

1 **Admixed Populations Improve Power for Variant Discovery and Portability in Genome-**
2 **wide Association Studies**

3

4

5 Meng Lin^{1,*}, Danny S. Park², Noah A. Zaitlen³, Brenna M. Henn⁴, Christopher R. Gignoux^{1,*}

6

7 ¹ Colorado Center for Personalized Medicine, University of Colorado Anschutz Medical

8 Campus, Aurora, CO

9 ² Department of Bioengineering and Therapeutic Sciences, University of California, San

10 Francisco, CA

11 ³ Departments of Neurology and Computational Medicine, University of California Los Angeles,

12 Los Angeles, CA

13 ⁴ Department of Anthropology, Center for Population Biology and the Genome Center,

14 University of California Davis, Davis, CA

15

16

17 *Correspondence:

18 meng.lin@cuanschutz.edu

19 chris.gignoux@cuanschutz.edu

20

21

22

23 **Abstract**

24

25 Genome-wide association studies (GWAS) are primarily conducted in single-ancestry settings.

26 The low transferability of results has limited our understanding of human genetic architecture

27 across a range of complex traits. In contrast to homogeneous populations, admixed populations

28 provide an opportunity to capture genetic architecture contributed from multiple source

29 populations and thus improve statistical power. Here, we provide a mechanistic simulation

30 framework to investigate the statistical power and transferability of GWAS under directional

31 polygenic selection or varying divergence. We focus on a two-way admixed population and

32 show that GWAS in admixed populations can be enriched for power in discovery by up to 2-fold

33 compared to the ancestral populations under similar sample size. Moreover, higher accuracy of

34 cross-population polygenic score estimates is also observed if variants and weights are trained in

35 the admixed group rather than in the ancestral groups. Common variant associations are also

36 more likely to replicate if first discovered in the admixed group and then transferred to an

37 ancestral population, than the other way around (across 50 iterations with 1,000 causal SNPs,

38 training on 10,000 individuals, testing on 1,000 in each population, $p=3.78e-6$, $6.19e-101$, ~ 0 for

39 $F_{ST} = 0.2, 0.5, 0.8$, respectively). While some of these F_{ST} values may appear extreme, we

40 demonstrate that they are found across the entire phenome in the GWAS catalog. This

41 framework demonstrates that investigation of admixed populations harbors significant

42 advantages over GWAS in single-ancestry cohorts for uncovering the genetic architecture of

43 traits and will improve downstream applications such as personalized medicine across diverse

44 populations.

45

46 **Keywords:** admixture, statistical power, complex trait genetics, polygenic score, genetic
47 architecture

48
49

50 **Introduction**

51

52 Genome-wide association studies (GWAS) have allowed for significant progress in the field of
53 human complex traits. However, groups with multiple ancestral origins have seldom been a
54 primary focus in large scale genetic studies because: (1) admixed groups, along with other non-
55 European populations, have largely been underrepresented in GWAS designs in the past
56 (Bustamante et al., 2011; Popejoy and Fullerton, 2016; Martin et al., 2017a), and (2) the
57 population structure from heterogeneous ancestries in an admixed group, if not properly
58 corrected, can result in spurious correlation signals and thus greater false positive rates
59 (Rosenberg et al., 2010). However this mixture of ancestries present in admixed populations
60 provides opportunities for novel discovery. Recent advancements in methodologies tailored for
61 genetic mapping in admixed populations include disentangling of ancestry principal components
62 and relatedness in the presence of admixture (Thornton et al., 2012; Conomos et al., 2015, 2016),
63 combining local ancestry and allelic information to improve quantitative trait locus (QTL)
64 mapping (Pasaniuc et al., 2011; Shriner et al., 2011; Atkinson et al., 2021), leveraging local
65 ancestries for detection of epistasis (Aschard et al., 2015), and better fine mapping from linkage
66 disequilibrium (LD) variability in diverse groups (Zaitlen et al., 2010; Asimit et al., 2016;
67 Wojcik et al., 2019; Shi et al., 2020). Despite the fast development and practicality of these
68 methods, they have not often been applied to sample sizes of hundreds of thousands to millions
69 because study design and data collection in mega-scale cohorts routinely prioritize recruitment of
70 participants of single ancestry (Atkinson et al., 2021). This greatly impedes downstream progress,
71 such as polygenic risk score application across populations, where much lower accuracy is
72 observed in non-European populations for many traits (Duncan et al., 2019; Martin et al., 2019;
73 Cavazos and Witte, 2021).

74

75 In addition, complex traits in admixed groups potentially harbor differing genetic architectures
76 and varying environmental exposures compared to most widely studied groups such as
77 Europeans. Some biomedical traits have higher risk prevalence in admixed groups, such as
78 prostate cancer in African Americans (Bhardwaj et al., 2017; Conti et al., 2021), asthma in
79 Puerto Ricans (Lara et al., 2006; Pino-Yanes et al., 2015), obesity and type II diabetes in Native
80 Hawaiians (Maskarinec et al., 2009), and active tuberculosis in a South African admixed
81 population (Chimusa et al., 2014), which are likely attributed to elevated ancestry-specific risk
82 allele frequency. Among anthropometric traits, skin pigmentation in groups with admixed
83 ancestry harbor greater phenotypic variance than those with single ancestries (Martin et al.,
84 2017b). Here, the larger phenotypic variance is likely caused by increased polygenicity in
85 admixed groups, where in contrast some causal variants are nearly fixed in the single ancestry
86 groups due strong directional selection of skin pigmentation (e.g., rs1426654 in *SLC24A5* (Lin et
87 al., 2018)). The increase in minor allele frequencies in admixed populations compared to the
88 populations of ancestral origin could be ubiquitous in traits that have been under differential
89 processes of selection among ancestral populations or simply among populations that are deeply
90 diverged. This would theoretically result in greater power of discovery in GWAS, as the analysis
91 is most powerful for variants with higher minor allele frequency.

92
93 While genetic epidemiologists have typically focused on homogeneous populations, there are
94 clear opportunities to improve discovery in admixed populations. For example, local ancestry
95 can be leveraged to improve power in certain scenarios (e.g. (Pasaniuc et al., 2011)). In addition,
96 Zhang and Stram observed a power gain in admixed individuals in dichotomous traits compared
97 to pooled ancestral populations with stratification without environmental confounding (Zhang
98 and Stram, 2014). Here, we extend a similar framework using a flexible mechanistic model to
99 address the question of power compared to a similar-sized each ancestral population on its own,
100 across a range of allelic differentiation (as measured by F_{ST} (Weir and Cockerham, 1984)) with
101 varying narrow-sense heritability, and across a range of ancestry–phenotype associations,
102 whether driven by genetics or environment. We further extend these insights to look at power for
103 replication, whether from admixed populations to ancestral populations or *vice versa*, as well as
104 opportunities to derive trans-ethnic polygenic scores.

106 107 **Methods**

108 109 *Simulation-based power estimate between a trait and global ancestry*

110
111 The first simulator we provide in this study builds phenotypes in admixed populations using only
112 global ancestries, without involving genotype. The aim is to assess if the sample size is adequate
113 for observing a dichotomous trait by ancestry correlation. The details are described in
114 *Supplementary Notes*.

115 116 *Genotype-mediated simulation framework*

117
118 The general simulation framework consists of two steps: first, we model ancestries and simulate
119 genotypes based on ancestry specific frequencies and phenotypes (Fig. 1); then, we test
120 associations between the phenotype and causal variants via a linear model for a quantitative trait,
121 or a logistic regression for a dichotomous trait, and summarize the statistical power. If the
122 population is admixed, global ancestry is supplied as a fixed effect to correct for population
123 structure.

124 125 *1) Ancestry modeling*

126
127 Global ancestry in a 2-way admixed population is modeled as a beta distribution. The i th
128 individual's ancestry θ is characterized as

$$129
130 \theta_i \sim \text{Beta} \left(m^2 \left(\frac{1-m}{v} - \frac{1}{m} \right), \frac{1-m}{m} \alpha \right)$$

131
132 where m and v are the mean and variance of the global ancestry (from a presumptive population
133 1 in this framework) in an admixed population of interest.

