

Computationally efficient demographic history inference from allele frequencies with supervised machine learning

Linh N. Tran,^{1,2} Connie K. Sun,² Travis J. Struck,² Mathews Sajan,² and Ryan N. Gutenkunst^{2,*}

¹*Genetics Graduate Interdisciplinary Program and* ²*Department of Molecular & Cellular Biology, University of Arizona, Tucson, AZ, USA.*

***Correspondence:** Ryan N. Gutenkunst, Department of Molecular & Cellular Biology, University of Arizona, Tucson, AZ, 85721, USA. E-mail: rgutenk@arizona.edu.

1 Abstract

2 Inferring past demographic history of natural populations from genomic data is of central concern
3 in many studies across research fields. Previously, our group had developed *dadi*, a widely used
4 demographic history inference method based on the allele frequency spectrum (AFS) and maximum
5 composite likelihood optimization. However, *dadi*'s optimization procedure can be computationally
6 expensive. Here, we developed *donni* (demography optimization via neural network inference), a
7 new inference method based on *dadi* that is more efficient while maintaining comparable inference
8 accuracy. For each *dadi*-supported demographic model, *donni* simulates the expected AFS for a
9 range of model parameters then trains a set of Mean Variance Estimation neural networks using the
10 simulated AFS. Trained networks can then be used to instantaneously infer the model parameters
11 from future input data AFS. We demonstrated that for many demographic models, *donni* can
12 infer some parameters, such as population size changes, very well and other parameters, such
13 as migration rates and times of demographic events, fairly well. Importantly, *donni* provides
14 both parameter and confidence interval estimates from input AFS with accuracy comparable to
15 parameters inferred by *dadi*'s likelihood optimization while bypassing its long and computationally
16 intensive evaluation process. *donni*'s performance demonstrates that supervised machine learning
17 algorithms may be a promising avenue for developing more sustainable and computationally
18 efficient demographic history inference methods.

19 INTRODUCTION

20 Inferring demographic history from genomic data has become routine in many research fields,
21 from elucidating the anthropological origins and migration patterns of modern and archaic human
22 populations (Gutenkunst et al. 2009; Bergström et al. 2020; Marchi et al. 2022; Gopalan et al. 2022),
23 to inferring the population genetic trajectories of endangered animals (Mays Jr et al. 2018; Miller-
24 Butterworth et al. 2021; Chavez et al. 2022). Accounting for demographic history is also essential in
25 setting the appropriate background for detecting signals of natural selection (Nielsen et al. 2005;
26 Boyko et al. 2008; Kim et al. 2017), disease associations (Mathieson & McVean 2012), and recomb-
27 nation hotspots (Johnston & Cutler 2012). Due to the wide range of possible demographic models
28 and high dimensionality of genome sequence data, such analysis often involves computationally
29 expensive modeling. As the size of genomic datasets rapidly grows to thousands of full genomes,
30 there is a great need for more efficient and scalable methods for extracting information from such
31 datasets.

32 One class of widely used methods infers demographic history from sequence data summarized
33 as an allele frequency spectrum (AFS). An AFS is a multidimensional array with dimension equal
34 to the number of populations being considered in a given demographic model. Each array entry
35 is the number of observed single nucleotide polymorphisms (SNP) with given frequencies in the
36 sampled populations. For example, the [1,2] entry would count SNPs that were singletons in the
37 first population and doubletons in the second. A major advantage of using the AFS as a summary
38 statistic is the ease of scaling to whole genome data (Marchi et al. 2021), as it efficiently reduces the
39 high dimensionality of population genomic data. AFS-based inference methods are, therefore, often
40 fast and suitable for exploring many demographic models (Spence et al. 2018). Given its wide use in
41 countless empirical studies, much progress has been made towards understanding the theoretical
42 guarantees and limitations of the AFS and AFS-based inference (Myers et al. 2008; Achaz 2009;
43 Bhaskar & Song 2014; Terhorst & Song 2015; Baharian & Gravel 2018).

44 Demographic inference methods based on the AFS often work by maximizing the composite
45 likelihood of the observed AFS under a user-specified demographic history model with parameters
46 such as population sizes, migration rates, and divergence times (Coffman et al. 2016). The expected

47 AFS can be computed via a wide range of approaches (Gutenkunst et al. 2009; Naduvilezhath et al.
48 2011; Lukić & Hey 2012; Excoffier et al. 2013; Kern & Hey 2017; Jouganous et al. 2017; Kamm et al.
49 2017) with varying degrees of computational expense, model flexibility, and scalability. Because
50 this is the most computationally intensive step in the procedure, new methods developed thus far
51 have focused on devising algorithms to speed up AFS calculation (Jouganous et al. 2017; Kamm
52 et al. 2017, 2020). However, not much attention has been given to optimizing how the computed
53 AFS is stored and used for inference. In a typical likelihood optimization procedure, hundreds to
54 thousands of expected AFS are computed and compared to the data to obtain the best-fit parameter
55 set. These generated AFS and their corresponding demographic parameters contain information
56 regarding the mapping between the AFS and demographic parameters but are discarded after each
57 optimization run. As there are often a few common demographic models regularly used across
58 studies, if these simulated data could be captured, stored, and distributed for future use, individual
59 groups as well as the research community as a whole could save a lot of time and computational
60 effort by avoiding unnecessary repetition.

61 The mapping between the AFS and its associated demographic history model parameters can
62 be efficiently captured by supervised machine learning (ML) algorithms. Given a training data
63 set with feature vectors (AFS — input) and labels (demographic history parameters — output),
64 these algorithms can learn the function mapping from the input to the output. While training
65 ML algorithms can be computationally intensive up front, subsequent inference from trained
66 models will have minimal cost (Schrider & Kern 2018). ML algorithms have been widely adopted in
67 population genetics over the past decade, thanks to their efficiency and flexibility. Several studies
68 have used supervised ML algorithms such as random forest (RF) and multilayer perceptron (MLP)
69 with AFS as training data for demographic model selection and demographic parameter inference
70 (Sheehan & Song 2016; Smith et al. 2017; Villanea & Schraiber 2019; Mondal et al. 2019; Lorente-
71 Galdos et al. 2019; Sanchez et al. 2021). In Smith et al. (2017) specifically, the RF algorithm was
72 used to replace the rejection step in the approximate Bayesian computation (ABC) framework,
73 significantly improving overall efficiency (Pudlo et al. 2016). This improvement in efficiency was
74 in part due to more efficient use of simulated data. Whereas in a typical ABC procedure, any

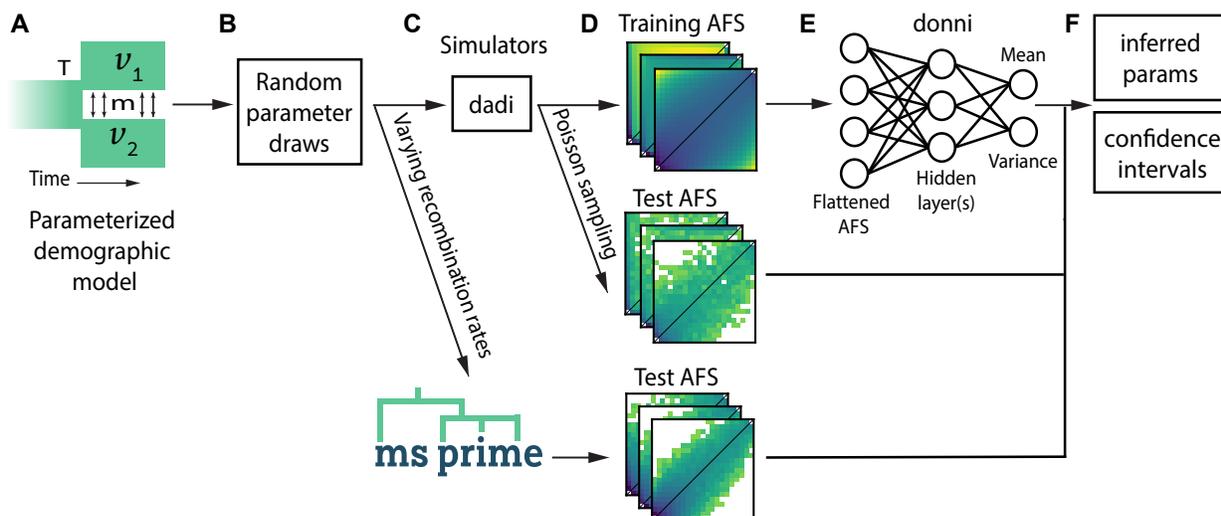


Figure 1: **Schematic of the workflow for training and testing donni.** For a given demographic model (A), we drew sets of model parameters (B) from a biologically relevant range (Table 1). Each parameter set represents a demographic history and corresponds to an expected AFS. These parameters were input into simulator programs (C) to generate training and test AFS (D). We use the expected AFS simulated with *dadi* and their corresponding parameters as training data for *donni*'s MVE networks (E). We generated test data either by Poisson sampling from *dadi*-simulated AFS or by varying recombination rates with *msprime*, resulting in a change in variance compared to training AFS. The output of *donni*'s trained networks includes both inferred parameters and their confidence intervals (F).

75 simulations beyond a threshold of difference to the data will be discarded, there all simulations
 76 were used as input for training the RF classification algorithm. The same principle can be applied
 77 in the maximum likelihood optimization and regression framework, where an ML algorithm can
 78 be trained by simulated AFS to provide estimates of demographic parameter values, bypassing
 79 likelihood optimization.

80 NEW APPROACHES

81 Here, we introduce *donni* (Demography Optimization via Neural Network Inference), a supervised
 82 ML extension to *dadi*, a widely used AFS-based method for inferring models of demographic
 83 history (Gutenkunst et al. 2009) and natural selection (Kim et al. 2017). *dadi* computes the ex-
 84 pected AFS by numerically solving a diffusion approximation to the Wright-Fisher model and uses
 85 composite-likelihood maximization to fit the model to the data. While the initial implementation

86 of the software could only handle up to three populations, a recent update supports up to five
87 populations (Gutenkunst 2021). donni uses dadi to generate AFS and demographic parameter labels
88 for training Mean Variance Estimation (MVE) networks (Nix & Weigend 1994) (Fig. 1). Researchers
89 can then use donni's trained MVE networks to instantaneously infer the parameter values and their
90 associated uncertainty from future AFS input data, obviating the need for likelihood optimization.
91 donni supports a wide range of common demographic parameters that dadi supports, including
92 population sizes, divergence times, continuous migration rates, inbreeding coefficients, and an-
93 cestral state misidentification ratios. We show that donni has inference accuracy comparable to
94 dadi but requires less computational resource, even after accounting for the cost of training the
95 MVE networks. Our library of trained networks currently includes all demographic models in the
96 dadi API as well as the models from Portik et al. (2017) pipeline. The supported sample sizes are
97 10, 20, 40, 80, and 160 haplotypes per population (up to 20 haplotypes only for three-population
98 models). For users who only need to use the trained networks for available demographic models,
99 almost no computation is required. For users who require custom models, we also provide our
100 command-line interface pipeline for generating trained models that can save time compared to
101 running likelihood optimization with dadi. Furthermore, the custom models produced can be
102 contributed to our growing library and shared with the community.

103 RESULTS

104 *Choice of MVE network for demographic history model parameter estimation with uncertainty*

105 We wanted to develop a supervised ML method that can infer not only the demographic history
106 parameters but also their associated uncertainties. Uncertainty estimation has not been the focus of
107 previous supervised neural-network-based approaches in demographic history inference (Sheehan
108 & Song 2016; Flagel et al. 2019). There are several techniques for constructing a prediction interval
109 from neural-network-based point estimation as reviewed by (Khosravi et al. 2011). Among them, the
110 MVE method is one of the most conceptually straightforward and least computationally demanding,
111 which are important factors for our goal of improving efficiency.

112 An MVE network is a feedforward neural network with two output nodes, one for the mean

113 and one for the variance (Fig. 1E). This approach provides an uncertainty estimate in a regression
114 setting by assuming that the errors are normally distributed around the mean estimation. For
115 demographic history inference, the mean is the value of the demographic history model parameter
116 we want to infer. We can construct confidence intervals using the normal distribution defined by
117 the output mean and variance estimates. There are different implementations of the feedforward
118 network architecture for MVE network (Sluijterman et al. 2023). Our implementation is a fully
119 connected network, similar to the MLP, in which all hidden layer weights are shared by the mean
120 and variance output nodes.

