

1 Inferring the landscape of 2 recombination using recurrent 3 neural networks

4 Jeffrey R. Adrion^{1,†}, Jared G. Galloway^{1,†}, Andrew D. Kern¹

*For correspondence:
jadrion@uoregon.edu

5 ¹Institute of Ecology and Evolution, University of Oregon

†These authors contributed equally
to this work

7 **Abstract** Accurately inferring the genome-wide landscape of recombination rates in natural
8 populations is a central aim in genomics, as patterns of linkage influence everything from genetic
9 mapping to understanding evolutionary history. Here we describe ReLERNN, a deep learning
10 method for accurately estimating a genome-wide recombination landscape using as few as four
11 samples. Rather than use summaries of linkage disequilibrium as its input, ReLERNN considers
12 columns from a genotype alignment, which are then modeled as a sequence across the genome
13 using a recurrent neural network. We demonstrate that ReLERNN improves accuracy and reduces
14 bias relative to existing methods and maintains high accuracy in the face of demographic model
15 misspecification. We apply ReLERNN to natural populations of African *Drosophila melanogaster* and
16 show that genome-wide recombination landscapes, while largely correlated among populations,
17 exhibit important population-specific differences. Lastly, we connect the inferred patterns of
18 recombination with the frequencies of major inversions segregating in natural *Drosophila*
19 populations.

21 Introduction

22 Recombination plays an essential role in the meiotic production of gametes in most sexual species,
23 and is often required for proper pairing and segregation of chromosomes (*Hunter et al., 2006*;
24 *Mather, 1938*; *Smith and Nicolas, 1998*). During meiotic recombination, double-strand breaks are
25 resolved as crossover or non-crossover recombination events along the chromosome, and as
26 such, homologous chromosomes can exchange genetic information (reviewed in *Kirkpatrick, 2010*;
27 *Zelkowski et al., 2019*). Thus while recombination is often critical to development and reproduction,
28 it also has profound effects on both evolutionary and population genomics (*Burt, 2000*; *Felsenstein,*
29 *1974*; *Haenel et al., 2018*; *Hartfield and Otto, 2011*; *Hill and Robertson, 1966*; *Kondrashov, 1982*).

30 Indeed, the population recombination rate $\rho = 4Nr$ is a central parameter in population and
31 statistical genetics (reviewed in *Hahn, 2018*), as ρ largely determines patterns of linkage disequi-
32 librium (LD) across the genome. In regions of the genome where ρ is relatively small we expect
33 increased levels of LD, and conversely in genomic compartments with high ρ we expect little LD.
34 Deviations from our expected levels of LD given the local recombination rate can be illustrative of
35 the influence of other evolutionary forces such as selection or migration. For example, selective
36 sweeps are expected to dramatically elevate LD near the target of selection (*Kim and Nielsen, 2004*;
37 *O'Reilly et al., 2008*; *Parsch et al., 2001*).

38 Structural variation itself is expected to modulate the landscape of recombination along the chro-
39 mosomes, as both crossovers and non-crossovers are predicated on the alignment of homologous
40 sequences, and structural rearrangements may directly impact those alignments. Chromosomal
41 inversions, long-known to suppress crossing over along a chromosome (e.g. *Sturtevant, 1921*), are

perhaps the most well-studied example of such structural variation. Inversion polymorphisms have been implicated in diverse evolutionary phenomena including local adaptation (*Ayala et al., 2013; Kirkpatrick and Barton, 2006; Lowry and Willis, 2010*), reproductive isolation (*Ayala et al., 2013; Noor et al., 2001; Rieseberg, 2001*), and the maintenance of meiotic drive complexes (*Jaenike, 2001; Presgraves et al., 2009*). As suppressors of recombination, we expect *a priori* that segregating inversions should show distinct histories of recombination in comparison to standard karyotype chromosomes.

While recombination plays a central role in meiosis and reproduction, the frequency and distribution of crossovers along the chromosomes are themselves phenotypes that can evolve (reviewed in *Kirkpatrick, 2010; Ritz et al., 2017*). Importantly, recombination rate variation exists between species, among sexes of the same species (males generally having shorter maps than females), and extends even between individuals of the same sex (*Kong et al., 2010; Singh et al., 2013; Winckler et al., 2005*). Yet while there is abundant variation in the rate of recombination within and between taxa, most methods for accurately measuring this variation involve painstaking experiments or large pedigrees. Thus genetics, as a field, would like to have a tool for directly estimating recombination rates from sequence data, without relying on pedigree genotyping or other ancillary information.

Accordingly, there is a rich history of estimating ρ in population genetics, including efforts to obtain minimum bounds on the number of recombination events (*Hudson and Kaplan, 1985; Myers and Griffiths, 2003; Wu, 2002*), methods of moments estimators (*Hudson, 1987; Wakeley, 1997*), composite likelihood estimators (*Chan et al., 2012; Hudson, 2002; McVean et al., 2002*), and summary likelihood estimators (*Li and Stephens, 2003; Wall, 2000*). Recently, supervised machine learning methods for estimating ρ have entered the fray (*Gao et al., 2016; Lin et al., 2013*), and have proven to be competitive in accuracy with state-of-the-art composite likelihood methods such as LDhat (*McVean et al., 2002*), often with far less computing effort.

To this end, we sought to develop a novel method for inferring rates of recombination directly from a sequence alignment through the use of deep learning. In recent years deep artificial neural networks (ANNs) have produced remarkable performance gains in computer vision (*Krizhevsky et al., 2012; Szegedy et al., 2015*), speech recognition (*Hinton et al., 2012*), natural language processing (*Sutskever et al., 2014*), and data preprocessing tasks such as denoising (*Vincent et al., 2008*). Perhaps most illustrative of the potential of deep learning is the remarkable success of convolutional neural networks (CNNs; *Lecun et al., 1998*) on problems in image analysis. For example, prior to the introduction of CNNs to the annual ImageNet Large Scale Visual Recognition Challenge (*Krizhevsky et al., 2012*), no method had achieved an error rate of less than 25% on the ImageNet data set. In the years that followed, CNNs succeeded in reducing this error rate below 5%, exceeding human accuracy on the same tasks (*Russakovsky et al., 2015*).

In this study we focus our efforts on recurrent neural networks (RNNs), a promising network architecture for population genomics, which has proven adept for analyzing sequential data of arbitrary lengths (*Graves et al., 2013*). Unlike other machine learning methods, deep learning approaches do not require a predefined feature vector. When fed labeled training data (e.g. a set of haplotypes simulated under a known recombination rate), these methods algorithmically create their own set of informative statistics that prove most effective for solving the specified problem. By training deep learning networks directly on sequence alignments, we allow the neural network to automatically extract informative features from the data without human supervision. Learning directly from a sequence alignment for population genetic inference has recently been shown to be possible using CNNs (*Chan et al., 2018; Flagel et al., 2018*), and as we show below, is also true for RNNs.

Here we introduce **Recombination Landscape Estimation using Recurrent Neural Networks**, an RNN-based method for estimating the genomic landscape of recombination rates directly from a phased genotype alignment. We found that ReLERNN is both highly accurate and out-performs competing methods at small sample sizes. We also show that ReLERNN retains its high accuracy in

93 the face of demographic model misspecification. We then apply ReLERNN to population genomic
94 data from African samples of *Drosophila melanogaster*. We demonstrate that the landscape of
95 recombination is largely conserved in this species, yet individual regions of the genome show
96 marked population-specific differences. Finally, we found that chromosomal inversion frequencies
97 directly impact the inferred rate of recombination, and we demonstrate that the role for inversions
98 in suppressing recombination extends far beyond the inversion breakpoints themselves.

99 Results

100 ReLERNN: an accurate method for estimating the genome-wide recombination 101 landscape

102 We developed ReLERNN, a new deep learning method for accurately predicting genome-wide
103 per-base recombination rates from as few as four phased chromosomes. Briefly, ReLERNN provides
104 an end-to-end inferential pipeline for estimating a recombination landscape from a population
105 sample: it takes as input a user-filtered Variant Call Format (VCF) file of phased genotypes, and from
106 this estimates a set of simulation parameters reflective of the input samples. ReLERNN then uses
107 the coalescent simulation program, msprime (Kelleher *et al.*, 2016), to simulate training, validation,
108 and test data sets under either a user-supplied or an inferred demographic history, seeking to
109 mimic population genetic properties of the empirical samples. ReLERNN trains a specific type
110 of RNN, known as a Gated Recurrent Unit (GRU), to predict the per-base recombination rate for
111 these simulations, using only the raw genotype matrix and a vector of genomic coordinates for
112 each simulation example (Figure 1). It then uses this trained network to estimate genome-wide
113 per-base recombination rates for empirical samples using a sliding-window approach. ReLERNN
114 can optionally estimate 95% confidence intervals around each prediction using a parametric boot-
115 strapping approach, and it uses these bootstrap estimates to correct for inherent biases in the
116 training process (see Materials and Methods; Figure 1–Figure Supplement 1).

117 A key feature of ReLERNN's network architecture is the bidirectional GRU layer (Figure 1 inlay),
118 which takes advantage of the sequential nature of genomic data. While vanilla (feed-forward)
119 networks use as input a full block of data for each example, recurrent layers break sequence
120 data into time steps, and iterate over them sequentially. This process allows the gradient descent
121 algorithm, known as backpropagation through time, to share parameters across time steps as well
122 as make inferences based on the ordering of SNPs—i.e. to have a memory of allelic associations. The
123 bidirectional attribute of the GRU layer simply means that each example is duplicated and reversed,
124 so the sequence data are analyzed from both directions and then merged by concatenation.

125 Performance on Simulated Chromosomes

126 As a proof of principle, we performed coalescent simulations using msprime (Kelleher *et al.*, 2016) to
127 generate whole chromosome samples using a fine scale genetic map estimated from *D. melanogaster*
128 (Comeron *et al.*, 2012). We then used ReLERNN to estimate the landscape of recombination
129 for these examples. ReLERNN is able to predict the per-base recombination landscape along a
130 simulated chromosome to a high degree of accuracy across a wide range of realistic parameter
131 values, assumptions, and sample sizes ($R^2 \geq 0.82$; Mean absolute error (MAE) $\leq 1.28 \times 10^{-8}$).
132 Importantly, the accuracy of ReLERNN is only modestly diminished when comparing predictions
133 based on 20 samples ($R^2 = 0.93$; $MAE = 3.72 \times 10^{-9}$; Figure 2) to those based on four samples
134 ($R^2 = 0.82$; $MAE = 6.66 \times 10^{-9}$; Figure 2–Figure Supplement 1). While ReLERNN retains accuracy at
135 small sample sizes, it exhibits somewhat greater sensitivity to both the assumed per-base mutation
136 rate and the assumed maximum ratio of ρ to the population mutation parameter, θ —two mandatory
137 assumptions.

138 To assess the degree of sensitivity to these mutation rate assumptions, we ran ReLERNN on
139 simulations using an assumed per-base mutation rate both 50% greater and 50% less than the
140 simulated (true) mutation rate. In both scenarios, ReLERNN predicts crossover rates that are highly

141 correlated with the simulated rates ($R^2 > 0.91$). However, in both scenarios MAE is inflated but still
142 modest, and the absolute rates of recombination are underpredicted ($R^2 = 0.91$; $MAE = 1.23 \times 10^{-8}$;
143 **Figure 2–Figure Supplement 2**) and overpredicted ($R^2 = 0.94$; $MAE = 1.28 \times 10^{-8}$; **Figure 2–Figure**
144 **Supplement 3**) when assuming a mutation rate less than or greater than the true per-base mutation
145 rate, respectively. Together these results suggest that ReLERNN is in fact learning information about
146 the ratio of crossovers to mutations, and while ReLERNN is highly robust to errant assumptions when
147 predicting relative recombination rates within a genome, caution must be taken when comparing
148 absolute rates between organisms with large differences in per-base mutation rate estimates.

149 **ReLERNN compares favorably to competing methods, especially for small sample** 150 **sizes and under model misspecification**

151 To assess the accuracy of ReLERNN relative to existing methods, we took a comparative approach
152 whereby we made predictions on the same set of simulated test chromosomes using methods
153 that differ broadly in their approaches. Specifically, we chose to compare ReLERNN against two
154 types of machine learning methods—a boosted regression method, FastEPRR (**Gao et al., 2016**),
155 and a convolutional neural network (CNN) recently described in **Flagel et al. (2018)**—and LDhat
156 (**McVean et al., 2002**), a widely cited approximate-likelihood method. We independently simu-
157 lated 10^5 chromosomes using msprime (**Kelleher et al., 2016**) (parameters: $n \in \{4, 8, 16, 32, 64\}$,
158 $priorLowsRho = 0.0$, $priorHighsRho = 5e^{-8} \times 1.25$, $priorLowsMu = 2.5e^{-8} \times 0.75$, $priorHighsMu =$
159 $2.5e^{-8} \times 1.25$, $ChromosomeLength = 3e^5$). Half of these were simulated under demographic equilib-
160 rium and half were simulated under a realistic demographic model (based on the out-of-Africa
161 expansion of European humans; see Materials and Methods). We show that ReLERNN outperforms
162 all other methods, exhibiting significantly reduced absolute error under both the demographic
163 model and under equilibrium assumptions ($T \leq -31$; $P < 10^{-16}$; *post hoc* Welch’s two sample *t*-tests
164 for all comparisons; **Figure 3**). Importantly, ReLERNN is also more accurate than all methods we
165 compared for each of the tested samples sizes, although all methods generally performed well with
166 larger sample sizes.

167 We also sought to assess the robustness of ReLERNN to demographic model misspecification,
168 whereby different generative models are used for simulating the training and test sets—e.g. training
169 on assumptions of demographic equilibrium when the test data was generated by a population
170 bottleneck. Methods robust to this type of misspecification are crucial, as the true demographic
171 history of a sample is often unknown and methods used to infer population size histories can
172 disagree or be unreliable (see **Figure 5–Figure Supplement 1**). Moreover, population size changes
173 alter the landscape of LD across the genome (e.g. **Slatkin, 1994; Rogers, 2014**), and thus have the
174 potential to reduce accuracy or produce biased recombination rate estimates.

175 To this end, we trained ReLERNN on examples generated under equilibrium and made pre-
176 dictions on 5000 chromosomes generated by the human demographic model specified above
177 (and also carried out the reciprocal experiment). We compared ReLERNN to both the CNN and
178 LDhat, whereby all methods were similarly misspecified (see Materials and Methods). We found
179 that ReLERNN outperforms both the CNN and LDhat, exhibiting significantly lower absolute er-
180 ror under both directions of demographic model misspecification ($T \leq -26$; $P_{WTT} < 10^{-16}$ for all
181 comparisons; **Figure 4**). Interestingly, we show that the error attributed to model misspecification
182 (termed marginal error; see Materials and Methods) is significantly greater when ReLERNN was
183 trained on equilibrium simulations and tested on demographic simulations than under the recip-
184 rocal misspecification ($T = 26.3$; $P_{WTT} < 10^{-16}$; **Figure 4–Figure Supplement 1**). While this is true,
185 it is important to note that marginal error is quite modest in both directions of misspecification
186 ($< 1.30 \times 10^{-9}$; **Figure 4–Figure Supplement 1**), suggesting that the additional information gleaned
187 from an informative demographic model is limited.

188 Differences in the ratio of homologous gene conversion events to crossovers can also bias
189 the inference of recombination rates, as conversion tracts break down LD within the prediction
190 window (**Gay et al., 2007; Przeworski and Wall, 2001**). We treated the effect of gene conversion as

191 another form of model misspecification by training on examples that lacked gene conversion and
192 testing on examples that included gene conversion. As ReLERNN uses msprime for all training
193 simulations, and msprime cannot currently simulate gene conversion, we generated all test set
194 simulations with ms (Hudson, 2002). We found that including gene conversion in our simulations
195 biased our predictions, resulting in an overestimate of the true recombination rate (Figure 4–Figure
196 Supplement 2). Moreover, the magnitude of this bias increased with the ratio of gene conversion
197 events to crossovers. As expected, we also observed a similar pattern of bias for LDhat, although the
198 magnitude of bias for LDhat was somewhat less than that exhibited by ReLERNN (Figure 4–Figure
199 Supplement 2).

