

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*) or LDhelmet (*Chan et al., 2012*), 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 genotypes 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 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 the face of

93 demographic model misspecification. We then apply ReLERNN to population genomic data from
94 African samples of *Drosophila melanogaster*. We demonstrate that the landscape of recombination
95 is largely conserved in this species, yet individual regions of the genome show marked population-
96 specific differences. Finally, we found that chromosomal inversion frequencies directly impact the
97 inferred rate of recombination, and we demonstrate that the role for inversions in suppressing
98 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 chromosomes. Briefly, ReLERNN provides an
104 end-to-end inferential pipeline for estimating a recombination landscape from a population sample:
105 it takes as input a user-filtered Variant Call Format (VCF) file of phased or unphased genotypes,
106 and from this estimates a set of simulation parameters reflective of the input samples. ReLERNN
107 then uses the coalescent simulation program, msprime (Kelleher *et al.*, 2016), to simulate training,
108 validation, and test data sets under either a user-supplied or an inferred demographic history,
109 seeking to mimic population genetic properties of the empirical samples. ReLERNN trains a specific
110 type of RNN, known as a Gated Recurrent Unit (GRU), to predict the per-base recombination rate
111 for 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 S1).

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 S2). While ReLERNN retains accuracy at small sample sizes, it
135 exhibits somewhat greater sensitivity to both the assumed per-base mutation rate and the assumed
136 maximum ratio of ρ to the population mutation parameter, θ —two mandatory assumptions.

137 To assess the degree of sensitivity to these mutation rate assumptions, we ran ReLERNN on
138 simulations using an assumed per-base mutation rate both 50% greater and 50% less than the
139 simulated (true) mutation rate. In both scenarios, ReLERNN predicts crossover rates that are highly
140 correlated with the simulated rates ($R^2 > 0.91$). However, in both scenarios MAE is inflated but still

141 modest, and the absolute rates of recombination are underpredicted ($R^2 = 0.91$; $MAE = 1.23 \times 10^{-8}$;
142 **Figure S3**) and overpredicted ($R^2 = 0.94$; $MAE = 1.28 \times 10^{-8}$; **Figure S4**) when assuming a mutation
143 rate less than or greater than the true per-base mutation rate, respectively. Together these results
144 suggest that ReLERNN is in fact learning information about the ratio of crossovers to mutations,
145 and while ReLERNN is highly robust to errant assumptions when predicting relative recombination
146 rates within a genome, caution must be taken when comparing absolute rates between organisms
147 with large differences in per-base mutation rate estimates. Crucially, we also show that ReLERNN
148 performs at least as well on unphased genotypes as it does on 100% correctly phased genotypes
149 ($W = 68.5$; $P = 0.17$; Mann-Whitney U test; **Figure S5**), suggesting that any effect of computational
150 phasing error can potentially be mitigated by unphasing the input genotypes.

151 **ReLERNN compares favorably to competing methods, especially for small sample 152 sizes and under model misspecification**

153 To assess the accuracy of ReLERNN relative to existing methods, we took a comparative approach
154 whereby we made predictions on the same set of simulated test chromosomes using methods that
155 differ broadly in their approaches. Specifically, we chose to compare ReLERNN against two types
156 of machine learning methods—a boosted regression method, FastEPRR (*Gao et al., 2016*), and
157 a convolutional neural network (CNN) recently described in *Flagel et al. (2018)*—and both LDhat
158 (*McVean et al., 2002*) and LDhelmet (*Chan et al., 2012*), two widely cited approximate-likelihood
159 methods. We independently simulated 10^5 chromosomes using msprime (*Kelleher et al., 2016*)
160 (parameters: $n \in \{4, 8, 16, 32, 64\}$, $priorLowsRho = 0.0$, $priorHighsRho = 5e^{-8} \times 1.25$, $priorLowsMu =$
161 $2.5e^{-8} \times 0.75$, $priorHighsMu = 2.5e^{-8} \times 1.25$, $ChromosomeLength = 3e^5$). Half of these were simulated
162 under demographic equilibrium and half were simulated under a realistic demographic model
163 (based on the out-of-Africa expansion of European humans; see Materials and Methods). We
164 show that ReLERNN outperforms all other methods, exhibiting significantly reduced absolute error
165 under both the demographic model and under equilibrium assumptions ($T \leq -31$; $P < 10^{-16}$; *post*
166 *hoc* Welch's two sample t -tests for all comparisons; **Figure 3**). Importantly, ReLERNN is also more
167 accurate than all methods we compared for each of the tested samples sizes, although all methods
168 generally performed well with larger sample sizes.

169 We also sought to assess the robustness of ReLERNN to demographic model misspecification,
170 whereby different generative models are used for simulating the training and test sets—e.g. training
171 on assumptions of demographic equilibrium when the test data was generated by a population
172 bottleneck. Methods robust to this type of misspecification are crucial, as the true demographic
173 history of a sample is often unknown and methods used to infer population size histories can
174 disagree or be unreliable (see **Figure S8**). Moreover, population size changes alter the landscape
175 of LD across the genome (e.g. *Slatkin, 1994*; *Rogers, 2014*), and thus have the potential to reduce
176 accuracy or produce biased recombination rate estimates.

177 To this end, we trained ReLERNN on examples generated under equilibrium and made pre-
178 dictions on 5000 chromosomes generated by the human demographic model specified above
179 (and also carried out the reciprocal experiment). We compared ReLERNN to the CNN, LDhat, and
180 LDhelmet, whereby all methods were similarly misspecified (see Materials and Methods). We found
181 that ReLERNN outperforms these methods under nearly all conditions, exhibiting significantly lower
182 absolute error under both directions of demographic model misspecification ($T \leq -26$; $P_{WTT} < 10^{-16}$
183 for all comparisons, with the exception of the comparison to LDhelmet using 16 chromosomes; **Fig-
184 ure 4**). Interestingly, we show that the error attributed to model misspecification (termed marginal
185 error; see Materials and Methods) is significantly greater when ReLERNN was trained on equilibrium
186 simulations and tested on demographic simulations than under the reciprocal misspecification
187 ($T = 26.3$; $P_{WTT} < 10^{-16}$; **Figure S6**). While this is true, it is important to note that marginal error is
188 quite modest in both directions of misspecification ($< 1.30 \times 10^{-9}$; **Figure S6**), suggesting that the
189 additional information gleaned from an informative demographic model is limited.

190 Differences in the ratio of homologous gene conversion events to crossovers can also bias the

191 inference of recombination rates, as conversion tracts break down LD within the prediction window
192 (*Gay et al., 2007; Przeworski and Wall, 2001*). We treated the effect of gene conversion as another
193 form of model misspecification by training on examples that lacked gene conversion and testing on
194 examples that included gene conversion. As ReLERNN uses msprime for all training simulations,
195 and msprime cannot currently simulate gene conversion, we generated all test set simulations
196 with ms (*Hudson, 2002*). We found that including gene conversion in our simulations biased our
197 predictions, resulting in an overestimate of the true recombination rate (*Figure S7*). Moreover,
198 the magnitude of this bias increased with the ratio of gene conversion events to crossovers. As
199 expected, we also observed a similar pattern of bias for LDhat, although the magnitude of bias for
200 LDhat was somewhat less than that exhibited by ReLERNN (*Figure S7*).

201 **Recombination landscapes are largely concordant among populations of African *D.*** 202 ***melanogaster***

203 Using our method, we characterized the genome-wide recombination landscapes of three popula-
204 tions of African *D. melanogaster* (sampled from Cameroon, Rwanda, and Zambia). Each population
205 was derived from the sequencing of 10 haploid embryos (detailed in *Lack et al., 2015; Pool et al.,*
206 *2012*), hence these data represent an excellent opportunity to exploit ReLERNN's high accuracy
207 on small sample sizes. We first sought to model the demographic history of each population, as
208 ReLERNN can simulate training data under demographic models inferred by three published soft-
209 ware methods—stairwayplot (*Liu and Fu, 2015*), SMC++ (*Terhorst et al., 2016*), and MSMC (*Schiffels*
210 *and Durbin, 2014*). Using all three methods, we show that inferred historical population sizes are
211 unreliable for these populations—no two methods recapitulate the same history, and the histories
212 generated by MSMC vary dramatically depending on the number of samples used (*Figure S8, Fig-*
213 *ure S9*). For these reasons, and because results from our simulations suggest that marginal error
214 due to demographic misspecification is quite low for our method (above; *Figure S6*), we decided to
215 simulate our training data under the 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 S10*), suggesting that we have sufficient power to discern fine-scale
229 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 mosomes 2L, 2R, 3L, 3R, and X, respectively (*Figure S11*; 1 Mb windows). These results differ from
234 those observed by *Chan et al. (2012)*, who compared 22 *D. melanogaster* sampled from the same
235 Rwandan population to the FlyBase map and found $R^2 = 0.55, 0.63, 0.45, 0.42, 0.41$ for the same chro-
236 mosomes. The minor differences we observed between methods for chromosomes 2L, 2R, 3L, and
237 the X chromosome can likely be attributed to the fact that we are comparing estimates from two
238 different methods, using different African flies, to a different experimentally derived map. However,
239 the larger differences found between methods for chromosome 3R seem less likely attributable
240 to methodological differences. Importantly, African *D. melanogaster* are known to harbor large

241 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.

