

1 **Tractor: A framework allowing for improved inclusion of admixed individuals in large-scale association**
2 **studies.**

3

4 Elizabeth G. Atkinson^{1,2,3,*}, Adam X. Maihofer⁴, Masahiro Kanai^{1,2,3,5,6}, Alicia R. Martin^{1,2,3}, Konrad J.
5 Karczewski^{1,2}, Marcos L. Santoro^{2,7}, Jacob C. Ulirsch^{1,2,3,8}, Yoichiro Kamatani⁹, Yukinori Okada^{6,10,11}, Hilary K.
6 Finucane^{1,2,3}, Karestan C. Koenen^{2,12}, Caroline M. Nievergelt^{4,†}, Mark J. Daly^{1,2,3,13,†}, Benjamin M. Neale^{1,2,†}

7

8 ¹Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston,
9 MA, USA

10 ²Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA

11 ³Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

12 ⁴Department of Psychiatry, University of California San Diego, La Jolla, CA, USA

13 ⁵Program in Bioinformatics and Integrative Genomics, Harvard Medical School, Boston, MA, USA

14 ⁶Department of Statistical Genetics, Osaka University Graduate School of Medicine, Suita, Japan

15 ⁷Departamento de Psiquiatria, Universidade Federal de São Paulo, São Paulo, Brazil

16 ⁸Program in Biological and Biomedical Science, Harvard Medical School, Boston, MA, USA

17 ⁹Laboratory of Complex Trait Genomics, Department of Computational Biology and Medical Sciences,
18 Graduate School of Frontier Sciences, The University of Tokyo, Tokyo, Japan

19 ¹⁰Laboratory of Statistical Immunology, Immunology Frontier Research Center (WPI-IFReC), Osaka University,
20 Suita, Japan

21 ¹¹Integrated Frontier Research for Medical Science Division, Institute for Open and Transdisciplinary Research
22 Initiatives, Osaka University, Suita, Japan

23 ¹²Department of Epidemiology, Harvard TH Chan School of Public Health, Boston, MA, USA

24 ¹³Institute for Molecular Medicine Finland, University of Helsinki, Helsinki, Finland

25

26 [†] These authors contributed equally to this work.

27

28 *Corresponding author contact information:

29 eatkinso@broadinstitute.org

30 (202) 246-0666

31 Analytic and Translational Genetics Unit

32 Richard B. Simches Building, 6th Floor

33 185 Cambridge Street

34 Boston, MA 02114

35

36

37

38 **Abstract**

39 Admixed populations are routinely excluded from medical genomic studies due to concerns over
40 population structure. Here, we present a statistical framework and software package, *Tractor*, to facilitate the
41 inclusion of admixed individuals in association studies by leveraging local ancestry. We test *Tractor* with
42 simulations and empirical data focused on admixed African-European individuals. *Tractor* generates ancestry-
43 specific effect size estimates, can boost GWAS power, and improves the resolution of association signals.
44 Using a local ancestry aware regression model, we replicate known hits for blood lipids in admixed
45 populations, discover novel hits missed by standard GWAS procedures, and localize signals closer to putative
46 causal variants.

47

48

49 Introduction

50 Admixed groups, whose genomes contain more than one ancestral population such as African
51 American and Hispanic/Latino individuals, make up more than a third of the US populace, and the population
52 is becoming increasingly mixed over time¹. Many common, heritable, diseases including prostate cancer²⁻⁵,
53 asthma⁶⁻⁹, and several cardiovascular disorders such as atherosclerosis^{10,11} are enriched in admixed
54 populations of the US. However, only a minute proportion of association studies address the genetic
55 architecture of complex traits in such groups^{12,13}; admixed individuals are systematically removed from many
56 studies due to the lack of methods and pipelines to effectively account for their ancestry such that population
57 substructure can infiltrate analyses and bias results¹⁴⁻²¹. Large-scale efforts to collect genetic data alongside
58 medically-relevant phenotypes are beginning to focus more on non-Eurasian ethnic groups that contain higher
59 amounts of admixture²²⁻²⁷, motivating the timely development of scalable methods to allow well-calibrated
60 statistical genomic work on these populations. If not addressed, this inability to analyze admixed people will
61 limit the clinical utility of large-scale data-collection efforts for minorities, exacerbating the concerning health
62 disparities that already exist²⁸⁻³².

63 In GWAS, the specific concern regarding including admixed participants is obtaining false positive hits due
64 to alleles being at different frequencies across populations. Most studies currently attempt to control for this by
65 using Principle Components (PCs) in a linear or linear mixed model framework. However, PCs capture broader
66 admixture fractions, and individuals' local ancestry makeup may differ between case and control cohorts even
67 if their global fractions are identical. Even including PCs as covariates, then, still leaves open the possibility for
68 false positive associations, as well as absorbing power.

69 Studying diverse populations in gene discovery efforts not only reduces disparities but also benefits
70 genetic analysis for individuals of all ancestries. Perhaps the most notable example of this is in multi-ethnic
71 fine-mapping, which can dramatically reduce the variant credible set by leveraging the differing LD structures
72 observed across populations³³⁻³⁸. This is particularly helpful in populations of African descent, where LD blocks
73 are the shortest and individuals have nearly a million more variants per person than individuals outside of the
74 continent³⁹. We find that with admixed populations we not only can utilize the LD patterns from multiple

75 ancestries, but have further disrupted LD blocks within each one, offering a more refined LD landscape with
76 which to localize GWAS signal.

77 To help ensure that advances in genomic medicine will apply globally, we have developed a scalable
78 framework that allows for the easy incorporation of admixed individuals into psychiatric genomics efforts by
79 using local ancestry inference (LAI). Our framework, distributed as a scalable software package named
80 *Tractor*, generates ancestry dosages at each site from input LAI calls, extracts painted haplotype segments for
81 correction at the genotype level and runs a local ancestry-aware regression model, producing ancestry-specific
82 effect size estimates and p values. Through testing in simulations and with empirical data on phenotypes with
83 differing levels of polygenicity, we demonstrate that *Tractor* produces accurate results in admixed cohorts and
84 boosts GWAS power across many genetic contexts. We further demonstrate improvements in association
85 signal localization from the higher resolution of haplotype breakpoints in admixed genomes. These efforts fill a
86 gap in existing resources and will improve our understanding of complex diseases across diverse populations.