200 **Recombination landscapes are largely concordant among populations of African *D.*** 201 ***melanogaster***

202 Using our method, we characterized the genome-wide recombination landscapes of three popula-
203 tions of African *D. melanogaster* (sampled from Cameroon, Rwanda, and Zambia). Each population
204 was derived from the sequencing of 10 haploid embryos (detailed in Lack et al., 2015; Pool et al.,
205 2012), hence these data represent an excellent opportunity to exploit ReLERNN's high accuracy
206 on small sample sizes. We first sought to model the demographic history of each population, as
207 ReLERNN can simulate training data under demographic models inferred by three published soft-
208 ware methods—stairwayplot (Liu and Fu, 2015), SMC++ (Terhorst et al., 2016), and MSMC (Schiffels
209 and Durbin, 2014). Using all three methods, we show that inferred historical population sizes are
210 unreliable for these populations—no two methods recapitulate the same history, and the histories
211 generated by MSMC vary dramatically depending on the number of samples used (Figure 5–Figure
212 Supplement 1, Figure 5–Figure Supplement 2). For these reasons, and because results from our
213 simulations suggest that marginal error due to demographic misspecification is quite low for our
214 method (above; Figure 4–Figure Supplement 1), we decided to simulate our training data under the
215 assumptions of demographic equilibrium.

216 Using ReLERNN, we discovered that the fine-scale recombination landscapes are highly corre-
217 lated among all three populations of *D. melanogaster* (genome-wide mean pairwise Spearman's
218 $\rho = 0.76$; $P < 10^{-16}$; 100 Kb windows; Figure 5). The genome-wide mean pairwise coefficient of
219 determination between populations was somewhat lower, $R^2 = 0.63$ ($P < 10^{-16}$; 100 Kb windows),
220 suggesting there may be important population-specific differences in the fine-scale drivers of
221 allelic association. These differences may also contribute to within-chromosome differences in
222 recombination rate between populations. Indeed, we estimate that mean recombination rates are
223 significantly different among populations for all chromosomes with the exception of chromosome
224 3L ($P \leq 3.78 \times 10^{-4}$; one-way analysis of variance). Post-hoc pairwise comparisons suggest that
225 this difference is largely driven by an elevated rate of recombination in Zambia, identified on all
226 chromosomes ($P \leq 8.21 \times 10^{-4}$; Tukey's HSD tests) except for 3L ($P_{HSD} \geq 0.15$). ReLERNN predicts
227 the recombination rate in simulated test sets to a high degree of accuracy for all three populations
228 ($R^2 \geq 0.93$; $P < 10^{-16}$; Figure 5–Figure Supplement 3), suggesting that we have sufficient power to
229 discern fine-scale differences in per-base recombination rates across the genome.

230 When comparing our recombination rate estimates to those derived from experimental crosses
231 of North American *D. melanogaster* (reported in Comeron et al., 2012), we find that the coefficients
232 of determination averaged over all three populations were $R^2 = 0.46, 0.70, 0.47, 0.08, 0.73$ for chro-
233 somes 2L, 2R, 3L, 3R, and X, respectively (Figure 5–Figure Supplement 4; 1 Mb windows). These
234 results differ from those observed by Chan et al. (2012), who compared 22 *D. melanogaster* sampled
235 from the same Rwandan population to the FlyBase map and found $R^2 = 0.55, 0.63, 0.45, 0.42, 0.41$ for
236 the same chromosomes. The minor differences we observed between methods for chromosomes
237 2L, 2R, and 3L can likely be attributed to the fact that we are comparing estimates from two different
238 methods, using different African flies, to a different experimentally derived map. However, the
239 larger differences found between methods for chromosomes 3R and the X seem less likely at-
240 tributable to methodological differences. Importantly, African *D. melanogaster* are known to harbor

241 large polymorphic inversions (*Corbett-Detig and Hartl, 2012; Lack et al., 2015*), often at appreciable
242 frequencies. For example, the inversion *In(3R)K* segregates in our Cameroon population at $p = 0.9$.
243 It is potentially these differences in inversion frequencies that contribute to the exceptionally weak
244 correlation observed using our method for chromosome 3R and the larger differences between
245 methods for chromosome X.

246 An important cause of population-specific differences in recombination landscapes might be
247 population-specific differences in the frequencies of chromosomal inversions, as recombination is
248 expected to be strongly suppressed between standard and inversion arrangements. Segregating
249 inversions in *D. melanogaster* have been shown to affect broad patterns of chromosomal varia-
250 tion, and are thought to have quite recent origins when taken together (*Corbett-Detig and Hartl,*
251 *2012*). To test for an effect of inversion frequency on our measurement of recombination rates, we
252 resampled haploid genomes from Zambia to create sampled populations with the cosmopolitan
253 inversion *In(2L)t* segregating at varying frequencies, $p \in \{0.0, 0.2, 0.6, 1.0\}$. In Zambia, *In(2L)t* segre-
254 gates at $p = 0.22$ (*Lack et al., 2015*), suggesting that recombination within the inversion breakpoints
255 may be strongly suppressed in individuals with the inverted arrangement relative to those with
256 the standard arrangement. Moreover, *In(2L)t* arose recently, likely within the past 100,000 years
257 (*Corbett-Detig and Hartl, 2012*). For these reasons, we predict that the inferred recombination rate
258 should decrease as the low-frequency inverted arrangement is increasingly overrepresented in the
259 set of sampled chromosomes (i.e. as more of the samples contain the high-LD inverted arrange-
260 ments). As predicted, we found a strong effect of the sample frequency of *In(2L)t* on estimated rates
261 of recombination for chromosome 2L in Zambia (*Figure 6*). Recombination rates are negatively
262 correlated with inversion frequency in our sample, not only within the inversion, but also in regions
263 3 Mb outside the inversion (flanking regions) ($\rho_{Spearman's} = -1$; $P = 0.04$ for both comparisons). We
264 also see a similar negative correlation outside the flanking regions, although this association is
265 weakened relative to that within or flanking the inversion (*Figure 6*). Importantly, varying the size of
266 the flanking regions (from 1-5 Mb) produces patterns that are qualitatively identical, suggesting that
267 the effect of inversions on recombination suppression extends far beyond the inversion breakpoints
268 themselves (*Figure 6–Figure Supplement 1*).

269 While the effect of inversion frequency on recombination rates may extend beyond the inver-
270 sion breakpoints, we expect that rates of recombination should be correlated with distance to the
271 inversion breakpoint on smaller spatial scales. To test this we looked at the recombination rates in
272 our African *D. melanogaster* populations, binned by distance to the nearest inversion breakpoints
273 segregating in these populations. Importantly, we curated the samples for our population com-
274 parisons by seeking to match the frequency of each inversion segregating in our samples with
275 its true population frequency, as measured in the whole of the DGN database (see Materials and
276 Methods). We show that recombination rates in the flanking regions are positively correlated with
277 distance to inversion breakpoints in both Rwanda and Zambia ($\rho_{Spearman's} = 1$; $P = 0.04$ for both
278 comparisons) but not in Cameroon ($\rho_{Spearman's} = 0.8$; $P = 0.17$; *Figure 7*). Likewise, recombination
279 rates in the inversion interior (> 2 Mb from the breakpoints) are expected to be higher than in
280 those regions immediately surrounding the breakpoints. However, with the exception of Cameroon
281 (Inversion interior compared to < 250 Kb from breakpoint; $P_{WTT} = 0.035$), we did not observe this
282 pattern ($P_{WTT} \geq 0.057$; *Figure 7*).

283 To further explore population-specific differences in recombination landscapes we took a statis-
284 tical outlier approach, whereby we define two types of recombination rate outliers—global outliers
285 and population-specific outliers (see Materials and Methods). Global outliers are characterized by
286 windows with exceptionally high variance in rates of recombination between all three populations
287 (*Figure 5*; red triangles) while population-specific outliers are those windows where the rate of re-
288 combination in one population is strongly differentiated from the rates in the other two populations
289 (*Figure 5*; population-colored triangles). We find that population-specific outliers, but not global
290 outliers, are significantly enriched within inversions ($P = 0.005$; randomization test; *Figure 5*; grey
291 boxes). Moreover, this enrichment remains significant when extending the inversion boundaries

292 by up to 250 Kb ($P_{rand} \leq 0.004$). However, extending the inversion boundaries beyond 250 Kb, or
293 restricting the overlap to windows surrounding only the breakpoints (250 Kb, 500Kb, 1 Mb, 2 Mb),
294 erodes this pattern ($P_{rand} \geq 0.055$ for all comparisons), suggesting that the role for inversions in
295 generating population-specific differences in recombination rates is complex, at least for these
296 populations.

297 Selection is another important factor that may confound the inference of recombination rates.
298 For instance selective sweeps generate localized patterns of high LD on either side of the sweep site
299 (*Kim and Nielsen, 2004; Schrider et al., 2015*), thus regions flanking selective sweeps may mimic
300 regions of reduced recombination. Inasmuch population-specific selective sweeps are expected to
301 contribute to population-specific differences in recombination rate estimates. We used diploS/HIC
302 (*Kern and Schrider, 2018*) to identify hard and soft selective sweeps in our African *D. melanogaster*
303 populations, and we tested for an excess of recombination rate outliers overlapping with windows
304 classified as sweeps. In total, diploS/HIC classified 27.4%, 28.1%, and 26.8%, of all genomic windows
305 as selective sweeps (either "hard" or "soft") for Cameroon, Rwanda, and Zambia, respectively, when
306 looking at 5kb, non-overlapping windows. The associated False Discovery Rates (FDR) for calling
307 sweeps in these populations were appreciable: 33.9%, 33.1% and 34.7%, respectively (**Figure 5–**
308 **Figure Supplement 5**). As expected, windows classified as sweeps had significantly lower rates of
309 recombination relative to neutral windows in all three populations ($P_{WTT} \leq 10^{-16}$ for all comparisons;
310 **Figure 7**). However, we found that neither global nor population-specific outliers were enriched
311 for selective sweeps ($P_{rand} \geq 0.246$ for both comparisons), suggesting that, when treated as a class,
312 recombination rate outliers are not likely driven by sweeps in these populations. When treated
313 separately (i.e. independent permutation tests for each recombination rate outlier window), we
314 identified 7 outliers enriched for sweeps at the $P \leq 0.05$ threshold, corresponding to an expected
315 FDR of 77%. However, given our FDR for calling sweeps in these populations, our measure of
316 the enrichment in overlap with recombination rate outliers is likely to be conservative. Two of
317 these outlier windows may represent potential true positives; an outlier in Cameroon contains 5
318 out of 6 non-overlapping 5 kb windows classified as "hard" sweeps, the second from Rwanda has
319 10 out of 12 windows classified as "hard" sweeps ($P_{rand} = 0.0$ for both comparisons). These two
320 recombination rate outlier windows are potentially ripe for future studies on selective sweeps in
321 these populations, and suggest that in at least some instances, selection contributes to observed
322 differences in estimates of recombination rates between *Drosophila* populations.

323 Discussion

324 We introduced a new method, ReLERNN, for predicting the genome-wide landscape of per-base
325 recombination rates from phased haplotypes from as few as four samples through the use of
326 deep neural networks. Population genomics, as a field, relies on estimates of recombination rates
327 to understand the effects of diverse phenomena ranging from the impacts of natural selection
328 (*Elyashiv et al., 2016*), to patterns of admixture and introgression (*Price et al., 2009; Brandvain*
329 *et al., 2014; Schumer et al., 2018*), to polygenic associations in genome-wide association studies
330 (*Bulik-Sullivan et al., 2015*). As befits this need, there has been a long tradition of development of
331 statistical methods for estimating the population recombination parameter, $\rho = 4Nr$ (*Chan et al.,*
332 *2012; Gao et al., 2016; Hudson and Kaplan, 1985; Hudson, 1987, 2002; Li and Stephens, 2003; Lin*
333 *et al., 2013; McVean et al., 2002; Myers and Griffiths, 2003; Wakeley, 1997; Wall, 2000; Wiuf, 2002*).

334 We sought to harness the power of deep learning, specifically deep recurrent neural networks, to
335 address the problem of estimating recombination rates, and in so doing, we developed a workflow
336 that reconstructs the genome-wide recombination landscape to a high degree of accuracy from
337 very small sample sizes—e.g. four phased haploid chromosomes. The use of deep learning has
338 recently revolutionized the fields of computer vision (*Krizhevsky et al., 2012; Szegedy et al., 2015*),
339 speech recognition (*Hinton et al., 2012*), and natural language processing (*Sutskever et al., 2014*),
340 and while its use in population genomics has only recently begun, it is anticipated to be similarly
341 fruitful (*Schrider and Kern, 2018*). The natural extension of deep learning to population genomic

342 analyses comes as a result of the ways in which ANNs learn abstract representations of their
343 inputs. In the case of population genomic analyses, the inputs can be naturally represented as
344 DNA sequence alignments, eliminating the need for human oversight (and potentially constraint)
345 in the form of statistical summaries (i.e. compression) of the raw data. ANNs can then learn
346 high-dimensional statistical associations directly from the sequence alignments, and use these to
347 return highly accurate predictions.

348 ReLERNN utilizes a variant of an ANN, known as a Gated Recurrent Unit (GRU), as its primary
349 technology. GRU networks excel at identifying temporal associations (*Jozefowicz et al., 2015*), and
350 therefore we model our sequence alignment as a bidirectional time series, whereby each ordered
351 SNP represents a new time step along the chromosome. We also model the distance between
352 SNPs using a separate input tensor, and these two inputs are concatenated after passing through
353 the initial layers of the network (see *Figure 1* inlay). We demonstrated that ReLERNN can predict
354 a simulated recombination landscape with a high degree of accuracy ($R^2 = 0.93$; *Figure 2*), and
355 that these predictions remain high, even when using small sample sizes ($R^2 = 0.82$; *Figure 2-Figure*
356 *Supplement 1*). Importantly, these predictions compared favorably to those made by a leading
357 composite likelihood method (LDhat; *McVean et al., 2002*), as well as other machine learning
358 methods (the CNN and FastEPRR; *Figure 3*). While the abstract nature of the data represented in
359 its internal layers constrains our ability to interpret the exact information ReLERNN relies on to
360 inform its predictions, our experiments using incorrect assumed mutation rates (*Figure 2-Figure*
361 *Supplement 3, Figure 2-Figure Supplement 2*) suggests that ReLERNN is potentially learning the
362 relative ratio of recombination rates to mutation rates. For these reasons, an extra caveat is
363 warranted—use caution when interpreting the results from ReLERNN as precise measures of the
364 per-base recombination rate unless precise mutation rate estimates are also known.

365 Demographic model misspecification is another potential source of error that should affect not
366 only deep learning methods targeted at estimating ρ , but also likelihood-based methods. Historical
367 demographic events (e.g. population bottlenecks, rapid expansions, etc.), because they may alter
368 the structure of LD genome-wide, can bias inference of recombination based on genetic variation
369 data. Our simulations demonstrated that while all the methods we tested had elevated error in
370 the context of demographic model misspecification, ReLERNN remained the most accurate across
371 all misspecification scenarios (*Figure 4*). While we caution against generalizing too much from
372 this experiment, the model misspecification tested here was extreme: we are replacing a human-
373 like demography of a bottleneck followed by exponential growth with a model of demographic
374 equilibrium. We suspect that ReLERNN, by using an RNN, is able to encode higher-order allelic
375 associations across the genome, for instance three-locus or four-locus linkage disequilibrium, and
376 in so doing capture more of the information available than traditional methods that use composite
377 likelihoods of two-locus LD summaries. Additionally, there are clear opportunities for future
378 improvements to ReLERNN. For instance, our simulation studies demonstrated that the RNN used
379 by ReLERNN is also sensitive to gene conversion events (*Figure 4-Figure Supplement 2*), thus the
380 joint estimation of rates of recombination and gene conversion may be quite feasible. Ultimately,
381 it remains far from clear what network architectures will be best suited for population genetic
382 inference, though we remain optimistic that ANNs will prove useful for a variety of applications in
383 the field.

