

1 **Characterizing features affecting local ancestry inference performance in admixed populations**

2 Jessica Honorato-Mauer<sup>1</sup>, Nirav N. Shah<sup>1</sup>, Adam X. Maihofer<sup>2</sup>, Clement C. Zai<sup>3</sup>, Sintia Belangero<sup>4</sup>,  
3 Caroline M. Nievergelt<sup>2</sup>, Psychiatric Genomics Consortium for PTSD Ancestry Working Group,  
4 Marcos Santoro<sup>5,\*</sup>, & Elizabeth Atkinson<sup>1,6,\*</sup>

5 **Affiliations**

6 1 Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, 77030,  
7 USA

8 2 Department of Psychiatry, School of Medicine, University of California at San Diego, La Jolla, CA  
9 92093, USA

10 3 Department of Psychiatry, Institute of Medical Science, Laboratory Medicine and Pathobiology,  
11 University of Toronto

12 4 Department of Morphology and Genetics, Universidade Federal de São Paulo, São Paulo, 04023-  
13 062, Brazil

14 5 Department of Biochemistry, Molecular Biology Division, Universidade Federal de São Paulo, São  
15 Paulo, 04023-062, Brazil

16 6 The Jan and Dan Duncan Neurological Research Institute, Texas Children's Hospital, Houston, TX  
17 77030, USA

18

19 **Correspondence emails:** [elizabeth.atkinson@bcm.edu](mailto:elizabeth.atkinson@bcm.edu) or [santoro@unifesp.br](mailto:santoro@unifesp.br)

20 \* Authors contributed equally to this work.

21

22 **Abstract (250 words)**

23 In recent years, significant efforts have been made to improve methods for genomic studies of  
24 admixed populations using Local Ancestry Inference (LAI). Accurate LAI is crucial to ensure  
25 downstream analyses reflect the genetic ancestry of research participants accurately. Here, we test  
26 analytic strategies for LAI to provide guidelines for optimal accuracy, focusing on admixed  
27 populations reflective of Latin America's primary continental ancestries – African (AFR),

28 Amerindigenous (AMR), and European (EUR). Simulating LD-informed admixed haplotypes under a  
29 variety of 2 and 3-way admixture models, we implemented a standard LAI pipeline, testing three  
30 reference panel compositions to quantify their overall and ancestry-specific accuracy. We examined  
31 LAI miscall frequencies and true positive rates (TPR) across simulation models and continental  
32 ancestries. AMR tracts have notably reduced LAI accuracy as compared to EUR and AFR tracts in all  
33 comparisons, with TPR means for AMR ranging from 88-94%, EUR from 96-99% and AFR 98-99%.  
34 When LAI miscalls occurred, they most frequently erroneously called European ancestry in true  
35 Amerindigenous sites. Using a reference panel well-matched to the target population, even with a  
36 lower sample size, LAI produced true-positive estimates that were not statistically different from a  
37 high sample size but mismatched reference, while being more computationally efficient. While  
38 directly responsive to admixed Latin American cohort compositions, these trends are broadly useful  
39 for informing best practices for LAI across other admixed populations. Our findings reinforce the  
40 need for inclusion of more underrepresented populations in sequencing efforts to improve reference  
41 panels.

42

### 43 **Introduction**

44 Admixed populations present a challenge in genome-wide analyses, as their genomes contain  
45 components from different continental ancestries, which vary from person to person along the  
46 genome, even if two people have the same overall ancestry proportions. This makes it statistically  
47 challenging to control for population structure, which can bias tests if left uncorrected<sup>1,2</sup>. Despite  
48 recent advances in complex trait genetics, limitations remain in our understanding of the architecture  
49 of genetic disorders in diverse populations due to their exclusion from many genomic studies<sup>3-5</sup>. For  
50 example, Latin American (LatAm) populations currently represent only 1.3% of all genome-wide  
51 association studies (GWAS) samples, despite accounting for 8.4% of the world population and  
52 contributing disproportionately to GWAS findings<sup>6</sup>. As large-scale efforts begin to focus more heavily  
53 on admixed groups, there is an unmet need for the design of well-suited pipelines to appropriately  
54 study these underrepresented populations<sup>7</sup>.

55           Local Ancestry Inference (LAI) is a machine learning approach for assigning each genomic  
56 region to a specific ancestry group by comparing phased genotype data to a reference sample  
57 containing representative phased whole genome sequence data. This method allows the researcher to  
58 assign the ancestral origin for each ancestry component for subsequent analyses, providing a better  
59 framework for control over population structure than only considering global admixture, as is  
60 accomplished with including principal components as covariates in statistical testing. There are  
61 several newly established analysis methods and tools that are tailored to admixed populations that  
62 implement LAI to deconvolute continental ancestry components in admixed samples. This work has  
63 shown that LAI can improve discovery power in genome-wide association studies for identifying  
64 ancestry-specific hits<sup>8</sup>, improve in polygenic risk scoring<sup>9</sup>, can be meaningful in evolutionary  
65 research<sup>10</sup>, assist in characterizing gene-gene interactions<sup>11</sup>, as well as provide more meaningful  
66 patient stratification in precision medicine<sup>12,13</sup>. To ensure and maximize the success of improving  
67 accuracy and statistical power in analyses involving LAI for admixed samples, it is paramount that  
68 local ancestry is correctly called in the individual haplotypes. However, limitations in available  
69 reference panels for many understudied populations hinder analysis, and the reference panel  
70 characteristics and algorithm parameters that result in optimal local ancestry inference accuracy across  
71 populations are still not firmly established.

72           One of the main features affecting LAI analysis is the reference panel used to infer local  
73 ancestry on the target sample. Reference panels are broadly required for genomic pipelines including  
74 LAI, yet are often sparse for admixed populations, particularly those who have some Amerindigenous  
75 (AMR) ancestry. Further, many reference samples are themselves admixed, which complicates  
76 assigning ancestral tracts, often resulting in admixed populations having diminished accuracy if  
77 reference panel homogeneity is assumed. As such, LAI may perform differentially for different  
78 ancestry components or populations such that there is an unmet need in establishing guidelines for  
79 best practice for diverse cohorts who may not have large, well-matched banks of reference samples to  
80 train algorithms on.

81           Here, we comprehensively test strategies for conducting LAI using existing reference  
82 resources with simulated "truth" genomic datasets reflective of the demographic history of Latin  
83 America to identify features that result in the best true positive rates. Latin American (LatAm)  
84 populations have a complex ancestry makeup resulting from past admixture events from multiple  
85 continental areas. Though the specific patterns vary between different geographic regions, historical  
86 admixture events generally involved substantial contributions from Amerindigenous (AMR),  
87 European (EUR) and/or African (AFR) populations<sup>14-16</sup>. Thus, the genomes of Latin American  
88 individuals are complex mosaics of different ancestral tracts that vary in length depending on the  
89 historical timing of when pulses of admixture occurred. We wish to highlight that we only describe  
90 inference of individuals genetic ancestry throughout this manuscript, rather than any metric of self-  
91 identification.