87 The incorporation of local ancestry into variant identification for admixed populations is a concept that has
88 been discussed previously^{40–50}, particularly with regard to ‘admixture mapping,’ whereby researchers associate
89 an elevation of a given ancestry at a locus in the genome with increased risk of a disease that is known to be
90 stratified in prevalence across ancestries^{51–55}. Admixture mapping has proven successful in diseases which are
91 highly stratified across populations, such as asthma and cardiovascular phenotypes^{56–61}. We build upon this
92 important work by modeling the local ancestry haplotype dosage for each person at each variant in a way that
93 allows for the generation of ancestry-specific effect size estimates while allowing for differences in minor allele
94 frequency (MAF) across populations without an increased false positive risk.

95 The statistical model built into *Tractor* for binary phenotypes tests each SNP for an association with the
96 phenotype using the logistic regression model: $\text{logit}(Y) = b_0 + b_1X_1 + b_2X_2 + b_3X_3 \dots + b_kX_k$

97 where X_1 is the number of haplotypes of the index ancestry present at that locus for each individual, X_2 is the
98 number of copies of the risk allele coming from the first ancestry, X_3 is the number of copies coming from the
99 second ancestry, and X_4 to X_k are other covariates such as PCs, age, sex, etc. The significance of the risk

100 allele is evaluated with a likelihood ratio test comparing the full model to a model fit without the risk allele, thus

101 allowing estimation of the aggregated effects in the presence of effect size heterogeneity. To further test if a
102 risk allele is ancestry-specific, we evaluate if the difference between b_2 and b_3 is non-zero using a Z-test. The
103 model presented is for a 2-way admixed scenario but can be scaled to an arbitrary number of ancestries. We
104 have built pipelines to implement this joint model as well as generate genotype files containing extracted
105 ancestry portions. These tools can be implemented in python and Hail either locally or on the cloud.

106

107 **Results**

108 ***LAI has high accuracy for African Americans***

109 We ran LAI using the program RFmix, a discriminative approach which estimates local ancestry by using
110 conditional random fields parameterized with random forests⁶². RFmix can run on multi-way admixture
111 populations, outperforms other local ancestry inference methods for minority populations, and leverages the
112 ancestry components in admixed reference panel individuals, highly important when there is a lack of
113 homogenous reference panels – often the case for understudied groups^{63,64}. As *Tractor* relies heavily on LAI
114 calls, we first ran simulations to ensure that RFmix called local ancestry accurately. LAI was highly accurate in
115 a realistic demographic model for African American (AA) individuals (one pulse of admixture 9 generations ago
116 with 84% contribution of haplotypes from Africa (AFR) and 16% from Europe (EUR); see *Methods*), assigning
117 the correct ancestry ~98% of the time (Table S1). To ensure that *Tractor* performed well across demographic
118 models, we varied demographic parameters including admixture fractions and pulse timings. Specifically, we
119 varied the pulse of admixture in time to 3 generations and 20 generations ago and changed the admixture
120 fractions to 30/70% and 50/50% EUR and AFR ancestry, respectively (**Figure S1**). We also checked the
121 ancestry-specific accuracy in the realistic demographic scenario to assess if there was a bias in calling
122 dependent on ancestry. Across all demographic models and ancestries, site-wise LAI calls were similarly
123 accurate, with the correct call being obtained ~98% of the time (Table S1). While we refer solely to continental
124 level ancestry here, we appreciate the high level of diversity and admixture within the continents and
125 particularly in Africa. As reference panels for diverse groups grow in size, we will have increased ability to
126 examine more geographically refined groups and deconvolve ever more specific haplotypes.

127

128 ***Recovery of long-range haplotypes disrupted by statistical phasing***

129 While errors in statistical phasing can lead to errors in LAI, we found that iterating between LAI and
130 statistical phasing improved the accuracy of both. Errors in statistical phasing are a major concern^{65,66}, but few
131 methods to recover disrupted haplotypes exist. Taking advantage of the unique ability to visualize tracts
132 offered by admixed individuals, we additionally improve long-range haplotype resolution by correcting
133 chromosome strand switch errors from phasing, which we find to be common in admixed cohorts. We
134 demonstrate that using local ancestry information, we are able to consistently correct switch errors and recover
135 disrupted haplotypes, making tract distributions look significantly more realistic (**Figures 1, 2**).

136 To replicate standard analytical procedures employed on cohort data, we statistically phased our truth
137 dataset using SHAPEIT2⁶⁷ software with a balanced reference panel composed of EUR and AFR continental
138 individuals from the 1000 Genomes Project³⁹. We then examined the distribution and lengths of the EUR
139 tracts. Analyzing the less common ancestry tracts allows for more precise quantification of tract counts
140 because it is less likely that recombination will mask their phase switch errors. The probability of obtaining the
141 observed number of tracts after phasing (131 after phasing vs 42 in the truth dataset) given the input
142 demographic model was extremely unlikely, $p=5.0 \times 10^{-26}$. After correcting phase switch errors, the likelihood of
143 the tract distribution improved, albeit still with significantly more switches than in the truth data ($p=2.7 \times 10^{-11}$, 96
144 tracts) – approximately half the excess tracts without phase error correction. After then implementing one
145 additional round of LAI on the corrected genotype files, the number of excess tracts was further reduced
146 ($p=0.009$, 62 tracts). Thus, our procedure for correcting phase switch errors successfully recovers long-range
147 haplotypes and better approximates the true tract length distributions (**Figure 2**).

148 To ensure that phase switch error correction performs well across different population histories, we ran
149 simulations to assess how closely the tract length distributions approximated the truth for a range of
150 demographic models. Specifically, we checked performance varying the timing of admixture pulses (including
151 3, 9, and 20 generations ago) and the admixture proportions in the simulation (70/30, and 50/50 AFR/EUR,
152 respectively). Under all scenarios, our tract recovery procedure improves strand flips in painted karyograms

153 **(Figures 1, S1)** and decreases the significance of the difference between the observed versus true tract length
154 distributions (**Figure 2, Table S2**). We note that running an additional round of LAI after recovering haplotypes
155 produced the most accurate tract length distributions. This is due to the improved ability of the model to
156 recognize ancestry switch points once more complete haplotypes have been recovered, resulting in smoothing
157 over previous short miscalls.