384 A natural application of ReLERNN, due in part to its high accuracy with small sample sizes, was
385 to characterize and compare the recombination landscapes for multiple populations of African *D.*
386 *melanogaster*, for which few populations with large sample sizes are currently available. Previous
387 estimates of genome-wide fine-scale recombination maps in flies have focused on characterizing
388 recombination in experimental crosses (*Comeron et al., 2012*), or by running LDhat (or the related
389 LDhelmet) on populations with relatively moderate sample sizes (i.e. ≥ 22 samples) (*Chan et al.,*
390 *2012; Langley et al., 2012*). Here, we applied ReLERNN to three populations for which at least ten
391 haploid embryos were sequenced: Cameroon, Rwanda, and Zambia (*Lack et al., 2015; Pool et al.,*
392 *2012*). Generally, recombination landscapes were well correlated among populations. Mean pair-

393 wise coefficients of determination among all three populations were $R^2 = 0.69, 0.61, 0.77, 0.43, 0.66$
394 for chromosomes 2L, 2R, 3L, 3R, and X, respectively. These correlations are notably lower than
395 those observed in humans (*Myers et al., 2005*) and mice (*Wang et al., 2017*), and one potential
396 biological cause for this large difference could be the cosmopolitan chromosomal inversions that
397 segregate in African *D. melanogaster* (*Corbett-Detig and Hartl, 2012; Lack et al., 2015*).

398 We demonstrated a significant negative association between inversion sample frequency and
399 recombination rate as inferred by ReLERNN through experimentally manipulating the frequency
400 of the inversion karyotype in our sample (*Figure 6*). Our results suggest that recombination sup-
401 pression extends well beyond the predicted breakpoints of the inversion (at least 5 Mb beyond in
402 the case of *In(2L)t*; *Figure 6–Figure Supplement 1*). This large-scale suppression of recombination
403 due to inversions in *Drosophila* has been observed both directly in experimental crosses (*Dobzhansky*
404 *and Epling, 1948; Novitski and Braver, 1954; Kulathinal et al., 2009; Miller et al., 2016; Fuller*
405 *et al., 2018*), and indirectly from patterns of variation surrounding known inversion breakpoints
406 (*Corbett-Detig and Hartl, 2012; Langley et al., 2012*). While it is true that the negative relationship
407 between inversion frequency and recombination should only exist for inversions segregating at low
408 frequencies (e.g. crossover suppression is not expected in inversion homozygotes), we predict a
409 negative relationship to dominate in these populations, as the majority of polymorphic inversions
410 are young, segregate at low frequencies, and show elevated LD along their lengths perhaps due to
411 the actions of natural selection (*Corbett-Detig and Hartl, 2012; Lack et al., 2015*).

412 While polymorphic inversions exert strong effects on recombination landscapes, support for
413 their role in explaining the most diverged regions among populations was mixed—we found that
414 population-specific recombination rate outliers, but not global outliers, were significantly enriched
415 within the inversions known to segregate in these populations (*Figure 5*). Moreover, our predictions
416 for the relative rates of recombination among populations, based on inversion frequencies per
417 chromosome, were largely not met—the inversions *In(2L)t*, *In(2R)NS*, and *In(3L)Ok* segregate at the
418 highest frequencies in Zambia, yet this population also has the highest average recombination
419 rate for these three chromosomes. Chromosome 3R, however, did match these predictions,
420 having inversions segregating at the highest frequencies of any chromosome (e.g. $p_{In(3R)K} = 0.9$ in
421 Cameroon) and also both the lowest coefficient of determination ($R^2 = 0.43$) and population-specific
422 recombination rates ranked in accordance with inversion frequencies (*Figure 5*).

423 Interestingly, while we identified two individual outlier regions characterized by numerous
424 selective sweeps, we did not observe a significant enrichment of sweeps overlapping either global
425 or population-specific outliers when these outliers were treated as a class of genomic elements.
426 This is perhaps surprising, given that selective sweeps are known to create characteristic elevations
427 of LD (*Kim and Nielsen, 2004*), and perhaps could mimic regions with very divergent levels of
428 recombination in a population-specific way. A number of other evolutionary forces might explain
429 the existence of our outlier regions as well. For example, mutation rate heterogeneity along
430 the chromosomes could, in principle, generate spurious peaks or troughs in our estimates of
431 recombination rate, as ReLERNN in effect scales its per-base recombination rate estimates by a
432 mutation rate that is assumed to be constant along the chromosome (*Figure 2–Figure Supplement 3*,
433 *Figure 2–Figure Supplement 2*). Moreover, introgression from diverged populations might affect
434 patterns of allelic association in a local way along the genome (*Schrider et al., 2018; Schumer*
435 *et al., 2018*). Taken together, our results suggest that while both inversions and selection can
436 influence population-specific differences in the landscape of recombination, the preponderance of
437 these differences likely have complex causes.

438 In this report we described ReLERNN, a novel deep learning method for inferring fine-scale rates
439 of recombination across the genome. While ReLERNN currently stands as a functional end-to-end
440 pipeline for measuring recombination rates, the modular design herein presents a number of
441 important opportunities for extension, with the potential to address myriad questions in population
442 genomics. For example, while ReLERNN is currently designed to use phased haplotype data as
443 its input, we see no reason why unphased, diploid genotypes couldn't be substituted (e.g. *Flagel*

444 *et al., 2018*). Moreover, the RNN structure we exploit here could be used for inference of the
445 distribution of selection coefficients and/or migration rates from natural populations. In addition,
446 ReLERNN presents an excellent opportunity for the implementation of transfer learning, whereby
447 ReLERNN could be trained in-house on an otherwise prohibitively extensive parameter space,
448 allowing end-users to make accurate predictions by generating only a small fraction of the current
449 number of simulations and training epochs presently required. The application of machine learning,
450 and deep learning in particular, to questions in population genomics is ripe with opportunity.
451 ReLERNN provides a platform for jumping off, that we hope to see advance our understanding of
452 both population genetics and adaptation itself.

453 **Materials and Methods**

454 **The ReLERNN workflow**

455 Here we briefly describe ReLERNN, a software package for accurately estimating a genome-wide re-
456 combination landscape from as few as four phased chromosomes. The ReLERNN workflow proceeds
457 by the use of four python modules—ReLERNN_SIMULATE, ReLERNN_TRAIN, ReLERNN_PREDICT,
458 and ReLERNN_BSCORRECT (*Figure 1*). The first three modules are mandatory, and include functions
459 to calculate Watterson's estimator and historical population sizes, functions for simulating the
460 training set, functions for training the neural network, and functions for reporting rates of recom-
461 bination along the chromosomes. The fourth module, ReLERNN_BSCORRECT, is optional (though
462 recommended) and includes functions for estimating 95% confidence intervals and implements a
463 correction function to reduce biases that may arise during training. The output from ReLERNN is a
464 list of genomic windows and their corresponding recombination rate prediction (reported as per-
465 base crossover events), along with 95% confidence intervals if the optional ReLERNN_BSCORRECT
466 module was used.

467 **Parameter estimation and coalescent simulation**

468 ReLERNN takes as input a VCF file of phased biallelic variants, which can either be coded as
469 nucleotides or ancestral/derived states (i.e. 0/1). A minimum of four sample chromosomes must
470 be included, and users should ensure proper filtering of the input file beforehand—e.g. excluding
471 low-coverage or low-quality sites, non-biallelic sites, and missing data. ReLERNN also requires the
472 user to provide an assumed per-base mutation rate and an assumed maximum value for the ratio
473 ρ/θ . These parameters are used to set an acceptable window size for prediction, by restricting the
474 total number of segregating sites in each window to remain below a critical threshold. ReLERNN
475 therefore uses a dynamic window size to reduce the probability of training failure due to having
476 too many, or too few, segregating sites present in a window (e.g. experimental trials showed that
477 the training loss function eventually returns NaNs when training on windows containing multiple
478 thousands of segregating sites). As a result, the output predictions file may return different window
479 sizes for different chromosomes, even within the same genome. For comparing rates between
480 populations, an optional script ("force_window_size_predictions.py") is provided to force rates to
481 conform to a given window size. This is accomplished by taking a weighted average of recombination
482 rates, whereby rates are weighted by the fraction of overlap between their original window positions
483 and the new forced window positions.

484 Once the appropriate window sizes have been estimated, ReLERNN_SIMULATE uses the coales-
485 cent simulation software, msprime (*Kelleher et al., 2016*), to independently generate 10^5 training
486 examples and 10^3 validation and test examples. By default, these simulations are generated under
487 assumptions of demographic equilibrium, using a range of per-base mutation and recombination
488 rates. However, ReLERNN can optionally simulate under a demographic history inferred by one of
489 three programs: stairwayplot (*Liu and Fu, 2015*), SMC++ (*Terhorst et al., 2016*), or MSMC (*Schiffels
490 and Durbin, 2014*), and the handling of output from these programs is fully integrated into ReL-
491 ERNN_SIMULATE. This provides users the ability to model a demographic history and to estimate
492 rates of recombination from different files (e.g. one that includes only intergenic sites). When each

493 simulation is completed, ReLERNN dumps both the genotype matrix and a vector of the positions
494 for every SNP into a temporary .npy file.

495 Sequence batch generation and network architecture

496 To reduce the large memory utilization common to the analysis of genomic sequence data, we took
497 a batch generation approach, whereby only small batches of simulations are called into memory
498 at any one time. Data normalization and padding occurs when a training batch is called, by which
499 the genotype and position arrays are read into memory. In ReLERNN, ancestral states are coded
500 as -1 , derived states are coded as 1 , and both genotype and positions arrays are padded with
501 0 s to the maximum number of segregating sites generated across all examples. In addition, a
502 framing pad of five 0 s is applied to both arrays, and the order of samples in each batch is randomly
503 shuffled. The targets for each training batch are the per-base recombination rates used by msprime
504 when simulating each example. These targets are z-score normalized across all training examples.
505 The normalized and padded genotype and position arrays form the input tensors for our neural
506 network.

507 ReLERNN trains a recurrent neural network with Keras (*Chollet et al., 2015*) using a Tensorflow
508 backend (*Abadi et al., 2015*). The complete details of our neural architecture can be found in the
509 python module "ReLERNN_networks.py", and a simplified flow diagram showing the connectivity
510 between layers can be found in *Figure 1*. Briefly, the ReLERNN neural network utilizes distinct input
511 layers for the genotype and position tensors, which are later merged using a concatenation layer
512 in Keras. The genotype tensor is first fed to a GRU layer, as implemented with the bidirectional
513 wrapper in Keras, and the output of this layer is passed to a dense layer followed by a dropout
514 layer. On the positions side of the network, the input positions tensor is fed directly to a dense
515 layer and then to a dropout layer. Dropout was used extensively in our network, as hypertuning
516 trials (below) demonstrated significantly improved accuracy when employing dropout relative to
517 networks without dropout. Once concatenated, output from the dropout layer is passed to a final
518 round of dense and dropout layers, and the final dense layer returns a single z-score normalized
519 prediction for each example, which is unnormalized back to units of crossovers per-base. ReLERNN
520 completes 250 training epochs and implements this training using the "Adam" optimizer and a
521 Mean Squared Error (*MSE*) loss function. Though the number of epochs is user-selectable, the
522 vast majority of networks are sufficiently trained within 250 epochs, largely due to how ReLERNN
523 handles the input tensor size and simulation parameters. Our hyper-tuning trials were completed
524 via a grid search over the set of parameters: Recurrent layer output dimensions (64, 82, 128), Loss
525 function (*MSE*, *MAE*), Input merge strategy (concatenate, average), and dense layer dimensionality
526 (64, 128), optimizing for *MSE*.

527 Parametric bootstrap analysis and prediction corrections

528 ReLERNN includes the option to both generate confidence intervals around each predicted re-
529 combination rate and to correct for potential biases generated during training using a parametric
530 bootstrapping approach. After the network has been trained and predictions have been gener-
531 ated, users can run ReLERNN_BSCORRECT, which resimulates 10^3 test examples for each of 100
532 recombination rate bins drawn from the distribution of recombination rates used to simulate
533 the original training set. Predictions are then generated for these 10^5 simulated test examples
534 using the previously trained network, generating a distribution of predictions for each respective
535 recombination rate bin. 95% confidence intervals are calculated from by taking the upper and lower
536 2.5% rate predictions from this distributions.

537 The distribution of test predictions can be biased in systematic ways, such as predictably under-
538 estimating rates of recombination for those examples with the highest simulated crossover events
539 (*Figure 1–Figure Supplement 1*). These biases may potentially be caused an inability to resolve very
540 high recombination rates with a limited number of informative SNPs. ReLERNN_BSCORRECT, esti-
541 mates the magnitude of this bias through bootstrapping, and applies a bias correction function to

542 the empirical predictions. The bias correction function takes each empirical prediction and identifies
543 the nearest median value in the bootstrap distribution. The correction function then adds to this
544 prediction the difference between this median value and the true recombination rate used to simu-
545 late the distribution of test examples at that recombination rate bin. This correction method has the
546 effect of elevating the empirical prediction in regions of parameter space where we are reasonably
547 confident that we are underestimating recombination rates and lowering the prediction in areas
548 where we are likely to be overestimating recombination rates. ReLERNN_BSCORRECT is provided
549 as an optional module in ReLERNN, as the resimulation of 10^5 test examples is computationally
550 expensive and may not be warranted in all circumstances.

551 **Testing the accuracy of ReLERNN on simulated recombination landscapes**

552 To test the accuracy of ReLERNN at recapitulating a dynamic recombination landscape, we ran our
553 complete ReLERNN workflow on simulation data replicating chromosome 2L of *D. melanogaster*.
554 Using crossover rates estimated by [Comeron et al. \(2012\)](#), we simulated varying numbers of samples
555 of *D. melanogaster* chromosome 2L with msprime using the RecombinationMap class. Simulated
556 samples were exported to a VCF file using ploidy = 1, and all simulations were generated under
557 demographic equilibrium. We used these simulated VCF files as the input to our ReLERNN pipeline,
558 and ran all ReLERNN modules with default parameters, with the exception of varying the assumed
559 per-base mutation rate and the assumed maximum ratio of ρ to θ . Assumed mutation rates were
560 varied from 50% less than the rate used in simulations (true rate) to 50% greater than the true
561 rate. Likewise, the ratio of ρ to θ was either held constant, resulting in the training set containing
562 on average higher or lower per-base recombination rates than the true rate, or was adjusted to
563 correctly reflect the true maximum per-base recombination rate used—i.e. approximately 1.2×10^{-7}
564 crossovers per base.

565 **Comparative methods**

566 We chose to compare ReLERNN to three published methods for estimating recombination rates—
567 FastEPRR ([Gao et al., 2016](#)), a 1-dimensional CNN recently described in [Flagel et al. \(2018\)](#) and
568 LDhat ([McVean et al., 2002](#)). We generated a training set (used by ReLERNN and the CNN) with
569 10^5 examples and tested each method on an identical set of 5×10^3 simulation examples for
570 testing. We generated two classes of simulations, one simulated under demographic equilibrium
571 and one using a demographic history derived from European humans (CEU model; detailed in
572 "ReLERNN_demographic_models.py"; [Tennessen et al., 2012](#); [Gravel et al., 2011](#)). Both classes
573 of simulations were generated for $n \in \{4, 8, 16, 32, 64\}$, where n is the number of chromosomes
574 sampled from the population. All simulations were generated in msprime with the common
575 set of parameters: $priorLowsRho = 0.0$, $priorHighsRho = 5e^{-8} \times 1.25$, $priorLowsMu = 2.5e^{-8} \times 0.75$,
576 $priorHighsMu = 2.5e^{-8} \times 1.25$, $ChromosomeLength = 3e^5$, whereby values for both per-base mutation
577 and recombination rates were drawn from a uniform distribution between the low and high priors.

