

Fast and accurate local ancestry inference with Recomb-Mix

Yuan Wei¹, Degui Zhi², and Shaojie Zhang¹

¹Department of Computer Science, University of Central Florida, Orlando, FL, USA

²McWilliams School of Biomedical Informatics, University of Texas Health Science Center at Houston, Houston, TX, USA

Corresponding authors:

shzhang@cs.ucf.edu

degui.zhi@uth.tmc.edu

Running title:

Efficient local ancestry inference

1

Abstract

2

The availability of large genotyped cohorts brings new opportunities for revealing the high-resolution genetic structure of admixed populations via local ancestry inference (LAI), the process of identifying the ancestry of each segment of an individual haplotype. Though current methods achieve high accuracy in standard cases, LAI is still challenging when reference populations are more similar (e.g., intra-continental), when the number of reference populations is too numerous, or when the admixture events are deep in time, all of which are increasingly unavoidable in large biobanks. Here, we present a new LAI method, Recomb-Mix. Recomb-Mix integrates the elements of existing methods of the site-based Li and Stephens model and introduces a new graph collapsing trick to simplify counting paths with the same ancestry label readout. Through comprehensive benchmarking on various simulated datasets, we show that Recomb-Mix is more accurate than existing methods in diverse sets of scenarios while being competitive in terms of resource efficiency. We expect that Recomb-Mix will be a useful method for advancing genetics studies of admixed populations.

3

4

5

6

7

8

9

10

11

12

13

14 Introduction

15 Local ancestry inference (LAI) is a process of assigning the ancestral population labels of each segment
16 on an individual's genome sequence. LAI is not only useful for better study of human demographic
17 history (Martin et al. 2017) but also can enable several downstream tasks, including admixture map-
18 ping (Reich et al. 2005), ancestry-aware genome-wide association studies (GWAS) (Pasaniuc et al.
19 2011), and ancestry-specific polygenic risk scores (Duncan et al. 2019). Recent studies show that local
20 ancestry information improves the resolution of association signals in GWAS (Atkinson et al. 2021),
21 helping to infer the high-resolution of genomic regions containing genes as under selection (Hamid
22 et al. 2023). Local ancestry calls contribute to understanding the impact of genetic variants that cause
23 disease (Hou et al. 2023) and the accuracy of polygenic scores of genetically based predictions (Ding
24 et al. 2023).

25 The recent availability of biobank-scale genotyped datasets (Bycroft et al. 2018; Kurki et al. 2023)
26 and the rising of enormous databases from direct-to-consumer genetic companies (Durand et al. 2021;
27 Wang et al. 2021) create new challenges and opportunities for LAI. Participants in biobanks may
28 be from highly imbalanced source populations. Inferring ancestral components underrepresented in
29 these reference populations is in great need. On the other hand, the admixture to be inferred in the
30 samples may be multi-ways and may be from both recent and distant admixture events. However,
31 opportunities coexist with such challenges. More diverse samples, e.g., the Human Genome Diversity
32 Project (HGDP) (Bergström et al. 2020), are becoming available as reference panels. The number of
33 participants in biobanks is much larger and more representative than in previous reference panels, and
34 thus, more potential sub-continental ancestral information becomes available from biobanks. Although
35 methods for revealing sub-continental or even sub-population clusters are available (e.g., (Lawson et al.
36 2012)), they are mostly non-LAI methods and only capture global ancestry. With the availability of
37 diverse samples in biobanks and the need for in-depth knowledge of admixed individuals, current LAI
38 methods are unable to capitalize on these opportunities fully.

39 Existing LAI methods fall into two categories: site-based and window-based. Originally, Hap-
40 Mix (Price et al. 2009), based on the extended Li and Stephens (LS) Hidden Markov model (Li and
41 Stephens 2003), is developed to model different transition probabilities for within population and be-

42 tween population jumps. To make it tolerate mismatches, emission probabilities are also introduced,
43 and the model is not very efficient and cannot be scaled up (Geza et al. 2018). Later, window-based
44 methods are gaining popularity (e.g., RFMix (Maples et al. 2013), G-Nomix (Hilmansson et al. 2021),
45 SALAI-Net (Oriol Sabat et al. 2022)). These methods take short stretches of sites as windows and
46 define window-based features. It first makes local ancestry prediction over each window and then
47 uses certain post-processing to smooth out the labels across all windows. For each window, a certain
48 machine-learning approach is typically used. However, the predefined window boundaries are not nec-
49 essarily optimized, and the noisy initial window labels can be difficult to correct by post-processing.
50 Loter (Dias-Alves et al. 2018) is a recent site-based method under the LS framework. It formulates
51 the LAI problem as a combinatorial best-path problem in a graph, which can be solved efficiently by
52 dynamic programming. However, its problem formulation is simplistic in that it does not take into
53 account the useful information encoded in the LS model, such as differential transition probabilities
54 for within and between populations and variable recombination rates across sites. Loter reported that
55 it underperformed RFMix (Maples et al. 2013) and LAMP-LD (Baran et al. 2012) for datasets with
56 recent admixture events (i.e. < 150 generations). Therefore, there is room for improvement over the
57 Loter approach by introducing an LS-inspired parametrization of its scoring function.

58 In this study, we developed Recomb-Mix, a site-based method that is both accurate and efficient.
59 Our main insight is that we do not have to have an exact LS formulation. The gist of the HapMix
60 LS model is the differential transition penalties for within and between populations and assigning the
61 population labels for a site by comparing the paths going through it versus the paths by-passing it.
62 We achieve the same spirit by setting the within-population transition penalty to zero and collapsing
63 the nodes representing the allele values at each site. This allows both run time and space efficiency
64 while achieving higher accuracy than Loter. Recomb-Mix is designed to have robustness, scalability,
65 and superior accuracy on LAI. Applications to real human datasets confirmed the genetic differences
66 among populations and provided potential explanations for how the evolutionary processes shape these
67 differences.

68 **Methods**

69 **The Li and Stephens framework for site-based local ancestry inference**

70 Recomb-Mix is inspired by HapMix (Price et al. 2009) and Loter (Dias-Alves et al. 2018), all of which
71 are under the Li and Stephens (LS) framework. These two methods extend the basic LS model (Li
72 and Stephens 2003) to capture the difference between inter-population and intra-population transition
73 probabilities. The LS model defines the conditional probability of any haplotype sequence given a set
74 of haplotypes in a panel as a Hidden Markov Model (HMM). The states of the HMM are individual
75 sequence positions in the panel. By treating all haplotypes equally, transition probability in the HMM
76 can be specified by just the probability of switching a haplotype template or staying at the same
77 template. In a non-probabilistic combinatorial formulation, transition probability can be modeled as
78 a template change penalty. HapMix extends the model to have population labels as augmented HMM
79 state labels and introduces two transition probability parameters, i.e., small-scale (between haplotypes
80 from within a reference population) and large-scale (between the reference populations). See equations
81 (0.1) and (0.2) in the HapMix paper for the detailed definition of the population-label-aware LS model.
82 Loter formulates LAI as a graph optimization problem that finds a best-scoring path over a site-level
83 graph. It can be viewed as an LS “copying model” with simplified non-probabilistic parameterization.
84 Loter applies the same penalty to haplotype template switches in both cases within or across reference
85 populations.

86 Through the unified view of the LS framework and the graph optimization formulation (Table 1),
87 Recomb-Mix introduces special parameterization to the LS model to induce graph simplification and
88 more biologically relevant scoring function. First, by assuming no template change penalty when
89 switching haplotype templates within a reference population, Recomb-Mix enables the collapsing of
90 the reference panel to a compact population graph. Generating a compact population graph greatly
91 reduces the size of reference populations and retains the ancestry information, as genetic markers
92 having the same allele values per population are collapsed in the compact population graph. Different
93 template change penalties are used when switching haplotype templates within a reference population
94 and between the reference populations. The template change penalty within a reference population
95 is set to zero, and recombination rates from a genetic map parameterize the template change penalty

Scoring Function	HapMix	Loter	Recomb-Mix
Within population transition probability / template change penalty	Recombination parameter	1	0
Across population transition probability / template change penalty	Recombination parameter and miscopying parameter	1	r , recombination penalty
Emission probability / mismatch penalty	Mutation parameter	Mismatch penalty	d , mismatch penalty
Parameterization of transition probability	$1 - e^p$, where p is the product of genetic distance and recombination parameter	Bootstrap aggregation (bagging)	w , weight

Table 1: Comparison of scoring functions in HapMix (Price et al. 2009), Loter (Dias-Alves et al. 2018), and Recomb-Mix (this study) under the Li and Stephens (LS) framework.

96 between the reference populations. Second, Recomb-Mix’s scoring function $\sum(d + w \cdot r)$ is a simplified
97 version of HapMix’s. Still, it is a richer version than that of Loter’s. In Recomb-Mix’s scoring function,
98 d is a mismatch penalty score regarding a site on the query haplotype and the corresponding site in the
99 reference panel, w is the weight for the relative importance of the mismatch cost and the recombination
100 cost, and r is a normalized recombination rate penalty score between two sites translated from the
101 genetic distance. Table 1 shows the differences in scoring functions between some LAI methods under
102 the LS framework. HapMix incorporates recombination, miscopying, mutation, and genetic distance
103 parameters into its scoring function. With such a number of required biological parameters, lacking
104 accurate population information may lead to biased inference results; that is, the estimated biological
105 parameters required for HapMix may not be the correct parameters for the given dataset, and the
106 ancestry inference result based on such parameters is influenced (Patin et al. 2014; Suarez-Pajes et al.
107 2021). In contrast to HapMix’s complex scoring function, Loter uses a simple one that does not adopt
108 any recombination information. Recomb-Mix takes Loter’s simplicity but adds the notion of genetic
109 distance to encode genetic information, to achieve high computability and accuracy simultaneously.

110 Besides the scoring function, Recomb-Mix has other differences from HapMix and Loter. Recomb-
111 Mix and Loter can handle multi-way admixture inference, as HapMix can only tackle two-way admix-
112 ture. Both Recomb-Mix and HapMix use genetic map (Church et al. 2011; Schneider et al. 2017) to
113 help out the inference, while Loter does not take any biological information as input. From the HMM
114 algorithm perspective, HapMix uses the forward-backward algorithm to update transition and emission
115 probabilities and then estimate the hidden ancestral states (Price et al. 2009; Wu et al. 2021). Loter’s
116 approach minimizes an objective function using dynamic programming, a Viterbi-like algorithm (Dias-
117 Alves et al. 2018; Oriol Sabat et al. 2022). Recomb-Mix takes advantage of the graph optimization
118 formulation as Loter does but keeps population-level information only in a compact population graph.
119 All possible paths in the graph can be viewed as a set of “combined” paths from the original graph
120 that emit the same population label readout. Thus, it has a flavor of a “forward-like” algorithm as an
121 ancestry label is assigned to an individual node according to “combined” paths passing through it.

122 **Recomb-Mix**

123 The Recomb-Mix method is inspired by the LS HMM and implemented using a graph optimization
124 approach. Like the LS model, it assumes that an admixed individual haplotype is modeled as a mosaic
125 of individual haplotypes from a reference panel. Recomb-Mix constructs a population graph from a
126 given reference panel to infer the ancestral label at each locus on a given admixed individual haplotype
127 by finding a threading path that resembles the admixed individual haplotype the most among all the
128 paths. In the population graph, the allele values of individual haplotypes are grouped by each site.
129 Then, the population graph is transformed into a compact population graph by collapsing the site nodes
130 with the same allele value and ancestral label into one node. A mismatch penalty at each site occurs
131 when there is a difference between the collapsed site nodes’ allele value and the corresponding site’s
132 allele value in the admixed individual haplotype. The collapsed site nodes are linked to population
133 emission nodes based on the ancestral label of each node, and the population emission nodes make
134 a cross-population connection to the population emission nodes of the next site. A template change
135 penalty regarding recombinations between each site occurs when two population emission nodes of
136 adjacent sites having different ancestral labels are connected. Then, the population emission nodes

137 are expanded to genotype emission nodes linked to the site nodes for the next site. This process is
138 similar to a “forward-like” approach as each node in the compact population graph can be viewed as a
139 bundle of nodes in the original population graph being consolidated as one, and their ancestral labels
140 are assigned by all low penalty threading paths passing through it. Recomb-Mix sums over mismatch
141 penalties and template change penalties through all possible threading paths to determine the one that
142 has the minimum penalty score. The ancestral label of each site in the admixed individual haplotype
143 is assigned the same ancestral label of each corresponding node on such threading path.