158

159 ***Evaluating the landscape of GWAS power gains from Tractor***

160 We simulated individuals' likelihoods of being cases as a function of AFR admixture fraction, the ancestral
161 haplotype of each copy of the risk allele, and the risk allele dosage (See *Methods, Supplementary*
162 *Information*). This framework is equivalent to modifying the marginal effect sizes due to a tag SNP for a shared
163 causative mutation being monomorphic in EUR but variable in AFR, which is plausible as individuals from
164 Africa contain almost a million more variants than other populations⁶⁸. This also incorporates the clinically
165 observed phenomenon of disease prevalence differing as a function of ancestry. We then ran association tests
166 and compared the power across the odds ratio spectrum under the traditional GWAS model and under our
167 model. Compared to the traditional model, there is a significant gain in power using the *Tractor* framework with
168 similar improvements across sample sizes and disease prevalences (**Figure 3**). Power increases further when
169 there is a difference in MAF across ancestries.

170 We ran similar sets of simulations varying the parameters of the effect size difference, the absolute
171 MAF, MAF difference across ancestries, and admixture fractions (**Figures 3, S2, S3**). The biggest power gain
172 comes if an allelic effect is present in the smaller fraction ancestry only. For example, in a realistic AA
173 demographic model, EUR ancestry makes up only ~20% of the sample. If we model an allele with an effect
174 only active in the EUR background (**Figure 3D**), analyzing the tracts together without LAI information will have
175 essentially no power to detect an association due to the higher noise relative to signal from uninformative
176 tracts. However, *Tractor* is able to recover the effective sample size and power that one would have had if
177 analyzing just the effect haplotypes, i.e. the EUR segments alone.

178 The scenario where *Tractor* is most powerful is when there is heterogeneity in the apparent effect size for
179 the same variant across ancestries. Such heterogeneity in effect sizes may be a consequence of the same
180 variant having different effects in different populations (e.g., in the context of gene-environment interactions) or
181 may arise from differences in the indirect association evidence of the variant (i.e., the contribution to the
182 estimated effect size from tagging other causal genetic variants). Differences in indirect association can come
183 from ancestry-specific variation or from different patterns of linkage disequilibrium. Moreover, under the
184 converse scenarios where there is predicted to be no benefit, the power loss is minimal—we lose a degree of
185 freedom, resulting in a less precise error estimate for each SNP effect (Figure S2,3). In no case does *Tractor*
186 dramatically underperform compared to the traditional GWAS model. See *Supplementary Information* for
187 additional simulations and detail about power results.

188

189 ***No increase in false positive rate with the Tractor model***

190 We quantified the false positive rate of the *Tractor* model by simulating a variant with no effect and counting
191 the spurious significant associations identified in a simulated realistic AA population given $\alpha = 0.05$. Across our
192 tests at various MAFs and between-ancestry MAF differences, we observe no clear difference in false positive
193 rate between *Tractor* and traditional GWAS (**Figure S4**). In addition, we calculated the genomic inflation factor,
194 λ_{GC} , of null phenotypes across GWAS permutations and confirmed no significant inflation using the *Tractor*
195 GWAS model (**Figure S5b**). Therefore, there does not appear to be an elevation in false positive rates with the
196 *Tractor* framework, suggesting that the observed power increases are from improved detection of true
197 biological signal.

198

199 ***Tractor replicates known associations and identifies hits for blood lipids in admixed individuals*** 200 ***missed by standard GWAS***

201 To ensure that our *Tractor* joint-analysis GWAS model also performs well on empirical data, we ran the
202 method for three well characterized blood lipid phenotypes which have been demonstrated to have ancestry-
203 specific effects: Total Cholesterol (TC), high-density lipoprotein cholesterol (HDL-C), and low-density

204 lipoprotein cholesterol (LDL-C). We constructed a pseudo-cohort of 4309 two-way African-European admixed
205 individuals from the UK Biobank (UKBB) with biomarker phenotype data to serve as our sample (**Figure S4**).
206 *Tractor* GWAS replicated previously implicated associations for blood lipids in this cohort^{45,69-71}, reaching the
207 standard genome-wide significance level of 5×10^{-8} at previous top associations, including in genes *PCSK9*,
208 *LDLR*, and *APOE* (**Figure 4**). In some cases, *Tractor* improved the observed top hit significance.

209 Our LAI-incorporating model was also able to identify hits in these admixed UKBB individuals that
210 standard GWAS was not when using the traditional genome-wide significance threshold (**Figure 4**). For
211 example, we identify a hit missed by standard GWAS on the same dataset that is present only on the AFR
212 background on chromosome 1 (rs12740374, $p=3.46 \times 10^{-8}$). This locus has previously been shown to affect
213 blood lipid levels, metabolic syndrome, and coronary heart disease risk in independent AA cohorts^{69,72-77}, and
214 was determined to be the causal variant for affecting LDL-C in a multi-ethnic fine-mapping study⁷⁸. Had we not
215 deconvolved ancestral tracts for our GWAS, we would have missed this site with a demonstrated effect on our
216 phenotype and population of interest.

217 We additionally identify a novel peak on chr15 that only reached significance in the AFR tracts in this
218 UKBB cohort. The lead SNP (rs12594517, $p=1.915 \times 10^{-8}$) lies in an intergenic area and is uncharacterized.
219 The closest gene neighboring it is *MEIS2*, lying ~70kb upstream, followed by *C15orf41*. While the precise role
220 and mechanism this locus plays in affecting blood lipids remains unclear, *MEIS2* has previously been found to
221 be associated to body mass index and waist circumference and *C15orf41* was a significant hit in a previous
222 GWAS of cholesterol^{79,80}. Though further follow-up is needed to clarify any direct relationship to TC, this
223 association highlights the utility of *Tractor* to identify signals that would be undetectable in admixed cohorts
224 without accounting for local ancestry.