92           In our tests, we modify key parameters affecting LAI performance, including: 1) how well  
93 matched the reference panel is to the sample, 2) the absolute size of the panel, 3) the presence of  
94 admixture in the reference sample, 4) genomic data type/the number of variants (i.e. genotyping  
95 arrays vs whole genome sequencing data), as well as 5) demographic features of the cohort (global  
96 admixture proportions and timing of admixture events), and 6) parameter selection in LAI models  
97 (e.g. window size, number of EM iterations) (Figure 1). This informs best practice for researchers  
98 when conducting LAI on LatAm and other admixed populations to produce the highest accuracy  
99 results.

100

## 101 **Methods**

### 102 *Dataset generation and quality control*

103 To generate both a simulated truth dataset and comparison reference panels, we used data from the  
104 jointly called dataset of 1000 genomes (1KG) and Human Genome Diversity Project (HGDP) on  
105 GRch38<sup>17-19</sup>. For our three-way admixed analyses, we used data from Amerindigenous populations  
106 from HGDP (Karitiana, Surui, Colombian, Maya, Pima) as well as the Peruvians from Lima, Peru  
107 (PEL) and, in one test, East Asian (EAS) populations from 1KG to capture AMR ancestry. We wish

108 to clarify that in this manuscript we use the term ‘AMR’ to refer to the Amerindigenous ancestry  
109 present in modern-day Latin America, rather than as a population label for admixed American  
110 samples, as has been occasionally done in prior efforts. Each of the AMR populations from HGDP  
111 was randomly split in half, with one half used for admixture simulations (N = 31) and the other used  
112 as reference sample for LAI (N=31). To keep sample sizes balanced between ancestries for  
113 simulations, we selected 30 Iberians in Spain (IBS) samples from 1KG to capture southern European  
114 ancestry and 30 Yoruba in Ibadan, Nigeria (YRI) samples from 1KG for western African ancestry.  
115 For the reference panel for LAI, we used the remaining samples from IBS and YRI populations (N=77  
116 each) and the other half of the AMR samples, to represent a common analytic scenario. We filtered to  
117 keep only unrelated individuals and excluded multiallelic or duplicated variants, as well as those with  
118 a missingness rate > 10% and minor allele frequency < 0.5%. Genomic phasing of the complete  
119 dataset was conducted using SHAPEIT4<sup>20</sup>, and after phasing we subset the populations of interest as  
120 described below for our various simulations, with some samples used to model truth individuals and  
121 some used as LAI reference. By using distinct samples for our sample generation and reference panels  
122 we have an unbiased estimate of accuracy, at the cost of reducing the reference sample size.

123

#### 124 *Simulating truth admixed haplotypes*

125 Because of the reference sample size limitations, we simulated 60 haplotypes for each admixed  
126 cohort. Sample sizes for the reference component ancestries were selected to be equivalent to avoid  
127 biases due to unbalanced representation, and the terminal node size flag (-n 5) was implemented in  
128 LAI runs to further account for any sample size differences.

129 Latin America is a highly diverse region, and cohort admixture proportions vary widely  
130 depending on the country and even within each country<sup>14,15,16</sup>. Here, we simulated cohorts with six  
131 global ancestry patterns based on common ancestry proportions observed across Latin America.  
132 Briefly, in these simulations, one pulse of admixture is simulated at a designated point in time with  
133 specified global ancestry proportions contributed from the relevant source populations, after which  
134 haplotypes taken from the reference dataset are copied from the previous generation until the present,

135 with tract switches informed by a recombination map. We used here the hg38 HapMap combined  
136 recombination map which includes representatives from relevant global populations<sup>21</sup>. This results in  
137 a simulated truth dataset that is highly similar to modern empirical LatAm cohorts but has known  
138 phase and local ancestry which can be used for method benchmarking.

139 We tested four two-way models of AMR/EUR admixture and two three-way models of  
140 AMR/EUR/AFR admixture. In the two-way models, we compared the effects of ancestry proportions  
141 in models with: average proportions for a two-way LatAm individual (70% AMR/30% EUR, termed  
142 ‘average two-way model’)<sup>15</sup>, even AMR/EUR proportions (‘even model’), and two models to analyze  
143 the effect of extreme ancestry proportions, each with 5% of one ancestry and 95% of the other  
144 (‘extreme models’). This allows us to assess the performance that may be expected in a typical two-  
145 way admixed empirical sample, as well as assess features of sample composition influencing accuracy  
146 performance.

147 For three-way models, we tested a model of average proportions for a three-way admixed  
148 Latin American individual (15% AMR/60% EUR/25% AFR - average proportions for a Brazilian  
149 individual, termed ‘average 3-way model’)<sup>22</sup> and an even-proportioned model. Simulations were  
150 conducted using the admix-simu tool<sup>23</sup>. We simulated the average 3-way model in three different  
151 admixture demographic scenarios, considering a single pulse of admixture at 9 generations, 12, or 17  
152 generations ago<sup>14</sup>. This allowed us to evaluate the impact of varying tract lengths on true positive LAI  
153 rates. All other models were simulated considering a single pulse of admixture at 9 generations ago  
154 for the sake of comparability. In the simulation of the admixture model that has 3-way average LatAm  
155 proportions and a pulse of admixture 12 generations ago, we used data for all autosomes to obtain the  
156 highest precision. This admixture model was landed upon as it is reflective of the intermediate  
157 admixture pulse in a population migration model for the Brazilian population according to Kehdy et  
158 al., 2015<sup>14</sup>. For all other simulations, we simulated only chromosome 1 for the sake of computational  
159 efficiency.

160 For comparisons of DNA data generation type, we created a pseudo-genotype array dataset by  
161 selecting all SNVs present in the Global Screening Array (Illumina GSA) from our WGS-density

162 simulation reference dataset, a genotyping array that has been regularly used for non-European  
163 datasets. To test the effect of imputation on LAI accuracy, given that imputation is a typical step in  
164 cohort data processing for genomic analyses such as GWAS, we imputed the simulated haplotypes  
165 with SNP array-density sites using the TOPMed panel<sup>24</sup> imputation server and filtered imputed sites  
166 with  $> 0.8$  INFO score and  $MAF > 0.005$ .