245 An important cause of population-specific differences in recombination landscapes might be
246 population-specific differences in the frequencies of chromosomal inversions, as recombination is
247 expected to be strongly suppressed between standard and inversion arrangements. Segregating
248 inversions in *D. melanogaster* have been shown to affect broad patterns of chromosomal varia-
249 tion, and are thought to have quite recent origins when taken together (*Corbett-Detig and Hartl,*
250 **2012**). To test for an effect of inversion frequency on our measurement of recombination rates, we
251 resampled haploid genomes from Zambia to create sampled populations with the cosmopolitan
252 inversion *In(2L)t* segregating at varying frequencies, $p \in \{0.0, 0.2, 0.6, 1.0\}$. In Zambia, *In(2L)t* segre-
253 gates at $p = 0.22$ (*Lack et al., 2015*), suggesting that recombination within the inversion breakpoints
254 may be strongly suppressed in individuals with the inverted arrangement relative to those with
255 the standard arrangement. Moreover, *In(2L)t* arose recently, likely within the past 100,000 years
256 (*Corbett-Detig and Hartl, 2012*). For these reasons, we predict that the inferred recombination rate
257 should decrease as the low-frequency inverted arrangement is increasingly overrepresented in the
258 set of sampled chromosomes (i.e. as more of the samples contain the high-LD inverted arrange-
259 ments). As predicted, we found a strong effect of the sample frequency of *In(2L)t* on estimated rates
260 of recombination for chromosome 2L in Zambia (**Figure 6**). Recombination rates are negatively
261 correlated with inversion frequency in our sample, not only within the inversion, but also in regions
262 3 Mb outside the inversion (flanking regions) ($\rho_{Spearman's} = -1$; $P = 0.04$ for both comparisons). We
263 also see a similar negative correlation outside the flanking regions, although this association is
264 weakened relative to that within or flanking the inversion (**Figure 6**). Importantly, varying the size of
265 the flanking regions (from 1-5 Mb) produces patterns that are qualitatively identical, suggesting that
266 the effect of inversions on recombination suppression extends far beyond the inversion breakpoints
267 themselves (**Figure S13**).

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

282 To further explore population-specific differences in recombination landscapes we took a statis-
283 tical outlier approach, whereby we define two types of recombination rate outliers—global outliers
284 and population-specific outliers (see Materials and Methods). Global outliers are characterized by
285 windows with exceptionally high variance in rates of recombination between all three populations
286 (**Figure 5**; red triangles) while population-specific outliers are those windows where the rate of re-
287 combination in one population is strongly differentiated from the rates in the other two populations
288 (**Figure 5**; population-colored triangles). We find that population-specific outliers, but not global
289 outliers, are significantly enriched within inversions ($P = 0.005$; randomization test; **Figure 5**; grey
290 boxes). Moreover, this enrichment remains significant when extending the inversion boundaries
291 by up to 250 Kb ($P_{rand} \leq 0.004$). However, extending the inversion boundaries beyond 250 Kb, or

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

296 Selection is another important factor that may confound the inference of recombination rates.
297 For instance selective sweeps generate localized patterns of high LD on either side of the sweep site
298 (*Kim and Nielsen, 2004; Schrider et al., 2015*), thus regions flanking selective sweeps may mimic
299 regions of reduced recombination. Inasmuch population-specific selective sweeps are expected to
300 contribute to population-specific differences in recombination rate estimates. We used diploS/HIC
301 (*Kern and Schrider, 2018*) to identify hard and soft selective sweeps in our African *D. melanogaster*
302 populations, and we tested for an excess of recombination rate outliers overlapping with windows
303 classified as sweeps. In total, diploS/HIC classified 27.4%, 28.1%, and 26.8%, of all genomic windows
304 as selective sweeps (either "hard" or "soft") for Cameroon, Rwanda, and Zambia, respectively, when
305 looking at 5kb, non-overlapping windows. The associated False Discovery Rates (FDR) for calling
306 sweeps in these populations were appreciable: 33.9%, 33.1% and 34.7%, respectively (*Figure S12*).
307 As expected, windows classified as sweeps had significantly lower rates of recombination relative
308 to neutral windows in all three populations ($P_{WTT} \leq 10^{-16}$ for all comparisons; *Figure 7*). However,
309 we found that neither global nor population-specific outliers were enriched for selective sweeps
310 ($P_{rand} \geq 0.246$ for both comparisons), suggesting that, when treated as a class, recombination
311 rate outliers are not likely driven by sweeps in these populations. When treated separately (i.e.
312 independent permutation tests for each recombination rate outlier window), we identified 7 outliers
313 enriched for sweeps at the $P \leq 0.05$ threshold, corresponding to an expected FDR of 77%. However,
314 given our FDR for calling sweeps in these populations, our measure of the enrichment in overlap
315 with recombination rate outliers is likely to be conservative. Two of these outlier windows may
316 represent potential true positives; an outlier in Cameroon contains 5 out of 6 non-overlapping 5 kb
317 windows classified as "hard" sweeps, the second from Rwanda has 10 out of 12 windows classified
318 as "hard" sweeps ($P_{rand} = 0.0$ for both comparisons). These two recombination rate outlier windows
319 are potentially ripe for future studies on selective sweeps in these populations, and suggest that in
320 at least some instances, selection contributes to observed differences in estimates of recombination
321 rates between *Drosophila* populations.

322 Discussion

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

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

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

347 ReLERNN utilizes a variant of an ANN, known as a Gated Recurrent Unit (GRU), as its primary
348 technology. GRU networks excel at identifying temporal associations (Jozefowicz *et al.*, 2015), and
349 therefore we model our sequence alignment as a bidirectional time series, whereby each ordered
350 SNP represents a new time step along the chromosome. We also model the distance between
351 SNPs using a separate input tensor, and these two inputs are concatenated after passing through
352 the initial layers of the network (see **Figure 1** inlay). We demonstrated that ReLERNN can predict a
353 simulated recombination landscape with a high degree of accuracy ($R^2 = 0.93$; **Figure 2**), and that
354 these predictions remain high, even when using small sample sizes ($R^2 = 0.82$; **Figure S2**). These
355 predictions compared favorably to those made by a leading composite likelihood methods (LDhat
356 and LDhelmet; McVean *et al.*, 2002; Chan *et al.*, 2012), as well as other machine learning methods
357 (the CNN and FastEPRR; **Figure 3**). While the abstract nature of the data represented in its internal
358 layers constrains our ability to interpret the exact information ReLERNN relies on to inform its
359 predictions, our experiments using incorrect assumed mutation rates (**Figure S4**, **Figure S3**) suggests
360 that ReLERNN is potentially learning the relative ratio of recombination rates to mutation rates.
361 For these reasons, an extra caveat is warranted—use caution when interpreting the results from
362 ReLERNN as precise measures of the per-base recombination rate unless precise mutation rate
363 estimates are also known. Importantly, we also demonstrate that ReLERNN is just as accurate when
364 given unphased input genotypes as it is when provided with perfectly phased genotypes (**Figure S5**).

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,
376 and in so doing capture more of the information available than traditional methods that use
377 composite likelihoods of two-locus LD summaries. Additionally, there are clear opportunities for
378 future improvements to ReLERNN. For instance, our simulation studies demonstrated that the RNN
379 used by ReLERNN is also sensitive to gene conversion events (**Figure S7**), thus the joint estimation of
380 rates of recombination and gene conversion may be quite feasible. Ultimately, it remains far from
381 clear what network architectures will be best suited for population genetic inference, though we
382 remain optimistic that ANNs will prove useful for a variety of applications in the field.

383 A natural application of ReLERNN, due in part to its high accuracy with small sample sizes, was
384 to characterize and compare the recombination landscapes for multiple populations of African *D.*
385 *melanogaster*, for which few populations with large samples sizes are currently available. Previous
386 estimates of genome-wide fine-scale recombination maps in flies have focused on characterizing
387 recombination in experimental crosses (Comeron *et al.*, 2012), or by running LDhat (or the related
388 LDhelmet) on populations with relatively moderate sample sizes (i.e. ≥ 22 samples) (Chan *et al.*,
389 2012; Langley *et al.*, 2012). Here, we applied ReLERNN to three populations for which at least ten
390 haploid embryos were sequenced: Cameroon, Rwanda, and Zambia (Lack *et al.*, 2015; Pool *et al.*,
391 2012). Generally, recombination landscapes were well correlated among populations. Mean pair-
392 wise coefficients of determination among all three populations were $R^2 = 0.69, 0.61, 0.77, 0.43, 0.66$

393 for chromosomes 2L, 2R, 3L, 3R, and X, respectively. These correlations are notably lower than
394 those observed in humans (*Myers et al., 2005*) and mice (*Wang et al., 2017*), and one potential
395 biological cause for this large difference could be the cosmopolitan chromosomal inversions that
396 segregate in African *D. melanogaster* (*Corbett-Detig and Hartl, 2012; Lack et al., 2015*).

397 We demonstrated a significant negative association between inversion sample frequency and
398 recombination rate as inferred by ReLERNN through experimentally manipulating the frequency
399 of the inversion karyotype in our sample (*Figure 6*). Our results suggest that recombination
400 suppression extends well beyond the predicted breakpoints of the inversion (at least 5 Mb beyond
401 in the case of *In(2L)t*; *Figure S13*). This large-scale suppression of recombination due to inversions
402 in *Drosophila* has been observed both directly in experimental crosses (*Dobzhansky and Epling,*
403 *1948; Novitski and Braver, 1954; Kulathinal et al., 2009; Miller et al., 2016; Fuller et al., 2018*),
404 and indirectly from patterns of variation surrounding known inversion breakpoints (*Corbett-Detig*
405 *and Hartl, 2012; Langley et al., 2012*). Moreover, the extension of recombination rate differences
406 outside the inversion breakpoints may in part be driven the interchromosomal effect (*Lucchesi*
407 *and Suzuki, 1968*), whereby recombination suppression in heterozygous inversions acts to enhance
408 crossing over across the remaining chromosomes. While it is true that the negative relationship
409 between inversion frequency and recombination should only exist for inversions segregating at low
410 frequencies (e.g. crossover suppression is not expected in inversion homozygotes), we predict a
411 negative relationship to dominate in these populations, as the majority of polymorphic inversions
412 are young, segregate at low frequencies, and show elevated LD along their lengths perhaps due to
413 the actions of natural selection (*Corbett-Detig and Hartl, 2012; Lack et al., 2015*).