134
135 The local ancestries, i.e. the source of ancestry of both the maternal and paternal copies of
136 haplotypes at any genomic position, can be obtained from a binomial sampling with the

137 probability equaling the global ancestry. The process is repeated independently for diploid
138 chromosomes over a presumptive number (n) of LD-independent loci to form a $n \times N$ local
139 ancestry matrix, where N is the proposed sample size.

140 2) *Genotype simulation*

141 We first draw allele frequencies in the two ancestries (Population 1 and 2) from a beta
142 distribution under the Balding–Nichols model (Balding and Nichols, 1995) with a given F_{ST}
143 distribution under the Balding–Nichols model (Balding and Nichols, 1995) with a given F_{ST}
144 distribution under the Balding–Nichols model (Balding and Nichols, 1995) with a given F_{ST}
145 distribution under the Balding–Nichols model (Balding and Nichols, 1995) with a given F_{ST}
146 distribution under the Balding–Nichols model (Balding and Nichols, 1995) with a given F_{ST}

$$147 p_{1s}, p_{2s} \sim \text{Beta} \left(\frac{p_s(1 - F_{ST})}{F_{ST}}, \frac{(1 - p_s)(1 - F_{ST})}{F_{ST}} \right)$$

148 where p_s is the allele frequency at an independent locus s in an ancestral population prior to the
149 divergence and drawn from a uniform distribution $unif(0.001, 0.999)$. Within the model we
150 additionally provide additional distributions if the focus is not on common variants as it is here.
151 We set F_{ST} as a flexible value to increase from a baseline genome-wide F_{ST} when the ancestral
152 allele becomes rare. This is to reflect that a rarer variant in the ancestral population is easier to
153 drift to different frequencies in diverged populations, especially if one population has undergone
154 a severe bottleneck (as would be expected to increase F_{ST}).
155 a severe bottleneck (as would be expected to increase F_{ST}).
156 a severe bottleneck (as would be expected to increase F_{ST}).
157 a severe bottleneck (as would be expected to increase F_{ST}).

$$158 F_{ST} = F_{ST_G} + (1 - MAF_a)\delta$$

159 where F_{ST_G} is the genome-wide background F_{ST} between the two populations of ancestral origin
160 and is considered lower than the F_{ST} between trait-causal loci because of the difference in
161 directional selection. MAF_a is the minor allele frequency of the variant in ancestral populations
162 and δ is the increment with regard to minor allele frequency (MAF) decrease, set as 0.3 in this
163 study. Alternatively, we also test for fixed F_{ST} under the genome-wide background value when
164 exploring the effect on power from various F_{ST} values ranging from 0.1 to 0.9. The genotypes are
165 then drawn from binomial sampling using the allele frequency corresponding to the local
166 ancestry assigned at the locus (i.e., p_{1s} or p_{2s}) across all loci.
167 ancestry assigned at the locus (i.e., p_{1s} or p_{2s}) across all loci.
168 ancestry assigned at the locus (i.e., p_{1s} or p_{2s}) across all loci.

169 3) *Genetic contribution to trait*

170 We randomly assign w out of the total n loci to be causal variants, where w is the proposed
171 polygenicity of the trait. The weights for the causal variants are drawn from a standard normal
172 distribution $N(0,1)$, and the signs of the weights are tied to the prevalence of the allele in the two
173 populations of ancestral origin: the direction of the weight, positive or negative, at a locus is
174 decided by the binary outcome of trial with probability $\frac{p_{1s}}{p_{2s}}$. In this way, a difference in
175 directional selection of the complex trait in the two populations is introduced to facilitate a
176 correlation between the trait and ancestry in the admixed group. Then, polygenic risk scores
177 (PRS) in samples can be calculated based on the weights and the genotypes at causal loci.
178 (PRS) in samples can be calculated based on the weights and the genotypes at causal loci.
179 (PRS) in samples can be calculated based on the weights and the genotypes at causal loci.

180 4) *Non-genetic contribution to trait*

181

182 The non-genetic component is treated as the sum of two parts in admixed populations: (1)
183 random environmental variation modeled as Gaussian noise and (2) environmental confounders
184 correlated with ancestry, such as socioeconomic status and education, modeled as ancestry by
185 environment interaction. Details are described in *Supplementary Notes*.

186

187 5) *Phenotype*

188

189 For quantitative traits, the phenotype is the direct sum of the genetic component (i.e. PRS) and
190 the non-genetic score. For dichotomous traits, the phenotype is converted from the sum of
191 genetic and non-genetic scores to binary case and control status based on the given liability
192 threshold of the case prevalence.

193

194 6) *Association testing*

195

196 Association between the trait and a variant is tested via a linear regression for a quantitative trait,
197 or a logistic regression for a dichotomous trait in all three populations. The global ancestry is
198 corrected in the admixed group. Power is defined as the proportions of causal variants with a
199 significant p value above a given stringency threshold.

200

201 *Estimation of false positives*

202

203 A false positive rate in associations is empirically verified against the association stringency by
204 calculating the proportions of non-causal variants being discovered with significant association p
205 values. This is calculated separately in each population.

206

207 *PRS estimation*

208

209 For training purposes, we obtained the “estimated” weights of causal variants by conducting
210 association analyses in 10,000 individuals in each population. We then used these weights to
211 estimate PRS in another 1,000 individuals in each population as a test. We tested the PRS
212 construction in two ways: firstly, we only used significant ($p < 0.05$) causal-variants from the
213 training set (true positives); secondly, we included all significant variants (all positives) over a
214 range of different stringency ($p = 0.05, 5e-4, 5e-6, 5e-8$, respectively). The true PRS of the
215 individuals in the test set were available through an intermediate step in the simulations (Fig. 1)
216 and were used to test the accuracy of the estimated PRS via correlation coefficients.

217

218 *Calculation of F_{ST} for traits from the GWAS catalog*

219

220 We used the full NHGRI-EBI GWAS catalog “All associations v1.0” (Buniello et al., 2018) to
221 extract variants that are significant genome-wide ($< 5e-8$). We restricted traits to 899 that have
222 more than 10 significantly associated variants that can be found in the 1000 Genome Project
223 Phase 3 (Consortium et al., 2015), and computed Weir and Cockerham’s F_{ST} (Weir and
224 Cockerham, 1984) between 99 Utah Residents (CEPH) with Northern and Western European
225 Ancestry (CEU) and 108 Yorubans in Ibadan, Nigeria (YRI) samples using PLINK v1.9

226 (www.cog-genomics.org/plink/1.9/) (Chang et al., 2015). The genomic background weighted F_{ST}
227 was calculated on common variants (MAF >5%) only.

228

229 **Results**

230

231 *Correlation between a trait and ancestry is common*

232

233 Complex trait studies in groups with heterogeneous ancestries usually require a correction for
234 population structure. The implicit assumption is often a correlation between global ancestry and
235 the trait that is commonly observed *a priori*. The estimated ancestries, or typically ancestry
236 informative principal components, are included as a fixed effect to adjust for phenotypic variance
237 from non-genetic confounders (e.g., social and cultural factors correlated with population
238 structure), and to avoid spurious associations (Price et al., 2006). The correlations between
239 ancestries and complex traits can also be due to changes in genetic architectures among ancestral
240 groups either due to differential selection or deep divergence among populations. This in turn
241 forms one of the basic motivations of multi-ancestry genetic studies, including admixture
242 mapping (loci with ancestry deviating from genome-wide expectation), and cross-population
243 transferability of genetic predictors. Therefore, we provide a power estimate for whether a
244 significant correlation with ancestries can be observed, within a given incidence rate and
245 ancestry distributions (Methods, Fig. S1).