167

### 168 ***Local Ancestry Inference***

169 Local Ancestry was deconvoluted using RFMix v1.5.4<sup>25</sup>. We used the TrioPhased option, with a base  
170 window size of 0.2 cM, terminal node size of 5, 2 Expectation–maximization (EM) iterations, with  
171 reference panels reanalyzed in EM to account for any admixture present in the reference (flags -w 0.2,  
172 -n 5, -e 2, and --use-reference-panels-in-EM, respectively), and the number of generations since  
173 admixture was specified depending on the simulation model (9, 12 or 17). For reference panel testing,  
174 we used three different reference panel combinations from HGDP and 1KG (AMR/EUR or  
175 AMR/EUR/AFR) that varied only in the AMR reference samples, given that this group has much less  
176 representation in reference panels relative to the other two ancestries. This was done to benchmark  
177 how variations in the reference for this ancestry impacts LAI accuracy with particular attention to  
178 improving AMR accuracy given the limitations of available reference resources. The three reference  
179 panels for LAI were constructed using, for EUR and AFR components respectively, the remaining  
180 IBS and YRI samples from 1KG not used in the simulations ( $N_{IBS} = 77$ ,  $N_{YRI} = 77$ ). For the AMR  
181 component, the three panels varied as follows: 1) *Well matched to the target but low sample size*:  
182 using the other half of HGDP-AMR samples ( $N = 30$ ); 2) *Medium sample size containing some*  
183 *admixture in AMR component*: using the 1KG sample from Lima, Peru (PEL) ( $N = 85$ ); 3) *Large*  
184 *sample size but AMR component poorly matched to the target*: 1KG-PEL ( $N = 85$ ) plus 1KG East  
185 Asian (EAS) populations ( $N = 505$ ). We included the EAS population on panel 3 to capture highly  
186 diverged AMR ancestry considering the demographic history of human migrations, since the  
187 ancestors of modern Amerindigenous peoples of the Americas migrated from East Asia across the  
188 Bering Strait around fifteen thousand years ago<sup>26</sup>. As such, AMR and EAS ancestry are less diverged

189 than other ancestral components, but even so, this composition makes for a poorly matched reference  
190 panel to the target data. Panel 2 represents the procedure most commonly conducted in current  
191 studies. We used the LAI reference panel containing the AMR samples described in panel 1 for all  
192 comparisons that did not involve reference panel testing.

193

#### 194 *Statistical Analysis*

195 We quantified the LAI true positive rates (TPR) for each run to assess their respective  
196 performance. We define TPR as follows: for each run we calculated the sums of genomic positions for  
197 which a given ancestry was correctly called compared to the simulated true ancestry for that position,  
198 divided by the total number of positions for that ancestry overall in the cohort, in each simulated  
199 haplotype (Supplementary Figure 1). We analyzed the best-guess ancestry calls output by RFMix,  
200 regardless of confidence (Forward-Backward) estimates. We computed ancestry-specific TPR to  
201 assess if there was differential LAI performance depending on the background truth ancestry and  
202 tested for statistically significant differences using the Wilcoxon rank-sum test. Significance was  
203 considered when the Bonferroni-adjusted p-value ( $p\text{-adj} < 0.05$ ).

204

#### 205 **Results**

##### 206 *Impact of demography and ancestry proportions on LAI performance*

207 We compared the effect of different demographic models on LAI performance. Specifically,  
208 we assessed the impact of varying component ancestry proportions in two and three-way models as  
209 well as different generation times since an admixture pulse occurred (considering a single pulse) by  
210 simulating 9, 12 and 17 generations since admixture.

211 In general, LAI accuracy for a given ancestry increased as that global ancestry proportion  
212 increased (Figure 1, Table 1). A low global ancestry percentage tended both to decrease the accuracy  
213 and result in larger standard deviations, as observed in the “extreme proportions” simulations (95%  
214 EUR/5% AMR and 5% EUR/95% AMR, Figure 1). In all tested models, we additionally observed  
215 significantly lower true positive rates for the AMR component ( $p\text{-adj} < 0.05$ ). This result was



216 consistent for all chromosomes and demographic models tested (Figure 1, Figure 2, Supplementary  
217 Table S1).

218 The number of generations since admixture had a slight impact on TPR, which modestly  
219 decreased in older admixture events (17 generations ago) compared to more recent ones (9  
220 generations ago). We observed statistically significant differences in TPR between the 9 and 12  
221 generations models in the AFR and EUR components, and 12 and 17 generations models in the AFR  
222 component ( $p\text{-adj} < 0.05$ ). In all models the 17 generations since admixture model had the lowest  
223 accuracy (Figure 1-B, Table 1, Supplementary Table S1). This is likely due to increased difficulty in  
224 painting short ancestry tracts; the further back in time a pulse occurred, the shorter the relative  
225 ancestral tracts will be in the current day as recombination breaks ancestral stretches down over  
226 time<sup>27</sup>. The general trends observed in relative TPR per ancestry proportion were the same regardless  
227 of admixture pulse generation times.

228

### 229 *Impact of reference panel compositions*

230 For benchmarking the performance of different reference panel compositions, we tested  
231 multiple AMR/EUR and AMR/EUR/AFR reference panel combinations comprising individuals from  
232 the Human Genome Diversity Project (HGDP) and the Thousand Genomes Project (1KG). We  
233 organized our test reference panels to reflect: 1) a very well-matched panel but with low sample size;  
234 2) a moderate sized panel that includes admixed individuals in the reference versus restricting to only  
235 homogeneous individuals; or 3) a very large reference but that is poorly matched.

236 The average TPRs were similar across the three tested reference panels, and had no  
237 statistically significant differences ( $p\text{-adj} < 0.05$ , Figure 2-A, Supplementary table S2), however the  
238 small but well-matched reference panel (N=184) resulted in considerably faster running time  
239 compared to the admixed AMR reference panel (N=239) and the large but unmatched reference panel  
240 (N=659), which took three and sixteen times longer to complete a LAI run for chromosome 1,  
241 respectively. Figure 2-A and Table 1 summarize the TPR results for each reference panel.

242

243 ***Impact of genetic data type***

244 We assessed the impact of genetic data type by comparing true positive rates (TPR) of LAI  
245 with simulations generated from WGS-density reference, a subset of SNP array variants, and a dataset  
246 of imputed variants (using the subset of SNP array variants as input), in the average 3-way (15%  
247 AMR/ 60% EUR/ 25% AFR) admixture model simulation considering 12 generations since  
248 admixture. We selected genomic variants targeted by the GSA chip for these tests, as this is a  
249 commonly used array for diverse datasets.

250 LAI run on the WGS-density simulated dataset achieved better TPR for AFR and EUR  
251 ancestry components than the SNP array-density dataset ( $p\text{-adj} < 0.05$ ), likely due to fuller haplotype  
252 coverage, but was roughly 6 times slower to complete for all autosomes. Imputation slightly improved  
253 LAI calls for these ancestry components compared to the SNP-array only runs, indicating increased  
254 SNP density improved performance. These trends were different in the AMR component, however, in  
255 which we observed no statistically significant differences between either WGS, SNP array and  
256 imputed datasets, and observed a decrease in TPR following imputation (the lowest TPR in this  
257 component), although not statistically significant (Table 1, Figure 2-B, Supplementary table S2).  
258 Additionally, we performed a validation analysis of the imputation accuracy results by selecting only  
259 the original SNP-array sites from the LAI results of the imputed dataset, to observe if imputation  
260 changed LAI on these sites which could lead to changes in accuracy performance. We observed no  
261 significant changes in TPR compared to the full imputed dataset (Supplementary Figure 2).

262

263 ***RFmix window size parameter changes do not improve LAI accuracy from the default value***

264 As RFmix calls LA with a sliding window approach, we tested whether halving or doubling  
265 the default window size improved calls, which could change results especially at the borders of  
266 chromosomes that may only have anchoring haplotype information on one side of the window.