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

425 Interestingly, while we identified two individual outlier regions characterized by numerous
426 selective sweeps, we did not observe a significant enrichment of sweeps overlapping either global
427 or population-specific outliers when these outliers were treated as a class of genomic elements.
428 This is perhaps surprising, given that selective sweeps are known to create characteristic elevations
429 of LD (*Kim and Nielsen, 2004*), and perhaps could mimic regions with very divergent levels of
430 recombination in a population-specific way. A number of other evolutionary forces might explain
431 the existence of our outlier regions as well. For example, mutation rate heterogeneity along
432 the chromosomes could, in principle, generate spurious peaks or troughs in our estimates of
433 recombination rate, as ReLERNN in effect scales its per-base recombination rate estimates by
434 a mutation rate that is assumed to be constant along the chromosome (*Figure S4, Figure S3*).
435 Moreover, introgression from diverged populations might affect patterns of allelic association in a
436 local way along the genome (*Schrider et al., 2018; Schumer et al., 2018*). Taken together, our results
437 suggest that while both inversions and selection can influence population-specific differences in the
438 landscape of recombination, the preponderance of these differences likely have complex causes.

439 In this report we described ReLERNN, a novel deep learning method for inferring fine-scale rates
440 of recombination across the genome. While ReLERNN currently stands as a functional end-to-end
441 pipeline for measuring recombination rates, the modular design herein presents a number of
442 important opportunities for extension, with the potential to address myriad questions in population
443 genomics. For example, while ReLERNN is currently designed to use phased or unphased genotypes

444 from sequenced individuals as input, we see no reason why allele counts from pool-seq experiments
445 couldn't be substituted. Moreover, the RNN structure we exploit here could be used for inference
446 of the distribution of selection coefficients and/or migration rates from natural populations. In
447 addition, ReLERNN presents an excellent opportunity for the implementation of transfer learning,
448 whereby ReLERNN could be trained in-house on an otherwise prohibitively extensive parameter
449 space, allowing end-users to make accurate predictions by generating only a small fraction of
450 the current number of simulations and training epochs presently required. The application of
451 machine learning, and deep learning in particular, to questions in population genomics is ripe
452 with opportunity. ReLERNN provides a platform for jumping off, that we hope to see advance our
453 understanding of both population genetics and adaptation itself.

454 **Materials and Methods**

455 **The ReLERNN workflow**

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

468 **Parameter estimation and coalescent simulation**

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

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

493 rates of recombination from different files (e.g. one that includes only intergenic sites). When each
494 simulation is completed, ReLERNN dumps both the genotype matrix and a vector of the positions
495 for every SNP into a temporary .npy file.

496 Sequence batch generation and network architecture

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

508 ReLERNN trains a recurrent neural network with Keras (*Chollet et al., 2015*) using a Tensorflow
509 backend (*Abadi et al., 2015*). The complete details of our neural architecture can be found in the
510 python module "ReLERNN_networks.py", and a simplified flow diagram showing the connectivity
511 between layers can be found in *Figure 1*. Briefly, the ReLERNN neural network utilizes distinct input
512 layers for the genotype and position tensors, which are later merged using a concatenation layer
513 in Keras. The genotype tensor is first fed to a GRU layer, as implemented with the bidirectional
514 wrapper in Keras, and the output of this layer is passed to a dense layer followed by a dropout
515 layer. On the positions side of the network, the input positions tensor is fed directly to a dense
516 layer and then to a dropout layer. Dropout was used extensively in our network, as hypertuning
517 trials (below) demonstrated significantly improved accuracy when employing dropout relative to
518 networks without dropout. Once concatenated, output from the dropout layer is passed to a final
519 round of dense and dropout layers, and the final dense layer returns a single z-score normalized
520 prediction for each example, which is unnormalized back to units of crossovers per-base. ReLERNN
521 completes 250 training epochs and implements this training using the "Adam" optimizer and a
522 Mean Squared Error (*MSE*) loss function. Though the number of epochs is user-selectable, the
523 vast majority of networks are sufficiently trained within 250 epochs, largely due to how ReLERNN
524 handles the input tensor size and simulation parameters. Our hyper-tuning trials were completed
525 via a grid search over the set of parameters: Recurrent layer output dimensions (64, 82, 128), Loss
526 function (*MSE*, *MAE*), Input merge strategy (concatenate, average), and dense layer dimensionality
527 (64, 128), optimizing for *MSE*.

528 Parametric bootstrap analysis and prediction corrections

529 ReLERNN includes the option to both generate confidence intervals around each predicted re-
530 combination rate and to correct for potential biases generated during training using a parametric
531 bootstrapping approach. After the network has been trained and predictions have been gener-
532 ated, users can run ReLERNN_BSCORRECT, which resimulates 10^3 test examples for each of 100
533 recombination rate bins drawn from the distribution of recombination rates used to simulate
534 the original training set. Predictions are then generated for these 10^5 simulated test examples
535 using the previously trained network, generating a distribution of predictions for each respective
536 recombination rate bin. 95% confidence intervals are calculated from by taking the upper and lower
537 2.5% rate predictions from this distributions.

538 The distribution of test predictions can be biased in systematic ways, such as predictably under-
539 estimating rates of recombination for those examples with the highest simulated crossover events
540 (*Figure S1*). These biases may potentially be caused an inability to resolve very high recombination
541 rates with a limited number of informative SNPs. ReLERNN_BSCORRECT, estimates the magnitude of

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

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

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

566 **Comparative methods**

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

580 For both ReLERNN and the CNN, the same training set consisting of 10^5 examples was used
581 to train each neural network, and the same test examples were used to compare the predictions
582 produced by each method. Comparisons with LDhat and LDhelmet were made using the above
583 training examples to parameterize the generation of independent coalescent likelihood lookup
584 tables. For each set of examples of sample size N , we calculated the maximum value of ρ from
585 the training set and the average per-base values for θ for the test examples, using Watterson's
586 estimator. These parameter values were given to the functions for the lookup table generation
587 in LDhat and LDhelmet, and the resulting tables was used to make predictions on our 5×10^3 test
588 examples using the *pairwise* function. Comparisons with FastEPRR were made by transforming the
589 genotype matrices resulting from our test simulations into fasta-formatted input files, and running
590 the FastEPRR_ALN function (using format = 1) in R. As LDhat, LDhelmet, and FastEPRR all predict ρ ,

591 the resulting predictions were transformed to per-base recombination rates for comparison with
592 ReLERNN using the function $r = \frac{\rho_{pred} \times \mu_{true}}{\theta_W}$, whereby ρ_{pred} is the prediction output by each method, and
593 θ_W and μ_{true} are Watterson's estimator and the true per-base mutation rate used in the simulation
594 example, respectively. To compare accuracy among methods we directly compared the distribution
595 of absolute errors ($|r_{predicted} - r_{true}|$) for each method for each set of examples of sample size N .

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

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

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

626 Before running ReLERNN, we first set out to model the demographic history for each population
627 using each of three methods: stairwayplot (Liu and Fu, 2015), SMC++ (Terhorst et al., 2016), and
628 MSMC (Schiffels and Durbin, 2014). With the exception of MSMC, all methods were run using default
629 parameters. For MSMC, the use of default parameters generated predictions that were unusable
630 (Figure S9). For these reasons, and after direct communication with MSMC's authors, we determined
631 that running MSMC with a sample size of two chromosomes would be the most appropriate.
632 Ultimately we decided to run our ReLERNN pipeline with simulations generated under demographic
633 equilibrium [options: `-estimateDemography False -assumedMu 3.27e-9 -upperRhoThetaRatio`
634 `35`], as estimates of historical population size were unreliable for these data—all three methods
635 produced significantly different demographic histories (Figure S8)—and tests on simulated data
636 suggest little effect of demographic model misspecification (Figure S6). All code required to run our
637 ReLERNN analysis is deposited on GitHub (<https://github.com/kern-lab/ReLERNN>).

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.

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/](https://www.biorxiv.org/content/early/2018/11/05/267211)
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 **Comeron 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](https://doi.org/10.1371/journal.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](https://doi.org/10.1371/journal.pgen.1003056).
- 713 **Dobzhansky 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 **Fligel 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.](https://dx.doi.org/10.1093/molbev/msy224)
720 [1093/molbev/msy224](https://dx.doi.org/10.1093/molbev/msy224), doi: 10.1093/molbev/msy224.
- 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 **Haanel 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/
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 **Lucchesi JC**, Suzuki DT. The interchromosomal control of recombination. *Annual review of genetics*. 1968;
796 2(1):53–86.
- 797 **Mather K**. Crossing-over. *Biological Reviews*. 1938; 13(3):252–292.
- 798 **McVean G**, Awadalla P, Fearnhead P. A coalescent-based method for detecting and estimating recombination
799 from gene sequences. *Genetics*. 2002; 160(3):1231–1241.
- 800 **Miller DE**, Cook KR, Arvanitakis AV, Hawley RS. Third Chromosome Balancer Inversions Disrupt Protein-Coding
801 Genes and Influence Distal Recombination Events in *Drosophila melanogaster*. *G3: Genes, Genomes, Genetics*.
802 2016; 6(7):1959–1967. <https://www.g3journal.org/content/6/7/1959>, doi: 10.1534/g3.116.029330.
- 803 **Myers S**, Bottolo L, Freeman C, McVean G, Donnelly P. A fine-scale map of recombination rates and hotspots
804 across the human genome. *Science*. 2005; 310(5746):321–324.
- 805 **Myers SR**, Griffiths RC. Bounds on the minimum number of recombination events in a sample history. *Genetics*.
806 2003; 163(1):375–394.
- 807 **Noor MA**, Grams KL, Bertucci LA, Reiland J. Chromosomal inversions and the reproductive isolation of species.
808 *Proceedings of the National Academy of Sciences*. 2001; 98(21):12084–12088.
- 809 **Novitski E**, Braver G. An analysis of crossing over within a heterozygous inversion in *Drosophila melanogaster*.
810 *Genetics*. 1954; 39(2):197.
- 811 **O'Reilly PF**, Birney E, Balding DJ. Confounding between recombination and selection, and the Ped/Pop method
812 for detecting selection. *Genome research*. 2008; 18(8):1304–1313.
- 813 **Parsch J**, Meiklejohn CD, Hartl DL. Patterns of DNA sequence variation suggest the recent action of positive
814 selection in the janus-ocnus region of *Drosophila simulans*. *Genetics*. 2001; 159(2):647–657.
- 815 **Pool JE**, Corbett-Detig RB, Sugino RP, Stevens KA, Cardeno CM, Crepeau MW, Duchon P, Emerson JJ, Saelao P,
816 Begun DJ, Langley CH. Population Genomics of Sub-Saharan *Drosophila melanogaster*: African Diversity and
817 Non-African Admixture. *PLOS Genetics*. 2012 12; 8(12):1–24. <https://doi.org/10.1371/journal.pgen.1003080>,
818 doi: 10.1371/journal.pgen.1003080.
- 819 **Presgraves DC**, Gérard PR, Cherukuri A, Lyttle TW. Large-scale selective sweep among segregation distorter
820 chromosomes in African populations of *Drosophila melanogaster*. *PLoS genetics*. 2009; 5(5):e1000463.
- 821 **Price AL**, Tandon A, Patterson N, Barnes KC, Rafaels N, Ruczinski I, Beaty TH, Mathias R, Reich D, Myers S.
822 Sensitive detection of chromosomal segments of distinct ancestry in admixed populations. *PLoS genetics*.
823 2009; 5(6):e1000519.
- 824 **Przeworski M**, Wall JD. Why is there so little intragenic linkage disequilibrium in humans? *Genetics Research*.
825 2001; 77(2):143–151.

