

SIA: Selection Inference Using the Ancestral Recombination Graph

Hussein A. Hejase^{1*}, Ziyi Mo^{1,2*}, Leonardo Campagna^{3,4}, Adam Siepel¹

¹Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

²School of Biological Sciences, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

³Fuller Evolutionary Biology Program, Cornell Lab of Ornithology, Ithaca, NY, USA

⁴Department of Ecology and Evolutionary Biology, Cornell University, Ithaca, NY, USA

*These authors contributed equally

Corresponding author: Adam Siepel (asiepel@cshl.edu)

Keywords

ancestral recombination graph, machine learning, positive selection, selective sweep

1 Abstract

2 Detecting signals of selection from genomic data is a central problem in population genetics.
3 Coupling the rich information in the ancestral recombination graph (ARG) with a powerful and
4 scalable deep learning framework, we developed a novel method to detect and quantify positive
5 selection: **Selection Inference using the Ancestral recombination graph (SIA)**. Built on a Long
6 Short-Term Memory (LSTM) architecture, a particular type of a Recurrent Neural Network (RNN),
7 SIA can be trained to explicitly infer a full range of selection coefficients, as well as the allele
8 frequency trajectory and time of selection onset. We benchmarked SIA extensively on simulations
9 under a European human demographic model, and found that it performs as well or better as
10 some of the best available methods, including state-of-the-art machine-learning and ARG-based
11 methods. In addition, we used SIA to estimate selection coefficients at several loci associated
12 with human phenotypes of interest. SIA detected novel signals of selection particular to the
13 European (CEU) population at the *MC1R* and *ABCC11* loci. In addition, it recapitulated signals of
14 selection at the *LCT* locus and several pigmentation-related genes. Finally, we reanalyzed
15 polymorphism data of a collection of recently radiated southern capuchino seedeater taxa in the
16 genus *Sporophila* to quantify the strength of selection and improved the power of our previous
17 methods to detect partial soft sweeps. Overall, SIA uses deep learning to leverage the ARG and
18 thereby provides new insight into how selective sweeps shape genomic diversity.
19

20 Introduction

21 The ability to accurately detect and quantify the influence of selection from genomic sequence
22 data enables a wide variety of insights, ranging from understanding historical evolutionary events
23 to characterizing the functional and disease relevance of observed or potential genetic variants.
24 Adaptive evolution is driven by increases in frequency of alleles that enhance reproductive fitness.
25 In addition, alleles experiencing such positive selection often provide insights into the functional
26 or mechanistic basis of phenotypes of interest. Examples of genetic determinants of important
27 phenotypic traits under selection in human populations include a family of mutations in the
28 hemoglobin- β cluster, which confer resistance to malaria and are at high frequencies in many
29 populations [1,2], loci controlling growth factor signaling pathways that contribute to short stature
30 in Western Central African hunter-gatherer populations [3,4], as well as mutations in several
31 genes involved in immunity, hair follicle development, and skin pigmentation [5] (reviewed in refs.
32 [6–9]).

33
34 Population genetic methods predominantly identify positive selection through the detection of
35 selective sweeps. As the frequency of an advantageous allele increases, linked variants in the
36 vicinity can “hitchhike” to high frequency, leading to local reductions in genetic diversity. Previous
37 approaches to detecting selective sweeps (such as traditional summary statistics [10],
38 approximate likelihood and Approximate Bayesian Computation (ABC) methods [11], or
39 supervised machine learning (ML) methods [12,13]) exploit the effect of genetic hitchhiking on the
40 spatial haplotype structure and site frequency spectrum (SFS). Summary statistics have the
41 advantage of being fast and easy to compute, but may confound the effects of selection on genetic
42 diversity with the effects of complex demographic histories including bottlenecks, population
43 expansions and structured populations. Besides, they cannot easily be used to estimate the value
44 of the selection coefficient. Approximate likelihood and ABC methods, on the other hand, can

45 provide an estimate of the strength of selection by aggregating multiple summary statistics [11],
46 but can be prohibitively computationally expensive when applied at a large scale. ML methods for
47 inferring selection can be more scalable, and can capture complex nonlinear relationships among
48 features. With the exception of a handful of recently developed methods that operate on the
49 multiple sequence alignment itself [14,15], however, the majority of ML approaches to selection
50 inference solely make use of traditional summary statistics as features for prediction. In short,
51 previous methods (including ABC and most ML methods) predominantly rely on low-dimensional
52 summary statistics, which, even in combination, capture only a small portion of the information in
53 the sequence data.

54

55 Recently, a new generation of inference methods have made it possible to go beyond summary
56 statistics and estimate or sample a full ancestral recombination graph (ARG) [16–18] for a
57 collection of sequences of interest. The ARG is a complex data structure that summarizes the
58 shared evolutionary history and recombination events that have occurred in a collection of DNA
59 sequences, and therefore contains highly informative features that can potentially be leveraged
60 to make accurate inferences about selection. The ARG representation is interchangeable with a
61 sequence of local genealogies along the genome and the recombination events that transform
62 each genealogy to the next. The influence of selection on each allele can be characterized from
63 the ARG, based on departures from the patterns of coalescence and recombination expected
64 under neutrality as reflected in the local genealogies. Traditional ARG inference methods [19–23]
65 were restricted in accuracy and scalability, limiting the practical application of ARGs. Recent
66 advances [24], however, have enabled scalable yet statistically rigorous genome-wide ARG
67 inference with dozens of genomes. Moreover, methods such as Relate [25] and tsinfer [26] have
68 further dramatically improved the scalability of ARG inference to accommodate thousands or even
69 hundreds of thousands of genomes. The latest progress in genealogical inference has paved the
70 way for ARG-based methods to address many different questions in population genetics [24–27].

71

72 One natural way to exploit the richness of the ARG representation in inference of selection would
73 be to extract features from inferred ARGs and feed them into a modern supervised machine-
74 learning framework. Deep-learning methods, in particular, have recently achieved unprecedented
75 success on a variety of challenging problems, including image recognition, machine translation,
76 and game-play [28]. Deep learning is also highly flexible, providing many opportunities for the
77 design of novel model architectures motivated by biological knowledge. An ARG-guided deep-
78 learning model could potentially provide new insight into how natural selection impacts the human
79 genome, human diseases and other phenotypes, and human evolution.

80

81 With these goals in mind, we developed a new method, called SIA (**S**election **I**nference using the
82 **A**ncestral recombination graph), that uses a Recurrent Neural Network (RNN) [29,30] to infer the
83 selection coefficient and allele frequency trajectory of a variant that maps to a gene tree
84 embedded in an ARG. Rather than relying on traditional sequence-based summary statistics, SIA
85 makes use of features based on the local genealogies extracted from the ARG. Based on these
86 local topological features, SIA learns to infer the selection coefficient and allele frequency
87 trajectory of a beneficial variant (see **Figure 1**). As described below, SIA performs well on
88 benchmarks and is reasonably robust to model misspecification. Applying SIA to data from the
89 1000 Genomes Northern and Western European (CEU) population, we identified new and known
90 loci under positive selection that are associated with a variety of phenotypes and estimated
91 selection coefficients at these loci. In addition, using SIA, we built on our previous work [31] on a
92 bird species-complex in the genus *Sporophila* by elucidating the strength and targets of selection
93 at specific loci tied to a collection of rapid speciation events. Overall, SIA is the first method that
94 couples ARG-based features with a machine-learning approach for population genetic inference.

95

96 Results

97 *Methodological overview.* SIA is based on an RNN that is trained to predict selection at a genomic
98 site from genealogical features at that site of interest and nearby sites (see **Methods** for detailed
99 descriptions, see **Figure 1** for a conceptual overview of SIA, and **Figure S1** for an illustration of
100 ARG features and the RNN architecture). Based on the demography of a particular population of
101 interest, training data including genomic regions under various strengths of selection are
102 simulated. The ARG is then inferred from each simulated data set. ARG-level statistics are
103 extracted at the site under selection (or a neutral site) as features to be used as input to the deep-
104 learning model. Specifically, we use lineage counts at a set of discrete time points as a fixed-
105 dimension encoding of a genealogy. The encoding of the genealogy at the focal site as well as
106 similar encodings of flanking genealogies constitute the feature vector for that site. SIA uses a
107 Long Short-Term Memory (LSTM) architecture, designed specifically to handle the temporal
108 nature of the feature set. The LSTM unrolls temporally such that the lineage counts at each time
109 point are fed to the network iteratively. Finally, the model trained on simulations is applied to
110 ARGs inferred from empirical data to identify sweeps, infer selection coefficients, and allele-
111 frequency trajectories.

112
113 *Classification of sweeps.* We first compared SIA with several existing methods, including the
114 Tajima's D [10] and H1 [32] summary statistics, iHS [33], a genealogy-based statistic [25] and a
115 summary-statistic-based machine-learning method [12,13] (see **Methods**), in the classification
116 task of distinguishing hard sweeps from neutrally evolving regions. Our performance comparison
117 was conducted across 16 combinations of selection coefficients and segregating allele
118 frequencies such that the beneficial site was subjected to selection ranging from weak to strong,
119 resulting in low to high derived allele frequencies (DAFs). Since *a priori* we expected sweep sites
120 with lower selection coefficients and lower DAFs to be harder to detect, we performed a stratified

121 analysis of SIA's performance by selection coefficient and DAF. **Figure 2** reports the Receiver
122 Operating Characteristic (ROC) curves using simulations based on the CEU demographic model
123 [34] where inferred genealogies were used as input to SIA to account for gene tree uncertainty.
124 As expected, all methods tended to perform better in a regime with higher selection coefficients
125 and DAFs, as indicated by increasing values of the area under the ROC curve (AUROC) statistic
126 from left to right (increasing selection) and from top to bottom (increasing DAF). SIA outperformed
127 the other methods across model conditions, with a more pronounced performance advantage for
128 sites under weaker selection and segregating at lower DAFs (**Figure 2**). For each given selection
129 coefficient, the AUROC of the Relate tree statistic (shown in red in **Figure 2**), which measures
130 how unlikely it is that the observed expansion of the derived lineages is purely due to genetic drift,
131 did not substantially improve as the DAF increased. Alleles at higher frequency tend to be older
132 and subjected to drift over longer periods, which may lead to reduced power for Relate to
133 distinguish lineage expansion under selection from the neutral expectation. Consequently, while
134 the ARG-based methods SIA and Relate both outperformed other methods at low DAFs, SIA was
135 alone in maintaining this advantage at higher DAFs.

136

137 In addition, we validated the ability of SIA to classify genomic regions with additional test sets
138 simulated under a demographic model for southern capuchinos, a group of songbirds in which we
139 previously identified and characterized many examples of sweeps [31], finding a predominance
140 of “soft” rather than “hard” sweeps (meaning that they tend to be based on standing genetic
141 variation rather than new mutations; see **Methods**). **Figure S2** reports the ROC curves for the
142 task of distinguishing partial soft sweeps from neutral regions. Despite soft sweeps being harder
143 to detect, the classifier achieved good performance in the moderate-to-strong selection regimes
144 ($s = 0.005$ and $s = 0.0075$) where the accuracy ranged between 82% and 96%, a substantial
145 improvement over the previous accuracy of 56% [31]. SIA performed particularly well in identifying
146 partial soft sweeps when the site under selection was at a high segregating frequency. For

147 example, at segregating frequencies of 0.75 and 0.9, the performance of SIA ranged between
148 80% and 96% across a variety of selection regimes ($s = 0.0025, 0.005, \text{ and } 0.0075$). The
149 performance of SIA degraded somewhat for weak selection ($s = 0.001$) with an accuracy ranging
150 between 63% and 74%.

151
152 *Selection coefficient inference using true gene trees.* We assessed the performance of SIA in
153 correctly predicting the selection coefficient and compared it to CLUES [35]. Like SIA, CLUES
154 uses local genealogies based on the ARG to infer a selection coefficient. However, CLUES
155 calculates the likelihood of the genealogy analytically using a hidden Markov model (HMM), and
156 does not rely on simulated training data. In addition, CLUES uses a single genealogy at the focal
157 site, whereas SIA additionally considers flanking trees.

