1  **The Unreasonable Effectiveness of Convolutional Neural Networks in Population**

2  **Genetic Inference**

3

4  Lex Flagel[1,2], Yaniv Brandvain[2], and Daniel R. Schrider[3,*]

5

6  [1]Monsanto Company, Chesterfield, MO.

7  [2]Department of Plant and Microbial Biology, University of Minnesota, St. Paul, MN.

8  [3]Department of Genetics, University of North Carolina, Chapel Hill, NC.

9

10  *Corresponding author. Email: drs@unc.edu

11 **ABSTRACT**

12 Population-scale genomic datasets have given researchers incredible amounts of information

13 from which to infer evolutionary histories. Concomitant with this flood of data, theoretical and

14 methodological advances have sought to extract information from genomic sequences to infer

15 demographic events such as population size changes and gene flow among closely related

16 populations/species, construct recombination maps, and uncover loci underlying recent

17 adaptation. To date most methods make use of only one or a few summaries of the input

18 sequences and therefore ignore potentially useful information encoded in the data. The most

19 sophisticated of these approaches involve likelihood calculations, which require theoretical

20 advances for each new problem, and often focus on a single aspect of the data (e.g. only allele

21 frequency information) in the interest of mathematical and computational tractability. Directly

22 interrogating the entirety of the input sequence data in a likelihood-free manner would thus offer

23 a fruitful alternative. Here we accomplish this by representing DNA sequence alignments as

24 images and using a class of deep learning methods called convolutional neural networks (CNNs)

25 to make population genetic inferences from these images. We apply CNNs to a number of

26 evolutionary questions and find that they frequently match or exceed the accuracy of current

27 methods. Importantly, we show that CNNs perform accurate evolutionary model selection and

28 parameter estimation, even on problems that have not received detailed theoretical treatments.

29 Thus, when applied to population genetic alignments, CNN are capable of outperforming

30 expert-derived statistical methods, and offer a new path forward in cases where no likelihood

31 approach exists.

## 32  INTRODUCTION

33  Using genetic data to make inferences about the natural histories of populations represents a
34  major goal of evolutionary research. As the ever-increasing throughput of DNA sequencing
35  technologies makes the generation of large population genomic data sets more routine,
36  researchers can leverage patterns of genetic variation across the genome to characterize the
37  evolutionary forces at play (Hahn 2018). For example, advances have been made in identifying
38  historical demographic events such as population size changes (Marth *et al.* 2004; Tennessen *et al.*
39  2012; Gazave *et al.* 2014) and genetic exchange between populations and species (Martin *et al.*
40  2013; Hellenthal *et al.* 2014; Sankararaman *et al.* 2014; Corbett-Detig and Nielsen 2017; Schrider
41  *et al.* 2018). Population genomic analyses have also revealed the pervasive impact of selection on
42  linked neutral polymorphism (Begun and Aquadro 1992; Begun *et al.* 2007; Langley *et al.* 2012;
43  Elyashiv *et al.* 2016), both through positive selection (Maynard Smith and Haigh 1974; Kaplan *et*
44  *al.* 1989) and purifying selection (Charlesworth *et al.* 1993). As the volume of population genomic
45  data sets has increased, so too has the demand for powerful computational methods capable of
46  using these data to learn about the fundamental evolutionary processes shaping genomic
47  variation.

48      To meet this need, myriad statistical and computational tools have been devised to
49  answer evolutionary questions using population genetic data. One particularly common
50  paradigm, which predates the high-throughput sequencing revolution, is that of the population
51  genetic summary statistic: a value (or sometimes a vector of values) designed to capture the
52  information present in a sequence alignment of individuals from one or more populations. When
53  a particular evolutionary phenomenon acts on a population it alters the shapes of genealogies,
54  and this effect is manifest in the observed sequence alignment. For example, a population
55  expansion will result in genealogies with longer branches near the leaves of the tree, which will
56  manifest as an excess of rare alleles. Many summary statistics seek to uncover the signature of
57  these genealogical skews through their effect on the alignment; e.g. Tajima's $D$ will be negative
58  following a recent expansion or recovery from a bottleneck (Tajima 1989; Simonsen *et al.* 1995).
59  Ideally a summary statistic will only detect the signal of the evolutionary process it is being used
60  to investigate, but in practice summary statistics are frequently confounded by other forces that
61  may have similar effects on the shapes and/or sizes of genealogies. For example, Tajima's $D$ is
62  sensitive to positive selection as well as population size changes (Simonsen *et al.* 1995). Moreover,

3

63    such summary statistics do not capture all of the information present in the alignment. Thus a

64    major challenge of population genetic inference is to create methods that utilize as much

65    information from the input data as possible in order to maximize our ability to distinguish among

66    the numerous evolutionary processes that can give rise to an observed signal.

67    One approach researchers have adopted to address this challenge is to incorporate a

68    larger number of observations from the data into likelihood-based inference methods. However,

69    calculating likelihoods of population genomic data sets is often mathematically and

70    computationally intractable, and therefore such approaches often use composite likelihoods

71    which ignore the non-independence of observations (e.g. Hudson 2001; Nielsen *et al.* 2005). For

72    example, Nielsen et al.'s SweepFinder (2005), which examines allele frequencies at

73    polymorphisms flanking a focal region to determine whether that region has experienced a recent

74    selective sweep (Maynard Smith and Haigh 1974), treats each allele frequency as an independent

75    observation despite the partially shared evolutionary histories linked alleles experience. Another

76    drawback of most likelihood-based methods is that they generally compute the likelihood of only

77    a few features of the data (often only one), and therefore additional information that could

78    improve accuracy is ignored. For example, SweepFinder examines allele frequencies but ignores

79    linkage disequilibrium (LD), which is elevated in areas flanking the selected site (Kim and Nielsen

80    2004). Hidden Markov models (Hobolth *et al.* 2007; Boitard *et al.* 2009; Dutheil *et al.* 2009; Kern

81    and Haussler 2010), including those based on the sequential Markov coalescent (Li and Durbin

82    2011; Schiffels and Durbin 2014), have also proved effective at using population genetic

83    observations along a recombining chromosome to make evolutionary inferences.

84    More recently, population geneticists have begun to explore an alternative strategy of

85    using a large set of complementary summary statistics for model selection and parameter

86    estimation, an approach that often results in more powerful and robust inference (e.g. Lin *et al.*

87    2011; Pybus *et al.* 2015; Gao *et al.* 2016; Schrider and Kern 2016; Sheehan and Song 2016). Each

88    summary statistic seeks to measure a particular attribute of the genealogy, and one can thus

89    design a customized set of summary statistics to more fully represent the genealogical information

90    present in the sequence alignment. This view deploys summary statistics less for their individual

91    links to underlying theory, and more for their collective ability to perform pattern recognition.

92    The challenge then becomes extracting information about the underlying evolutionary processes

93    from the set of summary statistics. Two exciting approaches for dealing with this challenge that

94    have garnered increasing attention in recent years are approximate Bayesian computation (ABC;

95    reviewed in Beaumont 2010) and supervised machine learning (reviewed in Schrider and Kern

96    2018). Both of these approaches make use of suites of user-defined summary statistics and

97    training data generated under known parameters to identify reasonable evolutionary models and

98    parameterizations that could have generated the observed data. Here we focus on the supervised

99    machine learning approach, as it sets the scene for the convolutional neural networks described

100    below.

101          In the terminology of supervised machine learning, each summary statistic is called a

102    feature, and the full set of statistics used is called a feature vector. To use supervised machine

103    learning, a researcher must first obtain training data (often referred to as "labeled" data)—a set of

104    data points each summarized by a feature vector (the explanatory variables) accompanied by a

105    known outcome (the response variable). Next, a supervised machine learning algorithm is trained

106    to predict the outcome given the feature vector using the labeled training data. Thus, the

107    supervised machine learning technique automates the process of extracting information and

108    constructing rules from a set of summary statistics. Across many areas of research, supervised

109    machine learning techniques are fast replacing rules developed by human experts because they

110    are often more accurate (LeCun *et al.* 2015).

111          Supervised machine learning methods are increasingly being applied to numerous

112    problems in population genetics (Schrider and Kern 2018). In this context, labeled training data

113    are usually generated via population genetic simulation, an endeavor that has grown

114    considerably more feasible given recent improvements in simulation flexibility and efficiency (e.g.

115    Thornton 2014; Kelleher *et al.* 2016; Haller and Messer 2017; Kelleher *et al.* 2018). To date,

116    population genetic applications of machine learning include demographic inference (Pudlo *et al.*

117    2016; Sheehan and Song 2016), local ancestry inference (Schrider *et al.* 2018), inferring

118    recombination rates (Lin *et al.* 2013; Gao *et al.* 2016), and detecting genomic regions experiencing

119    recent selective sweeps (Pavlidis *et al.* 2010; Lin *et al.* 2011; Ronen *et al.* 2013; Pybus *et al.* 2015;

120    Schrider and Kern 2016). While such methods have great promise, they still rely on a user-

121    defined set of summary statistics (ranging in number from dozens to hundreds). Moreover, it is

122    not known whether it is possible to construct a set of statistics that sufficiently captures all

123    relevant information in the input data.

124        Unlike other machine learning approaches, convolutional neural networks (CNN; LeCun

125    *et al.* 1998) are pattern recognition algorithms that do not require a predefined feature vector.

126    When fed labeled training data (e.g. a set of haplotypes simulated under a known biological

127    scenario), a CNN discovers meaningful features, in essence making a feature vector, and then

128    extracts information from these features in order to make inferences. CNNs have proved effective

129    in a number of fields (reviewed in LeCun *et al.* 2015), and particularly in the field of image

130    recognition, where they have achieved dramatic improvements over previous efforts (e.g.

131    Lawrence *et al.* 1997; Krizhevsky *et al.* 2012; Simonyan and Zisserman 2014). The application of

132    CNNs to population genomic inference is just beginning, and shows great promise (Chan *et al.*

133    2018). Population genetic questions may be particularly well suited for CNN-based learning

134    because they take matrices as inputs, and alignments of sequenced chromosomes are quite

135    naturally represented in this manner.

136        The goal of this paper is to assess the effectiveness of CNNs as a general strategy for

137    population genomic inference. We demonstrate that CNNs can be successfully applied to a

138    number of population genomic problems, in some cases achieving surprising accuracy. In

139    particular, we use simulation to show that CNNs can leverage images of aligned sequences to

140    accurately uncover regions experiencing gene flow between related populations/species, estimate

141    recombination rates, detect selective sweeps, and make demographic inferences. Indeed, in most

142    cases we observe performance that matches or exceeds that of current methods. We also use a

143    CNN to accurately infer recombination rates from read coverage data in a simulated

144    autotetraploid, demonstrating this approach's flexibility in handling noisy data while solving a

145    complex problem for which no theoretical solution exists. In light of these encouraging findings,

146    we argue that population genetics researchers should consider CNNs as a potential solution to a

147    variety of problems involving evolutionary inferences from sequence data. Because some readers

148    may have little background with this tool, we also provide an overview of the inner workings of

149    CNNs and explore several technical considerations that may impact performance.

150

151    **RESULTS**

152    Our goal is to use a CNN to make population genetic inferences from an alignment image, which

153    can be thought of as matrices where each entry represents the allele present in a given

154    chromosome at a given site. In particular, we focus on four distinct problems: identifying local

155 introgression, estimating the recombination rate, detecting selective sweeps, and inferring
156 population size changes. We chose these four tasks because each represents a different challenge
157 in population genetic inference, each with its own attendant branch of theory. To show the
158 ability of CNNs to solve problems for which no statistical approaches have been proposed, we
159 extended our recombination inference to infer recombination rates in autotetraploids with
160 tetrasomic inheritance.

161     Below, we address each of these problems in turn, providing a brief overview of the
162 phenomenon in question and existing methodology before describing our results using CNNs.
163 But prior to tackling these problems, we first give an overview of CNNs and discuss strategies for
164 reorganizing our input data that we found helpful in making CNNs work more efficiently with
165 population genetic alignments.

166

167 **Overview of convolutional neural networks**

168 Internally, a CNN is a type of artificial neural network − a collection of connected layers of
169 combinatorially linked mathematical functions (termed *artificial neurons*) that take an input and
170 transform it into an output value (Mitchell 1997). In a typical fully connected artificial neural
171 network, the input values are fed through a series of layers of artificial neurons (fig. 1A), termed
172 hidden layers, before reaching the output layer which transforms its inputs into a final prediction.
173 The output for the $j$th neuron within one of the hidden layers is given by the following:

$$f\left(\sum_i^n w_{ij}x_i + b_j\right)$$

174 In the expression above, $x_i$ is the neuron's $i$th input value (either an input value from the data or
175 from a neuron in the previous layer's output), and $w_{ij}$ is the *weight* attached to the connection
176 between that node ($i$) and the current node ($j$) and $b_j$ is the current node's *bias* term. That is, to
177 obtain the value of neuron $j$, we compute the linear combination of the vector containing all
178 values from the previous layer and the $j$th neuron's vector of weights; the results of this summation
179 are in turn added to neuron $j$'s bias term and then fed as input to some function $f$, termed the
180 *activation function* and which may be nonlinear. Thus, an artificial neural network is a
181 mathematical function.