144 To formally define the Recomb-Mix method, a reference panel having m individual haplotypes with
145 n sites can be transformed as a population graph $G = (V, E)$, representing the HMM of this set of
146 haplotypes. V is the union of the starting and the ending nodes $\{s, e\}$ and the site nodes S_j for
147 $j \in [1, n]$. $S_j = \{s_j^1, s_j^2 \dots s_j^m\}$ is the set of nodes representing alleles of haplotypes at position j . E is
148 the union of the edges from s_j^i to s_{j+1}^k for all $j \in [1, n - 1]$, $i, k \in [1, m]$, and the edges from s to s_1^i
149 and s_n^i to e for $i \in [1, m]$. s_j^i is the node that the i -th haplotype has at site j . Each node at every site
150 of every haplotype has an associated ancestral label. It is assumed there are p populations presented
151 in the reference panel, and each population has an ancestral label in $[1, p]$. The ancestral label of node
152 s_j^i is $l(s_j^i) \in [1, p]$. The allele value of node s_j^i is $a(s_j^i) \in [0, 1]$, assuming all sites are bi-allelic. The
153 population graph G is further transformed into a compact population graph $G' = (V', E')$ by collapsing
154 all nodes with the same allele value and ancestral label to one node in every site. V' is the union of
155 the starting and the ending nodes $\{s, e\}$ and the site nodes $S'_j = \{s'^1_j, s'^2_j, \dots, s'^{|S'_j|}_j\}$ for $j \in [1, n]$. S'_j
156 is a set of nodes representing all unique pairs of allele values and ancestral labels in S_j (i.e., there is a
157 node $s'^i_j \in S'_j$ if and only if there is a node $s \in S_j$ such that $a(s) = a(s'^i_j)$ and $l(s) = l(s'^i_j)$ and for all
158 $k \in [1, |S'_j|]$, $k \neq i$, $a(s'^i_j) \neq a(s'^k_j)$ or $l(s'^i_j) \neq l(s'^k_j)$). E' is the union of the edges from u_1 to u_2 for all
159 $u_1 \in S'_j$ and $u_2 \in S'_{j+1}$ for $j \in [1, n - 1]$, and the edges from s to u_1 and u_2 to e for all $u_1 \in S'_1$ and
160 $u_2 \in S'_n$.

161 To calculate all possible threading paths in G' for a query admixed individual haplotype Q with n
162 sites, $Q = (q_1, q_2, \dots, q_n)$ (the allele value of q_i is $a(q_i) \in [0, 1]$), Recomb-Mix incorporates the mismatch
163 penalty and the template change penalty into its objective function. The mismatch penalty function is
164 defined as $d(x_1, x_2)$, where x_1 and x_2 are allele values. The template change penalty function is defined
165 as $r(y_1, y_2)$, where y_1 and y_2 are ancestral label values. Then, the cost of a candidate threading path

166 $P = (u_1, u_2, \dots, u_n)$ is defined as:

$$f(P) = \sum_{j=1}^n d(a(q_j), a(u_j)) + w \sum_{j=1}^{n-1} r(l(u_j), l(u_{j+1})). \quad (1)$$

167 In Equation 1, w is a scale factor to balance the mismatch and recombination costs. Let P^* be the
168 threading path having the minimum penalty cost among all candidate threading paths in G' . The
169 ancestral labels of the nodes in P^* from s to e are the estimated ancestral labels of sites in Q , that is
170 $(l(u^*_1), l(u^*_2), \dots, l(u^*_n))$. Thus, LAI can be formulated as a problem to find P^* in G' .

171 Figure 1 is an example of local ancestry inference with Recomb-Mix. A reference panel having seven
172 individual haplotypes, eight sites, and two ancestral labels (shown in red and blue) is represented as
173 a population graph G . Nodes representing each site are fully connected to nodes representing their
174 adjacent site. A node s is connected to all nodes for the first site, and all nodes for the last site are
175 connected to a node e . Q is a query of an admixed individual haplotype. G is then transformed into
176 a compact population graph G' , and a threading path having the minimum penalty score is selected
177 from node s to node e in G' , to be used to paint the admixed individual haplotype query by assigning
178 the estimated ancestral label to each site in Q . Figure 1B is an example to show why G' is still an LS
179 model. It demonstrates how nodes representing sites in positions three and four in G are transformed
180 into the corresponding nodes in G' . First, the set of nodes representing alleles of haplotypes at position
181 three (the first column) is freely collapsed to genotyping emission nodes (the second column) with the
182 same allele value and ancestral label to one node. Second, the genotyping emission nodes are linked to
183 population emission nodes (filled red and blue nodes in the third column) according to their ancestral
184 labels. Then, the population emission nodes at position three make cross-population connections to
185 the population emission nodes at position four (filled red and blue nodes in the fourth column). A
186 penalty r is applied to the connections of population emission nodes when their ancestral labels differ.
187 Finally, the population emission nodes at position four are linked to the set of nodes representing alleles
188 of haplotypes at position four (the fifth column), and those nodes are freely collapsed to genotyping
189 emission nodes (the sixth column).

190 Recomb-Mix uses a simplified scoring function (Equation 1) like the one HapMix uses to calculate
191 the penalty score of each threading path. The mismatch penalty is determined by simply comparing

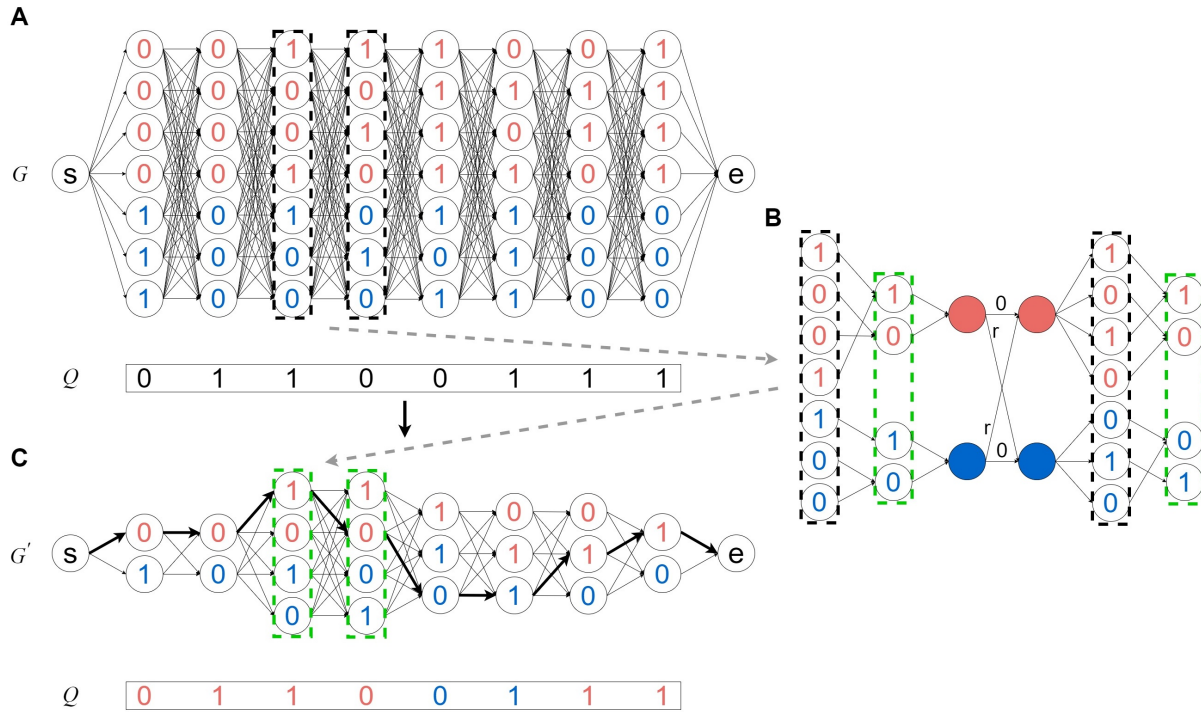


Figure 1: An example of local ancestry inference with Recomb-Mix. (A) G is a population graph representing the HMM in Recomb-Mix, constructed from a given reference panel. G contains seven haplotypes with eight sites belonging to two populations (shown in red and blue). Q is a query of an admixed individual haplotype. (B) A transformation process from nodes in sites three and four in G to nodes in the corresponding sites in G' . Nodes in black boxes correspond to the nodes in sites three and four in G . Nodes in green boxes correspond to the nodes in sites three and four in G' . The filled nodes in red and blue are population emission nodes in sites three and four. r is a cross-population penalty. (C) G' is a compact population graph transformed from G . Q is assigned with estimated ancestral labels for each site (shown in red and blue on allele values), according to a threading path selected with minimum penalty score (shown as bold edges) in G' .

192 the allele values of each site in Q to the corresponding nodes in G' . For each site, a mismatch penalty
193 is applied if the allele values are not the same. $d(\cdot)$ in Equation 1 is implemented as $d(a(q_j), a(u_j)) = 0$
194 if $a(q_j) = a(u_j)$; otherwise $d(a(q_j), a(u_j)) = 1$ for site j of u_j in G' and q_j in Q . The template
195 change penalty is determined by the recombination rate between site j and $j + 1$ and the ancestral
196 labels of u_j and u_{j+1} in G' . Since the recombination rate between two sites is inversely proportional
197 to the probability of the edge connecting these sites being a recombination breakpoint, the template
198 change penalty cost is determined by the reciprocal of the recombination rate between the two sites.
199 To leverage the linkage disequilibrium (LD) effect (haplotype information of allele correlations) from
200 the LS model, the template change penalty is applied to the edge connecting two adjacent nodes if
201 they have different ancestral labels. No template change penalty is applied if two adjacent nodes on a
202 threading path share the same ancestral label. This setting allows the representation of a more diverse
203 set of haplotypes than those explicitly listed in the panel. Not giving any penalties to them allows all
204 possible threading paths jumping between templates within a population in G to be treated equally
205 in G' . Of course, this setting is quite simplistic: it could risk allowing too much diversity. Also, it is
206 possible to set the within-population template change penalty to a value other than zero, or some more
207 sophisticated settings. However, setting this to zero captures the main idea of differentiating within-
208 versus across-population transition probability. $r(\cdot)$ in Equation 1 is implemented as $r(l(u_j), l(u_{j+1})) =$
209 0 if $l(u_j) = l(u_{j+1})$; otherwise $r(l(u_j), l(u_{j+1})) = R_{j,j+1}$ for site j and $j + 1$ of u_j and u_{j+1} in G' . $R_{j,j+1}$
210 is the normalized reciprocal of the recombination rate between site j and $j + 1$. To calculate $R_{j,j+1}$,
211 min-max normalization is used to scale the range of the recombination rates from the genetic map
212 into $[0, 1]$. The normalized recombination rates $\mathcal{R}_{j,j+1}$ are further processed by applying a reciprocal
213 function to obtain $R_{j,j+1}$, where $R_{j,j+1} = 2/(\mathcal{R}_{j,j+1} + 1)$. The range of the normalized reciprocal of
214 the recombination rates is $[1, 2]$. The normalization step is to ensure the template change penalty is in
215 the same order of magnitude as the mismatch penalty to prevent the domination of any penalties.