158
159 We began by supplying both methods with true genealogies, in order to later disentangle the error
160 deriving from the ARG inference step from other sources of error (see **Discussion**). We found
161 that SIA identified regions under neutrality with approximately no bias (median inferred $s = 7.5e-$
162 05 ; **Figure 3**). Similarly, SIA correctly inferred the selection coefficient for regions under moderate
163 to strong selection ($s \in \{0.0025, 0.005, 0.0075, 0.01\}$) with the median inferred s deviated from
164 the true s by at most 3%. On the other hand, SIA somewhat underestimated the selection
165 coefficient (median inferred $s = 0.00037$) for the weak selection regime (true $s = 0.001$), likely
166 owing to limits in the training set within that selection regime (see **Discussion**). We further binned
167 the results by segregating frequency and selection coefficient and found that, in general, the
168 variance in estimates of s for SIA (as well as CLUES) tended to decrease as the segregating
169 frequency of the beneficial allele increased (**Figure S3**).

170
171 CLUES performed roughly similarly to SIA in this experiment, but tended to slightly overestimate
172 s for the neutral regions (i.e., true $s = 0$) and underestimate s for the moderate to high selection

173 regimes (i.e., true $s = 0.005, 0.0075, \text{ and } 0.01$). Under these conditions, SIA's median predictions
174 of s were noticeably closer to the true values (**Figure 3A**). At the same time, CLUES performed
175 slightly better than SIA in weak selection regimes (i.e., true $s = 0.001$ and 0.0025) (**Figure 3**).
176 Overall, SIA (RMSE = $9.52e-4$) achieved a lower error in estimating s than CLUES (RMSE =
177 $1.44e-3$), when true genealogies were used as input to both methods (Wilcoxon signed-rank test
178 for difference in mean of squared error, $p = 1.25e-42$). This finding potentially reflects the benefit
179 of linkage information utilized by SIA through the additional flanking genealogies (see
180 **Discussion**).

181
182 *Selection coefficient inference using inferred gene trees.* To account for gene-tree uncertainty,
183 we next used ARGs inferred with Relate, which is scalable to the size of the training dataset for
184 SIA (see **Methods**), as input to SIA and CLUES and compared their performance on CEU
185 simulations. Furthermore, we compared both methods to a supervised machine learning method,
186 ImaGene (see **Figure S20**), that operates directly on an image of the alignment itself. ImaGene
187 does not require gene trees as input and instead uses a Convolutional Neural Network (CNN) to
188 perform dimensionality reduction of the sequence alignment, allowing for accurate and efficient
189 classification and regression.

190
191 Overall, we found that SIA and ImaGene outperformed CLUES in these experiments (**Figure 4**).
192 CLUES tended to underestimate selection coefficients for the moderate-to-strong selection
193 regimes, to a greater extent compared to the case where true genealogies were used for inference
194 (**Figures 3A & 4A**). This decrease in performance of CLUES evidently derives from error at the
195 ARG reconstruction step. SIA, on the other hand, appeared to be more robust to the same ARG
196 reconstruction error. ImaGene performed remarkably similarly to SIA, given that it relies solely on
197 the sequence alignment. SIA exhibited lower error at neutral sites and sites with low-to-moderate
198 values of s , whereas ImaGene prevailed at sites under strong selection (**Figure 4B**).

199 Nevertheless, SIA showed a slightly smaller overall RMSE ($2.75e-3$) compared to ImaGene
200 ($2.91e-3$) (Wilcoxon signed-rank test, $p = 6.18e-38$), and in particular, SIA produces estimates of
201 s much closer to 0 for neutral loci. Notably, in this case both SIA and ImaGene were trained with
202 simulations under the same uniform distribution of s values (see **Methods**). A different choice of
203 training distribution could impact their performance across selection regimes (see **Discussion**).
204 Furthermore, we binned the results of these methods by both the segregating frequency and the
205 selection coefficient (see **Figure S4**) and again found that in general they exhibit higher variance
206 under low segregating frequency of the beneficial allele. As before, we also tested our regression
207 framework on true and inferred gene trees of test sets simulated under the *S. hypoxantha*
208 demographic model (see **Figure S5**). We found that SIA was approximately unbiased for the
209 moderate ($s = 0.005$) and high ($s = 0.01$) selection regimes but appeared to overestimate the
210 selection coefficient for regions under weak selection ($s = 0.001$ and 0.0025), when both true and
211 inferred genealogies were used as input. Furthermore, SIA appeared to overestimate the
212 selection coefficient for neutral regions when inferred gene trees were used as input, whereas it
213 was approximately unbiased for true gene trees.

214
215 *Performance on selection coefficient prediction with different sample sizes.* To explore the
216 tradeoffs associated with the use of larger data sets, we examined the performance of SIA under
217 different sample sizes, assuming a constant-sized demographic model ($N_e=10,000$). **Figure S6**
218 shows the error in selection coefficient inference on a held-out test set, stratified by the age of the
219 allele (panels **A&B**) and present-day derived allele frequency (panels **C&D**) at the site of interest.
220 We observed that sites with low frequency ($AF < 0.33$) and more recent (onset $< 0.2 \times 2N_e$
221 generations) alleles experience the most significant reduction in error as sample size increases.
222 Notably, the performance of SIA on more ancient alleles (onset $> 0.2 \times 2N_e$ generations) had little
223 to no improvement as the sample size increased from 32 to 254. These observations are in line
224 with the expectation that having more samples improves the chance of capturing low-frequency

225 alleles, but provides limited information about more ancient events. The reason for this age-
226 dependency is that, looking backwards in time, most lineages coalesce rapidly and only a few
227 survive to more ancient epochs, in a manner that depends only weakly on the sample size. It may
228 be useful to consider these observations when choosing the sample size for use in studying
229 selection in a particular context (see **Discussion**).

230

231 *Inference of allele frequency trajectory.* We further adapted the deep-learning architecture of SIA
232 to model the allele frequency (AF) trajectory at a site by retaining the output of the LSTM at each
233 time point (**Figure S1**, see **Methods**). We then evaluated the performance of SIA in the inference
234 of the AF trajectory using simulations under the CEU demography across a range of selection
235 coefficients and current DAFs. SIA was largely able to capture the expected trend of more rapidly
236 increasing AF under stronger selection (**Figure S7** and **S9**). In addition, AF estimates by SIA
237 using both true and inferred genealogies were generally unbiased, although AF at more recent
238 time points tended to be slightly underestimated when data was simulated under weaker
239 selection. AF estimates also appeared to be more accurate in terms of variance for alleles under
240 stronger selection (**Figure S8** and **S10**). As expected, the variance of AF estimates tended to
241 increase going further back in time (**Figure S8** and **S10**).

242

243 *Model performance on simulations with misspecified demographic models.* To evaluate the
244 robustness of SIA to mismatches between the demographic parameters used for simulating
245 training data and the true underlying demography of real data, we tested the method on the
246 selection-coefficient inference task with datasets simulated under a range of alternative
247 parameters. Each aspect of this model misspecification was assessed independently of the
248 others. In particular, the misspecified datasets contained simulations under (i) combinations of
249 population mutation (θ) and recombination (ρ) rates sampled beyond the range used for the
250 training data (**Figures S11** and **S14**), (ii) various alternative demographic scenarios (**Figures S12**,

251 **S15**, and **S17**), and (iii) various effective population sizes (**Figures S13** and **S16**). We compared
252 the performance of SIA on these misspecified datasets to that of CLUES [35], supplying both
253 methods with the true genealogies. We consider CLUES the “silver standard” when it comes to
254 robustness because it is unsupervised and therefore should not be susceptible to misspecified
255 training data compared to supervised learning methods such as SIA. Overall, we found that both
256 CLUES and SIA were reasonably robust to model misspecification (**Figures S11-13**), although
257 the performance of both methods inevitably declined when tested on severely misspecified data
258 (**Figure S13**). Interestingly, SIA tended to overestimate selection coefficient when the true N_e was
259 much smaller than that used for training, and underestimate it when the true N_e was much larger,
260 whereas CLUES did the opposite (**Figure S13**). Since the CLUES likelihood model of allele
261 frequency transition is parameterized by the population-scaled selection coefficient ($\alpha = 2Ns$), a
262 larger N_e likely appears to CLUES as equivalent to a higher s . On the other hand, features used
263 by SIA capture broad information of coalescence and linkage in the ARG, and therefore can be
264 distorted by misspecified N_e in more subtle ways (see **Discussion**). Using the same misspecified
265 dataset, we also ran SIA with Relate-inferred genealogies and compared its performance to that
266 of the genotyped-based deep-learning model ImaGene [14,15]. In general, SIA appeared to be
267 more robust to model misspecifications, achieving an overall RMSE of 0.00362, 0.00318 and
268 0.00374 in the misspecified θ/ρ , demography, and N_e experiments, respectively, compared to
269 ImaGene, whose RMSE was 0.00416, 0.00330 and 0.00462 in the corresponding experiments
270 (**Figures S14-16**). The advantage of SIA was particularly noticeable in cases of misspecified
271 demographic parameters (**Figures S15 & S16**). Notably, SIA exhibited reduced bias when
272 working with inferred genealogies compared to true genealogies, under conditions of extremely
273 mismatched N_e (compare **Figures S13 & S16**).

274

275 *Model prediction at genomic loci of interest in CEU population.* We then applied the SIA model to
276 identify selective sweeps and infer selection coefficients at selected genomic loci in the 1000

277 Genomes CEU population. These loci included the canonical example of selection at the *MCM6*
278 gene, which regulates the neighboring *LCT* gene and contributes to the lactase persistence trait
279 [36], the *ABCC11* gene regulating earwax production, several pigmentation-related genes, as well
280 as genes associated with obesity, diabetes and addiction (**Table 1**).

281
282 For *LCT*, SIA detected a strong signal of selection at the nearby SNP that has been associated
283 with the lactase persistence trait (rs4988235). At this SNP, SIA inferred a sweep probability close
284 to 1 and a selection coefficient greater than 0.01, making this one of the strongest signals of
285 selection in the human genome. A close examination of the local genealogy at this site reveals a
286 clear pattern indicative of a selective sweep — a burst of recent coalescence among the derived
287 lineages (orange taxa are the lineages carrying the derived allele) is clearly visible from the tree
288 (**Figure 5**).

289
290 At a number of pigmentation genes [37–41], SIA detected signals of moderate selection, including
291 *MC1R* (rs1805007, $P(\text{sweep}) = 0.95$, $s \approx 0.0037$), *KITLG* (rs12821256, $P(\text{sweep}) = 0.87$, $s \approx$
292 0.0019), *ASIP* (rs619865, $P(\text{sweep}) = 0.78$, $s \approx 0.0019$), *OCA2* (rs12913832, $P(\text{sweep}) = 0.75$, s
293 ≈ 0.0056) and *TYR* (rs1393350, $P(\text{sweep}) = 0.62$, $s \approx 0.0011$). In addition, SIA identified a weak
294 signal of selection at a SNP in the *ABCC11* gene (rs17822931), which influences earwax and
295 sweat production [42], with a selection coefficient of around 0.00035. There are few other
296 estimates for these genes available for comparison, but, notably, our estimate for *LCT* of $s \approx 0.01$
297 is consistent with a previous estimate on the order of 0.01-0.1 [36], and with recent studies of
298 ancient DNA samples [43,44] suggesting a value closer to 0.01. Our estimates suggest that
299 selection at the pigmentation loci is considerably weaker than at *LCT*, in contrast to previous
300 estimates for these loci, which covered a wide range but were generally considerably larger
301 (ranging from 0.02-0.1) [45]. Interestingly, CLUES estimated s at the *OCA2* locus to be on the
302 order of 0.001 (roughly similar to SIA's estimate of 0.0056), but s at the *KITLG*, *ASIP*, *TYR* loci to

303 be greater than 0.01 (in comparison to SIA's considerably smaller estimates of 0.0019, 0.0019,
304 and 0.0011) [35]. The apparent discrepancy between the estimates may be partially due to the
305 fact that the two methods used samples from two different populations (CEU for SIA and
306 GBR/British for CLUES).