246

247 *The power of genetic discovery in an admixed population is higher than in ancestral populations*

248

249 We primarily focused on a genotype-mediated simulation framework to investigate the GWAS
250 setting in an admixed group. We started by modeling global ancestries, then generating LD-
251 independent genotypes based on population divergence, and subsequently the corresponding
252 phenotypes under an additive model (Methods, Fig. 1). We set up the model in a 2-way admixed
253 group with similar proportions to African Americans, here an average of ~75% West African
254 ancestry (denoted as Population 1 in simulations) and ~25% European ancestry (denoted as
255 Population 2) (Bryc et al., 2015; Baharian et al., 2016). We simulated a complex trait assuming
256 50% narrow sense heritability with 100 causal variants, either as a quantitative or a dichotomous
257 trait with a liability threshold of 5%, in both the admixed population of interest and the
258 homogenous populations of ancestral origin (N=1,000 each). To induce a difference between
259 ancestral phenotypic distributions and a correlation between a trait and global ancestries, we tied
260 the direction of effect sizes to the minor allele frequencies in the two populations of ancestral
261 origin (Methods, Fig. 2). In this study, each independent setting was repeated in 50 runs.

262

263 In the standard setting, we modeled the parameters described above, and F_{ST} across the 100
264 causal variants between Population 1 and 2 as a flexible value with a baseline equaling the
265 genome-wide F_{ST} of 0.2 and an increment associated with the rarity of the ancestral MAF. This is
266 referred to as $F_{ST}=0.2+$ in the text. The aim of this is to mirror the larger stochasticity in the
267 frequency change of an ancestrally-rare variant in diverged populations, especially when one of
268 the derived populations has experienced the severe out-of-Africa bottleneck. We then tested
269 associations between the phenotype and each locus while correcting for global ancestries over
270 various sample sizes ranging from 1,000 to 1,000,000. We found the power to discover a causal
271 variant at a canonical threshold of $p \leq 0.05$ significantly higher in an admixed population than

272 the average in either of the populations of ancestral origin (Wilcoxon $p=1.23e-20$ and $5.79e-10$
273 across the range of sample sizes for quantitative and dichotomous traits described in Fig. 2,
274 respectively).

275
276 In addition to the standard setting where an environment by ancestry effect (Env \times Anc) is
277 modeled as ancestry-weighted Gaussian noise, we explored an alternative where we model Env
278 \times Anc as linearly dependent on the ancestry percentages, which would explain a range of
279 proportions of phenotypic variance from 0% to $1-h^2$ (Fig. S2). The power advantage in admixed
280 populations remains consistent between the default Gaussian Env \times Anc and linear modeling,
281 where the latter was set as up to 10% of non-genetic components (Fig. S3).

282
283 The comparatively high power in admixed populations is more pronounced when the trait
284 distributions have greater distance between Population 1 and 2, or the two populations are more
285 deeply diverged, reflected by the larger F_{ST} at causal variants (Fig. 3). To relate to real-world
286 GWAS, we compared our levels of differentiation to the NHGRI-EBI GWAS catalog (Buniello
287 et al., 2018). Among the 899 traits that have more than 10 genome-wide significant hits found in
288 1000 Genomes Project, the majority ($N=877$) have at least one associated variant beyond the
289 background F_{ST} of 0.155 (Fig. 3), we provide a list of the most-differentiated traits between CEU
290 and YRI in Table S1. In contrast to the response to varying F_{ST} , the statistical power does not
291 obviously change when the narrow sense heritability of the trait differs (Fig. S4). When
292 increasing the overall stringency of the type I error rate up to a conventional genome-wide
293 significance of $5e-8$, the power advantage remains very similar across different thresholds,
294 despite the expected decrease in power value on the absolute scale in all populations (Fig. S5,
295 S6). Thus we picked the canonical threshold of $p \leq 0.05$ for the remaining analyses, as it can
296 represent all stringency levels when this study focuses on the relative power comparison, and this
297 more relaxed cutoff would include a larger number of causal variants for further discussion. The
298 actual false positive rate of associations, as calculated from the 900 non-causal variants from the
299 simulation, remained at approximately 5% in admixed samples across the full F_{ST} and h^2 range,
300 though it was lower in Population 1 and 2 as F_{ST} increases high enough to drive minor allele
301 frequencies towards zero (Fig. S7).

302
303 *Cross-population replication and transferability is asymmetric between the admixed group and*
304 *homogenous groups*

305
306 As GWAS is conventionally focused on common variants, to investigate replication and
307 transferability we then increased the polygenicity of a trait to 1,000 causal variants, and set MAF
308 filtering at 5% for each population's genotypes prior to testing associations. We compare
309 discovery in the major ancestral population (Population 1) relative to the admixed population.
310 The proportion of significant signals that replicate in the reciprocal group is asymmetric between
311 the two populations. Discovery in the admixed samples was more likely to replicate in
312 Population 1 than the other way around, and this trend becomes more exaggerated as trait F_{ST}
313 increases (one-way Wilcoxon $p=3.78e-6$, $<2.2e-16$, and $<2.2e-16$ for $F_{ST} = 0.2, 0.5,$ and $0.8,$
314 respectively; Fig. 4).

315
316 We tested the cross-population transferability of polygenic scores (or polygenic risk scores, PRS)
317 constructed from discovered loci with MAF $> 5\%$ in each population by increasing the sample

318 sizes of the training set to 10,000 per population and separately estimating the GWAS-based
319 PRS in an additional testing set of 1,000 samples in each population. We measured the PRS
320 accuracy as the correlation coefficient r between the estimated values and the true value in each
321 test group across 50 repeats of simulations. Interestingly, the prediction accuracy is also
322 asymmetric between admixed and homogenous samples. When we constructed PRS using only
323 true positive signals at an alpha of 0.05, the accuracy of estimating PRS in Population 1 or 2
324 using weights and loci trained from the admixed population is significantly higher than the other
325 way around. This holds true when using all (both true and false) positive signals at various
326 stringency levels (Fig. 4; Figure S8; Table S2), suggesting another advantage in conducting
327 GWAS in admixed populations.

328

329 Discussion

330

331 Our simulation framework, APRICOT, Admixed Population powerR Inference Computed for
332 phenotypic Traits, demonstrated that GWAS in admixed populations has greater power for
333 discovery than in the homogenous populations of ancestral origin, given the same sample sizes.
334 The difference in power increases when the trait is under more differentiated polygenic selection
335 in the two populations of ancestral origin, reflected by F_{ST} . This is because when a trait is driven
336 by more-differentiated variants, its causal variants are likely to be pushed to more extreme allele
337 frequencies, thus weakening the statistical power of discovery in that population. In contrast, the
338 frequency of the same causal loci in admixed populations likely have become more intermediate
339 due to variation in ancestries, making them much easier to detect. An extreme yet classic
340 example that echoes with the observation would be skin pigmentation, where selection is in the
341 opposite direction in populations at high latitude and those living near the equator. A non-
342 synonymous, skin-lightening mutation at rs1426654 is fixed in European descendants, with a
343 high F_{ST} of 0.985 between CEU and YRI. This mutation would never have been discovered
344 through GWAS if analyses were only conducted in European populations, but it is highly
345 detectable through association analyses in admixed populations (Martin et al., 2017b).

346

347 Additionally, the power advantage in admixed populations may persist even for traits that have
348 not been under such strong differentiation: for almost all the 899 traits we examined from the
349 GWAS catalog, some associated SNPs can have a much larger than background F_{ST} between
350 CEU and YRI, even when the traits themselves on average show limited differentiation (Fig. 3).
351 We note however that the high F_{ST} across these trait-associated variants could partially be
352 attributed to ascertainment bias, where the “tagging SNPs” by design are common in Europeans,
353 making the corresponding genetic component of these traits seemingly more differentiated across
354 populations (Novembre and Barton, 2018). The true causal variants that were tagged by these
355 signals could have moderately attenuated F_{ST} , yet the differences in allele frequency likely
356 remain larger than expected, as previously observed from GWAS on simulated whole genome
357 sequences between Africans and Europeans (Kim et al., 2018). Therefore, attempts to discover
358 variants similar to these “ F_{ST} outlier” signals would benefit from GWAS designed in admixed
359 samples.