182     Importantly, by changing the values of the weights and biases, an artificial neural network
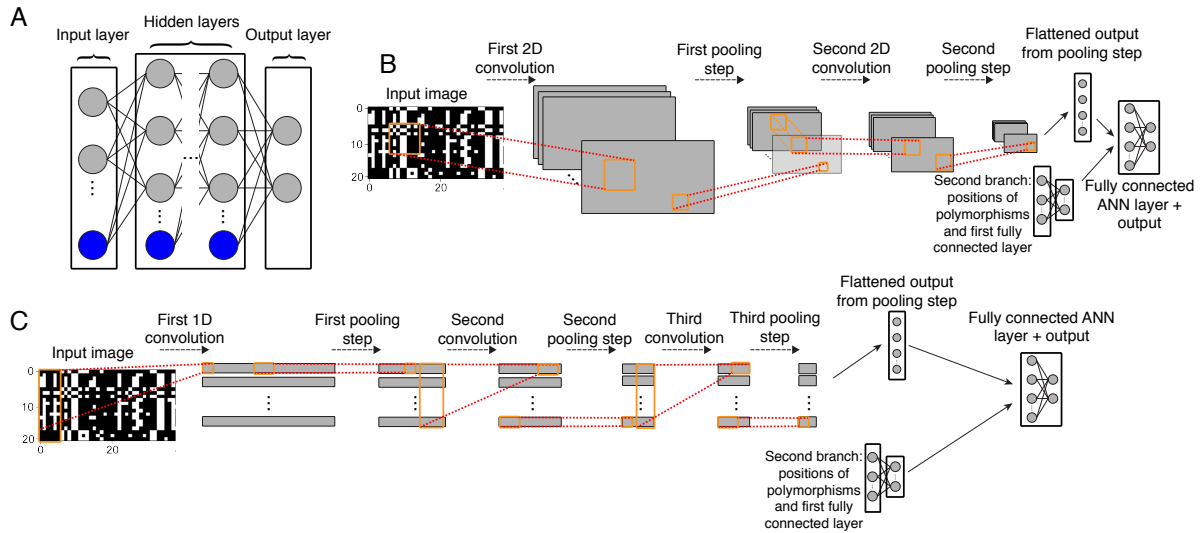183 can be tuned to detect informative patterns in the input data in order to produce the desired

7

**Fig. 1: Schematics of a standard feedforward neural network and two convolutional neural network designs used in this study.** A) Diagram of a fully connected feedforward neural network. Gray circles represent input (left side), output (right side), or hidden (center) neurons. Blue circles represent collections of bias terms. With the exception of the input layer, the value of any given neuron is a linear combination of values from the previous layer plus a bias term; this sum is then passed to an activation function (not shown). Each edge represents a distinct weighted input or bias term. Outputs may represent class membership posterior probabilities or estimates of continuous variables. B) A diagram of a 2D CNN similar to that used in this study to infer demographic parameters. The input is an alignment represented as an image which is passed through a first convolutional layer in order to create a set of feature maps. These feature maps are then downsized via a pooling step which replaces the values of a 1 or 2D matrix within a feature map with a single value summarizing it (e.g. the mean or maximum value of that matrix). For example, a 2D pooling operation of size 2 will reduce the size of a feature map by a factor of 4, as each adjacent 2×2 matrix within the input feature map is replaced by a single value (e.g. the maximum of those four values). These downsized feature maps are then passed through a second convolutional filter and pooling step, and the resulting output is flattened into a one dimensional vector and passed as input into a fully connected feedfoward layer (bias terms not shown). Also passed into this layer is output from a second branch of this network: the vector of positions of segregating sites in the alignment which have been passed through their own fully connected layer. Finally, the last fully connected neural network layer yields the predicted output values. C) Similar to panel B, but showing a 1D CNN with three convolutional layers (each followed by a pooling step), as used for our recombination rate estimator.

184
185     output. In the case of image recognition, an image is first represented numerically, typically as a

186     matrix of pixel intensities, and then transformed by the artificial neural network to produce an

187     output, for example a prediction of the type of object in the image. CNNs (Fig. 1B–C) differ from

8

188    standard artificial neural networks in that they begin with one or more convolutional layers, in
189    which a series of smaller weight matrices referred to as "filters" slide across the input image—
190    mimicking the manner in which animal cortical neurons each focus on input only from a small
191    receptive field—and perform a matrix convolution at each step until a series of filtered image
192    matrices are produced (LeCun *et al.* 1998). These filters are constructed during training (see
193    below). Each convolutional layer is often followed by a pooling layer (see Fig 1B and caption)
194    which reduces the size of these filtered image matrices while maintaining potentially important
195    discriminatory information obtained by the convolutional filters. Finally, these matrices are
196    flattened into one-dimensional vectors and then fed into a fully connected (or "dense") artificial
197    neural network (for an accessible overview see LeCun *et al.* 2015). Thus, salient features derived
198    from the image matrix by the convolutional and pooling layers are passed into one or more
199    layers of a fully connected neural network whose output layer then yields our predicted response
200    value.

201         CNNs allow for two types of convolutional layers: 1-dimensional and 2-dimensional,
202    which differ only with respect to the possible shapes that the convolutional filter can take (Fig.
203    1B–C). 1-dimensional (1D) convolutions are often used in the application to time-series data (e.g.
204    Dieleman and Schrauwen 2014; Kim 2014), and are thus applicable to sequence alignment
205    matrices. Despite its name, a 1D filter is not a vector but rather a rectangular matrix that spans a
206    user-defined number of entries (called the "kernel size") in one dimension in the input data (in
207    our case this dimension is that of the polymorphic sites in the alignment), and stretches entirely
208    across the other dimension (in our case across all chromosomes in the sample). A 2-dimensional
209    (2D) convolutional filter, which is more often used with image data, allows the user to specify
210    both dimensions of the filter matrix (often using a square matrix). Whether 1- or 2-dimensional,
211    the benefit of incorporating convolutions is that it allows the CNN to take advantage of structural
212    information in the input data. For example, from an image of a face, a CNN can learn to detect
213    the repeated pattern of the eye shape and the location of both eyes relative to one another and to
214    other features. When there is meaningful structural information such as this, CNNs tend to
215    outperform non-convolutional neural networks.

216         Here our input data is an alignment of linked segregating sites with partially shared
217    evolutionary histories. Our hope is that a CNN can discover structural information in these data
218    in order to make evolutionary inferences—for example, locating the valley in diversity at the

9

219    center of a sweep (Maynard Smith and Haigh 1974), the "shoulders" on the flanks of a sweep
220    where linkage disequilibrium and allele frequencies are both elevated (Schrider *et al.* 2015), or
221    even the spatial relationship between these patterns. We also note that neural networks such as
222    CNNs can have multiple "branches" each with separate architectures and input types—in some
223    of the cases discussed in this paper we incorporate an additional network branch whose input is
224    the vector of the positions of the segregating sites (Fig. 1B–C).

225        Like all supervised machine learning methods, a CNN must be trained on labeled
226    training data before it can make predictions on unlabeled data (i.e. data whose response variables
227    are unknown). Training is accomplished by tuning the weights and biases that control the
228    behavior of its artificial neurons so that together they maximize the accuracy of the outputs on
229    the training data. Note that the weights determined during the training process include the values
230    of the convolutional filter matrices, and thus different filters will be algorithmically created for
231    each task we address in this paper. This tuning occurs over a number of iterations using the
232    backpropagation algorithm (Rumelhart *et al.* 1986), which in modern implementations feeds a
233    small number of training examples (a "mini-batch") through the network and then estimates the
234    error gradient on the output vectors produced for these examples. The error gradient is then
235    propagated in reverse through the network—a given hidden neuron's contribution to the error is
236    proportional to the linear combination of its weight vector and the errors associated with each
237    neuron in the next layer. The weights are then updated using one of the many flavors of
238    stochastic gradient descent (e.g. Kingma and Ba 2014). This process repeats until each training
239    example has been fed through the network, marking the completion of a single training iteration.
240    Training continues for a number of these iterations (often called epochs) until a specified stopping
241    criterion is reached (e.g. a predefined number of iterations has been performed, accuracy on the
242    validation set has not improved relative to the previous iteration, etc.).

243        In the context of population genetics, the CNN's input could be a matrix of allelic states
244    at each polymorphic site (Fig. 2). For example, an alignment of haploid individuals $M$, where
245    $M_{ij}=0$ if the $i$th individual has the ancestral allele at the $j$th segregating site in the alignment, and 1
246    if this individual has the derived allele (an input format that can easily be altered to allow for
247    multiallelic polymorphisms); we adopt this approach and variants of it below. The output can be
248    a categorical indicator (e.g. whether or not the genomic window experienced a recent selective
249    sweep) in which the problem is referred to as a classification task in machine learning

10

250 terminology, a quantitative value (e.g. the population recombination rate) in which case the task

251 is referred to as regression, or a vector containing both categorical and quantitative values. Once

252 the CNN has been trained to produce the desired output, it can be applied to unlabeled data (e.g.
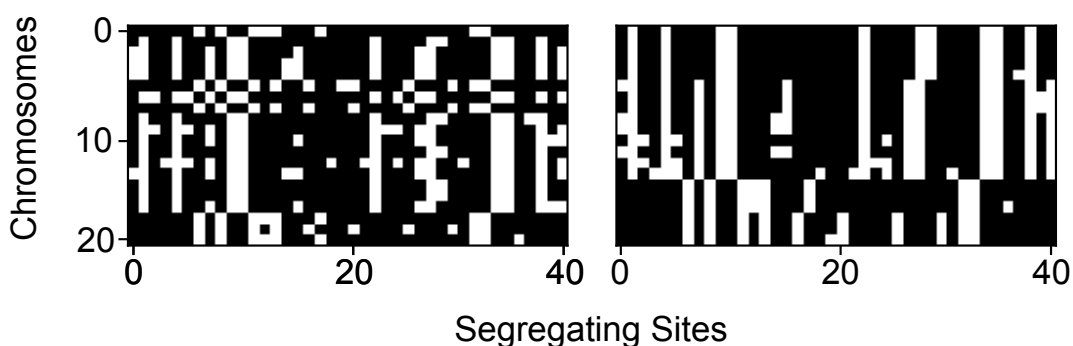
253 sequence from natural populations).

254



**Fig. 2: Example population genetic alignments visualized as black-and-white images.** An unsorted alignment matrix (left) and this same matrix sorted by genetic similarity among chromosomes (right) are shown. Each row represents one of twenty chromosomes in the sample and each column represents one of forty segregating sites. Derived and ancestral states are encoded as black and white, respectively.

255

256 Because supervised machine learning relies on predictive functions tuned algorithmically

257 from training data, CNNs can be applied to any problem for which a training set can be

258 obtained, and therefore our inference is not limited to problems for which appropriate likelihood

259 models or statistics have been derived and implemented. In a population genetics context,

260 coalescent simulations provide a versatile and computationally efficient (Hudson 2002; Teshima

261 and Innan 2009; Ewing and Hermisson 2010; Kelleher *et al.* 2016; Kern and Schrider 2016)

262 means to generate training data. In this paper we relied exclusively on coalescent simulations to

263 produce training data for the CNN. However, compute-intensive forward population simulations

264 may offer greater flexibility than coalescent simulations in some situations, and recent advances

265 are making them more computationally feasible (Kelleher *et al.* 2018).

266

267 **Using a CNN to make inferences from an alignment: a simple test case**

268 We evaluated the performance impact of transposing the alignment matrix (so that columns

269 rather than rows correspond to chromosomes) and sorting the chromosomes in the alignment

270 matrix by genetic similarity. We did this using a 1D CNN trained to estimate the population-

271  scaled mutation rate, $\theta$, in an equilibrium population. We found that both of these techniques
272  accelerate the decline in root-mean-square error (RMSE; Fig. 3), showing that they help the
273  network achieve better performance. Transposing the alignment matrix so that chromosomes are
274  represented by rows and polymorphisms by columns has a particularly notable effect (compare
275  blue and black lines in Fig. 3). Additionally, sorting the chromosomes by genetic similarity further
276  increases the accuracy of the CNN when combined with the matrix transposition above
277  (magenta line); alternatively, using a permutation-invariant network architecture would obviate
278  any need for this step (Chan *et al.* 2018). The effect of transposition should disappear when using
279  2D convolutions because in those cases we always used a square convolutional filter matrix
280  (Methods), but we found that 1D CNNs often performed as well as 2D CNNs (data not shown).
281  Thus, unless otherwise specified we use 1D convolutions for the tasks discussed below.
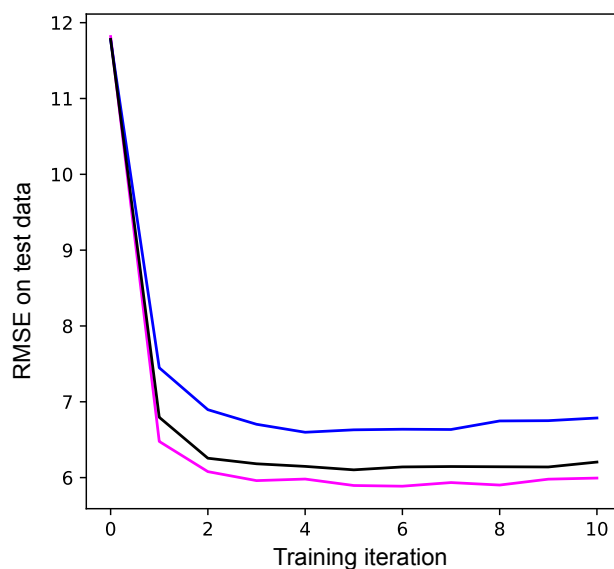282



**Fig. 3: The impact of input data reorganization on accuracy.** We show the root mean squared error (RMSE) of a 1D CNN's predictions of $\theta$ as assessed on 1,000 test alignments after a given number of training iterations. Each line is the average of 10 runs. The blue line shows accuracy after training using alignment matrices with each row representing one chromosome. The black line shows accuracy after transposing all matrices so that chromosomes correspond to columns; this makes 1D convolutional filters examine each individual at a group of adjacent segregating sites. The magenta line shows the impact of transposing matrices, and sorting the chromosomes in the alignment matrix by genetic similarity.