216 Representing the LS model as a compact population graph is efficient in terms of time and space.
217 Recomb-Mix uses a dynamic programming approach to solve the problem of finding P^* in G' . Using G'
218 instead of G to compute the minimum penalty score among all possible threading paths significantly
219 improves the computing time. The maximum indegree and outdegree of any node in G' is a constant
220 value $2p$, assuming all sites are bi-allelic and the number of populations presented in the reference panel

221 is small (i.e., $p \ll m$). For each site, each population’s minimum mismatch penalty score is tracked
222 alongside the minimum mismatch penalty score over all the populations. The candidate threading
223 paths of each node having the current minimum penalty score can be determined by comparing two
224 candidate penalty scores of the adjacent site, which are the minimum penalty score whose population
225 is the same and the minimum penalty score over all the populations. Thus, the time complexity of
226 computing the penalty scores on G' is $O(np)$. Using G' also substantially alleviates the demand for
227 spaces to store a large reference panel. The space complexity is reduced from $O(nm)$ to $O(np)$, as G'
228 stores at most $2p$ nodes per site. Using a compact population graph to reduce space usage is similar
229 to but different from existing approaches in phasing and imputation. In a popular phasing method
230 SHAPEIT (Delaneau et al. 2019), each node represents its allele value and each edge represents the
231 weight of the number of individual haplotypes in its reference panel. This approach makes SHAPEIT
232 have a space complexity of $O(nj)$ (j is the number of conditioning states for each marker), which
233 helps speed up its subsequent HMM calculation (Delaneau et al. 2012). Likewise, another popular
234 phasing and imputation method, Beagle (Browning et al. 2018), constructs its HMM state space from
235 its reference panel by leveraging composite reference haplotypes. The same haplotype compression
236 technique is later adopted by FLARE (Browning et al. 2023). For each query individual haplotype,
237 FLARE finds Identity-By-State (IBS) segments using Positional Burrows-Wheeler Transform (PBWT)
238 algorithm (Durbin 2014) from the reference haplotypes and then makes composite reference haplotypes
239 by stitching those IBS segments. Utilizing the composite reference haplotypes costs a $O(nm')$ space
240 complexity, where m' depends on the number and locations of IBS segments of the query individual
241 haplotype against the original reference panel. Usually, m' is relatively small, as it is expected there
242 exist many long IBS segments.

243 Another advantage of presenting a reference panel as a compact population graph is that Recomb-
244 Mix can process the reference panel regardless of whether the panel is phased or not. When converting
245 a reference panel into a compact population graph, the order of the sites from two haplotypes of an
246 individual is irrelevant, thanks to a property that the out-neighborhood of a node u in a graph is the
247 set of nodes adjacent to u . Thus, Recomb-Mix is flexible to handle both phased and unphased reference
248 panels.

249 Simulated datasets

250 To evaluate the performance of Recomb-Mix, several admixture datasets were simulated using SLiM
251 v4.0 (Haller and Messer 2019). The input data were individuals in Whole-Genome Sequencing (WGS)
252 form from various populations in the study of the 1000 Genomes Project (TGP) (Auton et al. 2015;
253 Clarke et al. 2016) and the Human Genome Diversity Project (HGDP) (Bergström et al. 2020). Each
254 input population was split into disjoint sets of founders and references. The admixture population was
255 simulated as the descendants of admixed founders from different populations. The admixed individuals
256 were sampled from the admixture population, and the referenced individuals were sampled from each
257 reference population. A set of three-way inter-continental datasets of Chromosome 18 were simulated
258 using YRI, CEU, and CHB individuals (representing African, European, and Asian populations; more
259 descriptions of the populations are available in Supplemental Table S1) from the TGP dataset. Various
260 sizes of reference panels (i.e., 100, 250, 500, and 1,000) and numbers of generations after the admixture
261 event (i.e., 15, 50, 100, and 200) were examined. The average recombination rate and mutation rate
262 used for the simulation was $1.46455e-08$ and $1.29e-08$ per base pair per generation, according to the
263 stdpopsim library (Adrion et al. 2020). A 0.02% genotyping error, following Browning et al. (2023), was
264 added to admixed and reference individuals. The ground truth ancestral labels of admixed individuals
265 were extracted from the SLiM output tree sequence (Haller et al. 2019). Additionally, a set of seven-way
266 inter-continental datasets of Chromosome 18 were simulated using AFR, EAS, EUR, NAT, OCE, SAS,
267 and WAS individuals (representing African, East Asian, European, American, Oceanian, Central and
268 South Asian, and West Asian populations) from the HGDP dataset. The HGDP dataset was phased
269 and imputed using Beagle 5.4 (Browning et al. 2018, 2021). The goal of simulating the seven-way
270 admixture is to explore how well LAI methods are able to distinguish local ancestral segments from
271 the admixture of a large number of ancestral populations.

272 To explore the power of LAI at the intra-continental level, a set of intra-continental datasets were
273 simulated using TSI, FIN, and GBR individuals (representing Italian, Finnish, and British populations)
274 using the same settings as the inter-continental datasets. To explore the influence of the uneven
275 proportion of individuals per population in founders and references, two variations of the three-way
276 15-generation intra-continental datasets with uneven founders (i.e., 68 Italian, 32 Finnish, and 100

277 British individuals) or uneven references (i.e., 170 Italian, 80 Finnish, and 250 British individuals) were
278 simulated. Both cases were 1/3, 1/6, and 1/2 individuals to the entire population panel. Additionally,
279 an experiment was conducted to test the case when the reference panel size was ultra-small, i.e., the
280 reference panel size was 20 and 50 (or only about 7 or 17 individuals per population in the reference
281 panel).

282 **Benchmark setup**

283 Two conventional measurements were used to evaluate the performance of LAI methods. The squared
284 Pearson's correlation coefficient r^2 value (used by FLARE, LAMP-LD (Baran et al. 2012), and MO-
285 SAIC (Salter-Townshend and Myers 2019)) and the accuracy rate of the correctly-predicted markers
286 (used by G-Nomix, Loter, RFMix, and SALAI-Net). The r^2 value was followed by LAMP-LD's defini-
287 tion (Baran et al. 2012), in which the r^2 value is defined as the one between the true and the inferred
288 number of alleles from each of the populations, averaged over all the populations. The criteria used
289 by FLARE that markers were filtered with minor allele frequency ≤ 0.005 and minor allele count \leq
290 50 (Browning et al. 2023) was also applied. r^2 values are mainly reported in the benchmarks but
291 accuracy rates are also available, mostly in the supplemental. It is found that the results of r^2 values
292 and accuracy rates are often consistent.

293 Recomb-Mix was tested against several datasets with the following LAI methods: FLARE (Brown-
294 ing et al. 2023), G-Nomix (Hilmarrsson et al. 2021), Loter (Dias-Alves et al. 2018), RFMix (Maples
295 et al. 2013), and SALAI-Net (Oriol Sabat et al. 2022). The weight parameter w was tuned to $w = 1.5$,
296 which provides the best performance of Recomb-Mix for the given datasets. Parameters used for other
297 methods are available in Supplemental Table S2. FLARE is the most recent proposed LAI method.
298 It is a site-based generative method under the LS framework, extended from HapMix (Price et al.
299 2009), which models the hidden local ancestry of each site. FLARE shows encouraging speed and
300 accuracy because its performance is optimized for computation resources, and its model is designed
301 to be flexible for optional parameters to deal with various situations. G-Nomix is a window-based
302 discriminative approach using two-layer prediction to perform LAI, in which it uses Logistic Regression
303 as a base model and XGBoost as a smoother model. It is currently the leading LAI method due to its

304 promising speed and accuracy. Loter frames the ancestry prediction problem as a graph optimization
305 problem. It is prioritized for LAI on distant admixture events and good for non-model species as no
306 biological information is required. RFMix, another window-based discriminative approach, is one of the
307 popular LAI methods. It applies a conditional random field model to LAI, particularly using random
308 forest classification. RFMix shows a robust performance on multi-way admixture datasets. SALAI-Net
309 is also a window-based discriminative approach developed from its predecessor LAI-Net (Montserrat
310 et al. 2020), the first neural network-based LAI. It uses a reference matching layer and a smoother layer
311 (i.e., a combination of cosine similarity score and neural network) to perform LAI. In addition to the
312 adoption of GPU hardware, SALAI-Net utilizes a pre-trained generalized model, making it free from
313 re-training and parameter tuning during the inference process; thus, it is considered to be very fast.
314 The above LAI methods were selected for benchmarking because FLARE, G-Nomix, and SALAI-NET
315 are the newly published ones, reporting promising results; RFMix is the most popular and widely
316 used for ancestry-related applications. Loter has the same problem formulation as Recomb-Mix does.
317 HapMix was not included because it cannot tackle multi-way admixture and only produces inference
318 results at the diploid level.

319 **Results**

320 **Local ancestry inference for three-way inter-continental admixed populations**

321 For three-way inter-continental simulated datasets, Recomb-Mix had the best r^2 values and accuracy
322 rates in reference panel sizes 100, 250, 500, and 1,000 with 15 generations and in generations 15, 50,
323 100, and 200 with 500 references. Figure 2 shows the r^2 values of the inference results on six LAI
324 methods, FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net. Supplemental Figures S1
325 and S2 show the accuracy rates (values are in Supplemental Tables S5 and S6). Overall, as the reference
326 panel size increases, the average r^2 value and the accuracy rate increase for all methods. The large
327 reference panel containing more individual samples than those in small ones helps improve the inference
328 result. When using a reference panel with 1,000 individuals, all methods had at least 0.99 r^2 value or
329 92% accuracy rate, while Recomb-Mix reached the best r^2 value 0.9989 or accuracy rate of 99.10%.

330 Recomb-Mix performed well for a small reference panel (i.e., for a 100-individual panel it achieved the
331 best performance, that is r^2 value 0.9919 or accuracy rate of 97.96%). G-Nomix and SALAI-Net had
332 the second and the third highest r^2 values, which were 0.9681 and 0.9480. SALAI-Net and G-Nomix
333 had the second and the third highest accuracy rates, which were 86.69% and 86.63%.

334 Recomb-Mix's performance on ultra-small reference panels (i.e., size 20-50) was tested. Such cases
335 are interesting because small reference panels can benefit low-resourced populations. Meanwhile, LAI
336 with small reference panels is challenging because allele frequencies and haplotype frequencies are
337 noisy. For a small reference panel size of 20, Recomb-Mix achieves the best accuracy rate of 62.85%;
338 for a 50-individual reference panel, Recomb-Mix achieves 94.45%, while other methods' accuracy rates
339 are around 60% to 70% (see Supplemental Figure S3 and Table S7). We include MOSAIC (Salter-
340 Townshend and Myers 2019) in this experiment as it reportedly performs well on small reference
341 panels (Browning et al. 2023). MOSAIC achieves better accuracy rates than FLARE, Loter, and
342 RFMix on reference panels of sizes 50 and 100, but its performance is worse than Recomb-Mix.