578 For both ReLERNN and the CNN, the same training set consisting of 10^5 examples was used
579 to train each neural network, and the same test examples were used to compare the predictions
580 produced by each method. Comparisons with LDhat were made using the above training examples
581 to parameterize the generation of independent coalescent likelihood lookup tables. For each set of
582 examples of sample size N , we calculated the maximum value of ρ from the training set and the
583 average per-base values for θ for the test examples, using Watterson's estimator. These parameter
584 values were given to LDhat's *complete* function for the lookup table generation, and the resulting
585 table was used to make predictions on our 5×10^3 test examples using the *pairwise* function.
586 Comparisons with FastEPRR were made by transforming the genotype matrices resulting from
587 our test simulations into fasta-formatted input files, and running the FastEPRR_ALN function (using
588 format = 1) in R. As both LDhat and FastEPRR predict ρ , the resulting predictions were transformed
589 to per-base recombination rates for comparison with ReLERNN using the function $r = \frac{\rho_{pred} \times \mu_{true}}{\theta_W}$,
590 whereby ρ_{pred} is the prediction output by each method, and θ_W and μ_{true} are Watterson's estimator

591 and the true per-base mutation rate used in the simulation example, respectively. To compare
592 accuracy among methods we directly compared the distribution of absolute errors ($|r_{predicted} - r_{true}|$)
593 for each method for each set of examples of sample size N .

594 To test the effects of model misspecification on predictions, we simply directed ReLERNN and
595 the CNN to use a training set generated under demographic equilibrium for making predictions
596 on a test set generated under the CEU model, and vice versa. To test for the effects of model
597 misspecification in LDhat, we generated a lookup table using parameter values estimated from
598 the misspecified training set (e.g. the lookup table used for predicting the CEU model test set was
599 generated by using parameter values directly inferred from training simulations under equilibrium.
600 We did not directly test the effect of model misspecification using FastEPRR, as this method takes
601 as input only a fasta sequence file, and therefore the internal training of the model was not able to
602 be separated from the input sequences. To address the effects of model misspecification, we also
603 directly compared the distribution of absolute errors ($|r_{predicted} - r_{true}|$). Additionally, we compared the
604 marginal error directly attributable to model misspecification among methods. We defined marginal
605 error as $\epsilon_m - \epsilon_c$, where ϵ_m and ϵ_c are equal to $|r_{predicted} - r_{true}|$ when the model is misspecified and
606 correctly specified, respectively. We simulated gene conversion test sets using ms (Hudson, 2002),
607 with a mean conversion tract length of 352 bp (corresponding to the mean empirically derived
608 tract length in *D. melanogaster* (Hilliker et al., 1994)) and simulated a ratio of conversion events to
609 crossover events of 0, 1, 2, 4, and 8.

610 **Recombination rate variation in *D. melanogaster***

611 We obtained *D. melanogaster* population sequence data from the *Drosophila* Genome Nexus (DGN;
612 <https://www.johnpool.net/genomes.html>; Lack et al., 2015; Pool et al., 2012). We converted DGN
613 "consensus sequence files" to VCF format using custom python scripts, excluding all non-biallelic
614 sites and sites containing missing data. We chose to analyze populations from Cameroon, Rwanda,
615 and Zambia, as these populations contained at least 10 haploid embryo sequences per population
616 and each population included multiple segregating chromosomal inversions (supplemental table
617 1). To ensure roughly equivalent power to compare rates among populations, we downsampled
618 both Rwanda and Zambia to 10 chromosomes. We selected individual haploid genomes for each
619 population by requiring that our sampled inversion frequencies for each of the six segregating
620 inversions—*In(1)Be*, *In(2L)t*, *In(2R)NS*, *In(3L)Ok*, *In(3R)K*, and *In(3R)P*—closely approximate their popu-
621 lation frequencies as measured in the complete set of haploid genomes for that population. All
622 sample accessions and their corresponding inversion frequencies are located in the supporting
623 materials.

624 Before running ReLERNN, we first set out to model the demographic history for each population
625 using each of three methods: stairwayplot (Liu and Fu, 2015), SMC++ (Terhorst et al., 2016), and
626 MSMC (Schiffels and Durbin, 2014). With the exception of MSMC, all methods were run using
627 default parameters. For MSMC, the use of default parameters generated predictions that were
628 unusable (Figure 5–Figure Supplement 2). For these reasons, and after direct communication with
629 MSMC's authors, we determined that running MSMC with a sample size of two chromosomes would
630 be the most appropriate. Ultimately we decided to run our ReLERNN pipeline with simulations
631 generated under demographic equilibrium [options: `-estimateDemography False -assumedMu`
632 `3.27e-9 -upperRhoThetaRatio 35`], as estimates of historical population size were unreliable for
633 these data—all three methods produced significantly different demographic histories (Figure 5–
634 Figure Supplement 1)—and tests on simulated data suggest little effect of demographic model
635 misspecification (Figure 4–Figure Supplement 1). All code required to run our ReLERNN analysis
636 is deposited on GitHub (<https://github.com/kern-lab/ReLERNN> and <https://github.com/kern-lab/ReLERNN-analysis>).

638 We measured the correlation in recombination rates between each African *D. melanogaster*
639 populations in 100 kb sliding windows, as ReLERNN will predict the rates of recombination in slightly
640 different window sizes, depending on θ for each chromosome. The recombination rate for each

641 sliding window was calculated by taking the average of all rate windows predicted by ReLERNN,
642 weighted by the fraction that each window overlapped the larger sliding window. Recombination
643 rate outliers were identified in two ways: as global outliers and population-specific outliers. Global
644 outliers were identified by first calculating the mean and standard deviation in recombination rates
645 for all three populations in each 100 kb sliding window. We then used the top 1% of outliers from
646 the distribution of residuals, after fitting a linear model to the standard deviation on the mean.
647 Population-specific outliers were identified by using a modification of the population branch statistic
648 (herein PBS*; *Yi et al., 2010*), whereby we replaced pairwise F_{ST} with the pairwise differences in
649 recombination rates. We then used the top 1% of all PBS* scores as our population-specific outliers,
650 with each outlier corresponding to a PBS* score for a single population.

651 To test the effect of inversion frequency on predicted recombination rates, we resampled
652 10 haploid chromosomes from the available set of haploid genomes from Zambia to generate
653 sampled populations containing *In(2L)t* at varying frequencies, $p \in \{0.0, 0.2, 0.6, 1.0\}$. We then ran
654 ReLERNN on chromosome 2L for each of these resampled Zambian populations. We classified
655 recombination windows by their overlap with the coordinates of *In(2L)t* (as defined in *Corbett-Detig
656 and Hartl, 2012*), defining windows within the breakpoints (inside), windows up to 3 Mb outside the
657 breakpoints (flanking), and windows > 3 Mb outside the breakpoints (outside).

658 To test the effect of genome-wide inversion breakpoints on differences in recombination land-
659 scapes between populations, we classified windows by their overlap with inversion interiors (> 2 Mb
660 inside the inversion breakpoints) and their overlap with windows within 200 Kb, 500 Kb, 1 Mb, and
661 2 Mb of inversion breakpoints. We tested for an enrichment of both global and population-specific
662 outliers within inversions by randomization tests, whereby we permuted the labels for outliers
663 10^4 times and counted the overlap with inversions for each permutation to calculate the empirical
664 p-values. We also tested for an effect of selection on recombination rates in these populations,
665 by running diploS/HIC (*Kern and Schrider, 2018*) to detect selective sweeps. We ran diploS/HIC
666 on each population, training on simulations generated under demographic equilibrium. For each
667 population we simulated 2000 training examples from each of the five classes of regions required
668 by diploS/HIC using the coalescent simulation software discoal (*Kern and Schrider, 2016*). For simu-
669 lations which included sweeps we drew the selection coefficient from a uniform distribution such
670 that $s \sim U(0.0001, 0.005)$, the time of completion of the sweep from $\tau \sim U(0, 0.05)$, and the frequency
671 at which a soft sweep first comes under selection as $f \sim U(0, 0.1)$. We drew θ from $U(65, 654)$ and
672 we drew ρ from an exponential distribution with mean 1799 and the upper bound truncated at triple
673 the mean. For the discoal simulations we simulated 605 kb of data with the goal of classification of
674 the central most 55 kb window. We looked at the overlap with "sweep" windows (those classified
675 as either "hard" or "soft") and those windows classified as "neutral" by diploS/HIC. Our complete
676 diploS/HIC pipeline for these samples is available in the supporting materials online. All statistical
677 tests were completed in R (*R Core Team, 2018*), with the exception of empirical randomization tests,
678 which were completed using Python.

679 **Data availability**

680 ReLERNN is currently available at <https://github.com/kern-lab/ReLERNN>. Supporting information,
681 tables, and figures will be deposited online at the publication journal.

682 **Acknowledgments**

683 The authors would like to gratefully acknowledge Matthew Hahn, Dan Schrider, and Peter Ralph
684 for their helpful comments and suggestions. This work benefited from access to the University of
685 Oregon high performance computer, Talapas. JRA, JGG, and ADK were supported by NIH award
686 R01GM117241 to ADK. We would also like to thank the Hearth for their fine coffee.

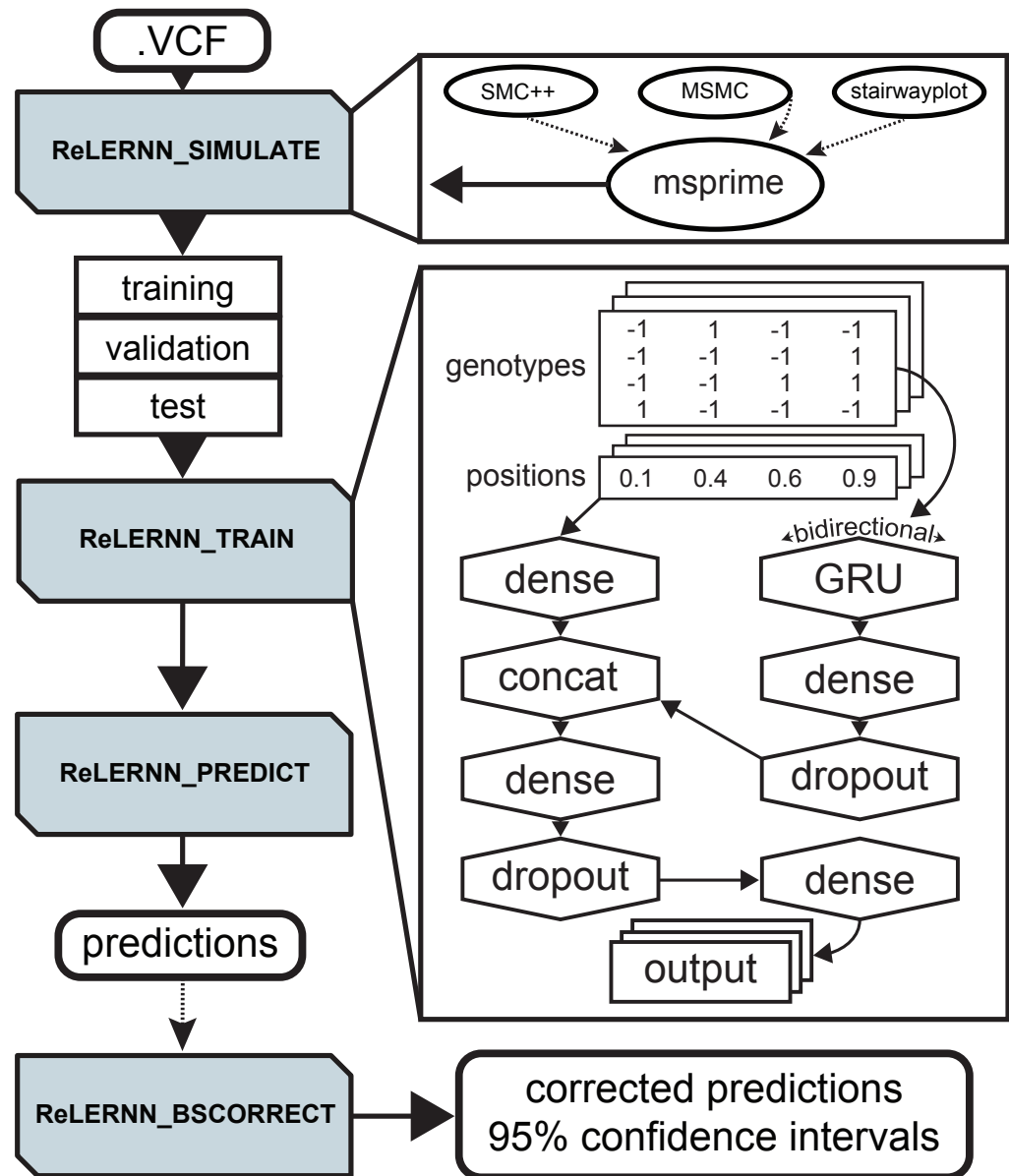


Figure 1. Diagram depicting a typical workflow using ReLERNN's four modules (shaded boxes). ReLERNN_SIMULATE can optionally (dotted lines) utilize output from stairwayplot, SMC++, MSMC to simulate under a demographic history in msprime. The breakout of ReLERNN_TRAIN depicts the GRU network architecture used for training. The input genotype matrix shows alleles encoded as ancestral (-1), derived (1), or padded (0; *not shown*), and the input position matrix shows variant position coded along the real number line (0-1).

Figure 1-Figure supplement 1. Parametric bootstrapping results as implemented by ReLERNN. Lines represent the minimum (blue), lower 5% (orange), lower 25% (green), median (red), upper 25% (purple), upper 95% (brown), and maximum (pink) bounds for each of 1000 replicate simulations and predictions (y-axis) across 100 recombination rate bins (x-axis)

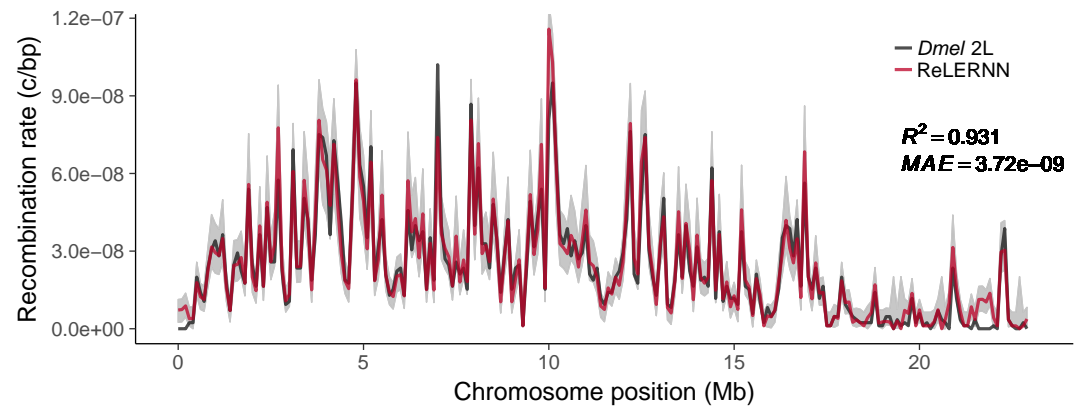


Figure 2. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

Figure 2-Figure supplement 1. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 4$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

Figure 2-Figure supplement 2. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Here the per-base mutation rate was assumed to be 50% less than the rate used for simulation. Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

Figure 2-Figure supplement 3. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Here the per-base mutation rate was assumed to be 50% greater than the rate used for simulation. Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

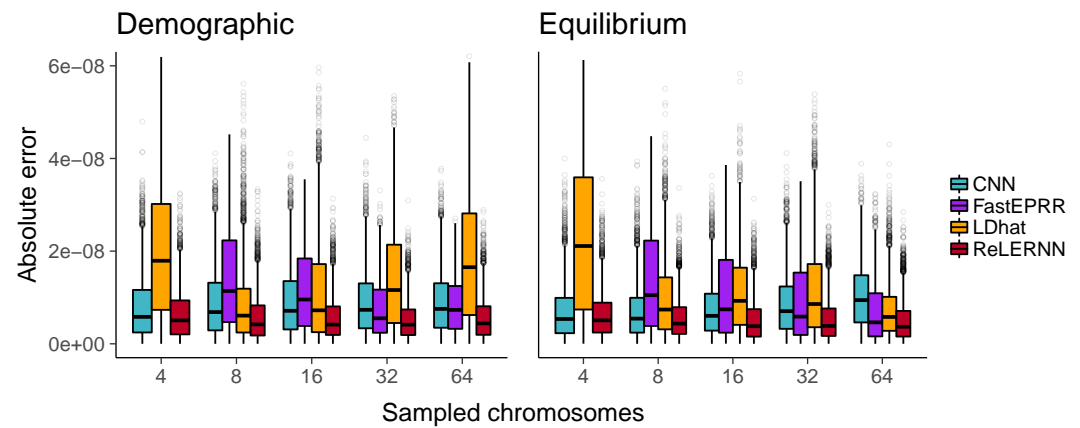


Figure 3. Distribution of absolute errors ($|r_{predicted} - r_{true}|$) for each method across 5000 simulated chromosomes (1000 for FastEPRR). Independent simulations were run under a known demographic history (left) or an assumption of demographic equilibrium (right). Sampled chromosomes indicate the number of independent sequences that were sampled from each msprime (Kelleher et al., 2016) coalescent simulation.