267 We found that halving the default window size to 0.1cM did not significantly change TPR for  
268 the AFR and AMR components, and significantly lowered TPR in the EUR component compared to  
269 the default 0.2 cM ( $p\text{-adj} < 0.0001$ ). Doubling the default window size to 0.4 cM significantly

270 decreased TPR for the AFR component ( $p\text{-adj} < 1e\text{-}5$ ) but did not significantly change for the other  
271 components compared to the default (Figure 2-C, Table 1, Supplementary table S2). As such, we  
272 recommend retaining the default 0.2cM window size for RFmix runs.

273 We additionally examined the ForwardBackward probability estimates from RFMix to check  
274 how confident the algorithm estimated the wrong calls, as such confidence estimates could be a  
275 readily implemented filter to remove poorly called loci. We observed, however, that miscalls had high  
276 confidence estimates, therefore setting a stringent filter for ForwardBackward probabilities in an  
277 attempt to reduce LAI miscalls would not be sufficient to improve results.

278

### 279 *LAI miscalls are more frequent in certain genomic locations, but vary between cohorts*

280 When miscalls in LAI occurred, we observed that although they may occur at any point in the  
281 genome, they were more frequent around telomere and centromere regions (Figure 3-A, 3-B,  
282 Supplementary figures 3-8). Considering sites with over 10% miscalls in both three-way model runs  
283 and considering a window of 1kb upstream and downstream, we observed that these regions can span  
284 or flank genes (Supplementary tables S3-S4), most of which have been previously associated in  
285 GWAS studies according to GWAS Catalog (Supplementary Tables S5-S6). We compared these sites  
286 with low complexity regions from the UCSC RepeatBrowser hg38 dataset and observed an overlap of  
287 >98% in both models. It is important to note, however that the sites/regions with over 10% miscalls  
288 varied between the admixture models in which we ran this analysis, therefore we do not supply a list  
289 of regions that will be more frequently miscalled, as this may vary between cohorts.

290

### 291 *LAI miscalls occur with a consistent error mode*

292 We summed miscall counts and divided them between "error modes", i.e.: how many truth  
293 sites of one ancestry are being miscalled as each of the other ancestries. This allowed us to  
294 characterize trends in miscall directions and observe whether one ancestry was systematically being  
295 over- or under-called than other. We observed that the most common direction for miscalls to occur

296 was for truth AMR sites to be incorrectly called EUR (Figure 3-C, 3-D, Supplementary figures 9-12).

297 The second most frequent error mode was EUR positions being miscalled AFR.

298

## 299 **Discussion**

300 In this study, we evaluated characteristics impacting the performance of LAI for a range of  
301 two and three-way admixed demographic models reflective of many Latin American populations.  
302 Specifically, we assessed the impact of reference panel composition, demographic features such as the  
303 proportions of major ancestry groups and number of generations since admixture in the cohort,  
304 genetic data technology (genotyping arrays versus whole genome sequencing), the impact of  
305 imputation, and LAI analytic thresholds in affecting performance for diverse cohorts. Given the high  
306 LAI accuracy observed in the literature for 2-way admixed AFR/EUR cohorts<sup>8</sup>, we focused our  
307 analyses in this manuscript on determining the best practices for cohorts involving AMR, as the  
308 smaller divergence time between EUR and AMR tracts poses a challenge for deconvolution, as does  
309 the particularly limited availability relevant reference individuals for AMR ancestry. Thus, we  
310 focused the construction of our reference panel tests in the service of optimizing AMR accuracy.

311 These benchmarks allow us to provide a set of recommendations for parameter and panel  
312 selection to achieve optimal LAI performance in LatAm populations. Specifically, comparing the  
313 performance of different reference panel compositions, we observed that there was not significant  
314 difference in accuracy across the three panels (well-matched but small sample size, medium size with  
315 some degree of admixture in the reference, and large but poorly matched to the target cohort),  
316 although we do observe a large difference in runtime, with the small but well-matched panel running  
317 substantially faster than the other panels. Given the high computational burden required by LAI,  
318 having a quicker runtime for analysis is an important point of consideration in practical use. As such,  
319 a curated reference panel reflective of the ancestries present in the target cohort appears to be the best  
320 option for LAI reference panel construction. Importantly, across all demographic and reference panel  
321 models tested, Amerindigenous ancestry tracts suffer from notably reduced accuracy as compared to  
322 European and African tracts. This is likely due to there being less representative (and less

323 homogeneous) reference data for AMR ancestry in existing reference resources. Moving forward, it  
324 will be vital for efforts to focus on ethical recruitment of more diverse and geographically distributed  
325 reference samples in large scale data collection efforts to maximize the performance of LAI across all  
326 ancestry backgrounds.

327       Regarding ancestry proportions in the admixed simulations, we observed that overall, having a  
328 higher proportion of an ancestry in the simulation improved the true positive rates for that ancestry in  
329 LAI. When ancestries represent a very small (e.g. 5%) global proportion, all ancestries suffer, with  
330 AMR suffering the most.

331       We investigated LAI miscalls in the realistic three-way simulation models to evaluate the typical  
332 error mode when wrong calls are produced. Specifically, we examined the rates of miscalls for each  
333 ancestry component to: characterize trends in the relative amount and direction of miscalls, see if a  
334 particular ancestry was being systematically over or under-called, document if there were genomic  
335 regions where miscalls were most frequent, to identify other factors that could be driving error modes,  
336 and to assess if alterations to RFMix parameters could improve miscall rates. Investigating the typical  
337 error modes when LAI miscalls occur, we observed a much higher frequency of miscalls in the  
338 direction of calling simulated AMR regions as EUR compared to other miscall directions. This may  
339 be explained by the smaller genetic divergence between AMR and EUR haplotype tracts than either is  
340 to AFR, resulting in closer haplotype similarity. This finding implies that reference panels will need to  
341 be grown substantially to confidently assess within-continental ancestral components for many  
342 geographic regions. The specific direction of AMR/EUR miscalls being dominated in the direction of  
343 AMR to EUR rather than vice versa can be explained by the substantially (2x) larger sample size  
344 available for EUR compared to AMR. Another point of consideration is that the AMR reference  
345 samples themselves have some degree of admixture with European ancestry, which adds uncertainty  
346 to the model, though we did implement EM procedures to attempt to correct for this. These results are  
347 consistent with miscall trends observed in other studies of diverse populations<sup>28</sup>. As this prior cited  
348 work was done with older LAI software than RFMix, we have confirmed that this error mode is

349 consistent between different LAI algorithms and therefore likely to be driven by the genetic data,  
350 rather than a feature specific to RFMix.