121 *Variance in allele frequencies affects donni training and performance*

122 Since the AFS is the key input data in our method, we first considered how different levels of
123 variance in the AFS might affect training and performance of the MVE networks underlying donni.
124 While the expected AFS computed by dadi under a given set of demographic model parameter gives
125 the mean value of each AFS entry, AFS summarized from observed data will have some variance.
126 We asked whether training the network on AFS with some level of variance or AFS with no variance
127 would lead to better overall performance. When generating AFS simulations, we modeled such
128 variance in the AFS by Poisson-sampling from the expected AFS (examples in Fig. S1A and S2B-D.)
129 We implemented four levels of AFS variance: none, low, moderate, and high in AFS used for training
130 and testing. We then surveyed the inference accuracy for all pairwise combinations for each type of
131 variance in training sets versus test sets.

132 Overall, we found that networks trained on AFS with no to moderate level of variance perform
133 similarly across all variance levels in test AFS (Fig. S3-S6 for the split-migration model). High
134 variance in training AFS led to substantially poorer performance in parameters that are more
135 difficult to infer, such as time and migration rate. The population size change and ancestral state
136 misidentification parameters were the least affected by AFS variance, and inference accuracy
137 remained similarly high across all variance scenarios. For the time parameter, training on AFS
138 with moderate variance produced the best-performing accuracy across all test cases (Fig. S4).
139 However, for the migration rate parameter, training on AFS with no variance produced the overall
140 best-performing accuracy (Fig. S5). We concluded that for subsequent analyses and model library

141 production for donni, we would train using AFS with no variance, since there was no clear benefit
142 from adding an extra variance simulation step in training. For test AFS, we would use AFS with
143 moderate level of variance to better match real data.

144 *donni is efficient and has comparable inference accuracy to dadi*

145 Since we built donni to be an alternative to dadi’s likelihood optimization, we compared with
146 dadi in our performance analysis. We validated the inference accuracy of donni for three models: a
147 two-population model with an ancestral population split and symmetric migration between the
148 populations (split-migration model, Fig. 2A), a one-population model with one size change event
149 (two epoch model, Fig. 3A), and a three-population model for human migration out of Africa (the
150 OOA model, Fig. 5A) from Gutenkunst et al. (2009). We also compared the computational efficiency
151 of donni and dadi for two different sample sizes of the split-migration model.

152 For the split-migration model, donni was able to infer all demographic history parameters
153 with accuracy comparable to dadi (Fig. 2B-I). The population size change parameters ν_1 and ν_2 were
154 inferred very well by both donni (Fig. 2B, C) and dadi (Fig. 2F, G). The time parameter T (Fig. 2D, H)
155 and migration rate m (Fig. 2E, I) were more difficult to accurately infer for both methods, with dadi
156 having trouble optimizing parameter values close to the specified parameter boundary (Fig. 2E). We
157 used Spearman’s correlation coefficient ρ to quantify the monotonic relationship between the true
158 and the inferred parameter values, similar to Flagel et al. (2019). For a more direct measurement of
159 inference accuracy, we also provide the RMSE scores for all models in Table S1.

160 To compare the efficiency of donni and dadi, we benchmarked the computational resources
161 required by each method to infer demographic parameters from the same 100 test AFS (Fig. 2J-K).
162 Since inferring parameters with donni’s trained networks is computationally trivial, we instead
163 measured the resources required by donni to generate trained networks. For both methods, compu-
164 tation was substantially more expensive as the sample size increased from 20 haplotypes to 160.
165 For dadi (Fig. 2J), there was a spread of optimization runtime among the 100 test AFS, with several
166 difficult spectra requiring more than 500 CPU hours to reach convergence for both sample sizes. By
167 comparison, the computation required for donni (Fig. 2K), including generating training data with

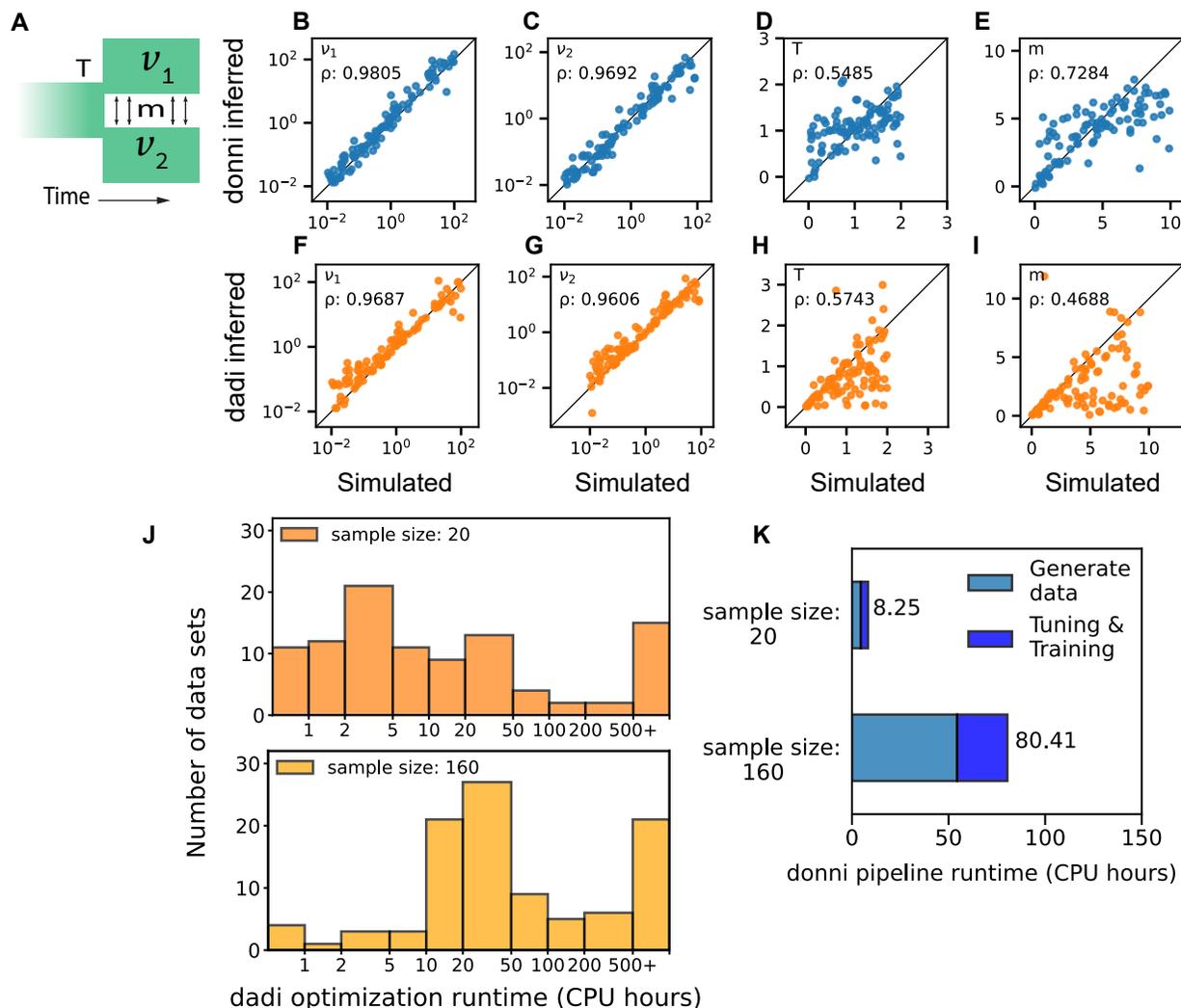


Figure 2: Inference accuracy and computing time of donni and dadi for a two-population model. (A) The two-population split-migration model with four parameters: v_1 and v_2 are relative sizes of each population to the ancestral, T is time of split, and m is the migration rate. (B-I) Inference accuracy by donni (B-E) and dadi (F-I) for the four parameters on 100 test AFS (sample size of 20 haplotypes). (J) Distribution of optimization times among test data sets for dadi. (K) Computing time required for generating donni's trained networks for two sample sizes. Generate data includes computing time for generating 5000 dadi-simulated AFS as training data. Tuning & training is the total computing time for hyperparameter tuning and training the MVE network using the simulated data.

168 dadi, hyperparameter tuning and training, was less than the average time required for running
169 dadi optimization on a single AFS. This result suggests that donni may benefit many cases where
170 dadi optimization can take a long time to reach convergence.

171 Fig. 2K also suggests that generating the expected AFS with dadi is computationally expensive,
172 often equivalent to if not more so than tuning and training a network. Such expensive operations
173 are indeed what we aimed to minimize with donni. During each dadi optimization, a large number
174 of expected AFS are also calculated. As opposed to discarding all these expensive calculations after
175 each dadi optimization, donni's trained network effectively captures the mapping between the
176 expected AFS and demographic history model parameter values in its network weights, which can
177 be reused instantaneously in the future.

178 *donni accurately estimates uncertainty of inferred parameter values*

179 Sometimes, demographic model parameters may be unidentifiable, because multiple parameter sets
180 generate nearly identical AFS. As a simple example, we considered the one-population two epoch
181 model (Fig. 3A), which is parameterized by the relative size ν of the contemporary population and
182 the time at which the population size changed T . For this model, donni inferences are inaccurate
183 when T/ν is large (Fig. 3B-C). In this parameter regime, over the time T after the size change, the
184 AFS relaxes back to that of an constant-sized equilibrium population. Therefore, in this case, the
185 true parameters are unrecoverable because the AFS itself does not have the appropriate signal to
186 infer them. While this problem may be avoided if users follow the best practice for model selection
187 of exploring simpler models before complex ones (Marchi et al. 2021), it also highlights the need
188 for uncertainty quantification, where a wide confidence interval would appropriately indicate
189 problematic inference.

190 Using the variance output from the trained MVE networks, donni can calculate any range of
191 confidence intervals specified by the user for each inferred parameter. We validated our uncertainty
192 quantification approach by measuring the observed coverage for six confidence intervals: 15, 30,
193 50, 75, 80, and 95% intervals (details in Materials & Methods). For the two-epoch model, our
194 approach provided well-calibrated confidence intervals (Fig. 3D). Considering individual test AFS,

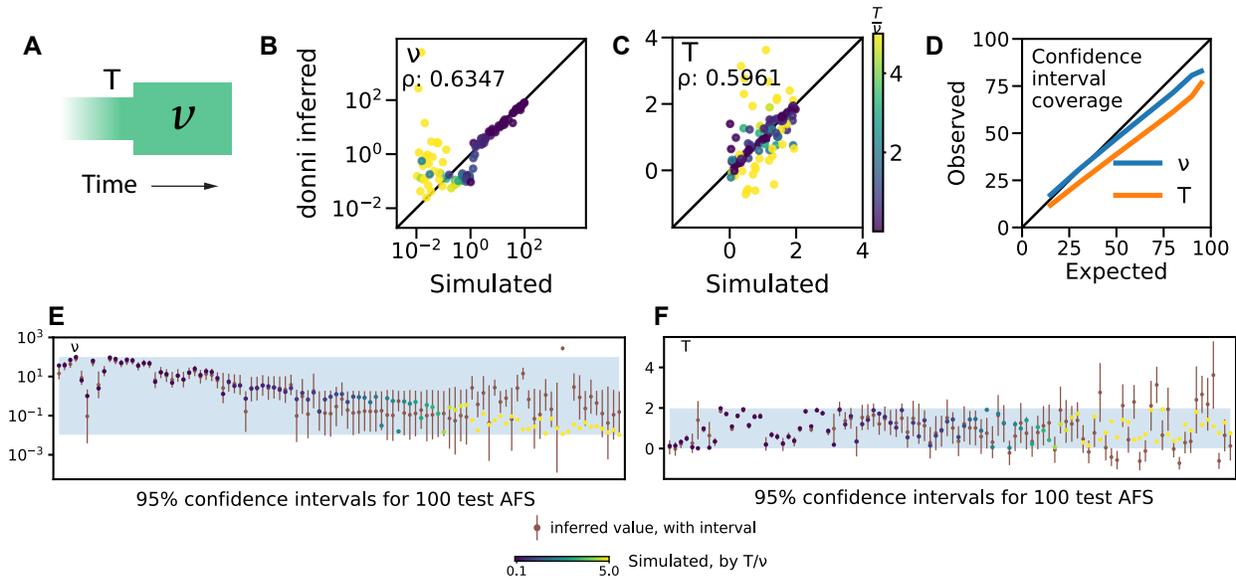


Figure 3: Uninformative AFS affecting inference accuracy and uncertainty quantification method validation. (A) The one-population two-epoch model with two parameters, ν for size change and T for time of size change. (B-C) Inference accuracy for ν and T by donni on 100 test AFS, colored by simulated $\frac{T}{\nu}$ values. (D) Confidence interval coverage for ν and T by donni. The observed coverage is the percentage of test AFS that have the simulated parameter values captured within the corresponding expected interval. (E-F) As an example, we show details of the 95% confidence interval data points from panel D for 100 test AFS. The simulated values for ν (E) and T (F) of these AFS are colored by their $\frac{T}{\nu}$ values, similar to panels B-C. donni's inferred parameter values and 95% confidence interval outputs are in brown. The percentage of simulated color dots lying within donni's inferred brown interval gives the observed coverage at 95%. The light shades are the simulated parameter range (Table S2) used in simulating training and test AFS. The 100 test AFS are sorted along the x-axis using true $\frac{T}{\nu}$ values.