307

308 On the other hand, SIA did not detect significant evidence of positive selection at several disease-
309 associated loci (rs7903146/*TCF7L2*, rs1800497/*ANKK1*, and rs9939609/*FTO*) or at several other
310 pigmentation loci (rs13289810/*TYRP1*, rs1003719/*TTC3*, and rs7495174/*OCA2*) (**Table 1**).
311 Notably, allele frequencies at these six loci tend to be similar in African and European populations
312 [46], suggesting that they are not likely to be under strong environment-dependent positive
313 selection, although it is possible that they have experienced very recent selective pressure that
314 SIA lacks the power to detect (see **Discussion**). Notably, *TYRP1* and *TTC3* also lacked signals
315 of selection in the CLUES analysis. Compared to the genealogies at sweep sites (**Figure 5**), the
316 trees at these putatively neutral loci lack the distinctive signature of recent bursts of coalescence
317 among derived lineages (**Figure 6**).

318

319 *Southern capuchino species analysis.* Our previous study of southern capuchino seedeaters
320 made use of the full ARG and machine learning to detect and characterize selective sweeps, and
321 suggested that soft sweeps are the dominant mode of adaptation in these species (see **Methods**
322 for more details). To further characterize the targets and strengths of positive selection in these
323 species, we applied SIA to polymorphism data [47] for *S. hypoxantha*, and adopted a conservative
324 approach by reporting only sites with DAF ≥ 0.5 , SIA-inferred $s \geq 0.0025$, and SIA-inferred sweep
325 probability ≥ 0.99 (see **Methods**). In addition to loci near top F_{ST} peaks and known pigmentation-
326 related genes (**Table 2**), we identified many more sites under positive selection located outside
327 the previously scanned F_{ST} peaks, amounting to a total of 15,551 putative partial soft sweep sites
328 across the 333 scanned scaffolds for *S. hypoxantha*. These sites can be prioritized for further

329 evaluation and downstream analysis. Notably, SIA enabled us to distinguish between selection at
330 regulatory and coding sequences, and we found that sweep loci near F_{ST} peaks and pigmentation
331 genes fall mostly in non-coding regions (**Table 2**). We additionally surveyed all putative sweep
332 sites identified by SIA and found that they are indeed enriched in non-coding regions (Fisher's
333 exact test, $p = 6.80 \times 10^{-5}$), particularly noticeable in the "near-coding" regions (**Figure S21**).
334 Consistent with the observation that the most highly differentiated SNPs among taxa are non-
335 coding [47,48] our finding suggests that positive selection may act on *cis*-regulatory regions to
336 drive differentiation and the subsequent speciation process. Furthermore, we examined many
337 individual predictions in detail, considering the local trees inferred by Relate at these high-
338 confidence predictions (**Figure 7**). We found, in numerous cases, that these sweeps had distinct
339 genealogical features, displaying evidence of a burst of coalescence events, corresponding to
340 unusually large and young clades. Prominent examples include predictions near pigmentation-
341 related genes *ASIP*, *KITL*, *SLC45A2*, and *TYRP1*.

342

343 Discussion

344 The ARG is useful for addressing a wide variety of biological questions ranging from inferring
345 demographic parameters to estimating allele ages. SIA exploits the particular utility of the ARG
346 for accurate inference of positive selection in a way that makes use of the full dataset, as opposed
347 to traditional summary statistics, which necessarily discard substantial information. Direct use of
348 the ARG improves upon traditional summary statistics in two key ways. First, it enables
349 consideration of the temporal distribution of coalescence and recombination events in the history
350 of the analyzed sequences, in contrast to traditional summary statistics that simply average over
351 these coalescence and/or recombination events. In addition, ARG-based methods provide better
352 spatial resolution by separately examining individual genealogies and the recombination
353 breakpoints between them, rather than averaging across windows containing unknown numbers

354 of genealogies. These detailed patterns of coalescences and linkage enable the ARG-based
355 approaches to capture a more localized and fine-grained picture of selection (e.g. infer selection
356 coefficient and allele frequency trajectory) as well as to achieve a better classification
357 performance. This performance advantage is particularly noticeable at lower DAFs and when
358 selection is weak, a regime where previous methods for selection inference fall short (**Figure 2**).

359

360 At the same time, the supervised machine-learning approach sets SIA apart from another ARG-
361 based method, CLUES, which approximates a full likelihood function for ARGs in the presence of
362 selection using importance sampling and a HMM. Although the accuracy of both SIA and CLUES
363 degraded when using inferred genealogies compared to true genealogies, reflecting the error and
364 uncertainty at the ARG inference step, SIA appeared to be more robust to gene tree uncertainty
365 (**Figures 3 and 4**). One possible reason for this observation is that CLUES effectively assumes
366 that the selection coefficient at the focal site is conditionally independent of the flanking trees
367 given the focal tree. This assumption should hold in the presence of fully specified genealogies,
368 but it may make CLUES more sensitive to errors in the inferred genealogies. In other words,
369 through its use of supervised learning, SIA may be able to compensate for the effects of
370 genealogy inference error on its estimation of the selection coefficient by also directly considering
371 the flanking trees and LD-related patterns among them. Still, the drop in accuracy observed
372 across methods underscores the dependency of ARG-based approaches on the ARG inference
373 method. For this reason, we anticipate that SIA may benefit substantially from further
374 improvement in ARG inference tools (see ref. [9]).

375

376 The ARG-based feature set distinguishes SIA from other supervised machine learning
377 approaches for characterizing selective sweeps. SIA uses local topological features of the ARG,
378 which are more informative than the SFS- or LD-based summary statistics employed by machine
379 learning methods such as S/HIC, SFselect, and evolBoosting. Using simulations, we

380 demonstrated that the SIA classifier outperformed a deep-learning method that aggregates these
381 traditional summary statistics (**Figure 2**). We also compared SIA with ImaGene, which represents
382 another flavor of supervised learning methods, inspired by the recent rise of CNNs for image
383 recognition. ImaGene encodes sequence alignments as images for powerful population genetic
384 inferences with CNNs and provides a state-of-the-art benchmark to compare against. We found
385 that ImaGene performs remarkably well across a wide range of simulations, but SIA does appear
386 to be somewhat less biased and more robust to model misspecification than ImaGene. The
387 evolutionary information in the ARG is implicit in the sequence alignment but some of this
388 information may be difficult for a brute-force machine learning model to discover directly.

389

390 We demonstrated that utilizing the ARG granted SIA considerably improved performance over
391 deep learning models solely employing traditional summary statistics. However, a possible
392 drawback of an ARG-based model is the potentially prohibitive computational overhead incurred
393 by ARG inference, especially as sample size grows. Picking a sample size when running SIA
394 involves a tradeoff between scalability (fewer samples, faster ARG inference) and performance
395 (more samples, slower ARG inference). We have found that SIA can infer selection coefficients
396 reasonably well with as few as 16 haplotypes. Including more samples did improve performance
397 but with a sublinear reduction in error (**Figure S6**). Therefore, a sample size from a few dozen to
398 a few hundreds — well within the capabilities of most modern ARG inference methods — strikes
399 a good balance between performance and scalability. Moreover, we found that larger sample
400 sizes improved prediction performance primarily for alleles at lower frequencies but had little
401 impact on the performance for more ancient alleles (as most lineages would have already
402 coalesced going further back in time) (**Figure S6**). This observation suggests that the choice of
403 the sample size when applying SIA should be guided by the biological question of interest —
404 ancient selection can be studied with just a handful of samples, whereas a larger sample size is
405 better suited to detect more recent sweeps.

406

407 Like other supervised learning methods, SIA relies on simulations to generate training data, and
408 therefore could be biased by subjective choices of simulation parameters. For example, SIA and
409 ImaGene cannot make accurate predictions of selection coefficients outside the range
410 represented in the training data (**Figure S18**), whereas unsupervised methods such as CLUES
411 are not limited to a pre-defined range (**Figure S19**). This problem could be circumvented by
412 training on an extended range of s . Similarly, the tendency of SIA to underestimate the selection
413 coefficient for sites under weak selection (**Figures 3, 4**) could be mitigated by augmenting the
414 training set with simulations densely sampled from the weak selection regime. A more subtle
415 issue, however, arises when the underlying generative process of the real data does not match
416 the assumptions made for the simulations of the training data, potentially compromising the
417 accuracy of the method when applied to real data. Thus, we tested SIA on simulations with
418 parameters mismatching those used in the training procedure. In general, we found that SIA was
419 fairly robust to alternative parameter values, although, as expected, performance did degrade
420 somewhat under severely misspecified models. Notably, SIA achieved a similar level of
421 robustness to model parameter misspecification as the unsupervised (i.e. not relying on training
422 data) likelihood method CLUES, yet outperformed the supervised deep learning method
423 ImaGene.

424

425 Applying SIA to the CEU panel from the 1000 Genomes Project yielded several noteworthy
426 findings at loci with known ties to phenotypes of interest. In addition to confirming the canonical
427 signal of selective sweep at the *LCT* locus, SIA detected a novel signal of selection at a GWAS
428 SNP in the *MC1R* gene associated with red hair color, contrasting a previous study that could not
429 find evidence of selection at *MC1R* in the European population [49]. The derived allele at this
430 locus segregates at around 10% in the CEU population but is nearly absent in non-European
431 populations [46]. In addition, at the *MC1R* locus the Relate test statistic for selection [25], which

432 tends to perform particularly well at low segregating frequencies (**Figure 2**), falls slightly below
433 the significance threshold of 0.05, supporting the evidence of positive selection at this locus. SIA
434 also detected evidence of selection at a SNP in the *ABCC11* gene reported to be the determinant
435 of wet versus dry earwax as well as sweat production, mirroring the signal of selection previously
436 found in the East Asian population [50], although selection in the CEU population appeared to be
437 much weaker. In addition, SIA identified selection at a few other pigmentation-related loci, yet
438 determined previously identified SNPs in the *TYRP1* and *TTC3* genes to be largely free from
439 selection (**Table 1**). These results were consistent with a previous study [35], which reported
440 similar results for these pigmentation-related loci, albeit in a slightly different population (GBR).
441 SIA notably did not detect positive selection at GWAS loci in the *TCF7L2* gene associated with
442 type-2 diabetes, the *ANKK1* gene implicated in addictive behaviors, and the *FTO* gene associated
443 with obesity. Overall, this empirical study with the 1000 Genomes CEU population has illustrated
444 how SIA can be applied to assess natural selection at the resolution of individual sites, suggesting
445 that it may be useful in prioritizing GWAS variants for further scrutiny.

446
447 In our previous work on southern capuchino seedeaters [31] (see **Methods**), we applied newly
448 developed statistical methods for ancestral recombination graph inference and machine-learning
449 for the prediction of selective sweeps. We found evidence suggesting that a substantial fraction
450 of soft sweeps are partial but had limited power to identify them (i.e. average accuracy of 56%).
451 SIA considerably improved our characterization of positive selection in the southern capuchino
452 species in two key ways. The SIA framework performs inference of selection directly from
453 genealogies instead of traditional summary statistics, and in doing so achieved an accuracy of up
454 to 96% in detecting partial soft sweeps. Consequently, we found abundant evidence of soft
455 sweeps beyond the previously scanned F_{ST} peaks, and additionally were able to estimate their
456 selection coefficients. Importantly, SIA also took the analysis of selection beyond broad genomic
457 windows containing sweeps to the identification of specific putative causal variants. We took

458 advantage of this substantial improvement in genomic resolution and analyzed the distribution of
459 these sweep sites, which revealed that positive selection on regions that likely contain *cis*-
460 regulatory elements plays a role in driving the differentiation and speciation of southern capuchino
461 seedeaters.

462

463 While we believe SIA represents an important step forward in the use of the ARG for machine-
464 learning-based selection inference, there remain several possible avenues for improvement. For
465 example, SIA currently uses a point-estimate of the ARG, rather than a distribution, and therefore
466 does not explicitly take gene-tree uncertainty into account. We plan to improve SIA by using
467 strategies for inferring approximate posterior distribution of ARGs (e.g., [24]), as well as designing
468 better algorithms for ARG reconstruction that balance accuracy with scalability and can handle
469 thousands of genomes. In addition, the SIA framework was applied in the context of single-locus
470 selective sweeps, but could be extended to study polygenic selection, by making use of summary
471 statistics from genome-wide association studies (as in [51]) and adapting the architecture of our
472 neural network to account for selection acting at multiple sites. Finally, the robustness of SIA to
473 model misspecifications can be further improved by ensuring the simulated data is generated
474 under a distribution that is compatible with the real target data set. We anticipate that the continual
475 advancement in ARG inference methods has the potential to open up many new applications for
476 this flexible and powerful model of ARG-based deep learning in population genetics.