360

361 In this study, we provide a mechanistic framework to explore the relationship between power
362 gain in single variant associations and variation in ancestries, mediated by the nature of
363 intermediate allele frequencies in admixed populations. A similar hypothesis of power increase

364 was also explored via simulations in Zhang and Stram (Zhang and Stram, 2014), though the
365 focus of their study was to explore the role of local ancestry in genetic associations; therefore,
366 the assumptions of architecture for comparison between admixed and ancestral populations were
367 simplified, where non-genetic components (such as environmental effect and environment by
368 ancestry interactions) and heritability were not considered in the model, and a constant effect
369 size was assumed for all causal variants. Under this model, Zhang and Stram observed a power
370 increase in admixed groups when compared to stratified analyses in the ancestral populations
371 pooled with a proportion identical to the mean global ancestry percentage. Our simulations
372 extended this framework to dive deeper into more-realistic scenarios across various ranges of
373 environmental effect, trait divergence, and heritability. Moreover, we modeled the ancestry-
374 phenotype association observed in many real-world traits under various different distributions.
375 With a more adjustable genetic architecture in the model, we were able to quantify power
376 advantage in admixed populations within different circumstances in order to investigate practical
377 applications as replication portability and PRS.

378
379 In realistic practice, some confounders and restrictions beyond the model assumptions exist: first,
380 some non-additive genetic components, such as genetic by environment interactions ($G \times E$) and
381 epistasis (Park et al., 2018; Rau et al., 2020), could potentially induce effect size heterogeneity at
382 causal loci with or among populations (Rosenberg et al., 2019), thus obscuring the prediction of
383 power advantage in admixed populations because the power of discovery would be variant-
384 specific and balanced by the gain *vs.* loss from the increase in frequency and change in effect
385 size. However, the increase in power is still expected to be substantial from additive components
386 that are usually considered major in a genetic architecture, with effect sizes highly similar across
387 populations (Wojcik et al., 2019). Additionally, currently the contribution from epistasis or $G \times$
388 E components to most trait variability is estimated to be relatively small (Wang et al., 2019; Dahl
389 et al., 2020; Hivert et al., 2020). For variants with heterogeneous effect sizes per ancestry, other
390 local ancestry-aware regression methods could potentially improve the power of detection in
391 admixed populations (Atkinson et al., 2021). Second, the observations in this study that admixed
392 populations harbor a greater power of discovery in GWAS than the ancestral populations is
393 credited to the existence of ancestry variance, independent of specific demographic history of
394 either the admixed or the ancestral population. It is possible that the demographic details or
395 specific assumptions of the genetic architecture would affect the absolute value of power
396 estimate on a finer scale, which has not been the focus of this study, yet is worth being further
397 explored through forward or coalescent simulations with additional details (including modeling
398 differential linkage disequilibrium patterns) in the future. Third, we focused on a single
399 admixture scenario, albeit one reflecting a realistic scenario. We would anticipate our observed
400 patterns to be exaggerated in populations with even contributions from Populations 1 and 2.
401 Further our framework could be extended to k -way admixed populations, but the interpretation
402 and degree of population-specific interpretation become far more complex to be described here.

403
404 Despite the underrepresentation of admixed groups in large GWAS, recent research has
405 highlighted the importance of conducting genetic research with more diversity. Our work joins
406 burgeoning efforts to quantify the statistical benefits of complex trait studies in diverse
407 populations, especially populations of mixed ancestry. Our work suggests another advantage for
408 conducting genetic studies in admixed populations, which comes from elevated allele
409 frequencies when traits are moderately to highly differentiated. Moreover, discoveries from such

410 studies aid improvement in cross-population PRS, which is critical in clinical prediction in
411 personalized medicine yet presently has suboptimal performance for many biomedical traits in
412 non-European populations (Martin et al., 2019; Rosenberg et al., 2019; Cavazos and Witte,
413 2021). We therefore highlight that insights gained from admixed populations provide improved
414 and appealing generalizable properties compared to homogeneous populations. As the field
415 increasingly moves towards personalized medicine applications we must be mindful of
416 opportunities to incentivize novel studies and analyses in diverse and, particularly as we
417 highlight here, populations of mixed ancestry.

418
419

420 **Acknowledgements:** This research was supported by NIH grants R35GM133531 (to BMH),
421 R01HG010297, U01HG009080/S1 (to NAZ and CRG) and R01HG011345 (to NAZ and CRG).
422 The content is solely the responsibility of the authors and does not necessarily represent the
423 official views of the National Institutes of Health.

424

425 **Conflicts of Interest to declare:** None.

426

427 **Software Availability:**

428

429 *APRICOT*: <https://github.com/menglin44/APRICOT>

430

431

432

433 **Reference**

434

435

436

437 Aschard, H., Gusev, A., Brown, R., and Pasaniuc, B. (2015). Leveraging local ancestry to detect
438 gene-gene interactions in genome-wide data. *Bmc Genet* 16, 124. doi:10.1186/s12863-015-
439 0283-z.

440 Asimit, J. L., Hatzikotoulas, K., McCarthy, M., Morris, A. P., and Zeggini, E. (2016). Trans-
441 ethnic study design approaches for fine-mapping. *Eur J Hum Genet* 24, 1330–1336.
442 doi:10.1038/ejhg.2016.1.

443 Atkinson, E. G., Maihofer, A. X., Kanai, M., Martin, A. R., Karczewski, K. J., Santoro, M. L., et
444 al. (2021). Tractor uses local ancestry to enable the inclusion of admixed individuals in
445 GWAS and to boost power. *Nat Genet* 53, 195–204. doi:10.1038/s41588-020-00766-y.

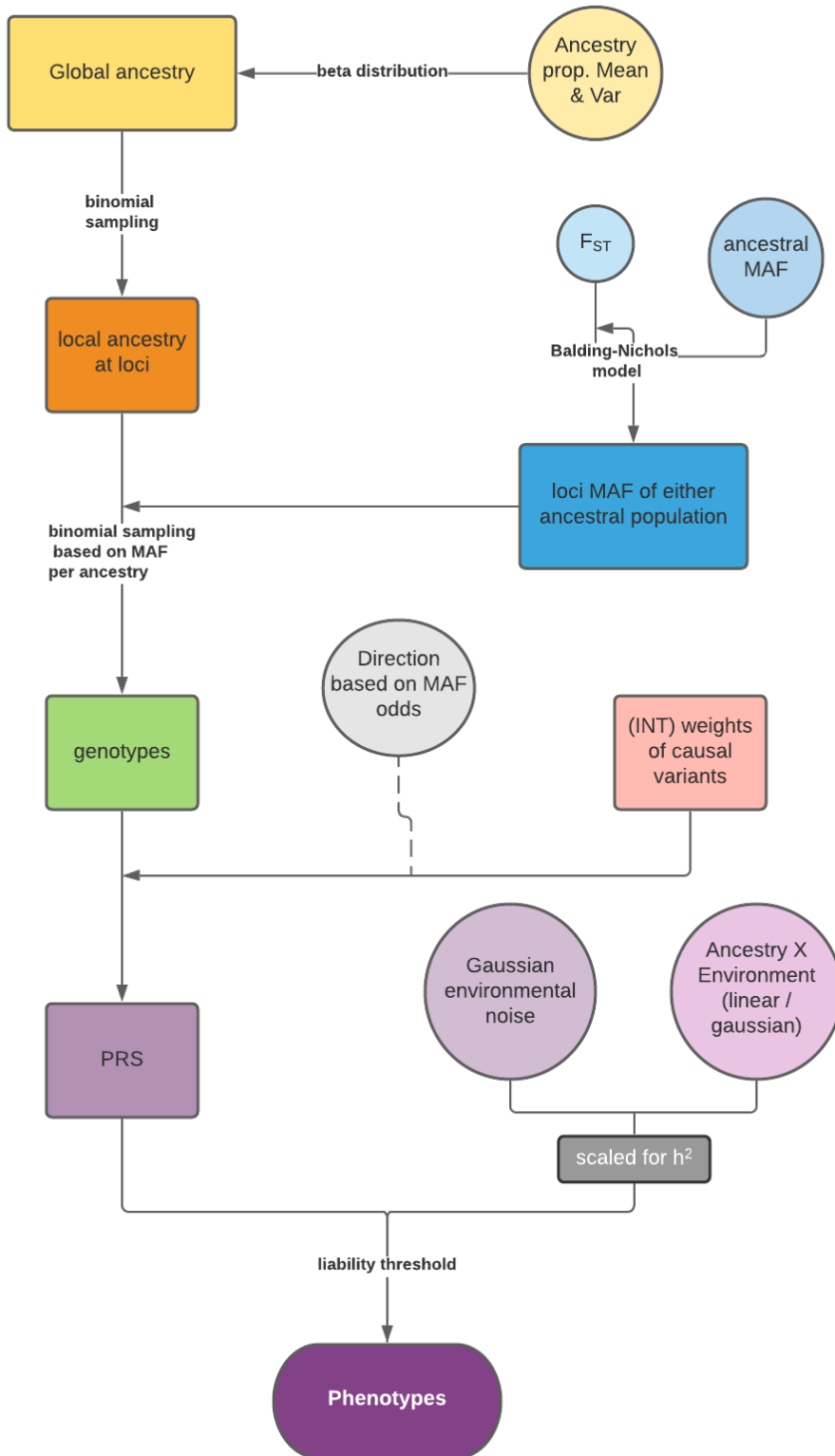
446 Baharian, S., Barakatt, M., Gignoux, C. R., Shringarpure, S., Errington, J., Blot, W. J., et al.
447 (2016). The Great Migration and African-American Genomic Diversity. *Plos Genet* 12,
448 e1006059. doi:10.1371/journal.pgen.1006059.