195 the uninformative AFS yielded appropriately wide confidence intervals (Fig. 3E-F, yellow points).
196 We found that confidence intervals were similarly well-calibrated for the split-migration model
197 (Fig. S7).

198 *donni is not biased by linkage between alleles*

199 The Poisson Random Field model underlying *dadi* (Sawyer & Hartl 1992) and thus *donni* assumes
200 independence of all genomic loci in the data, which is equivalent to assuming infinite recombination
201 between any pair of loci. But loci within close proximity on the same chromosome are likely sorted
202 together during recombination and therefore linked. To assess how linkage affects *donni* inference,
203 we tested *donni*'s networks that were trained on *dadi*-simulated AFS without linkage on test AFS
204 simulated with *msprime*, a coalescent simulator that includes linkage (Baumdicker et al. 2022). These
205 *msprime*-simulated test AFS (examples in Fig. S1B and S2E-G) represent demographic scenarios
206 similar to those in *dadi* but also include varying levels of linkage under a range of biologically
207 realistic recombination rates. Since smaller recombination rates lead to more linkage and further
208 departure from the training data assumption, we tested *donni* on AFS with decreasingly small
209 recombination rates down to $r = 10^{-10}$ crossover per base pair per generation, which is two orders
210 of magnitude smaller than the average recombination rate in humans.

211 Population size parameters ν were inferred well no matter the recombination rate, but the
212 inference accuracy for T and m decreased as the recombination rate decreased (Fig. 4). Confidence
213 intervals were well calibrated at the higher recombination rates (Fig. 4A&E), but too small at the
214 lowest recombination rate (Fig. 4I). These patterns are similar to those we found when testing
215 the effects of AFS variance by Poisson-sampling from expected AFS with *dadi* (Fig.S3-S7), where
216 accuracy decreased with higher variance, and confidence intervals were underestimated at the
217 highest variances. Note that at $r = 10^{-10}$, linkage disequilibrium often extends entirely across the
218 simulated test regions, so in this regime methods assuming zero recombination, such as IMA3 (Hey
219 et al. 2018), may be more appropriate. Importantly, even though more linkage did lead to higher
220 variance in the estimated parameter values, we did not observe bias in our inferences.

221 *Comparison with dadi for the Out-of-Africa model*

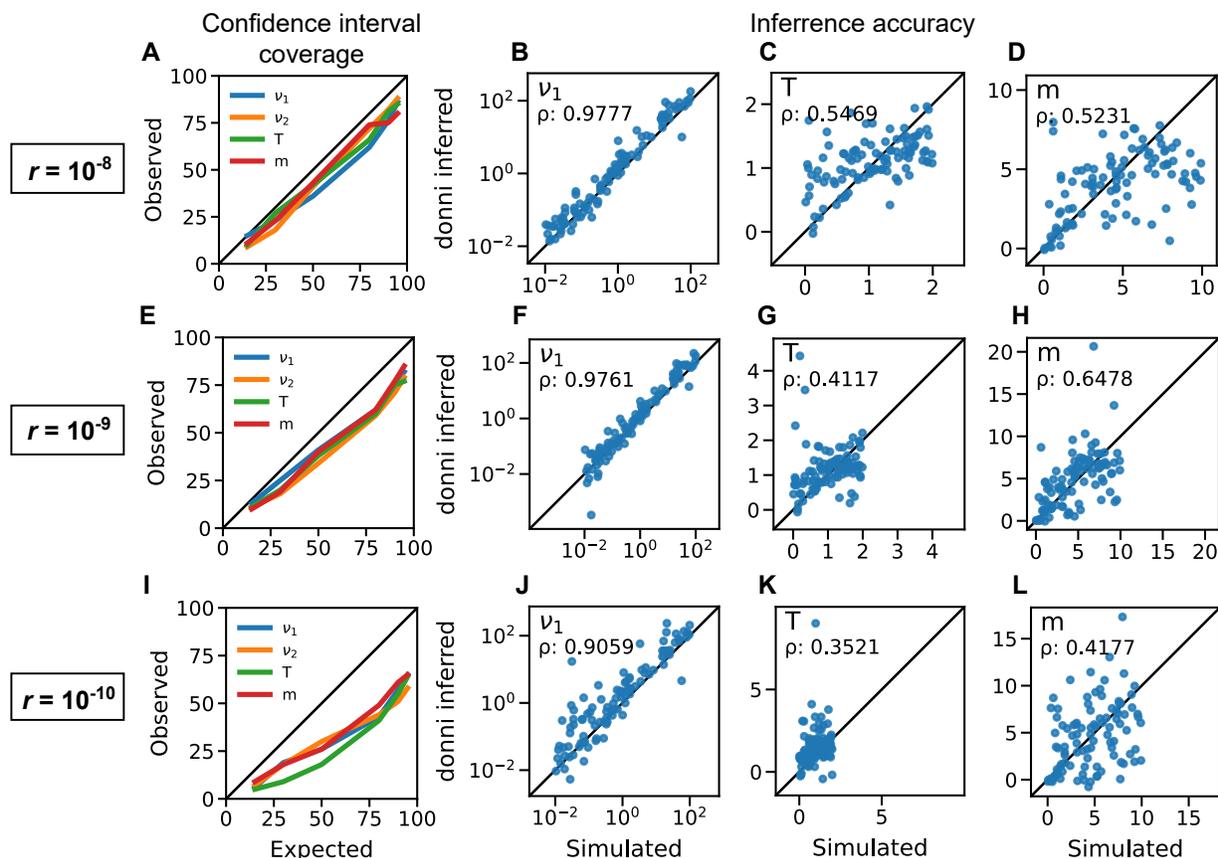


Figure 4: **donni's inference accuracy and uncertainty quantification coverage on msprime-simulated test AFS with linkage.** Each row shows the confidence interval coverage and inference accuracy for select parameters of the split-migration demographic model (Fig. 2A) at varying recombination rate. Recombination rate decreases from top to bottom row, corresponding to increased linkage and variance in the msprime-simulated test AFS. The same networks (train on dadi-simulated AFS) were used in this analysis as in Fig. 2F-I.

222 We tested donni on the three-population Out-of-Africa (OOA) model with 6 size change parameters,
223 4 migration rates, and 3 time parameters (Fig. 5A). In general, we observed a similar pattern
224 to previous models; size change parameters were often easier to infer than times or migration
225 rates (Fig. 5). For example, both donni and dadi showed near perfect inference accuracy for ν_{Af}
226 (Fig. 5B&G). They both also performed well for the for ν_{Eu} , ν_{As} and *misid* parameters (Fig. S8). But
227 several parameters were challenging for both methods, including some size change parameters,
228 such as ν_{As0} (Fig. 5C,H), ν_B and ν_{Eu0} (Fig. S8). The time parameters proved to be the most challenging
229 with relatively lower accuracy for both methods, with T_{Af} (Fig. 5D,I and T_B S8) being particularly
230 difficult. Overall, both methods agree on the parameters that are easy versus difficult to infer.

231 However, when inference accuracy is poor on difficult parameters, dadi and donni tend to
232 have different failure patterns. For instance with the m_{AfB} parameter, dadi tended to get stuck
233 at the parameter boundaries for many AFS (Fig. 5J), while donni essentially inferred the average
234 value for all test AFS (Fig. 5E). This indicates a failure by donni to learn any information from the
235 training AFS for this particular parameter. For all other migration rate parameters in the model,
236 donni performs well, matching dadi (Fig. S8).

237 While performance varied between the two methods among parameters, donni still had com-
238 parable accuracy to dadi in most cases. donni was also able to produce well calibrated confidence
239 intervals for all parameters (Fig. 5F). Due to the computational expense of dadi optimization for this
240 model, we only analyzed 30 test AFS for direct comparison between donni and dadi. Since donni
241 is not as computationally constrained, we also tested donni on all 1000 test AFS per our standard
242 procedure, finding similar results (Table S1).

243 Finally, we investigated the empirical AFS data from (Gutenkunst et al. 2009) using donni's
244 trained MVE networks for the Out-of-Africa demographic model (S3). We found that donni's
245 estimates differ from dadi's to varying degrees across the parameters. The similarity in accuracy
246 pattern between donni and dadi in Fig. 5 and Fig. S8 does not translate to similar inference values
247 between the two approaches on these data. For example, donni and dadi have similarly high
248 accuracy patterns for ν_{As} but have very different estimates on the empirical AFS data ($\nu_{As} = 7.29$
249 for dadi and $\nu_{As} = 1.276$ for donni). For this model, donni also tends to infer a stronger migration

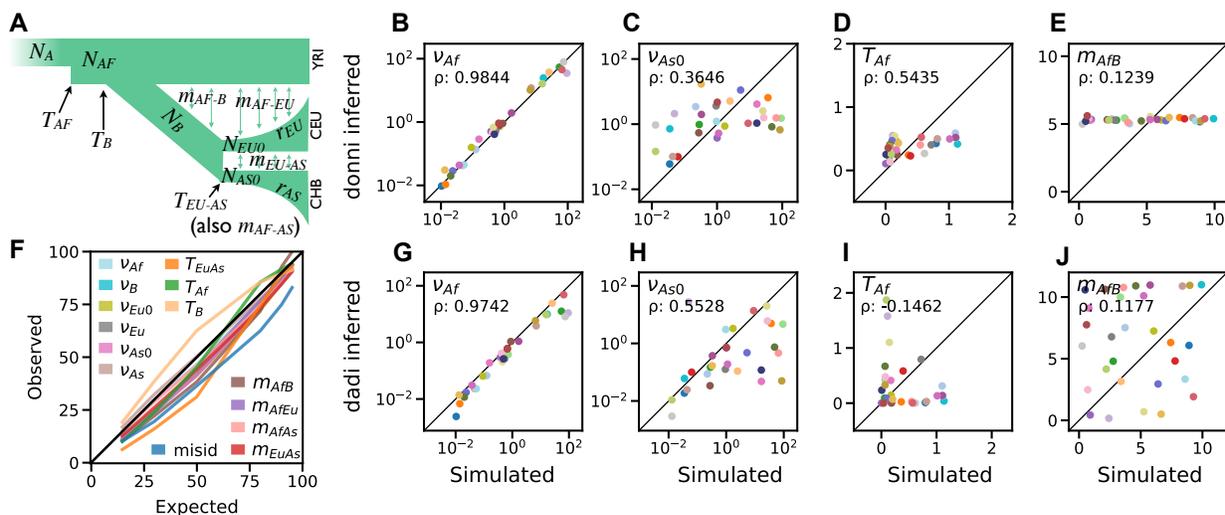


Figure 5: **Inference accuracy compared with dadi and confidence interval coverage by donni for the Out-of-Africa demographic model.** (A) The three-population Out-of-Africa model with 14 demographic history parameters. (B-E) Inference accuracy for representative parameters on 30 simulated test AFS inferred by donni. (G-J) Inference accuracy for the same parameters and 30 test AFS inferred by dadi. Each of the 30 test AFS is represented by a different color dot. For the accuracy of the rest of the parameters see Fig. S8. (F) donni confidence interval coverage for all model parameters.

250 rate than dadi does, with a higher estimate across all four migration rate parameters. Despite these
 251 differences in the estimated parameter values, dadi’s estimates are within donni’s 95% confidence
 252 intervals for all parameters.

253 *donni’s trained networks are accessible*

254 Given its speed, we expect that donni will be useful for quickly exploring many demographic
 255 scenarios given a user’s data set. To support this, we have produced trained networks for a large
 256 collection demographic history models. These include five one-population and eight two-population
 257 models from the current dadi API, plus the 34 two-population and 33 three-population models
 258 from Portik et al. (2017). For each of these models, we provide trained networks for unfolded
 259 and folded AFS for each of five sample sizes (only two sample sizes for three-population models).
 260 For large-scale production, we developed a comprehensive command-line interface pipeline for
 261 generating training data, tuning hyperparameters, and assessing the quality of the trained networks.
 262 donni’s pipeline is open-source and available on GitHub (<https://github.com/lntran26/donni>) for

263 users interested in training custom models. The trained network library is available on CyVerse
264 (Merchant et al. 2016; Center 2011) and donni's command-line interface will automatically download
265 appropriate networks. The library also included all accuracy and confidence interval coverage plots
266 for all supported demographic history models.