- 826 **R Core Team.** R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing,
827 Vienna, Austria; 2018. <https://www.R-project.org>.
- 828 **Rieseberg LH.** Chromosomal rearrangements and speciation. *Trends in ecology & evolution.* 2001; 16(7):351–
829 358.
- 830 **Ritz KR, Noor MA, Singh ND.** Variation in recombination rate: adaptive or not? *Trends in Genetics.* 2017;
831 33(5):364–374.
- 832 **Rogers AR.** How population growth affects linkage disequilibrium. *Genetics.* 2014; 197(4):1329–1341.
- 833 **Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg AC,
834 Fei-Fei L.** ImageNet Large Scale Visual Recognition Challenge. *Int J Comput Vision.* 2015 Dec; 115(3):211–252.
835 <http://dx.doi.org/10.1007/s11263-015-0816-y>, doi: 10.1007/s11263-015-0816-y.
- 836 **Schiffels S, Durbin R.** Inferring human population size and separation history from multiple genome sequences.
837 *Nature Genetics.* 2014 06; 46:919 EP -. <https://doi.org/10.1038/ng.3015>.
- 838 **Schrider DR, Ayroles J, Matute DR, Kern AD.** Supervised machine learning reveals introgressed loci in the
839 genomes of *Drosophila simulans* and *D. sechellia*. *PLoS genetics.* 2018; 14(4):e1007341.
- 840 **Schrider DR, Kern AD.** Supervised Machine Learning for Population Genetics: A New Paradigm. *Trends in
841 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.
- 842 **Schrider DR, Mendes FK, Hahn MW, Kern AD.** Soft shoulders ahead: spurious signatures of soft and partial
843 selective sweeps result from linked hard sweeps. *Genetics.* 2015; 200(1):267–284.
- 844 **Schumer M, Xu C, Powell DL, Durvasula A, Skov L, Holland C, Blazier JC, Sankararaman S, Andolfatto P, Rosen-
845 thal GG, Przeworski M.** Natural selection interacts with recombination to shape the evolution of hybrid
846 genomes. *Science.* 2018; 360(6389):656–660. <https://science.sciencemag.org/content/360/6389/656>, doi:
847 10.1126/science.aar3684.
- 848 **Singh ND, Stone EA, Aquadro CF, Clark AG.** Fine-scale heterogeneity in crossover rate in the garnet-scalloped
849 region of the *Drosophila melanogaster* X chromosome. *Genetics.* 2013; 194(2):375–387.
- 850 **Slatkin M.** Linkage disequilibrium in growing and stable populations. *Genetics.* 1994; 137(1):331–336.
- 851 **Smith KN, Nicolas A.** Recombination at work for meiosis. *Current opinion in genetics & development.* 1998;
852 8(2):200–211.
- 853 **Sturtevant A.** A case of rearrangement of genes in *Drosophila*. *Proceedings of the National Academy of
854 Sciences of the United States of America.* 1921; 7(8):235.
- 855 **Sutskever I, Vinyals O, Le QV.** Sequence to Sequence Learning with Neural Networks. In: *Proceedings of the 27th
856 International Conference on Neural Information Processing Systems - Volume 2 NIPS'14*, Cambridge, MA, USA: MIT
857 Press; 2014. p. 3104–3112. <http://dl.acm.org/citation.cfm?id=2969033.2969173>.
- 858 **Szegedy C, Liu W, Jia Y, Sermanet P, Reed SE, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A.** Going deeper
859 with convolutions. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA,
860 June 7-12, 2015*; 2015. p. 1–9. <https://doi.org/10.1109/CVPR.2015.7298594>, doi: 10.1109/CVPR.2015.7298594.
- 861 **Tennessen JA, Bigham AW, O'Connor TD, Fu W, Kenny EE, Gravel S, McGee S, Do R, Liu X, Jun G, Kang HM,
862 Jordan D, Leal SM, Gabriel S, Rieder MJ, Abecasis G, Altshuler D, Nickerson DA, Boerwinkle E, Sunyaev S, et al.**
863 Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes. *Science.*
864 2012; 337(6090):64–69. <https://science.sciencemag.org/content/337/6090/64>, doi: 10.1126/science.1219240.
- 865 **Terhorst J, Kamm JA, Song YS.** Robust and scalable inference of population history from hundreds of unphased
866 whole genomes. *Nature Genetics.* 2016 12; 49:303 EP -. <https://doi.org/10.1038/ng.3748>.
- 867 **Vincent P, Larochelle H, Bengio Y, Manzagol PA.** Extracting and Composing Robust Features with Denoising
868 Autoencoders. In: *Proceedings of the 25th International Conference on Machine Learning ICML '08*, New York, NY,
869 USA: ACM; 2008. p. 1096–1103. <http://doi.acm.org/10.1145/1390156.1390294>, doi: 10.1145/1390156.1390294.
- 870 **Wakeley J.** Using the variance of pairwise differences to estimate the recombination rate. *Genetics Research.*
871 1997; 69(1):45–48.
- 872 **Wall JD.** A comparison of estimators of the population recombination rate. *Molecular Biology and Evolution.*
873 2000; 17(1):156–163.

- 874 **Wang RJ**, Gray MM, Parmenter MD, Broman KW, Payseur BA. Recombination rate variation in mice from an
875 isolated island. *Molecular ecology*. 2017; 26(2):457–470.
- 876 **Winckler W**, Myers SR, Richter DJ, Onofrio RC, McDonald GJ, Bontrop RE, McVean GA, Gabriel SB, Reich D,
877 Donnelly P, et al. Comparison of fine-scale recombination rates in humans and chimpanzees. *Science*. 2005;
878 308(5718):107–111.
- 879 **Wiuf C**. On the minimum number of topologies explaining a sample of DNA sequences. *Theoretical population*
880 *biology*. 2002; 62(4):357–363.
- 881 **Yi X**, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, Zheng H,
882 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
883 adaptation to high altitude. *Science (New York, NY)*. 2010 07; 329(5987):75–78. [https://www.ncbi.nlm.nih.gov/
884 pubmed/20595611](https://www.ncbi.nlm.nih.gov/pubmed/20595611), doi: 10.1126/science.1190371.
- 885 **Zelkowski M**, Olson M, Wang M, P Pawlowski W. Diversity and Determinants of Meiotic Recombination
886 Landscapes. *Trends in Genetics*. 2019 04; doi: 10.1016/j.tig.2019.02.002.

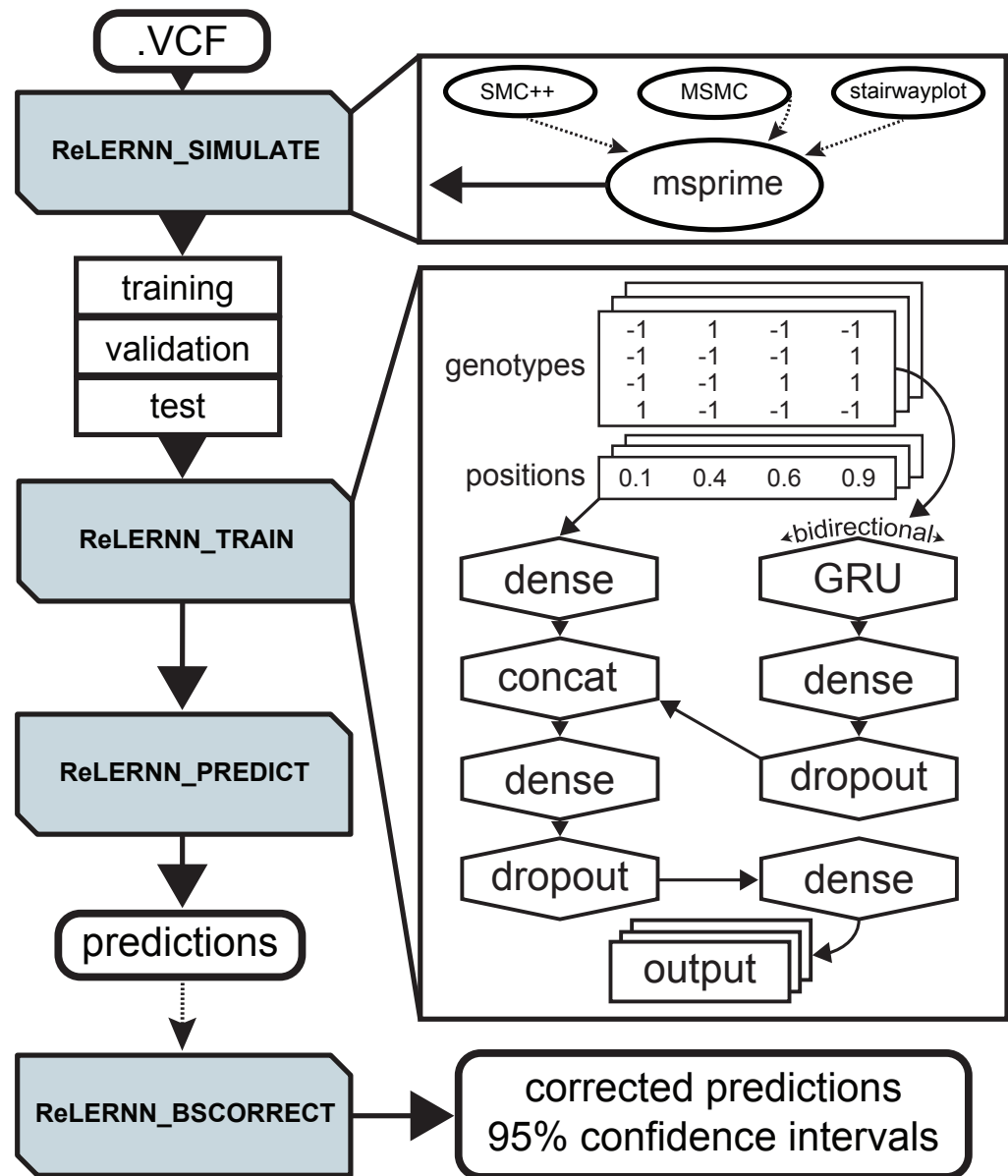


Figure 1 Diagram depicting a typical workflow using ReLERNN's four modules (shaded boxes). ReLERNN_SIMULATE can optionally (dotted lines) utilize output from stairwayplot, SMC++, MSMC to simulate under a demographic history in msprime. The breakout of ReLERNN_TRAIN depicts the GRU network architecture 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).

