

TITLE: Mutation rate variation shapes genome-wide diversity in *Drosophila melanogaster*

AUTHORS: Gustavo V. Barroso^{*1,2} and Julien Y. Dutheil¹

5 1) Max Planck Institute for Evolutionary Biology. Department of Evolutionary Genetics. August-Thienemann-Straße 2 24306 Plön – GERMANY

2) Current address: University of California, Los Angeles. Department of Ecology and Evolutionary Biology. 621 Charles E. Young Drive South 90095 – 1606 Los Angeles, CA – USA

* Corresponding author e-mail: gvbarroso@gmail.com

10

AUTHOR CONTRIBUTIONS: GVB and JYD conceived and designed the study, conducted the simulations, analysed the data and wrote the manuscript. GVB implemented the method.

ACKNOWLEDGMENTS: The authors thank Chris Kyriazis, Kai Zeng, Josep Comeron,

15 Kirk Lohmueller and Pier Palamara for discussions about this work; Ana Filipa Moutinho for providing organised data on *Drosophila*. JYD acknowledges funding from the Max Planck Society. This work was supported by a grant from the German Research Foundation (Deutsche Forschungsgemeinschaft) attributed to JYD, within the priority program (SPP) 1590 “probabilistic structures in evolution”.

20

The authors declare no competing financial interests.

The iSMC software package is freely available at <https://github.com/gvbarroso/iSMC>

Scripts used to reproduce results can be found in https://github.com/gvbarroso/ismc_dm_analyses

25 Data required to reproduce the results are deposited under DOI [10.6084/m9.figshare.13164320](https://doi.org/10.6084/m9.figshare.13164320)

ABSTRACT: What shapes the distribution of nucleotide diversity along the genome? Attempts to answer this question have sparked debate about the roles of neutral stochastic processes and natural selection in molecular evolution. However, the mechanisms of evolution do not act in isolation, and integrative models that simultaneously consider the influence of multiple factors on diversity are lacking; without them, confounding factors lurk in the estimates. Here we present a new statistical method that jointly infers the genomic landscapes of genealogies, recombination rates and mutation rates. In doing so, our model captures the effects of genetic drift, linked selection and local mutation rates on patterns of genomic variation. Guided by our causal model, we use linear regression to estimate the individual contributions of these micro-evolutionary forces to levels of nucleotide diversity. Our analyses reveal the signature of selection in *Drosophila melanogaster*, but we estimate that the mutation landscape is the major driver of the distribution of diversity in this species. Furthermore, our simulation study suggests that in many evolutionary scenarios the mutation landscape will be a crucial force shaping diversity, depending notably on the genomic window size used in the analysis. We argue that incorporating mutation rate variation into the null model of molecular evolution will lead to more realistic inference in population genomics.

INTRODUCTION:

Understanding how various evolutionary forces shape genetic diversity (π) is a major goal of population genomics (Charlesworth, 2009; Ellegren & Galtier, 2016), with a rich history of studies, both theoretical and empirical (Charlesworth & Charlesworth, 2016; Casillas & Barbadilla, 2017).

45 For many years, the debate was restricted to the relative importance of genetic drift and natural selection to the genome-wide average π (Kimura, 1968; Ohta, 1992). The observation that π does not scale linearly with population size across species (Lewontin, 1974) was termed “Lewontin’s Paradox”, and recent work has taken a new stab at this old problem (Buffalo, 2021; Galtier & Rousselle, 2020). Later on, with recognition that linkage and recombination wrap the genome in
50 regions of correlated evolutionary histories (Hudson, 1983; Kaplan & Hudson, 1985), focus shifted toward understanding how diversity levels vary along chromosomes of a single species (Pouyet & Gilbert, 2020). In 1992, Begun and Aquadro found a positive correlation between π and local recombination rate in *Drosophila melanogaster* (Begun & Aquadro, 1992), which they interpreted as the signature of linked selection (Hudson & Kaplan, 1988; Cutter & Payseur, 2013). In the three
55 decades since this seminal work, identifying the drivers of the distribution of diversity became a leading quest in the field. Nevertheless, this search has so far been incomplete. The literature has mostly focused on selection (Comeron, 2014; Elyashiv et al., 2016; McVicker et al., 2009; Stankowski et al., 2019) and introgression (Schrider et al., 2018; Stankowski et al., 2019; Hubisz et al., 2020), whereas variation in the *de novo* mutation rate (μ) along the genome has been largely
60 ignored, presumably due to challenges in its estimation (Jónsson et al., 2018; Besenbacher et al., 2019). Yet a recent study based on human trios suggests that the impact of the mutation landscape on polymorphism may be greater than previously recognised: up to 46% of the human-chimpanzee divergence, and up to 69% of within-human diversity, can potentially be explained by variation in *de novo* mutation rate at the 100 kb scale (Smith et al., 2018). It is unclear, however, how well these
65 results generalise to species with distinct genomic features and life history traits. The few studies conducted in non-human organisms relied on proxies of the local mutation rate such as synonymous

diversity or divergence with a closely-related outgroup (Castellano, et al., 2018a; Castellano et al., 2018b). Since these indirect measures of the mutation rate are susceptible to the confounding effect of selection (both direct, *e.g.* codon usage (Lawrie et al., 2013; Machado et al., 2020) and indirect, *e.g.* background selection (Phung et al., 2016)), developing dedicated statistical methods to infer mutation rate variation from polymorphism data is of high interest. Through simultaneous inference of the landscapes of genetic drift, selection, recombination and mutation, confounding factors can be better teased apart and, in a second step, the relative contribution of these micro-evolutionary forces to the distribution of diversity can be more meaningfully quantified.

Disentangling the effects of multiple forces shaping the evolution of DNA sequences is challenging. For example, a genomic region with reduced diversity can be explained by either linked selection, drift, low mutation rate or a combination thereof. Zeng and Jackson developed a likelihood-based framework that jointly infers the effective population size (N_e) (Charlesworth, 2009) and μ with high accuracy (Zeng & Jackson, 2018). However, since it relies on the single-site frequency spectrum, their method is restricted to unlinked loci. While this approach avoids the confounding effect of linkage disequilibrium (Slatkin, 2008), it discards sites in the genome where local variation in the mutation rate may be relevant. In this article we describe a new model to fill in this gap. We have previously described a statistical framework (the integrative sequentially Markovian coalescent, iSMC) that jointly infers the demographic history of the sample and variation in the recombination rate along the genome (Barroso et al., 2019) via a Markov-modulated Markov process. We now further extend this framework to account for sequential variation of the mutation rate. This integration allows statistical inference of variation along the genome in both recombination and mutation rates, as well as in Times to the Most Recent Common Ancestor, that is, the ancestral recombination graph of two haploids (Rosenberg & Nordborg, 2002). Because natural selection disturbs τ away from its distribution under neutrality around functionally constrained regions of the genome (Palamara et al., 2018; Rasmussen et al., 2014; Stern et al.,

2019), iSMC offers estimators of all micro-evolutionary forces and we can further use causal inference (Pearl & Mackenzie, 2018) to simultaneously estimate their effects on π . Our analyses of *Drosophila melanogaster* reveals the impact of linked selection, however, it suggests that the rate of *de novo* mutations is quantitatively the most important factor shaping genetic diversity in this species.

RESULTS:

100 The sequentially Markov coalescent with heterogeneous mutation and recombination

The sequentially Markovian Coalescent (SMC) frames the genealogical process as unfolding spatially along the genome (McVean & Cardin, 2005; Marjoram & Wall, 2006). Its first implementation derives the transition probabilities of genealogies between adjacent sites as a function of the historical variation of N_e and the genome-average scaled recombination rate $\rho = 4 \times N_e \times r$ (Li & Durbin, 2011). Model fitting is achieved by casting the SMC as a hidden Markov model (HMM) (Dutheil, 2017) and letting the emission probabilities be a function of the underlying Time to the Most Recent Common Ancestor (TMRCA, τ) and the scaled mutation rate $\theta = 4 \times N_e \times \mu$ (see Methods). The SMC has proven to be quite flexible and serves as the theoretical basis for several models of demographic inference (Schiffels & Durbin, 2014; Terhorst et al., 2017; Sellinger et al., 2020). We have previously extended this process to account for the variation of ρ along the genome, thereby allowing for a heterogeneous frequency of transitions between local genealogies along the genome (Barroso et al., 2019). In this more general process called iSMC (Dutheil, 2020), recombination rate heterogeneity is captured by an auto-correlation parameter, where the local ρ is taken from a discrete distribution and the transition between recombination rates along the genome follows a first-order Markov process.

In the general case, the iSMC process is a Markov-modulated Markov process that can be cast as a HMM where the hidden states are n -tuples storing all combinations of genealogies and discretised

values of each parameter that is allowed to vary along the genome (Dutheil, 2020). If one such
120 parameter contributes to either the transition or emission probabilities of the HMM, then the
parameters that shape its prior distribution can be optimised, *e.g.* by maximum likelihood (see
Methods). In the iSMC with heterogeneous recombination (ρ -iSMC) the hidden states are 2-tuples
containing pairs of genealogies and recombination rates (Barroso et al., 2019). Here, we extend this
model by allowing the mutation rate to also vary along the genome (**Figure 1**), following an
125 independent Markov process, *i.e.*, letting the hidden states of the HMM be $\{\theta$ -category, ρ -category,
genealogy} triplets. The signal that variation in ρ and θ leaves on the distribution of SNPs is
discernible because their contributions to the likelihood are orthogonal: the recombination and
mutation rates affect the transition and emission probabilities in the forward algorithm of the HMM,
respectively. Parameter optimisation and subsequent posterior decoding is performed as detailed by
130 (Barroso et al., 2019). Under strict neutrality, the inferred θ landscape is an approximation to the
landscape of *de novo* mutations (μ). iSMC can therefore be used to infer genome-wide variation in
mutation rates with single-nucleotide resolution, and statistical noise is reduced by averaging the
posterior estimates of θ within larger genomic windows.

135 In order to increase power, we further extend iSMC to accommodate multiple genomes. In this
augmented model, input genomes are combined in pairs such that the underlying genealogies have a
trivial topology which is reduced to their τ (**Figure 1**). Although under Kingman's Coalescent
(Kingman, 1982) the genealogies of multiple pairs of genomes are not independent, we approximate
and compute the composite log-likelihood of the entire dataset by summing over the "diploid" log-
140 likelihoods, similarly to MSMC2 (Malaspinas et al., 2016). Furthermore, iSMC enforces all
diploids to share their prior distributions of τ , ρ and θ so that multiple sequences provide aggregate
information to our parameter inference; it does not, however, explicitly enforce that they have
common genomic landscapes. Rather, iSMC uses posterior probabilities to reconstruct
recombination and mutation maps separately for each diploid.

