, Tilak M-K, Figuet 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(4):e1008668.
- Schlenke TA, Begun DJ. 2003. Natural selection drives drosophila immune system evolution. *Genetics* 164(4):1471–1480.
- Shiple B. 2016. Cause and correlation in biology: a user's guide to path analysis, structural equations and causal inference with R. 2nd ed. Cambridge: Cambridge University Press.
- Smith NGC, Eyre-Walker A. 2002. Adaptive protein evolution in *Drosophila*. *Nature* 415(6875):1022–1024.
- Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, et al. 2005. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A*. 102(43):15545–15550.
- Tataru P, Mollion M, Glémin S, Bataillon T. 2017. Inference of distribution of fitness effects and proportion of adaptive substitutions from polymorphism data. *Genetic* 207(3):1103–1119.
- Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, Sivertsson A, Kampf C, Sjostedt E, Asplund A, et al. 2015. Tissue-based map of the human proteome. *Science* 347(6220):1260419.
- Uricchio LH, Petrov DA, Enard D. 2019. Exploiting selection at linked sites to infer the rate and strength of adaptation. *Nat Ecol Evol*. 3(6):977–984.
- Vaser R, Adusumalli S, Leng SN, Sikic M, Ng PC. 2016. SIFT missense predictions for genomes. *Nat Protoc*. 11(1):1–9.
- Wong WC, Kim D, Carter H, Diekhans M, Ryan MC, Karchin R. 2011. CHASM and SNVBox: toolkit for detecting biologically important single nucleotide mutations in cancer. *Bioinformatics* 27(15):2147–2148.
- Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, Ophir R, Bar-Even A, Horn-Saban S, Safran M, Domany E, et al. 2005. Genome-wide midrange transcription profiles reveal expression level relationships in human tissue specification. *Bioinformatics* 21(5):650–659.
- Yang Z, Nielsen R, Goldman N, Pedersen A. 2000. Codon-substitution models for heterogeneous selection pressure at amino acid sites. *Genetics* 155(1):431–449.
- Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, Azov AG, Bennett R, et al. 2020. Ensembl 2020. *Nucleic Acids Res*. 48(D1):D682–D688.
- Zhang J, Yang J-R. 2015. Determinants of the rate of protein sequence evolution. *Nat Rev Genet*. 16(7):409–420.
- Zhang L, Li W-H. 2004. Mammalian housekeeping genes evolve more slowly than tissue-specific genes. *Mol Biol Evol*. 21(2):236–239.
- Zhu J, He F, Hu S, Yu J. 2008. On the nature of human housekeeping genes. *Trends Genet*. 24(10):481–484.