- 449 Balding, D. J., and Nichols, R. A. (1995). A method for quantifying differentiation between
450 populations at multi-allelic loci and its implications for investigating identity and paternity.
451 *Genetica* 96, 3–12. doi:10.1007/bf01441146.
- 452 Bhardwaj, A., Srivastava, S. K., Khan, M. A., Prajapati, V. K., Singh, S., Carter, J. E., et al.
453 (2017). Racial disparities in prostate cancer a molecular perspective. *Front Biosci* 22, 772–
454 782. doi:10.2741/4515.
- 455 Bryc, K., Durand, E. Y., Macpherson, J. M., Reich, D., and Mountain, J. L. (2015). The Genetic
456 Ancestry of African Americans, Latinos, and European Americans across the United States.
457 *Am J Hum Genetics* 96, 37–53. doi:10.1016/j.ajhg.2014.11.010.
- 458 Buniello, A., MacArthur, J. A. L., Cerezo, M., Harris, L. W., Hayhurst, J., Malangone, C., et al.
459 (2018). The NHGRI-EBI GWAS Catalog of published genome-wide association studies,
460 targeted arrays and summary statistics 2019. *Nucleic Acids Res* 47, D1005–D1012.
461 doi:10.1093/nar/gky1120.
- 462 Bustamante, C. D., Vega, F. M. D. L., and Burchard, E. G. (2011). Genomics for the world.
463 *Nature* 475, 163. doi:10.1038/475163a.
- 464 Cavazos, T. B., and Witte, J. S. (2021). Inclusion of variants discovered from diverse populations
465 improves polygenic risk score transferability. *Hum Genetics Genom Adv* 2, 100017.
466 doi:10.1016/j.xhgg.2020.100017.
- 467 Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., and Lee, J. J. (2015).
468 Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience*
469 4, 1–16. doi:10.1186/s13742-015-0047-8.
- 470 Chimusa, E. R., Zaitlen, N., Daya, M., Möller, M., Helden, P. D. van, Mulder, N. J., et al.
471 (2014). Genome-wide association study of ancestry-specific TB risk in the South African
472 Coloured population. *Hum Mol Genet* 23, 796–809. doi:10.1093/hmg/ddt462.
- 473 Conomos, M. P., Miller, M. B., and Thornton, T. A. (2015). Robust Inference of Population
474 Structure for Ancestry Prediction and Correction of Stratification in the Presence of
475 Relatedness. *Genet Epidemiol* 39, 276–293. doi:10.1002/gepi.21896.
- 476 Conomos, M. P., Reiner, A. P., Weir, B. S., and Thornton, T. A. (2016). Model-free Estimation
477 of Recent Genetic Relatedness. *Am J Hum Genetics* 98, 127–148.
478 doi:10.1016/j.ajhg.2015.11.022.
- 479 Consortium, 1000 Genomes Project, Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P.,
480 Kang, H. M., et al. (2015). A global reference for human genetic variation. *Nature* 526, 68.
481 doi:10.1038/nature15393.
- 482 Conti, D. V., Darst, B. F., Moss, L. C., Saunders, E. J., Sheng, X., Chou, A., et al. (2021). Trans-
483 ancestry genome-wide association meta-analysis of prostate cancer identifies new

- 484 susceptibility loci and informs genetic risk prediction. *Nat Genet* 53, 65–75.
485 doi:10.1038/s41588-020-00748-0.
- 486 Dahl, A., Nguyen, K., Cai, N., Gandal, M. J., Flint, J., and Zaitlen, N. (2020). A Robust Method
487 Uncovers Significant Context-Specific Heritability in Diverse Complex Traits. *Am J Hum*
488 *Genetics* 106, 71–91. doi:10.1016/j.ajhg.2019.11.015.
- 489 Duncan, L., Shen, H., Gelaye, B., Meijsen, J., Ressler, K., Feldman, M., et al. (2019). Analysis
490 of polygenic risk score usage and performance in diverse human populations. *Nat Commun*
491 10, 3328. doi:10.1038/s41467-019-11112-0.
- 492 Hivert, V., Sidorenko, J., Rohart, F., Goddard, M. E., Yang, J., Wray, N. R., et al. (2020).
493 Estimation of non-additive genetic variance in human complex traits from a large sample of
494 unrelated individuals. *Biorxiv*, 2020.11.09.375501. doi:10.1101/2020.11.09.375501.
- 495 Kim, M. S., Patel, K. P., Teng, A. K., Berens, A. J., and Lachance, J. (2018). Genetic disease
496 risks can be misestimated across global populations. *Genome Biol* 19, 179.
497 doi:10.1186/s13059-018-1561-7.
- 498 Lara, M., Akinbami, L., Flores, G., and Morgenstern, H. (2006). Heterogeneity of Childhood
499 Asthma Among Hispanic Children: Puerto Rican Children Bear a Disproportionate Burden.
500 *Pediatrics* 117, 43–53. doi:10.1542/peds.2004-1714.
- 501 Lin, M., Siford, R. L., Martin, A. R., Nakagome, S., Möller, M., Hoal, E. G., et al. (2018). Rapid
502 evolution of a skin-lightening allele in southern African Khoesans. *Proc National Acad Sci*
503 115, 201801948. doi:10.1073/pnas.1801948115.
- 504 Martin, A. R., Gignoux, C. R., Walters, R. K., Wojcik, G. L., Neale, B. M., Gravel, S., et al.
505 (2017a). Human Demographic History Impacts Genetic Risk Prediction across Diverse
506 Populations. *Am J Hum Genetics* 100, 635–649. doi:10.1016/j.ajhg.2017.03.004.
- 507 Martin, A. R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B. M., and Daly, M. J. (2019).
508 Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat Genet* 51,
509 584–591. doi:10.1038/s41588-019-0379-x.
- 510 Martin, A. R., Lin, M., Granka, J. M., Myrick, J. W., Liu, X., Sockell, A., et al. (2017b). An
511 Unexpectedly Complex Architecture for Skin Pigmentation in Africans. *Cell* 171, 1340-
512 1353.e14. doi:10.1016/j.cell.2017.11.015.
- 513 Maskarinec, G., Grandinetti, A., Matsuura, G., Sharma, S., Mau, M., Henderson, B. E., et al.
514 (2009). Diabetes prevalence and body mass index differ by ethnicity: the Multiethnic Cohort.
515 *Ethnic Dis* 19, 49–55.
- 516 Novembre, J., and Barton, N. H. (2018). Tread Lightly Interpreting Polygenic Tests of Selection.
517 *Genetics* 208, 1351–1355. doi:10.1534/genetics.118.300786.