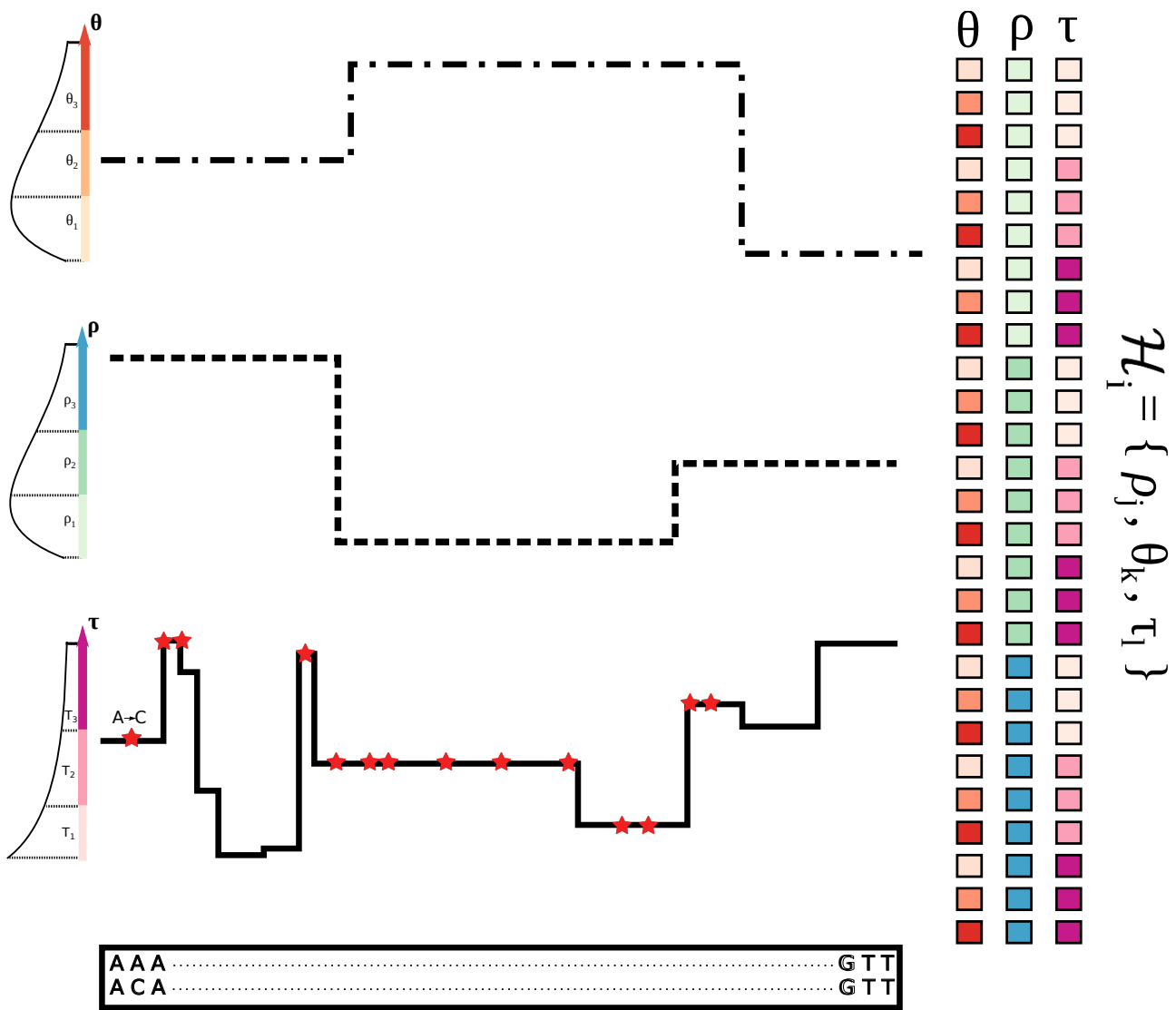


Figure 1. Schematic representation of p-θ-iSMC for one pair of genomes. This cartoon model has three time intervals, three recombination rate categories and three mutation rate categories. The genome-wide distribution of diversity depends on the mutation landscape (top) and on the τ landscape (bottom), which is modulated by the recombination landscape (middle). Discretised values of these distributions (left) are combined in triplets as the hidden states of our Hidden Markov Model (right).

The large variance of the coalescent process results in variation among the posterior landscapes inferred from each pair of genomes since their τ landscapes modulate the amount of information they contain. To obtain a consensus landscape of the whole sample in order to reduce noise, iSMC averages the individual posterior estimates for each site in the genome (see Methods). On the other hand, differences in the τ distribution between diploids primarily reflect variance in the coalescent. We average these individual landscapes to obtain a measure of drift in neutral simulations, noting,

160 however, that the average τ of the sample also contains information about natural selection
(Palamara et al., 2018) – a property we exploit in the analyses of *Drosophila* data.

Mutation rate variation impacts diversity more than linked selection in *Drosophila*

We sought to quantify the determinants of genome-wide diversity in *Drosophila melanogaster* using
165 10 haploid genome sequences from the Zambia population. To infer the genomic landscapes, we
employed a ρ - θ -iSMC model with five mutation rate classes, five recombination rate classes and 30
coalescence intervals, leading to 750 hidden states. Owing to the complexity of the model, which
led to large computation times, we proceeded in two steps: we first estimated model parameters on
a subset of the data (chromosome arm 2L), and then used the fitted model to infer the landscape of
170 mutation, recombination and TMRCA for all autosomes (see Methods). The estimated parameters
suggested an exponential-like distribution of recombination rates ($\hat{\alpha} = \hat{\beta} = 1.03$ for their
Gamma distribution) whereas the inferred distribution of mutation rates was more tightly centered
around the mean ($\hat{\alpha} = \hat{\beta} = 2.93$ for their Gamma distribution). iSMC also inferred that the
change in recombination rate across the genome was more frequent (auto-correlation parameter
175 $\hat{\delta}_\rho \sim 0.9999$, corresponding to a change of recombination rate on average every 10 kb) than the
change in mutation rate (auto-correlation parameter $\hat{\delta}_\theta \sim 0.99999$, corresponding to a change of
mutation rate on average every 100 kb). This suggests that our model is mostly sensitive to factors
that determine large-scale variation in the mutation rate instead of fine-scale sequence motifs such
as highly mutable triplets (DeWitt et al., 2021; Harris & Pritchard, 2017). Our estimated genome-
180 wide average $\hat{\rho}$ (0.036) is in line with previous estimates (Chan et al., 2012), and the coalescence
rates suggest a relatively recent ~ 4 -fold bottleneck followed by a fast recovery. We used these
parameters estimated from *D. melanogaster* to simulate 10 replicate datasets (see Methods). The
aims of these simulations are (1) to benchmark iSMC's accuracy in reconstructing the mutation
landscape; and (2) to understand how the genomic landscapes of ρ , θ and τ interact to influence

185 diversity levels under neutrality, thereby providing a measure of contrast for the analyses of real
data (where natural selection is present).

We first report strong correlations between simulated and inferred maps, ranging from 0.975 to
0.989 (**Table 1, Figure 2**), showing that our model is highly accurate under strict neutrality and
190 when mutation rate varies along the genome in Markovian fashion.

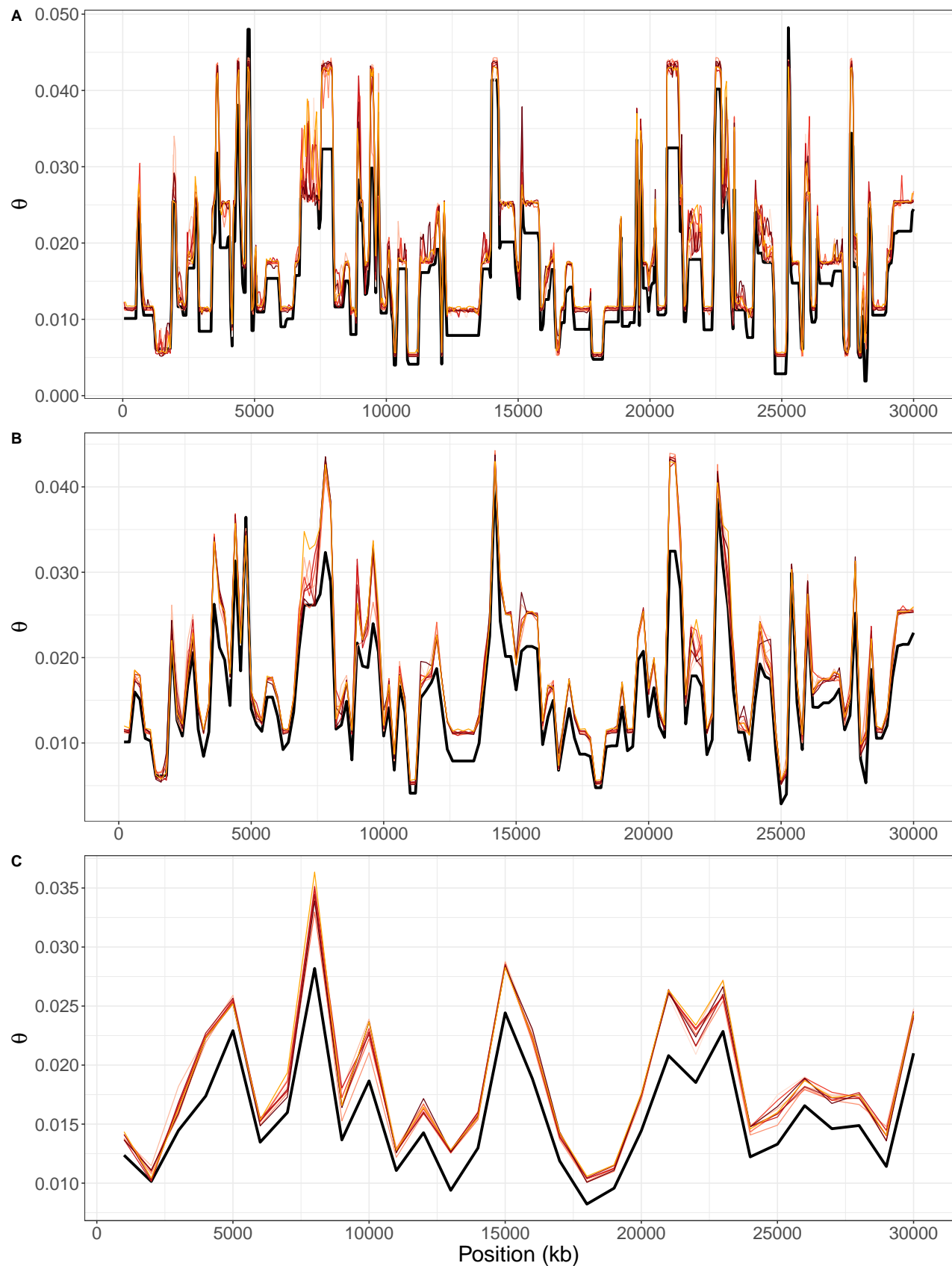
Table 1. Spearman correlations between simulated and inferred mutation maps. All p-values are smaller than $2.2e-16$.