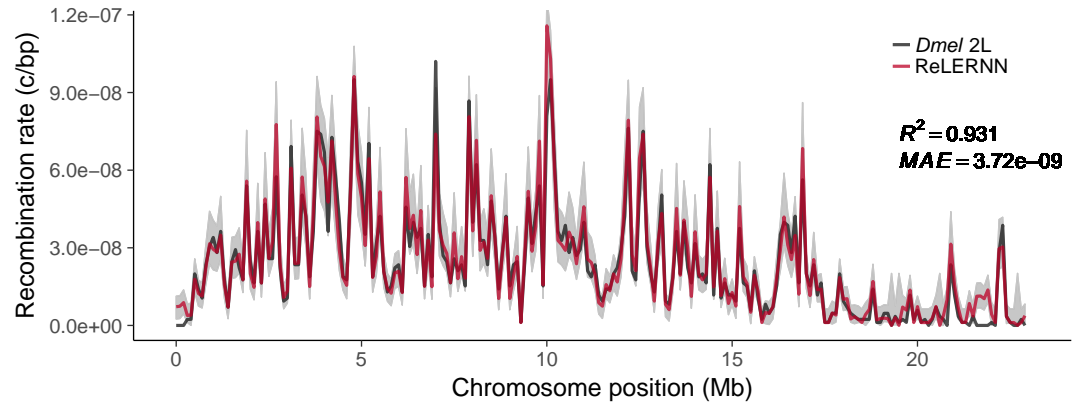


Figure 2 Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *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.

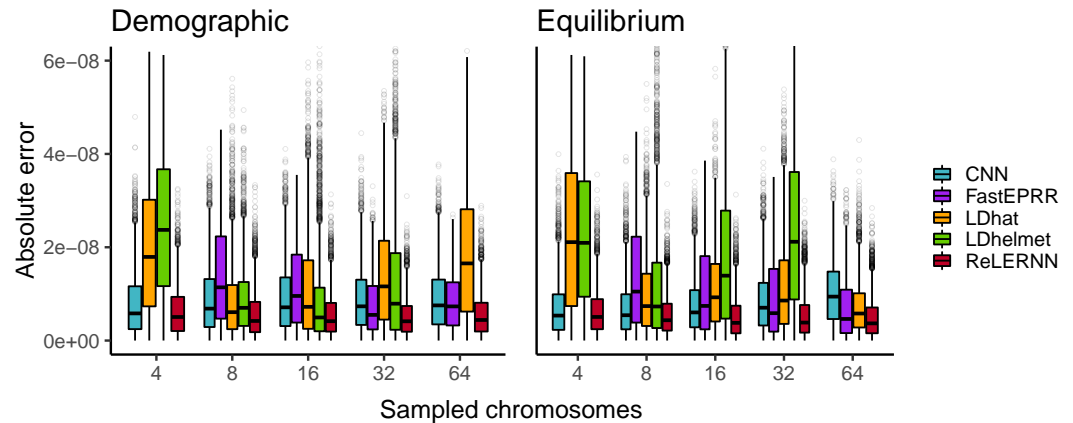


Figure 3 Distribution of absolute errors ($|r_{\text{predicted}} - r_{\text{true}}|$) for each method across 5000 simulated chromosomes (1000 for FastEPRR). Independent simulations were run under a known demographic history (left) or an assumption of demographic equilibrium (right). Sampled chromosomes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

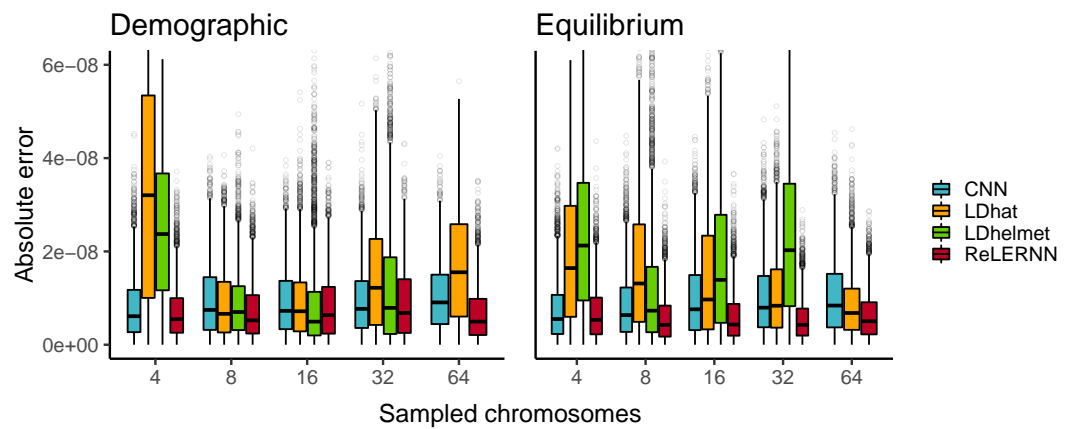


Figure 4 Distribution of absolute errors ($|r_{\text{predicted}} - r_{\text{true}}|$) for each method across 5000 simulated chromosomes after model misspecification. For the CNN and ReLERNN, predictions were made by training on equilibrium simulations and testing on sequences simulated under a demographic model (left) or training on demographic simulations and testing on sequences simulated under equilibrium (right). For LDhat and LDhelmet, the lookup tables were 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.

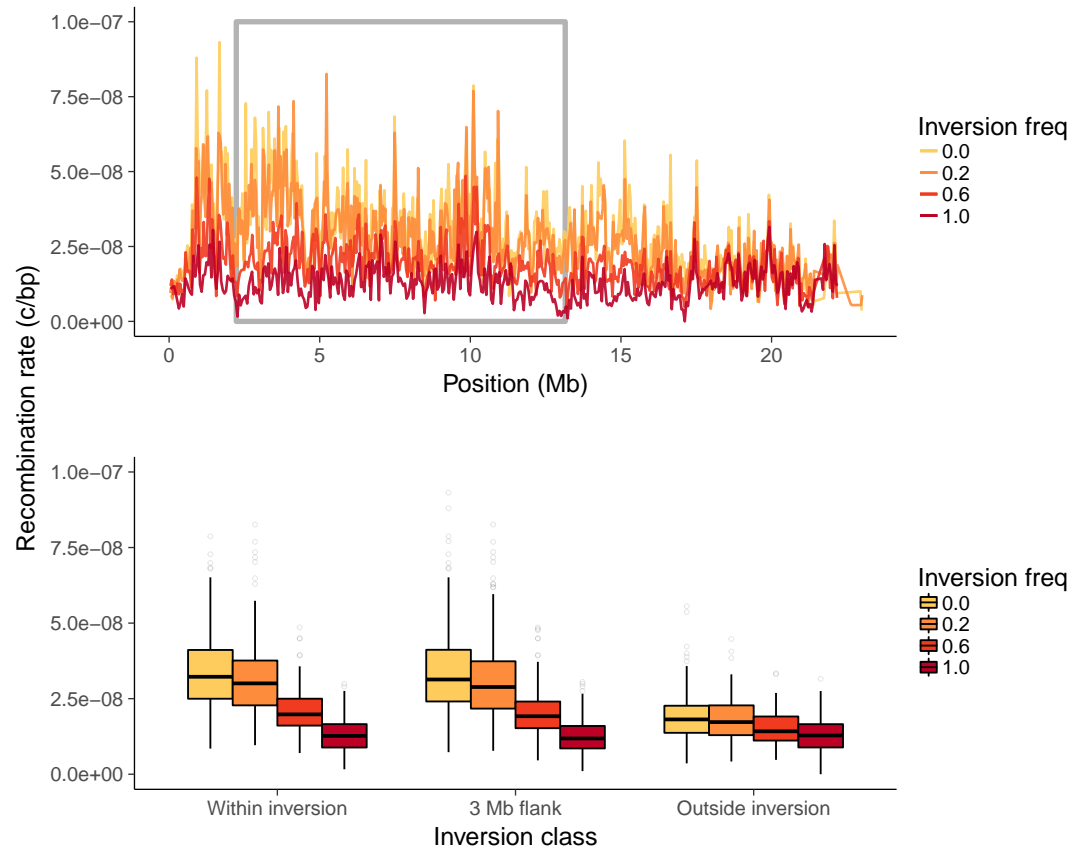


Figure 6 (Top) Recombination landscapes for Zambian *D. melanogaster* surrounding *In(2L)t*, sampled at different inversion frequencies. The grey box denotes the inversion boundaries of *In(2L)t* in *Drosophila* (Corbett-Detig and Hartl, 2012). (Bottom) Recombination rate estimates from genomic windows within the inversion, within a 3 Mb region flanking the inversion, and 3 Mb outside the inversion, sampled at different inversion frequencies.