267 **DISCUSSION**

268 We addressed dadi's computationally intensive optimization procedure by developing donni, a
269 new inference method based on a supervised machine learning algorithm, the MVE network. We
270 found that donni's trained MVE networks can instantaneously infer many demographic history
271 parameters with accuracy comparable to dadi on simulated data. Even when including comput-
272 ing time required for training the network networks, for many cases donni is faster than dadi's
273 maximum likelihood optimization. Users are also provided a confidence interval for each inferred
274 demographic history model parameter value from donni. Through examples of one-, two-, and
275 three-population demographic models, we demonstrated that donni's uncertainty quantification
276 method works well for a wide range of demographic parameters. We also showed that donni works
277 well for AFS simulated by msprime, which includes linkage.

278 Our approach of using supervised machine learning to reduce the computational expense of
279 the maximum likelihood optimization step is similar in spirit to Smith et al. (2017) using random
280 forests to improve the efficiency of the computationally intensive ABC procedure. While Smith
281 et al. (2017) developed a classification approach for demographic model selection, our method is a
282 regression approach, where we provide a suite of pre-trained regressors for many commonly used
283 demographic history models. Users can quickly explore many possible scenarios and get an estimate
284 for several demographic parameters based on their input AFS data. However, we caution users to
285 always start with simpler models first before trying more complex ones, to avoid exacerbating the
286 uninformative parameter space problem. While we have implemented an accompanying uncertainty
287 quantification tool to aid in identifying such problematic scenarios, best practices in model-based
288 inference should still be followed.

289 Our choice of AFS as input data for training the network algorithm has several limitations.

290 First, because the size of the AFS depends on the sample size but the network requires a fixed
291 input size, we have to train a different set of networks for different sample sizes within the same
292 demographic history model. Different sets of trained networks are also required for unfolded
293 versus folded AFS. We have limited our trained network library to sample sizes of 10, 20, 40, 80,
294 and 160 haplotypes per population. User data that don't match exactly these sample sizes will be
295 automatically down-projected (Marth et al. 2004) by donni to the closest available option, leading
296 to some data loss. It is, however, possible to use donni's pipeline to train custom models that can
297 support a different sample size. We also verified that donni still provides accurate inference and
298 well-calibrated confidence intervals on down-projected data (Fig. S9).

299 Second, for optimal network performance, we need to normalize the AFS data for training,
300 leading to the loss of information about the parameter $\theta = 4N_a\mu L$, where N_a is the ancestral effective
301 population size, μ is the mutation rate and L is the sequence length. Estimating θ is required for
302 converting all demographic parameters in genetic units to absolute population sizes and divergence
303 times. While donni can provide a point estimate for θ , it cannot provide the uncertainty, which
304 is necessary for estimating the uncertainty of absolute parameter values. This limitation can be
305 overcome with a hybrid approach between donni and dadi, where donni's inferred parameter
306 outputs become the starting point for dadi's optimization procedure and uncertainty estimation
307 (Coffman et al. 2016). While this approach requires running likelihood optimization, a good starting
308 value provided by donni should reduce overall computing time.

309 donni trains a separate MVE network for each parameter in a given demographic history
310 model, even though the model parameters are correlated. This is a limitation of our implementation,
311 because the canonical MVE network architecture includes only one node for the mean and one
312 node for the variance. It may be possible to add additional nodes to output means, variances,
313 and covariances from a single network, but we found that this often affects the overall inference
314 quality of the trained MVE network. Additionally, we tested an alternative multi-output regression
315 approach (the scikit-learn Multilayer-Perceptron Regressor) and found that our single-output
316 approach provided similarly accurate estimates (Fig. S10). To our knowledge, existing methods for
317 estimating uncertainties of multi-output neural network regressions are limited.

318 At its heart, the neural network approach of donni corresponds to a nonlinear regression of
319 model parameters on AFS entries, in contrast to existing approaches which typically maximize a
320 composite likelihood through optimization. Neural networks can be used to estimate likelihoods
321 (e.g., Tejero-Cantero et al. (2020)), which could then be optimized or sampled over, but here we
322 prefer the more direct regression approach. Although dadi and donni display comparable overall
323 accuracy (Fig. 2&5), they may differ when applied to any given data set (Table S3), reflecting
324 differences between regression and composite likelihood optimization.

325 In conclusion, our results indicate that using supervised machine learning algorithms trained
326 with AFS data is a computationally efficient approach for inferring demographic history from
327 genomic data. Despite implementation limitations discussed above, the AFS is fast to simulate
328 compared to other types of simulated data such as genomic sequence images (Flagel et al. 2019;
329 Sanchez et al. 2021) or coalescent trees (Baumdicker et al. 2022; Kelleher et al. 2016). Furthermore,
330 while ignoring linkage may be a weakness of AFS-based methods, it can also be a strength in that
331 it is more species-agnostic and therefore trained models are transferable among species. A major
332 challenge for AFS-based methods such as ours is the poor scaling to large sample sizes and number
333 of populations, where the AFS matrix becomes high dimensional and sparse, and simulation
334 becomes prohibitively expensive. While we limited this study to three-population models, there
335 have been major improvements in AFS-based methods that can handle more (Gutenkunst 2021;
336 Jouganous et al. 2017; Kamm et al. 2017, 2020). Given our results, a supervised machine learning
337 approach might be a promising next step to extend to such AFS-based methods to further improve
338 their computational efficiency.

339 **MATERIALS AND METHODS**

340 *Simulations with dadi*

341 We used dadi v.2.3.0 (Gutenkunst et al. 2009) to simulate AFS for training and testing the networks.
342 For each demographic model, we uniformly drew parameter sets from a biologically relevant range
343 of parameters (Table S2). We then generated each expected AFS by specifying the demographic
344 model and parameters in dadi. We calculated the extrapolation grid points used for dadi inte-

345 gration based on the number of haplotypes per population according to Gutenkunst (2021) for
346 one-population models. For models with more than one population, we used the same formula but
347 also increased the grid points by a factor of 1.5 for each additional population. The demographic
348 model parameter values are used as labels for the generated AFS data. To simulate AFS with
349 different levels of variance, we started with the original expected AFS set (no variance). We then
350 Poisson-sampled from the expected AFS to generate a new AFS with variance. We controlled the
351 level of variance by the parameter θ , by which we multiplied the expected AFS before sampling. We
352 used $\theta = 10000, 1000$, and 100 corresponding to low, moderate, and high levels of variance, respec-
353 tively (Fig. S3-S7.) Intuitively, modifying $\theta = 4N_a\mu L$ is equivalent to altering the effective number
354 of sites surveyed L . Assuming $\mu \sim 10^{-8}$ and $N_a \sim 10^4$, $\theta = 1000$ is equivalent to $L \sim 2.5 \times 10^6$ sites.
355 Smaller θ is equivalent to fewer sites surveyed, hence noisier AFS. Finally, we normalized both
356 expected and Poisson-sampled AFS for training and testing. The results shown in Fig. 2,3,5, and S8
357 are based on unfolded AFS with sample size of 20 haplotypes per population.

358 *Simulations with msprime*

359 We used msprime v1.2.0 (Baumdicker et al. 2022) to simulate AFS from demographic history models
360 while including linkage. We first specified dadi-equivalent demography in msprime for the two
361 epoch and split-migration models. This included the population size change ratio ν and time of
362 change T parameters for the two epoch model, and population size change ratios ν_1 and ν_2 , time T
363 of split, and migration rate m for the split-migration model. We then specified additional parameters
364 required for msprime to yield $\theta = 4N_A L \mu = 40,000$, with ancestral population size $N_A = 10,000$,
365 sequence length $L = 10^8$ base pairs, and mutation rate $\mu = 10^{-8}$ per base pair per generation.
366 We used three recombination rates $10^{-8}, 10^{-9}$, and 10^{-10} per base pair per generation to simulate
367 different levels of linkage and variance in the AFS. We then generated tree-sequence data with
368 msprime before converting to the corresponding unfolded AFS of sample size 20 haplotypes per
369 population and normalizing for testing with trained networks.

370 *Network architecture and hyperparameter optimization*

371 We used TensorFlow v2.12.1 and Keras v2.12.0 to generate all trained MVE networks for donni.

372 These networks have two fully connected hidden layers containing between 4 and 64 neurons. The
373 exact number of neurons in each hidden layer are hyperparameters that were automatically selected
374 via our tuning procedure described below. The input layer is a flattened AFS with varying sizes
375 depending on the sample size and whether it is a folded or unfolded AFS. The output layer has two
376 nodes for the mean and variance of one demographic history parameter. For tuning and training the
377 network, we implemented a custom loss function based on the negative log-likelihood of a normal
378 distribution:

$$L(\theta) = \sum_{i=1}^N \frac{1}{2} \log(\sigma_{\theta}^2(x_i)) + \frac{1}{2} \frac{(y_i - \mu_{\theta}(x_i))^2}{\sigma_{\theta}^2(x_i)}$$

379 For automatic hyperparameter tuning, we used the HyperBand and RandomSearch tuning
380 algorithms available in keras-tuner v.1.4.6. The 5000 AFS training data set was split 80% for training
381 and 20% for validation. For a given network, we first used HyperBand to optimize both the
382 hidden layer size and learning rate. We then kept the MVE network from HyperBand with the
383 best performance on the validation data, froze the hidden layer size, and then continued tuning
384 only the learning rate using RandomSearch. The MVE network with the best performance on the
385 validation data after RandomSearch is then selected for subsequent training on the full training
386 data. All hyperparameter configurations and non-default settings for the tuning algorithms are
387 listed in Table S4.

388 *Uncertainty quantification coverage*

389 For uncertainty quantification, the trained MVE network outputs a variance for each inferred
390 demographic history parameter. donni pipeline converts this variance to confidence intervals using
391 the normal distribution. To validate our uncertainty quantification method, we first obtained the
392 method's estimation for six confidence intervals, 15, 30, 50, 80, 90, and 95% on all test AFS. We then
393 get the observed coverage by calculating the percentage of test AFS that have their corresponding
394 simulated parameter value captured within the estimated interval. The expected versus observed
395 percentages are plotted in our confidence interval coverage plots.

396 *donni training and testing pipeline*

397 We used 5,000 AFS (no variance) for training and tuning and 1,000 AFS (moderate variance,

398 $\theta = 1000$) for accuracy and uncertainty coverage validation. For visualization, only 100 test AFS
399 (30 AFS for the out-of-Africa model) are shown to compare with dadi. However, accuracy scores
400 by donni on all 1000 test AFS are provided in Table S1. Our pipeline tunes and trains one network
401 for each demographic model parameter and sample size. For example, the two epoch model with
402 two parameters ν and T has 20 independently trained networks: 2 networks for ν and T times 5
403 supported sample sizes times 2 polarization states.

404 *Likelihood optimization with dadi-cli*

405 To infer demographic parameters for a large number of test AFS in parallel (100 AFS for the
406 split-migration model and 30 AFS for the out-of-Africa model), we used dadi's command-line
407 interface (Huang 2023). We specified the upper and lower bound values for optimization based
408 on the parameter range provided in Table S2. Optimization ran until convergence, as defined by
409 $\delta \log(L) = 0.0005$ for the Out-of-Africa model and $\delta \log(L) = 0.001$ for the split-migration model.

410 *Benchmarking dadi optimization and donni pipeline*

411 To benchmark the computational expense required for dadi optimization versus for training the
412 networks, we used 10 CPUs on a single computing node for each task. For donni, the tasks are
413 generating training AFS, hyperparameter tuning with HyperBand, and training using the tuned hy-
414 perparameters. Estimating demographic parameters for 100 test AFS with donni's trained networks
415 is nearly instantaneous. For dadi, each test AFS is a task that was optimized until convergence, at
416 which time was recorded, or until the specified cut-off time (50 hours \times 10 CPUs = 500 CPU hours).

417 *Acknowledgments* This work was supported by the National Institute of General Medical Sciences
418 of the National Institutes of Health (R01GM127348 and R35GM149235 to R.N.G.).