343 Multi-way admixture

344 Besides the experiments on three-way admixed individuals, a case study on seven-way admixed indi-
345 viduals was investigated. The goal is to find out how LAI methods perform on individuals admixed
346 from a large number of founder populations, as in a real case scenario, human individuals are involved
347 in multiple population admixture events. Seven-way inter-continental datasets with various reference
348 panel sizes and generations were simulated, and for such challenging datasets, the r^2 values and the
349 accuracy rates of LAI methods dropped, but Recomb-Mix kept performing well. Figure 3 shows the
350 r^2 values of the inference results on six LAI methods, FLARE, G-Nomix, Loter, Recomb-Mix, RFMix,
351 and SALAI-Net. The average accuracy rates are in Supplemental Figures S4 and S5 (values are in
352 Supplemental Tables S10 and S11). Figure 4B illustrates an inferred haplotype sample of an admixed
353 individual from the methods. Compared to Figure 4A, more population labels were mistakenly as-
354 signed, as inferring seven-way admixed individuals is a harder task than inferring three-way ones. In
355 general, window-based LAI methods performed better than site-based ones, except Recomb-Mix. Using
356 a window as the smallest unit of inference helps tolerate errors within the window since the population

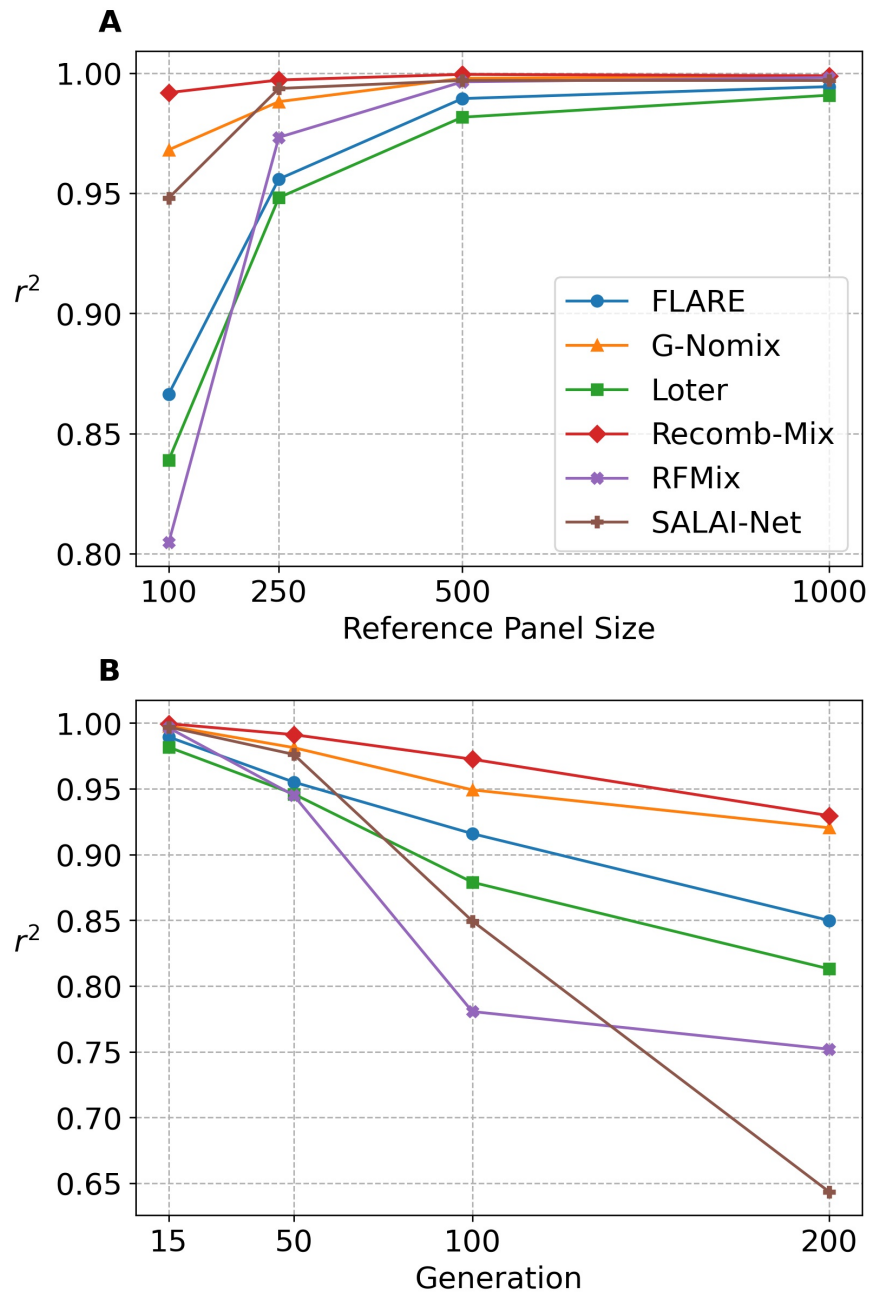


Figure 2: The squared Pearson's correlation coefficient r^2 of three-way inter-continental simulated datasets on FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net. Markers were filtered with minor allele frequency ≤ 0.005 and minor allele count ≤ 50 . (A) The three-way 15-generation datasets with the reference panel sizes 100, 250, 500, and 1,000 (values are in Supplemental Table S3). (B) The three-way 500-reference datasets with the generations 15, 50, 100, and 200 (values are in Supplemental Table S4).

357 label having the highest estimated probability determines the inference result for the entire window.
358 On the other hand, site-based methods may focus more on a single site's label, which may affect the
359 inference result of one's surrounding region when making incorrect inferences, especially if the number
360 of potential population labels is large and the number of reference haplotype templates is limited.
361 Though Recomb-Mix uses a site-based approach, it achieved high accuracy. Allowing individual vari-
362 ations within the same population helps inflate the panel so more reference haplotype templates (e.g.,
363 relatives) are available for local site inference.

364 For a seven-way 200-generation 500-individual inter-continental WGS dataset, FLARE, G-Nomix,
365 and Loter had better r^2 values than Recomb-Mix. They were 0.1005, 0.0292, and 0.0017 higher
366 than that of Recomb-Mix, respectively (see Figure 3B). FLARE, G-Nomix, and Loter are claimed to
367 be the LAI methods good for identifying distant admixture events, as demonstrated high-resolution
368 accuracy when the admixture event occurs over 100 generations (Browning et al. 2023; Hilmarsson
369 et al. 2021; Dias-Alves et al. 2018). FLARE incorporates the number of generations as a parameter
370 in their model and its value is updated using an iterative expectation maximization (EM) approach
371 to calculate the probabilities of a change of ancestry state for each marker and haplotype. The longer
372 the admixture event occurs, the higher the probability that ancestral segments or tracts having a
373 large length difference appear. This information helps FLARE to update their generation parameter
374 better. On top of G-Nomix's base module's classifier, a smoother module is added to refine the
375 inference result. The smoother is a data-driven approach, which learns to capture the distribution
376 of recombination breakpoints. Usually, the distant admixture event has richer information on the
377 distribution of recombinations, which helps G-Nomix's smoother module improve the accuracy. Loter
378 adopts the bagging technique to generate the averaged result, which avoids putting a strong prior on
379 a particular length of ancestry segment. This helps improve the inference accuracy since the ancestry
380 segments appearing in distant admixture events are not the same length.

381 **Intra-continental admixture**

382 Compared to the inter-continental admixture, the LAI on the intra-continental is relatively less stud-
383 ied. The same benchmarks were set up and evaluated as the inter-continental ones. Similar to the

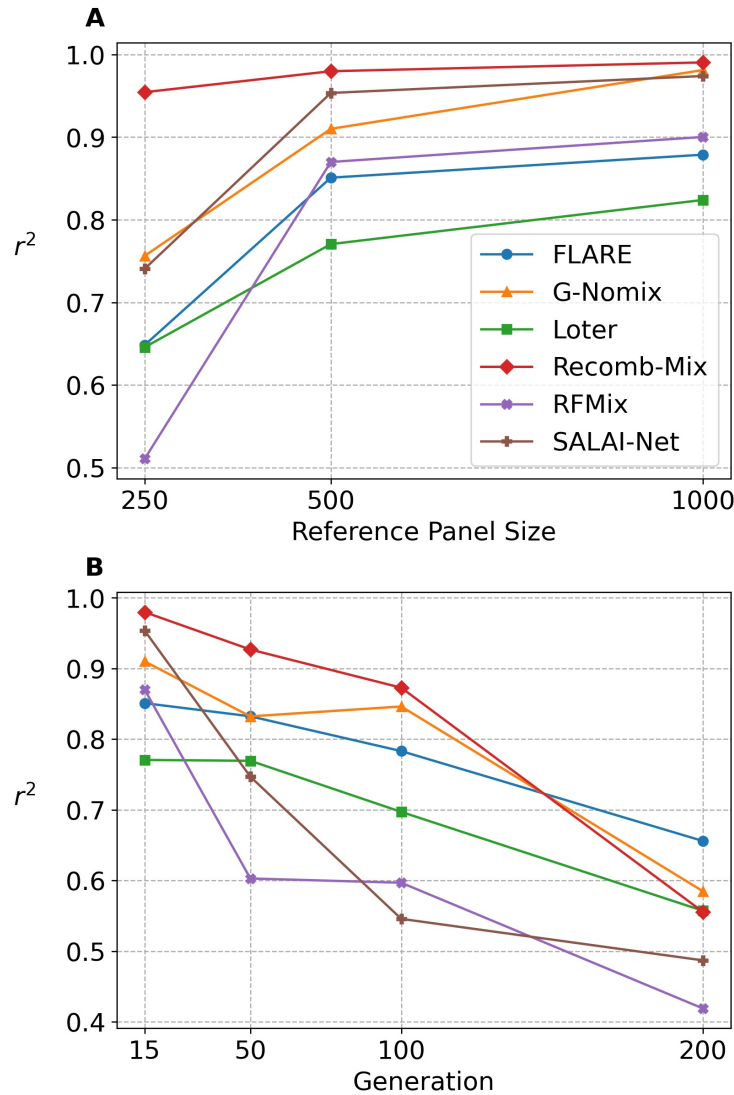


Figure 3: The squared Pearson's correlation coefficient r^2 of seven-way inter-continental simulated datasets on FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net. Markers were filtered with minor allele frequency ≤ 0.005 and minor allele count ≤ 50 . (A) The seven-way 15-generation datasets with the reference panel sizes 250, 500, and 1,000 (values are in Supplemental Table S8). The reference panel size 100 case was not included because the number of markers was too small and may have influenced the outcome after the filtering. (B) The seven-way 500-reference datasets with the generations 15, 50, 100, and 200 (values are in Supplemental Table S9).

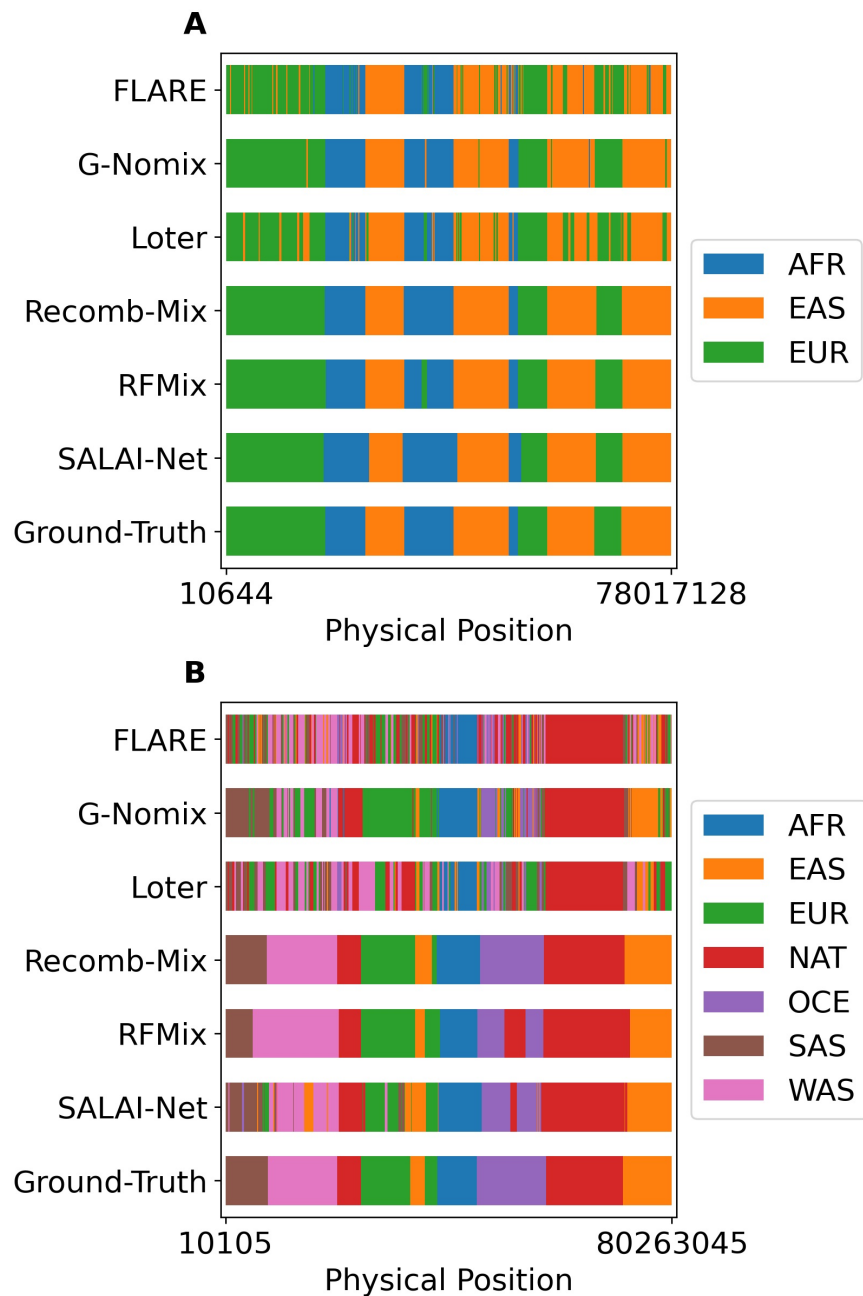


Figure 4: Sample haplotypes inferred by FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net with the ground truth of ancestry labels. (A) An inferred sample haplotype from a three-way 15-generation 500-reference inter-continental simulated dataset. (B) An inferred sample haplotype from a seven-way 15-generation 500-reference inter-continental simulated dataset.