Replicate/Scale	50 kb	200 kb	1 Mb
Replicate 1	0.9780	0.9850	0.9790
Replicate 2	0.9850	0.9890	0.9750
Replicate 3	0.9780	0.9830	0.9810
Replicate 4	0.9780	0.9820	0.9820
Replicate 5	0.9820	0.9870	0.9830
Replicate 6	0.9840	0.9870	0.9810
Replicate 7	0.9800	0.9870	0.9750
Replicate 8	0.9810	0.9880	0.9820
Replicate 9	0.9830	0.9870	0.9810
Replicate 10	0.9840	0.9890	0.9840

195

We then used the genomic landscapes from these simulated datasets to investigate how evolutionary
forces shape the distribution of genetic diversity along the genome, measured as π , the average
heterozygosity of the sample. The structure of our hypothesized causal model of diversity (**Figure**
200 **3**) suggests the absence of "backdoor paths" creating non-causal associations between ρ , θ , τ and π
(Pearl & Mackenzie, 2018). Therefore, we represented our model as an ordinary least squares
regression (OLS) that seeks to explain π as a linear combination of the centered variables ρ , θ and τ ,

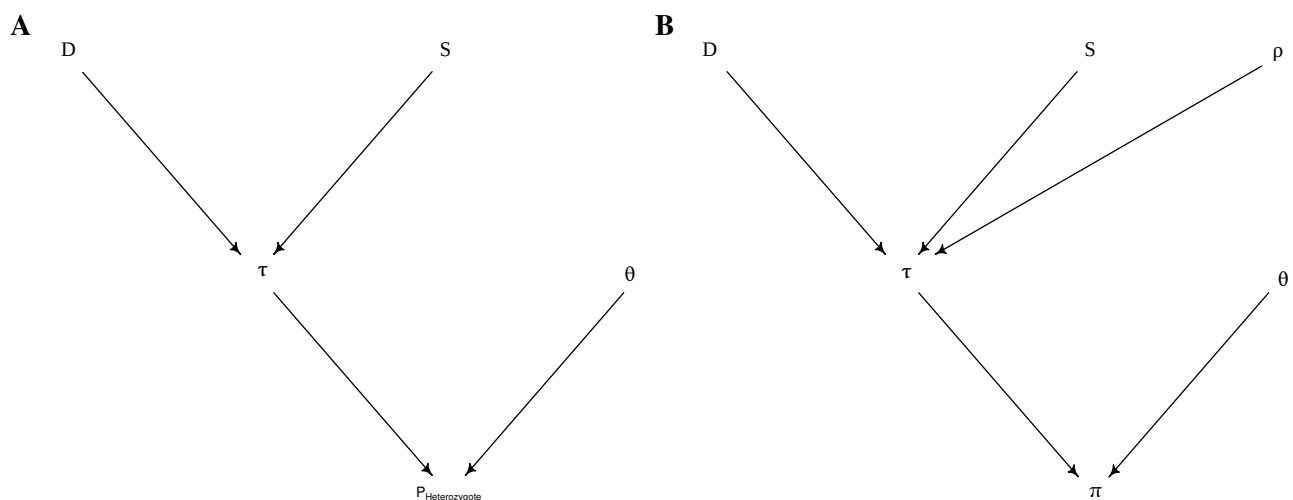
Figure 2. iSMC recovers the mutation landscape in simulations. Simulated values are shown by the thick black line whereas inferred values are shown, for each replicate, by thin lines in shades of red. **A**, 50 kb scale. **B**, 200 kb scale. **C**, 1 Mb scale.



binned in windows of 50 kb, 200 kb and 1 Mb. Since simulations grant direct access to the true genomic landscapes, the ensuing linear model is free of estimation noise in the explanatory variables and serves as a ground truth assessment of how neutral evolutionary forces influence genetic diversity. In all replicates, we found that model selection using Akaike's information criterion favours a regression with an interaction term between the two variables that directly influence π $\pi_i = \beta_1 \cdot \tau_i + \beta_2 \cdot \theta_i + \beta_3 \cdot \rho_i + \beta_4 \cdot \theta : \tau + \epsilon_i$ over the simpler model

$$\pi_i = \beta_1 \cdot \tau_i + \beta_2 \cdot \theta_i + \beta_3 \cdot \rho_i + \epsilon_i .$$

Figure 3. Directed acyclic graphs depicting our causal model for the determinants of genetic diversity. **A:** for a single, hypothetical nucleotide that is independent of any neighbors, its probability of heterozygosity is directly influenced by the local mutation rate (θ) and TMRCA (τ), which in turn is affected by drift (D) and selection (S). **B:** when contiguous sites are grouped into genomic windows, their correlated histories imply that the local recombination rate (ρ) plays a role in determining the variance in drift and the breadth of linked selection via τ , which interacts with local θ to influence π .



Fitting this linear model at the 50 kb, 200 kb and 1 Mb scales shows significant and positive effects of θ and τ , but not of ρ , on π . This is expected since both deeper ancestry and higher mutation rate lead to increased diversity, and the effect of recombination rate on π not only is mediated by τ (thus disappearing after conditioning on it in the linear model), but also should not be felt under neutrality (**Figure 3B**). There is also a significant and positive effect of the interaction between θ and τ , where

the effect of the mutation rate on diversity can only be fully manifested if ancestry is deep enough (reciprocally, ancestry can only be seen if the local mutation rate is high enough). Strikingly, the total variance explained by the model is $> 99\%$ at all scales, suggesting that these three landscapes
235 are sufficient to describe the genome-wide distribution of diversity. To understand the relative contributions of drift, mutation and recombination on local diversity levels, we used type II ANOVA to partition the R^2 contributed by each explanatory variable of our linear model. Our estimates show that the θ landscape explains most of the variance in π in our simulated scenario and that its contribution increases with the genomic scale (96.3% at 50 kb, 98.6% at 200 kb and 99.3% at 1 Mb,
240 yellow squares in **Figure 4A**). On the other hand, the contribution of the τ landscape decreases with the genomic scale (2.7% at 50 kb, 1% at 200 kb and 0.54% at 1 Mb, blue squares in **Figure 4**). These trends stem from the very fine scale of variation in τ (changing on average every 48.42 base pairs due to recombination events in our coalescent simulations, median = 19 base pairs), which smooths out more rapidly when averaged within larger windows. Conversely, the broader scale of
245 heterogeneity in θ (changing every ~ 100 kb) makes it comparatively more relevant at larger windows. We then fitted such OLS models to the same simulated data except using the genomic landscapes as inferred by iSMC. When using the inferred landscapes instead of the true ones the sign and significance of the estimated coefficients of the OLS models remained unchanged, but in some replicates the residuals of the model were found to be correlated and/or with heterogeneous
250 variance. As this violation of the OLS assumption could bias our estimates of the linear coefficients, we also fitted Generalized Least Squares (GLS) models accounting for both effects, which reassuringly produced consistent results. Although co-linearity between $\hat{\tau}$ and $\hat{\theta}$ arises due to confounding in their estimation by iSMC, the variance inflation factors are always < 5 , indicating that the coefficients are robust to this effect (Ferré, 2009). Type II ANOVA using the inferred
255 landscapes shows that the contribution of $\hat{\tau}$ is slightly higher than using the true landscapes (5.1%, 2.9% and 1.4%, increasing window size) whereas the contribution of $\hat{\theta}$ is slightly lower (92.5%, 95.4% and 97.5%, increasing window size), but the variance explained by

each variable largely agrees between the two cases (**Figure 4**). We conclude that the joint-inference approach of iSMC can infer the genomic landscapes of τ , ρ and θ and that our linear model can
260 quantify their effect on the distribution of nucleotide diversity.

We finally employed the landscapes obtained with ρ - θ -iSMC to quantify the determinants of genome-wide diversity in *D. melanogaster*. We used our inferred maps to fit an OLS regression of the form $\pi_i = \beta_1 \cdot \hat{\tau}_i + \beta_2 \cdot \hat{\theta}_i + \beta_3 \cdot \hat{\rho}_i + \beta_4 \cdot \hat{\theta}_i + \epsilon_i$. As in our simulations, the regression model shows
265 positive effects of both τ and θ , but not of ρ , on π across all scales (**Table 2**). Likewise, a GLS model yields the same trends, and the variance inflation factors are < 5 , indicating that the estimated coefficients are robust to co-linearity (Ferré, 2009). Partitioning of variance shows a small contribution of $\hat{\tau}$ that decreases with increasing genomic scale (5.9% at 50 kb, 2.1% at 200 kb and 2.1% at 1 Mb) whereas the opposite is true for $\hat{\theta}$ (91.7% at 50 kb, 96.7% at 200 kb and
270 96.8% at 1 Mb, **Figure 4**). Our linear model explains $> 99\%$ of the variation in π along *D. melanogaster* autosomes, and the effect of our inferred landscapes on diversity are remarkably close to those from our neutral simulations (**Figure 4**), suggesting that iSMC is robust to the occurrence of selection in this system. Unlike in our neutral simulations, however, in real *Drosophila* data the simple correlation test between $\hat{\rho}$ and π yields a positive and significant estimate (Spearman's rho
275 = 0.20, p-value = $2e-13$ at the 50 kb scale; Spearman's rho = 0.15, p-value = 0.0025 at the 200 kb scale; Spearman's rho = 0.20, p-value = 0.07 at the 1 Mb scale), recapitulating the classic result of Begun and Aquadro (Begun & Aquadro, 1992) and indicating the presence of linked selection. We also found a positive correlation between $\hat{\rho}$ and $\hat{\tau}$ (Spearman's rho = 0.48, p-value $< 2.2e-16$ at 50 kb; Spearman's rho = 0.45, p-value $< 2.2e-16$ at 200 kb; Spearman's rho = 0.48, p-value $< 2.2e-16$ at 1 Mb), once again contrasting the results under neutrality and suggesting that the effect of
280 linked selection is indeed captured by the distribution of genealogies and modulated by the recombination rate (Cutter & Payseur, 2013). Although in SMC-based methods τ is primarily influenced by demography (fluctuating population sizes explicitly by use of a Coalescent prior

taming the transition probabilities of the HMM (Li & Durbin, 2011; Schiffels & Durbin, 2014) and
 285 population structure implicitly (Beichman et al., 2018)), it has also been demonstrated to carry the
 signature of selection (Palamara et al., 2018) due to local changes in coalescence rates. We tested
 the sensitivity of our framework to this effect by fitting linear models without $\hat{\tau}$ as an explanatory
 variable $\pi_i = \beta_1 \cdot \hat{\theta}_i + \beta_2 \cdot \hat{\rho}_i + \epsilon_i$, hypothesizing that in this case the recombination rate would have a
 significant and positive effect on diversity. Indeed, this is what we found at all genomic scales
 290 (Table 2), corroborating our interpretation of the causal relationships among ρ , N_e and π in the
 presence of selection (Figure 3B). In summary, our results show that recombination shapes
 diversity via the τ distribution and linked selection, but that in *D. melanogaster*, the impact of
 selection on the diversity landscape is smaller than that of mutation rate variation.