225 *Tractor* is also able to refine the location of GWAS signals to closer to the causal variant than is
226 possible using standard GWAS procedures. TC has previously been mapped to the gene *DOCK6* in AA
227 cohorts⁶⁹, a finding we replicate for the suggestive GWAS threshold with standard GWAS on UKBB admixed
228 individuals in the same intronic area as previously found. However, when we run the *Tractor* model, we identify
229 a lead *DOCK6* SNP 20kb downstream in the AFR samples, as well as in a meta-analysis of hits from

230 deconvolved AFR and EUR tracts. This new lead SNP (rs2278426) is a missense mutation spanning both
231 *DOCK6* as well as *ANGPTL8*, a gene which may play a key role in blood lipid regulation (see *Discussion*,
232 **Figure 5**). To assess whether our improved ability to localize to this variant was due to a true effect size
233 differences between the EUR and AFR or to a marginal effect size difference driven by MAF or LD differences
234 across the ancestries, we further attempted validating its association by fine-mapping TC in other large-scale
235 populations: 345,235 white British individuals from UKBB and 135,808 Japanese individuals from Biobank
236 Japan (⁸¹, Ulirsch, JC., Kanai, M. et al., in prep., Kanai, M. et al., in prep). rs2278426 was successfully fine-
237 mapped to a 95% credible set in both populations, with maximum posterior inclusion probability of 0.993 in
238 Biobank Japan. This variant is at 26% frequency in the gnomAD ⁸² East Asian ancestry individuals, 18% in
239 African, and 4% frequency in the non-Finish Europeans. These frequency patterns suggest higher power in
240 non-European population to localize a causal variant compared to Europeans. Though below the traditional
241 genome-wide significance level in our sample of ~4300 individuals, this locus highlights the improved ability to
242 localize GWAS signal thanks to leveraging additional breakpoints in admixed genomes.

243

244 **Discussion**

245 Despite the recent advances in understanding the genetics of complex diseases, major limitations
246 remain in our knowledge of the architecture of such disorders in minority and admixed populations. Here, we
247 present an analytical framework and statistical gene discovery method distributed as a scalable software
248 package named *Tractor*, which allows admixed samples to be appropriately included alongside homogenous
249 ones in a well calibrated manner in statistical genomics efforts. We test our framework in a simulation model
250 designed to emulate real AA cohorts. We also apply it to empirical data from admixed African-descent
251 individuals of the UKBB. We observe a gain in power to detect risk loci across sample sizes, demographic
252 models, and disease prevalences using the *Tractor* framework, particularly when effect sizes are
253 heterogeneous across populations. Our approach incorporates a local-ancestry aware GWAS method that can
254 extend the traditional GWAS model. *Tractor* generates ancestry specific *p* values and effect size estimates,
255 which admixture mapping cannot, and which can be extremely helpful in post-GWAS efforts such as

256 constructing genetic risk scores for understudied populations. We demonstrate that our framework also gives
257 increased precision in localizing GWAS signal by leveraging the disrupted LD blocks visible with ancestral
258 chromosome painting in recently admixed groups. This reduces the credible set of SNPs and aids in the
259 prioritization of variants for subsequent functional testing.

260 The *Tractor* pipeline requires several inputs, most importantly accurate local ancestry calls. Users
261 should ensure good LAI performance in their target cohorts. A major determinant of accurate LAI calls is a
262 comprehensive and well-matched reference panel^{49,83}. Of relevance is that reference panels are more plentiful
263 for Eurasian populations than for other groups, underscoring the need to expand sequencing efforts in more
264 global populations. To ensure LAI was unbiased across regions of the genome in our GWAS, we examined the
265 distribution of local ancestry across the genome. Local ancestry inference appears relatively evenly distributed
266 across the genome and proportional to global admixture proportions (**Figures 1, S6**). However, we caution that
267 calls around centromeres and at the ends of chromosomes are most likely to include error, as these genomic
268 regions do not have anchor points on one edge. We similarly recommend that LAI ideally to be conducted on
269 whole genome sequencing data to avoid the introduction of biases. We also note that we have thoroughly
270 tested *Tractor* here in the two-way admixture model that reflects AA demographic history. The analytic
271 infrastructure, however, can currently also run on a three-way admixed model and our statistical model can
272 scale to an arbitrary number of ancestries. Future work will test power and optimize the code in a variety of
273 multi-way admixed demographic scenarios. A final consideration is to ensure consistent phenotyping across
274 ancestry groups, as is standard in multi-ethnic GWAS.

275 We thoroughly evaluated the landscape of when *Tractor* does and does not add power to association
276 studies in simulated data modeled after AA cohorts of the PGC-PTSD (**Figures 3, S2, S3**). In situations where
277 there are differences across ancestries, *Tractor* recovers the power that would be lost from analyzing
278 populations together. In particular, power gains are greatest when there is an effect size difference at a locus
279 between ancestries coupled with differing MAF. For example, our simulated case of a 20% MAF difference for
280 an allele with an effect only in the AFR genetic background would allow for identification of risk variants with an
281 odds ratio ~0.1 smaller than would be possible with traditional GWAS. This allows for detection of additional

282 loci that would have been undetected without modeling local ancestry. *Tractor* can also boost power when
283 there is an effect only on one haplotype background or an allele only present in one ancestry, again most
284 dramatically when that ancestry is less frequent in the dataset. In such an instance in a standard GWAS
285 setting, the signal would be dramatically reduced due to noise from the uninformative majority haplotypes,
286 resulting in extremely low power to detect the locus (**Figure 3D**). Power can be recovered, however, by
287 deconvolving local ancestry and analyzing genotypes on ancestry-specific haplotypes, thus controlling for
288 population structure as well as identifying risk variants that would otherwise be undetectable.

289 Conversely, we find that it is not generally necessary to include local ancestry in a GWAS model when
290 there is no effect size difference between groups. Notably, we are referring here to detection of marginal effect
291 sizes as well as true effect sizes. There is evidence suggesting that in most cases (with some notable
292 exceptions^{8,69,84}), the true effect sizes of causative variants are likely to be equal across ancestries^{33,38,85–90}.
293 However, the marginal effect size of a tag SNP might routinely be different across ancestries due to
294 differences in ancestral MAF and the LD patterns resulting from each ancestral population's demographic
295 history^{91,92}. We underscore that power gains appear to be from true biological signal rather than false positives,
296 as we quantified the *Tractor* false positive rate to be no higher than standard GWAS (Figures S4, S5).