283
284  **CNN's can accurately detect introgressed loci**

12

285  Recent studies indicate that closely related species often exchange genes (Kulathinal *et al.* 2009;
286  Martin *et al.* 2013; Brandvain *et al.* 2014; Fontaine *et al.* 2015). There are several motivations for
287  locating genomic segments introgressed from one species into another. For one, the occurrence
288  of cross-species gene flow raises the possibility of adaptive introgression, wherein a beneficial
289  allele enters a population via migration from a related species (reviewed in Hedrick 2013).
290  Discovering introgressed loci can therefore identify alleles underlying rapid ecological adaptation
291  as well as the source of these alleles. In addition, uncovering genomic regions that are and are not
292  porous to cross species gene flow may help to illuminate the genomic basis of reproductive
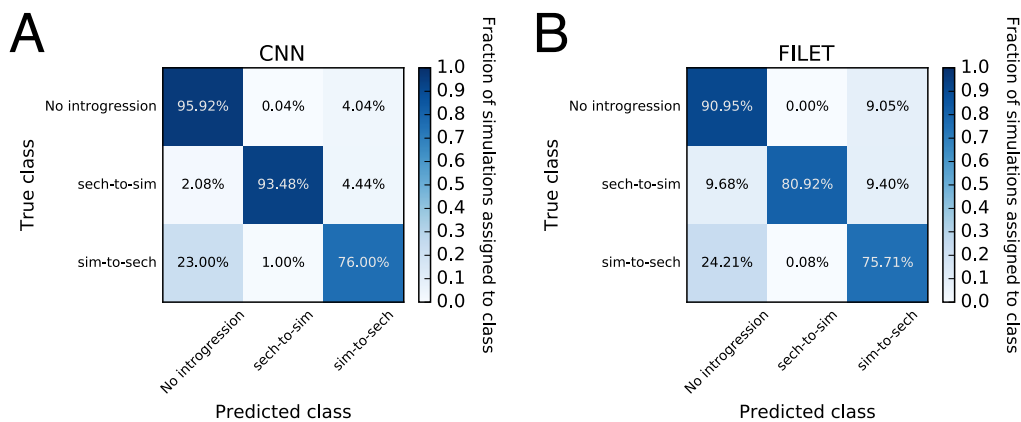293  isolation (Turner *et al.* 2005).



**Fig. 4: Performance of classifiers for detecting introgression.** We use confusion matrices to show the performance of a CNN trained to detect genomic regions of introgression between two closely related species (panel A), and a competing method that uses a vector of summary statistics to the same end (FILET; panel B). These classifiers were both trained and tested on the same data sets which were simulated under a joint demographic history inferred from a sample of *Drosophila simulans* and *D. sechellia* individuals (as described in the Methods) with and without introgression. The classifiers seek to discriminate among three classes: no introgression in the genomic window being examined, introgression from *D. sechellia* to *D. simulans*, and introgression from *D. simulans* into *D. sechellia*. Each entry in the matrix shows the fraction of test examples belonging to the class specified on the *y*-axis that were inferred by the method to belong to the class specified on the *x*-axis. Correct classifications are those found along the diagonals, while all off-diagonal entries represent incorrect classifications.

294

295  Researchers have thus sought to devise methods capable of detecting introgressed regions
296  from multispecies population genomic data sets. These include methods that attempt to infer the
297  ancestry for each individual at each site (e.g. Price *et al.* 2009; Lawson *et al.* 2012; Sohn *et al.* 2012)
298  and those that explicitly seek to discriminate between introgressed and non-introgressed loci

299    (Sankararaman *et al.* 2014; Geneva *et al.* 2015; Rosenzweig *et al.* 2016; Schrider *et al.* 2018). We

300    trained a CNN to identify introgression in a scenario modeled after the demographic history of

301    the *Drosophila simulans-D. sechellia* species pair (Methods), for which there is evidence for recent

302    gene flow (Garrigan *et al.* 2012).

303         Fig. 4A displays the results of these tests in the form of confusion matrices, which show

304    the fraction of test examples correctly predicted for each class (diagonal values) as well as the

305    fractions incorrectly assigned (off-diagonal values). To compare the performance of our CNN to

306    competing approaches, Fig. 4B displays the confusion matrix for FILET, a method previously

307    shown to outperform several methods, including two statistics for detecting introgression (Joly *et*

308    *al.* 2009; Geneva *et al.* 2015), and a tool that infers local ancestry tracks for each individual

309    (Lawson *et al.* 2012). Overall, this CNN classified 88.5% of test simulations correctly (95%

310    confidence interval: 87.7–89.2%). The most difficult scenario for the CNN was introgression

311    from *D. simulans* into *D. sechellia*, which it misclassified as "no introgression" 23% of the time. For

312    the other two classes the CNN accuracy was >95%. Importantly, for every class this CNN

313    achieved greater accuracy than FILET (overall accuracy of 82.5%; 95% confidence interval:

314    81.7%–83.4%), a machine learning approach that leverages a vector of 31 summary statistics

315    (Schrider *et al.* 2018). Thus, it is a useful measuring stick for assessing the CNN's accuracy, and

316    the CNN's success in this comparison is encouraging.

317

318    **Estimating historical recombination rates**

319    Recombination creates new combinations of alleles, and the degree of linkage between selected

320    sites affects the efficiency with which natural selection can act on each individual site (Hill and

321    Robertson 1966). The interplay of selection and recombination also influences the landscape of

322    diversity across the genome (Begun and Aquadro 1992). Knowledge of recombination rates is

323    thus key to population genetics research. As a more practical alternative to estimating rates

324    directly (e.g. from pedigrees; Kong *et al.* 2010), one can infer recombination rates from

325    population genetic data by examining associations among alleles at different sites. A number of

326    methods have been proposed to solve this problem, including summary statistic estimation

327    approaches (e.g. Hudson and Kaplan 1985; Hudson 1987; Hey and Wakeley 1997), composite

328    likelihood-based methods (e.g. Hudson 2001; McVean *et al.* 2004; Chan *et al.* 2012), and machine

329    learning tools using a vector of statistics (Lin *et al.* 2013; Gao *et al.* 2016). We sought to determine
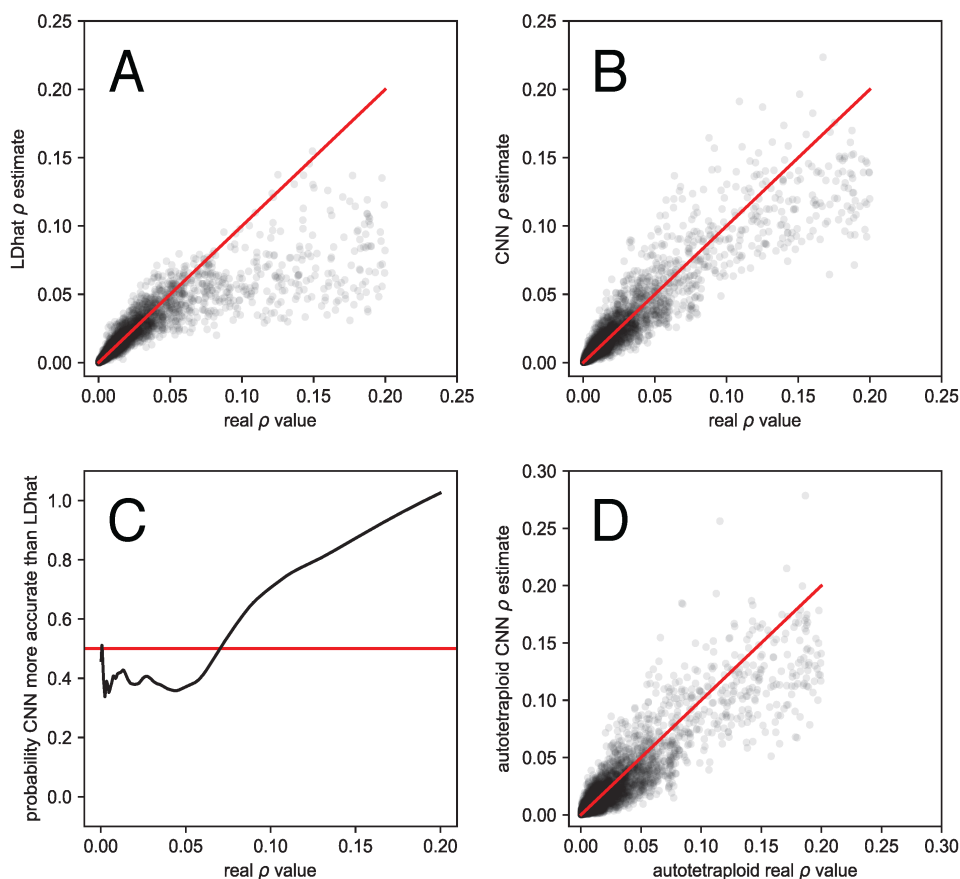
14

**Fig. 5: Accuracy of recombination rate estimates from LDhat and our CNN.** Panels A and B show the real $\rho$ values per base pair on the *x*-axes and LDhat's (A) and the CNN's (B) predictions on the *y*-axes. Panel C again shows the real $\rho$ values on the *x*-axis, and the probability that the CNN was more accurate than LDhat (black line) on the *y*-axis. This probability was calculated by scoring estimates where the CNN outperformed LDhat as one and the reciprocal as zero, and then smoothing these values with a lowess curve with a span of 15%. The red line represents the expectation if both methods had identical accuracy. Panel D shows the results from the simulated autotraploid model, with the real $\rho$ values on the *x*-axes and the CNN prediction on the *y*-axes.

330
331  whether a CNN taking an alignment image as input could be trained to tackle this task. To

332  address this problem, we first trained a CNN to estimate the historical population recombination

333  rate $\rho = 4Nr$ (where $r$ is the crossover rate per base pair per meiosis) from phased chromosomes.

334  This is the simplest scenario, as the arrangement of alleles on chromosomes is completely

335  resolved. Following training, we compared the CNN's performance to that of `LDhat` (McVean

336  *et al.* 2004), a widely used composite likelihood method, on the same testing data (Fig. 5). We

337  generated a test set of alignments whose values of $\rho$ spanned three orders of magnitude, from

338  0.0002 to 0.2 (expressed per bp). Overall, both approaches performed well at predicting the true

339  value of $\rho$. LDhat had an $R^2 = 0.77$ and an RMSE = 0.016, whereas the CNN had a $R^2 = 0.86$

340  and an RMSE = 0.011 (Fig. 5A,B). LDhat appears to estimate $\rho$ slightly better than the CNN

341  for lower recombination rates, whereas the CNN performs better at the higher values of $\rho$ (Fig.

342  5C). Additionally, the CNN appears to provide a roughly unbiased estimator of $\rho$, while LDhat's

343  estimates appear downwardly biased.

344      Because the CNN was capable of estimating $\rho$ independent of $\theta$, we were interested to see

345  how well it could interpolate between the $\theta$ values it was trained with. The CNN was trained with

346  a large gap between $N = 20,000$ and $N = 50,000$ (and thus a large gap in $\theta$; see Methods), so we

347  used coalescent simulations to generate an additional test set with $N$ values drawn uniformly

348  among 30,000, 35,000, 40,000, and 45,000. When tested on these data the CNN's predictions

349  had an $R^2 = 0.82$ and an RMSE = 0.017. This represents a slight decrease in accuracy from the

350  values obtained when tested on the same $N$ values used in training, but nonetheless shows that

351  the CNN can interpolate between training parameters without a dramatic loss in accuracy. This

352  could be a useful property, for example in cases where $N$ (or $\theta$) is unknown, but where one can

353  generate coalescent simulations across a range of plausible values.

354      Further complications arise when estimating $\rho$ from unphased data. Under this scenario

355  the arrangement of alleles on chromosomes is not known. One work-around is to first phase the

356  alleles and then infer $\rho$ as above, but not all data sources are easily phased, and phasing errors

357  will, of course, reduce accuracy. Another approach is to analyze the unphased data directly. The

358  relevant theory required to tackle this problem in a probabilistic manner has been worked out for

359  unphased diploids (Auton and McVean 2007), but expanding this theory to higher ploidies would

360  require a substantial effort. Take for example an autotetraploid with tetrasomic inheritance,

361  where there are five possible genotypes (*AAAA*, *AAAa*, *AAaa*, *Aaaa*, and *aaaa*). To further

362  complicate things, after sequencing an autotetraploid genome to a moderate depth of coverage

363  and identifying polymorphisms, the true underlying genotype may be uncertain. For example,

364  given a site with 10 reads supporting *A* and 10 supporting *a*, the true genotype could be *AAAa*,

365  *AAaa*, or *Aaaa*. To show the utility of CNNs in addressing novel population genomic inference

366  problems, we designed a CNN capable of inferring $\rho$ from a simulated set of sequence reads from

367  an unphased autotetraploid population sample.

16

368        We used a simple simulation scheme to produce read counts for each allele at each site

369    for each individual in a sample of 12 autotetraploids, each with approximately 25X expected

370    genome-wide coverage (see Methods). Rather than allelic assignments, the input matrix for this

371    CNN contains for every site in each individual the fraction of reads bearing the $a$ allele. Deriving

372    a likelihood function for $\rho$ under this formulation may be challenging, and such a solution has not

373    yet been attempted. However, appropriately designed artificial neural networks are universal

374    approximators, meaning that they have the potential to approximate any continuous function

375    over a compact input space (Hornik 1991). Thus it is possible for a CNN to approximate the

376    desired likelihood function, even in its absence. To this end we trained a CNN with a similar

377    architecture to the one used above on phased haploid chromosomes (see Methods). We evaluated

378    the performance of this CNN on a set of simulations where $\rho$ again ranged from 0.0002 to 0.2

379    (still scaling by $4N$, rather than $8N$ which would be appropriate for tetraploids, so the result can

380    be compared to those above). The CNN's predictions had an $R^2 = 0.83$ and an RMSE = 0.012

381    (Fig. 5D). As before, the estimate of $\rho$ was made independent of $\theta$, which varied over an order of

382    magnitude. The fact that this autotetraploid network performed only slightly worse than the

383    haploid version demonstrates that a CNN can solve problems for which no model-based

384    likelihood (or even composite likelihood) approach has been obtained, empowering empiricists

385    untrained in methods development to address questions specific to their biological system.

386

387    **CNNs can accurately detect and categorize signatures of recent positive selection**

388    When a new mutation is immediately favored by positive selection, it rapidly increases in

389    frequency until it fixes (i.e. completely replaces all other alleles at that site). This phenomenon,

390    referred to as a hard selective sweep, drastically reduces the amount of linked neutral variation

391    (Maynard Smith and Haigh 1974), and produces characteristic skews in the allele frequency

392    spectrum (Fay and Wu 2000) and linkage disequilibrium at linked sites (Kim and Nielsen 2004).