477

478 Methods

479 *Simulated datasets used for training and testing the selective sweep model.* Training and testing
480 data sets were generated using discoal [52] by simulating 1,000,000 regions of length 100 kb for
481 each model we considered (i.e., “neutral” or “hard sweep”). Aside from these regions, 2,000 were
482 simulated for validation and 5,000 were simulated for testing. The number of sampled sequences

483 was selected to match the number of individuals in the CEU population in the 1000 Genomes
484 dataset. Thus, a total of 198 haploid sequences were sampled. Simulations used a demographic
485 model based on European demography [34]. In non-neutral simulations, selection was applied to
486 a single focal site located in the middle of the simulated region. We sampled each of the main
487 demographic and selection parameters from a uniform distribution: (1) mutation rate $\mu \sim U(1.25e-$
488 $08, 2.5e-08)$, (2) recombination rate $\rho \sim U(1.25e-08, 2.5e-08)$, (3) selection coefficient $s \sim$
489 $U(0.0001, 0.02)$, and (4) segregating frequency of the site under selection $f \sim U(0.01, 0.99)$.

490

491 *ARG Feature Extraction.* For each target variant, we extracted the corresponding gene tree from
492 the ARG, then overlaid it with 100 discrete timepoints. These timepoints were fixed across all
493 trees in an approximately log-uniform manner that resulted in finer discretization of more recent
494 time scales (as in [24]). We considered biallelic sites only and assumed no recurrent mutations;
495 thus each mutation was assumed to occur on the branch of the tree where the ancestral allele
496 switches to the derived. For each timepoint, we calculated the number of active ancestral and
497 derived lineages. Furthermore, we computed the number of all active lineages (not distinguishing
498 between ancestral and derived) at the same set of predefined timepoints in the two left and right
499 flanking gene trees to account for linkage disequilibrium. Together, these features were
500 summarized in a 600-dimensional feature vector, which was then used as input to an RNN. The
501 feature of a simulated sweep region was extracted from the sweep site (by default at the center
502 in all simulations) whereas the feature of a simulated neutral region was extracted from a variant
503 site (randomly chosen) with a pre-defined matched derived allele frequency. The features for each
504 genomic locus of interest in the CEU population were extracted from all variant sites at that locus
505 having a derived allele frequency of >0.05 .

506

507 *Training an RNN to predict different modes of selection.* An RNN was applied to the simulated
508 training data sets to learn a classification or regression model for the task at hand. We used a

509 Long Short-Term Memory (LSTM), a particular form of RNN, to accommodate the temporal nature
510 of our features and account for long-term dependencies and the vanishing gradient problem
511 observed in traditional RNNs. Our model had 100 timepoints with the final target output depending
512 on the use of classification or regression. For the classification task, the final target output is a
513 label for a binary classification problem predicting whether a region is under selection or neutrality.
514 For the regression task, the final target output is a continuous value, representing the selection
515 coefficient or the time of selection onset. We also took a many-to-many approach to model the
516 allele-frequency trajectory for the site under selection. The *Keras* software was used to train and
517 test the model. We used a two-stacked LSTM to account for greater model complexity where the
518 number of units in each stack was set to 100 and the hyperbolic tangent (*tanh*) was used as an
519 activation function. The *Adam* optimization method with its default operating parameters was used
520 to update the network weights. For the classification task, the *Softmax* activation function was
521 applied on the final dense layer and the *binary_crossentropy* was used to compute the cross-
522 entropy loss between true labels and predicted labels. For the regression task, the *linear*
523 activation function was applied on the final dense layer and the *mean_squared_error* function was
524 used.

525

526 *Estimation of Confidence Intervals.* To turn our single-valued regression model into one capable
527 of returning a distribution of predictions of s , we reused the dropout technique that is typically
528 used during training. Dropout enables a fraction of nodes to be randomly “turned off” in a certain
529 layer, which assists in the regularization of the model and helps prevent overfitting. We applied
530 dropout during inference, enabling us to sample a “thinned” network to generate a sample
531 prediction. By repeatedly sampling thinned networks, we generated a distribution of predictions
532 and then computed confidence intervals based on this distribution [53].

533

534 *ARG Inference*. Relate [25] (v1.0.17) was used for inferring ARGs underlying simulated genomic
535 samples as well as the CEU population in the 1000 Genomes dataset. For simulations under the
536 Tennesen *et al.* demography [34], Relate was run with the true simulation parameters (μ , ρ and
537 N_e) specified; whereas for genomic loci of the CEU population, Relate was run with a mutation
538 rate of 2.5×10^{-8} /base/generation (-m 2.5e-8), a constant recombination map of 1.25×10^{-8}
539 /base/generation and a diploid effective population size of 188,088 (-N 376176). The choice of
540 mutation rate follows [35] based on estimates from [54]. Although some more recent estimates
541 have been lower [55], these differences in mutation rate are unlikely to have a major effect on our
542 selection inference since SIA appears to be fairly robust to misspecification of mutation rate
543 (**Figures S11 & S14**). For simulations and genomic loci of the *S. hypoxantha* population, Relate
544 was run with $\mu=\rho=1 \times 10^{-9}$ /base/generation and a diploid N_e of 130,000. The branch lengths of
545 Relate-inferred genealogies were estimated iteratively with the `EstimatePopulationSize.sh` script
546 in the Relate package. Specifically, population size history was inferred from the ARG, the branch
547 lengths are then updated for the estimated population size history and these steps are repeated
548 until convergence. This was done for a default of 5 iterations (--num_iter 5).

549

550 *Alternative methods for selection inference*. To benchmark the performance of SIA for
551 classification of sites under neutrality versus selective sweep, we ran the following methods:
552 Tajima's D [10], H1 [32], iHS [33], a summary statistics-based deep learning model, and a tree-
553 based statistic that is part of the Relate [25] program. Tajima's D, H1 and iHS were calculated
554 with the *scikit-allele* package. Haplotypes of the entire 100kb simulated genomic segment were
555 used for Tajima's D and H1 calculations. The unstandardized iHS was computed at every site
556 with minor allele frequency > 5%, with respect to all other sites in the genomic segment
557 (min_maf=0.05, include_edges=True). iHS scores of all sites were then standardized in 50 allele-
558 frequency bins. Finally, the iHS score of a genomic region was taken to be the mean of the iHS
559 scores of all of its variant sites. For the summary statistics-based deep learning model, we made

560 use of the summary statistics used by S/HIC [12,13] as features for our deep learning architecture.
561 These included 11 sequence-based summary statistics (see **Figure 3** in [56]) which were used
562 as features for our deep learning model to distinguish among the two classes at hand (selective
563 sweep versus neutral drift). All statistics were collected along five consecutive 20-kb windows with
564 the objective of identifying possible sweeps induced by a positively selected mutation in the third
565 (middle) window. Some of these summary statistics corresponded to standard measures of
566 diversity, such as ss (the number of segregating sites), π [57], Tajima's D [10], θ_w [58], θ_H [59],
567 the number of distinct haplotypes [60], $H1$, $H12$, $H2/H1$ [32], Z_{ns} [61], and maximum value of ω
568 [62]. For each of these statistics, we computed an average value for each of the five 20 kb
569 windows for the simulated population. Finally, each summary statistic was normalized by dividing
570 the value recorded for a given window by the sum of values across all five windows. The Relate
571 tree-based selection test was performed with an add-on module (*DetectSelection.sh*) using the
572 inferred genealogy with calibrated branch lengths at a site of interest, yielding a \log_{10} p-value for
573 each site. We also compared the performance of SIA for selection coefficient inference to that of
574 CLUES [35] and a genotype-based convolutional neural network (CNN) framework [14,15].
575 Selection coefficient inference from true genealogies was performed with *clues-v0*
576 (<https://github.com/35ajstern/clues-v0>). Transition probability matrices were built on a range of
577 selection coefficients [0, 0.05] at increments of 0.0001 and present-day allele frequencies [0.01,
578 0.99] at increments of 0.01. Selection-coefficient inference from Relate inferred genealogies was
579 performed with CLUES (<https://github.com/35ajstern/clues>). Branch lengths of the genealogy at
580 the site of interest were resampled with Relate for 600 MCMC iterations, and CLUES was run
581 with the following arguments: `--tCutoff 10000 --burnin 100 --thin 5`. For the genotype-based CNN
582 model, each simulated genomic segment was preprocessed by first sorting the haplotypes and
583 then converting the segment to a fixed-size genotype matrix. Haplotype sorting was performed
584 by 1) calculating the pairwise manhattan distances between haplotypes, 2) setting the haplotype
585 with the smallest total distance to all other haplotypes as the first haplotype, and 3) sorting the

586 remaining haplotypes in increasing distance to the first haplotype. To convert the sorted
587 haplotypes to a fixed-size genotype matrix, centered on the middle variant of a simulated region,
588 up to 180 variants on each side were retained. Variants beyond 180 were discarded and if there
589 were fewer than 180, the missing variants were padded with zeros. Ancestral and derived alleles
590 were coded with 0's and 1's, respectively. Consequently, each simulated genomic region was
591 encoded as a (198 x 360) binary matrix, along with a real-valued vector encoding the genomic
592 positions of the variants in the matrix. The CNN model had a branched architecture — one branch
593 with five 1D convolution layers taking the genotype matrix as input and another branch with a fully
594 connected layer taking the vector of variant positions as input. The output of the two branches
595 was flattened, concatenated and fed into 3 fully connected layers, followed by a linear output layer
596 to predict selection coefficient (**Figure S20**).

597

598 *Evaluation metrics.* To evaluate the performance of SIA's classification model and alternative
599 methods, we computed a receiver operating characteristic (ROC) curve for the binary class at
600 hand ("neutral" or "sweep"), to provide a more complete summary of the behavior of different
601 types of errors. We further assessed the performance of SIA and alternative methods in terms
602 of correctly predicting the selection coefficient numerically using mean absolute error (mae),
603 root mean square error (rmse), coefficient of determination (r^2), and visually using a box plot that
604 compares the simulated ground truth against the predictions by the method at hand.

605

606 *Robustness study.* We carried out an extensive analysis of the robustness of our approach,
607 considering not only alternative demographic parameters (such as population size), but also
608 alternative parameters for recombination rate, mutation rate, time of selection onset, and selection
609 coefficients. In all cases, we took care to test our prediction methods under parameters well
610 outside the range used in training.

611

612 *Analysis of CEU population in 1000 Genomes data.* We applied SIA to infer selection coefficients
613 and allele frequency trajectories in the 1000 Genomes [63] CEU population at 13 genomic loci
614 with known association to phenotypes, some of which were previously identified as likely targets
615 of positive selection (**Table 1**). For each gene of interest, the ARG was inferred with Relate from
616 SNPs within a 2Mb window centered at the gene. Once the ARG was inferred, only SNPs with
617 valid ancestral allele ('AA' INFO field in the vcf file) were retained for downstream analysis.
618 Following the aforementioned protocol (see *ARG feature extraction*), features at all variant sites
619 in the 2Mb window above a derived allele frequency threshold of 0.05 were extracted. Lastly, the
620 SIA model was applied to classify neutrality versus selection, and infer selection coefficient and
621 allele frequency trajectory at each site.