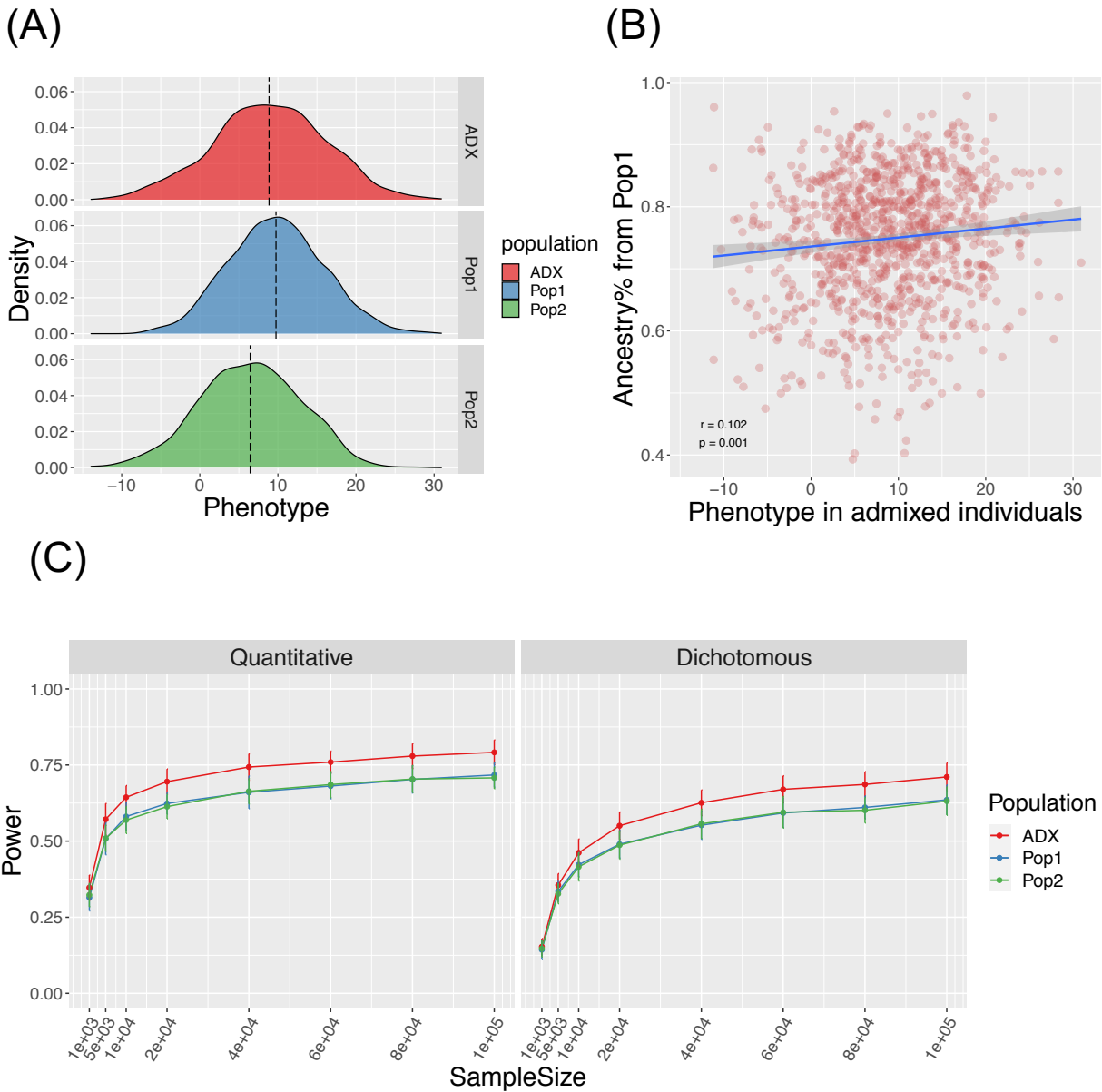
- 518 Park, D. S., Eskin, I., Kang, E. Y., Gamazon, E. R., Eng, C., Gignoux, C. R., et al. (2018). An
519 ancestry-based approach for detecting interactions. *Genet Epidemiol* 42, 49–63.
520 doi:10.1002/gepi.22087.
- 521 Pasaniuc, B., Zaitlen, N., Lettre, G., Chen, G. K., Tandon, A., Kao, W. H. L., et al. (2011).
522 Enhanced Statistical Tests for GWAS in Admixed Populations: Assessment using African
523 Americans from CARE and a Breast Cancer Consortium. *Plos Genet* 7, e1001371.
524 doi:10.1371/journal.pgen.1001371.
- 525 Pino-Yanes, M., Thakur, N., Gignoux, C. R., Galanter, J. M., Roth, L. A., Eng, C., et al. (2015).
526 Genetic ancestry influences asthma susceptibility and lung function among Latinos. *J Allergy*
527 *Clin Immun* 135, 228–235. doi:10.1016/j.jaci.2014.07.053.
- 528 Popejoy, A. B., and Fullerton, S. M. (2016). Genomics is failing on diversity. *Nat News* 538,
529 161. doi:10.1038/538161a.
- 530 Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A., and Reich, D.
531 (2006). Principal components analysis corrects for stratification in genome-wide association
532 studies. *Nat Genet* 38, 904–909. doi:10.1038/ng1847.
- 533 Rau, C. D., Gonzales, N. M., Bloom, J. S., Park, D., Ayroles, J., Palmer, A. A., et al. (2020).
534 Modeling epistasis in mice and yeast using the proportion of two or more distinct genetic
535 backgrounds: Evidence for “polygenic epistasis.” *Plos Genet* 16, e1009165.
536 doi:10.1371/journal.pgen.1009165.
- 537 Rosenberg, N. A., Edge, M. D., Pritchard, J. K., and Feldman, M. W. (2019). Interpreting
538 polygenic scores, polygenic adaptation, and human phenotypic differences. *Evol Medicine*
539 *Public Heal* 2019, 26–34. doi:10.1093/emph/eoy036.
- 540 Rosenberg, N. A., Huang, L., Jewett, E. M., Szpiech, Z. A., Jankovic, I., and Boehnke, M.
541 (2010). Genome-wide association studies in diverse populations. *Nat Rev Genet* 11, 356–366.
542 doi:10.1038/nrg2760.
- 543 Shi, H., Burch, K. S., Johnson, R., Freund, M. K., Kichaev, G., Mancuso, N., et al. (2020).
544 Localizing Components of Shared Transethnic Genetic Architecture of Complex Traits from
545 GWAS Summary Data. *Am J Hum Genetics* 106, 805–817. doi:10.1016/j.ajhg.2020.04.012.
- 546 Shriner, D., Adeyemo, A., and Rotimi, C. N. (2011). Joint Ancestry and Association Testing in
547 Admixed Individuals. *Plos Comput Biol* 7, e1002325. doi:10.1371/journal.pcbi.1002325.
- 548 Thornton, T., Tang, H., Hoffmann, T. J., Ochs-Balcom, H. M., Caan, B. J., and Risch, N. (2012).
549 Estimating Kinship in Admixed Populations. *Am J Hum Genetics* 91, 122–138.
550 doi:10.1016/j.ajhg.2012.05.024.