393    Alternatively, in a process known as a "soft sweep" populations may adapt via selection on a

394    polymorphism that has been segregating for some time, such that the adaptive allele exists on

395    numerous haplotypes (Hermisson and Pennings 2005). To uncover the mode of recent

396    adaptation and the genomic regions underlying recent adaptation, a large number of methods

397    have been devised to detect and characterize selective sweeps. These include summary statistics

398    (Kelly 1997; Fay and Wu 2000; Kim and Nielsen 2004; Voight *et al.* 2006; Garud *et al.* 2015),

399     composite likelihood-based approaches (Kim and Stephan 2002; Kim and Nielsen 2004; Nielsen
400     *et al.* 2005; Vy and Kim 2015), and supervised machine learning approaches using a vector of
401     statistics to obtain greater power than individual tests/statistics (Lin *et al.* 2011; Pybus *et al.* 2015;
402     Schrider and Kern 2016; Sheehan and Song 2016; Sugden *et al.* 2018). Although these efforts
403     have led to considerable progress, detecting and distinguishing between hard and soft sweeps
404     remains a major challenge.

405     We built a CNN to detect selective sweeps and to discriminate between hard sweeps and
406     soft sweeps. This CNN follows the S/HIC method of Schrider and Kern (2016) by casting the
407     problem as a classification task where the genomic region being examined is assigned to one of
408     five disjoint classes: a recent classic "hard" sweep, a recent "soft" sweep, a region linked to a
409     nearby hard sweep, a region linked to a nearby soft sweep, or a neutrally evolving region.
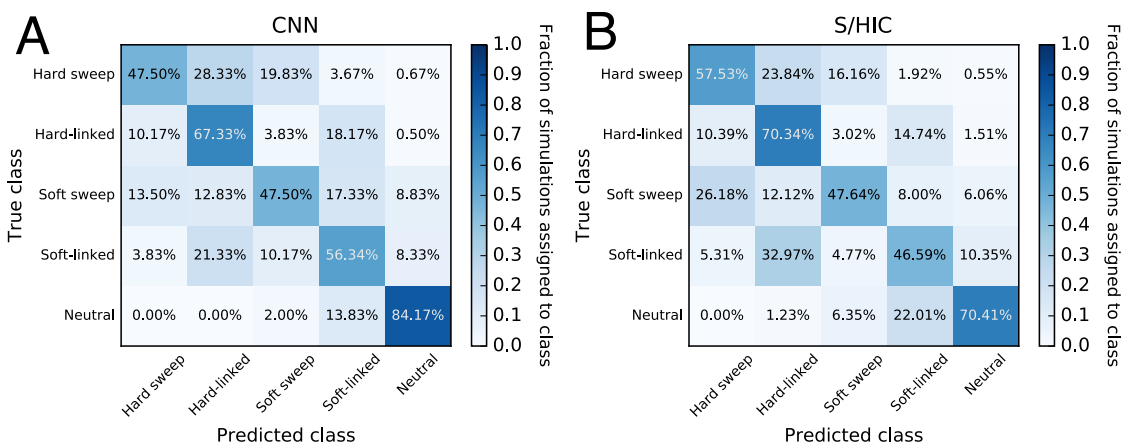


**Fig. 6: Confusion matrices showing accuracies of two methods that seek to detect recent positive selection by discriminating among hard sweeps, soft sweeps, unselected regions closely linked to hard and soft sweeps, and neutrally evolving regions.** (A) Confusion matrix summarizing the performance of our CNN, which uses an alignment image as input. (B) Performance of S/HIC, which uses a vector of summary statistics each measured in windows surrounding the region to be classified. These two classifiers were both trained and tested on the same data sets described in the Methods.

410

411     Like FILET for the problem of detecting introgression, comparing the CNN's accuracy to
412     that of S/HIC is informative because S/HIC was previously shown under a variety of simulated
413     scenarios to have greater power than a number of competing methods (Schrider and Kern 2016).
414     Rather than adopting S/HIC's approach of using a large vector of statistics, the CNN takes an
415     alignment image as input. We tested both methods against data simulated under a challenging
416     demographic history estimated from human population data (Methods). As evidenced by the

18

417   confusion matrices in Fig. 6, the CNN has slightly higher overall accuracy than S/HIC (60.6%
418   with 95% confidence interval: 58.8–62.3% for the CNN; versus 58.5% with 95% confidence
419   interval: 56.7%–60.2% for S/HIC). While S/HIC appears to be somewhat more sensitive to
420   sweeps, the CNN is achieves a more than 3-fold reduction in false positive rate: 2% of neutral
421   simulations are classified as sweeps by the CNN, versus 6.35% for S/HIC; all of these false
422   positives are classified as soft sweeps. This quality may be particularly desirable when scanning
423   genomes where sweeps are relatively rare and thus a high degree of specificity is required to
424   maintain a low false discovery rate, although the proclivity of either classifier to produce false
425   positives versus false negatives can be adjusted by imposing a posterior probability cutoff. Note
426   that these classifiers were both trained under the same demographic history from which the test
427   data were generated. We would not expect this CNN to match S/HIC's robustness to
428   demographic misspecification given that S/HIC's feature vector was designed with this in mind,
429   though we did not test this. Nonetheless, the fact that the CNN has similar accuracy to S/HIC
430   under this difficult test scenario is highly encouraging.

431

432   **CNNs can extract demographic information from alignments**
433   A major focus of population genetics research is to use genomic data to elucidate species'
434   demographic histories—the extent and timing of population size changes, and the history of
435   population splits and migration events. For example, a host of population genetic approaches
436   have been devised to infer the times and intensities of population contractions and expansions
437   over the course of a species' recent history (e.g. Marth *et al.* 2004; Schiffels and Durbin 2014; Liu
438   and Fu 2015), and to elucidate the history of population splits and subsequent gene flow (Nielsen
439   and Wakeley 2001; Hey 2009), and population merging events (e.g. Lipson *et al.* 2013; Loh *et al.*
440   2013). We asked whether CNNs can effectively extract demographic information from alignment
441   images, focusing on the task of inferring population size histories. In particular, we attempted to
442   train a CNN to estimate the parameters of a three-epoch model of instantaneous effective
443   population size changes. There are five such parameters: the ancestral population size ($N_2$), the
444   time of the more ancient population size change ($T_2$), the population size after this change ($N_1$),
445   the time of the more recent change ($T_1$), and the present-day population size ($N_0$); our response
446   variable is the vector of these 5 real-valued parameters. Thus this analysis also allows us to assess
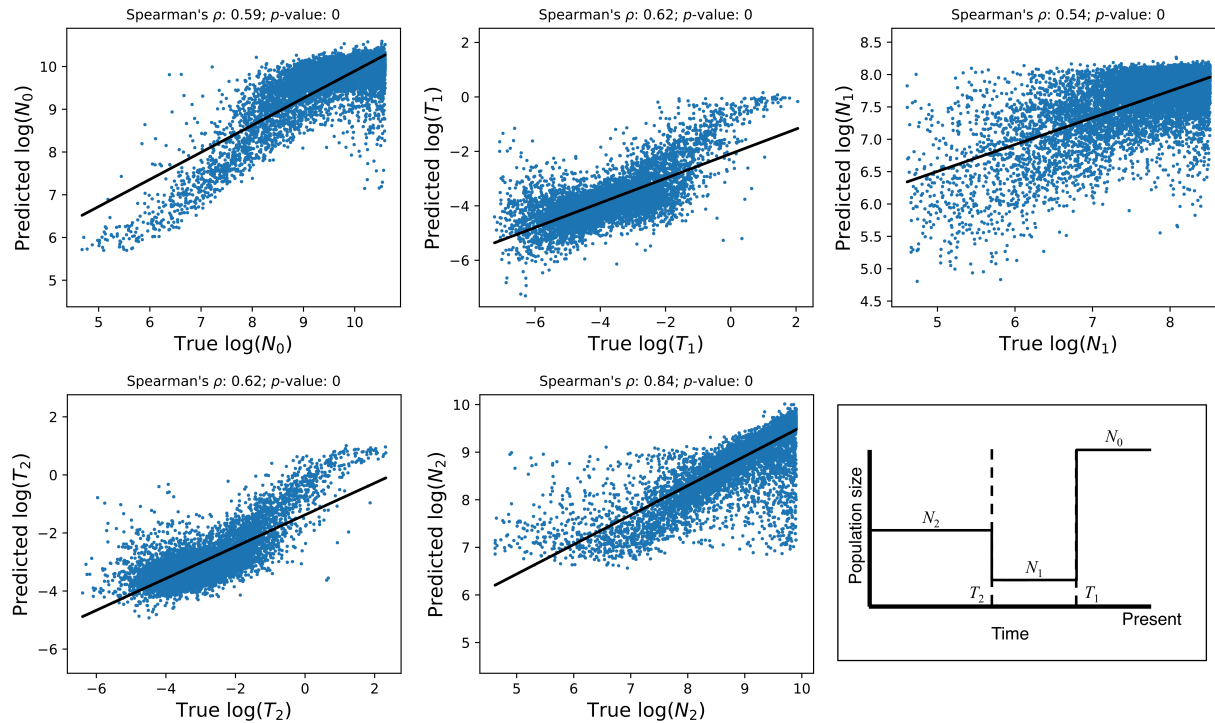447   the ability of CNNs to predict multiple population parameters simultaneously.

19

**Fig. 7: Accuracy of demographic inference CNN.** The scatterplots show the correlation between true and predicted demographic parameter values using our best-performing CNN for this task when applied to an independent test set. Note that there may be some monotonicity in the relationship between the true and predicted values of some of these parameters, which may affect calculations of the Spearman correlation coefficients shown above each scatterplot. These estimates should thus be viewed as a rough summary of this relationship, while the RMSE values reported in the text better summarize our accuracy. The inset on the bottom right shows the demographic model and its five parameters.

449

450        We simulated 50 haploid chromosomes under a variety of randomly selected population

451    size histories, and trained a CNN to estimate the demographic model parameters. The simulated

452    region was roughly equivalent in length to 1.5 Mbp of the human genome (Methods). Because

453    we found this problem to be comparatively difficult, we experimented with a variety of

454    hyperparameters governing the neural network structure and input/output format. In

455    supplementary table S1 we show the optimal RMSE (i.e. the minimum RMSE across training

456    iterations) for each hyperparameter combination examined. This experiment revealed several

457    general trends. First, 1D convolutional networks tended to fare slightly better than their 2D

458    counterparts (median RMSE of 0.52 across all hyperparameter combinations with 1D

459    convolutional filters, and median RMSE 0.54 for 2D convolutions; $p=1.1\times10^{-4}$; Mann-Whitney

20

460    *U* test); however several 2D networks performed nearly as well as the best 1D network, achieving

461    an RMSE of <0.5 while the best score obtained overall was 0.43. Second, smaller convolutional

462    filters tended to perform slightly better than larger ones—we observed a positive correlation of

463    kernel size with RMSE across hyperparameter combinations ($\rho=0.26$; $p=6.9\times10^{-4}$; Mann-

464    Whitney *U* test). For example, the median validation RMSE was 0.51 for a kernel size of 2 versus

465    0.56 for a kernel size of 10. Third, log-scaling the demographic parameters to be estimated

466    increased accuracy (RMSE decreased from 0.55 to 0.52; $p=0.020$; Mann-Whitney *U* test). For

467    this problem sorting chromosomes by relatedness resulted in a small improvement (RMSE

468    decreased from 0.54 to 0.53; $p=0.034$). Encoding ancestral and derived alleles as '0' and '255'

469    (i.e. black and white in a grayscale image), respectively, versus '-1' and '1' had a significant

470    influence on accuracy, with the former yielding better performance than the latter (RMSE of

471    0.51 vs. 0.60; $p=1.5\times10^{-15}$). Finally, using dropout resulted in a slight decrease in accuracy

472    (median RMSE increased from 0.52 to 0.55) though this was not statistically significant

473    ($p=0.092$). We note that these trends may change if the amount of training data is increased or

474    decreased, and may not necessarily hold for other tasks.

475         In Fig. 7, we show the correlation between the true and inferred values for each of these 5

476    parameters for the best performing network. For $N_0$ and $T_0$, these correlations are quite high,

477    implying that our CNN can recover the true values reasonably well. However, for the remaining

478    parameters, the correlation is lower (though still highly significant), and our CNN produces

479    downwardly biased estimates when the values of these parameters are larger. Although our

480    accuracy is far from perfect, we consider these results fairly encouraging because we are only

481    examining a single moderately sized genomic region, while other modern demographic inference

482    methods use data from across the genome. For example, ∂a∂i (Gutenkunst *et al.* 2009) uses allele

483    frequencies measured at a large number of polymorphisms (e.g. those found in all distal

484    intergenic regions across the genome; Gazave *et al.* 2014). PSMC and MSMC (Li and Durbin

485    2011; Schiffels and Durbin 2014) take data from a single very large recombining region such as

486    an entire chromosome. In essence, we are currently only able to utilize information about the

487    coalescent histories of the region in question—and this collection of histories may not match that

488    of the entire population, which would be more accurately reflected in genome-wide data. In the

489    Discussion, we address prospects for incorporating genome-scale data in demographic inference.

490

21

491 **DISCUSSION**

492 **Convolutional neural networks are well suited for population genetic problems**

493 Population geneticists have devised a wide array of computational methods to make evolutionary

494 inferences from genomic data. Typically the goal of these efforts is to aggregate information

495 across genomic sites in order to make an accurate inference. These methods include likelihood-

496 based approaches (e.g. Kim and Stephan 2002; Nielsen *et al.* 2005; Gutenkunst *et al.* 2009; Liu

497 and Fu 2015), probabilistic graphical models such as hidden Markov models (e.g. Turner *et al.*

498 2005; Boitard *et al.* 2009; Lawson *et al.* 2012), and those that rely on the use one or more

499 summary statistics designed to characterize patterns of variation within a genomic region (e.g.