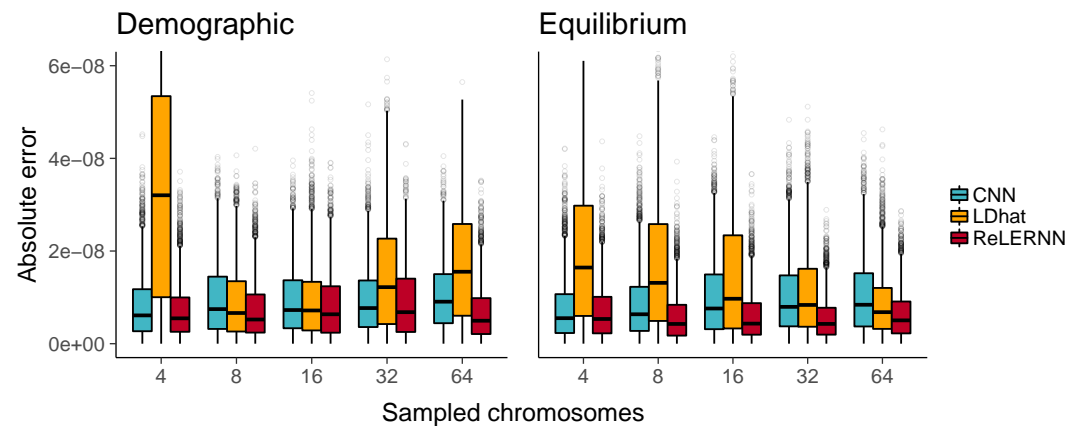


Figure 4. Distribution of absolute errors ($|r_{predicted} - r_{true}|$) for each method across 5000 simulated chromosomes after model misspecification. For the CNN and ReLERNN, predictions were made by training on equilibrium simulations and testing on sequences simulated under a demographic model (left) or training on demographic simulations and testing on sequences simulated under equilibrium (right). For LDhat, the lookup table was generated using parameters values that were estimated from simulations where the model was misspecified in the same way as described for the CNN and ReLERNN above. Sampled chromosomes indicate the number of independent sequences that were sampled from each msprime (Kelleher et al., 2016) coalescent simulation.

Figure 4–Figure supplement 1. Distribution of marginal errors attributed to model misspecification across 5000 simulated chromosomes. Predictions were made by training on equilibrium simulations and testing on sequences simulated under a demographic model (left) or training on demographic simulations and testing on sequences simulated under equilibrium (right). Here, marginal errors are represented as $\epsilon_m - \epsilon_c$, where ϵ_m and ϵ_c are equal to $|r_{predicted} - r_{true}|$ when the model is misspecified and correctly specified, respectively. Sampled chromosomes indicate the number of independent sequences that were sampled from each msprime (Kelleher et al., 2016) coalescent simulation.

Figure 4–Figure supplement 2. Distribution of predicted rates of recombination over true rates for 5000 examples simulated with gene conversion and $n = 8$. The ratio of gene conversion to crossovers was drawn from $U(0, c)$, with $c \in \{0, 1, 2, 4, 8\}$. Gene conversion tract lengths were fixed at 352 bp, and all simulations were completed in ms (Hudson, 2002).

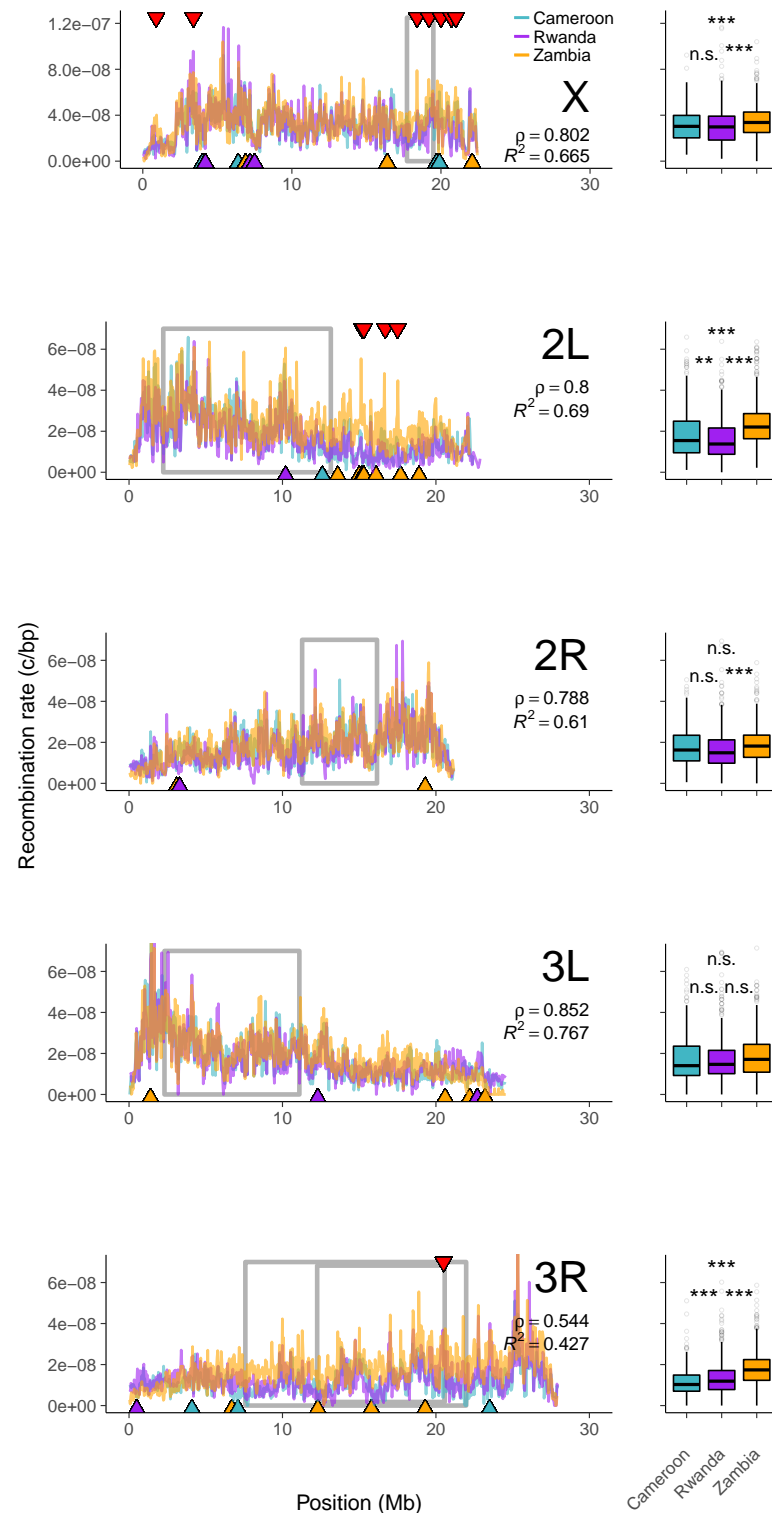


Figure 5. (Left) Genome-wide recombination landscapes for *D. melanogaster* populations from Cameroon (teal lines), Rwanda (purple lines), and Zambia (orange lines). Grey boxes denote the inversion boundaries predicted to be segregating in these samples (Pool et al., 2012; Corbett-Detig and Hartl, 2012). Red triangles mark the top 1% of global outlier windows for recombination rate. Blue, purple, and orange triangles mark the top 1% of population-specific outlier windows for recombination rate, with triangle color indicating the outlier population (see Materials and Methods). (Right) Per-chromosome recombination rates for each population. Spearman's ρ and R^2 are reported as the mean of pairwise estimates between populations for each chromosome. *** $P < 0.001$ and n.s. $P > 0.05$ are based on Tukey HSD tests for all pairwise comparisons.

Figure 5-Figure supplement 1. Historical population size estimates were inferred for Cameroon, Rwanda, and Zambia using three separate methods, all of which disagree with one another. Inferences are based on 10 samples for both stairwayplot (grey line) and SMC++ (orange line), and 2 samples for MSMC (purple line).

Figure 5-Figure supplement 2. Historical population size estimates were inferred for Cameroon, Rwanda, and

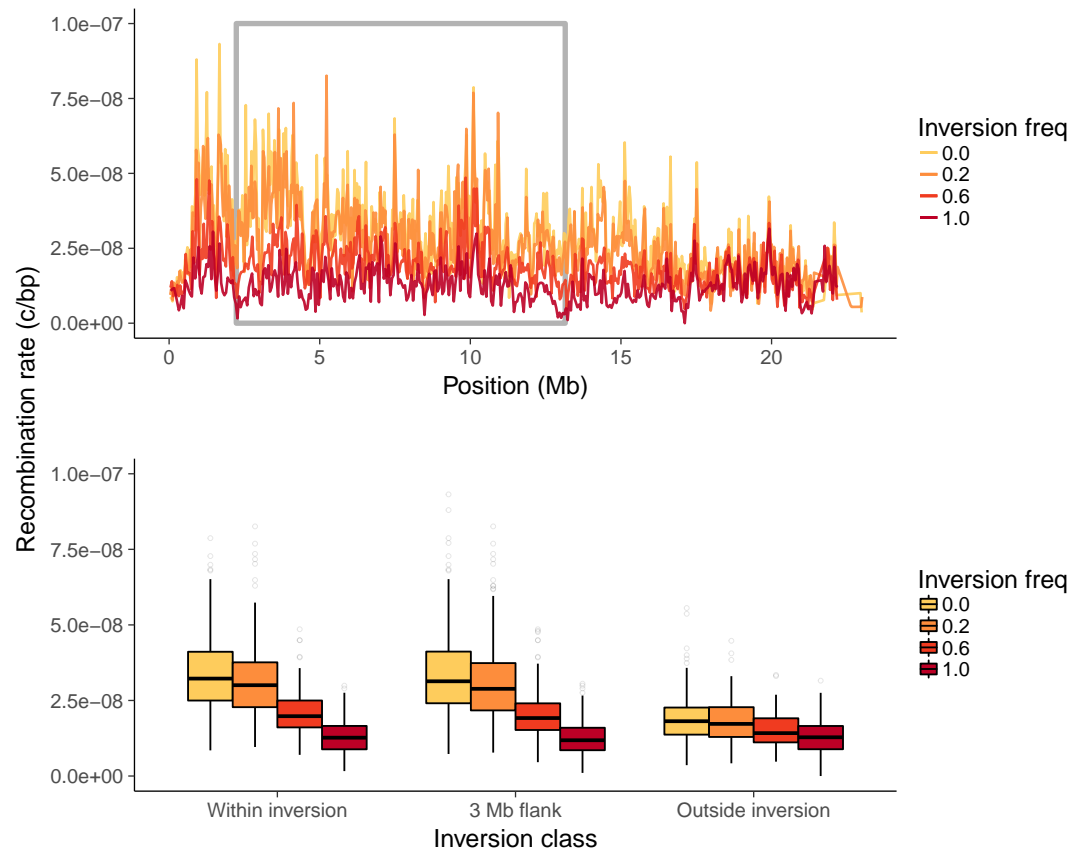


Figure 6. (Top) Recombination landscapes for Zambian *D. melanogaster* surrounding *In(2L)t*, sampled at different inversion frequencies. The grey box denotes the inversion boundaries of *In(2L)t* in *Drosophila* (Corbett-Detig and Hartl, 2012). (Bottom) Recombination rate estimates from genomic windows within the inversion, within a 3 Mb region flanking the inversion, and 3 Mb outside the inversion, sampled at different inversion frequencies. **Figure 6-Figure supplement 1.** Recombination rate estimates using flanking window sizes from 1-5 Mb. Rates are shown for genomic windows within the inversion, within regions flanking the inversion, and for regions outside both the inversion and flanking regions. All estimates are from chromosome 2L with *In(2L)t* sampled at different inversion frequencies

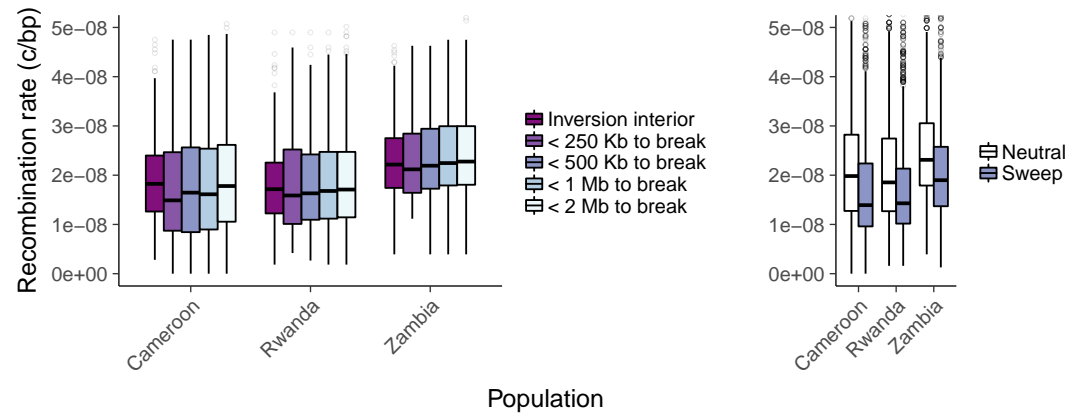


Figure 7. (Left) Recombination rate estimates for genomic windows > 2 Mb inside, < 250 kb surrounding, < 500 kb surrounding, < 1 Mb surrounding, and < 2 Mb surrounding all inversion breakpoints. (Right) Recombination rate estimates for all genomic windows overlapping windows predicted as either hard/soft sweeps (purple) or as neutral (white) by diploS/HIC (*Kern and Schrider, 2018*).