- 551 Wang, H., Zhang, F., Zeng, J., Wu, Y., Kemper, K. E., Xue, A., et al. (2019). Genotype-by-
552 environment interactions inferred from genetic effects on phenotypic variability in the UK
553 Biobank. *Sci Adv* 5, eaaw3538. doi:10.1126/sciadv.aaw3538.
- 554 Weir, B. S., and Cockerham, C. C. (1984). ESTIMATING F-STATISTICS FOR THE
555 ANALYSIS OF POPULATION STRUCTURE. *Evolution* 38, 1358–1370.
556 doi:10.1111/j.1558-5646.1984.tb05657.x.
- 557 Wojcik, G. L., Graff, M., Nishimura, K. K., Tao, R., Haessler, J., Gignoux, C. R., et al. (2019).
558 Genetic analyses of diverse populations improves discovery for complex traits. *Nature* 570,
559 514–518. doi:10.1038/s41586-019-1310-4.
- 560 Zaitlen, N., Paşaniuc, B., Gur, T., Ziv, E., and Halperin, E. (2010). Leveraging Genetic
561 Variability across Populations for the Identification of Causal Variants. *Am J Hum Genetics*
562 86, 23–33. doi:10.1016/j.ajhg.2009.11.016.
- 563 Zhang, J., and Stram, D. O. (2014). The Role of Local Ancestry Adjustment in Association
564 Studies Using Admixed Populations. *Genet Epidemiol* 38, 502–515. doi:10.1002/gepi.21835.
- 565
566
567
568



569
570
571
572

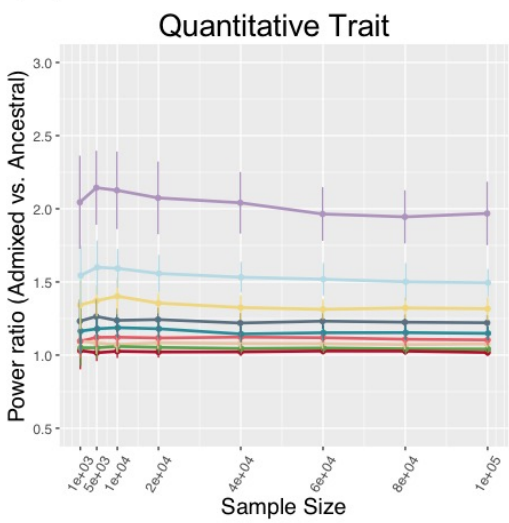
Figure 1. Flow chart of the genotype-mediated simulation framework (prior to association testing).



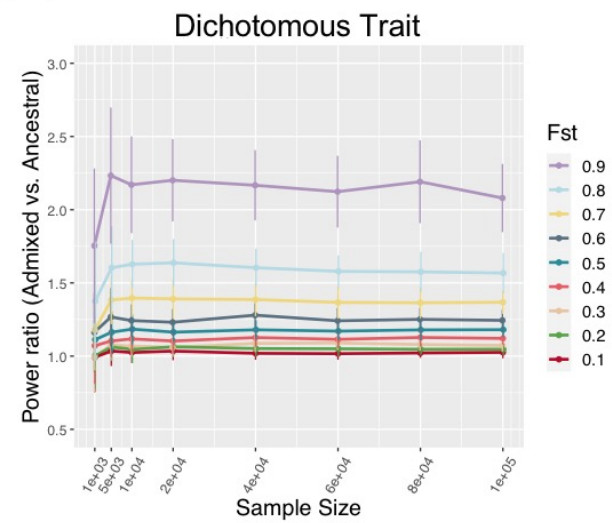
573
574
575
576
577
578
579
580
581
582
583
584

Figure 2. Genotype mediated simulation under an example condition. The simulated trait has 100 causal variants with a narrow sense $h^2=0.5$, F_{ST} at causal variants 0.2+ (explained in Result), and environment by ancestry effect modeled as the sum of ancestral Gaussian environmental noise proportional to global ancestry. (A) Simulated quantitative phenotype distribution of populations of ancestral origin (Pop1, Pop2, blue and green respectively) and admixed population (ADX, red) of 1,000 samples each. (B) Correlation between simulated phenotype in admixed population and the global ancestry from Population 1. (C) Power to discover a causal variant over a range of sample sizes in a quantitative and dichotomous trait. Data point and error bars represent the mean and standard deviation across 50 repeats, respectively.

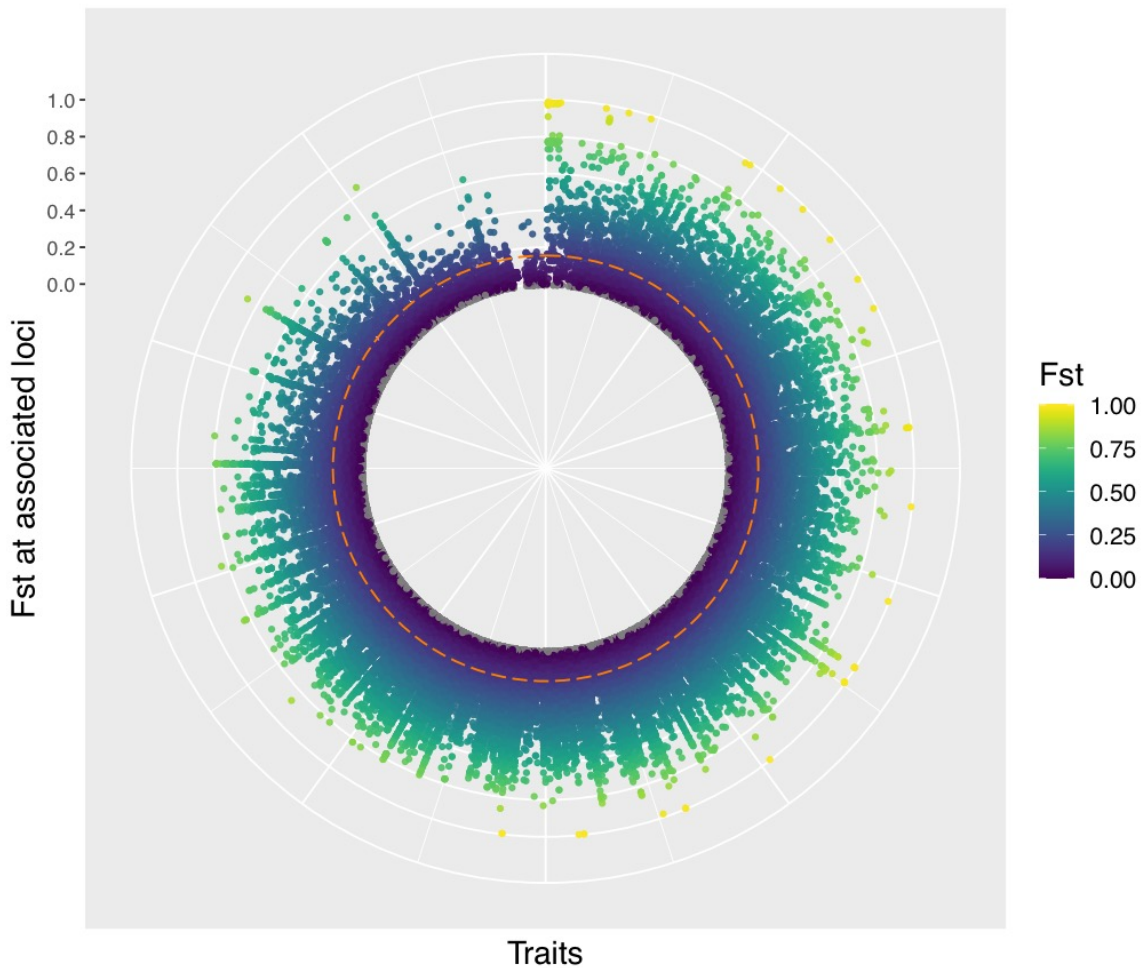
(A)