500 Tajima 1989; Fu and Li 1993; Kelly 1997; Fay and Wu 2000; Kim and Nielsen 2004; Voight *et*

501 *al.* 2006; Ferrer-Admetlla *et al.* 2014). While these approaches differ substantially from one

502 another, they all have one thing in common: they make use of population genomic theory to

503 connect the features of a data set to the underlying evolutionary process. Here we have

504 demonstrated the potential of an alternative approach: treating population genetic inference as

505 an image recognition problem where the "image" is the population genetic alignment, which is

506 directly fed as input to a CNN. In contrast to most mainstream approaches, this CNN approach

507 makes use of the entirety of the data, rather than using theoretically derived estimators or closed-

508 form likelihood functions to connect a small number of features of the data to an evolutionary

509 process.

510 Here we have shown that CNNs perform remarkably well on a number of problems in

511 population genetics. We developed CNNs with comparable if not greater power to detect

512 selective sweeps, identify introgressed loci, and infer local recombination rates when compared to

513 current methods on simulated data sets. The CNNs for detecting sweeps and introgression

514 demonstrate the ability to use an alignment image to distinguish among multiple evolutionary

515 models, while the recombination rate estimator demonstrates that continuous parameters can

516 also be inferred. Finally, although our demographic parameter estimates were fairly imprecise,

517 they were only based on a short stretch of the genome, and nonetheless demonstrate that CNNs

518 have the potential to infer multiple parameters from a sequence alignment. While we were in the

519 process of preparing this manuscript, Chan et al. completed an important study demonstrating

520 that a CNN can accurately detect recombination hotspots (Chan *et al.* 2018). Taken together

22

521   these results suggest that CNNs have enormous potential as a general paradigm for population
522   genetic inference.

523   The effectiveness and generality of CNNs in population genetic inference should not be
524   surprising. CNNs offer a number of intrinsic advantages that make them particularly amenable
525   to population genetic data. First, there have been a number of efforts to move in the direction of
526   making inferences on the basis of the full complement of data present in an alignment rather
527   than one or more summary statistics (Li and Stephens 2003; Lawson *et al.* 2012; Smith *et al.*
528   2018). CNNs represent a natural way of examining the entirety of an alignment in order to
529   increase inferential power. The development of novel CNN architectures to better handle spatial
530   associations in the data across multiple scales (Yu and Koltun 2015) has the potential to improve
531   CNN-driven population genetic inference even further. For example, improved ability to detect
532   both the localized reduction in diversity at a sweep (Maynard Smith and Haigh 1974) as well as
533   the potentially confounding skews in patterns of diversity produced in its flanking regions
534   (Schrider *et al.* 2015) would be beneficial in sweep detection.

535   Another desirable property of CNNs is that they effectively perform automated feature
536   detection (LeCun *et al.* 2015). Because they discover discriminatory information directly from the
537   image, there is no need to manually construct an optimal set of features. CNNs may thus
538   outperform methods based on a set of manually curated features as observed here, although this
539   may not be the case for all tasks (e.g. Bellot *et al.* 2018). This brings up perhaps the strongest
540   quality of CNNs in the context of evolutionary inference: because CNNs can make inference in
541   the absence of statistics or a likelihood function, they can make predictions for phenomena for
542   which there exists no analytical expectation.

543   Indeed, CNNs can tackle problems for which no relevant summary statistics have been
544   devised—vectors of such statistics are required for other likelihood-free methods such as ABC
545   (Beaumont 2010) or traditional supervised machine learning techniques (Schrider and Kern
546   2018). On a related note, neural networks are particularly amenable to the incorporation of
547   disparate data types with no prior knowledge of their relationships. For example, here we have
548   included both genotype information and positional information for segregating sites as branches
549   to our networks, allowing both to be used together in prediction despite the fact that our network
550   isn't instructed how these two pieces of information relate to one another. All that is required is
551   appropriate training data. Thus, we may not have to wait for theoretical advances in order to

23

552  draw inferences from data, provided we are concerned with evolutionary models for which
553  training data can be obtained from simulation—including the wide range of scenarios that could
554  potentially be investigated via increasingly flexible and efficient forward simulators (Thornton
555  2014; Haller and Messer 2017; Kelleher *et al.* 2018).

556  This point is driven home by the success of our CNN for estimating recombination rates
557  in autotetraploids from read pileup information alone—despite the input's lack of genotype calls,
558  let alone phased haplotypes, these inferences are nearly as accurate as those that we obtained
559  from haplotype alignments. This result also suggests that CNNs may be well suited for other
560  inferences where genotype calls are unreliable (e.g. low coverage sequencing data; Korneliussen *et*
561  *al.* 2014) or unobtainable (e.g. pooled population sequencing; Schlötterer *et al.* 2014). Given
562  CNNs' flexibility, future studies should evaluate their potential to tackle not only those problems
563  examined in this paper, but the myriad additional important challenges in evolutionary genetics
564  to which they could be readily applied, including but not limited to uncovering adaptive
565  introgression (Racimo *et al.* 2016), joint inference of selective and demographic histories (Sheehan
566  and Song 2016), and even inferring structured outputs such as ancestral recombination graphs
567  (Rasmussen *et al.* 2014).

568

569  **To what extent are CNNs robust to model misspecification?**
570  Another particularly encouraging result of our recombination rate estimation analysis is that we
571  were able to infer rates for data generated from a range of parameter values to which the CNN
572  had not been exposed during training with very little decrease in accuracy. This ability to
573  interpolate between training values is a particularly desirable property. First, it implies that
574  CNNs can be used to create flexible inference tools using a modest training data set, and second
575  that researchers can focus training between reasonable parameter bounds, without knowing the
576  true (and often unknowable) underlying parameters; future efforts must explore the possibility of
577  training networks to be robust to more extreme cases of model misspecification.

578  One illustrative example of the potential pitfalls of model misspecification is the problem
579  of detecting selective sweeps without accounting for confounding demographic events. For
580  example, population bottlenecks will skew genealogies in a manner similar to sweeps (Simonsen *et*
581  *al.* 1995), and thus may result in a large fraction of false positives (Jensen *et al.* 2005; Nielsen *et al.*
582  2005). Schrider and Kern (2016) were able to mitigate this problem by designing a feature vector

24

583    that is sensitive to the spatial skews in patterns of variation created by a sweep but insensitive to

584    genome-wide skews produced by demographic events. Although this strategy is not possible with

585    CNNs because they perform automated feature extraction, it may be that incorporating training

586    examples generated under potentially confounding scenarios could alleviate this issue.

587          Therefore, future work must thoroughly 1) assess how CNNs trained on data simulated

588    under one range of evolutionary parameters fare when applied to different parameterizations,

589    and 2) determine whether robustness to such misspecification might be achieved by training a

590    CNN under a wide range of parameter values that are likely to encapsulate the correct values—

591    the recombination rate estimator's successful interpolation suggests that this may be a possibility.

592    Model misspecification is not a concern for tasks where training data may be obtained without

593    simulation (e.g. detecting selective constraint; Schrider and Kern 2015), though in such cases one

594    must take care to prevent dependencies between training and test examples because of shared

595    evolutionary histories due to physical linkage or paralogy/orthology relationships (Washburn *et*

596    *al.* 2018).

597

## **Outstanding practical challenges associated with the application of CNNs to sequence data**

600    Although the CNN approach outlined above has great potential, there are several outstanding

601    challenges with applying CNNs to a wider spectrum of problems. One important obstacle is the

602    large amount of training data required by CNNs, which makes applications requiring alignments

603    of large regions (e.g. entire chromosomes) more difficult. This challenge includes both the

604    generation of large labeled training examples, and time- and memory-efficient training with these

605    large examples given limited computational resources. Fortunately, continued improvements in

606    simulation speed (Kelleher *et al.* 2016; Kelleher *et al.* 2018) and the efficiency of CNN training

607    (Chilimbi *et al.* 2014; Yu and Koltun 2015; Jouppi *et al.* 2017; Köster *et al.* 2017) is mitigating this

608    problem. Such advances would be a boon for efforts to infer demographic parameters, which

609    require simultaneously examining data sampled from across the genome or along an entire

610    chromosome, unlike scans to infer locus-by-locus histories of

611    selection/recombination/introgression. Advances in handling large or high-resolution images

612    may also prove fruitful. For example, CNN-based strategies that simultaneously examine a

613    number of smaller "patches", each covering a portion of the image rather than the entirety of the

614    image (e.g. Lu *et al.* 2015), may aid efforts to extract demographic information from genome-

615    scale data.

616           Another challenge with the application of CNNs is that their performance can be

617    sensitive to network architecture (Szegedy *et al.* 2015). There is no underlying theory for selecting

618    optimal network architecture, though improved architectures are sure to continue to arise, and

619    automated methods exist for optimizing the many hyperparameters of a given architecture (e.g.

620    Snoek *et al.* 2012). Though we uncover some promising CNN architectures for population

621    genetic inference, we suspect that substantial improvements can still be made.

622           We have also demonstrated that CNNs are sensitive to the input format of the population

623    genetic alignment, and our work has yielded several insights along this front. First, we found that

624    the ordering of haplotypes within the alignment can impact accuracy, and our results suggest that

625    it is often beneficial to reorder haplotypes so that more similar chromosomes appear next to one

626    another. This may be a suboptimal solution, and more creative approaches may be required to

627    provide a more general strategy. To this end, research into permutation-invariant neural

628    networks (Zaheer *et al.* 2017) may prove promising when dealing with sequence alignments. This

629    is evidenced by Chan et al.'s recent findings that a permutation-invariant architecture improves

630    both training speed and final accuracy of their CNN for detecting recombination rate hotspots

631    (Chan *et al.* 2018). Chan et al.'s network avoids any convolution or pooling operations that

632    combine information across individuals until an operation that collapses each column of the

633    (filtered) alignment matrix down to a single value in an order-invariant manner (e.g. site-wise

634    maximum). This design choice means that permuting the order of individuals within the

635    alignment will have no impact on their network's output. We also observed that 1D convolutions

636    in the proper orientation perform as well as the more widely used 2D convolutions in many

637    cases. Also, scaling response variables for regression problems (both log-scaling and

638    standardization) may also affect accuracy. We therefore recommend that users experiment with

639    these different ways of representing their data, as well as different CNN architectures, in order to

640    find the design that works best for the task at hand.

641           Another important consideration of CNNs is that once trained, they are specialized to a

642    particular problem as defined by the training set. That is, a CNN trained to infer recombination

643    rates under a European demographic history may have reduced accuracy when applied to an

644    African sample. Training under a variety of demographic scenarios may make a CNN more

26

645    robust to this problem, but a question for further study is whether this can be accomplished

646    without a loss in power relative to a more specialized CNN. Even a change as subtle as adding

647    another chromosome to a dataset will make one of our previously trained CNNs inapplicable, as

648    the input matrix would no longer be the proper size and either a new CNN must be trained or

649    the data subsampled. Importantly, Chan et al. (2018) describe an architecture that can allow for

650    variation in the number individuals in the input matrix. In spite of these limitations, recent

651    advances have greatly simplified training CNNs, and it will often be practical—or even

652    preferable—for a researcher to create a CNN tailored to their specific data set.

653

654    **Are CNNs a black box?**

655    Artificial neural networks are algorithms that seek to maximize their predictive accuracy by

656    optimizing their internal mathematical operations on training data and CNNs are an extremely

657    flexible subclass of these methods because they can act directly on the input data matrix (in our

658    case a sequence alignment). However, one consequence of this is that CNNs are in some ways a

659    "black box". For example, a CNN cannot "explain" why it made a particular prediction given its

660    input. Supervised machine learning algorithms in general have perhaps been unfairly maligned

661    with this "black box" label. These methods can in principle reveal much about underlying

662    processes by determining which features are most informative under certain scenarios (i.e. feature

663    ranking; see Breiman 2001). For example, the observation that certain features are highly

664    informative for recent but not ancient introgression (Schrider *et al.* 2018) suggests some key

665    differences between the genealogies produced under these two scenarios. Due to their complex

666    inner workings, less progress has been made in breaking through the CNN "black box" as

667    compared to more traditional supervised machine learning techniques. However, some successful

668    explanatory tools are available for CNNs (Ribeiro *et al.* 2016), and there is ongoing research in

669    this area. Moreover, because the CNN framework we adopt here works on images, it may be

670    possible to translate future breakthroughs in CNN interpretation from other fields (e.g. image

671    recognition) into population genetic inference. Thus a more optimistic view is that as CNNs and

672    related methods become more interpretable, these likelihood-free image recognition approaches

673    may help to reveal theoretical insights into evolutionary processes.

674        In the near-term, CNNs may remain useful only as a predictive tool, and we will continue

675    to rely on theoretical advances to improve our understanding of population genetic processes. In

676 spite of the shortcomings noted above, the highly encouraging results that we have laid out here
677 suggest that CNNs are able to discover information about the underlying genealogies from
678 alignment images and to use this information to more accurately elucidate the evolutionary
679 phenomena that have shaped these genealogies. CNNs have enormous potential for population
680 genomic inference. We believe that progress on a host of problems could accelerate appreciably
681 were this technology to be embraced by the field. Indeed, when it comes to the business-end of
682 population genetics—drawing accurate evolutionary inferences from data—we predict that
683 increasingly, likelihood-free approaches such as the ones we have describe here will prove most
684 effective at solving existing problems, and expand the universe of problems that researchers can
685 investigate.

686

687 **MATERIALS AND METHODS**

688

689 **Computational environment for training CNNs**

690 All CNNs used in this study were developed using two open source Python packages: Keras
691 (version 2.0.6; https://keras.io/) to define neural network architecture and orchestrate training
692 and testing, and TensorFlow (version 1.1.0; https://www.tensorflow.org/) as the backend (i.e.
693 TensorFlow performs the computation during training/testing). CNN training is computationally
694 intensive, but cloud-based GPU resources have made it affordable. As an example, our network
695 for detecting selective sweeps was trained on a cloud-based system with one Nvidia K80 GPU. It
696 took 6.6 hrs to train, and at $0.90 US dollars per hour the total cost was under $7. All code used
697 for training is available online (https://github.com/flag0010/pop_gen_cnn).