622

623 *Localizing sweeps in southern capuchino seedeaters.* We recently applied a combination of ARG
624 inference and machine-learning methods for identifying selective sweeps to study previously
625 identified “islands of differentiation” in southern capuchino seedeaters and distinguish among
626 possible evolutionary scenarios leading to their formation [31]. Taking advantage of its improved
627 power and genomic resolution, we applied SIA to sequence data for the species for which we
628 have the most samples, *Sporophila hypoxantha*. We simulated training (250,000 neutral; 250,000
629 soft sweeps), validation (1000 neutral; 1000 soft sweeps), and testing (2,500 neutral; 2,500 soft
630 sweeps) data sets for SIA under a demographic model inferred by G-PhoCS [64]. Simulations
631 were performed using discoal with the following parameters: (1) mutation rate $\mu = 1e-9$, (2)
632 recombination rate $\rho = 1e-9$, (3) derived $N_e = 130,000$, (4) root divergence time = 1,850,000
633 generations ago, (5) root $N_e = 1,450,000$, (6) ancestral divergence time = 44,000 generations ago,
634 (7) ancestral $N_e = 14,380,000$, (8) selection coefficient $s \sim U(0.001, 0.02)$, (9) initial frequency at
635 which selection starts acting on the allele $f_{init} \sim U(0.01, 0.05)$, and (10) segregating frequency of
636 the site under selection $f \sim U(0.25, 0.99)$. A total of 56 haploid sequences were sampled from
637 each simulation, matching the number of *S. hypoxantha* individuals (28) in the real data. The SIA

638 model for *S. hypoxantha* was built, trained and evaluated in an otherwise similar fashion to that
639 for the CEU population as outlined above.

640

641 Using a subset of polymorphism data from [47] of 28 *S. hypoxantha* and 2 *S. minuta* individuals,
642 we applied our trained model to localize selective sweeps in *S. hypoxantha* on 19 scaffolds that
643 contain top F_{ST} peaks in at least one pairwise species comparison [48] and/or harbor known
644 pigmentation-related genes such as *ASIP* (located on scaffold 252; induces melanocytes to
645 synthesize pheomelanin instead of eumelanin), *KITL* (located on scaffold 412; stimulates
646 melanocyte proliferation), *SLC45A2* (located on scaffold 404; transports substances needed for
647 melanin synthesis), and *CAMK2D* (located on scaffold 1717; cell communication), as well as 316
648 scaffolds that i) are longer than 100kb, ii) contain more than 1,000 variants, and iii) where more
649 than 95% of sites have a consensus ancestral allele, as determined by four identical haplotypes
650 for two individuals from the outgroup species *S. minuta*. The ARG was inferred with Relate for
651 each scaffold independently. Once the ARG was inferred, the SIA model was applied to sites with
652 consensus ancestral allele for classification and selection coefficient inference.

653

654 Acknowledgments

655 The authors would like to acknowledge Noah Dukler for help with Figure 1 preparation. This
656 research was supported by US National Science Foundation grant (NSFDEB) 1555769, US
657 National Institutes of Health grant R35-GM127070, the CSHL School of Biological Sciences
658 Gladys & Roland Harriman Fellowship, and the Simons Center for Quantitative Biology at Cold
659 Spring Harbor Laboratory. The content is solely the responsibility of the authors and does not
660 necessarily represent the official views of the US National Institutes of Health or the US National
661 Science Foundation.

662

663 Availability of data and materials

664 The scripts used for analyses in this study are available at [github.com/CshSiepellLab/arg-](https://github.com/CshSiepellLab/arg-selection)
665 [selection](https://github.com/CshSiepellLab/arg-selection) under a GNU GPLv3 license.

References

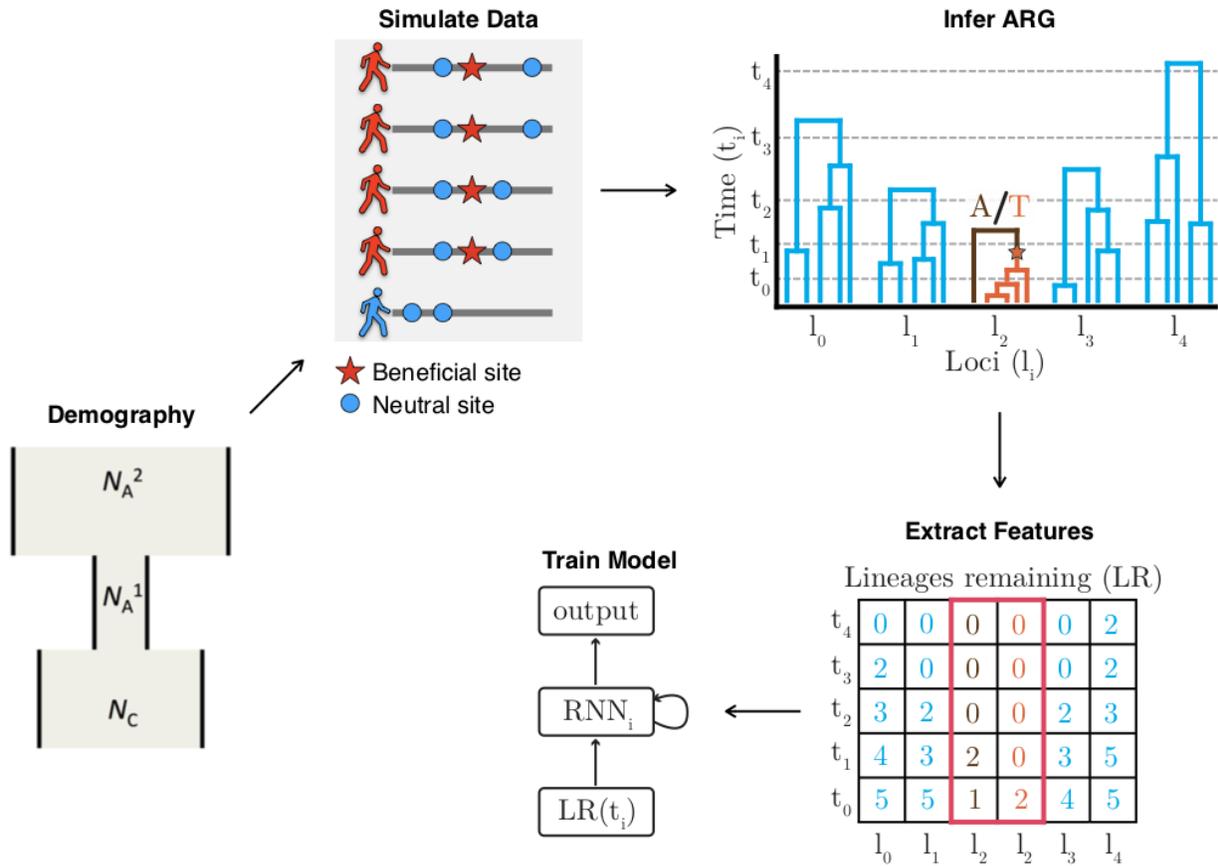
1. Ohashi J, Naka I, Patarapotikul J, Hananantachai H, Brittenham G, Looareesuwan S, et al. Extended Linkage Disequilibrium Surrounding the Hemoglobin E Variant Due to Malarial Selection. *Am J Hum Genet.* 2004;74: 1198–1208. doi:10.1086/421330
2. Currat M, Trabuchet G, Rees D, Perrin P, Harding RM, Clegg JB, et al. Molecular Analysis of the β -Globin Gene Cluster in the Niokholo Mandenka Population Reveals a Recent Origin of the β S Senegal Mutation. *Am J Hum Genet.* 2002;70: 207–223. doi:10.1086/338304
3. Jarvis JP, Scheinfeldt LB, Soi S, Lambert C, Omberg L, Ferwerda B, et al. Patterns of Ancestry, Signatures of Natural Selection, and Genetic Association with Stature in Western African Pygmies. *PLOS Genet.* 2012;8: e1002641. doi:10.1371/journal.pgen.1002641
4. Lachance J, Vernot B, Elbers CC, Ferwerda B, Froment A, Bodo J-M, et al. Evolutionary History and Adaptation from High-Coverage Whole-Genome Sequences of Diverse African Hunter-Gatherers. *Cell.* 2012;150: 457–469. doi:10.1016/j.cell.2012.07.009
5. Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, et al. Genome-wide detection and characterization of positive selection in human populations. *Nature.* 2007;449: 913–918. doi:10.1038/nature06250
6. Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, et al. Positive Natural Selection in the Human Lineage. *Science.* 2006;312: 1614–1620. doi:10.1126/science.1124309
7. Kelley JL, Swanson WJ. Positive Selection in the Human Genome: From Genome Scans to Biological Significance. *Annu Rev Genomics Hum Genet.* 2008;9: 143–160. doi:10.1146/annurev.genom.9.081307.164411
8. Fu W, Akey JM. Selection and Adaptation in the Human Genome. *Annu Rev Genomics Hum Genet.* 2013;14: 467–489. doi:10.1146/annurev-genom-091212-153509
9. Hejase HA, Dukler N, Siepel A. From Summary Statistics to Gene Trees: Methods for Inferring Positive Selection. *Trends Genet.* 2020;36: 243–258. doi:10.1016/j.tig.2019.12.008
10. Tajima F. Statistical Method for Testing the Neutral Mutation Hypothesis by DNA Polymorphism. *Genetics.* 1989;123: 585–595.
11. Peter BM, Huerta-Sanchez E, Nielsen R. Distinguishing between Selective Sweeps from Standing Variation and from a De Novo Mutation. *PLOS Genet.* 2012;8: e1003011. doi:10.1371/journal.pgen.1003011
12. Kern AD, Schrider DR. diploS/HIC: An Updated Approach to Classifying Selective Sweeps. *G3 Genes Genomes Genet.* 2018;8: 1959–1970. doi:10.1534/g3.118.200262
13. Schrider DR, Kern AD. S/HIC: Robust Identification of Soft and Hard Sweeps Using Machine Learning. *PLOS Genet.* 2016;12: e1005928. doi:10.1371/journal.pgen.1005928
14. Flagel L, Brandvain Y, Schrider DR. The Unreasonable Effectiveness of Convolutional Neural Networks in Population Genetic Inference. Kim Y, editor. *Mol Biol Evol.* 2019;36:

- 220–238. doi:10.1093/molbev/msy224
15. Torada L, Lorenzon L, Beddis A, Isildak U, Pattini L, Mathieson S, et al. ImaGene: a convolutional neural network to quantify natural selection from genomic data. *BMC Bioinformatics*. 2019;20: 337. doi:10.1186/s12859-019-2927-x
 16. Griffiths RC, Marjoram P. Ancestral inference from samples of DNA sequences with recombination. *J Comput Biol J Comput Mol Cell Biol*. 1996;3: 479–502. doi:10.1089/cmb.1996.3.479
 17. Hudson RR. Gene genealogies and the coalescent process. *Oxf Surv Evol Biol*. 1990;7: 1–44.
 18. Wiuf C, Hein J. Recombination as a Point Process along Sequences. *Theor Popul Biol*. 1999;55: 248–259. doi:10.1006/tpbi.1998.1403
 19. Hein J. A heuristic method to reconstruct the history of sequences subject to recombination. *J Mol Evol*. 1993;36: 396–405. doi:10.1007/BF00182187
 20. Song YS, Hein J. Constructing Minimal Ancestral Recombination Graphs. *J Comput Biol*. 2005;12: 147–169. doi:10.1089/cmb.2005.12.147
 21. Minichiello MJ, Durbin R. Mapping Trait Loci by Use of Inferred Ancestral Recombination Graphs. *Am J Hum Genet*. 2006;79: 910–922. doi:10.1086/508901
 22. Kuhner MK. LAMARC 2.0: maximum likelihood and Bayesian estimation of population parameters. *Bioinformatics*. 2006;22: 768–770. doi:10.1093/bioinformatics/btk051
 23. O’Fallon BD. ACG: rapid inference of population history from recombining nucleotide sequences. *BMC Bioinformatics*. 2013;14: 40. doi:10.1186/1471-2105-14-40
 24. Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. Genome-Wide Inference of Ancestral Recombination Graphs. *PLoS Genet*. 2014;10: e1004342. doi:10.1371/journal.pgen.1004342
 25. Speidel L, Forest M, Shi S, Myers SR. A method for genome-wide genealogy estimation for thousands of samples. *Nat Genet*. 2019;51: 1321–1329. doi:10.1038/s41588-019-0484-x
 26. Kelleher J, Wong Y, Wohns AW, Fadil C, Albers PK, McVean G. Inferring whole-genome histories in large population datasets. *Nat Genet*. 2019;51: 1330–1338. doi:10.1038/s41588-019-0483-y
 27. Arenas M. The importance and application of the ancestral recombination graph. *Front Genet*. 2013;4. doi:10.3389/fgene.2013.00206
 28. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521: 436–444. doi:10.1038/nature14539
 29. Hochreiter S, Schmidhuber J. Long Short-Term Memory. *Neural Comput*. 1997;9: 1735–1780. doi:10.1162/neco.1997.9.8.1735
 30. Maas AL, Daly RE, Pham PT, Huang D, Ng AY, Potts C. Learning Word Vectors for Sentiment Analysis. *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics; 2011. pp. 142–150. Available: <https://www.aclweb.org/anthology/P11-1015>
 31. Hejase HA, Salman-Minkov A, Campagna L, Hubisz MJ, Lovette IJ, Gronau I, et al. Genomic islands of differentiation in a rapid avian radiation have been driven by recent selective sweeps. *Proc Natl Acad Sci*. 2020;117: 30554–30565. doi:10.1073/pnas.2015987117
 32. Garud NR, Messer PW, Buzbas EO, Petrov DA. Recent Selective Sweeps in North American *Drosophila melanogaster* Show Signatures of Soft Sweeps. *PLoS Genet*. 2015;11: e1005004. doi:10.1371/journal.pgen.1005004
 33. Voight BF, Kudaravalli S, Wen X, Pritchard JK. A Map of Recent Positive Selection in the Human Genome. *Hurst L, editor. PLoS Biol*. 2006;4: e72. doi:10.1371/journal.pbio.0040072