384 results of the inter-continental datasets, Recomb-Mix performed well on the intra-continental datasets.
385 Figure 5A shows the r^2 values of the local ancestry inference of six LAI methods, FLARE, G-Nomix,
386 Loter, Recomb-Mix, RFMix, and SALAI-Net on a three-way 15-generation intra-continental simulated
387 dataset. (Supplemental Figure S6 shows the average accuracy rates (values are in Supplemental Tables
388 S14)). Overall, the r^2 values of each method were worse than the ones in the inter-continental datasets.
389 This is expected as the admixture occurring at the intra-continental level generates individuals who
390 resemble each other. Thus, performing LAI on such datasets is more challenging than at the inter-
391 continental level. Recomb-Mix had the best r^2 value in reference panel sizes 250, 500, and 1,000 with
392 15 generations. For a 250-individual reference panel, the r^2 value of Recomb-Mix was 0.9299, and the
393 second-best method, G-Nomix, only achieved 0.8560. For a 1,000-individual reference panel, the r^2
394 values of Recomb-Mix and G-Nomix were close (0.9800 and 0.9820).

395 The impact of the number of generations on LAI at the intra-continental level was also investigated.
396 Four three-way 500-reference intra-continental simulated datasets with generations 15, 50, 100, and 200
397 were tested, and the results show both the r^2 values and accuracy rates are inversely proportional to the
398 number of generations (see Figure 5B and Supplemental Figure S7 (values in Supplemental Table S15)).
399 For a 15-generation dataset, Recomb-Mix had the highest r^2 value, 0.9625. The second and third high
400 r^2 values were 0.9235 (G-Nomix) and 0.9081 (SALAI-Net). For a 200-generation dataset, G-Nomix
401 and FLARE had better r^2 values than Recomb-Mix. Also, Loter's performance increased with the
402 increasing number of generations. This result is consistent and observed in other simulated datasets,
403 as G-Nomix, FLARE, and Loter do well in ancestry inference on the cases of distant admixture events.

404 **Robustness against admixture with uneven proportions of founders and ref-** 405 **erences**

406 To verify the robustness of Recomb-Mix handling cases on uneven founder populations and reference
407 panels, LAI was experimented with using uneven founders for the imbalanced admixture simulation
408 and uneven reference panels for the inference. Being uneven means the group consists of one-third of
409 individuals from the first population, one-sixth from the second population, and half of individuals from
410 the third population. Being even means the numbers of individuals from the populations in the group

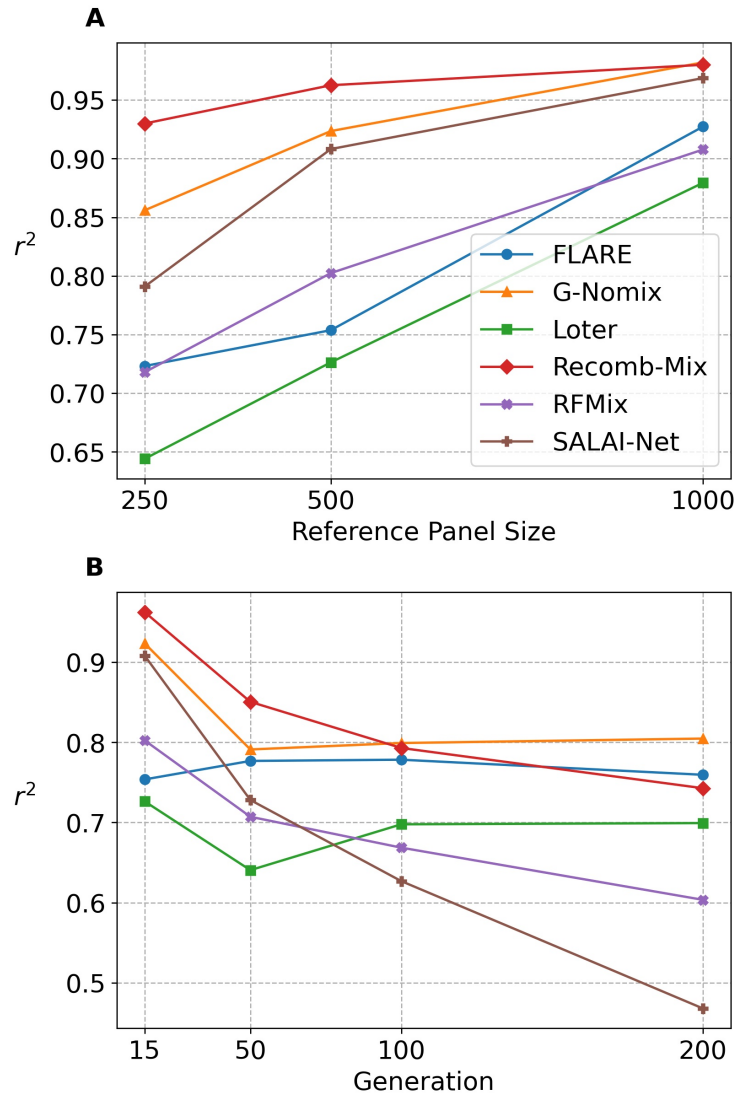


Figure 5: The squared Pearson's correlation coefficient r^2 of three-way intra-continental simulated datasets on FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net. Markers were filtered with minor allele frequency ≤ 0.005 and minor allele count ≤ 50 . (A) The three-way 15-generation datasets with the reference panel sizes 250, 500, and 1,000 (values are in Supplemental Table S12). The reference panel size 100 case was not included because the number of markers was too small and may have influenced the outcome after the filtering. (B) The three-way 500-reference datasets with the generations 15, 50, 100, and 200 (values are in Supplemental Table S13).

Method	Even Founders and References	Uneven Founders	Uneven References
Loter	0.7264	0.7263	0.7097
FLARE	0.7538	0.7187	0.7773
RFMix	0.8024	0.8112	0.7887
SALAI-Net	0.9081	0.8747	0.8692
G-Nomix	0.9235	0.9022	0.8811
Recomb-Mix	0.9625	0.9426	0.8944

Table 2: The squared Pearson’s correlation coefficient r^2 of FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net performing LAI on three-way 15-generation 500-reference intra-continental simulated datasets with even or uneven number of individuals per population in founder or reference panel. Markers were filtered with minor allele frequency ≤ 0.005 and minor allele count ≤ 50 .

411 are divided equally. Three sets of experiments were performed. One three-way admixture dataset
412 was simulated using even founders and inferred using even references, another was simulated using
413 uneven founders and inferred using even references, and the other was simulated using even founders
414 but inferred using uneven references.

415 The r^2 values and the accuracy rates in Table 2 and Supplemental Table S16 indicate that admixed
416 individuals with uneven founders and uneven reference panel slightly impact the performance across
417 all LAI methods. Among all LAI methods, Recomb-Mix had the highest r^2 values and accuracy rates
418 in both cases (0.9426 or 89.20% and 0.8944 or 83.61%, respectively). The process of Recomb-Mix
419 generating a collapsed graph helps convert the unbalanced reference populations into balanced ones.
420 Thus, Recomb-Mix keeps a high accuracy of inference results on the unbalanced reference populations.

421 Additionally, Recomb-Mix was tested on a modern Latino population admixture model that in-
422 volves uneven founders, which is a popular realistic model used as a study case for the local ancestry
423 inference (Maples et al. 2013; Wang et al. 2021). We used SLiM v4.0 to simulate the modern Latino
424 population dataset on Chromosome 1, using the same settings from the RFMix paper (Maples et al.
425 2013). Ten Latino genomes with 45% Native American (NAT), 50% European (CEU), and 5% African
426 ancestry (YRI) were simulated, originating from 400 individuals and 12 generations after the admixture
427 event. 30 individuals from each population were used to form the reference panel. We used Beagle 5.4

428 to phase the source data. The average LAI accuracy rate using Recomb-Mix and RFMix was 99.36%
429 and 93.79%, respectively. This shows that Recomb-Mix excels in the ancestry inference on the modern
430 Latino population admixture model derived from the uneven founders.

431 **Robustness against ancestry misspecification panel**

432 When performing real data analysis, the concern of data imperfection may be raised. Some populations
433 may be less studied and underrepresented in available reference panels. Furthermore, the existing
434 reference populations may contain a small fraction of admixture which may not make them the ideal
435 proxies for the labeled populations. Thus, it is necessary to investigate the impact of the ancestry
436 population misspecification on LAI.

437 An experiment was conducted by replacing the African reference population in a three-way inter-
438 continental admixed dataset with an imperfect reference panel. The imperfect version of the African
439 reference panel contains individuals who were Africans mixed with Europeans five generations from the
440 start of the simulation. This approach has similar effects as the one MOSAIC had, where their imperfect
441 reference panel contained admixed Sub-Saharan Africans and Europeans (Salter-Townshend and Myers
442 2019). We did not follow their process because the sampled individuals they used for the simulation were
443 from the extended HGDP dataset, whose data density is only at the single nucleotide polymorphism
444 (SNP) array level (Hellenthal et al. 2014). The ancestry misspecification experiments were repeated
445 for 15, 50, 100, and 200 generations since the admixture event, and FLARE, G-Nomix, Loter, Recomb-
446 Mix, RFMix, and SALAI-Net were tested. We did not include MOSAIC as it was designed for the case
447 when the source population lacked the availability of WGS data (Salter-Townshend and Myers 2019).

448 All LAI methods were impacted by the misspecification reference panel but still performed well,
449 as shown in Figure 6. Under the r^2 criteria with markers having minor allele values filtered, Recomb-
450 Mix performed the best in the cases of 50 and 100 generations since the admixture event. RFMix
451 performed the best for the most recent admixture case, and FLARE performed the best for the most
452 distant admixture case. Without filtering out any markers, Recomb-Mix had the highest accuracy rate
453 for most cases except the 15-generation case where RFMix performed the best.