698

699 **CNN validation strategy**

700 For each task, we divided our simulated inputs into three sets: a training set, a validation set, and
701 a test set. The training set was used to optimize the weights and biases of the CNN. The
702 validation set was used during training to determine how well the CNN generalizes to unseen
703 data, and adjustments were made to the CNN to improve its performance on the validation data.
704 We also used the validation set to terminate training once accuracy on this set appeared to
705 plateau—this process took different numbers of iterations for different tasks. Finally, the test set
706 was used to obtain a performance assessment of the final trained network. Importantly, this test

707 set was previously unseen by the CNN and therefore yields an unbiased evaluation of its

708 accuracy. We used binom.test in R to estimate 95% confidence intervals for classification

709 accuracies.

710

711 **Evaluating techniques for rescaling and reordering inputs to improve CNN**

712 **accuracy**

713 To evaluate the impact of alternative data preparation techniques, we developed a simple CNN

714 that estimates the locus-wide population mutation rate $\theta=4N\mu L$ where $\mu$ is the mutation rate per

715 base pair per generation and $L$ is the physical length of the locus being examined. This CNN is

716 trained using alignment images with forty chromosomes and $\theta$ drawn uniformly between 10 and

717 50 as simulated for a panmictic, constant sized population by **ms** (Hudson 2002). We trained this

718 CNN to minimize the root mean squared error (RMSE) between its prediction and the true value

719 of $\theta$ using 4,000 training matrices. Then its accuracy was scored on 1,000 test matrices that the

720 CNN was never trained on. These values were compared under different data preparation

721 approaches described below.

722 First, the matrices output by most coalescent simulation software, including **ms**, encode

723 ancestral and derived alleles for bialleleic sites as 0 and 1, respectively, and present the matrix

724 with phased haploid chromosomes as rows and sites as columns. When doing 1D convolutions,

725 we sought to use row-wise convolutional filters (Fig. 1C), i.e. those that examine each

726 chromosome in our sample across a small number of contiguous segregating sites (specified by

727 the "kernel_size" parameter in Keras) before sliding the filter forward one site (our stride length,

728 "strides" in Keras, was always set to 1). At present Keras does not allow for row-wise 1D

729 convolutions, so we accomplished this by transposing the alignment matrix and performing

730 column-wise convolutions.

731 We also assessed the impact on accuracy of sorting the chromosomes in the alignment by

732 genetic similarity. For example, the matrices in Fig. 2 contain identical information, but

733 chromosomes in the matrix on the left are randomized, while on the right they are sorted by

734 genetic similarity. We offer a fast algorithm for sorting matrices by genetic similarity

735 (https://github.com/flag0010/pop_gen_cnn/blob/master/sort.min.diff.py).

736

737 **Introgression detection**

738   To detect introgression, we simulated phased haploid training and test examples with msmove

739   (https://github.com/geneva/msmove) from the same demographic model that Schrider et al.

740   (2018) used to train the FILET classifier for detecting introgression between *Drosophila simulans*

741   and *D. sechellia*. In total we produced 237,500 coalescent simulations from 3 classes: 112,500

742   without no migration between species (No Introgression), 112,500 with gene flow from *D.*

743   *simulans* into *D. sechellia* (*sim→sech*), and 12,500 with gene flow from *D. sechellia* into *D. simulans*

744   (*sech→sim*). We used fewer *sech→sim* examples because test runs on smaller training sets suggested

745   that the network could detect this class fairly accurately, which allowed us to increase the

746   sampling of the other two more challenging classes by simulating more examples from them. To

747   our knowledge this approach of intentionally inflating the number and proportion of training

748   examples from the more challenging classes is unusual, as typically a balanced training set is

749   preferred. However we found that including additional examples from classes into our data set

750   substantially improved our ability to correctly them. The simulations were randomly assigned to

751   training and validation sets so that the training set included 107,500 examples each from the No

752   Introgression and *sim→sech* classes, and 7,500 examples from the *sech→sim* class. Both the

753   validation set and the test set contained 2,500 of each class (i.e. 7,500 total). Importantly, because

754   our test and validation sets were evenly balanced, they provided unbiased estimates of our

755   accuracy.

756        As in the *Drosophila* data set to which Schrider et al. applied FILET, each of our

757   coalescent simulations generated 34 chromosomes (14 *D. sechellia* and 20 *D. simulans*). Each

758   column in the alignment corresponded to a biallelic polymorphism, which was encoded as "0"

759   (ancestral allele) or "1" (derived allele) for each chromosome. In practice, the ancestral and

760   derived states may not be known with 100% certainty, and one may instead use major/minor

761   alleles, or randomly mispolarize a fraction of sites in the training data if one has an estimate of

762   the fraction of mispolarized sites in the true data. The effects of these design choices on

763   performance may then be evaluated on test data. Each matrix was organized so that individual

764   chromosomes were grouped by species. Each coalescent simulation produced a different number

765   of segregating sites (with the largest containing 1201 polymorphisms). Because the CNN's input

766   matrices must all have the same dimensions, we padded the right side of all matrices with fewer

767   than 1201 polymorphisms with columns containing only "0" until the total number of columns

768   reached 1201. Finally we transposed this matrix resulting in a 1201×34 matrix for each

30

769   coalescent simulation. In practice, one will have to set the image width to the largest number of

770   SNPs encountered across all training, test/validation, or real data examples included in the

771   analysis. Alternatively, one may select a fixed number of segregating sites to include in the

772   analysis, in which case each example may correspond to a different physical size (creating

773   additional variance in total recombination rates). Thus, when using this alternative approach,

774   one should adjust the lengths of simulated examples accordingly.

775   We trained a CNN architecture with three 1D-convolutional layers (kernel size = 2), each

776   followed by average-pooling, and finally two densely connected layers (i.e. the same network

777   architecture as the main network branch illustrated in Fig. 1C, but with one additional dense

778   layer). These layers contained 256, 128, 128, 128, and 128 neurons, respectively. To avoid

779   overfitting during training, each layer used dropout regularization (randomly removing 25% of

780   neurons between convolutional layers during each training iteration, and 50% between densely

781   connected layers) and rectified linear unit activation functions (i.e. ReLUs; Hahnloser *et al.* 2000;

782   Nair and Hinton 2010). Dropout regularization encourages the CNN to learn redundant

783   representations of the data, thereby reducing the network's dependence on individual weights

784   (Srivastava *et al.* 2014). The last layer was a sigmoid output layer with 3 neurons, each

785   corresponding to the 3 classes given above. The CNN was trained using the Adam optimization

786   procedure (Kingma and Ba 2014), a categorical cross-entropy loss function, and a mini-batch size

787   of 256. The CNN was run for 19 training iterations through the training data.

788

789   **Recombination rate: phased haplotype version**

790   For the recombination rate estimator we used `ms` (Hudson 2002) to simulate 50 phased

791   chromosomes, each with a target length of 20kb. To do so, we drew a population size ($N$) from

792   the following values: 5,000, 10,000, 15,000, 20,000, and 50,000, and set the population-scaled

793   mutation rate parameter $\theta = 4N\mu L$ (letting $\mu=1.5 \times 10^{-8}$ and $L=20$kb). We also set a population-

794   scaled recombination rate, $\rho = 4NrL$, where $r$ is the per bp crossover rate per meiosis, by drawing

795   $r$ from a bounded exponential distribution raging from $10^{-8}$ to $10^{-6}$. This yields a range of $\rho$ per

796   base pair of $2 \times 10^{-4}$ to $2 \times 10^{-1}$. These values roughly encompass the range of recombination rates

797   experienced in humans and *Drosophila*. Following this procedure, we generated 156,275

798   coalescent simulations. ~92% were used to train the CNN, and ~4% each were set aside for

799   validation and testing. To assess our CNNs ability to interpolate to unseen population sizes, we

31

800 also created 5,000 additional test matrices using the procedures above, but with $\mathcal{N}$ drawn

801 uniformly from the following: 30,000, 35,000, 40,000, and 45,000.

802    Each simulation was represented by a matrix of 50 rows, one for each chromosome, and

803 418 columns (the largest number of segregating sites). As before, we encoded the ancestral allele

804 with "0" and the derived allele with "1". Because not all simulations resulted in the same number

805 of polymorphisms, we padded both the genotype matrix and the position vector in the same

806 manner as for the introgression CNN, bringing the total size of each matrix to 50×418. Next, we

807 sorted each matrix by genetic similarity among chromosomes as described above and then

808 transposed the matrix to 418×50. We also extracted the segregating site positions vector from the

809 **ms** output which represents each position as a real number between zero (the leftmost position

810 on the simulated chromosome) and one (the rightmost position). For simulations with fewer than

811 418 segregating sites, we padded the positions vector with "-1"s.

812    We transformed the $\rho$ values for the training, validation, and test sets by taking the natural

813 log of each value and centering them on the mean of the training set. By using the mean from the

814 training set for all transformations, we ensure that there is no leakage of information between

815 training and validation/testing.

816    We trained a CNN with two input branches. The first branch took the haplotype

817 matrices as input and included three 1D-convolutional layers (kernel size = 2), each followed by

818 average-pooling. These layers contained 1250, 256, and 256 neurons, respectively. Each of these

819 layers uses dropout normalization (25%), L2-regularization of the weights ($\lambda$ = 0.0001), and

820 ReLU activation functions. The second branch took the position vector as input and contains

821 one densely connected layer with 64 neurons, again using dropout normalization (10%) and a

822 ReLU activation function. The two branches are then merged into another densely connected

823 layer of 256 neurons with ReLU activation functions. Finally, the output layer is a single neuron

824 with a simple linear activation function that predicts the continuous $\rho$ value. The CNN was

825 trained using the Adam optimization algorithm, using mean-squared error as our loss function,

826 and a mini-batch size of 32. The CNN was trained for 16 iterations.

827    We    compared    our    CNN's    results    to    those    of    **LDhat**    version    2.2a

828 (https://github.com/auton1/LDhat). We chose **LDhat** because it is widely used to estimate

829 historical recombination rates, and because it can be efficiently run on large data sets. LDhat will

830 estimate $\rho$ only for a specified population mutation rate ($\theta = 4\mathcal{N}\mu$), and we supplied it with the

32

831   exact $\theta$ value used for each coalescent simulation. This was done by creating five likelihood

832   lookup tables using the `complete` program, all set for 50 haploid chromosomes, for the

833   following $\theta$ values: 6, 12, 18, 24, and 60. Respectively, these correspond to $N$ = 5,000, 10,000,

834   15,000, 20,000, and 50,000 (the same values we used for training our CNNs). `LDhat` only

835   predicts values within the bounds of the lookup table. Therefore, to facilitate a fair comparison to

836   results from our CNN, which is unbounded, we selected the maximum $\rho$ value in the likelihood

837   lookup table to be 133.3% of the true maximum for each $\theta$. We then set the grid size of $\rho$ equal 1,

838   and estimated $\rho$ on the test set using `LDhat`'s `pairwise` program.

839       In contrast, the CNN was not provided information about $\theta$, and instead had to infer $\rho$

840   independent of $\theta$. This ability would be a desirable property for an estimator, as $\theta$ is likely to vary

841   considerably across the genome and outside of simulated data sets one may never know $\theta$

842   precisely. On the other hand, the CNN was provided with the physical distance between

843   segregating sites, information `LDhat` does not utilize but which will generally be available when

844   making inferences on real data. Both of these factors make our direct comparison of the CNN

845   with `LDhat` imperfect because each had access to information the other lacked when producing

846   its estimate. Nonetheless we consider this example a useful illustration of the CNN's

847   performance.

848

849   **Recombination rate: autotetraploid version**

850   We sought to train a CNN to estimate a locus-wide recombination rate in autotetraploid

851   genomes. To add a level of methodological realism to this problem, we did so from a matrix

852   storing a simple summary of read pileup information at each site for each individual.

853       To this end, we generated new coalescent simulations with 48 chromosomes each

854   following the procedure outlined above for the haploid CNN. This approach is reasonable

855   because it has been shown that the standard coalescent approximates the appropriate coalescent

856   for autotetraploids as long as $N$ is larger than a few hundred (Arnold *et al.* 2012). We generated

857   217,500 coalescent simulations, and randomly assigned 200,000 to the training set, 10,000 to the

858   validation set, and 7,500 to the test set. Next, within each coalescent simulation, we randomly

859   partitioned our 48 chromosomes into twelve sets of four. Each set represents one synthetic

860   autotetraploid genome and every site has five possible genotypes (*AAAA*, *AAAa*, *AAaa*, *Aaaa*, and

861   *aaaa*). For each autotetraploid genome $i$ and each site $j$ we simulated the number of reads

33

862    covering the site ($C_{ij}$) by drawing a random sample from a Poisson distribution with $\lambda = 25$. Then

863    we selected the number of reads representing the *a* allele $R_{ij} \sim Binom(n=C_{ij}, p=x_{ij})$, where $x_{ij}$

864    represents the frequency of the *a* allele in the tetraploid genotype (i.e. 0, 0.25, 0.5, 0.75, and 1 for

865    the five genotypes listed above). For each individual *i* at site *j*, the corresponding entry in the

866    input matrix was the fraction $R_{ij}/C_{ij}$, i.e. the fraction of reads supporting the derived allele. The

867    *AAAA* and *aaaa* genotypes were always 0 and 1, respectively. For the three heterozygous

868    genotypes (*AAAa*, *AAaa*, and *Aaaa*), $R_{ij}/C_{ij}$ varied based on sampling error but had expected values

869    of 0.25, 0.5, and 0.75, respectively. Thus at each site the original 48 chromosomes were reduced

870    to a set of 12 values corresponding to the fractions of reads supporting the *a* allele in a pool of

871    sequence reads from an autotetraploid sequenced at ~25X coverage. Note that this scheme

872    includes neither sequencing error, nor the site-specific depth which would be necessary to

873    calculate a likelihood, but is nonetheless adequate for our proof of concept.

874          As above, we sorted the rows of this matrix by genetic similarity and padded each matrix

875    with zeros to a length of 460 (the most segregating sites of any of the simulated matrices) before

876    transposing, yielding a 460×12 matrix. Again, we recorded the padded vector of positions from

877    the simulation output. Our CNN architecture was identical to the one given above for the phased

878    haplotype version, except for the dimensionality of the input changed to 460×12, and we

879    reduced the first convolutional layer from 1250 to 256 because of the smaller second dimension

880    of the input. The CNN was trained for 9 iterations.