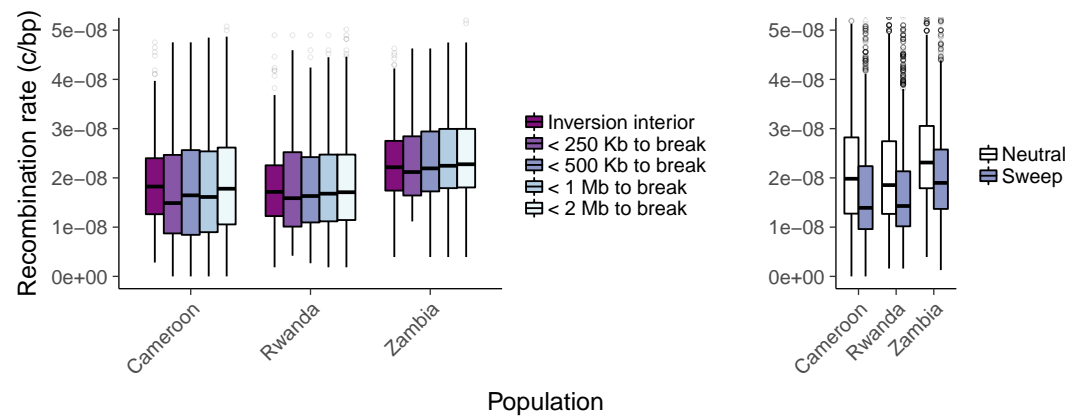


Figure 7 (Left) Recombination rate estimates for genomic windows > 2 Mb inside, < 250 kb surrounding, < 500 kb surrounding, < 1 Mb surrounding, and < 2 Mb surrounding all inversion breakpoints. (Right) Recombination rate estimates for all genomic windows overlapping windows predicted as either hard/soft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

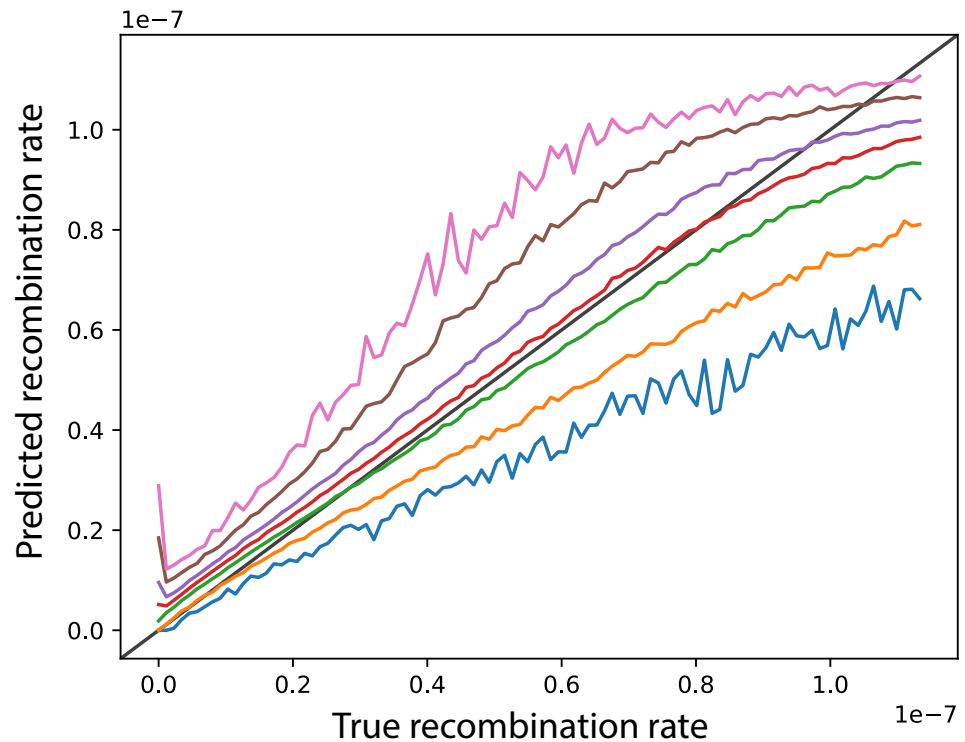


Figure S1 Parametric bootstrapping results as implemented by ReLERNN. Lines represent the minimum (blue), lower 5% (orange), lower 25% (green), median (red), upper 25% (purple), upper 95% (brown), and maximum (pink) bounds for each of 1000 replicate simulations and predictions (y-axis) across 100 recombination rate bins (x-axis)

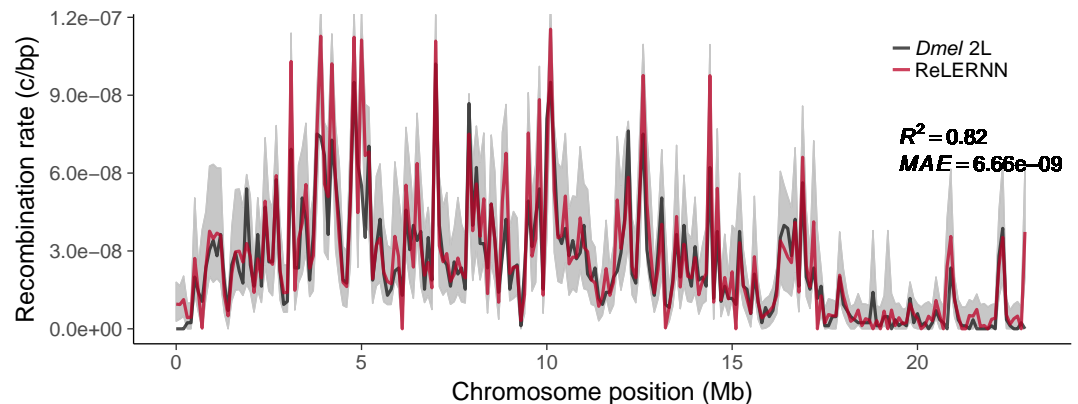


Figure S2 Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 4$ chromosomes under mutation-drift equilibrium using msprime (Kelleher *et al.*, 2016), with per-base crossover rates derived from *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.

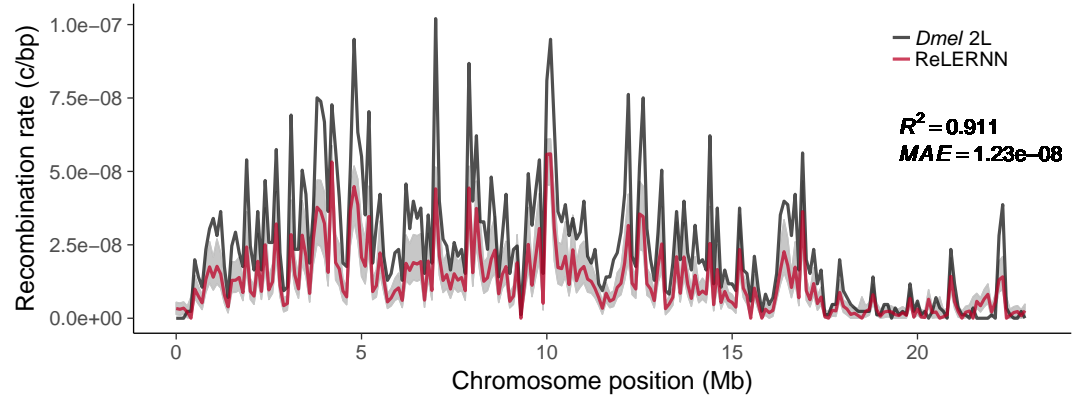


Figure S3 Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Cameron 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.

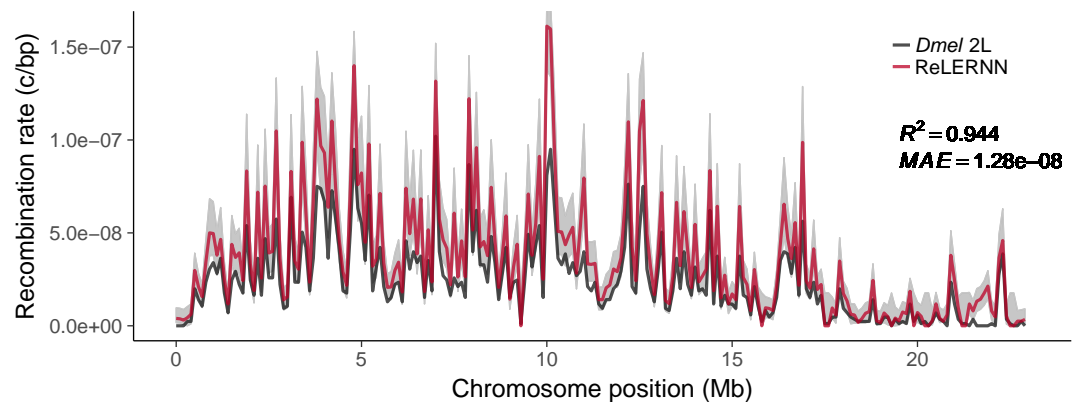


Figure S4 Recombination rate predictions for a simulated *Drosophila* chromosome (black line) using ReLERNN (red line). The recombination landscape was simulated for $n = 20$ chromosomes under mutation-drift equilibrium using msprime (Kelleher et al., 2016), with per-base crossover rates derived from *D. melanogaster* chromosome 2L (Cameron 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.