419 * References

- 420 Achaz G (2009) Frequency spectrum neutrality tests: one for all and all for one. *Genetics* 183:249.
- 421 Baharian S, Gravel S (2018) On the decidability of population size histories from finite allele
422 frequency spectra. *Theoretical Population Biology* 120:42.
- 423 Baumdicker F, Bisschop G, Goldstein D, Gower G, Ragsdale AP, Tsambos G, Zhu S, Eldon B,
424 Ellerman EC, Galloway JG, et al. (2022) Efficient ancestry and mutation simulation with msprime
425 1.0. *Genetics* 220:iyab229.
- 426 Bergström A, McCarthy SA, Hui R, Almarri MA, Ayub Q, Danecek P, Chen Y, Felkel S, Hallast
427 P, Kamm J, et al. (2020) Insights into human genetic variation and population history from 929
428 diverse genomes. *Science* 367:eaay5012.
- 429 Bhaskar A, Song YS (2014) Descartes' rule of signs and the identifiability of population demographic
430 models from genomic variation data. *Annals of statistics* 42:2469.
- 431 Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams MD,
432 Schmidt S, Sninsky JJ, Sunyaev SR, et al. (2008) Assessing the evolutionary impact of amino acid
433 mutations in the human genome. *PLoS genetics* 4:e1000083.
- 434 Center D (2011) The iplant collaborative: cyberinfrastructure for plant biology. Chardon, M, and
435 Vandewalle, P(1991) Acoustico-lateralis system Cyprinid Fishes .
- 436 Chavez DE, Gronau I, Hains T, Dikow RB, Frandsen PB, Figueiró HV, Garcez FS, Tchaicka L,
437 de Paula RC, Rodrigues FH, et al. (2022) Comparative genomics uncovers the evolutionary
438 history, demography, and molecular adaptations of south american canids. *Proceedings of the*
439 *National Academy of Sciences* 119:e2205986119.
- 440 Coffman AJ, Hsieh PH, Gravel S, Gutenkunst RN (2016) Computationally efficient composite
441 likelihood statistics for demographic inference. *Molecular biology and evolution* 33:591.
- 442 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M (2013) Robust demographic inference
443 from genomic and snp data. *PLoS genetics* 9:e1003905.
- 444 Flagel L, Brandvain Y, Schrider DR (2019) The unreasonable effectiveness of convolutional neural
445 networks in population genetic inference. *Molecular biology and evolution* 36:220.
- 446 Gopalan S, Berl RE, Myrick JW, Garfield ZH, Reynolds AW, Bafens BK, Belbin G, Mastoras M,
447 Williams C, Daya M, et al. (2022) Hunter-gatherer genomes reveal diverse demographic trajec-
448 tories during the rise of farming in eastern africa. *Current Biology* 32:1852.
- 449 Gutenkunst RN (2021) Dadi. cuda: accelerating population genetics inference with graphics pro-
450 cessing units. *Molecular biology and evolution* 38:2177.
- 451 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD (2009) Inferring the joint de-
452 mographic history of multiple populations from multidimensional snp frequency data. *PLoS*
453 *genetics* 5:e1000695.

- 454 Hernandez RD, Williamson SH, Bustamante CD (2007) Context Dependence, Ancestral Misidentifi-
455 cation, and Spurious Signatures of Natural Selection. *Molecular Biology and Evolution* 24:1792.
- 456 Hey J, Chung Y, Sethuraman A, Lachance J, Tishkoff S, Sousa VC, Wang Y (2018) Phylogeny
457 Estimation by Integration over Isolation with Migration Models. *Molecular Biology and Evolution*
458 35:2805.
- 459 Huang X (2023) dadi-cli. <https://github.com/xin-huang/dadi-cli>.
- 460 Johnston HR, Cutler DJ (2012) Population demographic history can cause the appearance of recom-
461 bination hotspots. *The American Journal of Human Genetics* 90:774.
- 462 Jouganous J, Long W, Ragsdale AP, Gravel S (2017) Inferring the joint demographic history of
463 multiple populations: beyond the diffusion approximation. *Genetics* 206:1549.
- 464 Kamm J, Terhorst J, Durbin R, Song YS (2020) Efficiently inferring the demographic history of many
465 populations with allele count data. *Journal of the American Statistical Association* 115:1472.
- 466 Kamm JA, Terhorst J, Song YS (2017) Efficient computation of the joint sample frequency spectra for
467 multiple populations. *Journal of Computational and Graphical Statistics* 26:182.
- 468 Kelleher J, Etheridge AM, McVean G (2016) Efficient coalescent simulation and genealogical analysis
469 for large sample sizes. *PLoS computational biology* 12:e1004842.
- 470 Kern AD, Hey J (2017) Exact calculation of the joint allele frequency spectrum for isolation with
471 migration models. *Genetics* 207:241.
- 472 Khosravi A, Nahavandi S, Creighton D, Atiya AF (2011) Comprehensive review of neural network-
473 based prediction intervals and new advances. *IEEE Transactions on neural networks* 22:1341.
- 474 Kim BY, Huber CD, Lohmueller KE (2017) Inference of the distribution of selection coefficients for
475 new nonsynonymous mutations using large samples. *Genetics* 206:345.
- 476 Lorente-Galdos B, Lao O, Serra-Vidal G, Santpere G, Kuderna LF, Arauna LR, Fadhlouli-Zid K,
477 Pimenoff VN, Soodyall H, Zalloua P, et al. (2019) Whole-genome sequence analysis of a pan
478 african set of samples reveals archaic gene flow from an extinct basal population of modern
479 humans into sub-saharan populations. *Genome biology* 20:1.
- 480 Lukić S, Hey J (2012) Demographic inference using spectral methods on snp data, with an analysis
481 of the human out-of-africa expansion. *Genetics* 192:619.
- 482 Marchi N, Schlichta F, Excoffier L (2021) Demographic inference. *Current Biology* 31:R276.
- 483 Marchi N, Winkelbach L, Schulz I, Brami M, Hofmanová Z, Blöcher J, Reyna-Blanco CS, Diekmann
484 Y, Thiéry A, Kapopoulou A, et al. (2022) The genomic origins of the world's first farmers. *Cell* .
- 485 Marth GT, Czabarka E, Murvai J, Sherry ST (2004) The allele frequency spectrum in genome-wide
486 human variation data reveals signals of differential demographic history in three large world
487 populations. *Genetics* 166:351.
- 488 Mathieson I, McVean G (2012) Differential confounding of rare and common variants in spatially
489 structured populations. *Nature genetics* 44:243.

- 490 Mays Jr HL, Hung CM, Shaner PJ, Denvir J, Justice M, Yang SF, Roth TL, Oehler DA, Fan J, Rekula-
491 pally S, et al. (2018) Genomic analysis of demographic history and ecological niche modeling in
492 the endangered sumatran rhinoceros *dicerorhinus sumatrensis*. *Current Biology* 28:70.
- 493 Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, Antin P (2016) The iplant collaborative:
494 cyberinfrastructure for enabling data to discovery for the life sciences. *PLoS biology* 14:e1002342.
- 495 Miller-Butterworth CM, Diefenbach DR, Edson JE, Hansen LA, Jordan JD, Gingery TM, Russell AL
496 (2021) Demographic changes and loss of genetic diversity in two insular populations of bobcats
497 (*lynx rufus*). *Global Ecology and Conservation* 26:e01457.
- 498 Mondal M, Bertranpetit J, Lao O (2019) Approximate bayesian computation with deep learning
499 supports a third archaic introgression in asia and oecania. *Nature communications* 10:1.
- 500 Myers S, Fefferman C, Patterson N (2008) Can one learn history from the allelic spectrum? *Theoretic-
501 al population biology* 73:342.
- 502 Naduvilezhath L, Rose LE, Metzler D (2011) Jaatha: a fast composite-likelihood approach to estimate
503 demographic parameters. *Molecular Ecology* 20:2709.
- 504 Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, Bustamante C (2005) Genomic scans for
505 selective sweeps using snp data. *Genome research* 15:1566.
- 506 Nix DA, Weigend AS (1994) Estimating the mean and variance of the target probability distribution.
507 In *Proceedings of 1994 ieee international conference on neural networks (ICNN'94)*, volume 1, pages
508 55–60. IEEE.
- 509 Portik DM, Leaché AD, Rivera D, Barej MF, Burger M, Hirschfeld M, Rödel MO, Blackburn DC,
510 Fujita MK (2017) Evaluating mechanisms of diversification in a guineo-congolian tropical forest
511 frog using demographic model selection. *Molecular ecology* 26:5245.
- 512 Pudlo P, Marin JM, Estoup A, Cornuet JM, Gautier M, Robert CP (2016) Reliable abc model choice
513 via random forests. *Bioinformatics* 32:859.
- 514 Sanchez T, Cury J, Charpiat G, Jay F (2021) Deep learning for population size history inference:
515 Design, comparison and combination with approximate bayesian computation. *Molecular
516 Ecology Resources* 21:2645.
- 517 Sawyer SA, Hartl DL (1992) Population genetics of polymorphism and divergence. *Genetics*
518 132:1161.
- 519 Schrider DR, Kern AD (2018) Supervised machine learning for population genetics: a new paradigm.
520 *Trends in Genetics* 34:301.
- 521 Sheehan S, Song YS (2016) Deep learning for population genetic inference. *PLoS computational
522 biology* 12:e1004845.
- 523 Sluijterman L, Cator E, Heskes T (2023) Optimal training of mean variance estimation neural
524 networks. *arXiv preprint arXiv:230208875* .
- 525 Smith ML, Ruffley M, Espíndola A, Tank DC, Sullivan J, Carstens BC (2017) Demographic model
526 selection using random forests and the site frequency spectrum. *Molecular Ecology* 26:4562.

- 527 Spence JP, Steinrücken M, Terhorst J, Song YS (2018) Inference of population history using coalescent
528 hmms: review and outlook. *Current opinion in genetics & development* 53:70.
- 529 Tejero-Cantero A, Boelts J, Deistler M, Lueckmann JM, Durkan C, Gonçalves PJ, Greenberg DS,
530 Macke JH (2020) sbi: A toolkit for simulation-based inference. *Journal of Open Source Software*
531 5:2505.
- 532 Terhorst J, Song YS (2015) Fundamental limits on the accuracy of demographic inference based on
533 the sample frequency spectrum. *Proceedings of the National Academy of Sciences* 112:7677.
- 534 Villanea FA, Schraiber JG (2019) Multiple episodes of interbreeding between neanderthal and
535 modern humans. *Nature ecology & evolution* 3:39.

536 **Supporting Information**

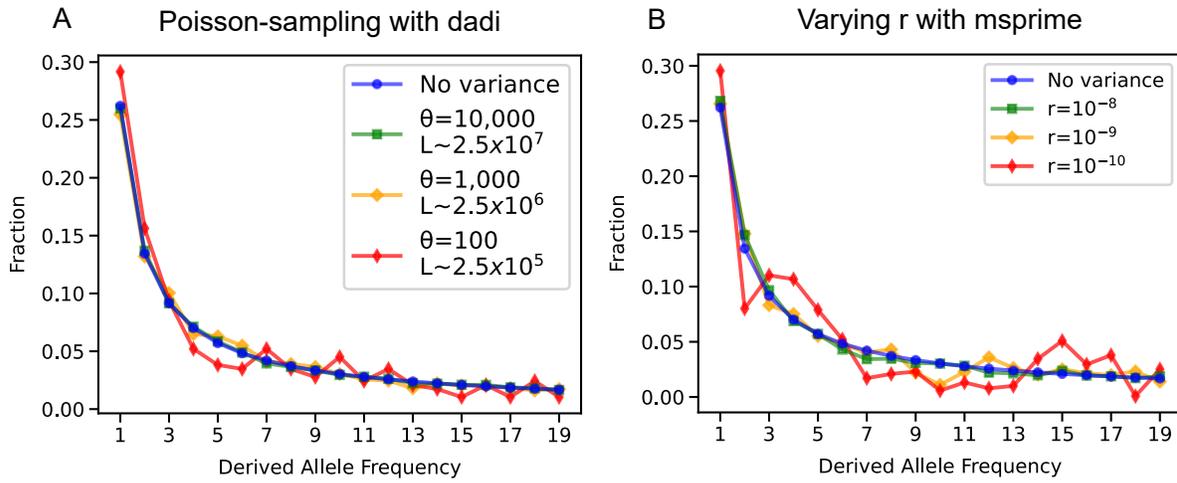


Figure S1: **Simulated AFS examples with different variance for the two epoch model.** All AFS are normalized and plotted on the same scale. The "No variance" line in both panels is the expected AFS generated by dadi with $\nu = 0.8$, $T = 0.5$ (A) AFS with different variance by Poisson-sampling from the "No variance" AFS. (B) msprime-simulated AFS with equivalent demography to (A) but with varying recombination rates. Here $\theta = 4N_a\mu L = 4 \times 10^3$ (with $\mu = 10^{-8}$ per nucleotide per generation and $L = 10^8$ base pairs).