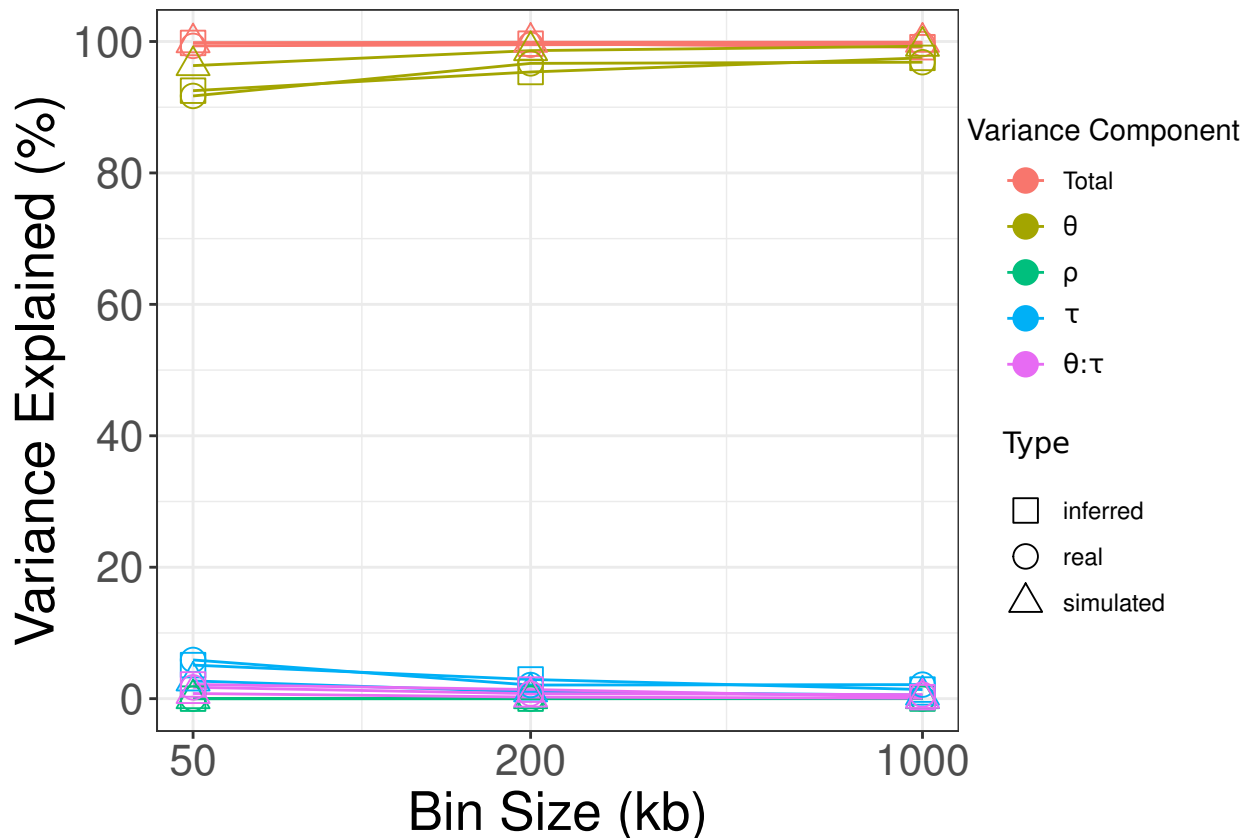
295 **Table 2. Coefficient estimates of linear regression models to explain the genome-wide distribution of diversity in *Drosophila melanogaster*.**

	50 kb scale		200 kb scale		1 Mb scale	
	OLS		OLS		OLS	
Variable	Coefficient	p-value	Coefficient	p-value	Coefficient	p-value
Theta	0.9746	<2.2e-16	0.9767	<2.2e-16	0.9903	<2.2e-16
Tau	0.0114	<2.2e-16	0.0151	<2.2e-16	0.0098	<2.2e-16
Rho	0.0018	0.0237	0.0016	0.2580	0.0049	0.0748
Theta:Tau	1.0510	<2.2e-16	0.9481	<2.2e-16	0.5321	<2.2e-16
	GLS		GLS		GLS	
Variable	Coefficient	p-value	Coefficient	p-value	Coefficient	p-value
Theta	0.9730	<1e-4	0.9729	<1e-4	0.9816	<1e-4
Tau	0.0115	<1e-4	0.0114	<1e-4	0.0106	<1e-4
Rho	0.0014	0.0700	0.0018	0.2300	0.0005	0.8658
Theta:Tau	1.0777	<1e-4	0.9323	<1e-4	0.6656	<1e-4
	GLS		GLS		GLS	
Variable	Coefficient	p-value	Coefficient	p-value	Coefficient	p-value
Theta	1.0991	<1e-4	1.0796	<1e-4	1.0840	<1e-4
Rho	0.0529	<1e-4	0.0601	<1e-4	0.0444	<1e-4

To investigate the signature of selection, we analysed the relationship between the local mutation rate and the levels of synonymous (π_s) and non-synonymous (π_N) diversity across *D. melanogaster* genes (see Methods). We computed these summary statistics across exons and matched their coordinates with our 50 kb-scale genomic landscapes to increase resolution (*i.e.*, to maximize variation in mutation and recombination rates among genes). We observed a stronger relationship between $\hat{\theta}$ and π_s (Spearman's rho = 0.68, 95% CI after 10,000 bootstrap replicates = [0.64, 0.72], partial correlation accounting for $\hat{\tau}$) than between $\hat{\theta}$ and π_N (Spearman's rho = 0.27, 95% CI after 10,000 bootstrap replicates = [0.22, 0.32], partial correlation accounting for $\hat{\tau}$) indicating that selection partially purges the excess of deleterious variants in genes with elevated mutation rate, whereas synonymous variants segregate more freely either because they are not directly affected by selection (but are still linked to selected sites) or because selection on codon usage (Lawrie et al., 2013; Machado et al., 2020) is not as strong as selection on protein function. Since synonymous sites are interdigitated with non-synonymous sites, the contrast between their correlation tests cannot be explained by a bias in $\hat{\theta}$ in functionally constrained regions of the genome. Furthermore, a correlation test between $\hat{\theta}$ and the proportion of exonic sites in the same 50 kb windows (Spearman's rho = -0.037, p-value = 0.19, partial correlation accounting for $\hat{\tau}$) fails to reveal a bias in our inference of the mutation rate in regions under stronger background selection. Conversely, we observed a negative and significant correlation between $\hat{\tau}$ and the proportion of exonic sites (Spearman's rho = -0.158, p-value = 2.03e-12, partial correlation accounting for $\hat{\theta}$), as expected since stronger background selection in coding regions should reduce the TMRCA more sharply (Palamara et al., 2018). To estimate the effect of these factors on more heavily constrained diversity, we fitted linear models considering only 50 kb windows with more than 20,000 coding sites. Once again, there were significant and positive effects of both $\hat{\tau}$ and $\hat{\theta}$, but not of $\hat{\rho}$, on π . Moreover, the mutation landscape remains the most important factor, explaining 93.2% of the distribution of diversity in gene-rich regions.

Figure 4. Variance in the distribution of diversity explained by each genomic landscape.

Partitioning of variance according to bin size (x-axis, shown in log10 scale), using either simulated data (*true* landscapes: triangles; *inferred* landscapes: squares) or real *Drosophila* data (circles). Colors represent explanatory variables in the linear model: θ (yellow), ρ (green), τ (blue), $\theta:\tau$ interaction (pink) and the total variance explained by the model (red).



Mutation rate variation shapes genome-wide diversity in several evolutionary scenarios

Our analyses of *Drosophila* data and *Drosophila*-inspired simulations suggest that the mutation

landscape is by far the most important factor influencing levels of diversity along the genome. But

are there scenarios where τ has a more pronounced effect on π ? We addressed this question by

exploring the parameter space of our neutral simulations. For fixed values of the long-term average

population size ($N_e = 100,000$), the average mutation rate per site per generation ($\mu = 2e-09$), the

Gamma distribution of scaling factors of θ ($\alpha = \beta = 2.5$) and the Gamma distribution of scaling

factors of ρ ($\alpha = \beta = 1.0$), we varied the demographic history (flat N_e ; 10-fold bottleneck 0.5

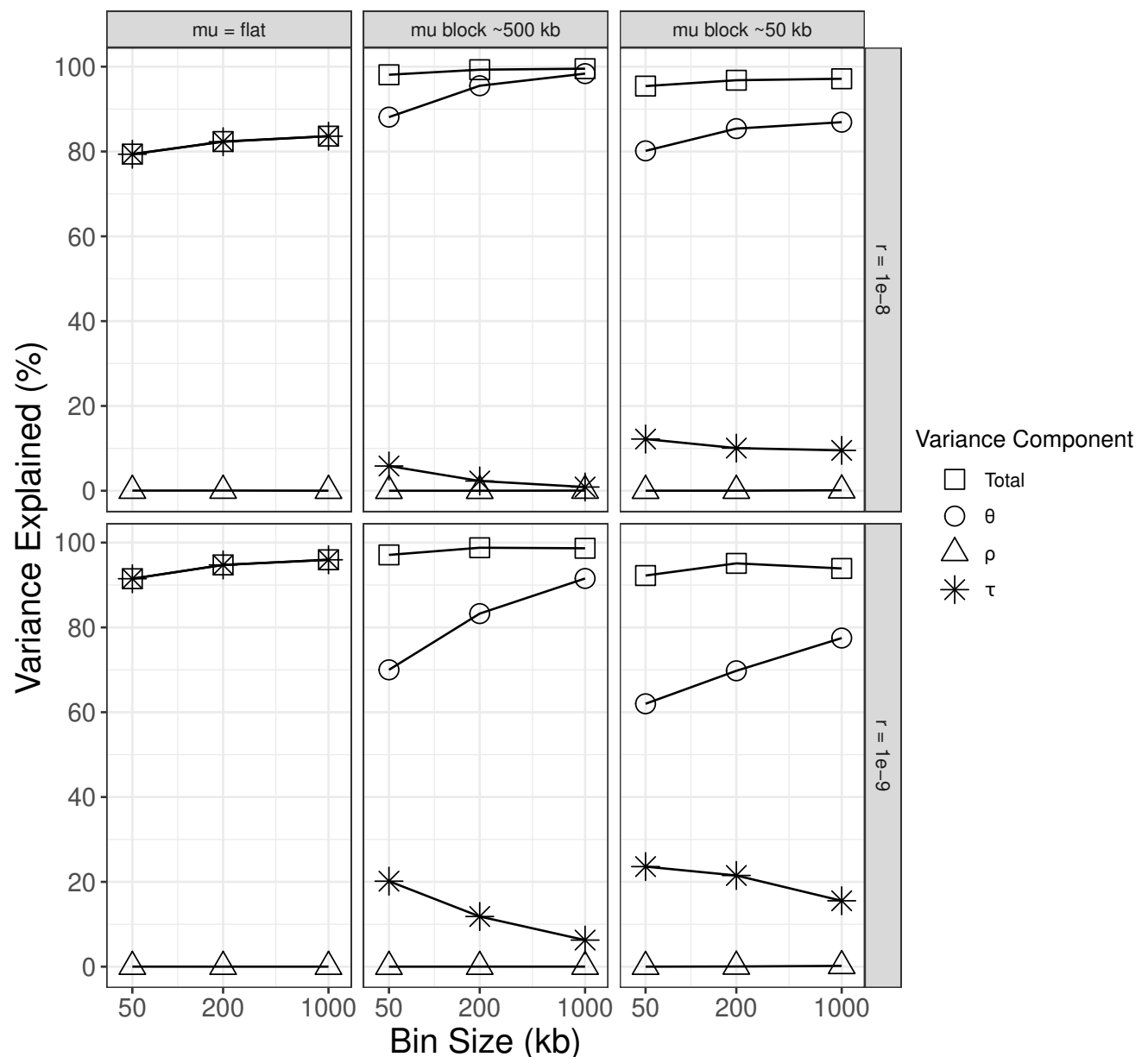
coalescent time units ago), the average recombination rate per site per generation ($r = 1e-08$; $1e-09$)

and the average length of genomic blocks of constant θ (50 kb; 500 kb; flat mutation landscape).