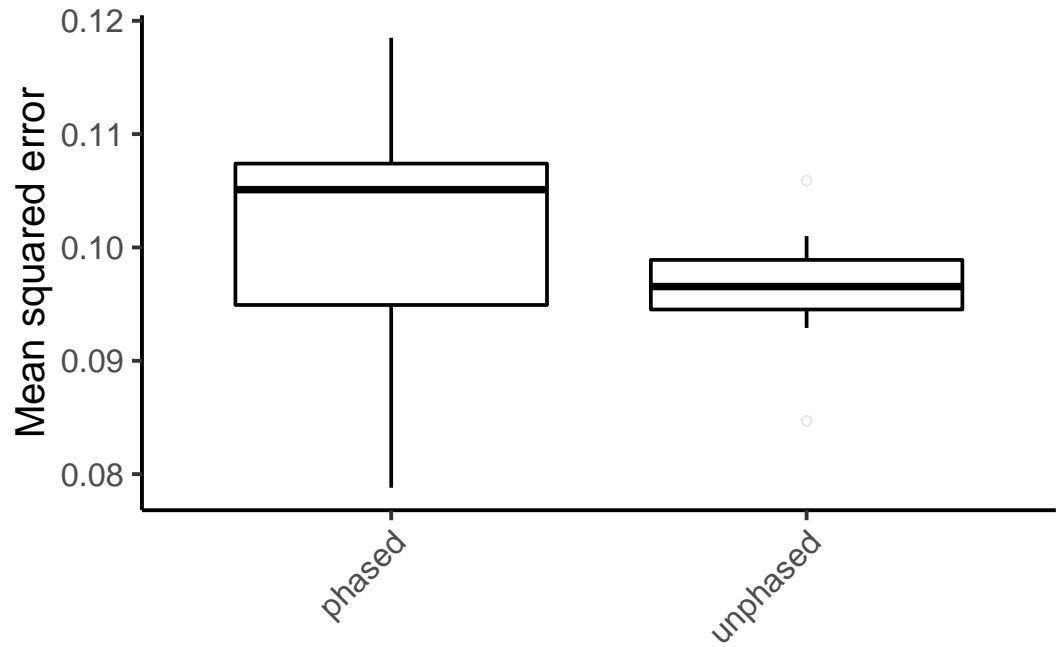


Figure S5 Mean squared error for ReLERNN predictions on 10 replicates of 1000 test simulations using 100% correctly phased input genotypes and completely unphased genotypes. All simulations used the recombination map derived from *D. melanogaster* chromosome 2L (Comeron *et al.*, 2012).

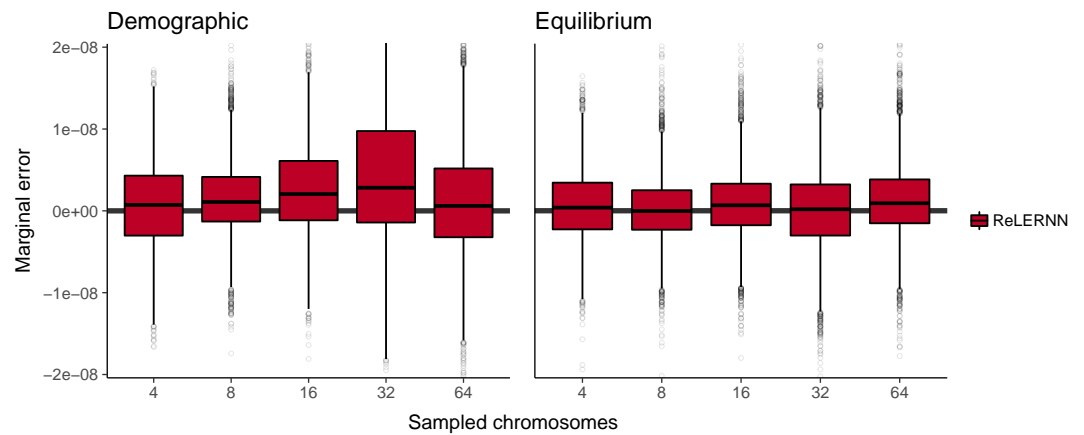


Figure S6 Distribution of marginal errors attributed to model misspecification across 5000 simulated chromosomes. Predictions were made by training on equilibrium simulations and testing on sequences simulated under a demographic model (left) or training on demographic simulations and testing on sequences 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.

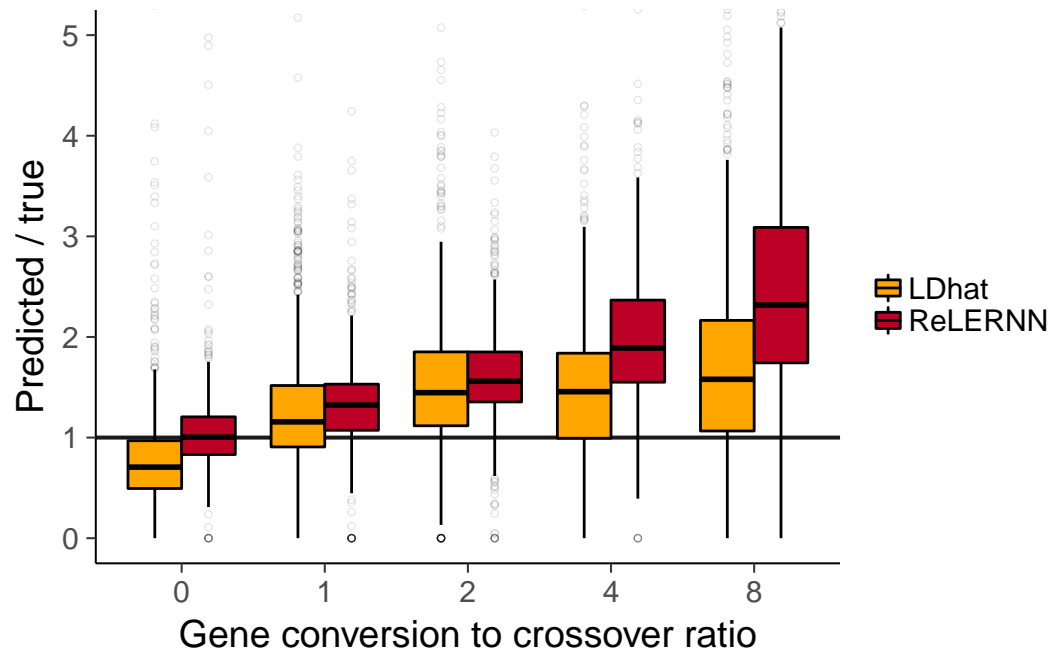


Figure S7 Distribution of predicted rates of recombination over true rates for 5000 examples simulated with gene conversion and $n = 8$. The ratio of gene conversion to crossovers was drawn from $U(0, c)$, with $c \in \{0, 1, 2, 4, 8\}$. Gene conversion tract lengths were fixed at 352 bp, and all simulations were completed in ms (Hudson, 2002).

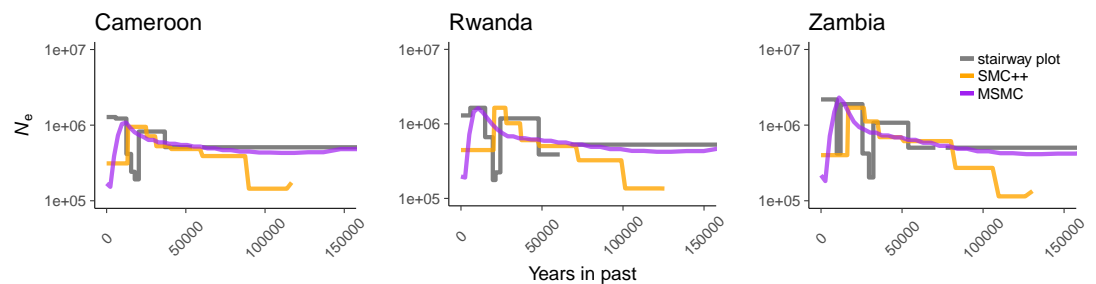


Figure S8 Historical population size estimates were inferred for Cameroon, Rwanda, and Zambia using three separate methods, all of which disagree with one another. Inferences are based on 10 samples for both stairwayplot (grey line) and SMC++ (orange line), and 2 samples for MSMC (purple line).

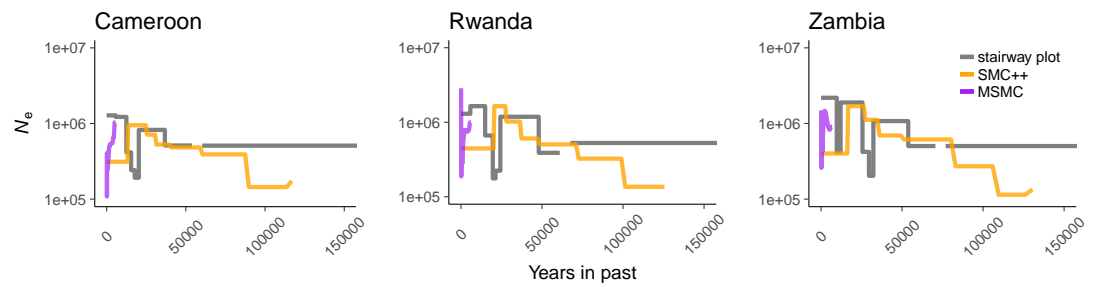


Figure S9 Historical population size estimates were inferred for Cameroon, Rwanda, and Zambia using three separate methods, all of which disagree with one another. Here, inferences are based on 10 samples for both stairwayplot (grey line) and SMC++ (orange line), and 10 samples for MSMC (purple line).

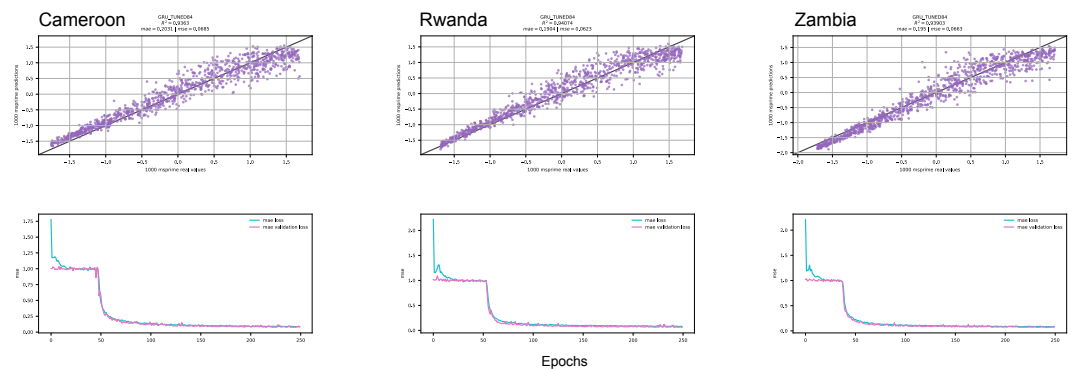


Figure S10 ReLERNN test results for Cameroon, Rwanda, and Zambia when trained under assumptions of mutation-drift equilibrium. Scatter plots (top) show raw (unnormalized) predictions for per-base recombination rates for 1000 test examples. Mean absolute error and mean squared error are calculated for 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).