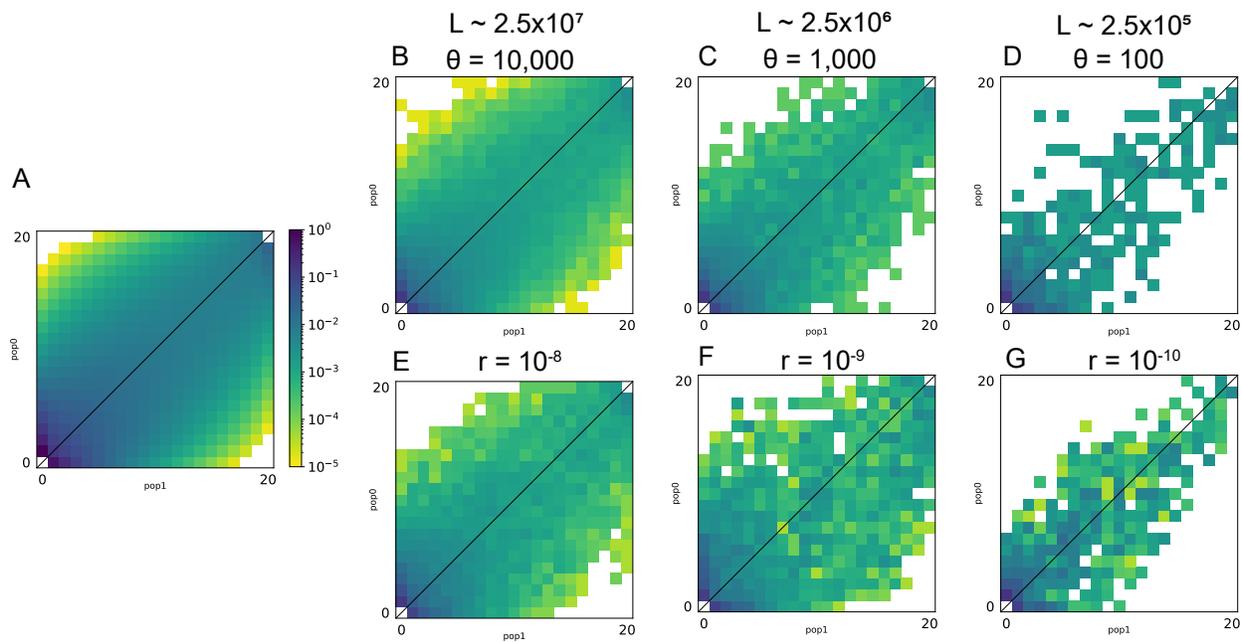


Figure S2: **Simulated AFS examples with different variance for the split-migration model.** All AFS are normalized and plotted on the same scale. (A) Expected AFS generated by dadi with $\nu_1 = 1$, $\nu_2 = 0.5$, $T = 2$, $m = 5$. (B-D) AFS with different variance by Poisson-sampling from (A). (E-G) msprime-simulated AFS with equivalent demography as in (A) under varying recombination rates. Entries below 10^{-5} are masked.

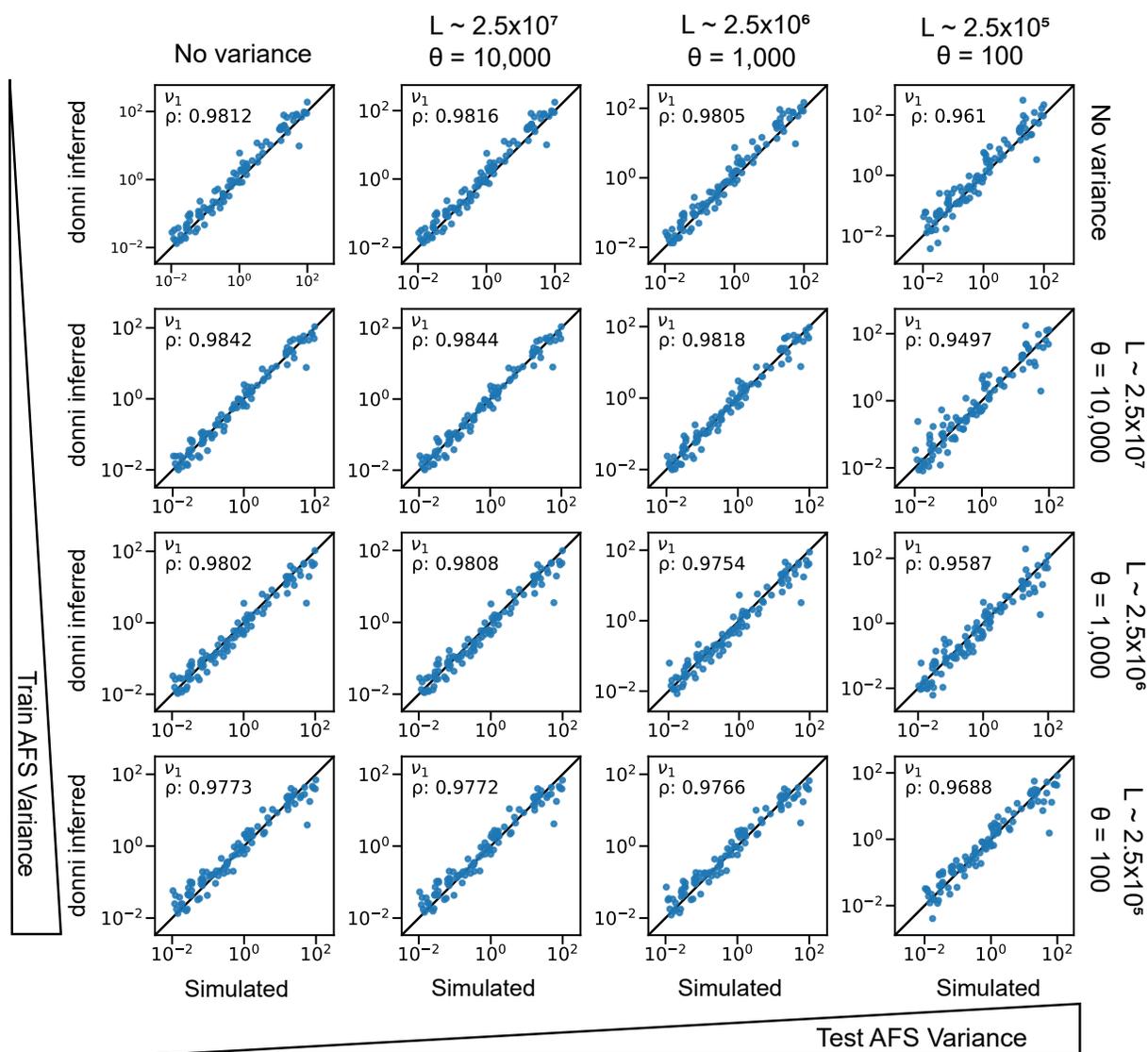


Figure S3: The effects of AFS variance on donni training and performance for the split-migration model size-change parameter v_1 . Each row corresponds to different levels of variance in training AFS, and each column corresponds to different levels of variance in test AFS. For example, the third panel from the left in the top row is the inference accuracy of a network trained on AFS with no variance tested on AFS with moderate levels of variance ($\theta = 1000$ or $L \sim 2.5 \times 10^6$ sites surveyed).

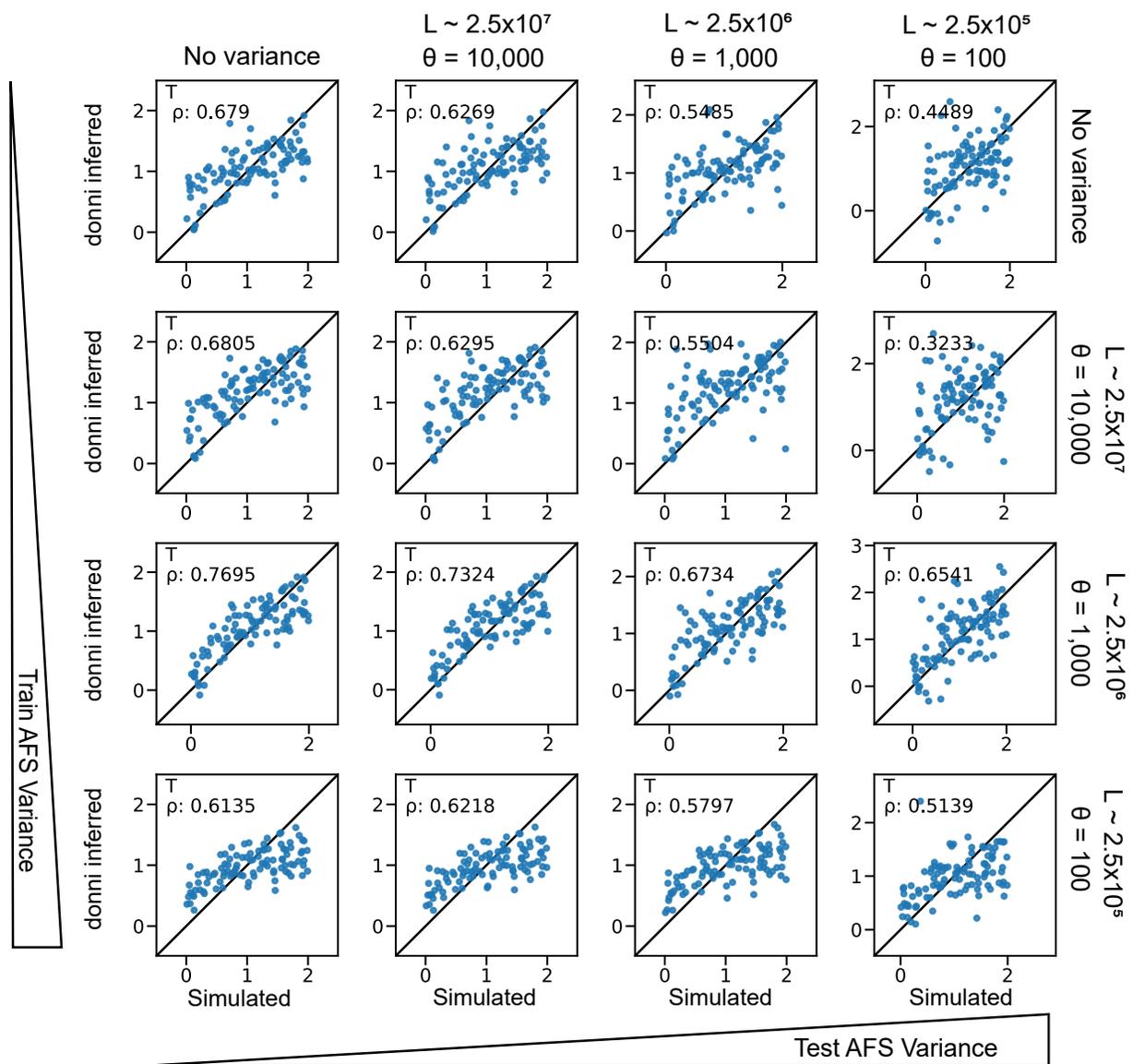


Figure S4: The effects of AFS variance on donni training and performance for the split-migration model time parameter T . Panels are as in Fig. S3.