We reasoned that the extent of the variation in τ along the genome relative to that of θ should modulate their relative influence on π . We fitted OLS models to explain π using the true simulated landscapes as explanatory variables, and computed their average R^2 over all replicates for each evolutionary scenario (**Figure 5, 6**). The OLS models included an interaction term between θ and τ but its individual R^2 was excluded from the plots because it is overall low ($\sim 1\%$) and of no direct interest. We observed clear trends emerging from these simulated data. First, for a given demographic history and pattern of variation in the mutation rate, increasing r reduces the influence of τ on π . This happens because with high recombination rates the genealogies change more often along the genome, thus displaying more homogeneous maps when averaged within windows. Second, for a given r and pattern of variation in the mutation rate, τ has a larger impact on π in the bottleneck scenario compared to the scenario of constant population size. This happens because when N_e varies in time, the distribution of coalescence times may become multi-modal (Hein et al., 2004) and therefore more heterogeneous along the genome. Third, for a given demographic history and r , frequent changes in θ along the genome (on average every 50 kb) reduce its impact on π relative to rare changes in θ (on average every 500 kb). This happens because frequent changes in θ lead to it being more homogeneous along the genome, for the window sizes used in our analyses. Finally, if the mutation landscape is flat, then as expected the variance explained by our linear model is entirely attributed to τ . Note that although in our neutral simulations τ varies along the genome as a result of genetic drift alone, it still has a non-negligible effect on the distribution of diversity in most scenarios (*i.e.*, binning does not lead to completely flat landscapes of genetic drift). This is in agreement with an observation that heterogeneous recombination rates lead to outliers in genome-wide F_{ST} scans, even under neutrality (Booker et al., 2020), which happens because the recombination landscape enlarges the variance of the τ distribution by making the frequency of genealogy transitions change according to the local ρ . From a practical standpoint, it means that drift should not be neglected as an explanation for the distribution of π , especially at narrow window sizes (≤ 10 kb).

Figure 5. Variance explained by each genomic landscape under constant Ne scenarios.

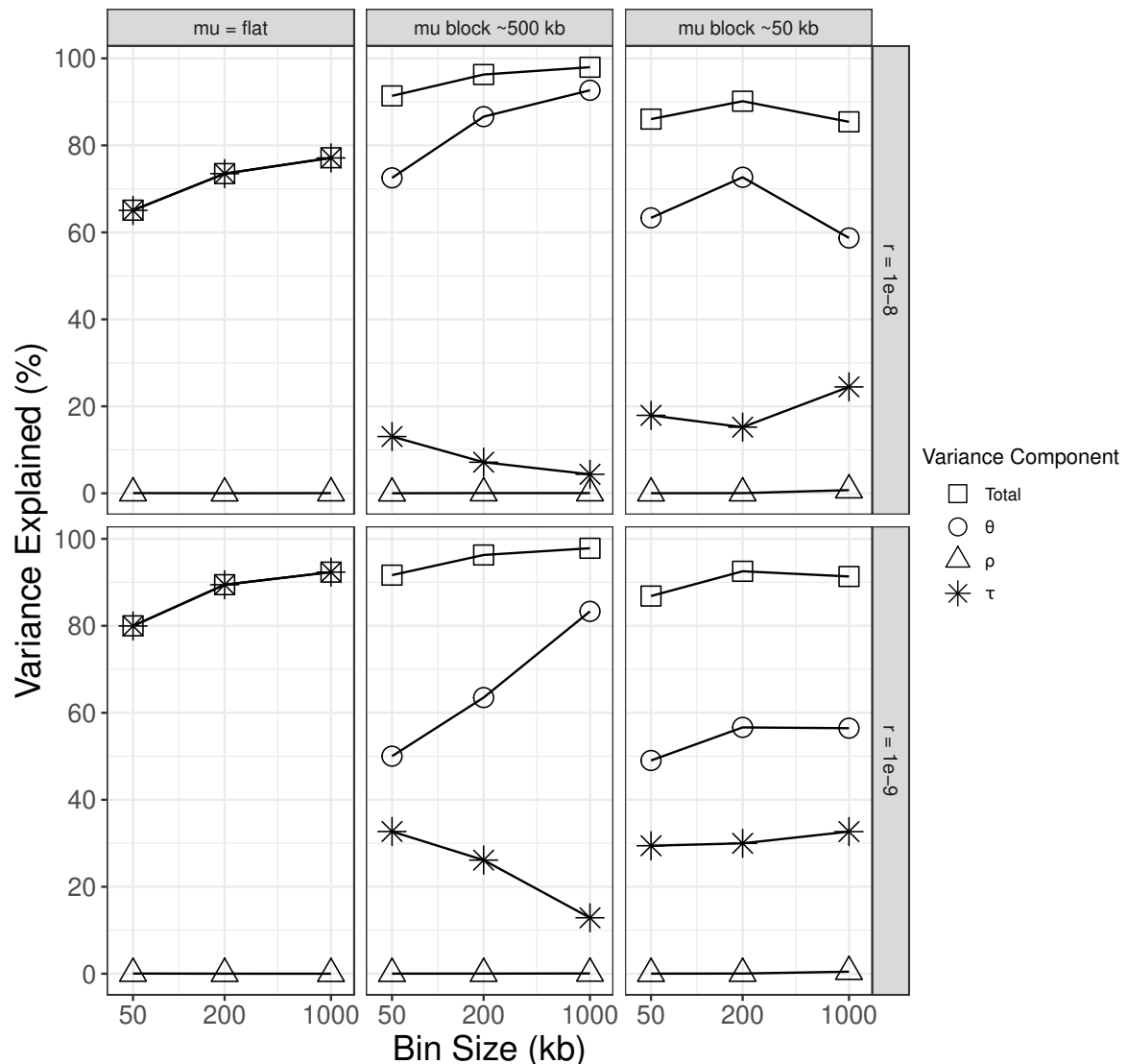
370 Partitioning of variance according to bin size and genomic parameters (rows = recombination rate, columns = scale of mutation rate variation). X-axis is shown in log(10) scale.



375 More generally, our simulation study shows that the relative impacts of evolutionary forces on π depend primarily on (1) the joint patterns of variation of τ and θ along the genome; and (2) the scale of the analysis, due to averaging parameter estimates within genomic windows. In light of these results, the genome of *D. melanogaster* – with its high effective recombination rate, broad pattern of variation in the mutation rate and the relatively smooth demographic history of the Zambia population – seems to be particularly susceptible to the effect of the mutation landscape on its

distribution of diversity. Yet, since the mutation landscape stood out as the most relevant factor in all of the explored scenarios where it was allowed to vary even slightly (**Figure 5, 6**), we predict that it is very likely to shape genome-wide diversity patterns in other species as well.

Figure 6. Variance explained by each genomic landscape under population bottleneck scenarios. Partitioning of variance according to bin size and genomic parameters (rows = recombination rate, columns = scale of mutation rate variation). X-axis is shown in log10 scale.



DISCUSSION:

The relative strengths of selection and drift in shaping patterns of nucleotide diversity have been debated for several decades (reviewed in Hey, 1999; and, more recently, Jensen et al., 2019; Kern & Hahn, 2018), with the contribution of local mutation rate only recently brought to light (Castellano, Eyre-Walker, et al., 2018; Harpak et al., 2016; Smith et al., 2018). We were able to employ our

extended iSMC model to jointly infer mutation, recombination and TMRCA landscapes and to use causal inference to estimate their impact on π along the genome. Our analyses revealed that these combined landscapes explained >99% of the distribution of diversity along the *Drosophila* genome; when looking into the detailed patterns, we found the footprints of linked selection, but the major driver of genome-wide diversity in this species seems to be the mutation landscape. This does not imply that selection cannot extend beyond the 18.3% of the *Drosophila* genome that is exonic (Alexander et al., 2010), but rather that variation in the mutation rate is strong enough to contribute relatively more to the variation in π , in the genomic scales of analyses that we conducted. Our findings, however, contrast with an estimate by Comeron that 70% of the distribution of diversity in *Drosophila* can be explained by background selection at the 100 kb scale (Comeron, 2014), where the author further argued that many regions of increased diversity may be experiencing balancing selection. Instead, we propose that mutation rate variation is responsible for at least some of these effects. We believe that such discrepancy between the results is partially due to the B-value statistic computed in (Comeron, 2014) not accounting for heterogeneity in the mutation rate along the genome, contrary to its original application in humans (McVicker et al., 2009). It is also conceivable that selection is not only manifested as distortions in the distribution of genealogies (the τ landscape) but also biases our estimate of the mutation landscape. However, based on the high similarity between real *Drosophila* data and our neutral simulations (Figure 4), we argue that the bias induced by selection is unlikely to overturn our conclusion of a major impact of the mutation landscape on the distribution of diversity. We also note that selection should have a stronger impact on π when binning is performed at smaller genomic scales (≤ 10 kb), which we have not explored because estimation noise with iSMC is typically high at such small window sizes. At larger scales, a putative explanation for our results is that the reduction in τ caused by linked selection could be relatively uniform across the chromosomes of species with a compact genome and high effective recombination rate. In summary, our results provide evidence that similarly to humans (Harpak et al., 2016; Smith et al., 2018), the mutation landscape is a crucial driver of the

distribution of diversity in the fruit fly. Our simulation study (**Figure 5, 6**) further suggests that in many evolutionary scenarios the mutation landscape will be the most relevant factor shaping π along the genome. Future studies using integrative models like the ones we introduced here and
425 applied to species with distinct genomic features and life-history traits will help elucidate how often – and by how much – the mutation landscape stands out as the main driver of the diversity landscape.

It is important to note that our results do not directly argue in favor of either genetic drift or natural
430 selection in the classic population genetics debate. Instead, they highlight the importance of a third element – the mutation landscape – in shaping patterns of genetic diversity. Nevertheless, the mutation landscape may impact the dynamics of natural selection by modulating the rate of mutations in genes according to their position in the genome. Consequently, levels of selective interference, genetic load and rates of adaptation should vary accordingly (Castellano et al., 2016).
435 In *D. melanogaster*, we inferred a ~10-fold change in mutation rate at the 50 kb scale, meaning that the impact of mutation rate variation on selective processes can be substantial. These results open intriguing lines of inquiry. First, under what conditions can the shape of the mutation landscape itself be selected for? For example, it has been shown that modifiers of the global mutation rate are under selection to reduce genetic load (Lynch, 2008; Lynch et al., 2016). It remains to be seen
440 whether the position of genes or genomic features correlated with local mutation rate (*e.g.*, replication timing (Francioli et al., 2015)) can likewise be optimised. Second, how conserved is the mutation landscape across species? Analogous work on the recombination landscape has revealed overall fast evolution in mammals and has helped uncover the molecular architecture responsible for the placement of double-strand breaks (Berg et al., 2011; Jabbari et al., 2019); it will be
445 interesting to test whether mutation events follow a similar pattern, now that the impact of various sequence motifs on μ is being more thoroughly investigated (DeWitt et al., 2021; Kim et al., 2021).