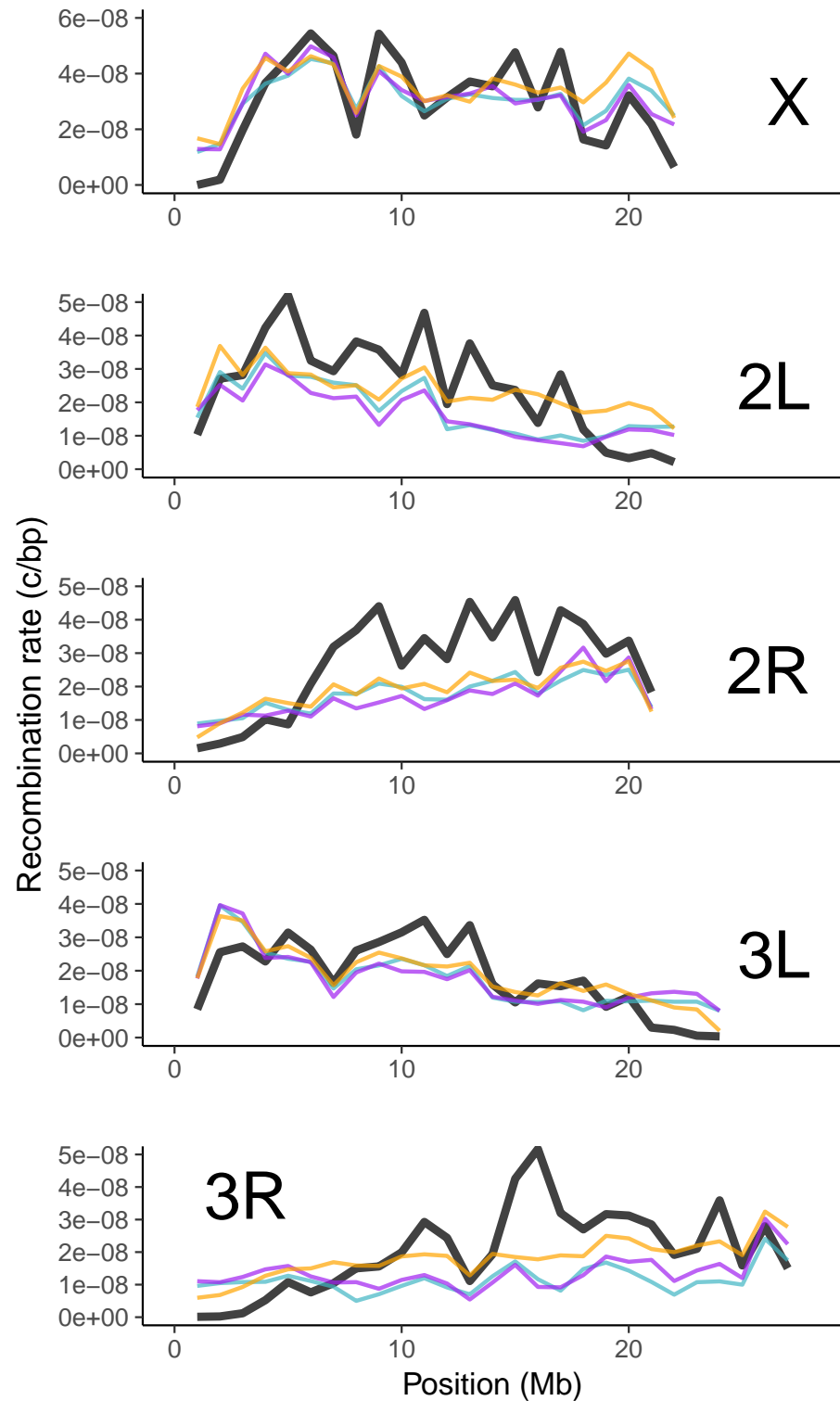


Figure S11 Genome-wide recombination landscapes for *D. melanogaster* populations from Cameroon (teal lines), Rwanda (purple lines), and Zambia (orange lines). Rates are compared to those experimentally derived by *Cameron et al. (2012)* (black lines). All rates have been scaled to 1 Mb windows by using a weighted average (see Materials and Methods).

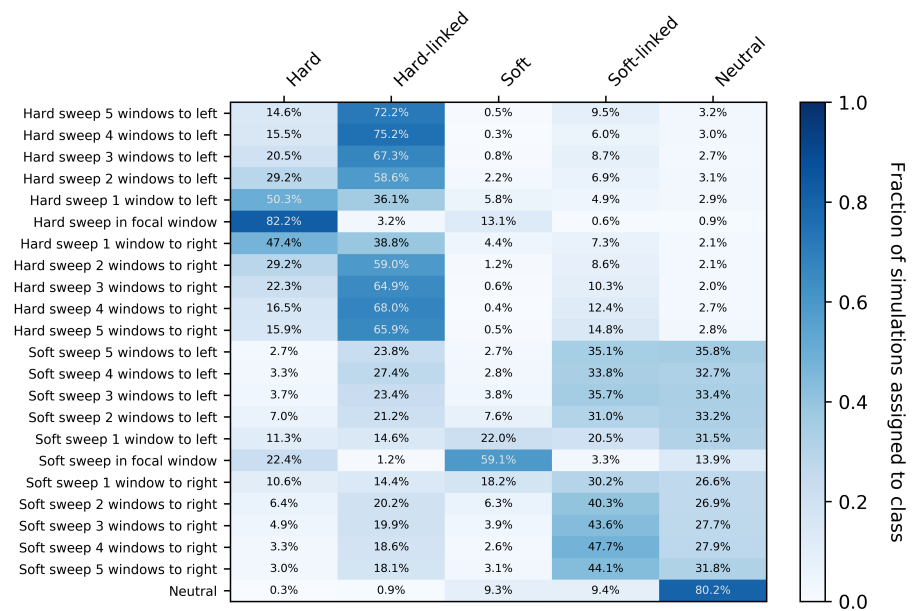


Figure S12 Confusion matrix showing the fraction of test simulation windows assigned to each of five prediction categories by diploS/HIC (*Kern and Schrider, 2018*): hard, hard-linked, soft, soft-linked, and neutral. The y-axis shows the location of the window being classified relative to the selected window.

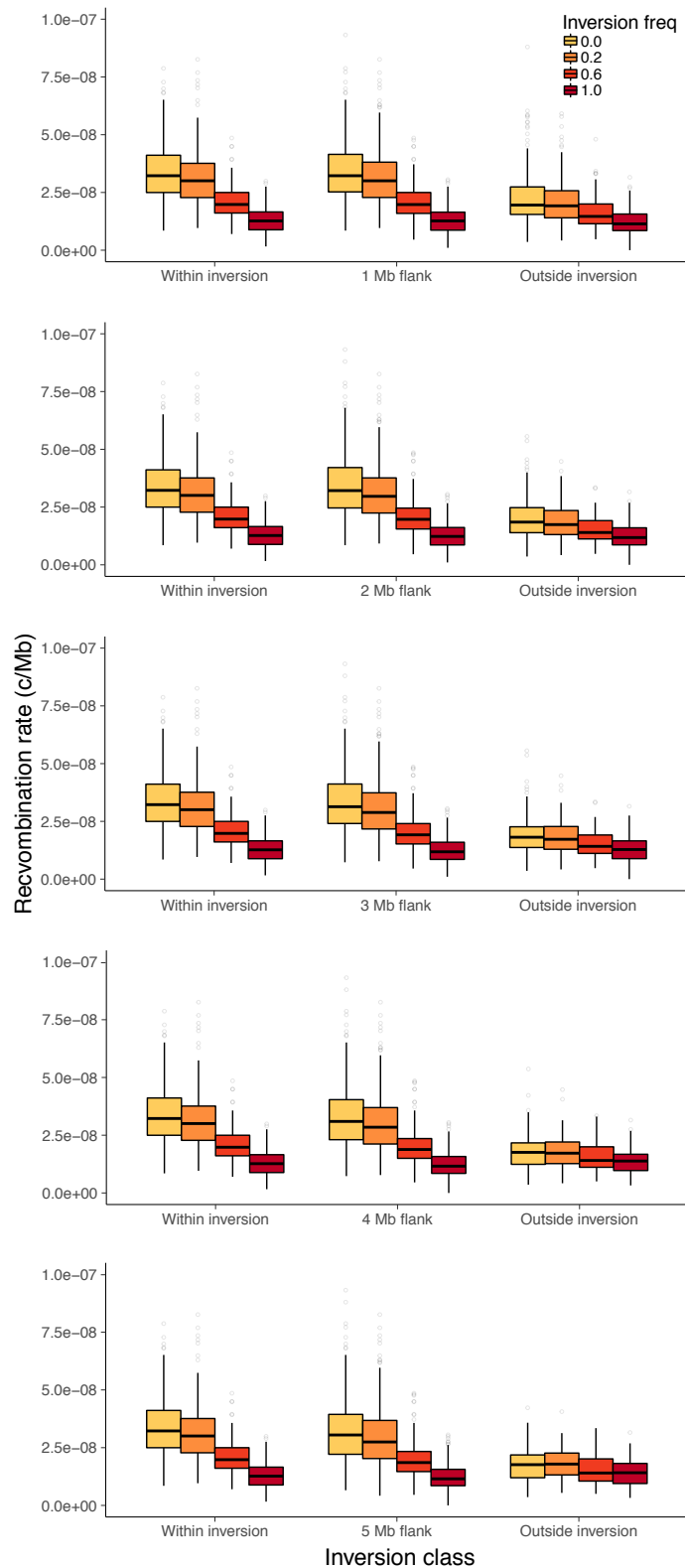


Figure S13 Recombination rate estimates using flanking window sizes from 1-5 Mb. Rates are shown for genomic windows within the inversion, within regions flanking the inversion, and for regions outside both the inversion and flanking regions. All estimates are from chromosome 2L with *In(2L)t* sampled at different inversion frequencies