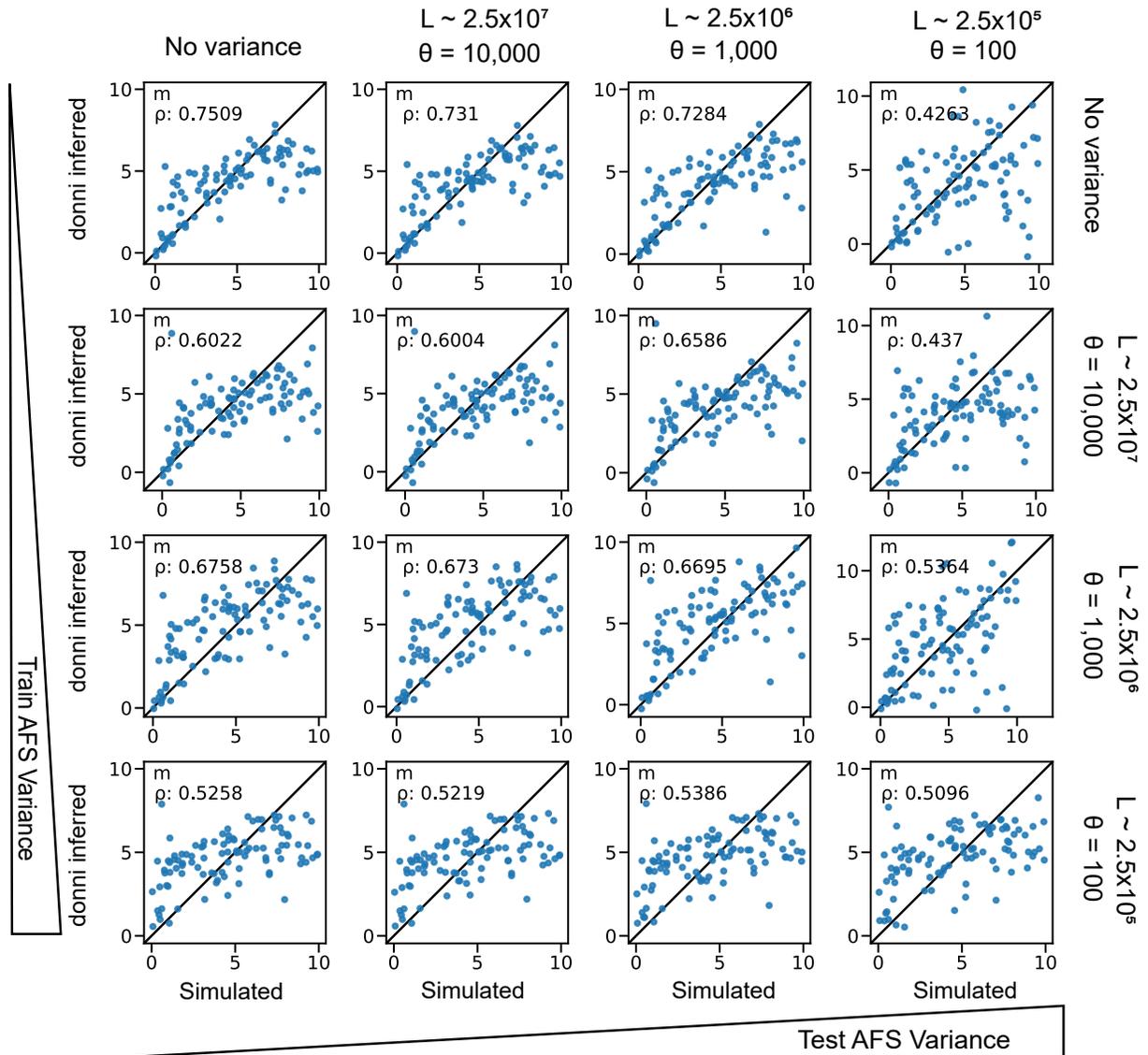


Figure S5: The effects of AFS variance on donni training and performance for the split-migration model migration rate parameter m . Panels are as in Fig. S3.

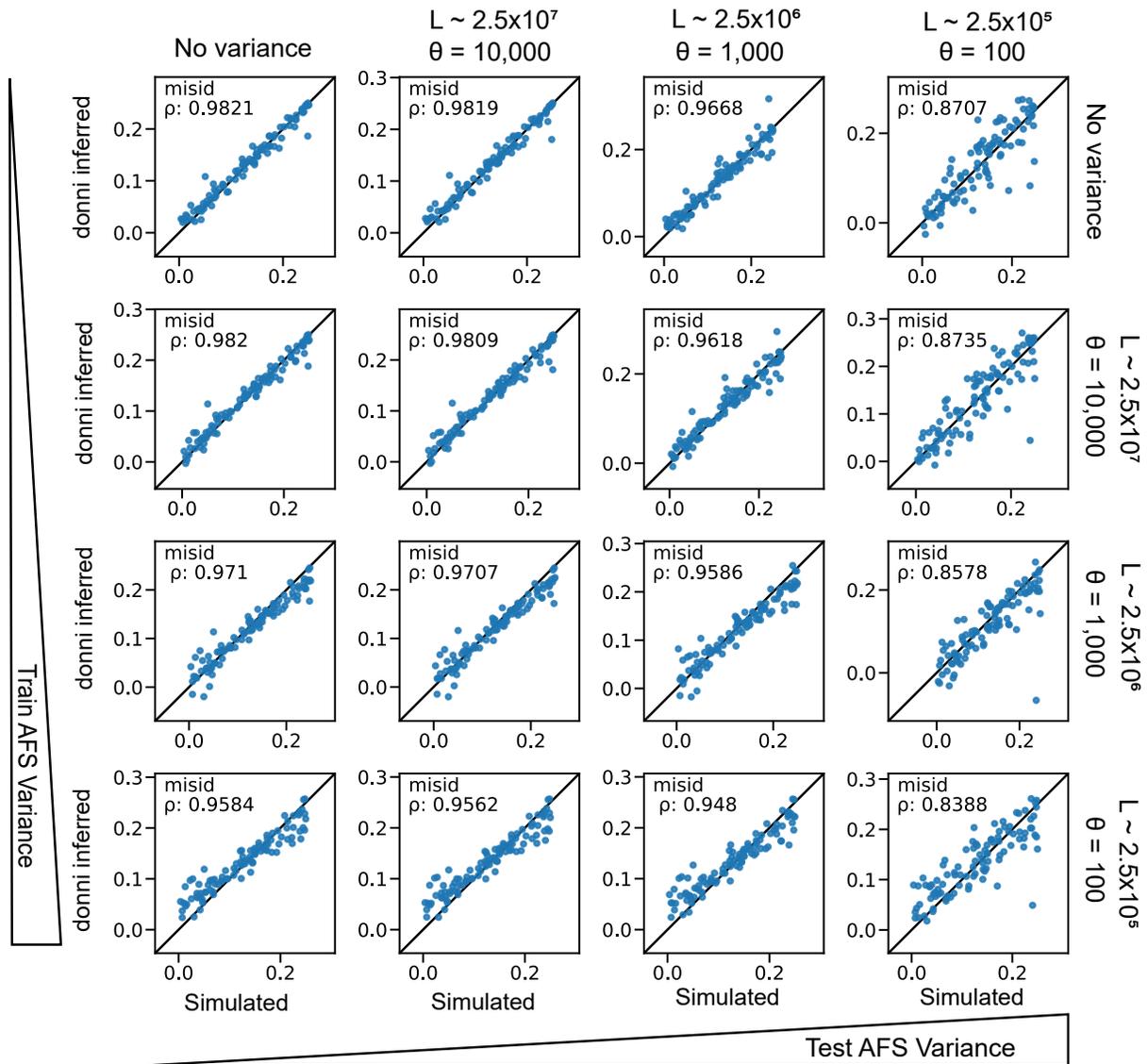


Figure S6: The effects of AFS variance on donni training and performance for the split-migration model ancestral state misidentification parameter. Panels are as in Fig. S3.

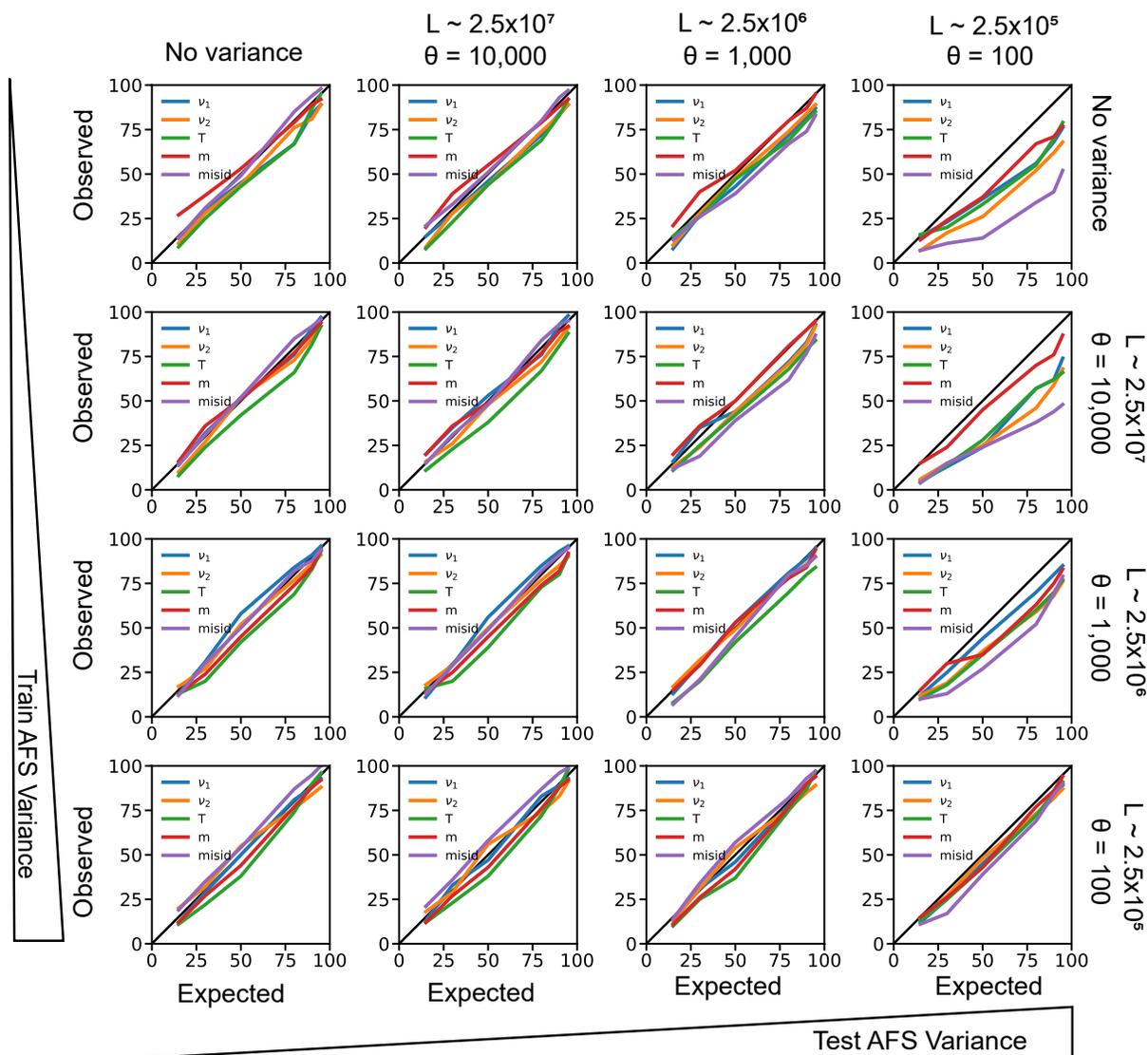


Figure S7: The effects of AFS variance on donni's uncertainty quantification method for the split-migration model. Panels are as in Fig. S3.

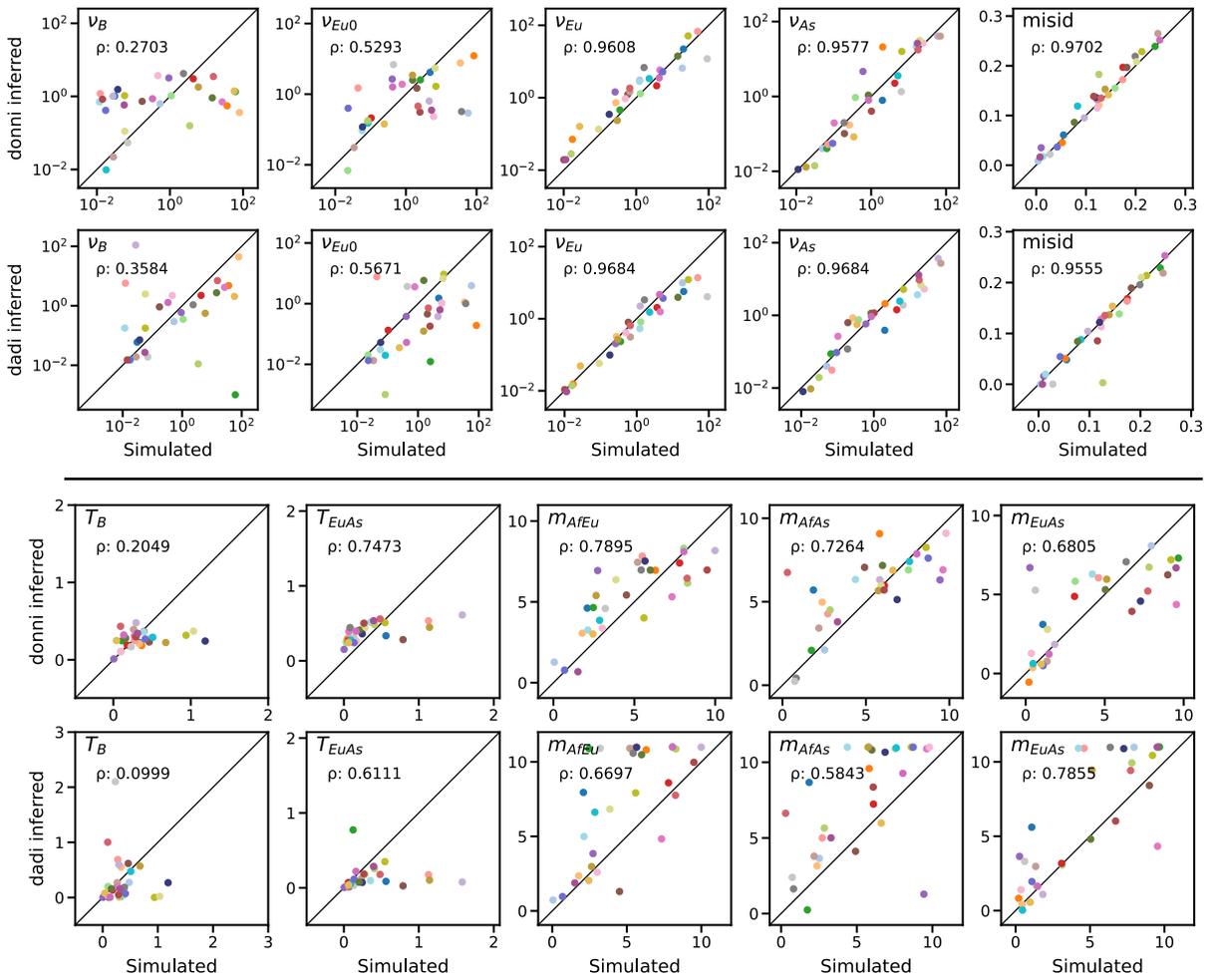


Figure S8: Inference accuracy of *dadi* and *donni* on the rest of Out-of-Africa model parameters. Each of the 30 test AFS is represented by a different color dot as in Fig. 5.