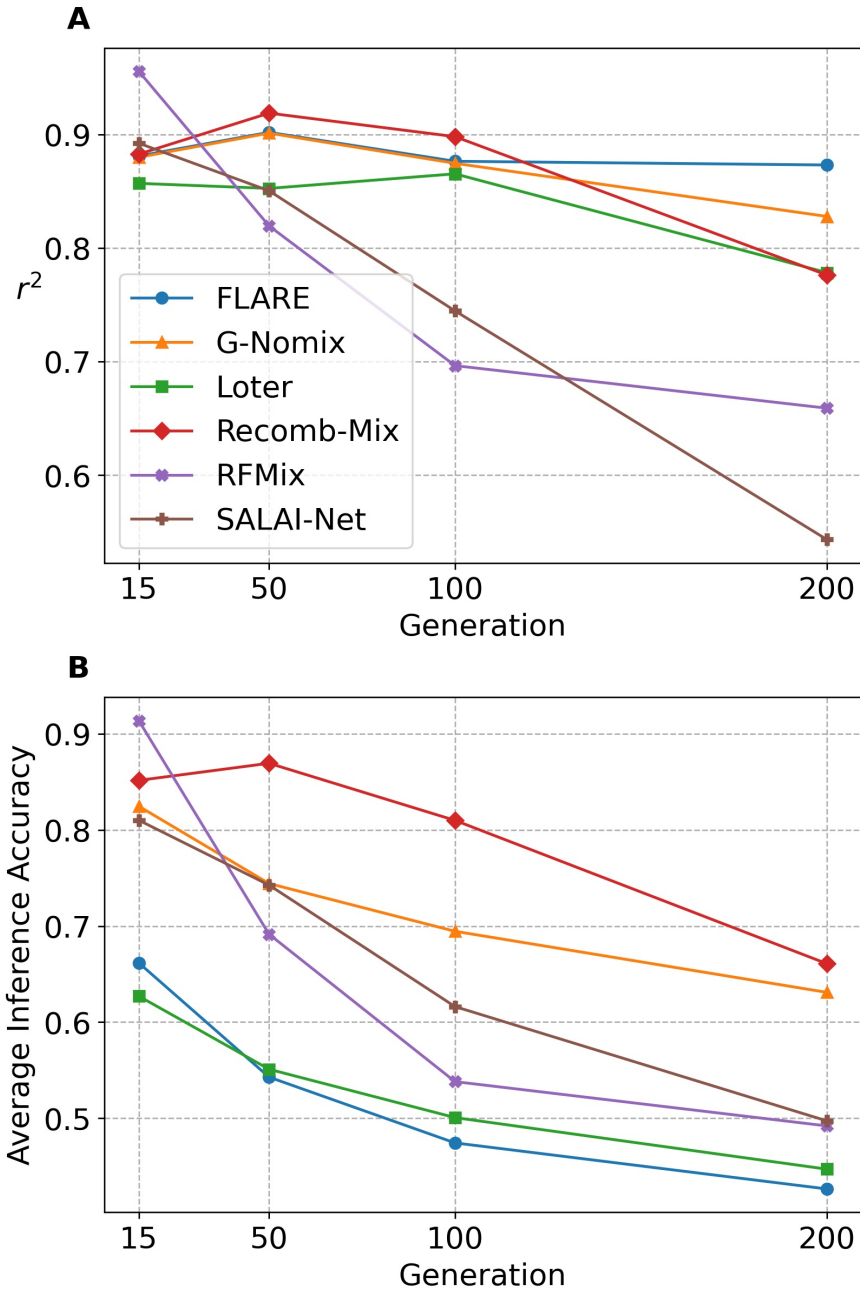


Figure 6: The performance of local ancestry inference with generations 15, 50, 100, and 200 of the three-way 500-misspecified-reference inter-continental simulated datasets on FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net. (A) The squared Pearson's correlation coefficient r^2 (values are in Supplemental Table S17). Markers were filtered with minor allele frequency ≤ 0.005 and minor allele count ≤ 50 . (B) The average accuracy rates (values are in Supplemental Table S18).

454 **Robustness against phasing error**

455 To investigate the impact of phasing error on local ancestry inference, two cases of phasing errors on
456 either the target panel or the reference panel were explored. The three-way 15-generation 100-reference
457 inter-continental simulated dataset was used as the baseline, and Beagle 5.4 was applied to phase the
458 panels. The phasing error rate was 0.58% for the target panel and 1.33% for the reference panel, verified
459 by VCFtools (Danecek et al. 2011). FLARE, G-Nomix, Loter, Recomb-Mix, RFMix, and SALAI-Net
460 were tested on the panels and diploid accuracy rates were used for the performance measurement as
461 RFMix did (Maples et al. 2013). The results in Supplemental Table S19 show that the diploid accuracy
462 rates did not fluctuate much when there were phasing errors on the panels, indicating that the low rate
463 of phasing errors may not have a substantial impact on the local ancestry inference.

464 **Recomb-Mix is efficient in memory, space, and run time**

465 We examined the run time and maximum amount of memory LAI methods used for their performance
466 on admixed individual haplotypes. Supplemental Figures S8, S9, and Table S20 show the average CPU
467 run time and maximum amount of physical memory that all six LAI methods consumed across different
468 experimental runs. In general, for the same method, inference on three-way admixed individuals was
469 faster than those on seven-way. This is expected as a seven-way admixture has many more local
470 ancestral segments across the chromosome than ones in a three-way, which costs more time for the
471 inference. All methods showed reasonable run time for an LAI query of an admixed individual haplotype
472 except Loter, which was about 10 or 100 times slower than other methods. SALAI-Net was the fastest
473 method and Recomb-Mix was the runner-up but only took 0.31 and 2.04 more seconds than SALAI-Net
474 in three-way and seven-way datasets. From the memory-consuming perspective, all methods' memory
475 usage was acceptable, and Recomb-Mix required the smallest amount of memory, 2.44 and 4.13 GB in
476 three-way and seven-way datasets, respectively.

477 Recomb-Mix has a feature that converts the compact population graph into a Variant Call Format
478 (VCF) file (Danecek et al. 2011). Later, they can be reused by Recomb-Mix to save processing time.
479 A compact VCF file is much smaller than the original one since it only contains individual haplotype
480 templates with population-level information. For example, the disk space needed to store a 3-way

481 inter-continental 500-reference panel was decreased from 665 to 13.3 MB (and 1.4 MB for a compressed
482 VCF file). Similarly, for a 7-way panel, the disk space was decreased from 787 to 21.8 MB (and 1.9
483 MB for a compressed VCF file).

484 **The 1000 Genomes Project and the Human Genome Diversity Project an-** 485 **cestry analysis**

486 To show the scalability and robustness of Recomb-Mix, we estimated the ancestry proportions from
487 the inferred local ancestries for the populations in the 1000 Genomes Project (TGP) data (Byrskä-
488 Bishop et al. 2022) using the four founder populations (Africans, Admixed Americans, East Asians, and
489 Europeans) as the reference panel from the Human Genome Diversity Project (HGDP) data (Bergström
490 et al. 2020). Similarly, We estimated the ancestry proportions for the populations in the HGDP data
491 using the four founder populations (Africans, East Asians, Europeans, and Native Americans) from the
492 TGP data. We merged two Chromosome 18 datasets (TGP with 3,457,645 markers and HGDP with
493 2,127,412 markers), yielding 1,165,399 intersected markers. Then the merged dataset was phased using
494 Beagle 5.4 and the individuals were assigned the population labels provided by their original datasets.

495 Like FLARE (Browning et al. 2023), we calculated the global ancestry composition by averaging the
496 estimated local ancestry proportions across the genome. Figure 7 is Recomb-Mix’s ancestry inference
497 result on the TGP dataset that is generally consistent with exceptions. Similar results using other
498 LAI methods are available in Supplemental Figure S10. For African individuals who reside in the
499 African continent, at least 97% segment was labeled as African on average. For ACB and ASW
500 individuals (located in the Caribbean and America), a small portion of the segment was labeled as
501 non-African due to their admixed backgrounds. Most segments of American individuals were labeled
502 as mixed percentages of Europeans and Native Americans. South Asian individuals were labeled as
503 mixed percentages of East Asian and European. East Asian and European individuals had at least
504 99% and 94% segments labeled as East Asian and European, respectively.

505 The ancestry inference result on the HGDP dataset is generally anticipated. Figure 8 shows that
506 Recomb-Mix predicted African, East Asian, European, and native American individuals have 99%,
507 94%, 88%, and 98% segments matched expected ancestries. For the Oceanian individuals, the segments

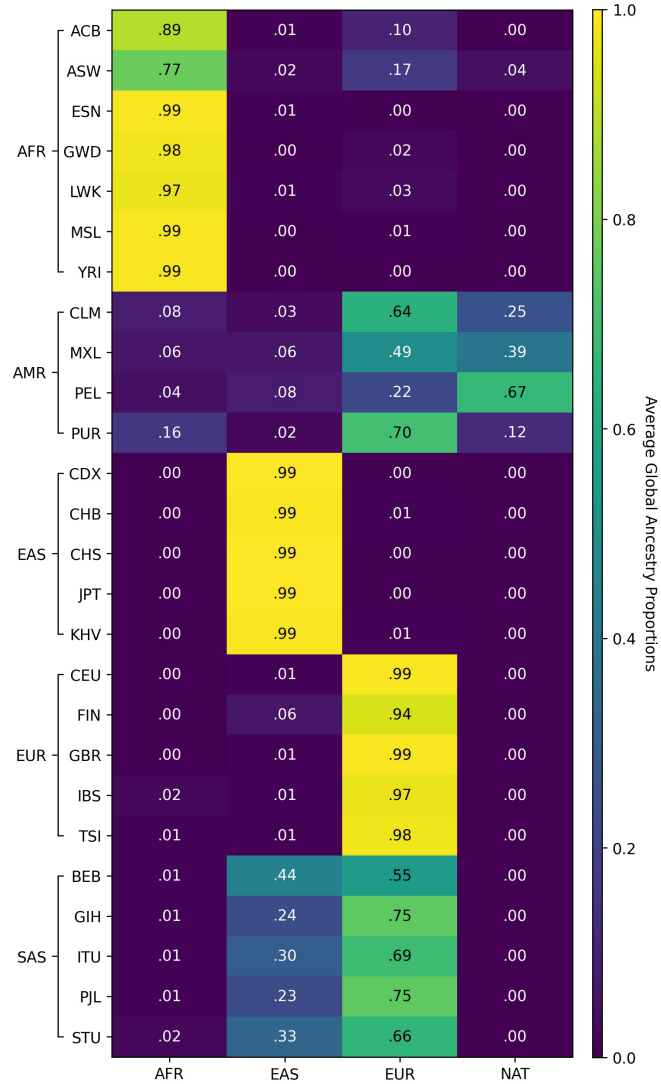


Figure 7: The average global ancestry proportions in the TGP Chromosome 18 data using four reference ancestries from the HGDP data. Descriptions of the populations are in Supplemental Table S1.

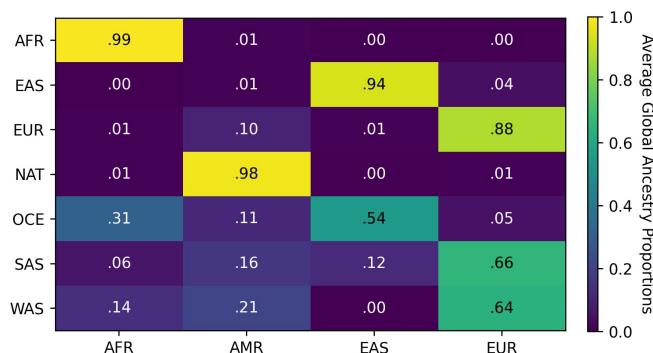


Figure 8: The average global ancestry proportions in the HGDP Chromosome 18 data using four reference ancestries from the TGP data. Descriptions of the populations are in Supplemental Table S1.

508 were decomposed into mixed ancestries, primarily East Asian and African. For South Asian and West
 509 Asian individuals, the segments were inferred as mixed and mainly consisted of European ancestry.
 510 Interpretation of ancestry inference results when the founder populations in the reference panel were
 511 more complex. For example, recent genetic evidence suggests that EUR, WAS, and SAS may share
 512 some Yamnaya DNA (Lazaridis et al. 2022; Narasimhan et al. 2019), which might be part of the causes
 513 of our results of possible shared (about 10%) ancient ancestry among AMR, EUR, OCE, SAS, and
 514 WAS. Similar behaviors were observed on other LAI methods, such as G-Nomix and SALAI-Net (see
 515 Supplemental Figure S11). Recomb-Mix was forced to give a single LAI call for these regions since the
 516 inference results were based on the given reference panels. If the ancient population were not in the
 517 reference panel, the segment would be labeled as the population closest to the ancient one.