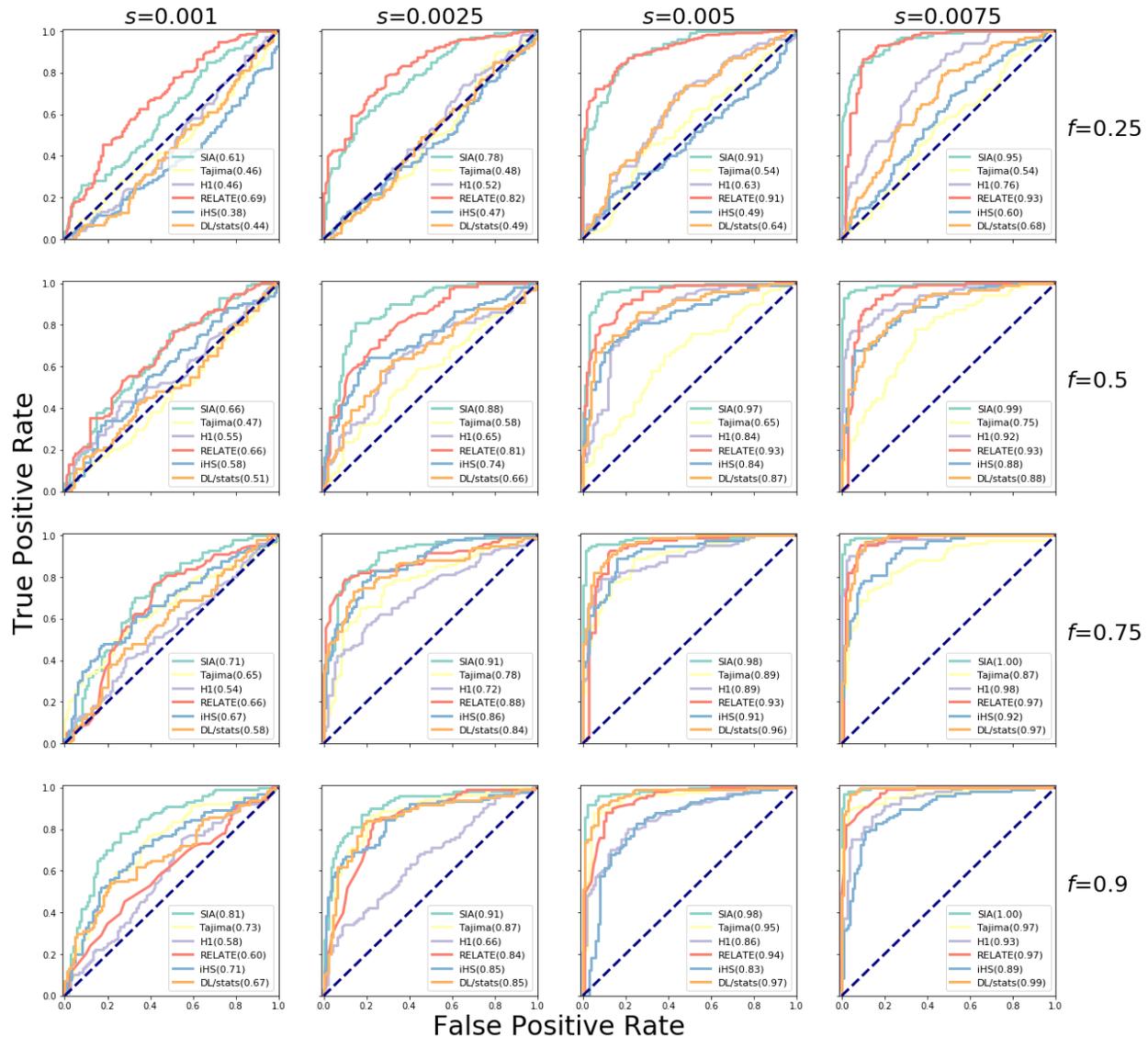
34. Tennesen JA, Bigham AW, O'Connor TD, Fu W, Kenny EE, Gravel S, et al. Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes. *Science*. 2012;337: 64–69. doi:10.1126/science.1219240
35. Stern AJ, Wilton PR, Nielsen R. An approximate full-likelihood method for inferring selection and allele frequency trajectories from DNA sequence data. Hernandez RD, editor. *PLOS Genet*. 2019;15: e1008384. doi:10.1371/journal.pgen.1008384
36. Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake JA, et al. Genetic Signatures of Strong Recent Positive Selection at the Lactase Gene. *Am J Hum Genet*. 2004;74: 1111–1120. doi:10.1086/421051
37. Han J, Kraft P, Nan H, Guo Q, Chen C, Qureshi A, et al. A Genome-Wide Association Study Identifies Novel Alleles Associated with Hair Color and Skin Pigmentation. *PLOS Genet*. 2008;4: e1000074. doi:10.1371/journal.pgen.1000074
38. Sturm RA, Duffy DL, Zhao ZZ, Leite FPN, Stark MS, Hayward NK, et al. A Single SNP in an Evolutionary Conserved Region within Intron 86 of the *HERC2* Gene Determines Human Blue-Brown Eye Color. *Am J Hum Genet*. 2008;82: 424–431. doi:10.1016/j.ajhg.2007.11.005
39. Sulem P, Gudbjartsson DF, Stacey SN, Helgason A, Rafnar T, Magnusson KP, et al. Genetic determinants of hair, eye and skin pigmentation in Europeans. *Nat Genet*. 2007;39: 1443–1452. doi:10.1038/ng.2007.13
40. Kenny EE, Timpson NJ, Sikora M, Yee M-C, Moreno-Estrada A, Eng C, et al. Melanesian blond hair is caused by an amino acid change in *TYRP1*. *Science*. 2012;336: 554. doi:10.1126/science.1217849
41. Liu F, Wollstein A, Hysi PG, Ankra-Badu GA, Spector TD, Park D, et al. Digital Quantification of Human Eye Color Highlights Genetic Association of Three New Loci. *PLOS Genet*. 2010;6: e1000934. doi:10.1371/journal.pgen.1000934
42. Yoshiura K, Kinoshita A, Ishida T, Ninokata A, Ishikawa T, Kaname T, et al. A SNP in the *ABCC11* gene is the determinant of human earwax type. *Nat Genet*. 2006;38: 324–330. doi:10.1038/ng1733
43. Mathieson S, Mathieson I. *FADS1* and the Timing of Human Adaptation to Agriculture. *Mol Biol Evol*. 2018;35: 2957–2970. doi:10.1093/molbev/msy180
44. Mathieson I. Estimating time-varying selection coefficients from time series data of allele frequencies. *bioRxiv*. 2020; 2020.11.17.387761. doi:10.1101/2020.11.17.387761
45. Wilde S, Timpson A, Kirsanow K, Kaiser E, Kayser M, Unterländer M, et al. Direct evidence for positive selection of skin, hair, and eye pigmentation in Europeans during the last 5,000 y. *Proc Natl Acad Sci*. 2014;111: 4832–4837. doi:10.1073/pnas.1316513111
46. Marcus JH, Novembre J. Visualizing the geography of genetic variants. *Bioinformatics*. 2017;33: 594–595. doi:10.1093/bioinformatics/btw643
47. Turbek SP, Browne M, Giacomo ASD, Kopuchian C, Hochachka WM, Estalles C, et al. Rapid speciation via the evolution of pre-mating isolation in the Iberá Seedeater. *Science*. 2021;371. doi:10.1126/science.abc0256
48. Campagna L, Repenning M, Silveira LF, Fontana CS, Tubaro PL, Lovette IJ. Repeated divergent selection on pigmentation genes in a rapid finch radiation. *Sci Adv*. 2017;3: e1602404. doi:10.1126/sciadv.1602404
49. Harding RM, Healy E, Ray AJ, Ellis NS, Flanagan N, Todd C, et al. Evidence for Variable Selective Pressures at *MC1R*. *Am J Hum Genet*. 2000;66: 1351–1361. doi:10.1086/302863
50. Ohashi J, Naka I, Tsuchiya N. The Impact of Natural Selection on an *ABCC11* SNP Determining Earwax Type. *Mol Biol Evol*. 2011;28: 849–857. doi:10.1093/molbev/msq264
51. Stern AJ, Speidel L, Zaitlen NA, Nielsen R. Disentangling selection on genetically correlated polygenic traits via whole-genome genealogies. *Am J Hum Genet*. 2021;108: 219–239. doi:10.1016/j.ajhg.2020.12.005

52. Kern AD, Schrider DR. Discoal: flexible coalescent simulations with selection. *Bioinformatics*. 2016;32: 3839–3841. doi:10.1093/bioinformatics/btw556
53. Gal Y, Ghahramani Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *International Conference on Machine Learning*. PMLR; 2016. pp. 1050–1059. Available: <http://proceedings.mlr.press/v48/gal16.html>
54. Nachman MW, Crowell SL. Estimate of the mutation rate per nucleotide in humans. *Genetics*. 2000;156: 297–304.
55. Scally A, Durbin R. Revising the human mutation rate: implications for understanding human evolution. *Nat Rev Genet*. 2012;13: 745–753. doi:10.1038/nrg3295
56. Schrider DR, Kern AD. Supervised Machine Learning for Population Genetics: A New Paradigm. *Trends Genet*. 2018;34: 301–312. doi:10.1016/j.tig.2017.12.005
57. Nei M, Li WH. Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proc Natl Acad Sci U S A*. 1979;76: 5269–5273.
58. Watterson GA. On the number of segregating sites in genetical models without recombination. *Theor Popul Biol*. 1975;7: 256–276. doi:10.1016/0040-5809(75)90020-9
59. Fay JC, Wu C-I. Hitchhiking Under Positive Darwinian Selection. *Genetics*. 2000;155: 1405–1413.
60. Messer PW, Petrov DA. Population genomics of rapid adaptation by soft selective sweeps. *Trends Ecol Evol*. 2013;28: 659–669. doi:10.1016/j.tree.2013.08.003
61. Kelly JK. A Test of Neutrality Based on Interlocus Associations. *Genetics*. 1997;146: 1197–1206.
62. Kim Y, Nielsen R. Linkage Disequilibrium as a Signature of Selective Sweeps. *Genetics*. 2004;167: 1513–1524. doi:10.1534/genetics.103.025387
63. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global reference for human genetic variation. *Nature*. 2015;526: 68–74. doi:10.1038/nature15393
64. Campagna L, Gronau I, Silveira LF, Siepel A, Lovette IJ. Distinguishing noise from signal in patterns of genomic divergence in a highly polymorphic avian radiation. *Mol Ecol*. 2015;24: 4238–4251. doi:10.1111/mec.13314
65. Eriksson N, Macpherson JM, Tung JY, Hon LS, Naughton B, Saxonov S, et al. Web-Based, Participant-Driven Studies Yield Novel Genetic Associations for Common Traits. *PLOS Genet*. 2010;6: e1000993. doi:10.1371/journal.pgen.1000993
66. Lyssenko V, Lupi R, Marchetti P, Guerra SD, Orho-Melander M, Almgren P, et al. Mechanisms by which common variants in the *TCF7L2* gene increase risk of type 2 diabetes. *J Clin Invest*. 2007;117: 2155–2163. doi:10.1172/JCI30706
67. Spellacy CJ, Harding MJ, Hamon SC, Mahoney JJ, Reyes JA, Kosten TR, et al. A variant in *ANKK1* modulates acute subjective effects of cocaine: a preliminary study. *Genes Brain Behav*. 2014;13: 559–564. doi:10.1111/gbb.12121
68. Frayling TM, Timpson NJ, Weedon MN, Zeggini E, Freathy RM, Lindgren CM, et al. A common variant in the *FTO* gene is associated with body mass index and predisposes to childhood and adult obesity. *Science*. 2007;316: 889–894. doi:10.1126/science.1141634