881

882    **Detecting selective sweeps and discriminating between modes of selection**

883    For detecting selective sweeps, we used the same coalescent simulations that Schrider and Kern

884    (2017) used to train a classifier to detect sweeps in the JPT population (Japanese individuals from

885    Tokyo) from Phase 3 of the 1000 Genomes dataset (Auton *et al.* 2015). The JPT demographic

886    scenario is one where detecting selective sweeps is fairly difficult (see Figure S1 from Schrider and

887    Kern 2017), as expected for bottlenecked populations (Jensen *et al.* 2005). For this CNN, we

888    began with a set of 269,000 simulated genomic windows with the 5 following classes: a recent

889    hard sweep (i.e. fixation of a *de novo* beneficial mutation), a recent soft sweep (i.e. fixation of a

890    beneficial but previously neutral segregating polymorphism), a region linked to a nearby hard

891    sweep, a region linked to a nearby soft sweep, and a neutrally evolving region. Each simulated

892    alignment contained 208 chromosomes and we kept only coalescent simulations that contained ≤

893    5,000 segregating sites, and again padded with zeros so that all matrices were 208×5000. This

894    left 238,655 simulations, and from those we constructed a training set of 233,655 simulations. In

895    trial runs, we found that regions flanking hard and soft sweeps were the most difficult classes to

896    predict, so we again simulated additional examples from these more challenging classes. This

897    shifted the balance of our training set so that is was comprised of approximately 13% neutral

898    regions, 17% each for hard and soft sweeps, and 26.5% each for regions linked to nearby hard

899    and soft sweeps windows. We then set aside an evenly balanced set of 2,000 simulations for

900    validation and 3,000 for testing.

901        As before, we sorted each matrix by genetic similarity among chromosomes and then

902    transposed the matrix to 5000×208. We also extracted the segregating site positions vector from

903    these simulations which were generated by `discoal` (Kern and Schrider 2016), which like `ms`

904    represents each position as a real number between zero and one.

905        As above, we trained a CNN with two input branches. The first branch took the

906    haplotype matrices as input and included five 1D-convolutional layers (kernel size = 2), each

907    followed by average-pooling. These layers each contained 256 neurons and used dropout

908    normalization (20%). The second branch took the position vector as input and contained one

909    densely connected layer with 64 neurons, again using dropout normalization (10%). The two

910    branches were then merged into another densely connected layer of 256 neurons with 25%

911    dropout. Each hidden layer of the network used L2-regularization of the weights ($\lambda = 0.0001$)

912    and ReLU as the activation function. Finally, the output of this layer was fed to a five neuron

913    layer with softmax activation functions that predicts the five classes given above. The CNN was

914    trained using the Adam optimization algorithm, the categorical cross-entropy loss function, and a

915    mini-batch size of 32. The CNN was trained for 3 iterations.

916

917    **Inferring population size histories**

918    To show how CNNs can be used to infer species' demographic histories, and how CNN

919    architecture can impact this inference, we experimented with a variety of CNN approaches to

920    infer the 5 parameters of a 3-epoch model of instantaneous population size changes (i.e. 3

921    population sizes and 2 times of size change). We also use this challenging problem as an

922    opportunity to evaluate how alternative approaches to building a CNN can influence its

923    performance. In effect, we conducted a full grid search of the following attributes of both our

924 CNN architecture and input/output format: the dimensionality of our convolutions (1D or 2D),

925 the kernel size (i.e. the width of our 1D convolutional filters and both the height and width of our

926 square 2D filters; we tried each multiple of 2 raging from 2 to 10), whether to include dropout

927 (yes or no) following max pooling steps or dense layers, whether to sort our rows based on

928 similarity (yes or no), whether to log-scale our response variables (yes or no), and whether to

929 represent ancestral and derived alleles as -1/1 or as 0/255. When included, our dropout layers

930 immediately followed both max pooling steps, the dense layer following the distance input layer,

931 and the final dense layer. Each of these dropout steps randomly removed 25% of neurons. Each

932 response variable was transformed to a $z$-score according to the sample mean and variance for

933 that variable across all simulated examples.

934 The network we used for this task had two branches: a standard CNN like that depicted

935 in Fig. 1B–C but with more convolutional layers (four CNN layers each producing 128 filters and

936 each followed by a max pooling layer with a kernel size of 2), and a dense neural network layer

937 (consisting of 32 nodes) taking positional information as its input, and concatenating its output

938 with that of the final max pooling layer of the CNN prior to being fed into the final dense layer

939 (256 nodes). The positional information was a vector, $\boldsymbol{d}$, whose length was the maximum of the

940 number of segregating sites observed across all simulated examples minus one. Each value in the

941 vector $d_i$ was simply the distance (scaled between zero and one where one is the total length of the

942 simulated region) between segregating site $i$ and site $i$-1.

943 In total, we simulated 100,000 alignments of phased chromosomes using `ms`. 10,000 each

944 were set aside for testing and validation, while the remaining 80,000 were used for training. The

945 simulated population size histories were generated randomly—each demographic model

946 parameter was drawn uniformly from a range listed in supplementary table S2. Each simulated

947 region was roughly equivalent 1.5 Mbp in the human genome, assuming per base pair mutation

948 and recombination rates of $1.2 \times 10^{-8}$ and $1 \times 10^{-8}$, respectively. However, in order to make the size

949 of the simulation output more tractable for processing in a CNN we divided the mutation rate by

950 10 (equivalent to randomly downsampling the number of polymorphisms included in the input

951 by a factor of 10). During training we used a batch size of 200, trained our networks for up to 10

952 iterations, and retained the best performing CNN as assessed on the validation set. Often the best

953 CNN was obtained prior to completing all 10 training iterations. We then evaluated the

954 performance of the best CNN for each network architecture and input format on the test set by

36

955 calculating total RMSE (our loss function for this task); we also calculated Spearman correlation

956 coefficients between the true and predicted values for each of the five demographic model

957 parameters.

958

**ACKNOWLEDGMENTS**

963

**REFERENCES**

965 Arnold, B., K. Bomblies and J. Wakeley, 2012 Extending coalescent theory to autotetraploids.
966         *Genetics* **192:** 195-204.
967 Auton, A., L. D. Brooks, R. M. Durbin, E. P. Garrison, H. M. Kang *et al.*, 2015 A global
968         reference for human genetic variation. *Nature* **526:** 68-74.
969 Auton, A., and G. McVean, 2007 Recombination rate estimation in the presence of hotspots.
970         *Genome Res.* **17:** 1219-1227.
971 Beaumont, M. A., 2010 Approximate Bayesian computation in evolution and ecology. *Annual*
972         *review of ecology, evolution, and systematics* **41:** 379-406.
973 Begun, D. J., and C. F. Aquadro, 1992 Levels of naturally occurring DNA polymorphism
974         correlate with recombination rates in *D. melanogaster*. *Nature* **356:** 519-520.
975 Begun, D. J., A. K. Holloway, K. Stevens, L. W. Hillier, Y.-P. Poh *et al.*, 2007 Population
976         genomics: whole-genome analysis of polymorphism and divergence in Drosophila
977         simulans. *PLoS Biol.* **5:** e310.
978 Bellot, P., G. de los Campos and M. Pérez-Enciso, 2018 Can Deep Learning Improve Genomic
979         Prediction of Complex Human Traits? *Genetics***:** genetics. 301298.302018.
980 Boitard, S., C. Schlötterer and A. Futschik, 2009 Detecting selective sweeps: a new approach
981         based on hidden Markov models. *Genetics* **181:** 1567-1578.
982 Brandvain, Y., A. M. Kenney, L. Flagel, G. Coop and A. L. Sweigart, 2014 Speciation and
983         introgression between Mimulus nasutus and Mimulus guttatus. *PLoS Genet.* **10:** e1004410.
984 Breiman, L., 2001 Statistical modeling: The two cultures (with comments and a rejoinder by the
985         author). *Statistical science* **16:** 199-231.
986 Chan, A. H., P. A. Jenkins and Y. S. Song, 2012 Genome-wide fine-scale recombination rate
987         variation in Drosophila melanogaster. *PLoS Genet.* **8:** e1003090.
988 Chan, J., V. Perrone, J. P. Spence, P. A. Jenkins, S. Mathieson *et al.*, 2018 A Likelihood-Free
989         Inference Framework for Population Genetic Data using Exchangeable Neural Networks.
990         *bioRxiv*.
991 Charlesworth, B., M. Morgan and D. Charlesworth, 1993 The effect of deleterious mutations on
992         neutral molecular variation. *Genetics* **134:** 1289-1303.
993 Chilimbi, T. M., Y. Suzue, J. Apacible and K. Kalyanaraman, 2014 Project Adam: Building an
994         Efficient and Scalable Deep Learning Training System, pp. 571-582 in *OSDI*.

Corbett-Detig, R., and R. Nielsen, 2017 A hidden Markov model approach for simultaneously estimating local ancestry and admixture time using next generation sequence data in samples of arbitrary ploidy. *PLoS Genet.* **13:** e1006529.

Dieleman, S., and B. Schrauwen, 2014 End-to-end learning for music audio, pp. 6964-6968 in *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. IEEE.

Dutheil, J. Y., G. Ganapathy, A. Hobolth, T. Mailund, M. K. Uyenoyama *et al.*, 2009 Ancestral population genomics: the coalescent hidden Markov model approach. *Genetics* **183:** 259-274.

Elyashiv, E., S. Sattath, T. T. Hu, A. Strutsovsky, G. McVicker *et al.*, 2016 A genomic map of the effects of linked selection in Drosophila. *PLoS Genet.* **12:** e1006130.

Ewing, G., and J. Hermisson, 2010 MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics* **26:** 2064-2065.

Fay, J. C., and C.-I. Wu, 2000 Hitchhiking under positive Darwinian selection. *Genetics* **155:** 1405-1413.

Ferrer-Admetlla, A., M. Liang, T. Korneliussen and R. Nielsen, 2014 On detecting incomplete soft or hard selective sweeps using haplotype structure. *Mol. Biol. Evol.* **31:** 1275-1291.

Fontaine, M. C., J. B. Pease, A. Steele, R. M. Waterhouse, D. E. Neafsey *et al.*, 2015 Extensive introgression in a malaria vector species complex revealed by phylogenomics. *Science* **347:** 1258524.

Fu, Y.-X., and W.-H. Li, 1993 Statistical tests of neutrality of mutations. *Genetics* **133:** 693-709.

Gao, F., C. Ming, W. Hu and H. Li, 2016 New software for the fast estimation of population recombination rates (FastEPRR) in the genomic era. *G3: Genes, Genomes, Genetics* **6:** 1563-1571.

Garrigan, D., S. B. Kingan, A. J. Geneva, P. Andolfatto, A. G. Clark *et al.*, 2012 Genome sequencing reveals complex speciation in the Drosophila simulans clade. *Genome Res.* **22:** 1499-1511.

Garud, N. R., P. W. Messer, E. O. Buzbas and D. A. Petrov, 2015 Recent selective sweeps in North American Drosophila melanogaster show signatures of soft sweeps. *PLoS Genet.* **11:** e1005004.

Gazave, E., L. Ma, D. Chang, A. Coventry, F. Gao *et al.*, 2014 Neutral genomic regions refine models of recent rapid human population growth. *Proceedings of the National Academy of Sciences* **111:** 757-762.

Geneva, A. J., C. A. Muirhead, S. B. Kingan and D. Garrigan, 2015 A new method to scan genomes for introgression in a secondary contact model. *PLoS ONE* **10:** e0118621.

Gutenkunst, R. N., R. D. Hernandez, S. H. Williamson and C. D. Bustamante, 2009 Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.* **5:** e1000695.

Hahn, M. W., 2018 *Molecular Population Genetics*. Oxford University Press.

Hahnloser, R. H., R. Sarpeshkar, M. A. Mahowald, R. J. Douglas and H. S. Seung, 2000 Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature* **405:** 947.

Haller, B., and P. Messer, 2017 SLiM 2: Flexible, Interactive Forward Genetic Simulations. *Mol. Biol. Evol.* **34:** 230.

Hedrick, P. W., 2013 Adaptive introgression in animals: examples and comparison to new mutation and standing variation as sources of adaptive variation. *Mol. Ecol.* **22:** 4606-4618.

Hellenthal, G., G. B. Busby, G. Band, J. F. Wilson, C. Capelli *et al.*, 2014 A genetic atlas of human admixture history. *Science* **343:** 747-751.

Hermisson, J., and P. S. Pennings, 2005 Soft sweeps molecular population genetics of adaptation from standing genetic variation. *Genetics* **169:** 2335-2352.

Hey, J., 2009 Isolation with migration models for more than two populations. *Mol. Biol. Evol.* **27:** 905-920.

Hey, J., and J. Wakeley, 1997 A coalescent estimator of the population recombination rate. *Genetics* **145:** 833-846.

Hill, W. G., and A. Robertson, 1966 The effect of linkage on limits to artificial selection. *Genetics Research* **8:** 269-294.

Hobolth, A., O. F. Christensen, T. Mailund and M. H. Schierup, 2007 Genomic relationships and speciation times of human, chimpanzee, and gorilla inferred from a coalescent hidden Markov model. *PLoS Genet.* **3:** e7.

Hornik, K., 1991 Approximation capabilities of multilayer feedforward networks. *Neural networks* **4:** 251-257.

Hudson, R. R., 1987 Estimating the recombination parameter of a finite population model without selection. *Genetics Research* **50:** 245-250.

Hudson, R. R., 2001 Two-locus sampling distributions and their application. *Genetics* **159:** 1805-1817.

Hudson, R. R., 2002 Generating samples under a Wright–Fisher neutral model of genetic variation. *Bioinformatics* **18:** 337-338.

Hudson, R. R., and N. L. Kaplan, 1985 Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics* **111:** 147-164.

Jensen, J. D., Y. Kim, V. B. DuMont, C. F. Aquadro and C. D. Bustamante, 2005 Distinguishing between selective sweeps and demography using DNA polymorphism data. *Genetics* **170:** 1401-1410.

Joly, S., P. A. McLenachan and P. J. Lockhart, 2009 A statistical approach for distinguishing hybridization and incomplete lineage sorting. *The American Naturalist* **174:** E54-E70.