518 Discrete ancestry informative markers

519 Local ancestry inference may benefit from ancestry informative markers (AIMs), which are genetic
 520 markers with significantly different allele frequencies in various populations (Parra et al. 1998). AIMs
 521 provide information regarding ancestry and can be determined in a panel by measuring marker in-
 522 formativeness for ancestry (Ding et al. 2011). Rather than relying on a selected set of markers, i.e.,
 523 AIMs, the proposed compact population graph keeps all markers, and the only differential information
 524 between populations is that some alleles might be missing in one or several populations. These markers

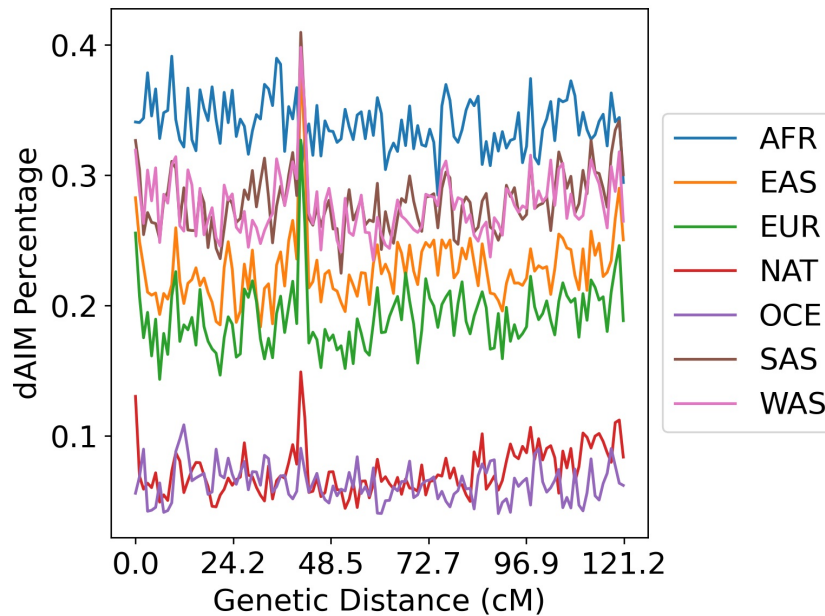


Figure 9: The discrete AIM (dAIM) density in the HGDP dataset per population on Chromosome 18. Each bin is 1 centiMorgan (cM), showing the markers' dAIM percentage.

525 are dubbed discrete AIMS (dAIMs), whose allele values in one population that at least one of the other
526 populations does not have. dAIMs are operationally defined and depend on some random chance of
527 whether an allele is present in the reference panel for the population or not. However, as it is shown
528 in Figure 9 (the dAIM densities (i.e., percentages of markers in the dataset being dAIMs) of Chromo-
529 some 18 in the HGDP data (Bergström et al. 2020)), dAIMs are densely available on a typical panel.
530 Therefore, even though collapsing nodes reduced the information in the original panel, the remaining
531 information in dAIMs might be sufficient for making good-quality ancestry calls. For example, there
532 is a dAIM density peak occurring around the 18q21 region in the HGDP dataset (see Figure 9). In
533 a previous admixture mapping study, a genome-wide significant admixture mapping peak contributed
534 from multiple ancestry signals was identified in the same region (Gignoux et al. 2019). This correlation
535 suggests that dAIM density has the potential to identify ancestry-specific selection.

536 Ablation study

537 We want to understand which component contributes the most to Recomb-Mix’s ancestry inference pro-
538 cess. Experiments were designed for Recomb-Mix to make inferences on a three-way inter-continental
539 15-generation 100-reference simulated dataset by not setting the within-population template change
540 penalty to zero or not using the recombination rates. We observed a slight decrease in performance
541 when the recombination rate was not used. However, the performance dropped significantly when
542 the within-population template change penalty was applied (see Supplemental Figure S12 and Ta-
543 ble S21). Supplemental Figure S12 shows LAI accuracies are significantly improved when setting the
544 within-population template change penalty to zero, especially for small reference panel cases. We
545 also calculated the average number of threading path changes across populations and the standard
546 deviations of the dataset with 228,503 markers for Recomb-Mix and the version that used the within-
547 population template change penalty. Recomb-Mix had 8.47 ± 2.54 , while the latter version had a much
548 larger number, 198.46 ± 46.94 . The LAI calling became less effective when using the template change
549 penalty within each population. This may be due to the numerous local optimal threading paths to be
550 explored within each population, which can lead to noise and deviation from finding the path with the
551 minimal global penalty score. By setting the within-population template change penalty to zero, the
552 number of explorations between the paths within a population is significantly reduced, and the focus
553 is shifted to only a few consolidated paths representing diverse haplotype templates.

554 The dAIMs are usually evenly distributed alongside the chromosome, as we showed in Figure 9.
555 To illustrate dAIM’s important role in LAI, we designed an experiment to engineer a new dataset
556 based on the simulated one by taking out all the dAIMs for certain regions. We tested Recomb-Mix
557 on the engineered dataset, and the result showed a strong correlation between the dAIM density and
558 the accuracy of the inference. Supplemental Figures S13 and S14 show the dAIM density and local
559 ancestry inference accuracy rate of the original dataset and the engineered dataset. In Supplemental
560 Figure S14, there are five instances where the low LAI accurate rates correspond with areas lacking
561 dAIMs in the engineered dataset. The Pearson correlation coefficient for this dataset’s dAIM density
562 and local ancestry inference accuracy rate is 0.79, demonstrating a strong correlation between dAIMs
563 and LAI accuracies.

564 Single nucleotide polymorphism array data analysis

565 We created an SNP array dataset to verify if Recomb-Mix works on a panel with a limited number of
566 dAIMs and rare variants. First, we followed Tang et al.'s pipeline (2022) to down-sample the three-way
567 inter-continental 15-generation 100-reference simulated sequencing dataset. The number of markers
568 was decreased from 228,503 to 15,584. Then, we further filtered out markers having minor allele
569 frequency (MAF) less than 5%, as typically genotyping array data contains common variants whose
570 MAF > 5% (Bomba et al. 2017; Verlouw et al. 2021). Eventually, the dataset had 8,744 markers that
571 resembled an imputed SNP array panel. The number of dAIMs in the panel also decreased compared
572 to the one in the sequencing panel (see Supplemental Figure S15 and S16); however, the dAIM density
573 in the SNP array data did not change as much as the one in the sequencing data (see Supplemental
574 Figure S17 and S18). We tested Recomb-Mix on the SNP array dataset, and it achieved 0.9949 r^2
575 value and 96.34% accuracy rate, compared to 0.9986 and 97.96% on the sequencing dataset without
576 filtering any sites. This result suggests Recomb-Mix has the same performance on the SNP array data
577 as on the sequencing data, possibly due to the dAIM density in the SNP array data being on the same
578 order of magnitude as the one in the sequencing data.

579 Discussion

580 We presented a new LAI method named Recomb-Mix, based on a simplified LS model formulating
581 LAI as a graph optimization problem. By not considering recombination penalties within populations,
582 Recomb-Mix shows promising LAI results under various circumstances. A compact population graph
583 also helps Recomb-Mix process LAI effectively and efficiently. Furthermore, it is convenient to store the
584 reference panel as a compact population graph on disk, which takes up little space for future ancestry
585 inference without a re-transformation process. Recomb-Mix is competitive with other state-of-the-art
586 LAI methods in accuracy and computational performance and is applicable to real genomic datasets.

587 We introduced the concept of dAIM, where dAIMs are determined by the allele values of each
588 population. We showed that dAIMs can have marker informativeness for ancestry. Of course, this
589 selection of markers is simplistic and mainly captures the differentially present or absent markers in
590 the reference panel across populations. In future studies, other ways of collapsing the graph might be

591 explored to retain more relevant ancestry information from the reference panel. It might be optimized
592 to allow the selection of markers present in multiple populations but with different allele frequencies.
593 This could further enhance the performance and enable our model to provide uncertainty estimates for
594 ancestry inference results.

595 As a site-based LAI method, Recomb-Mix is designed to exploit the site-level information to achieve
596 superior accuracy, especially at the intra-continental level. The number of dAIMs in intra-continental
597 admixed individuals is less than those in inter-continental admixed individuals, as intra-dAIMs are
598 a subset of inter-dAIMs. Window-based LAI methods may find it difficult to achieve high accuracy
599 during the intra-level ancestry inference. There is a higher probability for each window containing mul-
600 tiple intra-dAIMs than that for inter-dAIMs. Since the window is the smallest unit representing one
601 ancestral source, windows having intra-dAIMs representing different populations may easily misrepre-
602 sent the inference result. Decreasing the window size may mitigate the situation, with the potential
603 computational burden. However, its lower bound is one site per window, i.e., site-based.

604 Despite the high accuracy rates demonstrated, Recomb-Mix has limitations of not considering the
605 disparate genetic maps across populations, the allele frequencies, or genotyping and phasing errors.
606 Thus, other complementary methods may be useful for a well-specified model. FLARE takes optional
607 parameters such as minor allele frequency and number of generations since admixture. It may perform
608 well if these biological parameters are correctly estimated for the model. G-Nomix has a few pre-
609 trained models available which may be in handy if one was pre-trained specifically for the given model.
610 SALAI-Net may be a good choice as it employs a pre-trained model that is generalized and applicable
611 to any species and any set of ancestries. Though Recomb-Mix was not designed to handle erroneous
612 panels, the genotyping error and phasing error seem to have no large impact on the inference results
613 (see the Results section). If the error rate is high, a pre-processing step may be needed to correct the
614 noisy data panel before making the inference.

615 **Software availability**

616 The Recomb-Mix code is available at <https://github.com/ucfcb/Recomb-Mix>.

617 **Competing interest statement**

618 The authors declare no competing interests.

619 **Acknowledgments**

620 The authors thank Dr. Ahsan Sanaullah for the helpful suggestions on the manuscript. This work was
621 supported by the National Institutes of Health under grants R01HG010086 and R56HG011509.

622 **References**

623 Adrion JR, Cole CB, Dukler N, Galloway JG, Gladstein AL, Gower G, Kyriazis CC, Ragsdale AP,
624 Tsambos G, Baumdicker F, et al.. 2020. A community-maintained standard library of population
625 genetic models. *eLife* **9**: e54967.

626 Atkinson EG, Maihofer AX, Kanai M, Martin AR, Karczewski KJ, Santoro ML, Ulirsch JC, Kamatani
627 Y, Okada Y, Finucane HK, et al.. 2021. Tractor uses local ancestry to enable the inclusion of admixed
628 individuals in GWAS and to boost power. *Nature Genetics* **53**: 195–204.

629 Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, Clark AG, Donnelly
630 P, Eichler EE, Flicek P, et al.. 2015. A global reference for human genetic variation. *Nature* **526**:
631 68–74.

632 Baran Y, Pasaniuc B, Sankararaman S, Torgerson DG, Gignoux C, Eng C, Rodriguez-Cintron W,
633 Chapela R, Ford JG, Avila PC, et al.. 2012. Fast and accurate inference of local ancestry in Latino
634 populations. *Bioinformatics* **28**: 1359–1367.

635 Bergström A, McCarthy SA, Hui R, Almarri MA, Ayub Q, Danecek P, Chen Y, Felkel S, Hallast P,
636 Kamm J, et al.. 2020. Insights into human genetic variation and population history from 929 diverse
637 genomes. *Science* **367**: eaay5012.

638 Bomba L, Walter K, and Soranzo N. 2017. The impact of rare and low-frequency genetic variants in
639 common disease. *Genome Biology* **18**: 77.