297 We would like to highlight that *Tractor* benefits from power gains to detect the marginal beta (pertaining to,
298 for example, an allele which is only present in one population and tags a nearby causal variant), in addition to
299 the rarer case of variants with true effects only on one haplotypic backbone. Our framework therefore will
300 improve power at substantially more locations across the genome than only at sites which have ancestry-
301 specific causal effect differences. *Tractor* also benefits from increased power in cases where functionally
302 important (and likely rare) alleles only present in one population are missed by genotyping or imputation. In
303 such situations the common alleles in LD, despite being shared across populations, would be associated to the
304 phenotype as a function of which haplotypic background they are found on, and thus would have a haplotype-
305 specific effect. Another relevant scenario to consider would be the presence of LD in regions where there are
306 ancestry-specific markers intermingled with shared ones. This would affect the univariate scan results such

307 that considering the haplotypic background on which alleles fall would particularly aid in localizing signal
308 through improved marginal beta estimates, even when the causal effect is the same in both ancestries.

309 *Tractor* was also able to replicate established GWAS hits, discover new ones, and aid in the
310 localization of GWAS signal with empirical data. We replicated known hits for well-characterized blood lipid
311 phenotypes when testing *Tractor* on a dataset consisting of ~4300 2-way admixed African-European
312 individuals from the UKBB (**Figure 4**). We further demonstrate an improved ability to localize GWAS signal to
313 putative causal SNPs previously identified in another diverse collection, Biobank Japan⁸¹ (**Figure 5**).

314 Specifically, previous analyses of blood lipid phenotypes in an admixed AA cohort had pinpointed the TC top
315 hit to lie within the gene *DOCK6*⁶⁹, nearby the lead SNP for this region in standard GWAS on the admixed
316 UKBB individuals (rs4804576). The *Tractor* AFR ancestry, as well as results from a meta-analysis of summary
317 statistics from AFR and EUR deconvolved genotype files, identified a different top association ~20kb
318 downstream (rs2278426) which additionally lies over the *ANGPTL8* gene on the positive strand while spanning
319 an intronic area of *DOCK6* on the minus strand. *ANGPTL8*, also known as lipasin and betatrophin, has been
320 shown to regulate plasma lipid levels in mice by inhibiting the enzyme lipoprotein lipase⁹³⁻⁹⁶. In humans,
321 *ANGPTL8* levels correlate with metabolic phenotypes including type 2 diabetes and obesity⁹⁷⁻¹⁰⁰ and HDL-C
322 expression levels across diverse ancestry groups have been demonstrated to better correlate with than
323 *DOCK6*¹⁰¹. Together these make *ANGPTL8* a more promising candidate gene than *DOCK6*, which has no
324 clear tie to blood lipid phenotypes. Intriguingly, the *Tractor* lead SNP, rs2278426, is a missense mutation in
325 *ANGPTL8* (p.Arg59Trp) that is predicted to be possibly damaging and deleterious by PolyPhen and SIFT,
326 respectively^{102,103}. This variant is at 18% frequency in gnomAD⁸² AFR individuals but at 4% frequency in the
327 non-Finnish Europeans. These frequency differences highlight how leveraging different diverse populations
328 allows for the improved identification of risk variants as well as how employing multi-ethnic mapping methods
329 aids in the resolution of association signals.

330 The *Tractor* infrastructure released here may be helpful in multiple statistical genetics use cases
331 beyond GWAS. For example, correcting for population structure should be a key early step in evolutionary
332 genomic studies on admixed populations running analyses such as genome-wide scans of selection to avoid

333 bias in selection statistics^{104–106}. Within medical genetics, accounting for the ancestral background on which an
334 allele appears will be key in admixed populations, particularly in studies of rarer variants which are more
335 population specific^{107,108}. For example, because AFR and EUR haplotypes have different rates of background
336 variation⁶⁸, controlling for the local ancestral background may help pinpoint the differences between cases and
337 controls in burden testing that would previously have been overwhelmed by uninformative markers.

338 In sum, *Tractor* allows users to account for genotype-level ancestry in a precise manner, allowing for
339 the well-calibrated inclusion of admixed individuals in large-scale gene discovery efforts. This approach
340 provides a number of benefits over traditional GWAS, including the production of ancestry-specific effect size
341 estimates a p values, improved localization of GWAS signal, and power boosts in genetic scenarios such as
342 when there are effect size or MAF differences across ancestries. This infrastructure is designed as a series of
343 steps to be flexible and easily ported into other statistical genomics activities. We freely provide *Tractor* code in
344 python and Hail¹⁰⁹, a scalable cloud-compatible framework, as well as examples of implementation in a Jupyter
345 notebook¹¹⁰. *Tractor* advances the existing methodologies for studying the genetics of complex disorders in
346 admixed populations.

347

348 **Online Methods**

349 ***QC and LAI Pipeline***

350 The core feature of the *Tractor* framework relies on accounting for fine-scale population structure as
351 informed by local ancestry (i.e. ancestral chromosome painting). *Tractor* then uses this information to (i)
352 correct for individuals' ancestral dosage at all variant sites, (ii) recover long-range tracts in admixed individuals;
353 and (iii) extract the tracts and ancestry dosage counts from each ancestry component for use in ancestry-
354 specific association tests. We have tested and built this framework around LAI calls from RFmix versions 1 and
355 2⁶², and have built an automated pipeline (<https://github.com/eatkinson/Post-QC>) to perform all necessary
356 post-genotyping QC, data harmonization, phasing, and LA inference to consistently prepare the data for
357 downstream analysis. The main code is in bash, subscripts are written in python (See *Supplementary*
358 *Information* for additional details).

359 In all tests shown here, we ran RFmix_v2 with 1 EM iteration and a window size of 0.2 cM. We used the
360 HapMap b37 recombination map¹¹¹ to inform switch locations. The -n 5 flag (terminal node size for random
361 forest trees) was included to account for an unequal number of reference individuals per reference population.
362 We additionally used the --reanalyze-reference flag, which recalculates admixture in the reference samples
363 themselves for improved ability to distinguish ancestries. This is especially important when the reference
364 samples are themselves admixed. As a reference panel for our 2-way admixed simulated African-European
365 cohorts, we used relevant populations of the 1000G reference panel given *a priori* knowledge of AA's
366 demographic history¹¹²⁻¹¹⁴ consisting of 108 YRI and 99 CEU. Painted karyogram plots were produced using a
367 modified version of publicly available code (https://github.com/armartin/ancestry_pipeline). We have optimized
368 this pipeline under the two-way admixed AA demographic scenario. *Tractor* additionally supports 3-way
369 admixture calls with an expanded set of scripts (also at <https://github.com/eatkinson/Tractor>). In all cases we
370 recommend conducting tests of LAI accuracy to ensure reliability, as accurate LAI calls are required for good
371 performance.