351 Beyond error modes, we observe that miscall regions do not appear randomly across the genome,  
352 but are most likely to fall in areas that mark the edges of haplotypes, like centromere and telomeres  
353 (Figure 3A, 3-B, Supplementary Figures 3-8). We note several areas that had elevated miscall rates  
354 (higher than 10% miscalls). ForwardBackward probabilities for the LAI algorithm were still confident  
355 in such areas and tweaking RFmix parameters was insufficient to correct them. As these regions may  
356 vary between cohorts given that different models resulted in different regions with elevated false  
357 positive rates, we do not recommend blanket masking of those observed in this study. This highlights  
358 the importance of ensuring good LAI accuracy for gene discovery and other statistical genomics  
359 efforts, as misclassification both soaks up power in LAI-informed GWAS as well as can lead to false  
360 positive associations due to technical ancestry miscalls<sup>8</sup>. Importantly, miscall regions may contain  
361 genes of interest, so care should be taken to validate, for example, GWAS hits in border haplotype  
362 areas that show elevated miscall rates. Inflation of miscalls at particular regions could also impact the  
363 interpretation of other statistical genetics efforts, such as admixture mapping or evolutionary scans of  
364 selection that utilize local ancestry enrichment. We observed an inflation of miscalls in low  
365 complexity regions, such as short and long interspersed nuclear elements (SINE/LINE), DNA repeats  
366 and micro-satellites, therefore additional care should be considered when analyzing these regions  
367 and/or genes in close proximity. This could be due to the fact that low complexity regions are usually  
368 more challenging to map<sup>29,30</sup>, and/or are evolutionarily conserved<sup>31,32</sup>, with little variation across  
369 ancestry groups, and therefore these regions would be more prone to error in LAI. The development  
370 of methods that incorporate repeat polymorphisms, multi-allelic variants and other complex forms of  
371 genetic variation in genome-wide analyses may help improve LAI accuracy.

372 Examining how different DNA data types impact LAI performance, we observe that WGS  
373 and SNP array simulated data resulted in similar TPR estimates for genotyped sites, although having  
374 more variants in the dataset improved estimates. We also observe that LAI performs nearly as well on  
375 imputed data as directly genotyped data when a large and diverse reference panel is used. This

376 suggests that, provided imputation can be performed with a representative reference panel, LAI calls  
377 on imputed data may be confidently utilized for downstream efforts. Out of an abundance of caution,  
378 we recommend setting a stringent INFO threshold (e.g. 0.8) for imputed sites to ensure high  
379 confidence calls. We note, additionally, that non-significant differences in TPR in the context of this  
380 work do not mean that the differences that we observe are not relevant since small differences in LAI  
381 accuracy can impact statistical power in downstream applications<sup>8</sup> and may represent a difference  
382 observed in a large number of sites in the genome. Expanding available reference samples to contain  
383 representative haplotypes from diverse and understudied populations would improve the quality of  
384 imputation as well as LAI.

385         Of course, this work has some important limitations which must be considered. As the focus  
386 of the present study is LatAm populations, we limited our demographic models to those involving 2-  
387 way admixture between AMR and EUR or 3-way admixture between AMR, EUR, and AFR, which  
388 represents the majority of Latinx populations. We note, however, that some LatAm populations have  
389 other patterns than those directly benchmarked here. Despite this, the broader trends in LAI  
390 performance identified in this work should hold across demographic models beyond the specific use  
391 cases simulated in this manuscript. We also note that while we appreciate that there is a high level of  
392 diversity within continental regions<sup>14,33</sup>, only continental-level ancestry was able to be assessed here  
393 due to limitations in available reference panel geographic coverage. Similarly, having a small number  
394 of available reference AMR samples limited the number of individuals available for simulating and  
395 running LAI, which limits variability in the data for this component in comparison to EUR and AFR.  
396 Improved LAI call rates and finer scale LAI resolution would be possible in the future if reference  
397 panels are expanded. Regarding software, we have benchmarked only RFMix v1 in this work, as prior  
398 work has demonstrated that RFMixv1 performed the best in comparison to other methods for multi-  
399 way admixed samples<sup>34</sup>. We expect the trends observed here to be consistent across LAI software,  
400 though further benchmarking would be needed to confirm this.

401 In conclusion, in our reference panel benchmarking, the best cost-benefit in terms of LAI  
402 accuracy and speed is to use a well-matched reference even if it has a lower sample size. Examining  
403 the ancestry-specific performance of LAI across reference panels, we observed consistently lower  
404 performance for the AMR ancestry component across all simulation settings compared to EUR and  
405 AFR. Unfortunately, this inequity could not be overcome by any of the tested modifications to  
406 reference panel, LAI software parameters, or features of genetic data. The best way to improve AMR  
407 performance would be to increase the well-matched reference panel's sample size, underscoring the  
408 importance of furthering recruitment of larger and more representative reference samples for  
409 understudied populations. Given the high proportion of the global population that contains admixed  
410 ancestry and the fact that populations are getting increasingly admixed over time<sup>35</sup>, it is timely to  
411 establish the optimal methods for well-calibrated genomic analyses in admixed populations.

#### 412 **Acknowledgments**

413 This study was supported by Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP;  
414 fellowship number: 2021/09584-1). Financial support for the PTSD-PGC Ancestry Working Group  
415 was provided by the National Institute of Mental Health (NIMH; R01MH106595). EGA was  
416 supported by the National Institute of Mental Health (K01 MH121659, R01 HG012869), the Caroline  
417 Wiess Law Fund for Research in Molecular Medicine, and the ARCO Foundation Young Teacher-  
418 Investigator Fund at Baylor College of Medicine.

419

#### 420 **Author Contributions**

421 J.M. conducted analysis and wrote the manuscript. A.X.M. and N.N.S. assisted with software. C.Z.  
422 reviewed the manuscript. C.M.N. and S.B. advised on the project. E.G.A. and M.S. conceptualized,  
423 supervised, and funded the project as well as contributed to the manuscript. All authors reviewed and  
424 approved the final manuscript.

425

#### 426 **Declaration of interests**



427 The authors declare no competing interests.

428

#### 429 **Web Resources**

430 Admix-simu: <https://github.com/williamslab/admix-simu/>

431 RFMix V.1 <https://github.com/indraniel/rfmix>

432 Shapeit v4 <https://odelaneau.github.io/shapeit4/>

433 Pipeline used to prepare data for RFMix v1: [https://github.com/armartin/ancestry\\_pipeline/](https://github.com/armartin/ancestry_pipeline/)

434 R package utilized to create TPR figures: <https://phanstiellab.github.io/plotgardener/>

435 UCSC RepeatBrowser: <https://repeatbrowser.ucsc.edu/data/>

436 TOPMed Imputation Server: <https://imputation.biodatacatalyst.nhlbi.nih.gov/>

437 Thousand Genomes Project: <http://ftp.1000genomes.ebi.ac.uk/>.

438 The Human Genome Diversity Project:

439 [ftp://ngs.sanger.ac.uk/production/hgdp/hgdp\\_wgs.20190516/statphase/](ftp://ngs.sanger.ac.uk/production/hgdp/hgdp_wgs.20190516/statphase/).

440 Jointly called HGDP+1kG: <https://gnomad.broadinstitute.org/downloads#v3-hgdp-1kg>

441 HapMap GRCh38 recombination map:

442 [http://bochet.gcc.biostat.washington.edu/beagle/genetic\\_maps/plink.GRCh38.map.zip](http://bochet.gcc.biostat.washington.edu/beagle/genetic_maps/plink.GRCh38.map.zip)

443

#### 444 **Data/Code Availability**

445 Code generated in this project for simulating admixed data and quantifying LAI true positive rates is

446 freely available on github at <https://github.com/Atkinson-Lab/LAI-sims-accuracy>.