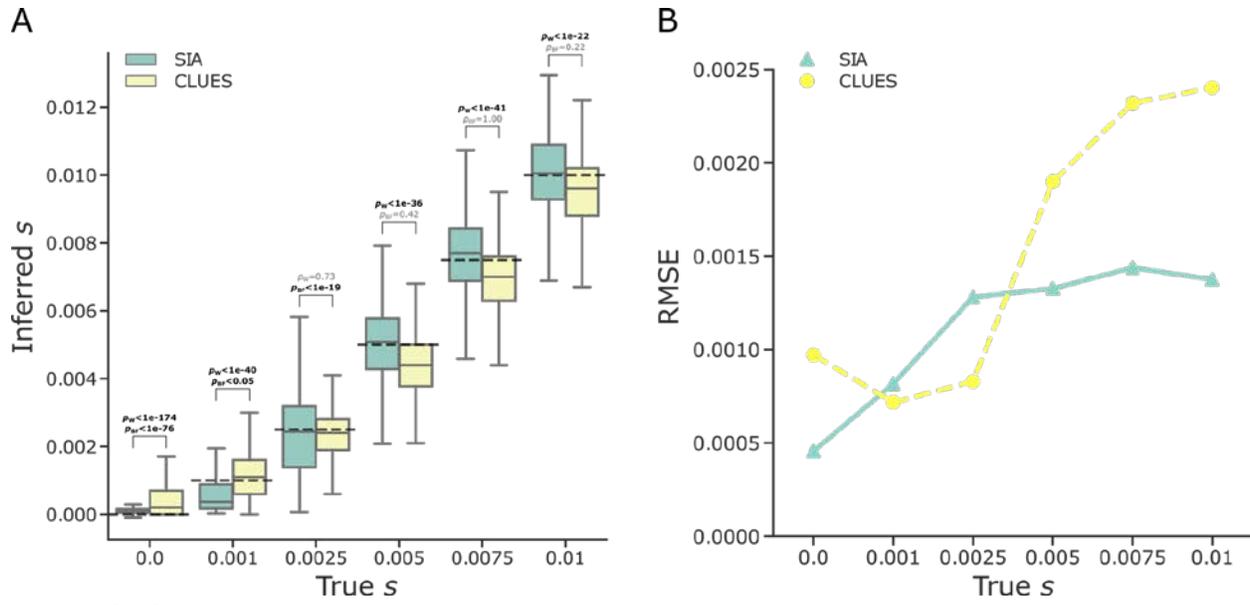
666 Figures



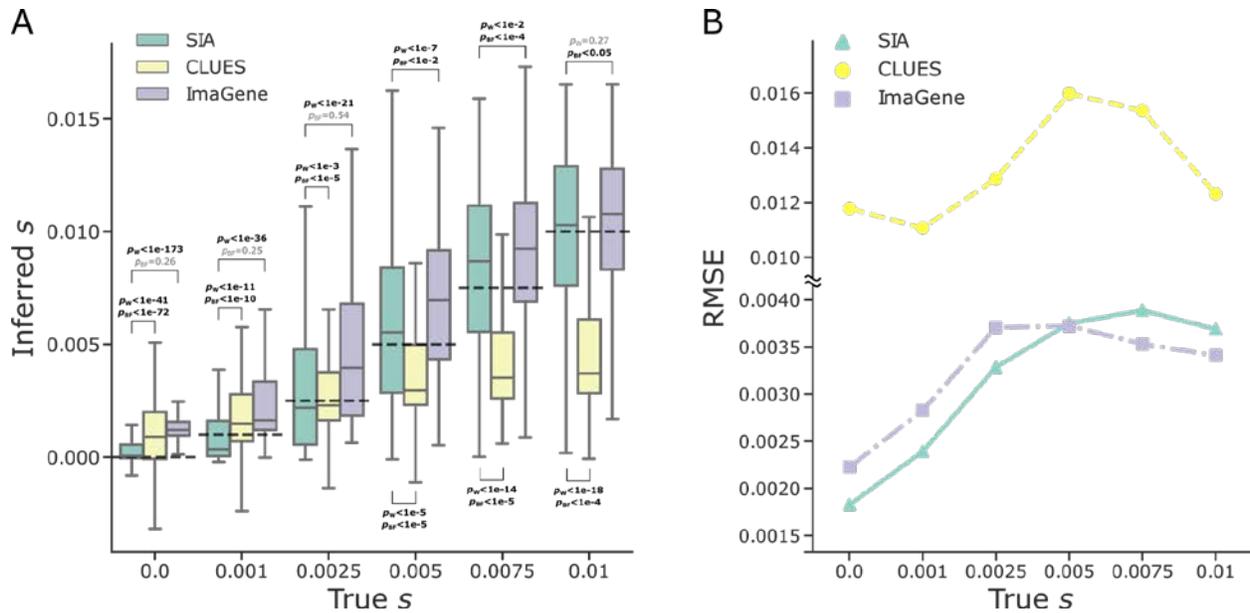
667 **Figure 1: A high-level framework for automating the detection of selective sweeps.** We
 668 first estimate the demographic history for the population of interest, then based on the estimated
 669 demographic history, we simulate neutral regions and sweeps using the discoal simulator [52].
 670 We proceed with ARG inference and then extract ARG-level statistics from each simulated
 671 region. The ARG-level statistics were used as features for a deep-learning Recurrent Neural
 672 Network (RNN) model. Finally, the learned model was applied to the empirical data to infer
 673 sweeps, selection coefficients, and allele-frequency trajectories.



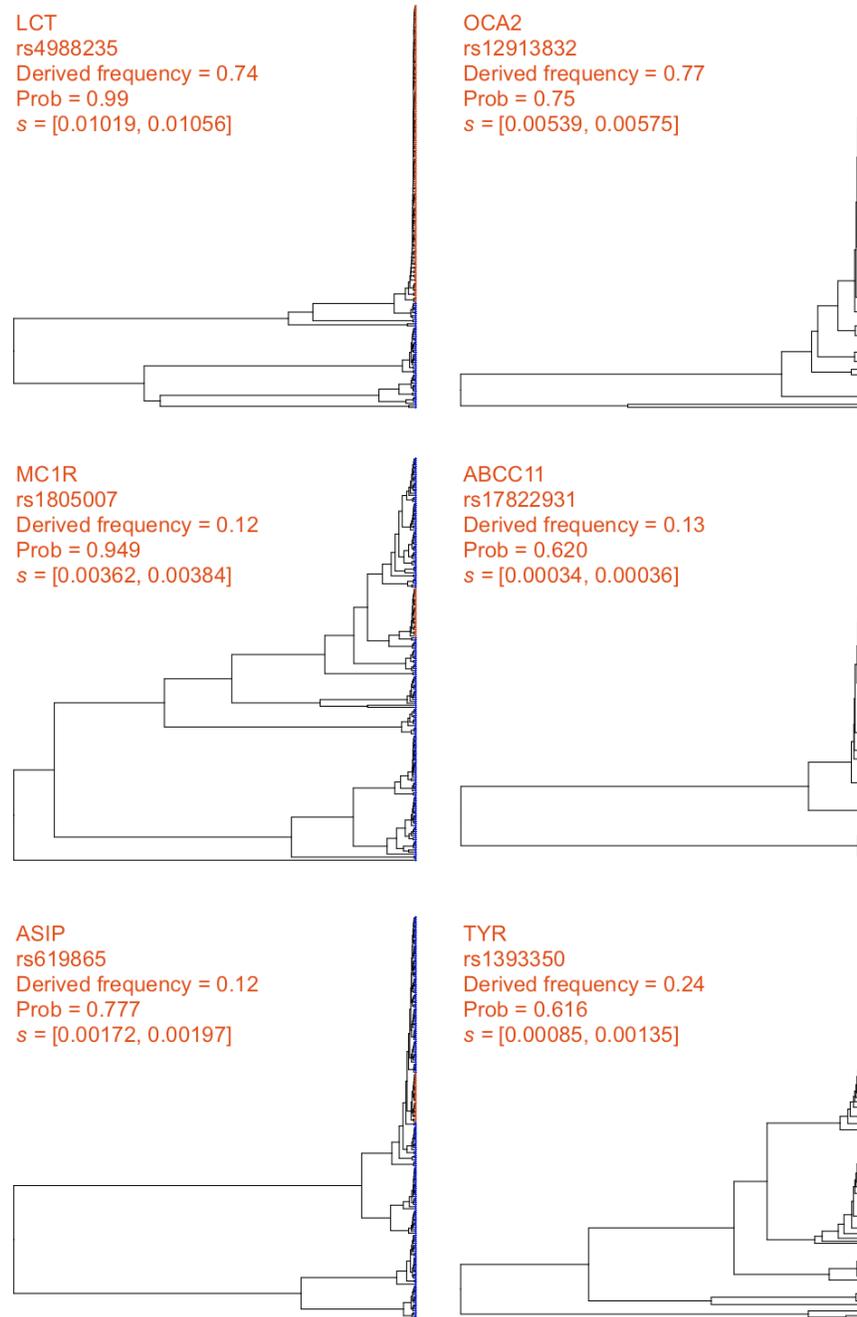
674 **Figure 2: Classification performance of SIA and other methods on simulated data.**
 675 Sequence data were simulated under a variety of selection regimes (s , shown horizontally) and
 676 derived allele frequencies (DAFs) for the beneficial mutation under selection (f , shown vertically)
 677 (see **Methods** for more details). The prediction task distinguished neutral regions and sweeps.
 678 The methods were tested on a set of 200 regions per panel (100 per class), and the receiver
 679 operating characteristic (ROC) curve records the true positive rate (TPR) as a function of the
 680 false positive rate (FPR). The curve is obtained by varying the prediction threshold from 0 to 1
 681 and recording for each threshold the number of regions correctly assigned (TPs) or misassigned
 682 (FPs) as positives (with prediction probability above the threshold). The performance of each
 683 method was evaluated based on the area under its ROC curve, or AUROC. We report each
 684 method's AUROC as an average across 200 replicate datasets for each model condition. Note
 685 that inferred genealogies were used as input to SIA.