372

373 **LAI Accuracy**

374 We validated that LAI was performing well in the AA use case. To do this, we generated a truth dataset by
375 simulating individuals with known phase and LA from empirical data. Our simulation reference panel consisted
376 of haplotypes from homogenous PGC-PTSD individuals who had $\geq 95\%$ EUR or AFR ancestry as inferred by
377 SNPweights¹¹⁵. We simulated admixture between these reference individuals with admix-simu¹¹⁶ using a
378 realistic demographic scenario for the AA population^{113,114} of 1 pulse of admixture 9 generations ago with 84%
379 contribution from Africa and 16% from Europe. The resultant population mixes amongst itself until the present
380 day, copying haplotypes from the previous generation with break points informed by the HapMap combined
381 recombination map¹¹¹. This retains the LD structure and genetic variation present in real genomic data and
382 ensures that the truth dataset resembles cohort data as closely as possible. We then called LA with the 1000
383 Genomes⁶⁸ AFR and EUR superpopulations as our reference panel, and calculated LAI accuracy as how often
384 the ancestry call was correct in the simulated truth data.

385

386 ***Correcting Switch Errors from Statistical Phasing using Local Ancestry***

387 Despite LAI calling ancestry dosage accurately, frequent chromosomal switches were visible in painted
388 karyograms (**Figure 1**), which we determined were due to phasing errors. It is important to retain complete
389 tracts, as spurious breakpoints will reduce the accuracy haplotype-based test. *Tractor* detects and fixes phase
390 switches using the most likely ancestry assignment of subpopulations as determined by a conditional random
391 field from RFmix. We define phase switches as a swap of ancestry across a chromosome within a 1 cM
392 window at a region with heterozygous ancestry dosage.

393 To ensure that correcting phase switch errors improved results compared to the truth expectations for
394 the input demographic scenario, we modeled the expected distributions of EUR tract lengths within AA
395 individuals using a Poisson process with rate=9, the number of generations ago when the pulse of admixture
396 occurred (**Figure 2**). The waiting time until a recombination event disrupts a tract is expected to follow this
397 distribution, with a slight shortening of tracts proportional to the percent admixture due to the inability to
398 visualize tract switches that occur across regions of the same ancestry. The overall proportion of the genome
399 in the realistic scenarios was within range of expectations given the simulation model of 16% European, 84%
400 African ancestry (15.1% and 84.9%, respectively).

401

402 ***GWAS power simulations incorporating local ancestry***

403 We assessed the improvements in GWAS power from using *Tractor* through simulations. We
404 formulated our simulation framework on the suggestions of Skotte et al (2019)¹⁷. Power calculations were
405 based on a simulation framework that initially models an AA population assuming a bi-allelic disease risk allele
406 with a 20% overall MAF and an additive effect in the AFR genetic background but not in the EUR. Specifically,
407 the overall admixture proportions were drawn from a beta distribution with shape parameters 7.76 and 2.17,
408 the fitted parameters to this distribution for AFR ancestry proportions observed in the PGC-PTSD Freeze 2 AA
409 cohorts. The genotype of each copy of the allele was drawn from a binomial distribution with the probability of
410 having the minor allele set to the MAF. We simulated a disease phenotype with individuals' risk drawn from a

411 binomial distribution assuming a 10% disease prevalence. Risk of developing the phenotype was modified on
412 a log-additive scale according to the admixture proportions and the presence of the minor allele on an AFR
413 background using a logit model. In this model, the probability of disease was set to $-2.19 + \log$ of allelic risk
414 effect size*number of copies of the minor allele coming an AFR ancestral background + $0.5*AFR$ admixture
415 proportion. -2.19 was chosen as it represents a 10% probability of disease given no AFR admixture or copies
416 of the minor allele from either ancestral background. The 0.5 value in $0.5*AFR$ Admixture was set in order to
417 induce stratification in the simulated population, as is observed in empirical data. In other words, all of our
418 simulations modeled increasing disease prevalence with admixture fractions, reflective of clinical observation.
419 With this simulation design, individuals with higher AFR ancestry proportions are more likely to be cases
420 whereas those with higher EUR ancestry proportions are more likely to be controls. Subjects' disease status
421 was then drawn from a binomial distribution with the probability parameterized to their individual disease risk
422 according to the logit model. Cases and controls were sampled at random from the simulated population at a
423 2.5:1 control to case ratio, the approximate ratio of controls to cases in PGC-PTSD freeze 2.

424 Under each simulation, we fit three logistic regression models of disease status that included: M1)
425 admixture only, M2) number of copies of the risk allele only, and M3) admixture + number of copies of the risk
426 allele on a EUR background + number of copies on an AFR background. M1 serves as a null comparison to
427 evaluate the significance of including the SNP as a predictor. The significance of M2 and M3 are evaluated by
428 likelihood ratio tests comparing them to M1. For each 100 simulations at a given effect size and sample size,
429 for both M2 and M3 we estimated power as the proportion of the time that the likelihood ratio test was
430 significant ($p < 5e-8$). We performed 100 rounds of simulation with this model at each level of allelic effect size
431 ranging from Odds Ratio (OR) 1.05 to 1.3 and case sample size $N=4000$ and 12000 .

432

433 ***Characterizing the landscape of Tractor power across genomic and disease contexts***

434 To evaluate *Tractor* power gains, we ran similar sets of simulations varying effect size differences across
435 ancestries, MAF differences, admixture fractions, and disease prevalence (**Figures 3, S2, S3**).