References

- 687
688 **Abadi M**, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin M, Ghemawat S,
689 Goodfellow I, Harp A, Irving G, Isard M, Jia Y, Jozefowicz R, Kaiser L, Kudlur M, Levenberg J, et al., TensorFlow:
690 Large-Scale Machine Learning on Heterogeneous Systems; 2015. <https://www.tensorflow.org/>, software
691 available from tensorflow.org.
- 692 **Ayala D**, Guerrero RF, Kirkpatrick M. Reproductive isolation and local adaptation quantified for a chromosome
693 inversion in a malaria mosquito. *Evolution: International Journal of Organic Evolution*. 2013; 67(4):946–958.
- 694 **Brandvain Y**, Kenney AM, Fligel L, Coop G, Sweigart AL. Speciation and introgression between *Mimulus nasutus*
695 and *Mimulus guttatus*. *PLoS genetics*. 2014; 10(6):e1004410.
- 696 **Bulik-Sullivan BK**, Loh PR, Finucane HK, Ripke S, Yang J, Patterson N, Daly MJ, Price AL, Neale BM, of the
697 Psychiatric Genomics Consortium SWG, et al. LD Score regression distinguishes confounding from polygenicity
698 in genome-wide association studies. *Nature genetics*. 2015; 47(3):291.
- 699 **Burt A**. Perspective: sex, recombination, and the efficacy of selection—was Weismann right? *Evolution*. 2000;
700 54(2):337–351.
- 701 **Chan AH**, Jenkins PA, Song YS. Genome-wide fine-scale recombination rate variation in *Drosophila melanogaster*.
702 *PLoS genetics*. 2012; 8(12):e1003090.
- 703 **Chan J**, Perrone V, Spence JP, Jenkins PA, Mathieson S, Song YS. A Likelihood-Free Inference Framework for
704 Population Genetic Data using Exchangeable Neural Networks. *bioRxiv*. 2018; [https://www.biorxiv.org/
705 content/early/2018/11/05/267211](https://www.biorxiv.org/content/early/2018/11/05/267211), doi: 10.1101/267211.
- 706 **Chollet F**, et al., Keras. GitHub; 2015. <https://github.com/fchollet/keras>.
- 707 **Cameron JM**, Ratnappan R, Bailin S. The Many Landscapes of Recombination in *Drosophila melanogaster*.
708 *PLOS Genetics*. 2012 10; 8(10):1–21. <https://doi.org/10.1371/journal.pgen.1002905>, doi: 10.1371/jour-
709 nal.pgen.1002905.
- 710 **Corbett-Detig RB**, Hartl DL. Population Genomics of Inversion Polymorphisms in *Drosophila melanogaster*.
711 *PLOS Genetics*. 2012 12; 8(12):1–15. <https://doi.org/10.1371/journal.pgen.1003056>, doi: 10.1371/jour-
712 nal.pgen.1003056.
- 713 **Dozhzhansky T**, Epling C. The suppression of crossing over in inversion heterozygotes of *Drosophila pseudoob-*
714 *scura*. *Proceedings of the National Academy of Sciences of the United States of America*. 1948; 34(4):137.
- 715 **Elyashiv E**, Sattath S, Hu TT, Strutsovsky A, McVicker G, Andolfatto P, Coop G, Sella G. A genomic map of the
716 effects of linked selection in *Drosophila*. *PLoS genetics*. 2016; 12(8):e1006130.

- 717 **Felsenstein J.** The evolutionary advantage of recombination. *Genetics*. 1974; 78(2):737–756.
- 718 **Flagel L, Brandvain Y, Schrider DR.** The Unreasonable Effectiveness of Convolutional Neural Networks in
719 Population Genetic Inference. *Molecular Biology and Evolution*. 2018 12; 36(2):220–238. <https://dx.doi.org/10.1093/molbev/msy224>, doi: 10.1093/molbev/msy224.
720
- 721 **Fuller ZL, Koury SA, Leonard CJ, Young RE, Ikegami K, Westlake J, Richards S, Schaeffer SW, Phadnis N.** Extensive
722 recombination suppression and chromosome-wide differentiation of a segregation distorter in *Drosophila*.
723 *bioRxiv*. 2018; <https://www.biorxiv.org/content/early/2018/12/21/504126>, doi: 10.1101/504126.
- 724 **Gao F, Ming C, Hu W, Li H.** New software for the fast estimation of population recombination rates (FastEPRR) in
725 the genomic era. *G3: Genes, Genomes, Genetics*. 2016; 6(6):1563–1571.
- 726 **Gay J, Myers S, McVean G.** Estimating meiotic gene conversion rates from population genetic data. *Genetics*.
727 2007; 177(2):881–894.
- 728 **Gravel S, Henn BM, Gutenkunst RN, Indap AR, Marth GT, Clark AG, Yu F, Gibbs RA, Bustamante CD.** Demographic
729 history and rare allele sharing among human populations. *Proceedings of the National Academy of Sciences*.
730 2011; 108(29):11983–11988. <https://www.pnas.org/content/108/29/11983>, doi: 10.1073/pnas.1019276108.
- 731 **Graves A, Jaitly N, Mohamed A.** Hybrid speech recognition with Deep Bidirectional LSTM. In: *2013 IEEE Workshop*
732 *on Automatic Speech Recognition and Understanding, Olomouc, Czech Republic, December 8-12, 2013*; 2013. p.
733 273–278. <https://doi.org/10.1109/ASRU.2013.6707742>, doi: 10.1109/ASRU.2013.6707742.
- 734 **Haenel Q, Laurentino TG, Roesti M, Berner D.** Meta-analysis of chromosome-scale crossover rate variation
735 in eukaryotes and its significance to evolutionary genomics. *Molecular Ecology*. 2018; 27(11):2477–2497.
736 <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.14699>, doi: 10.1111/mec.14699.
- 737 **Hahn MW.** *Molecular population genetics*. Sinauer Associates; 2018.
- 738 **Hartfield M, Otto SP.** Recombination and hitchhiking of deleterious alleles. *Evolution: International Journal of*
739 *Organic Evolution*. 2011; 65(9):2421–2434.
- 740 **Hill WG, Robertson A.** The effect of linkage on limits to artificial selection. *Genetics Research*. 1966; 8(3):269–294.
- 741 **Hilliker AJ, Harauz G, Reaume AG, Gray M, Clark SH, Chovnick A.** Meiotic gene conversion tract length distribution
742 within the rosy locus of *Drosophila melanogaster*. *Genetics*. 1994; 137(4):1019–1026.
- 743 **Hinton G, Deng L, Yu D, Dahl G, rahman Mohamed A, Jaitly N, Senior A, Vanhoucke V, Nguyen P, Sainath T,**
744 **Kingsbury B.** Deep Neural Networks for Acoustic Modeling in Speech Recognition. *Signal Processing Magazine*.
745 2012; .
- 746 **Hudson RR.** Estimation the recombination parameter of a finite population model without selection. *Genetical*
747 *Research*. 1987; 50:245–250.
- 748 **Hudson RR.** Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*. 2002
749 Feb; 18(2):337–338.
- 750 **Hudson RR, Kaplan NL.** Statistical properties of the number of recombination events in the history of a sample
751 of DNA sequences. *Genetics*. 1985; 111(1):147–164.
- 752 **Hunter N, Aguilera A, Rothstein R.** *Molecular Genetics of Recombination*. . 2006; .
- 753 **Jaenike J.** Sex chromosome meiotic drive. *Annual Review of Ecology and Systematics*. 2001; 32(1):25–49.
- 754 **Jozefowicz R, Zaremba W, Sutskever I.** An empirical exploration of recurrent network architectures. In: *Internat-*
755 *ional Conference on Machine Learning*; 2015. p. 2342–2350.
- 756 **Kelleher J, Etheridge AM, McVean G.** Efficient Coalescent Simulation and Genealogical Analysis for Large Sample
757 Sizes. *PLoS Computational Biology*. 2016 May; 12(5):e1004842. <https://doi.org/10.1371/journal.pcbi.1004842>,
758 doi: 10.1371/journal.pcbi.1004842.
- 759 **Kern AD, Schrider DR.** Discoal: flexible coalescent simulations with selection. *Bioinformatics*. 2016; 32(24):3839–
760 3841.
- 761 **Kern AD, Schrider DR.** diploS/HIC: An Updated Approach to Classifying Selective Sweeps. *G3: Genes, Genomes,*
762 *Genetics*. 2018; 8(6):1959–1970. <http://www.g3journal.org/content/8/6/1959>, doi: 10.1534/g3.118.200262.

- 763 **Kim Y, Nielsen R.** Linkage disequilibrium as a signature of selective sweeps. *Genetics*. 2004; 167(3):1513–1524.
- 764 **Kirkpatrick M.** How and why chromosome inversions evolve. *PLoS biology*. 2010; 8(9):e1000501.
- 765 **Kirkpatrick M, Barton N.** Chromosome inversions, local adaptation and speciation. *Genetics*. 2006; 173(1):419–
766 434.
- 767 **Kondrashov AS.** Selection against harmful mutations in large sexual and asexual populations. *Genetics*
768 *Research*. 1982; 40(3):325–332.
- 769 **Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, Jonasdottir A, Walters GB, Jonasdottir A,**
770 **Gylfason A, Kristinsson KT, et al.** Fine-scale recombination rate differences between sexes, populations and
771 individuals. *Nature*. 2010; 467(7319):1099.
- 772 **Krizhevsky A, Sutskever I, Hinton GE.** ImageNet Classification with Deep Convolutional Neural Net-
773 works. In: Pereira F, Burges CJC, Bottou L, Weinberger KQ, editors. *Advances in Neural Informa-*
774 *tion Processing Systems 25* Curran Associates, Inc.; 2012.p. 1097–1105. [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)
775 [4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf).
- 776 **Kulathinal RJ, Stevison LS, Noor MA.** The genomics of speciation in *Drosophila*: diversity, divergence, and
777 introgression estimated using low-coverage genome sequencing. *PLoS genetics*. 2009; 5(7):e1000550.
- 778 **Lack JB, Cardeno CM, Crepeau MW, Taylor W, Corbett-Detig RB, Stevens KA, Langley CH, Pool JE.** The *Drosophila*
779 *Genome Nexus*: A Population Genomic Resource of 623 *Drosophila melanogaster* Genomes, Including 197
780 from a Single Ancestral Range Population. *Genetics*. 2015; 199(4):1229–1241. [http://www.genetics.org/](http://www.genetics.org/content/199/4/1229)
781 [content/199/4/1229](http://www.genetics.org/content/199/4/1229), doi: 10.1534/genetics.115.174664.
- 782 **Langley CH, Stevens K, Cardeno C, Lee YCG, Schrider DR, Pool JE, Langley SA, Suarez C, Corbett-Detig RB,**
783 **Kolaczowski B, et al.** Genomic variation in natural populations of *Drosophila melanogaster*. *Genetics*. 2012;
784 192(2):533–598.
- 785 **Lecun Y, Bottou L, Bengio Y, Haffner P.** Gradient-based learning applied to document recognition. In: *Proceedings*
786 *of the IEEE*; 1998. p. 2278–2324.
- 787 **Li N, Stephens M.** Modeling linkage disequilibrium and identifying recombination hotspots using single-
788 nucleotide polymorphism data. *Genetics*. 2003; 165(4):2213–2233.
- 789 **Lin K, Futschik A, Li H.** A fast estimate for the population recombination rate based on regression. *Genetics*.
790 2013; p. genetics–113.
- 791 **Liu X, Fu YX.** Exploring population size changes using SNP frequency spectra. *Nature Genetics*. 2015 04; 47:555
792 EP –. <https://doi.org/10.1038/ng.3254>.
- 793 **Lowry DB, Willis JH.** A widespread chromosomal inversion polymorphism contributes to a major life-history
794 transition, local adaptation, and reproductive isolation. *PLoS biology*. 2010; 8(9):e1000500.
- 795 **Mather K.** Crossing-over. *Biological Reviews*. 1938; 13(3):252–292.
- 796 **McVean G, Awadalla P, Fearnhead P.** A coalescent-based method for detecting and estimating recombination
797 from gene sequences. *Genetics*. 2002; 160(3):1231–1241.
- 798 **Miller DE, Cook KR, Arvanitakis AV, Hawley RS.** Third Chromosome Balancer Inversions Disrupt Protein-Coding
799 Genes and Influence Distal Recombination Events in *Drosophila melanogaster*. *G3: Genes, Genomes, Genetics*.
800 2016; 6(7):1959–1967. <https://www.g3journal.org/content/6/7/1959>, doi: 10.1534/g3.116.029330.
- 801 **Myers S, Bottolo L, Freeman C, McVean G, Donnelly P.** A fine-scale map of recombination rates and hotspots
802 across the human genome. *Science*. 2005; 310(5746):321–324.
- 803 **Myers SR, Griffiths RC.** Bounds on the minimum number of recombination events in a sample history. *Genetics*.
804 2003; 163(1):375–394.
- 805 **Noor MA, Grams KL, Bertucci LA, Reiland J.** Chromosomal inversions and the reproductive isolation of species.
806 *Proceedings of the National Academy of Sciences*. 2001; 98(21):12084–12088.
- 807 **Novitski E, Braver G.** An analysis of crossing over within a heterozygous inversion in *Drosophila melanogaster*.
808 *Genetics*. 1954; 39(2):197.

- 809 **O'Reilly PF**, Birney E, Balding DJ. Confounding between recombination and selection, and the Ped/Pop method
810 for detecting selection. *Genome research*. 2008; 18(8):1304–1313.
- 811 **Parsch J**, Meiklejohn CD, Hartl DL. Patterns of DNA sequence variation suggest the recent action of positive
812 selection in the janus-ocnus region of *Drosophila simulans*. *Genetics*. 2001; 159(2):647–657.
- 813 **Pool JE**, Corbett-Detig RB, Sugino RP, Stevens KA, Cardeno CM, Crepeau MW, Duchon P, Emerson JJ, Saelao P,
814 Begun DJ, Langley CH. Population Genomics of Sub-Saharan *Drosophila melanogaster*: African Diversity and
815 Non-African Admixture. *PLOS Genetics*. 2012 12; 8(12):1–24. <https://doi.org/10.1371/journal.pgen.1003080>,
816 doi: 10.1371/journal.pgen.1003080.
- 817 **Presgraves DC**, Gérard PR, Cherukuri A, Lyttle TW. Large-scale selective sweep among segregation distorter
818 chromosomes in African populations of *Drosophila melanogaster*. *PLoS genetics*. 2009; 5(5):e1000463.
- 819 **Price AL**, Tandon A, Patterson N, Barnes KC, Rafaels N, Ruczinski I, Beaty TH, Mathias R, Reich D, Myers S.
820 Sensitive detection of chromosomal segments of distinct ancestry in admixed populations. *PLoS genetics*.
821 2009; 5(6):e1000519.
- 822 **Przeworski M**, Wall JD. Why is there so little intragenic linkage disequilibrium in humans? *Genetics Research*.
823 2001; 77(2):143–151.
- 824 **R Core Team**. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing,
825 Vienna, Austria; 2018, <https://www.R-project.org>.
- 826 **Rieseberg LH**. Chromosomal rearrangements and speciation. *Trends in ecology & evolution*. 2001; 16(7):351–
827 358.
- 828 **Ritz KR**, Noor MA, Singh ND. Variation in recombination rate: adaptive or not? *Trends in Genetics*. 2017;
829 33(5):364–374.
- 830 **Rogers AR**. How population growth affects linkage disequilibrium. *Genetics*. 2014; 197(4):1329–1341.
- 831 **Russakovsky O**, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC,
832 Fei-Fei L. ImageNet Large Scale Visual Recognition Challenge. *Int J Comput Vision*. 2015 Dec; 115(3):211–252.
833 <http://dx.doi.org/10.1007/s11263-015-0816-y>, doi: 10.1007/s11263-015-0816-y.
- 834 **Schiffels S**, Durbin R. Inferring human population size and separation history from multiple genome sequences.
835 *Nature Genetics*. 2014 06; 46:919 EP –. <https://doi.org/10.1038/ng.3015>.
- 836 **Schrider DR**, Ayroles J, Matute DR, Kern AD. Supervised machine learning reveals introgressed loci in the
837 genomes of *Drosophila simulans* and *D. sechellia*. *PLoS genetics*. 2018; 14(4):e1007341.
- 838 **Schrider DR**, Kern AD. Supervised Machine Learning for Population Genetics: A New Paradigm. *Trends in*
839 *Genetics*. 2018 Apr; 34(4):301–312. <https://doi.org/10.1016/j.tig.2017.12.005>, doi: 10.1016/j.tig.2017.12.005.
- 840 **Schrider DR**, Mendes FK, Hahn MW, Kern AD. Soft shoulders ahead: spurious signatures of soft and partial
841 selective sweeps result from linked hard sweeps. *Genetics*. 2015; 200(1):267–284.
- 842 **Schumer M**, Xu C, Powell DL, Durvasula A, Skov L, Holland C, Blazier JC, Sankararaman S, Andolfatto P, Rosen-
843 thal GG, Przeworski M. Natural selection interacts with recombination to shape the evolution of hybrid
844 genomes. *Science*. 2018; 360(6389):656–660. <https://science.sciencemag.org/content/360/6389/656>, doi:
845 10.1126/science.aar3684.
- 846 **Singh ND**, Stone EA, Aquadro CF, Clark AG. Fine-scale heterogeneity in crossover rate in the garnet-scalloped
847 region of the *Drosophila melanogaster* X chromosome. *Genetics*. 2013; 194(2):375–387.
- 848 **Slatkin M**. Linkage disequilibrium in growing and stable populations. *Genetics*. 1994; 137(1):331–336.
- 849 **Smith KN**, Nicolas A. Recombination at work for meiosis. *Current opinion in genetics & development*. 1998;
850 8(2):200–211.
- 851 **Sturtevant A**. A case of rearrangement of genes in *Drosophila*. *Proceedings of the National Academy of*
852 *Sciences of the United States of America*. 1921; 7(8):235.
- 853 **Sutskever I**, Vinyals O, Le QV. Sequence to Sequence Learning with Neural Networks. In: *Proceedings of the 27th*
854 *International Conference on Neural Information Processing Systems - Volume 2 NIPS'14*, Cambridge, MA, USA: MIT
855 Press; 2014. p. 3104–3112. <http://dl.acm.org/citation.cfm?id=2969033.2969173>.