There is an ongoing discussion about incorporating complex demography and background selection into the null model of molecular evolution (Comeron, 2014, 2017; Johri et al., 2020), motivated by the goal of providing more sensible expectations for rigorously testing alternative scenarios. Our results suggest that a more realistic null model should also include variation in the mutation rate along the genome. By doing so, genome-wide scans (*e.g.*, looking for regions with reduced diversity as candidates for selective sweeps) may become less susceptible to false negatives (in regions of high mutation rate) and positives (in regions of low mutation rate), paving the way to more robust inference (Booker et al., 2017; Haasl & Payseur, 2016).

METHODS:

Modelling spatial variation in θ

We now introduce our approach to modelling the mutation landscape starting from the original pair-wise SMC process. Because iSMC models pairs of genomes, the genealogies underlying each orthologous site can be conveniently summarized by τ , the time to their most recent common ancestor (Li & Durbin, 2011; Schiffels & Wang, 2020). The pair of DNA sequences is described as a binary sequence where 0 represents homozygous states and 1 represents heterozygous states (thus, once genome pairs are established, phasing information is discarded). The probability of observing 0 or 1 at any given position of the genome depends only on τ and the population mutation rate θ . If the hidden state configuration of the model excludes variation in the mutation rate, then θ is assumed to be a global parameter such that the emission probabilities of homozygous and heterozygous states can be compute for every site as $P(0|\tau) = e^{(-\theta \times \tau)}$, and $P(1|\tau) = 1 - e^{(-\theta \times \tau)}$ respectively, as presented by Li & Durbin (2011).

470

To incorporate spatial heterogeneity in the mutation rate, we set the genome-wide average θ as the average number of pair-wise differences between haplotypes, and modulate it by drawing scaling factors from a discretised Gamma distribution with mean equal to 1. The parameters shaping this

prior distribution are estimated by maximum likelihood together with other parameters of the
 475 model. We model the changes in mutation rate along the genome as a Markov process with a single
 parameter δ_θ , the transition probability between any class of mutation rate, independent of the
 genealogical process. The justification for the Markov model is that sites in close proximity are
 expected to have similar mutation rates, for example, as is the case when the efficiency of the
 replication machinery decreases with increasing distance from the start of the replication fork
 480 (Francioli et al., 2015). Let $n^{(\tau)}$ be the number of discretised τ intervals, and $n^{(\theta)}$ be the number
 of discretised categories of the prior distribution of scaling factors of θ . The ensuing Markov-
 modulated HMM has $n = n^{(\tau)} \times n^{(\theta)}$ hidden states. The transition matrix for spatial variation in θ is:

$$Q_\theta = \begin{bmatrix} P_{11} & P_{12} & \cdots & P_{1n^{(\theta)}} \\ P_{21} & P_{22} & \cdots & P_{2n^{(\theta)}} \\ \vdots & \vdots & \ddots & \vdots \\ P_{n^{(\theta)}1} & P_{n^{(\theta)}2} & \cdots & P_{n^{(\theta)}n^{(\theta)}} \end{bmatrix} = \begin{bmatrix} 1-\delta & \frac{\delta_\theta}{n^{(\theta)}-1} & \cdots & \frac{\delta_\theta}{n^{(\theta)}-1} \\ \frac{\delta_\theta}{n^{(\theta)}-1} & 1-\delta_\theta & \cdots & \frac{\delta_\theta}{n^{(\theta)}-1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\delta_\theta}{n^{(\theta)}-1} & \frac{\delta_\theta}{n^{(\theta)}-1} & \cdots & 1-\delta_\theta \end{bmatrix}$$

485 where δ_θ is the auto-correlation parameter. The resulting process is a combination of the SMC and
 the mutation Markov model, so that its transition probabilities are functions of the parameters from
 both processes, that is, the coalescence rates (parameterized by splines), δ_θ and the global
 recombination rate ρ (Barroso et al., 2019). The forward recursion for this model at genomic
 490 position i can be written as:

$$F_i(\tau_t, \theta_m) = \left(\sum_{k=1}^{n^{(\theta)}} \left(\sum_{j=1}^{n^{(\tau)}} F_{i-1}(\tau_j, \theta_k) \cdot Pr(\tau_j \rightarrow \tau_t) \cdot Pr(\theta_k \rightarrow \theta_m) \right) \right) \cdot Pr(\tau_t \rightarrow S_i | \theta_m) \quad (1)$$

where θ_m is the product of the genome-wide average mutation rate and the value of the m -th
 495 discretised category drawn from its prior Gamma distribution. The emission probability of binary

state S_i depends on the height of the t -th genealogy and the focal mutation rate θ_m . More

specifically, the emission probabilities of θ -iSMC are $P(0|\tau_t, \theta_m) = e^{(-\theta_m \times \tau_t)}$, and

$P(1|\tau_t, \theta_m) = 1 - e^{(-\theta_m \times \tau_t)}$. The forward recursion integrates over all $n^{(\theta)}$ categories of θ and over

all $n^{(\tau)}$ intervals of τ . In the double-modulated model, where both mutation and recombination are

500 allowed to vary along the genome, this integration is performed over θ , τ as well as ρ (giving a total

of $n^{(\tau)} \times n^{(\theta)} \times n^{(\rho)}$ hidden states). Since spatial variation in ρ contributes to the transition

probability between genealogies, the forward recursion is now given by:

$$F_i(\tau_t, \theta_m, \rho_r) = \left(\sum_{l=1}^{n^{(\theta)}} \left(\sum_{k=1}^{n^{(\rho)}} \left(\sum_{j=1}^{n^{(\tau)}} F_{i-1}(\tau_j, \rho_k, \theta_l) \cdot Pr(\tau_j \rightarrow \tau_t | \rho_k) \cdot Pr(\theta_l \rightarrow \theta_m) \cdot Pr(\rho_k \rightarrow \rho_r) \right) \right) \right) \cdot Pr(\tau_t \rightarrow S_i | \theta_m)$$

505 (2)

To obtain the single-nucleotide landscapes for a given diploid individual, we first compute the

posterior probability of each hidden state for every site i in the genome using standard HMM

procedures (Durbin et al., 1998). Afterwards, since in ρ - θ -iSMC the hidden states are triplets

510 (**Figure 1**), computing the posterior average of each variable of interest – ρ , θ or τ – amounts to first

marginalising the probability distribution of its categories and then using it to weight the

corresponding discretised values (Barroso et al., 2019). Let \hat{M} be the inferred discretised Gamma

distribution shaping mutation rate variation, and $\hat{\theta}_l$ be the product of the estimated genome-wide

average mutation rate $\hat{\theta}_0$ and the value of \hat{M} inside category l . Similarly, let \hat{R} be the inferred

515 discretised Gamma distribution shaping recombination rate variation, and $\hat{\rho}_k$ be the product of the

estimated genome-wide average recombination rate $\hat{\rho}_0$ and the value of \hat{R} inside category k .

Then the posterior average $\hat{\theta}$ at position i is given by:

$$\hat{\theta}_i = \hat{\theta}_0 \cdot \sum_{l=1}^{n^{(\theta)}} m_l \cdot \left(\sum_{k=1}^{n^{(\rho)}} \sum_{j=1}^{n^{(\tau)}} P_i(\theta_l, \rho_k, \tau_j) \right) \quad (3)$$

520

where $P_i(\theta_l, \rho_k, \tau_j)$ is the probability of the triplet $\{\theta_l, \rho_k, \tau_j\}$ underlying the i -th site. Likewise, the posterior average $\hat{\rho}$ at position i is given by:

$$\hat{\rho}_i = \hat{\rho}_0 \cdot \sum_{k=1}^{n^{(\rho)}} r_k \cdot \left(\sum_{l=1}^{n^{(\theta)}} \sum_{j=1}^{n^{(\tau)}} P_i(\theta_l, \rho_k, \tau_j) \right) \quad (4)$$

525

Finally, the posterior average $\hat{\tau}$ at position i is presented in units of $4 \times Ne$ generations and obtained with:

$$\hat{\tau}_i = \sum_{j=1}^{n^{(\tau)}} \hat{\tau}_j \cdot \left(\sum_{k=1}^{n^{(\rho)}} \sum_{l=1}^{n^{(\theta)}} P_i(m_l, r_k, \tau_j) \right) \quad (5)$$

530

We can then bin the inferred single-nucleotide landscapes into windows of length L by averaging our site-specific estimates over all sites within each window. A consensus map of the population is obtained by further averaging over all n individual (binned) maps in our sample, *i.e.*:

$$\hat{\theta}_{pop}^L = \frac{1}{n} \sum_{j=1}^n \left(\frac{1}{L} \sum_{i=1}^L \hat{\theta}_{i,j} \right) \quad (6)$$

535

is our estimate of the consensus mutation rate in a single genomic window of length L , where n is the number of pairs of genomes analysed by iSMC, and likewise for ρ and τ :

$$\hat{\rho}_{pop}^L = \frac{1}{n} \sum_{j=1}^n \left(\frac{1}{L} \sum_{i=1}^L \hat{\rho}_{i,j} \right) \quad (7)$$

540

$$\hat{\tau}_{pop}^L = \frac{1}{n} \sum_{j=1}^n \left(\frac{1}{L} \sum_{i=1}^L \hat{\tau}_{i,j} \right) \quad (8)$$

Simulation study

Using SCRM (Staab et al., 2015), we simulated 10 haploid sequences of length 30 Mb with
 545 parameters identical to those inferred from ρ - θ -iSMC in *Drosophila melanogaster*: $\theta = 0.0112$; $\rho = 0.036$; α (continuous Gamma distribution used as mutation rate prior) = 3.0; α (continuous Gamma distribution used as recombination rate prior) = 1.0; δ_θ (mutation rate transition probability) = 1e-05; δ_ρ (recombination rate transition probability) = 1e-04. **Figures 5 and 6** display the mean R^2 value over 10 replicates, but the standard deviation of these estimates are very small and confidence
 550 intervals were therefore omitted. All scripts necessary to reproduce the analyses and figures, as well as supplementary figures are available at https://github.com/gvbarroso/ismc_dm_analyses

Data analyses

Model fitting and posterior decoding by ρ - θ -iSMC in *Drosophila melanogaster* data was performed
 555 using a hidden-states configuration of 30 τ intervals, five ρ categories and five θ categories. We used publicly available data – haplotypes ZI103, ZI117, ZI161, ZI170, ZI179, ZI191, ZI129, ZI138, ZI198 and ZI206 coming from the Zambia population in the *Drosophila* Population Genomics Project Phase 3 (Lack et al., 2015). Gaps and unknown nucleotides in these FASTA sequences were assigned as missing data. For each scale in which the landscapes were binned (50 kb, 200 kb and 1
 560 Mb), we filtered out windows with more than 10% missing data in the resulting maps. To optimize computational time, ρ - θ -iSMC was first fitted to chromosome 2L only. Maximum likelihood estimates from this model were then used to perform posterior decoding on all other autosomes. Genomic coordinates for coding sequences and their summary statistics (π_N , and π_S) were taken from (Moutinho et al., 2019) . Data analyses procedures (starting from the inferred iSMC
 565 landscapes) such as building linear models and testing Gauss-Markov assumptions are detailed in the script `dm_analyses.Rmd`, available in the GitHub repository:

https://github.com/gvbarroso/ismc_dm_analyses/r_scripts/analyses. Intermediate tables are provided in the FigShare repository: [10.6084/m9.figshare.13164320](https://figshare.com/10.6084/m9.figshare.13164320).