436 *Varying effect size across populations:* To examine the effect of modifying the effect sizes, we introduced
437 an effect on EUR haplotypes as well, rather than just on AFR. All these simulations assumed 80% admixture,
438 10% disease prevalence, and 20% MAF in both groups. We modeled cases across the OR spectrum where
439 there was an effect of equal size in both ancestries, a 30% larger effect size on the EUR background, an effect
440 size 30% larger on the AFR haplotype, and an effect size only in the EUR.

441 *Varying absolute MAF:* We next fixed all other parameters and modified the absolute MAF of the
442 simulated risk allele, with the relative difference in MAF between ancestries remaining constant. We changed
443 our MAF from 20% to 10% and 40% under both the models of an effect only in the AFR background and with
444 matching effect sizes between EUR and AFR.

445 *MAF differences between groups:* To see if having a difference in the MAF between the two ancestral
446 groups affected GWAS power, we varied the MAF in the EUR background to be 10, 20, and 30% while
447 keeping the AFR MAF set to 20%.

448 *False positive rate:* We quantified the false positive rate by simulating a variant with no effect and
449 counting significant associations identified in a simulated realistic AA population given $\alpha = 0.05$.

450

451 ***Selection of two-way admixed African-European empirical individuals***

452 To select individuals with 2-way admixture with European and West African ancestry, we took a two-
453 pronged approach. First, we combined genetic reference data from the 1000 Genomes Project³⁹ and Human
454 Genome Diversity Panel¹¹⁸, then harmonized meta-data according to consistent continental ancestries. We
455 then ran PCA on unrelated individuals from the reference dataset. To partition individuals in the UKBB based
456 on their continental ancestry, we used the PC loadings from the reference dataset to project UK Biobank
457 individuals into the same PC space. We trained a random forest classifier given continental ancestry meta-data
458 based on the top 6 PCs from the reference training data. We applied this random forest to the projected UK
459 Biobank PCA data and assigned AFR ancestries if the random forest probability was >50%, otherwise
460 individuals were dropped from further analysis.

461 For those individuals classified by their genetic data to have AFR ancestry, we then combined the 1000

462 Genomes and Human Genome Diversity Panel reference data with genetic data from the African Genome
463 Variation Project as well as these UKBB individuals. To restrict to only two-way admixed West African-
464 European ancestry individuals, we restricted to individuals with at least 12.5% European ancestry, at least 10%
465 African ancestry, and who did not deviate more than 1 standard deviation from the AFR-EUR cline (**Figure**
466 S6A, B). This resulted in approximately 4300 individuals per blood lipid trait. Global ancestry fraction estimates
467 were obtained from running ADMIXTURE¹¹⁹ with $k=2$ (which was the best fit k value to this dataset based on 5-
468 fold cross-validation) on these individuals with 1000 Genomes Project³⁹ EUR and AFR superpopulation
469 individuals as reference data (**Figure S6C**). To ensure there were no major areas of the genome where local
470 ancestry inference was skewing significantly from the expected global fractions, we also assessed the
471 cumulative local ancestry calls across the genome for the UKBB admixed subset (**Figure S6D**).

472

473 **Software implementations**

474 We developed separate scripts to deconvolve ancestry tracts and calculate haplotype dosages, correct
475 phase switch errors, and run a *Tractor* GWAS to obtain ancestry-specific effect size estimates and p values.
476 Pre-GWAS steps are available as independent python scripts. We separated steps to allow for maximum
477 flexibility when using *Tractor*. To implement the joint modeling GWAS approach with the novel linear
478 regression model described here, we have built a scalable pipeline in Hail¹⁰⁹ which can be implemented locally
479 or on the Google Cloud Platform¹²⁰. Descriptions of the steps and an example Jupyter notebook¹¹⁰
480 demonstrating analytical steps and visualization of results of the *Tractor* joint-analysis GWAS are freely
481 available on github (<https://github.com/eatkinson/Tractor>).

482 An alternative pipeline designed for use across environments where Hail may not be as readily
483 implemented involves running the separate/meta-analysis GWAS version of *Tractor*. This pipeline requires the
484 initial processing steps to optionally correct phase switch errors and deconvolve ancestry tracts into their own
485 VCF files. Next, GWAS can be run for the deconvolved files containing different ancestral components with the
486 user's preferred GWAS software, such as plink¹²¹. In this implementation, a standard GWAS model can be run
487 on each ancestral component separately using the ancestry-specific VCF output by *Tractor*, which contains

488 fully or partially missing data including only haplotypes from the ancestry in question. Results from the different
489 ancestry runs could then be meta-analyzed to increase power by incorporating summary statistics from both
490 populations, though we recommend preferentially using the joint-analysis method described in this manuscript
491 to avoid any potential bias from combining multiple ancestral portions of the genome of the same individuals.
492 This implementation is also compatible in large-scale collections where there are large numbers of
493 homogenous individuals, for example many Europeans, but too limited a number of admixed individuals to be
494 run in a GWAS alone. The EUR sections of the admixed cohorts could be analyzed alongside the
495 homogenous European cohorts, making better use of the admixed samples even if other ancestry portions are
496 not utilized, and increasing the effective sample size.

497

498 ***Empirical test of Tractor on blood lipid phenotypes in European-African admixed UKBB individuals***

499 To ensure that *Tractor* replicated well-established associations, we ran standard GWAS, the *Tractor*
500 joint-analysis model, and a meta-analysis of summary statistics from EUR and AFR deconvolved tracts on
501 ~4300 admixed African-European individuals from the UKBB on the biomarker blood lipid traits of Total
502 Cholesterol (TC), high-density lipoprotein cholesterol (HDL), and low-density lipoprotein cholesterol (LDL).
503 We included covariates capturing global ancestry, age, sex, and blood dilution factor in all runs. We assessed
504 meta-analysis performance using different metrics to capture global ancestry, namely the first 20 PCs versus
505 the AFR fraction as determined by ADMIXTURE, which did not have substantive differences. In the joint-
506 analysis framework, we used the measure of global AFR ancestry fraction to more directly capture global
507 ancestry and avoid any potential collinearity with local ancestry from PCs. We generated QQ plots alongside
508 each trait and compared the inflation of test statistics in each GWAS case by looking at the genomic inflation
509 factor, λ_{GC} . We then compared results to those obtained from the same individuals using a standard GWAS
510 approach. No individuals overlap between the previous study of interest (Natarajan et al. 2018, which includes
511 diverse TOPMed²⁴ individuals) and the UKBB individuals included here. As expected, near-identical results
512 were obtained from the meta- and joint-approaches. Gene visualizations were produced with LocusZoom¹²²,
513 Manhattan and QQ plots with bokeh¹²³.