- 856 **Szegedy C**, Liu W, Jia Y, Sermanet P, Reed SE, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A. Going deeper
857 with convolutions. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA,*
858 *June 7-12, 2015*; 2015. p. 1–9. <https://doi.org/10.1109/CVPR.2015.7298594>, doi: 10.1109/CVPR.2015.7298594.
- 859 **Tennessen JA**, Bigham AW, O'Connor TD, Fu W, Kenny EE, Gravel S, McGee S, Do R, Liu X, Jun G, Kang HM,
860 Jordan D, Leal SM, Gabriel S, Rieder MJ, Abecasis G, Altshuler D, Nickerson DA, Boerwinkle E, Sunyaev S, et al.
861 Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes. *Science*.
862 2012; 337(6090):64–69. <https://science.sciencemag.org/content/337/6090/64>, doi: 10.1126/science.1219240.
- 863 **Terhorst J**, Kamm JA, Song YS. Robust and scalable inference of population history from hundreds of unphased
864 whole genomes. *Nature Genetics*. 2016 12; 49:303 EP -. <https://doi.org/10.1038/ng.3748>.
- 865 **Vincent P**, Larochelle H, Bengio Y, Manzagol PA. Extracting and Composing Robust Features with Denoising
866 Autoencoders. In: *Proceedings of the 25th International Conference on Machine Learning ICML '08, New York, NY,*
867 *USA: ACM; 2008*. p. 1096–1103. <http://doi.acm.org/10.1145/1390156.1390294>, doi: 10.1145/1390156.1390294.
- 868 **Wakeley J**. Using the variance of pairwise differences to estimate the recombination rate. *Genetics Research*.
869 1997; 69(1):45–48.
- 870 **Wall JD**. A comparison of estimators of the population recombination rate. *Molecular Biology and Evolution*.
871 2000; 17(1):156–163.
- 872 **Wang RJ**, Gray MM, Parmenter MD, Broman KW, Payseur BA. Recombination rate variation in mice from an
873 isolated island. *Molecular ecology*. 2017; 26(2):457–470.
- 874 **Winckler W**, Myers SR, Richter DJ, Onofrio RC, McDonald GJ, Bontrop RE, McVean GA, Gabriel SB, Reich D,
875 Donnelly P, et al. Comparison of fine-scale recombination rates in humans and chimpanzees. *Science*. 2005;
876 308(5718):107–111.
- 877 **Wiuf C**. On the minimum number of topologies explaining a sample of DNA sequences. *Theoretical population*
878 *biology*. 2002; 62(4):357–363.
- 879 **Yi X**, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, Zheng H,
880 Liu T, He W, Li K, Luo R, Nie X, Wu H, Zhao M, Cao H, Zou J, et al. Sequencing of 50 human exomes reveals
881 adaptation to high altitude. *Science (New York, NY)*. 2010 07; 329(5987):75–78. <https://www.ncbi.nlm.nih.gov/pubmed/20595611>, doi: 10.1126/science.1190371.
- 882
883 **Zelkowski M**, Olson M, Wang M, P Pawlowski W. Diversity and Determinants of Meiotic Recombination
884 Landscapes. *Trends in Genetics*. 2019 04; doi: 10.1016/j.tig.2019.02.002.

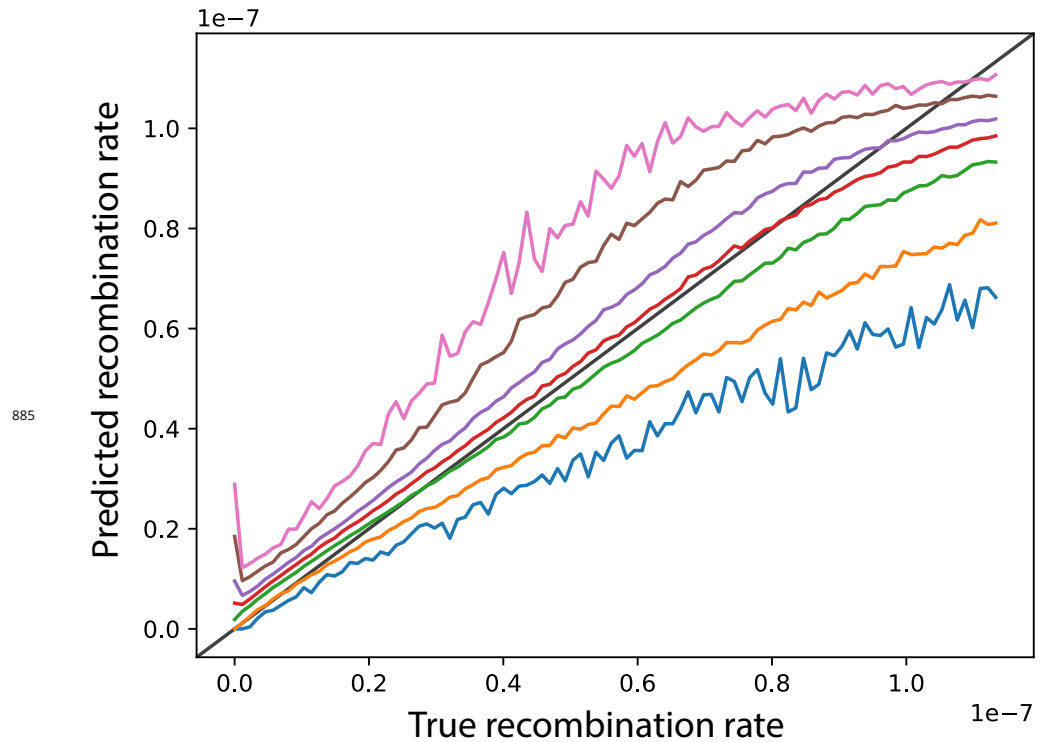


Figure 1-Figure supplement 1. Parametric bootstrapping results as implemented by ReLERNN. Lines represent the minimum (blue), lower 5% (orange), lower 25% (green), median (red), upper 25% (purple), upper 95% (brown), and maximum (pink) bounds for each of 1000 replicate simulations and predictions (y-axis) across 100 recombination rate bins (x-axis)

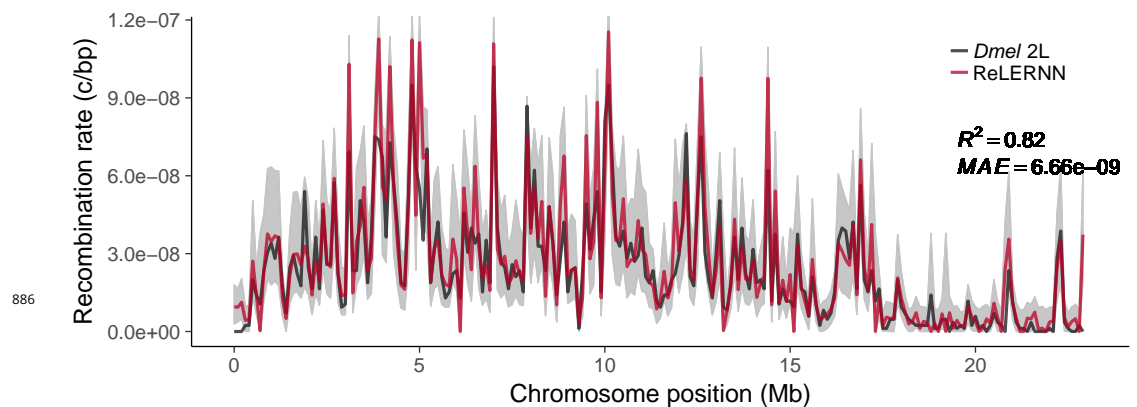


Figure 2-Figure supplement 1. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 4$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

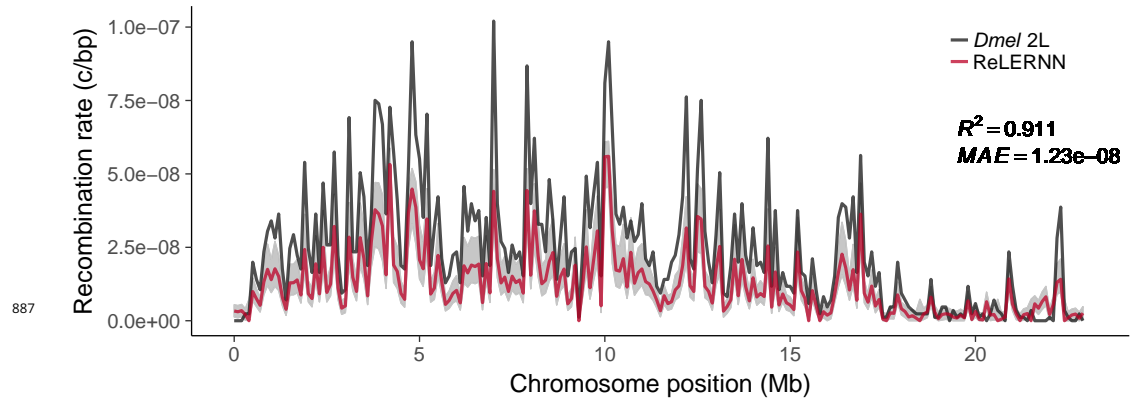


Figure 2-Figure supplement 2. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Here the per-base mutation rate was assumed to be 50% less than the rate used for simulation. Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

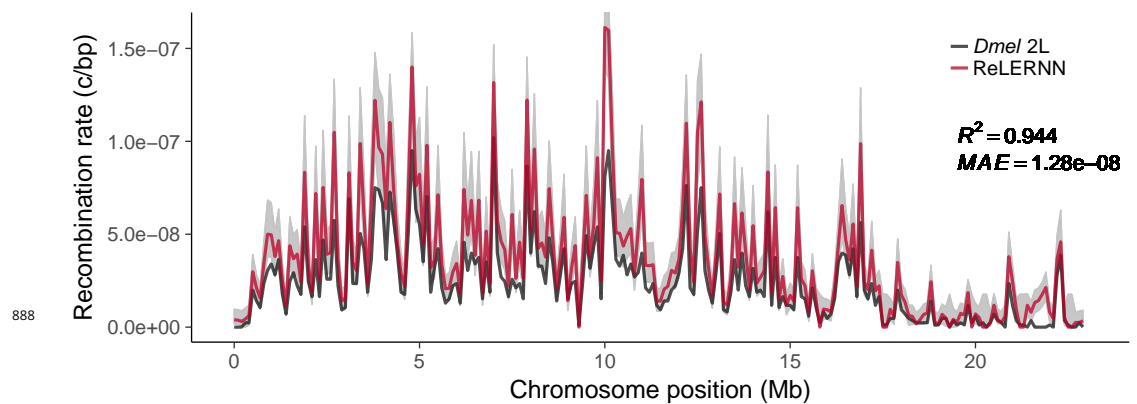


Figure 2-Figure supplement 3. Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Comeron et al., 2012). Here the per-base mutation rate was assumed to be 50% greater than the rate used for simulation. Gray ribbons represent 95% confidence intervals. R^2 is reported for the general linear model of predicted rates on true rates and mean absolute error was calculated across all 100 kb windows.

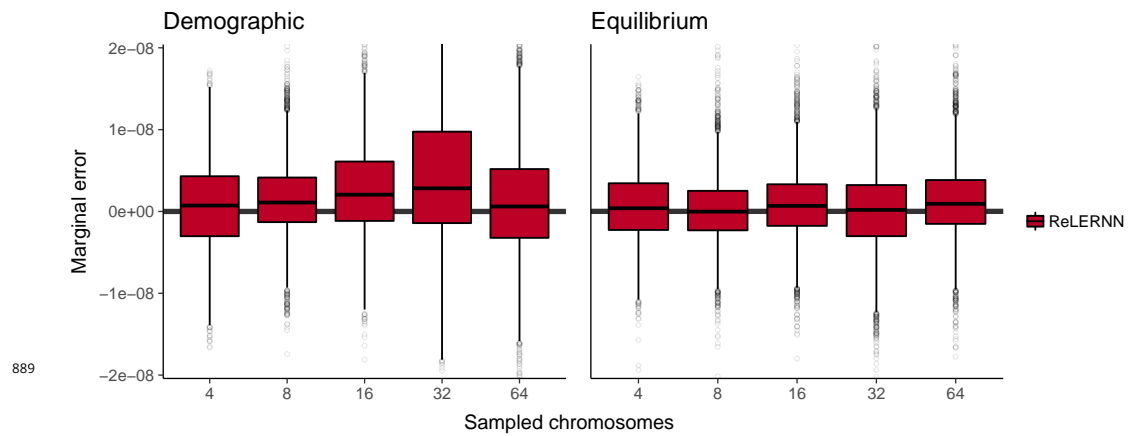


Figure 4-Figure supplement 1. Distribution of marginal errors attributed to model misspecification across 5000 simulated chromosomes. Predictions were made by training on equilibrium simulations and testing on sequences simulated under a demographic model (left) or training on demographic simulations and testing on sequences simulated under equilibrium (right). Here, marginal errors are represented as $\epsilon_m - \epsilon_c$, where ϵ_m and ϵ_c are equal to $|r_{predicted} - r_{true}|$ when the model is misspecified and correctly specified, respectively. Sampled chromosomes indicate the number of independent sequences that were sampled from each msprime (Kelleher *et al.*, 2016) coalescent simulation.

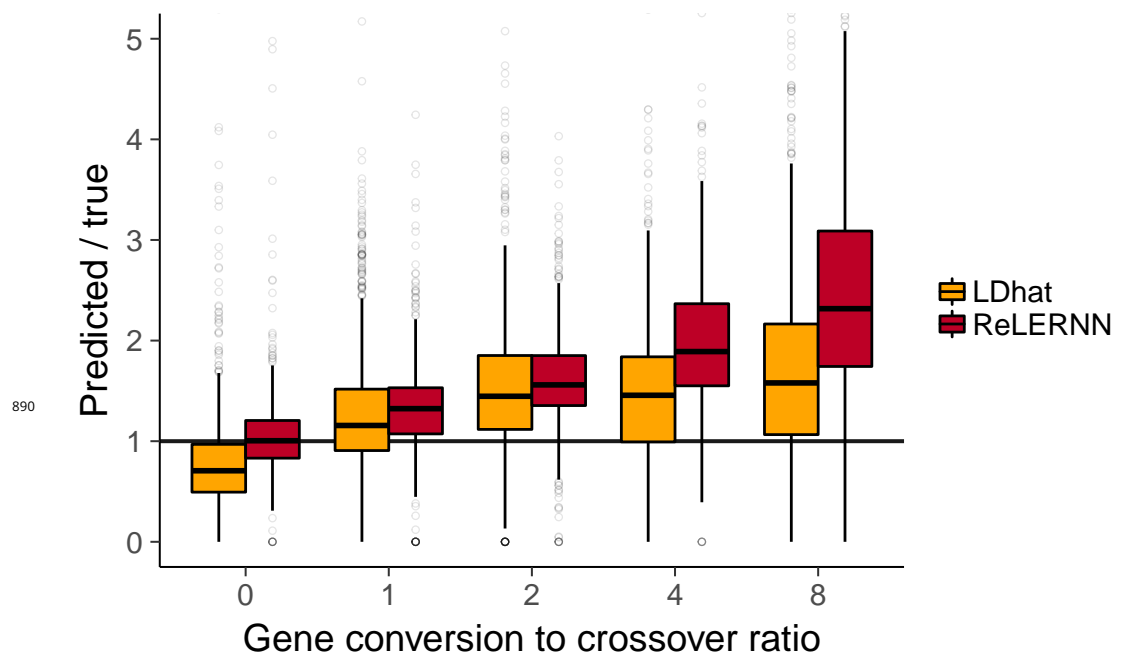


Figure 4-Figure supplement 2. Distribution of predicted rates of recombination over true rates for 5000 examples simulated with gene conversion and $n = 8$. The ratio of gene conversion to crossovers was drawn from $U(0, c)$, with $c \in \{0, 1, 2, 4, 8\}$. Gene conversion tract lengths were fixed at 352 bp, and all simulations were completed in ms (Hudson, 2002).