- 640 Browning BL, Tian X, Zhou Y, and Browning SR. 2021. Fast two-stage phasing of large-scale sequence
641 data. *The American Journal of Human Genetics* **108**: 1880–1890.
- 642 Browning BL, Zhou Y, and Browning SR. 2018. A One-Penny imputed genome from Next-Generation
643 reference panels. *The American Journal of Human Genetics* **103**: 338–348.
- 644 Browning SR, Waples RK, and Browning BL. 2023. Fast, accurate local ancestry inference with FLARE.
645 *The American Journal of Human Genetics* **110**: 326–335.
- 646 Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau
647 O, O’Connell J, et al.. 2018. The UK Biobank resource with deep phenotyping and genomic data.
648 *Nature* **562**: 203–209.
- 649 Byrska-Bishop M, Evani US, Zhao X, Basile AO, Abel HJ, Regier AA, Corvelo A, Clarke WE, Musunuri
650 R, Nagulapalli K, et al.. 2022. High-coverage whole-genome sequencing of the expanded 1000
651 Genomes Project cohort including 602 trios. *Cell* **185**: 3426–3440.e19.
- 652 Church DM, Schneider VA, Graves T, Auger K, Cunningham F, Bouk N, Chen HC, Agarwala R,
653 McLaren WM, Ritchie GR, et al.. 2011. Modernizing reference genome assemblies. *PLoS Biology* **9**:
654 1–5.
- 655 Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, Tassé AM, and Flicek P. 2016.
656 The international Genome sample resource (IGSR): A worldwide collection of genome variation
657 incorporating the 1000 Genomes Project data. *Nucleic Acids Research* **45**: D854–D859.
- 658 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth
659 GT, Sherry ST, et al.. 2011. The variant call format and VCFtools. *Bioinformatics* **27**: 2156–2158.
- 660 Delaneau O, Marchini J, and Zagury JF. 2012. A linear complexity phasing method for thousands of
661 genomes. *Nature Methods* **9**: 179–181.
- 662 Delaneau O, Zagury JF, Robinson MR, Marchini JL, and Dermitzakis ET. 2019. Accurate, scalable
663 and integrative haplotype estimation. *Nature Communications* **10**: 5436.

- 664 Dias-Alves T, Mairal J, and Blum MGB. 2018. Loter: A software package to infer local ancestry for a
665 wide range of species. *Molecular Biology and Evolution* **35**: 2318–2326.
- 666 Ding L, Wiener H, Abebe T, Altaye M, Go RC, Kercksmar C, Grabowski G, Martin LJ, Khurana Hershey
667 GK, Chakorborty R, et al.. 2011. Comparison of measures of marker informativeness for ancestry
668 and admixture mapping. *BMC Genomics* **12**: 622.
- 669 Ding Y, Hou K, Xu Z, Pimplaskar A, Petter E, Boulier K, Privé F, Vilhjálmsson BJ, Olde Loohuis
670 LM, and Pasaniuc B. 2023. Polygenic scoring accuracy varies across the genetic ancestry continuum.
671 *Nature* **618**: 774–781.
- 672 Duncan L, Shen H, Gelaye B, Meijssen J, Ressler K, Feldman M, Peterson R, and Domingue B. 2019.
673 Analysis of polygenic risk score usage and performance in diverse human populations. *Nature Com-*
674 *munications* **10**: 3328.
- 675 Durand EY, Do CB, Wilton PR, Mountain JL, Auton A, Poznik GD, and Macpherson JM. 2021. A
676 scalable pipeline for local ancestry inference using tens of thousands of reference haplotypes. *bioRxiv*
677 doi: 10.1101/2021.01.19.427308.
- 678 Durbin R. 2014. Efficient haplotype matching and storage using the positional Burrows-Wheeler trans-
679 form (PBWT). *Bioinformatics* **30**: 1266–1272.
- 680 Geza E, Mugo J, Mulder NJ, Wonkam A, Chimusa ER, and Mazandu GK. 2018. A comprehensive
681 survey of models for dissecting local ancestry deconvolution in human genome. *Briefings in Bioin-*
682 *formatics* **20**: 1709–1724.
- 683 Gignoux CR, Torgerson DG, Pino-Yanes M, Uricchio LH, Galanter J, Roth LA, Eng C, Hu D, Nguyen
684 EA, Huntsman S, et al.. 2019. An admixture mapping meta-analysis implicates genetic variation
685 at 18q21 with asthma susceptibility in Latinos. *Journal of Allergy and Clinical Immunology* **143**:
686 957–969.
- 687 Haller BC, Galloway J, Kelleher J, Messer PW, and Ralph PL. 2019. Tree-sequence recording in SLiM
688 opens new horizons for forward-time simulation of whole genomes. *Molecular Ecology Resources* **19**:
689 552–566.

- 690 Haller BC and Messer PW. 2019. SLiM 3: Forward genetic simulations beyond the Wright–Fisher
691 model. *Molecular Biology and Evolution* **36**: 632–637.
- 692 Hamid I, Korunes KL, Schrider DR, and Goldberg A. 2023. Localizing post-admixture adaptive vari-
693 ants with object detection on ancestry-painted chromosomes. *Molecular Biology and Evolution* **40**:
694 msad074.
- 695 Hellenthal G, Busby GBJ, Band G, Wilson JF, Capelli C, Falush D, and Myers S. 2014. A genetic
696 atlas of human admixture history. *Science* **343**: 747–751.
- 697 Hilmarsson H, Kumar AS, Rastogi R, Bustamante CD, Mas Montserrat D, and Ioannidis AG.
698 2021. High resolution ancestry deconvolution for next generation genomic data. *bioRxiv* doi:
699 10.1101/2021.09.19.460980.
- 700 Hou K, Ding Y, Xu Z, Wu Y, Bhattacharya A, Mester R, Belbin GM, Buyske S, Conti DV, Darst BF,
701 et al.. 2023. Causal effects on complex traits are similar for common variants across segments of
702 different continental ancestries within admixed individuals. *Nature Genetics* **55**: 549–558.
- 703 Kurki MI, Karjalainen J, Palta P, Sipilä TP, Kristiansson K, Donner KM, Reeve MP, Laivuori H,
704 Aavikko M, Kaunisto MA, et al.. 2023. Finngen provides genetic insights from a well-phenotyped
705 isolated population. *Nature* **613**: 508–518.
- 706 Lawson DJ, Hellenthal G, Myers S, and Falush D. 2012. Inference of population structure using dense
707 haplotype data. *PLoS Genetics* **8**: 1–16.
- 708 Lazaridis I, Alpaslan-Roodenberg S, Acar A, Açıkkol A, Agelarakis A, Aghikyan L, Akyüz U, Andreeva
709 D, Andrijašević G, Antonović D, et al.. 2022. The genetic history of the Southern Arc: A bridge
710 between West Asia and Europe. *Science* **377**: eabm4247.
- 711 Li N and Stephens M. 2003. Modeling linkage disequilibrium and identifying recombination hotspots
712 using single-nucleotide polymorphism data. *Genetics* **165**: 2213–2233.
- 713 Maples BK, Gravel S, Kenny EE, and Bustamante CD. 2013. RFMix: A discriminative modeling

- 714 approach for rapid and robust local-ancestry inference. *The American Journal of Human Genetics*
715 **93**: 278–288.
- 716 Martin AR, Gignoux CR, Walters RK, Wojcik GL, Neale BM, Gravel S, Daly MJ, Bustamante CD,
717 and Kenny EE. 2017. Human demographic history impacts genetic risk prediction across diverse
718 populations. *The American Journal of Human Genetics* **100**: 635–649.
- 719 Montserrat DM, Bustamante C, and Ioannidis A. 2020. Lai-Net: Local-ancestry inference with neural
720 networks. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal*
721 *Processing (ICASSP)*, pp. 1314–1318.
- 722 Narasimhan VM, Patterson N, Moorjani P, Rohland N, Bernardos R, Mallick S, Lazaridis I, Nakatsuka
723 N, Olalde I, Lipson M, et al.. 2019. The formation of human populations in South and Central Asia.
724 *Science* **365**: eaat7487.
- 725 Oriol Sabat B, Mas Montserrat D, Giro-i-Nieto X, and Ioannidis AG. 2022. SALAI-Net: Species-
726 agnostic local ancestry inference network. *Bioinformatics* **38**: ii27–ii33.
- 727 Parra EJ, Marcini A, Akey J, Martinson J, Batzer MA, Cooper R, Forrester T, Allison DB, Deka R,
728 Ferrell RE, et al.. 1998. Estimating African American admixture proportions by use of population-
729 specific alleles. *The American Journal of Human Genetics* **63**: 1839–1851.
- 730 Pasaniuc B, Zaitlen N, Lettre G, Chen GK, Tandon A, Kao WHL, Ruczinski I, Fornage M, Siscovick
731 DS, Zhu X, et al.. 2011. Enhanced statistical tests for GWAS in admixed populations: Assessment
732 using African Americans from CARE and a breast cancer consortium. *PLoS Genetics* **7**: 1–15.
- 733 Patin E, Siddle KJ, Laval G, Quach H, Harmant C, Becker N, Froment A, Régnault B, Lemée L, Gravel
734 S, et al.. 2014. The impact of agricultural emergence on the genetic history of african rainforest
735 hunter-gatherers and agriculturalists. *Nature Communications* **5**: 3163.
- 736 Price AL, Tandon A, Patterson N, Barnes KC, Rafaels N, Ruczinski I, Beaty TH, Mathias R, Reich
737 D, and Myers S. 2009. Sensitive detection of chromosomal segments of distinct ancestry in admixed
738 populations. *PLoS Genetics* **5**: 1–18.

- 739 Reich D, Patterson N, Jager PLD, McDonald GJ, Waliszewska A, Tandon A, Lincoln RR, DeLoa
740 C, Fruhan SA, Cabre P, et al.. 2005. A whole-genome admixture scan finds a candidate locus for
741 multiple sclerosis susceptibility. *Nature Genetics* **37**: 1113–1118.
- 742 Salter-Townshend M and Myers S. 2019. Fine-scale inference of ancestry segments without prior
743 knowledge of admixing groups. *Genetics* **212**: 869–889.
- 744 Schneider VA, Graves-Lindsay T, Howe K, Bouk N, Chen HC, Kitts PA, Murphy TD, Pruitt KD,
745 Thibaud-Nissen F, Albracht D, et al.. 2017. Evaluation of GRCh38 and de novo haploid genome
746 assemblies demonstrates the enduring quality of the reference assembly. *Genome Research* **27**: 849–
747 864.
- 748 Suarez-Pajes E, Díaz-de Usera A, Marcelino-Rodríguez I, Guillen-Guio B, and Flores C. 2021. Genetic
749 Ancestry Inference and Its Application for the Genetic Mapping of Human Diseases. *International*
750 *Journal of Molecular Sciences* **22**: 6962.
- 751 Tang K, Naseri A, Wei Y, Zhang S, and Zhi D. 2022. Open-source benchmarking of IBD segment
752 detection methods for biobank-scale cohorts. *GigaScience* **11**. Giac111.
- 753 Verlouw JAM, Clemens E, de Vries JH, Zolk O, Verkerk AJMH, am Zehnhoff-Dinnesen A, Medina-
754 Gomez C, Lanvers-Kaminsky C, Rivadeneira F, Langer T, et al.. 2021. A comparison of genotyping
755 arrays. *European Journal of Human Genetics* **29**: 1611–1624.
- 756 Wang Y, Song S, Schraiber JG, Sedghifar A, Byrnes JK, Turissini DA, Hong EL, Ball CA, and Noto
757 K. 2021. Ancestry inference using reference labeled clusters of haplotypes. *BMC Bioinformatics* **22**:
758 459.
- 759 Wu J, Liu Y, and Zhao Y. 2021. Systematic review on local ancestor inference from a mathematical
760 and algorithmic perspective. *Frontiers in Genetics* **12**: 698.