447

#### 448 **References**

449 1. Sul, J.H., Martin, L.S., and Eskin, E. (2018). Population structure in genetic studies: Confounding  
450 factors and mixed models. *PLoS Genet.* 14, e1007309.

451 2. Sohail, M., Maier, R.M., Ganna, A., Bloemendal, A., Martin, A.R., Turchin, M.C., Chiang, C.W.,  
452 Hirschhorn, J., Daly, M.J., Patterson, N., et al. (2019). Polygenic adaptation on height is  
453 overestimated due to uncorrected stratification in genome-wide association studies. *Elife* 8,.

454 3. Popejoy, A.B., and Fullerton, S.M. (2016). Genomics is failing on diversity. *Nature* 538, 161–164.

- 455 4. Sirugo, G., Williams, S.M., and Tishkoff, S.A. (2019). The Missing Diversity in Human Genetic  
456 Studies. *Cell* *177*, 1080.
- 457 5. Martin, A.R., Gignoux, C.R., Walters, R.K., Wojcik, G.L., Neale, B.M., Gravel, S., Daly, M.J.,  
458 Bustamante, C.D., and Kenny, E.E. (2020). Human Demographic History Impacts Genetic Risk  
459 Prediction across Diverse Populations. *Am. J. Hum. Genet.* *107*, 788–789.
- 460 6. Gurdasani, D., Barroso, I., Zeggini, E., and Sandhu, M.S. (2019). Genomics of disease risk in  
461 globally diverse populations. *Nat. Rev. Genet.* *20*, 520–535.
- 462 7. Peterson, R.E., Kuchenbaecker, K., Walters, R.K., Chen, C.-Y., Popejoy, A.B., Periyasamy, S.,  
463 Lam, M., Iyegbe, C., Strawbridge, R.J., Brick, L., et al. (2019). Genome-wide Association Studies in  
464 Ancestrally Diverse Populations: Opportunities, Methods, Pitfalls, and Recommendations. *Cell* *179*,  
465 589–603.
- 466 8. Atkinson, E.G., Maihofer, A.X., Kanai, M., Martin, A.R., Karczewski, K.J., Santoro, M.L., Ulirsch,  
467 J.C., Kamatani, Y., Okada, Y., Finucane, H.K., et al. (2021). Tractor uses local ancestry to enable the  
468 inclusion of admixed individuals in GWAS and to boost power. *Nat. Genet.* *53*, 195–204.
- 469 9. Marnetto, D., Pärna, K., Läll, K., Molinaro, L., Montinaro, F., Haller, T., Metspalu, M., Mägi, R.,  
470 Fischer, K., and Pagani, L. (2020). Ancestry deconvolution and partial polygenic score can improve  
471 susceptibility predictions in recently admixed individuals. *Nature Communications* *11*,.
- 472 10. Padhukasahasram, B. (2014). Inferring ancestry from population genomic data and its  
473 applications. *Front. Genet.* *5*, 204.
- 474 11. Aschard, H., Gusev, A., Brown, R., and Pasaniuc, B. (2015). Leveraging local ancestry to detect  
475 gene-gene interactions in genome-wide data. *BMC Genet.* *16*, 124.
- 476 12. Duconge, J., and Ruaño, G. (2010). The Emerging Role of Admixture in the Pharmacogenetics of  
477 Puerto Rican Hispanics. *J. Pharmacogenomics Pharmacoproteomics* *1*,.
- 478 13. Goetz, L.H., Uribe-Bruce, L., Quarless, D., Libiger, O., and Schork, N.J. (2014). Admixture and  
479 clinical phenotypic variation. *Hum. Hered.* *77*, 73–86.
- 480 14. Kehdy, F.S.G., Gouveia, M.H., Machado, M., Magalhães, W.C.S., Horimoto, A.R., Horta, B.L.,  
481 Moreira, R.G., Leal, T.P., Scliar, M.O., Soares-Souza, G.B., et al. (2015). Origin and dynamics of  
482 admixture in Brazilians and its effect on the pattern of deleterious mutations. *Proc. Natl. Acad. Sci. U.*  
483 *S. A.* *112*, 8696–8701.
- 484 15. Ruiz-Linares, A., Adhikari, K., Acuña-Alonzo, V., Quinto-Sanchez, M., Jaramillo, C., Arias, W.,  
485 Fuentes, M., Pizarro, M., Everardo, P., de Avila, F., et al. (2014). Admixture in Latin America:  
486 geographic structure, phenotypic diversity and self-perception of ancestry based on 7,342 individuals.  
487 *PLoS Genet.* *10*, e1004572.
- 488 16. Homburger, J.R., Moreno-Estrada, A., Gignoux, C.R., Nelson, D., Sanchez, E., Ortiz-Tello, P.,  
489 Pons-Estel, B.A., Acevedo-Vasquez, E., Miranda, P., Langefeld, C.D., et al. (2015). Genomic insights  
490 into the ancestry and demographic history of South America. *PLoS Genet.* *11*, e1005602.
- 491 17. Byrska-Bishop, M., Evani, U.S., Zhao, X., Basile, A.O., Abel, H.J., Regier, A.A., Corvelo, A.,  
492 Clarke, W.E., Musunuri, R., Nagulapalli, K., et al. (2021). High coverage whole genome sequencing  
493 of the expanded 1000 Genomes Project cohort including 602 trios.
- 494 18. Bergström, A., McCarthy, S.A., Hui, R., Almarri, M.A., Ayub, Q., Danecek, P., Chen, Y., Felkel,  
495 S., Hallast, P., Kamm, J., et al. (2020). Insights into human genetic variation and population history