Jouppi, N. P., C. Young, N. Patil, D. Patterson, G. Agrawal *et al.*, 2017 In-datacenter performance analysis of a tensor processing unit, pp. 1-12 in *Proceedings of the 44th Annual International Symposium on Computer Architecture*. ACM.

Kaplan, N. L., R. Hudson and C. Langley, 1989 The" hitchhiking effect" revisited. *Genetics* **123:** 887-899.

Kelleher, J., A. M. Etheridge and G. McVean, 2016 Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Comput. Biol.* **12:** e1004842.

Kelleher, J., K. Thornton, J. Ashander and P. Ralph, 2018 Efficient pedigree recording for fast population genetics simulation. *bioRxiv*: 248500.

Kelly, J. K., 1997 A test of neutrality based on interlocus associations. *Genetics* **146:** 1197-1206.

Kern, A. D., and D. Haussler, 2010 A population genetic hidden Markov model for detecting genomic regions under selection. *Mol. Biol. Evol.* **27:** 1673-1685.

Kern, A. D., and D. R. Schrider, 2016 discoal: flexible coalescent simulations with selection. *Bioinformatics* **32:** btw556.

Kim, Y., 2014 Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.

Kim, Y., and R. Nielsen, 2004 Linkage disequilibrium as a signature of selective sweeps. *Genetics* **167:** 1513-1524.

1088   Kim, Y., and W. Stephan, 2002 Detecting a local signature of genetic hitchhiking along a
1089          recombining chromosome. *Genetics* **160:** 765-777.
1090   Kingma, D. P., and J. Ba, 2014 Adam: A method for stochastic optimization. *arXiv preprint*
1091          *arXiv:1412.6980*.
1092   Kong, A., G. Thorleifsson, D. F. Gudbjartsson, G. Masson, A. Sigurdsson *et al.*, 2010 Fine-scale
1093          recombination rate differences between sexes, populations and individuals. *Nature* **467:**
1094          1099-1103.
1095   Korneliussen, T. S., A. Albrechtsen and R. Nielsen, 2014 ANGSD: analysis of next generation
1096          sequencing data. *BMC Bioinformatics* **15:** 356.
1097   Köster, U., T. Webb, X. Wang, M. Nassar, A. K. Bansal *et al.*, 2017 Flexpoint: An adaptive
1098          numerical format for efficient training of deep neural networks, pp. 1742-1752 in *Advances*
1099          *in Neural Information Processing Systems*.
1100   Krizhevsky, A., I. Sutskever and G. E. Hinton, 2012 Imagenet classification with deep
1101          convolutional neural networks, pp. 1097-1105 in *Advances in neural information processing*
1102          *systems*.
1103   Kulathinal, R. J., L. S. Stevison and M. A. Noor, 2009 The genomics of speciation in
1104          Drosophila: diversity, divergence, and introgression estimated using low-coverage
1105          genome sequencing. *PLoS Genet.* **5:** e1000550.
1106   Langley, C. H., K. Stevens, C. Cardeno, Y. C. G. Lee, D. R. Schrider *et al.*, 2012 Genomic
1107          variation in natural populations of *Drosophila melanogaster*. *Genetics* **192:** 533-598.
1108   Lawrence, S., C. L. Giles, A. C. Tsoi and A. D. Back, 1997 Face recognition: A convolutional
1109          neural-network approach. *IEEE transactions on neural networks* **8:** 98-113.
1110   Lawson, D. J., G. Hellenthal, S. Myers and D. Falush, 2012 Inference of population structure
1111          using dense haplotype data. *PLoS Genet.* **8:** e1002453.
1112   LeCun, Y., Y. Bengio and G. Hinton, 2015 Deep learning. *Nature* **521:** 436-444.
1113   LeCun, Y., L. Bottou, Y. Bengio and P. Haffner, 1998 Gradient-based learning applied to
1114          document recognition. *Proceedings of the IEEE* **86:** 2278-2324.
1115   Li, H., and R. Durbin, 2011 Inference of human population history from individual whole-
1116          genome sequences. *Nature* **475:** 493-496.
1117   Li, N., and M. Stephens, 2003 Modeling linkage disequilibrium and identifying recombination
1118          hotspots using single-nucleotide polymorphism data. *Genetics* **165:** 2213-2233.
1119   Lin, K., A. Futschik and H. Li, 2013 A fast estimate for the population recombination rate based
1120          on regression. *Genetics* **194:** 473-484.
1121   Lin, K., H. Li, C. Schlötterer and A. Futschik, 2011 Distinguishing positive selection from
1122          neutral evolution: boosting the performance of summary statistics. *Genetics* **187:** 229-244.
1123   Lipson, M., P.-R. Loh, A. Levin, D. Reich, N. Patterson *et al.*, 2013 Efficient moment-based
1124          inference of admixture parameters and sources of gene flow. *Mol. Biol. Evol.* **30:** 1788-
1125          1802.
1126   Liu, X., and Y.-X. Fu, 2015 Exploring population size changes using SNP frequency spectra.
1127          *Nat. Genet.* **47:** 555-559.
1128   Loh, P.-R., M. Lipson, N. Patterson, P. Moorjani, J. K. Pickrell *et al.*, 2013 Inferring admixture
1129          histories of human populations using linkage disequilibrium. *Genetics* **193:** 1233-1254.
1130   Lu, X., Z. Lin, X. Shen, R. Mech and J. Z. Wang, 2015 Deep multi-patch aggregation network
1131          for image style, aesthetics, and quality estimation, pp. 990-998 in *Proceedings of the IEEE*
1132          *International Conference on Computer Vision*.

Marth, G. T., E. Czabarka, J. Murvai and S. T. Sherry, 2004 The allele frequency spectrum in genome-wide human variation data reveals signals of differential demographic history in three large world populations. *Genetics* **166:** 351-372.

Martin, S. H., K. K. Dasmahapatra, N. J. Nadeau, C. Salazar, J. R. Walters *et al.*, 2013 Genome-wide evidence for speciation with gene flow in Heliconius butterflies. *Genome Res.* **23:** 1817-1828.

Maynard Smith, J., and J. Haigh, 1974 The hitch-hiking effect of a favourable gene. *Genet. Res.* **23:** 23-35.

McVean, G. A., S. R. Myers, S. Hunt, P. Deloukas, D. R. Bentley *et al.*, 2004 The fine-scale structure of recombination rate variation in the human genome. *Science* **304:** 581-584.

Mitchell, T. M., 1997 Artificial neural networks. *Machine Learning* **45:** 81-127.

Nair, V., and G. E. Hinton, 2010 Rectified linear units improve restricted boltzmann machines, pp. 807-814 in *Proceedings of the 27th international conference on machine learning (ICML-10)*.

Nielsen, R., and J. Wakeley, 2001 Distinguishing migration from isolation: a Markov chain Monte Carlo approach. *Genetics* **158:** 885-896.

Nielsen, R., S. Williamson, Y. Kim, M. J. Hubisz, A. G. Clark *et al.*, 2005 Genomic scans for selective sweeps using SNP data. *Genome Res.* **15:** 1566-1575.

Pavlidis, P., J. D. Jensen and W. Stephan, 2010 Searching for footprints of positive selection in whole-genome SNP data from nonequilibrium populations. *Genetics* **185:** 907-922.

Price, A. L., A. Tandon, N. Patterson, K. C. Barnes, N. Rafaels *et al.*, 2009 Sensitive detection of chromosomal segments of distinct ancestry in admixed populations. *PLoS Genet.* **5:** e1000519.

Pudlo, P., J.-M. Marin, A. Estoup, J.-M. Cornuet, M. Gautier *et al.*, 2016 Reliable ABC model choice via random forests. *Bioinformatics* **32:** 859-866.

Pybus, M., P. Luisi, G. M. Dall'Olio, M. Uzkudun, H. Laayouni *et al.*, 2015 Hierarchical boosting: a machine-learning framework to detect and classify hard selective sweeps in human populations. *Bioinformatics* **31:** 3946-3952.

Racimo, F., D. Marnetto and E. Huerta-Sanchez, 2016 Signatures of archaic adaptive introgression in present-day human populations. *Mol. Biol. Evol.* **34:** 296-317.

Rasmussen, M. D., M. J. Hubisz, I. Gronau and A. Siepel, 2014 Genome-wide inference of ancestral recombination graphs.

Ribeiro, M. T., S. Singh and C. Guestrin, 2016 Why should i trust you?: Explaining the predictions of any classifier, pp. 1135-1144 in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM.

Ronen, R., N. Udpa, E. Halperin and V. Bafna, 2013 Learning natural selection from the site frequency spectrum. *Genetics* **195:** 181-193.

Rosenzweig, B. K., J. B. Pease, N. J. Besansky and M. W. Hahn, 2016 Powerful methods for detecting introgressed regions from population genomic data. *Mol. Ecol.* **25:** 2387-2397.

Rumelhart, D. E., G. E. Hinton and R. J. Williams, 1986 Learning representations by back-propagating errors. *Nature* **323:** 533.

Sankararaman, S., S. Mallick, M. Dannemann, K. Prüfer, J. Kelso *et al.*, 2014 The genomic landscape of Neanderthal ancestry in present-day humans. *Nature* **507:** 354-357.

Schiffels, S., and R. Durbin, 2014 Inferring human population size and separation history from multiple genome sequences. *Nat. Genet.* **46:** 919-925.

Schlötterer, C., R. Tobler, R. Kofler and V. Nolte, 2014 Sequencing pools of individuals—mining genome-wide polymorphism data without big funding. *Nature Reviews Genetics* **15:** 749.

Schrider, D., J. Ayroles, D. R. Matute and A. D. Kern, 2018 Supervised machine learning reveals introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLoS Genet.* **14:** e1007341.

Schrider, D. R., and A. D. Kern, 2015 Inferring selective constraint from population genomic data suggests recent regulatory turnover in the human brain. *Genome Biol. Evol.* **7:** 3511-3528.

Schrider, D. R., and A. D. Kern, 2016 S/HIC: Robust Identification of Soft and Hard Sweeps Using Machine Learning. *PLoS Genet.* **12:** e1005928.

Schrider, D. R., and A. D. Kern, 2017 Soft sweeps are the dominant mode of adaptation in the human genome. *Mol. Biol. Evol.* **34:** 1863–1877.

Schrider, D. R., and A. D. Kern, 2018 Supervised Machine Learning for Population Genetics: A New Paradigm. *Trends Genet.* **34:** 301-312.

Schrider, D. R., F. K. Mendes, M. W. Hahn and A. D. Kern, 2015 Soft shoulders ahead: spurious signatures of soft and partial selective sweeps result from linked hard sweeps. *Genetics* **200:** 267-284.

Sheehan, S., and Y. S. Song, 2016 Deep learning for population genetic inference. *PLoS Comput. Biol.* **12:** e1004845.

Simonsen, K. L., G. A. Churchill and C. F. Aquadro, 1995 Properties of statistical tests of neutrality for DNA polymorphism data. *Genetics* **141:** 413-429.

Simonyan, K., and A. Zisserman, 2014 Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.

Smith, J., G. Coop, M. Stephens and J. Novembre, 2018 Estimating time to the common ancestor for a beneficial allele. *Mol. Biol. Evol.*

Snoek, J., H. Larochelle and R. P. Adams, 2012 Practical bayesian optimization of machine learning algorithms, pp. 2951-2959 in *Advances in neural information processing systems*.

Sohn, K.-A., Z. Ghahramani and E. P. Xing, 2012 Robust estimation of local genetic ancestry in admixed populations using a nonparametric Bayesian approach. *Genetics* **191:** 1295-1308.

Srivastava, N., G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, 2014 Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* **15:** 1929-1958.

Sugden, L. A., E. G. Atkinson, A. P. Fischer, S. Rong, B. M. Henn *et al.*, 2018 Localization of adaptive variants in human genomes using averaged one-dependence estimation. *Nature Communications* **9:** 703.

Szegedy, C., W. Liu, Y. Jia, P. Sermanet, S. Reed *et al.*, 2015 Going deeper with convolutions, pp. in *CVPR*.

Tajima, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123:** 585-595.

Tennessen, J. A., A. W. Bigham, T. D. O'Connor, W. Fu, E. E. Kenny *et al.*, 2012 Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science* **337:** 64-69.

Teshima, K. M., and H. Innan, 2009 mbs: modifying Hudson's ms software to generate samples of DNA sequences with a biallelic site under selection. *BMC Bioinformatics* **10:** 166.

Thornton, K. R., 2014 A C++ template library for efficient forward-time population genetic simulation of large populations. *Genetics* **198:** 157-166.

Turner, T. L., M. W. Hahn and S. V. Nuzhdin, 2005 Genomic islands of speciation in Anopheles gambiae. *PLoS Biol.* **3:** e285.

1226  Voight, B. F., S. Kudaravalli, X. Wen and J. K. Pritchard, 2006 A map of recent positive
1227       selection in the human genome. *PLoS Biol.* **4:** e72.
1228  Vy, H. M. T., and Y. Kim, 2015 A Composite-Likelihood Method for Detecting Incomplete
1229       Selective Sweep from Population Genomic Data. *Genetics* **200:** 633-649.
1230  Washburn, J. D., M. K. M. Guerra, G. Ramstein, K. A. Kremling, R. Valluru *et al.*, 2018
1231       Evolutionarily informed deep learning methods: Predicting transcript abundance from
1232       DNA sequence. *bioRxiv*: 372367.
1233  Yu, F., and V. Koltun, 2015 Multi-scale context aggregation by dilated convolutions. *arXiv*
1234       *preprint arXiv:1511.07122*.
1235  Zaheer, M., S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov *et al.*, 2017 Deep sets,
1236       pp. 3394-3404 in *Advances in Neural Information Processing Systems*.

1237

1238

1239  **SUPPLEMENTARY TABLE LEGENDS**

1240

1241  **Supplementary table S1: The effect of different neural network**

1242  **input/output/architecture hyperparameters on demographic inference error.**

1243

1244  **Supplementary table S2: Demographic parameter ranges used to simulate 3-epoch**

1245  **population size histories.**