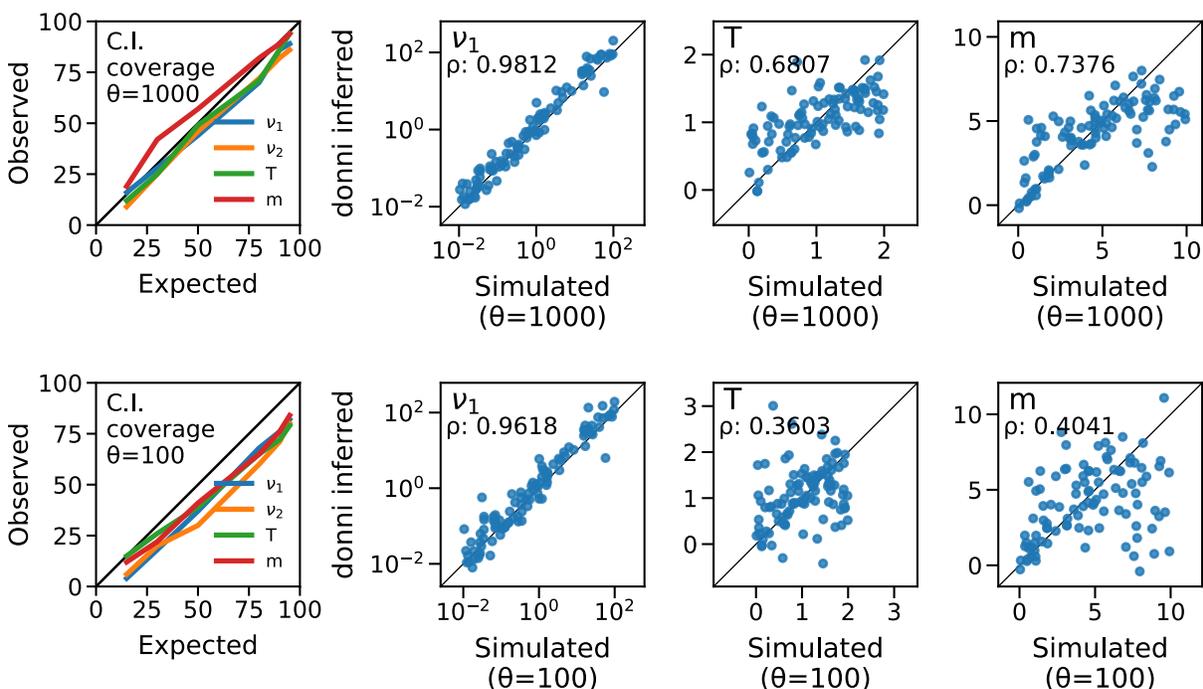


Figure S9: **Inference accuracy and confidence interval calibration by donni on down-projected test AFS for the split-migration model.** We simulated 100 test AFS with sample size 39 haplotypes per population then projected them to sample size 20 haplotypes per population. Top row is the result for test AFS projected from moderate variance ($\theta = 1000$) AFS and bottom row is for test AFS projected from high variance ($\theta = 100$) AFS.

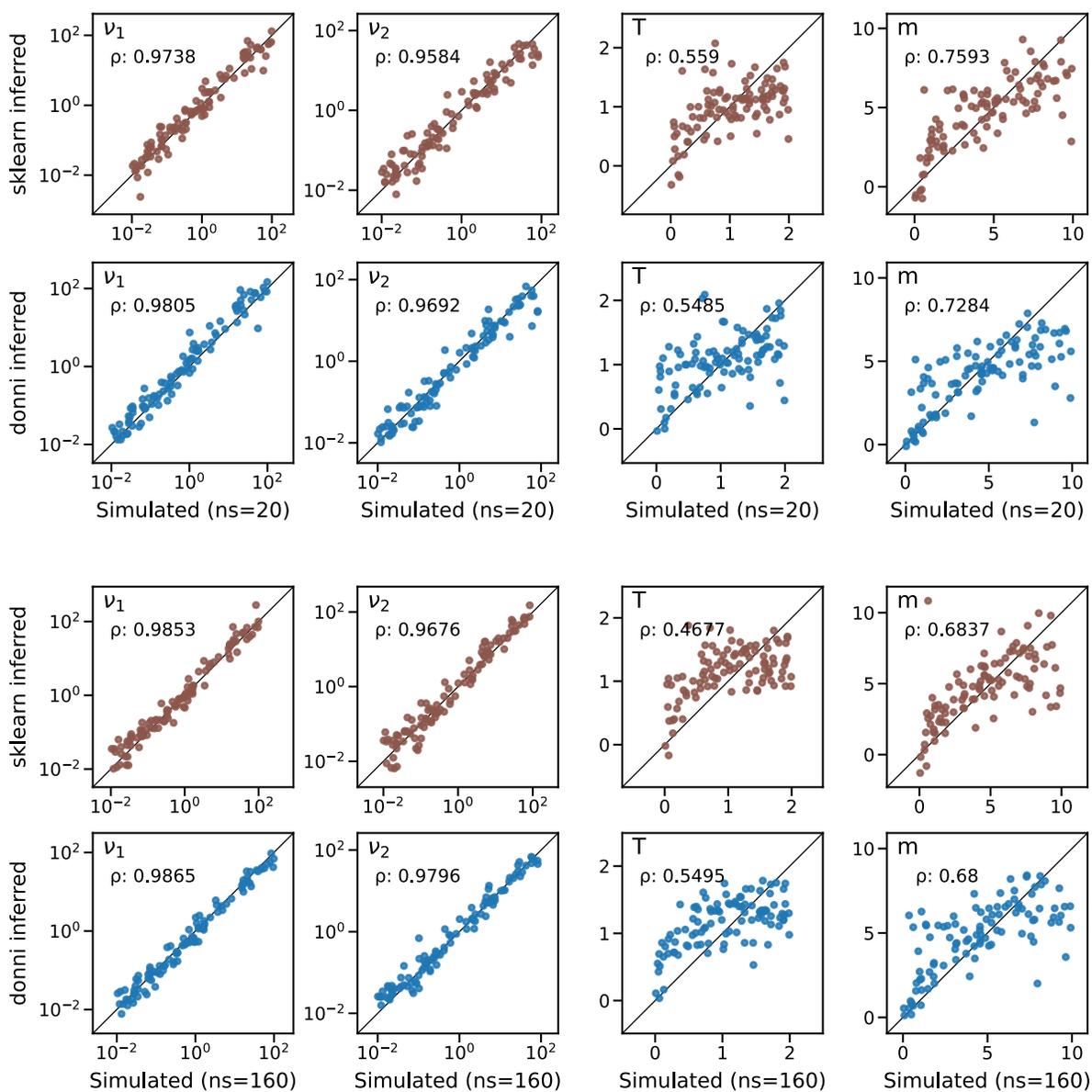


Figure S10: Inference accuracy by scikit-learn multi-output network compared with donni's single-output network for the split-migration model with sample sizes 20 and 160 haplotypes per population. scikit-learn multi-output network is one network network trained to predict all parameters in a demographic model, whereas donni trains a single network for each parameter. We used the same test AFS simulated with moderate variance ($\theta = 1000$) for sklearn and donni.

Table S1: donni inference accuracy on 1000 test AFS with moderate variance ($\theta = 1000$) and the best hidden layer sizes (HLS) architecture for demographic history models in this study.

Model	Parameter	RMSE	Spearman's ρ	network architecture (HLS)
1-population	ν	201.9	0.714	(64, 12)
Two epoch	T	0.72	0.514	(48, 9)
($ns = 20$ haplotypes)	$misid$	0.037	0.905	(64, 12)
2-population	ν_1	16.223	0.981	(48, 12)
Split Migration	ν_2	10.713	0.977	(64, 12)
($ns = 20$ haplotypes	T	0.444	0.65	(32, 16)
per population)	m	2.008	0.714	(48, 12)
	$misid$	0.019	0.968	(48, 8)
2-population	ν_1	7.616	0.987	(64, 16)
Split Migration	ν_2	5.828	0.985	(48, 12)
($ns = 160$ haplotypes	T	0.455	0.609	(48, 8)
per population)	m	1.927	0.736	(64, 16)
	$misid$	0.012	0.987	(64, 8)
3-population	ν_{Af}	11.504	0.985	(64,16)
Out of Africa	ν_B	23.7371	0.495	(64, 16)
($ns = 20$ haplotypes	ν_{Eu0}	22.609	0.451	(48, 8)
per population)	ν_{Eu}	17.67	0.921	(48, 16)
	ν_{As0}	24.06	0.405	(64, 16)
	ν_{As}	14.435	0.934	(64, 16)
	T_{Af}	0.299	0.435	(64, 4)
	T_B	0.322	0.366	(48, 12)
	T_{EuAs}	0.303	0.515	(64, 12)
	m_{AfB}	2.919	0.092	(64, 16)
	m_{AfEu}	2.038	0.733	(16, 16)
	m_{AfAs}	2.075	0.712	(64, 8)
	m_{EuAs}	2.171	0.676	(64, 12)
	$misid$	0.014	0.987	(48, 16)

Table S2: dadi demographic parameter range used for simulation in this study.

Parameter	Symbol	Lower bound	Upper bound
Population size change	ν	0.01	100
Time of event(*)	T	0.01	2
Migration rate	m	0	10
Ancestral state misidentification	$misid$	0	0.25

* For models with more than one T parameter, the range specified for time T applies to the sum of all T parameters (T_{sum}). For each demographic model, we drew different T_{sum} values according to the desired number of data sets. For each data set, we then drew a set of T parameters that sum to T_{sum} by sampling from the Dirichlet distribution.

Table S3: donni inferred compared to dadi inferred parameter values in genetic units for the Out-of-Africa model using data from (Gutenkunst et al. 2009).

Parameter	dadi	donni	donni 95% C.I.
θ	2788.2	2644.9*	n.a.
ν_{Af}	1.68	2.029	0.877 - 4.695
ν_B	0.287	0.192	0.011 - 3.500
ν_{Eu0}	0.129	0.145	0.004 - 5.627
ν_{Eu}	3.74	1.068	0.149 - 7.658
ν_{As0}	0.070	0.045	0.001 - 2.433
ν_{As}	7.29	1.276	0.200 - 8.135
m_{AfB}	3.65	5.089	-0.378 - 10.556
m_{AfEu}	0.44	1.673	-1.082 - 4.428
m_{AfAs}	0.28	0.373	-1.591 - 2.337
m_{EuAs}	1.40	4.871	-0.262 - 10.004
T_{Af}	0.607	0.432	-0.124 - 0.989
T_B	0.396	0.22	-0.310 - 0.749
T_{EuAs}	0.058	0.119	-0.084 - 0.321

* donni infers a slightly negative value for the probability of ancestral state misidentification ($misid = -0.00048$). These data were previously corrected for ancestral state misidentification using the approach of Hernandez et al. (2007). We thus rounded to $misid = 0$ when calculating θ .

Table S4: Hyperparameters tuned with KerasTuner for each demographic model parameter.

Tuner	Hyperparameter	Value range
Hyperband	First hidden layer	16, 32, 48, 64
	Second hidden layer	4, 8, 12, 16
	Learning rate	log sampling, [0.0001, 0.01]
	Max epochs	100
RandomSearch	Learning rate	log sampling, [0.0001, 0.01]
	Max trials	100