- 496 from 929 diverse genomes. *Science* 367,.
- 497 19. Koenig, Z., Yohannes, M.T., Nkambule, L.L., Zhao, X., Goodrich, J.K., Kim, H.A., Wilson,  
498 M.W., Tiao, G., Hao, S.P., Sahakian, N., et al. (2024). A harmonized public resource of deeply  
499 sequenced diverse human genomes. *Genome Res.* 34, 796–809.
- 500 20. Delaneau, O., Zagury, J.-F., Robinson, M.R., Marchini, J.L., and Dermitzakis, E.T. (2019).  
501 Accurate, scalable and integrative haplotype estimation. *Nat. Commun.* 10, 5436.
- 502 21. Consortium, †the International Hapmap, and †The International HapMap Consortium (2003). The  
503 International HapMap Project. *Nature* 426, 789–796.
- 504 22. Souza, A.M. de, Resende, S.S., Sousa, T.N. de, and Brito, C.F.A. de (2019). A systematic scoping  
505 review of the genetic ancestry of the Brazilian population. *Genet. Mol. Biol.* 42, 495–508.
- 506 23. Williams, A. (2016). Admix-simu: Admix-simu: Program to simulate admixture between multiple  
507 populations (Zenodo).
- 508 24. Taliun, D., Harris, D.N., Kessler, M.D., Carlson, J., Szpiech, Z.A., Torres, R., Taliun, S.A.G.,  
509 Corvelo, A., Gogarten, S.M., Kang, H.M., et al. (2021). Sequencing of 53,831 diverse genomes from  
510 the NHLBI TOPMed Program. *Nature* 590, 290–299.
- 511 25. Maples, B.K., Gravel, S., Kenny, E.E., and Bustamante, C.D. (2013). RFMix: a discriminative  
512 modeling approach for rapid and robust local-ancestry inference. *Am. J. Hum. Genet.* 93, 278–288.
- 513 26. Raghavan, M., Steinrücken, M., Harris, K., Schiffels, S., Rasmussen, S., DeGiorgio, M.,  
514 Albrechtsen, A., Valdiosera, C., Ávila-Arcos, M.C., Malaspina, A.-S., et al. (2015). POPULATION  
515 GENETICS. Genomic evidence for the Pleistocene and recent population history of Native  
516 Americans. *Science* 349, aab3884.
- 517 27. Gravel, S. (2012). Population genetics models of local ancestry. *Genetics* 191, 607–619.
- 518 28. Seldin, M.F., Pasaniuc, B., and Price, A.L. (2011). New approaches to disease mapping in  
519 admixed populations. *Nat. Rev. Genet.* 12, 523–528.
- 520 29. Tørresen, O.K., Star, B., Mier, P., Andrade-Navarro, M.A., Bateman, A., Jarnot, P., Gruca, A.,  
521 Grynberg, M., Kajava, A.V., Promponas, V.J., et al. (2019). Tandem repeats lead to sequence  
522 assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids*  
523 *Res.* 47, 10994–11006.
- 524 30. Treangen, T.J., and Salzberg, S.L. (2011). Repetitive DNA and next-generation sequencing:  
525 computational challenges and solutions. *Nat. Rev. Genet.* 13, 36–46.
- 526 31. Enright, J.M., Dickson, Z.W., and Golding, G.B. (2023). Low Complexity Regions in Proteins and  
527 DNA are Poorly Correlated. *Mol. Biol. Evol.* 40,.
- 528 32. Lenz, C., Haerty, W., and Golding, G.B. (2014). Increased substitution rates surrounding low-  
529 complexity regions within primate proteins. *Genome Biol. Evol.* 6, 655–665.
- 530 33. Atkinson, E.G., Dalvie, S., Pichkar, Y., Kalungi, A., Majara, L., Stevenson, A., Abebe, T., Akena,  
531 D., Alemayehu, M., Ashaba, F.K., et al. (2022). Genetic structure correlates with ethnolinguistic  
532 diversity in eastern and southern Africa. *Am. J. Hum. Genet.* 109, 1667–1679.
- 533 34. Schubert, R., Andaleon, A., and Wheeler, H.E. (2020). Comparing local ancestry inference  
534 models in populations of two- and three-way admixture. *PeerJ* 8, e10090.

535 35. Institute of Medicine, Board on Health Care Services, and Committee on Future Directions for the  
536 National Healthcare Quality and Disparities Reports (2010). Future Directions for the National  
537 Healthcare Quality and Disparities Reports (National Academies Press).

538

539

## 540 **Figure Legends**

541

542 **Figure 1:** A) True positive rates for LAI in six simulated cohorts with varying proportions of 2 or 3-  
543 way admixture between AFR/EUR/AMR (displayed in order of decreasing mean TPR). These  
544 simulated haplotypes consist of chromosome 1 and considered a pulse of admixture at 9 generations  
545 ago. B) True positive rates for LAI in varying generations since admixture models for the simulated  
546 haplotype data. The haplotypes in this comparison had 15% AMR/ 60% EUR/ 25% AFR proportions  
547 of admixture in all autosomes. Significance level: \*  $\leq 0.05$ , \*\*  $\leq 0.01$ , \*\*\*  $\leq 0.001$ , \*\*\*\*  $\leq$   
548 0.0001.

549 **Figure 2:** A) True positive rates for LAI in three reference panel comparisons that vary in the AMR  
550 component, separated by ancestry component. Benchmarking was run on the model reflecting a pulse  
551 of admixture at 9 generations ago with 15% AMR/ 60% EUR/ 25% AFR proportions in chromosome  
552 1. B) True positive rates for LAI in WGS vs. SNP array (GSA) vs. Imputed data. Benchmarking was  
553 run on the model reflecting a pulse of admixture at 12 generations ago with 15% AMR/ 60% EUR/  
554 25% AFR proportions in all autosomes. C) True positive rates for LAI runs varying the RFMix  
555 window size parameter in centimorgans (cM). Benchmarking was run on the model reflecting a pulse  
556 of admixture at 12 generations ago with 15% AMR/ 60% EUR/ 25% AFR proportions in all  
557 autosomes. Significance level: \*  $\leq 0.05$ , \*\*  $\leq 0.01$ , \*\*\*  $\leq 0.001$ , \*\*\*\*  $\leq 0.0001$ .

558 **Figure 3:** A) Percentage of wrong calls per site on chromosome 1, total and separated by error mode  
559 for LAI ran on the model reflecting a pulse of admixture at 12 generations ago with 15% AMR/ 60%  
560 EUR/ 25% AFR proportions. B) Percentage of wrong calls per site on chromosome 1, total and  
561 separated by error mode for LAI ran on the model reflecting a pulse of admixture at 12 generations  
562 ago with 33% AMR/ 33% EUR/ 34% AFR proportions. C) Miscall counts separated by error mode  
563 summing all autosomes for LAI ran on the model reflecting a pulse of admixture at 12 generations  
564 ago with 15% AMR/ 60% EUR/ 25% AFR proportions. D) Miscall counts separated by error mode  
565 summing all autosomes for LAI ran on the model reflecting a pulse of admixture at 12 generations  
566 ago with 33% AMR/ 33% EUR/ 34% AFR proportions.