570 **Linear modelling**

When building linear models from real data, we first fitted GLS models independently to each autosome arm (2L, 2R, 3L, 3R), correcting for both auto-correlation of and heteroscedasticity of the residuals. After using Bonferroni correction for multiple testing, we observed (across the autosome arms and for different window sizes) significant and positive effects of θ and τ on π , whereas the effect of ρ was only significant for chromosome 3L at the 200 kb scale, and the interaction between θ and τ is positive and significant except for arms 2R and 3L at the 1 Mb scale (**Table S1**). Since the trends in coefficients are overall consistent, we pulled the autosome arms and in the Results section we presented linear models fitted to the entire genome, for ease of exposition. Because we cannot rely on the GLS to partition the variance explained by each variable using type II ANOVA, we used OLS models to compute R^2 and restricted the GLS to assess the sign and significance of variables. We centered all explanatory variables before fitting the regression models to aid in both computation of variance inflation factors and interpretation of the coefficients (due to the interaction term between θ and τ). The procedures and intermediate results outlined here are detailed in the script `dm_analyses.Rmd`, available in the GitHub repository:

585 https://github.com/gvbarroso/ismc_dm_analyses/r_scripts/analyses.

Data and software availability

The iSMC software package and source code is freely available at

<https://github.com/gvbarroso/iSMC>

590 Scripts used to generate our results can be found in https://github.com/gvbarroso/ismc_dm_analyses

Data required to reproduce the results are deposited in FigShare under the DOI

[10.6084/m9.figshare.13164320](https://figshare.com/10.6084/m9.figshare.13164320)

REFERENCES:

- Alexander, R. P., Fang, G., Rozowsky, J., Snyder, M., & Gerstein, M. B. (2010). Annotating non-coding regions of the genome. *Nature Reviews Genetics*, 11(8), 559–571.
<https://doi.org/10.1038/nrg2814>
- Barroso, G. V., Puzović, N., & Dutheil, J. Y. (2019). Inference of recombination maps from a single pair of genomes and its application to ancient samples. *PLOS Genetics*, 15(11), e1008449.
<https://doi.org/10.1371/journal.pgen.1008449>
- Begun, D. J., & Aquadro, C. F. (1992). Levels of naturally occurring DNA polymorphism correlate with recombination rates in *D. melanogaster*. *Nature*, 356(6369), 519–520.
<https://doi.org/10.1038/356519a0>
- Beichman, A. C., Huerta-Sanchez, E., & Lohmueller, K. E. (2018). Using Genomic Data to Infer Historic Population Dynamics of Nonmodel Organisms. *Annual Review of Ecology, Evolution, and Systematics*, 49(1), 433–456. <https://doi.org/10.1146/annurev-ecolsys-110617-062431>
- Berg, I. L., Neumann, R., Sarbajna, S., Odenthal-Hesse, L., Butler, N. J., & Jeffreys, A. J. (2011). Variants of the protein PRDM9 differentially regulate a set of human meiotic recombination hotspots highly active in African populations. *Proceedings of the National Academy of Sciences of the United States of America*, 108(30), 12378–12383.
<https://doi.org/10.1073/pnas.1109531108>
- Besenbacher, S., Hvilsom, C., Marques-Bonet, T., Mailund, T., & Schierup, M. H. (2019). Direct estimation of mutations in great apes reconciles phylogenetic dating. *Nature Ecology & Evolution*, 1. <https://doi.org/10.1038/s41559-018-0778-x>
- Booker, T. R., Jackson, B. C., & Keightley, P. D. (2017). Detecting positive selection in the genome. *BMC Biology*, 15(1), 98. <https://doi.org/10.1186/s12915-017-0434-y>
- Booker, T. R., Yeaman, S., & Whitlock, M. C. (2020). Variation in recombination rate affects detection of outliers in genome scans under neutrality. *Molecular Ecology*, mec.15501.
<https://doi.org/10.1111/mec.15501>
- Buffalo, V. (2021). Why do species get a thin slice of π ? Revisiting Lewontin’s Paradox of Variation. *BioRxiv*, 2021.02.03.429633. <https://doi.org/10.1101/2021.02.03.429633>

- Casillas, S., & Barbadilla, A. (2017). Molecular Population Genetics. *Genetics*, 205(3), 1003–1035.
<https://doi.org/10.1534/genetics.116.196493>
- Castellano, D., Coronado-Zamora, M., Campos, J. L., Barbadilla, A., & Eyre-Walker, A. (2016). Adaptive Evolution Is Substantially Impeded by Hill–Robertson Interference in *Drosophila*. *Molecular Biology and Evolution*, 33(2), 442–455. <https://doi.org/10.1093/molbev/msv236>
- Castellano, D., Eyre-Walker, A., & Munch, K. (2018). Impact of mutation rate and selection at linked sites on fine-scale DNA variation across the homininae genome. *BioRxiv*, 452201.
<https://doi.org/10.1101/452201>
- Castellano, D., James, J., & Eyre-Walker, A. (2018). Nearly neutral evolution across the *Drosophila melanogaster* genome. *Molecular Biology and Evolution*, 35, 2685–2694.
<https://doi.org/10.1093/molbev/msy164>
- Chan, A. H., Jenkins, P. A., & Song, Y. S. (2012). Genome-Wide Fine-Scale Recombination Rate Variation in *Drosophila melanogaster*. *PLOS Genetics*, 8(12), e1003090.
<https://doi.org/10.1371/journal.pgen.1003090>
- Charlesworth, B. (2009). Effective population size and patterns of molecular evolution and variation. *Nature Reviews Genetics*, 10(3), 195–205. <https://doi.org/10.1038/nrg2526>
- Charlesworth, B., & Charlesworth, D. (2016). Population genetics from 1966 to 2016. *Heredity*, February, 1–8. <https://doi.org/10.1038/hdy.2016.55>
- Comeron, J. M. (2014). Background Selection as Baseline for Nucleotide Variation across the *Drosophila* Genome. *PLOS Genetics*, 10(6), e1004434.
<https://doi.org/10.1371/journal.pgen.1004434>
- Comeron, J. M. (2017). Background selection as null hypothesis in population genomics: Insights and challenges from *Drosophila* studies. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 372(1736), 20160471. <https://doi.org/10.1098/rstb.2016.0471>
- Cutter, A. D., & Payseur, B. A. (2013). Genomic signatures of selection at linked sites: Unifying the disparity among species. *Nature Reviews. Genetics*, 14(4), 262–274.
<https://doi.org/10.1038/nrg3425>
- DeWitt, W. S., Harris, K. D., Ragsdale, A. P., & Harris, K. (2021). Nonparametric coalescent inference of mutation spectrum history and demography. *Proceedings of the National Academy of Sciences*, 118(21). <https://doi.org/10.1073/pnas.2013798118>

- Durbin, R., Eddy, S. R., Krogh, A., & Mitchison, G. (1998). *Biological sequence analysis: Probabilistic models of proteins and nucleic acids*. Cambridge University Press.
<https://doi.org/10.1017/CBO9780511790492>
- Dutheil, J. Y. (2017). Hidden Markov Models in Population Genomics. *Methods in Molecular Biology (Clifton, N.J.)*, 1552, 149–164. https://doi.org/10.1007/978-1-4939-6753-7_11
- Dutheil, J. Y. (2020). Towards more realistic models of genomes in populations: The Markov-modulated sequentially Markov coalescent. *ArXiv:2010.08359 [q-Bio]*.
<http://arxiv.org/abs/2010.08359>
- Ellegren, H., & Galtier, N. (2016). Determinants of genetic diversity. *Nature Reviews. Genetics*, 17(7), 422–433. <https://doi.org/10.1038/nrg.2016.58>
- Elyashiv, E., Sattath, S., Hu, T. T., Strutsovsky, A., McVicker, G., Andolfatto, P., Coop, G., & Sella, G. (2016). A Genomic Map of the Effects of Linked Selection in *Drosophila*. *PLOS Genetics*, 12(8), e1006130. <https://doi.org/10.1371/journal.pgen.1006130>
- Ferré, J. (2009). Regression Diagnostics. In *Comprehensive Chemometrics* (Vol. 3, pp. 33–89). Elsevier. <https://www.sciencedirect.com/topics/mathematics/variance-inflation-factor>
- Francioli, L. C., Polak, P. P., Koren, A., Menelaou, A., Chun, S., Renkens, I., Consortium, G. of the N., Duijn, C. M. van, Swertz, M., Wijmenga, C., Ommen, G. van, Slagboom, P. E., Boomsma, D. I., Ye, K., Guryev, V., Arndt, P. F., Kloosterman, W. P., Bakker, P. I. W. de, & Sunyaev, S. R. (2015). Genome-wide patterns and properties of *de novo* mutations in humans. *Nature Genetics*, 47(7), 822–826. <https://doi.org/10.1038/ng.3292>
- Galtier, N., & Rousselle, M. (2020). How Much Does Ne Vary Among Species? *Genetics*, 216(2), 559–572. <https://doi.org/10.1534/genetics.120.303622>
- Haasl, R. J., & Payseur, B. A. (2016). Fifteen years of genomewide scans for selection: Trends, lessons and unaddressed genetic sources of complication. *Molecular Ecology*, 25(1), 5–23. <https://doi.org/10.1111/mec.13339>
- Harpak, A., Bhaskar, A., & Pritchard, J. K. (2016). Mutation Rate Variation is a Primary Determinant of the Distribution of Allele Frequencies in Humans. *PLOS Genetics*, 12(12), e1006489. <https://doi.org/10.1371/journal.pgen.1006489>
- Harris, K., & Pritchard, J. K. (2017). Rapid evolution of the human mutation spectrum. *ELife*, 6, e24284. <https://doi.org/10.7554/eLife.24284>