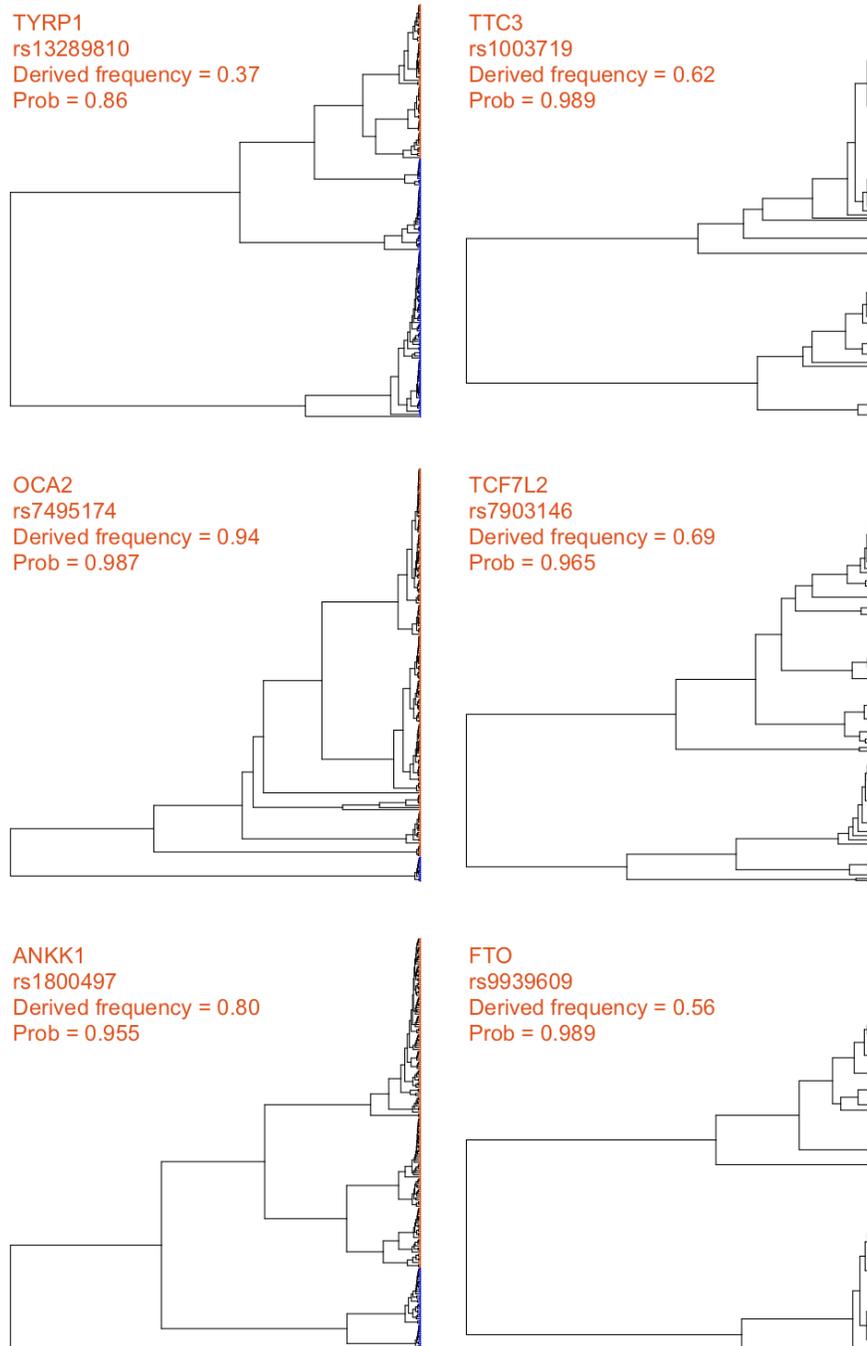
686 **Figure 3: Predictions of selection coefficients for simulated regions using SIA and**
 687 **CLUES based on true genealogies. (A)** The distribution of inferred selection coefficients for
 688 each method under each model condition are reported using a box plot. The box plot for each
 689 method reports these five statistics (from bottom to top): minimum, first quartile, median, third
 690 quartile, and maximum. The y-axis shows the inferred selection coefficient while the x-axis
 691 shows the true selection coefficient. The dashed-black line indicates the true selection
 692 coefficient for each model condition. The simulations are based on the CEU demographic model
 693 and true genealogies were used as input to both methods. Each model condition (i.e. box plot)
 694 represents a set of 400 independent simulations. The mean ranks and variances of the
 695 distributions of inferred s were compared using the Wilcoxon signed-rank test (p_W) and the
 696 Brown-Forsythe test (p_{BF}), respectively. **(B)** The root mean square error (RMSE) for each
 697 method under each model condition evaluated on 400 independent simulations.



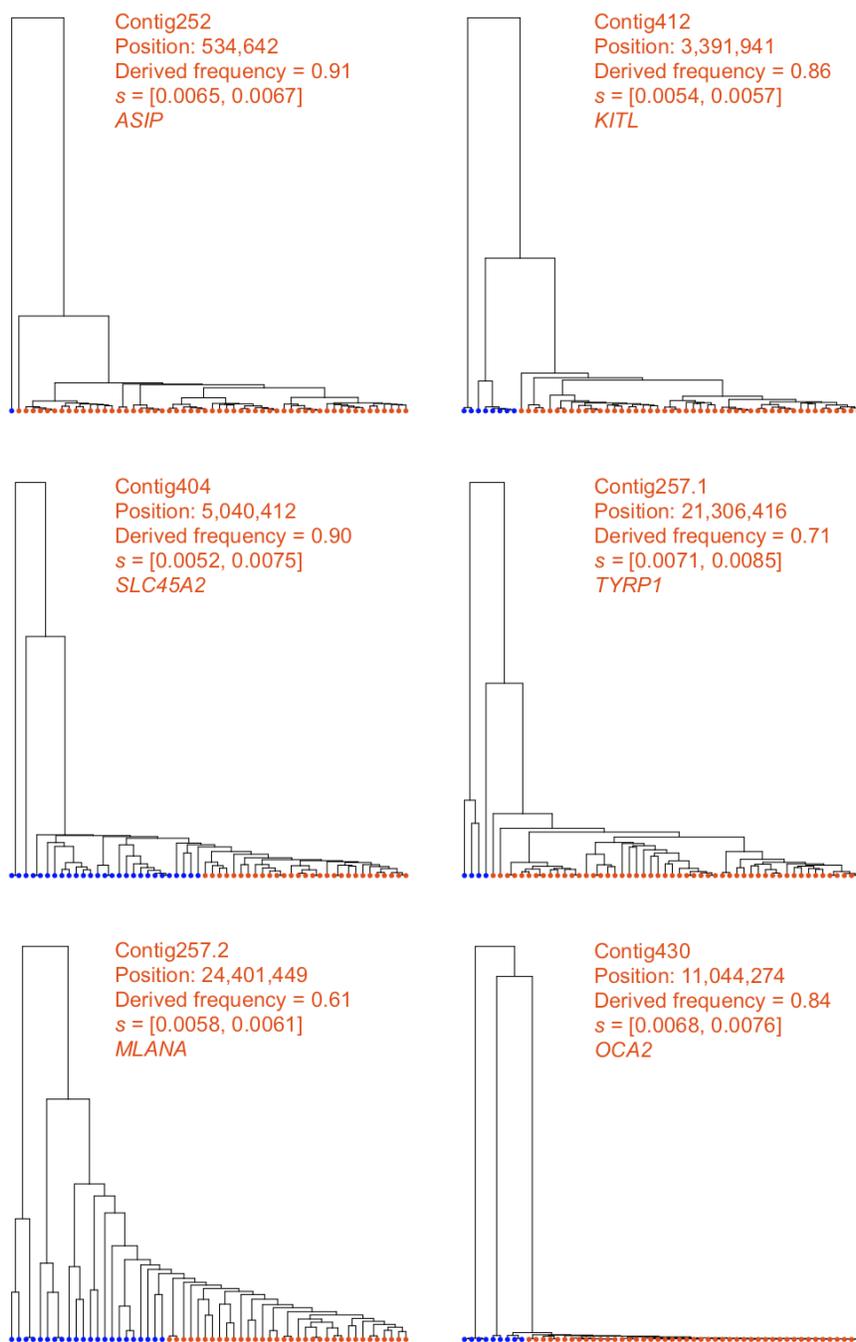
698 **Figure 4: Predictions of selection coefficient on simulated regions using SIA and CLUES**
 699 **based on inferred genealogies, and ImaGene. (A)** The distribution of inferred selection
 700 coefficients and **(B)** root mean square error (RMSE) for each method under each model
 701 condition. The simulations are based on the CEU demographic model where inferred
 702 genealogies were used as input to SIA and CLUES, whereas sequence alignments were used
 703 as input to ImaGene. Figure layout and description are otherwise similar to **Figure 3**.



704 **Figure 5: Local genealogies at six loci inferred to be under positive selection in the 1000**
705 **Genomes CEU population.** Gene name, RefSNP number, derived allele frequency, SIA-
706 inferred sweep probability and SIA-inferred selection coefficient range for each locus are
707 indicated at the top of each panel (see **Table 1** for more details). Taxa carrying the ancestral
708 and derived alleles are colored in blue and orange, respectively.



709 **Figure 6: Local genealogies at six loci lacking signal of positive selection in the 1000**
710 **Genomes CEU population.** Gene name, RefSNP number, derived allele frequency and
711 probability of neutrality inferred by SIA for each locus are indicated at the top of each panel (see
712 **Table 1** for more details). Taxa carrying the ancestral and derived alleles are colored in blue
713 and orange, respectively.



714 **Figure 7: Local genealogies at five loci inferred to be under positive selection in *S.***
715 ***hypoxantha*.** Contig name, position of SNP, derived allele frequency, SIA-inferred selection
716 coefficient range, and the pigmentation gene closest to the locus in question are indicated at the
717 top of each panel. Haplotype genomes carrying the ancestral and derived alleles are colored in
718 blue and orange, respectively.

719 **Tables**

720

721 **Table 1: List of genomic loci of interest along with their derived allele frequencies (DAF),**
 722 **sweep probabilities, and selection coefficients inferred by SIA in the 1000 Genomes CEU**
 723 **population.**

724

Gene	SNP ID	Chr	Position*	DAF	P(sweep)	Selection coefficient (95% CI)
<i>LCT</i> [36]	rs4988235	2	136608646	0.74	0.999	[0.01019, 0.01056]
<i>OCA2</i> [37,38]	rs12913832	15	28365618	0.77	0.750	[0.00539, 0.00575]
<i>MC1R</i> [37,39]	rs1805007	16	89986117	0.12	0.949	[0.00362, 0.00384]
<i>ABCC11</i> [42]	rs17822931	16	48258198	0.13	0.620	[0.00034, 0.00036]
<i>ASIP</i> [65]	rs619865	20	33867697	0.12	0.777	[0.00172, 0.00197]
<i>TYR</i> [39,65]	rs1393350	11	89011046	0.24	0.616	[0.00085, 0.00135]
<i>KITLG</i> [39]	rs12821256	12	89328335	0.13	0.869	[0.00183, 0.002]
<i>TYRP1</i> [40]	rs13289810	9	12396731	0.37	0.144	[0.00004, 0.00006]
<i>TTC3</i> [41]	rs1003719	21	38491095	0.62	0.011	[0, 0]
<i>OCA2</i>	rs7495174	15	28344238	0.94	0.013	[0, 0.00005]
<i>TCF7L2</i> [66]	rs7903146	10	114758349	0.69	0.035	[0, 0]
<i>ANKK1</i> [67]	rs1800497	11	113270828	0.80	0.045	[0, 0]
<i>FTO</i> [68]	rs9939609	16	53820527	0.56	0.011	[0, 0]

725 **Note: *Genomic coordinates in GRCh37 (hg19) assembly**

726 **Table 2: The top 25 F_{ST} peaks identified in [31] along with the number of partial soft sites**
 727 **in *S. hypoxantha* identified for each scaffold using SIA.** To avoid cases with limited power,
 728 we focused on sites with segregating frequency ≥ 0.5 , SIA-inferred $s > 0.0025$, and SIA-inferred
 729 sweep probability ≥ 0.99 .
 730

Scaffold	Start position (Mb)	End position (Mb)	Length (kb)	# of partial soft sites*
59	5.74	5.86	120	11
118	7.16	7.22	60	5
252	0.40	0.54	140	3
257.1	21.24	21.78	540	26
257.2	24.40	24.84	440	43
257.3	28.66	28.96	300	10
257.4	31.30	31.38	80	8
257.5	5.78	6.20	420	25 (1)
263	0.00	0.58	580	31
308	0.04	0.20	160	0
404.1	5.04	5.84	800	115 (7)
404.2	10.76	10.96	200	30
412	3.38	3.62	240	15
430	10.98	11.10	120	24
567	2.50	2.80	300	0
637.1	6.00	6.32	320	2
637.2	6.84	6.92	80	4
762	1.65	1.73	80	30
766	1.98	2.10	120	1
791	9.90	9.98	80	15
1717	0.92	0.98	60	7
3622	0.96	1.36	400	8
1635	3.71	3.75	40	4
1954	2.8	2.9	100	17
579	0.1	0.16	60	0

731 **Note: *The number of sweep sites in coding regions is shown in parenthesis.**