567

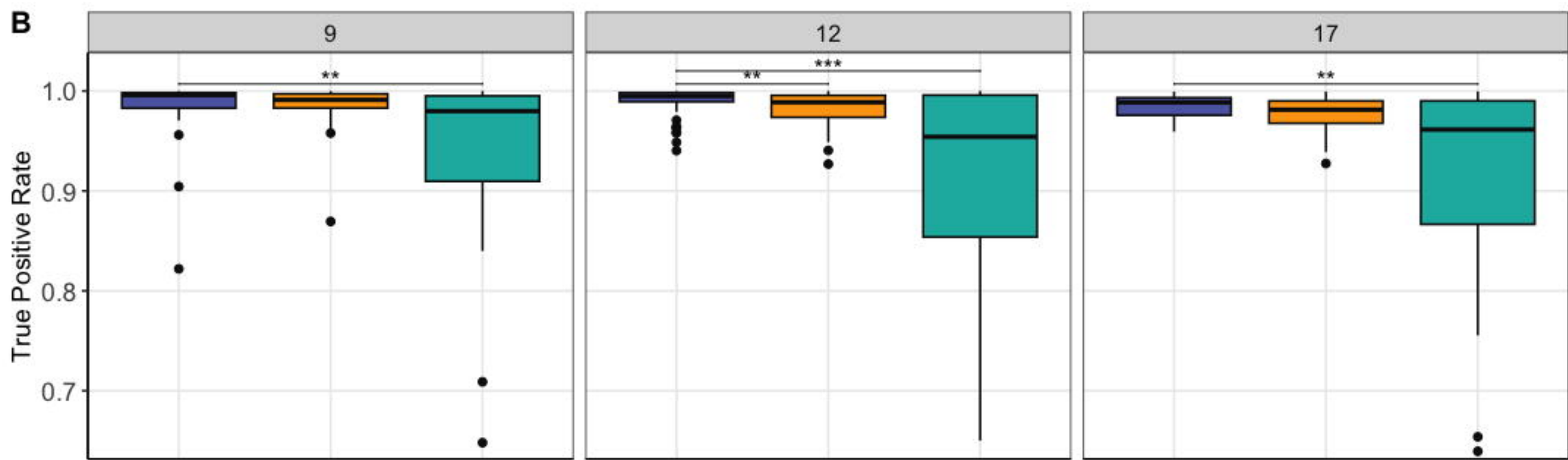
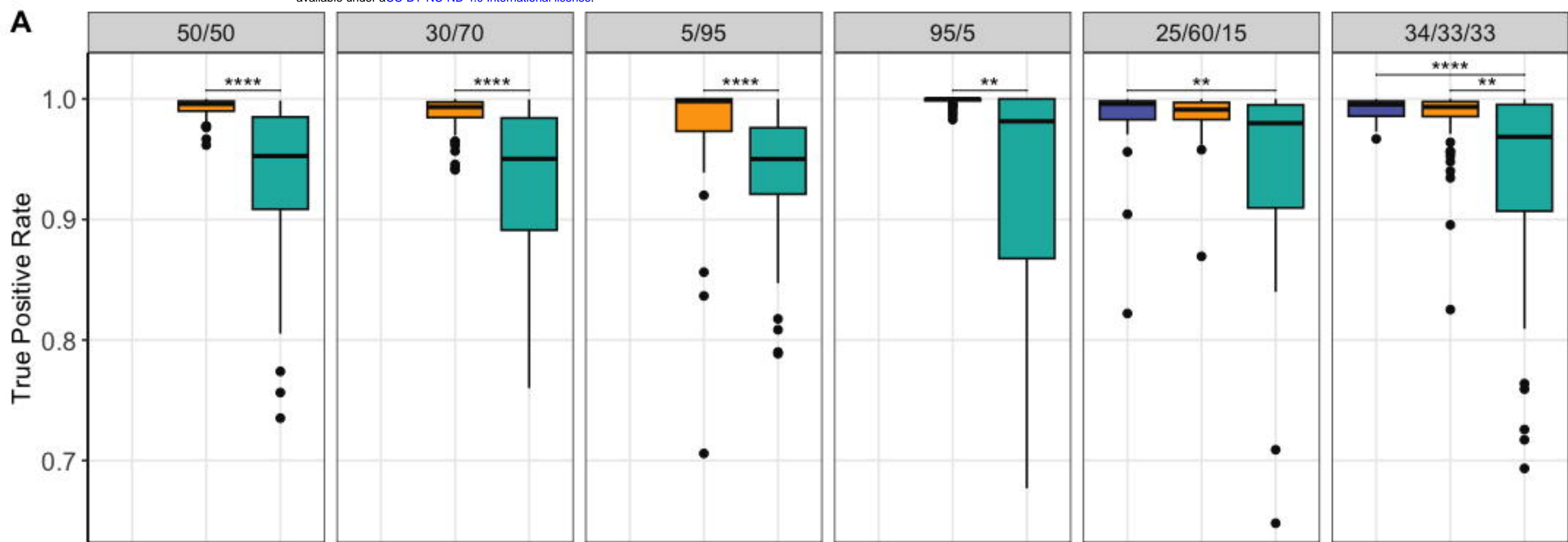
## 568 **Table Title and Legend**

569 **Table 1:** LAI true positive rate estimates per ancestry per comparison

Test			Mean TPR (SD) per ancestry			
Aim of comparison	Simulation model name	Analysis parameters*	AMR	EUR	AFR	
<b>Impact of Demography</b>	Proportions	Average 3-way	15%/60%/25%, 9gen, chr1	0.884 (0.224)	0.986 (0.018)	0.986 (0.026)
		Even 3-way	33%/33%/34%, 9gen, chr1	0.933 (0.082)	0.966 (0.130)	0.991 (0.008)
		Even 2-way	50%/50%/0, 9gen, chr1	0.918 (0.120)	0.992 (0.008)	N/A
		Average 2-way	70%/30%/0, 9gen, chr1	0.935 (0.058)	0.987 (0.014)	N/A
		Extreme proportions, high AMR	95%/5%/0, 9gen, chr1	0.937 (0.054)	0.961 (0.134)	N/A
		Extreme proportions, high EUR	5%/95%/0, 9gen, chr1	0.845 (0.280)	0.998 (0.003)	N/A
	Generations since admixture	Average 3-way	15%/60%/25%, 9gen, chr1	0.884 (0.224)	0.986 (0.018)	0.986 (0.026)
		Average 3-way	15%/60%/25%, 12gen, chr 1	0.906 (0.117)	0.982 (0.016)	0.989 (0.013)
		Average 3-way	15%/60%/25%, 17gen, chr1	0.896 (0.162)	0.977 (0.017)	0.984 (0.011)
<b>Features of data/analysis</b>	Reference Panel	Average 3-way, low N - well matched AMR reference	15%/60%/25%, 9gen chr1. HGDP AMR reference	0.884 (0.224)	0.986 (0.018)	0.986 (0.026)
		Average 3-way, medium N - admixed AMR reference	15%/60%/25%, 9gen chr1. 1KG PEL as AMR reference	0.878 (0.224)	0.986 (0.018)	0.989 (0.013)

	Average 3-way, high N - admixed + unmatched AMR reference	15%/60%/25%, 9gen chr1. 1KG PEL + EAS as AMR reference	0.875 (0.223)	0.983 (0.019)	0.991 (0.009)
Data type	Average 3-way	15%/60%/25%, 12gen, all autosomes. WGS density	0.935 (0.025)	0.983 (0.005)	0.989 (0.004)
	Average 3-way	15%/60%/25%, 12 gen, all autosomes. SNP array density	0.938 (0.029)	0.977 (0.008)	0.982 (0.004)
	Average 3-way	15%/60%/25%, 12 gen, all autosomes. Imputed SNP array density	0.927 (0.026)	0.982 (0.004)	0.987 (0.003)
Window size	Average 3-way	15%/60%/25%, 12gen, all autosomes, 0.1 cM	0.939 (0.023)	0.979 (0.005)	0.991 (0.003)
	Average 3-way	15%/60%/25%, 12gen, all autosomes, 0.2 cM (default)	0.935 (0.025)	0.983 (0.005)	0.989 (0.004)
	Average 3-way	15%/60%/25%, 12gen, all autosomes, 0.4 cM	0.923 (0.026)	0.984 (0.005)	0.980 (0.005)

570 **Table 1 Legend:** \*Analysis parameters: Admixture proportions of simulated cohort (AMR/EUR/AFR), number of generations since admixture, simulated  
 571 chromosome, other parameters.



Ancestry ■ AFR ■ EUR ■ AMR