- Hein, J., Schierup, M., & Wiuf, C. (2004). *Gene Genealogies, Variation and Evolution: A primer in coalescent theory*. Oxford University Press.
- Hey, J. (1999). The neutralist, the fly and the selectionist. *Trends in Ecology & Evolution*, 14(1), 35–38. [https://doi.org/10.1016/s0169-5347\(98\)01497-9](https://doi.org/10.1016/s0169-5347(98)01497-9)
- Hubisz, M. J., Williams, A. L., & Siepel, A. (2020). Mapping gene flow between ancient hominins through demography-aware inference of the ancestral recombination graph. *PLOS Genetics*, 16(8), e1008895. <https://doi.org/10.1371/journal.pgen.1008895>
- Hudson, R. R. (1983). Properties of a neutral allele model with intragenic recombination. *Theoretical Population Biology*, 23(2), 183–201. [https://doi.org/10.1016/0040-5809\(83\)90013-8](https://doi.org/10.1016/0040-5809(83)90013-8)
- Hudson, R. R., & Kaplan, N. L. (1988). The coalescent process in models with selection and recombination. *Genetics*, 120(3), 831–840.
- Jabbari, K., Wirtz, J., Rauscher, M., & Wiehe, T. (2019). A common genomic code for chromatin architecture and recombination landscape. *PLOS ONE*, 14(3), e0213278. <https://doi.org/10.1371/journal.pone.0213278>
- Jensen, J. D., Payseur, B. A., Stephan, W., Aquadro, C. F., Lynch, M., Charlesworth, D., & Charlesworth, B. (2019). The importance of the Neutral Theory in 1968 and 50 years on: A response to Kern and Hahn 2018. *Evolution*, 73(1), 111–114. <https://doi.org/10.1111/evo.13650>
- Johri, P., Charlesworth, B., & Jensen, J. D. (2020). Toward an Evolutionarily Appropriate Null Model: Jointly Inferring Demography and Purifying Selection. *Genetics*, 215(1), 173–192. <https://doi.org/10.1534/genetics.119.303002>
- Jónsson, H., Sulem, P., Arnadottir, G. A., Pálsson, G., Eggertsson, H. P., Kristmundsdottir, S., Zink, F., Kehr, B., Hjorleifsson, K. E., Jensson, B. Ö., Jonsdottir, I., Marelsson, S. E., Gudjonsson, S. A., Gylfason, A., Jonasdottir, A., Jonasdottir, A., Stacey, S. N., Magnusson, O. T., Thorsteinsdottir, U., ... Stefansson, K. (2018). Multiple transmissions of de novo mutations in families. *Nature Genetics*, 50(12), 1674–1680. <https://doi.org/10.1038/s41588-018-0259-9>

- Kaplan, N., & Hudson, R. R. (1985). The use of sample genealogies for studying a selectively neutral m-loci model with recombination. *Theoretical Population Biology*, 28(3), 382–396. [https://doi.org/10.1016/0040-5809\(85\)90036-X](https://doi.org/10.1016/0040-5809(85)90036-X)
- Kern, A. D., & Hahn, M. W. (2018). The Neutral Theory in Light of Natural Selection. *Molecular Biology and Evolution*, 35(6), 1366–1371. <https://doi.org/10.1093/molbev/msy092>
- Kim, Y.-A., Leiserson, M. D. M., Moorjani, P., Sharan, R., Wojtowicz, D., & Przytycka, T. M. (2021). Mutational Signatures: From Methods to Mechanisms. *Annual Review of Biomedical Data Science*, 4(1), 189–206. <https://doi.org/10.1146/annurev-biodatasci-122320-120920>
- Kimura, M. (1968). Evolutionary Rate at the Molecular Level. *Nature*, 217(5129), 624–626. <https://doi.org/10.1038/217624a0>
- Kingman, J. (1982). Genealogy Populations. *Journal of Applied Probability*, 19(1982), 27–43.
- Lack, J. B., Cardeno, C. M., Crepeau, M. W., Taylor, W., Corbett-Detig, R. B., Stevens, K. A., Langley, C. H., & Pool, J. E. (2015). The Drosophila Genome Nexus: A Population Genomic Resource of 623 Drosophila melanogaster Genomes, Including 197 from a Single Ancestral Range Population. *Genetics*, 199(4), 1229–1241. <https://doi.org/10.1534/genetics.115.174664>
- Lawrie, D. S., Messer, P. W., Hershberg, R., & Petrov, D. A. (2013). Strong Purifying Selection at Synonymous Sites in D. melanogaster. *PLOS Genetics*, 9(5), e1003527. <https://doi.org/10.1371/journal.pgen.1003527>
- Lewontin, R. C. (1974). *The Genetic Basis of Evolutionary Change*. Columbia University Press.
- Li, H., & Durbin, R. (2011). Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357), 493–496. <https://doi.org/10.1038/nature10231>
- Lynch, M. (2008). The Cellular, Developmental and Population-Genetic Determinants of Mutation-Rate Evolution. *Genetics*, 180(2), 933–943. <https://doi.org/10.1534/genetics.108.090456>
- Lynch, M., Ackerman, M. S., Gout, J.-F., Long, H., Sung, W., Thomas, W. K., & Foster, P. L. (2016). Genetic drift, selection and the evolution of the mutation rate. *Nature Reviews Genetics*, 17(11), 704–714. <https://doi.org/10.1038/nrg.2016.104>
- Machado, H. E., Lawrie, D. S., & Petrov, D. A. (2020). Pervasive Strong Selection at the Level of Codon Usage Bias in Drosophila melanogaster. *Genetics*, 214(2), 511–528. <https://doi.org/10.1534/genetics.119.302542>

- Malaspinas, A.-S., Westaway, M. C., Muller, C., Sousa, V. C., Lao, O., Alves, I., Bergström, A., Athanasiadis, G., Cheng, J. Y., Crawford, J. E., Heupink, T. H., Macholdt, E., Peischl, S., Rasmussen, S., Schiffels, S., Subramanian, S., Wright, J. L., Albrechtsen, A., Barbieri, C., ... Willerslev, E. (2016). A Genomic History of Aboriginal Australia. *Nature*, 1–20.
<https://doi.org/10.1038/nature18299>
- Marjoram, P., & Wall, J. D. (2006). Fast “coalescent” simulation. *BMC Genetics*, 7(1), 16–16.
<https://doi.org/10.1186/1471-2156-7-16>
- McVean, G. A. T., & Cardin, N. J. (2005). Approximating the coalescent with recombination. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1459), 1387–1393. <https://doi.org/10.1098/rstb.2005.1673>
- McVicker, G., Gordon, D., Davis, C., & Green, P. (2009). Widespread Genomic Signatures of Natural Selection in Hominid Evolution. *PLOS Genetics*, 5(5), e1000471.
<https://doi.org/10.1371/journal.pgen.1000471>
- Moutinho, A. F., Trancoso, F. F., & Dutheil, J. Y. (2019). The Impact of Protein Architecture on Adaptive Evolution. *Molecular Biology and Evolution*, 36(9), 2013–2028.
<https://doi.org/10.1093/molbev/msz134>
- Ohta, T. (1992). The Nearly Neutral Theory of Molecular Evolution. *Annual Review of Ecology and Systematics*, 23(1), 263–286. <https://doi.org/10.1146/annurev.es.23.110192.001403>
- Palamara, P. F., Terhorst, J., Song, Y. S., & Price, A. L. (2018). High-throughput inference of pairwise coalescence times identifies signals of selection and enriched disease heritability. *Nature Genetics*, 50(9), 1311–1317. <https://doi.org/10.1038/s41588-018-0177-x>
- Pearl, J., & Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect* (1st ed.). Basic Books, Inc.
- Phung, T. N., Huber, C. D., & Lohmueller, K. E. (2016). Determining the Effect of Natural Selection on Linked Neutral Divergence across Species. *PLOS Genetics*, 12(8), e1006199.
<https://doi.org/10.1371/journal.pgen.1006199>
- Pouyet, F., & Gilbert, K. J. (2020). Towards an improved understanding of molecular evolution: The relative roles of selection, drift, and everything in between. *PEER COMMUNITY IN EVOLUTIONARY BIOLOGY*, 22.

- Rasmussen, M. D., Hubisz, M. J., Gronau, I., & Siepel, A. (2014). Genome-Wide Inference of Ancestral Recombination Graphs. *PLoS Genetics*, 10(5).
<https://doi.org/10.1371/journal.pgen.1004342>
- Rosenberg, N. A., & Nordborg, M. (2002). Genealogical Trees, Coalescent Theory and the Analysis of Genetic Polymorphisms. *Nature Reviews Genetics*, 3(5), 380–390.
<https://doi.org/10.1038/nrg795>
- Schiffels, S., & Durbin, R. (2014). Inferring human population size and separation history from multiple genome sequences. *Nature Genetics*, 46(8), 919–925.
<https://doi.org/10.1038/ng.3015>
- Schiffels, S., & Wang, K. (2020). MSMC and MSMC2: The Multiple Sequentially Markovian Coalescent. In J. Y. Dutheil (Ed.), *Statistical Population Genomics* (pp. 147–166). Springer US. https://doi.org/10.1007/978-1-0716-0199-0_7
- Schrider, D. R., Ayroles, J., Matute, D. R., & Kern, A. D. (2018). Supervised machine learning reveals introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLOS Genetics*, 14(4), e1007341. <https://doi.org/10.1371/journal.pgen.1007341>
- Sellinger, T. P. P., Awad, D. A., Moest, M., & Tellier, A. (2020). Inference of past demography, dormancy and self-fertilization rates from whole genome sequence data. *PLOS Genetics*, 16(4), e1008698. <https://doi.org/10.1371/journal.pgen.1008698>
- Slatkin, M. (2008). Linkage disequilibrium—Understanding the evolutionary past and mapping the medical future. *Nature Reviews Genetics*, 9(6), 477–485. <https://doi.org/10.1038/nrg2361>
- Smith, T. C. A., Arndt, P. F., & Eyre-Walker, A. (2018). Large scale variation in the rate of germ-line de novo mutation, base composition, divergence and diversity in humans. *PLoS Genetics*, 14(3), e1007254. <https://doi.org/10.1371/journal.pgen.1007254>
- Staab, P. R., Zhu, S., Metzler, D., & Lunter, G. (2015). Scrm: Efficiently simulating long sequences using the approximated coalescent with recombination. *Bioinformatics*, 31(10), 1680–1682. <https://doi.org/10.1093/bioinformatics/btu861>
- Stankowski, S., Chase, M. A., Fuiten, A. M., Rodrigues, M. F., Ralph, P. L., & Streisfeld, M. A. (2019). Widespread selection and gene flow shape the genomic landscape during a radiation of monkeyflowers. *PLOS Biology*, 17(7), e3000391.
<https://doi.org/10.1371/journal.pbio.3000391>

- Stern, A. J., Wilton, P. R., & Nielsen, R. (2019). An approximate full-likelihood method for inferring selection and allele frequency trajectories from DNA sequence data. *PLOS Genetics*, 15(9), e1008384. <https://doi.org/10.1371/journal.pgen.1008384>
- Terhorst, J., Kamm, J. A., & Song, Y. S. (2017). Robust and scalable inference of population history from hundreds of unphased whole-genomes. *Nature Genetics*, 49(2), 303–309. <https://doi.org/10.1038/ng.3748>
- Zeng, K., & Jackson, B. C. (n.d.). Methods for estimating demography and detecting between-locus differences in the effective population size and mutation rate. *Molecular Biology and Evolution*. <https://doi.org/10.1093/molbev/msy212>