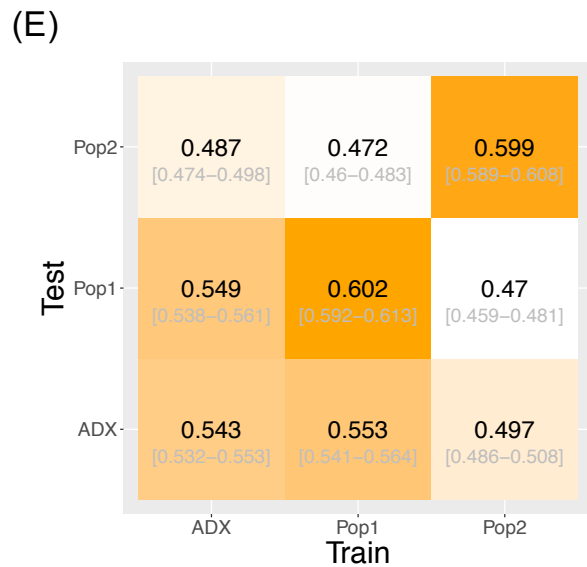
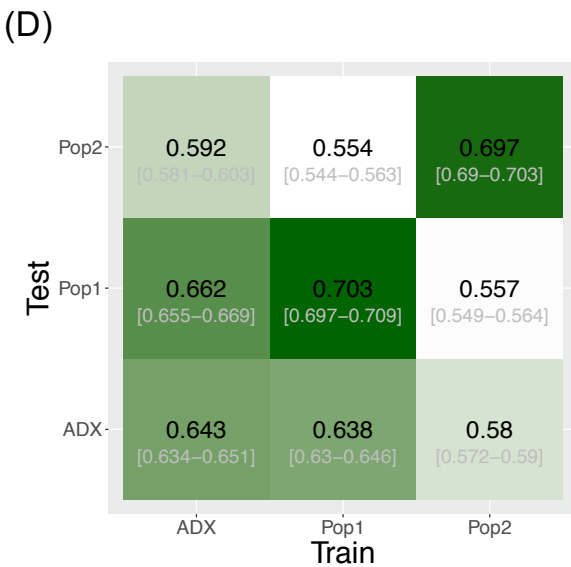
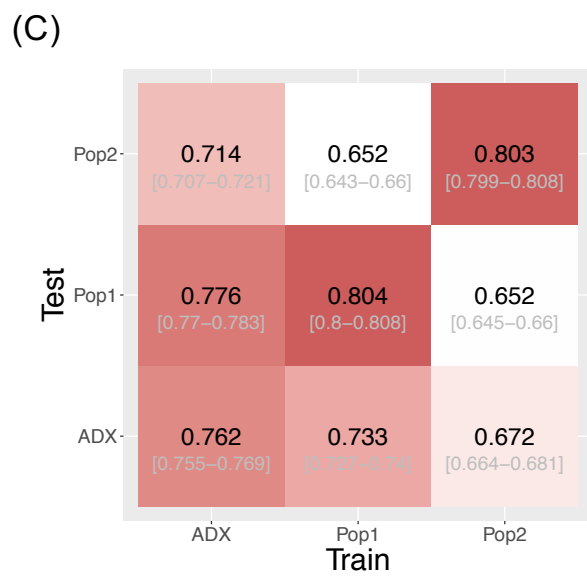
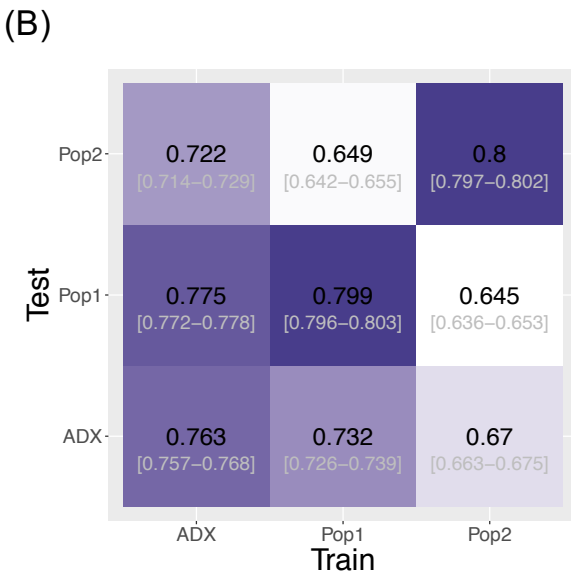
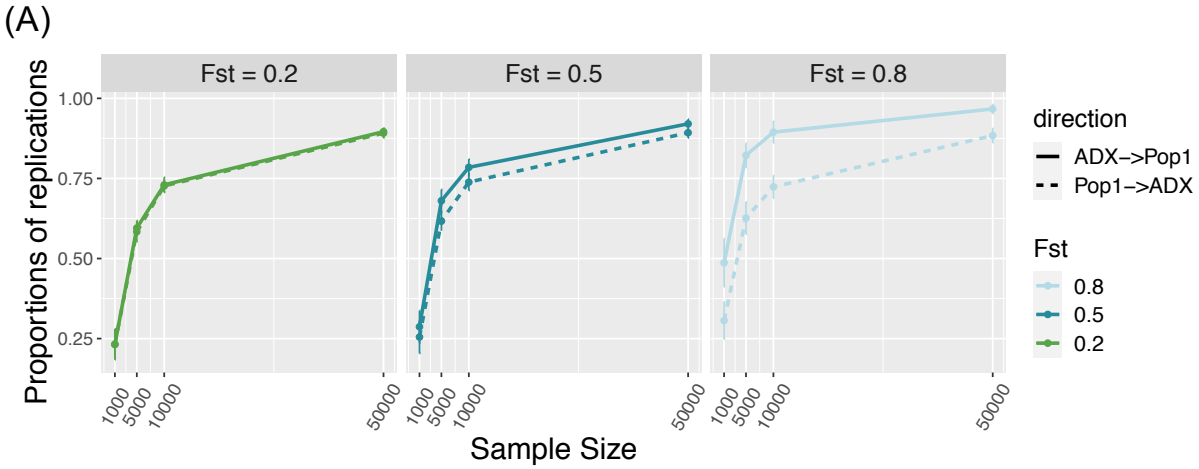
(B)



(C)



586 **Figure 3. Varying F_{ST} at trait associated loci.** Ratio of power in admixed population over the
587 average in the two populations of ancestral origin, with different F_{ST} at causal loci in (A) a
588 quantitative trait and (B) a dichotomous trait. F_{ST} was set to constant during simulations per a
589 specified value. The trait was assumed to have 100 causal loci and a narrow sense heritability of
590 0.5, with environment by ancestry effect modeled as a sum of ancestral Gaussian noise
591 proportional to the global ancestry. Data points and error bars represent the mean and standard
592 deviation across 50 repeats, respectively. (C) F_{ST} at genome-wide significant hits for 899 traits
593 from the GWAS catalog, between CEU and YRI from the 1000 Genomes Project Phase 3. Traits
594 are spread along the radian (x-axis), with variant F_{ST} shown along the radius (y-axis). The dashed
595 line represents the genomic background F_{ST} .
596



597
598

599 **Figure 4. Transferability of GWAS variants across populations.** (A) Replication of
600 individual signals that are common in both ancestral Population 1 and the admixed group.
601 Direction of replication is shown as a solid or dashed line: the former indicates loci are
602 discovered in an admixed population and replicated in Population 1; the latter loci are discovered
603 in Population 1 and replicated in the admixed population. Data point and error bars represent the
604 mean and standard deviation across 50 repetitions. (B), (C), (D), and (E) Heat map of accuracy
605 of PRS using signals above different stringency of significance level at 0.05, $5e-4$, $5e-6$, and $5e-$
606 8 , respectively. The accuracy is measured as the correlation coefficient between the estimated
607 PRS against the true PRS. The training population where the weights and variants were
608 identified, and the test population in which to construct PRS, are specified on the x- and y-axes.
609 Central numbers in black within each cell are the average correlation coefficient across 50
610 independent simulations, with the 95% confidence interval of the mean acquired from
611 bootstrapping ($n=1,000$).
612