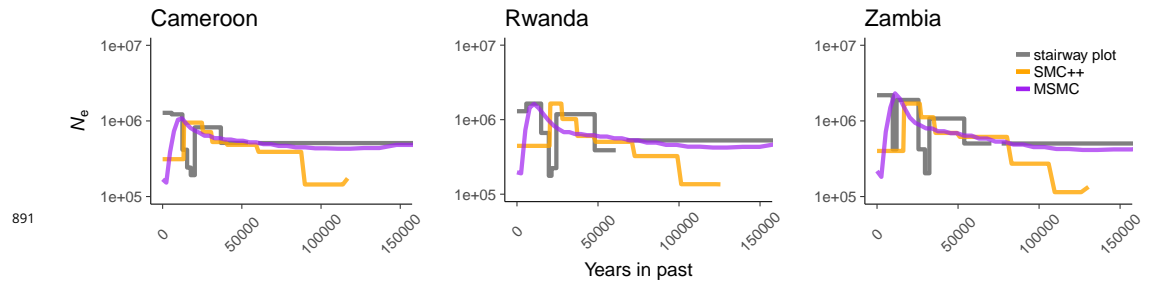


Figure 5-Figure supplement 1. Historical population size estimates were inferred for Cameroon, Rwanda, and Zambia using three separate methods, all of which disagree with one another. Inferences are based on 10 samples for both stairwayplot (grey line) and SMC++ (orange line), and 2 samples for MSMC (purple line).

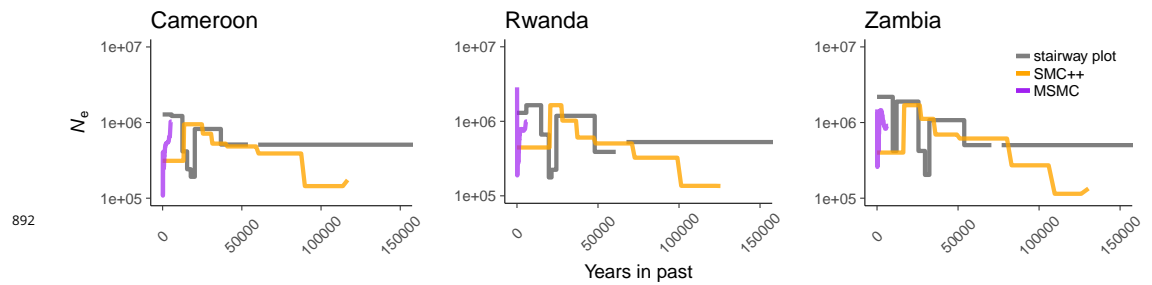


Figure 5-Figure supplement 2. Historical population size estimates were inferred for Cameroon, Rwanda, and Zambia using three separate methods, all of which disagree with one another. Here, inferences are based on 10 samples for both stairwayplot (grey line) and SMC++ (orange line), and 10 samples for MSMC (purple line).

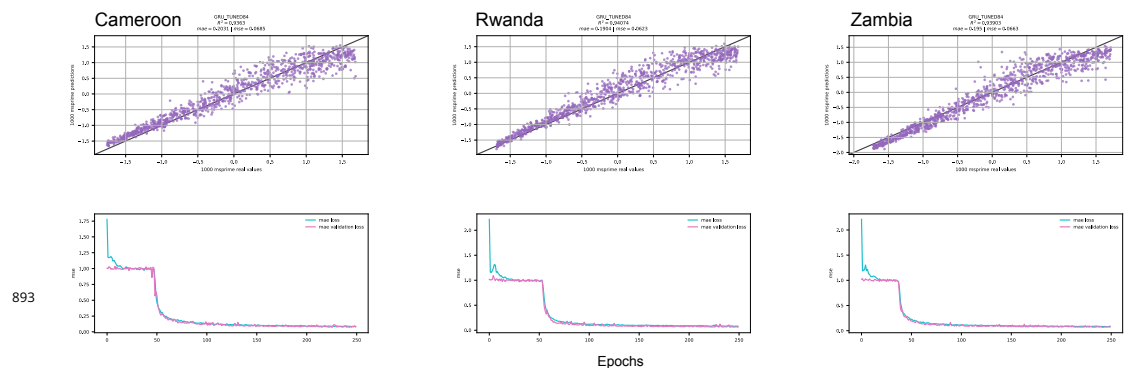


Figure 5-Figure supplement 3. ReLERNN test results for Cameroon, Rwanda, and Zambia when trained under assumptions of mutation-drift equilibrium. Scatter plots (top) show raw (unnormalized) predictions for per-base recombination rates for 1000 test examples. Mean absolute error and mean squared error are calculated for each population. Line graphs (bottom) show the decrease in the mean absolute error over time (epochs) for both the training set (blue lines) and the validation set (purple lines).

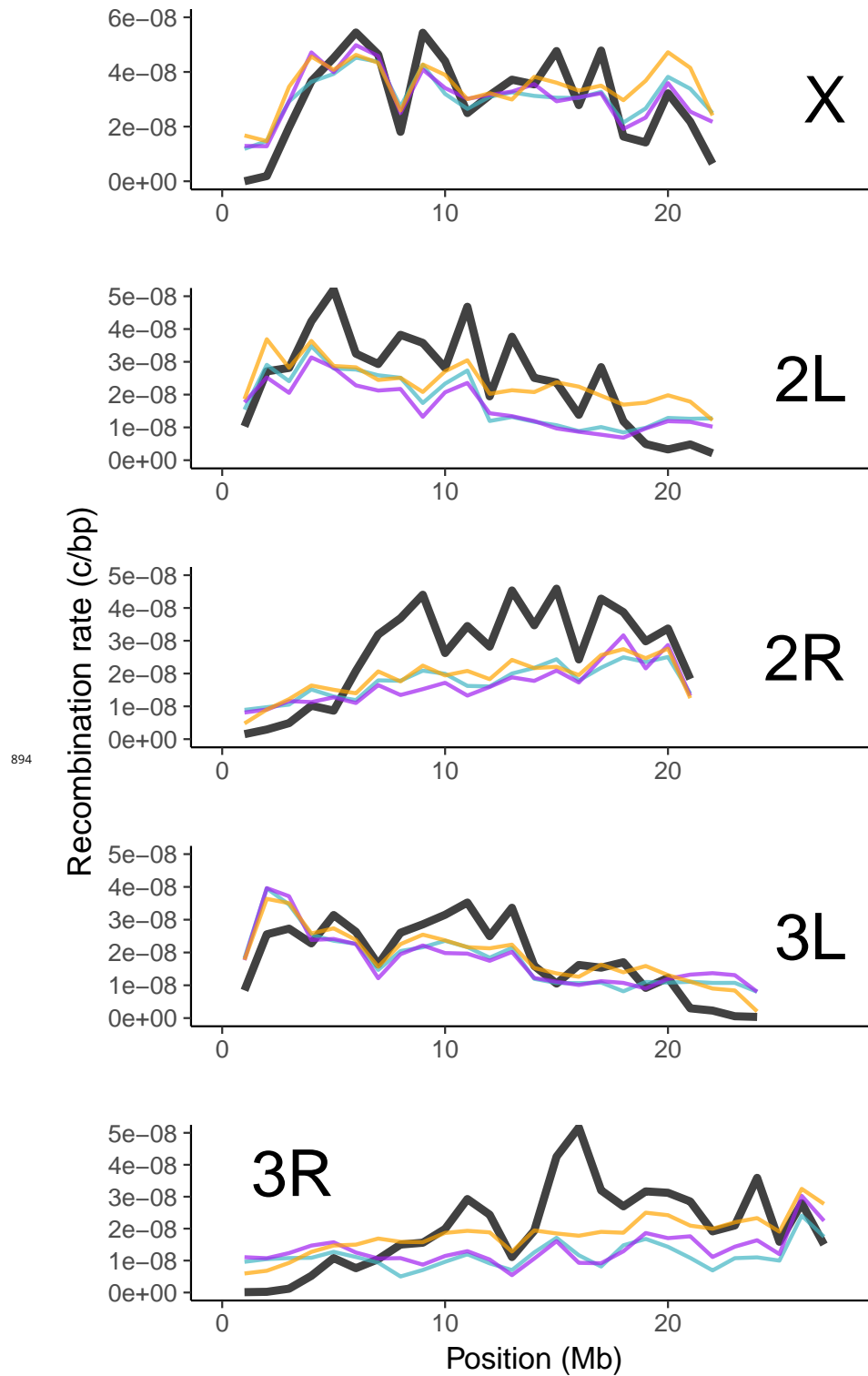


Figure 5-Figure supplement 4. Genome-wide recombination landscapes for *D. melanogaster* populations from Cameroon (teal lines), Rwanda (purple lines), and Zambia (orange lines). Rates are compared to those experimentally derived by *Comeron et al. (2012)* (black lines). All rates have been scaled to 1 Mb windows by using a weighted average (see Materials and Methods).

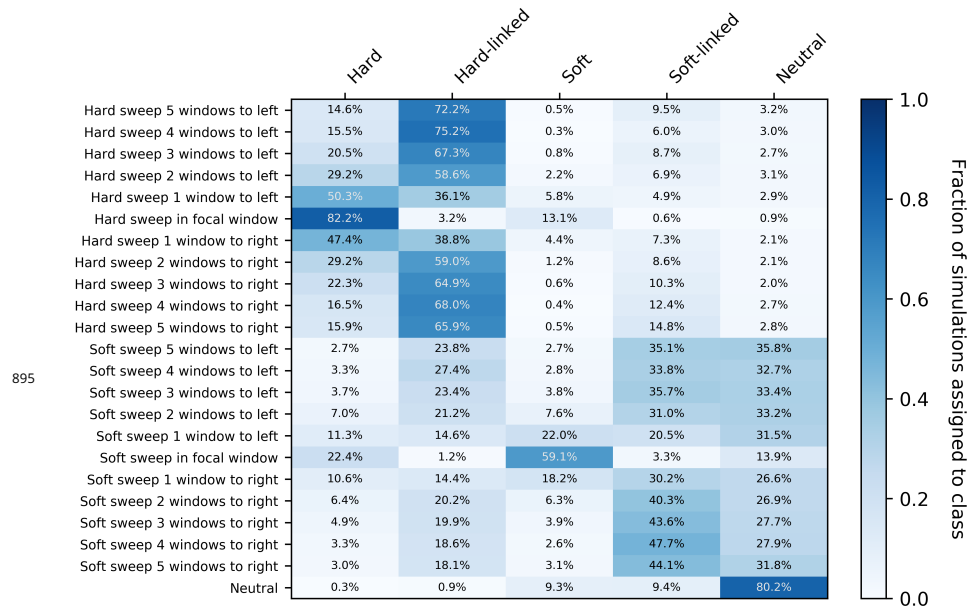


Figure 5—Figure supplement 5. Confusion matrix showing the fraction of test simulation windows assigned to each of five prediction categories by diploS/HIC (*Kern and Schrider, 2018*): hard, hard-linked, soft, soft-linked, and neutral. The y-axis shows the location of the window being classified relative to the selected window.

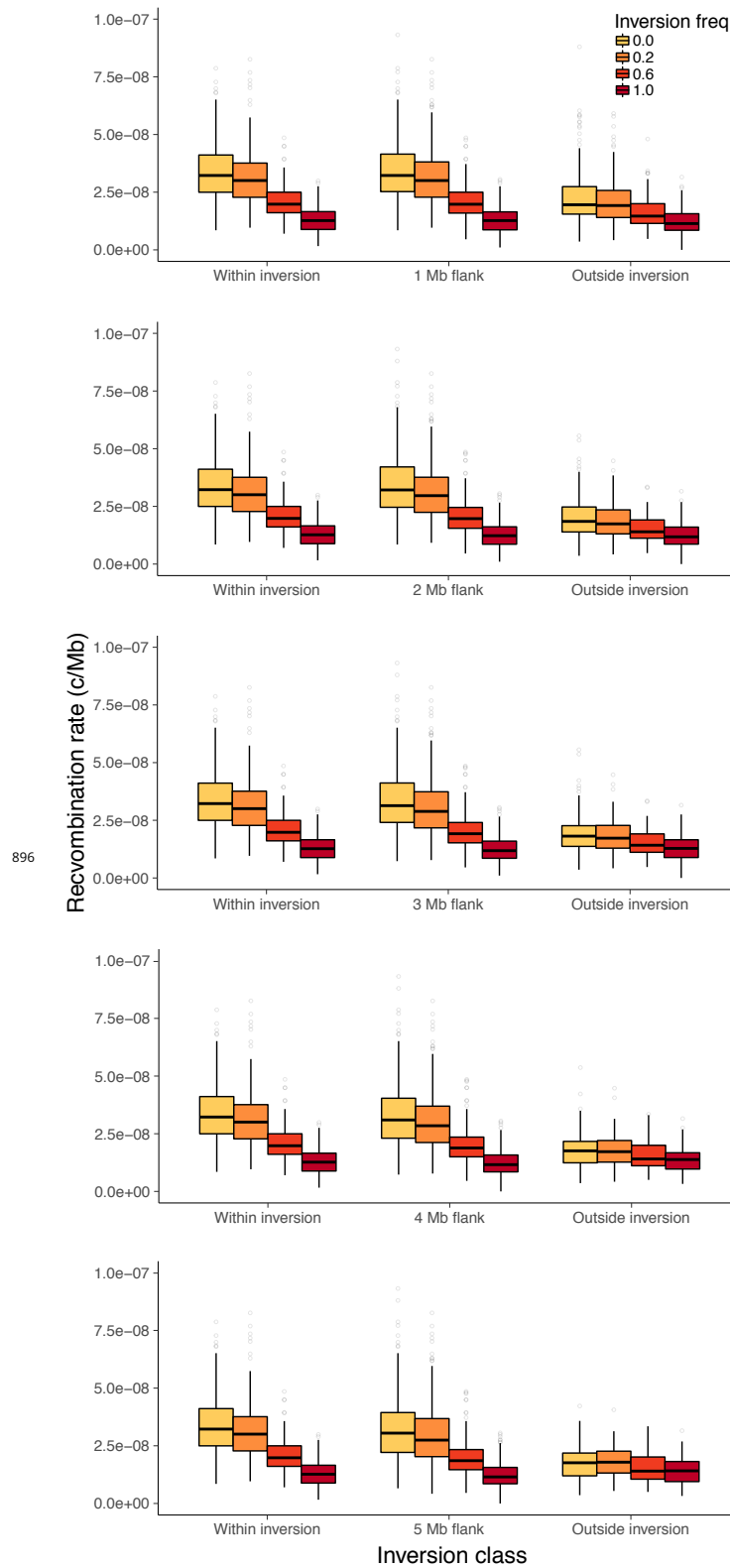


Figure 6-Figure supplement 1. Recombination rate estimates using flanking window sizes from 1-5 Mb. Rates are shown for genomic windows within the inversion, within regions flanking the inversion, and for regions outside both the inversion and flanking regions. All estimates are from chromosome 2L with *In(2L)t* sampled at different inversion frequencies