514

515 **Statistical fine-mapping of top hits in independent cohorts**

516 We conducted GWAS and statistical fine-mapping in two additional large-scale cohorts, 345,235 white
517 British individuals from UKBB and 135,808 Japanese individuals from BBJ. For the UKBB white British, we
518 used previously conducted fine-mapping results for TC (<https://www.finucanelab.org/data>). Briefly, we
519 computed association statistics for the variants with INFO > 0.8, MAF > 0.01% (except for rare coding variants
520 with MAC > 0), and HWE p-value > 1e-10 using BOLT-LMM¹²⁴ with covariates including the top 20 PCs, sex,
521 age, age², sex * age, sex * age², and blood dilution factor. We used FINEMAP v1.3.1^{125,126} and susieR
522 v0.8.1.0521¹²⁷ for fine-mapping using the GWAS summary statistics and in-sample dosage LD matrices
523 computed by LDstore v2.0b. We defined regions based on 3 Mb window surrounding lead variants and
524 merged them if overlapped. The maximum number of causal variants in a region was specified as 10. For BBJ,
525 we additionally conducted fine-mapping using the same pipeline as we did for UKBB. The GWAS summary
526 statistics of TC was computed for the variants with Rsq > 0.7 and MAF > 0.01% using BOLT-LMM with the
527 covariates including top 20 PCs, sex, age, age², sex * age, sex * age², and disease status (affected versus
528 non-affected) for the 47 target diseases in the BBJ. The details about genotyping and imputation was
529 extensively described previously^{81,128}.

530

531 **Assessment of the correct empirical p value for admixed individuals**

532 To evaluate the appropriate p value threshold for *Tractor* associations, we estimated ancestry-specific
533 empirical null p value distributions via permutation. Although the genome-wide significance threshold ($p < 5 \times$
534 10^{-8}) is widely adopted in the current literature, previous work has shown that different ancestry groups have
535 different numbers of independent variants⁶⁸. Here, we permuted a null continuous phenotype 1,000 times
536 using the same admixed African-European individuals from UKBB as in the *Tractor* cholesterol GWAS to
537 assess the correct p value threshold for the admixed individuals in this study. We measured the minimum p
538 values of associations (p_{\min}) for each ancestry and derived an ancestry-specific empirical genome-wide
539 significance threshold as the fifth percentile ($\alpha = 0.05$) of p_{\min} across permutations as previously described¹²⁹.

540 We calculated this percentile using the Harrell–Davis distribution-free quantile estimator¹³⁰ and calculated the
541 95% confidence interval via bootstrapping. Based on the permutation results (**Figure S5a**), we defined a study-
542 wide significance threshold at a conservative level of $p = 1 \times 10^{-8}$ for both AFR- and EUR-specific
543 associations. In addition, we calculated the genomic inflation factor, λ_{GC} , of null phenotypes across
544 permutations and confirmed no significant inflation using the *Tractor* GWAS model (**Figure S5b**).

545

546 **Acknowledgements**

547 We thank Pradeep Natarajan, Sarah Gagliano Taliun, and many other scientists within and beyond Boston for
548 their intellectual contributions to this work. This project was supported by the National Institute of Mental Health
549 (K01 MH121659 and T32 MH017119 to E.G.A.; K99MH117229 to A.R.M.; 2R01MH106595 to C.M.N. and K.C.
550 K.). M.K. was supported by a Nakajima Foundation Fellowship and the Masason Foundation. M.L.S. was
551 supported by Fundacao de Amparo a Pesquisa do Estado de Sao Paulo (#2018/09328-2). The BioBank Japan
552 Project was supported by the Tailor-Made Medical Treatment Program of the Ministry of Education, Culture,
553 Sports, Science, and Technology (MEXT), the Japan Agency for Medical Research and Development (AMED).

554

555 **Author Contributions**

556 E.G.A. designed and implemented the pipeline, ran analyses, and drafted the primary manuscript. A.X.M.
557 designed and ran analyses. M.K. designed and ran analyses with the aid of J.C.U., Y.K., Y.O., and H.K.F.
558 A.R.M. contributed code and aided in writing the manuscript. K.J.K. and M.S. aided in code implementation.
559 K.C.K, C.M. N., B.M.N. and M.J.D. supervised and advised on the project. All authors reviewed and approved
560 the final draft.

561

562 **Competing Interests statement**

563 M.J.D. is a founder of Maze Therapeutics. A.R.M. serves as a consultant for 23andMe and is a member of the
564 Precise.ly Scientific Advisory Board. B.M.N. is a member of the Deep Genomics Scientific Advisory Board and

565 serves as a consultant for the Camp4 Therapeutics Corporation, Takeda Pharmaceutical and Biogen. The

566 remaining authors declare no competing interests.

567

568 **Code availability**

569 All code is freely available. *Tractor* scripts, in both cloud executable Hail format and python formats, can be

570 found at <https://github.com/eatkinson/Tractor>. The automated QC pipeline to prepare datasets for *Tractor* and

571 run LAI is located at <https://github.com/eatkinson/Post-QC>.

572

573 **References Cited**

574

- 575 1. Parker, K., Morin, R., Juliana Menasce Horowitz & Rohal, M. *Multiracial in America: Proud, Diverse and*
576 *Growing in Numbers*. (2015).
- 577 2. Bhardwaj, A. *et al.* Racial disparities in prostate cancer a molecular perspective. *Front. Biosci.* **22**, 4515
578 (2017).
- 579 3. Grizzle, W. E. *et al.* Self-Identified African Americans and prostate cancer risk: West African genetic
580 ancestry is associated with prostate cancer diagnosis and with higher Gleason sum on biopsy. *Cancer*
581 *Med.* **8**, 6915–6922 (2019).
- 582 4. Duggan, M. A., Anderson, W. F., Altekruse, S., Penberthy, L. & Sherman, M. E. The Surveillance,
583 Epidemiology, and End Results (SEER) Program and Pathology: Toward Strengthening the Critical
584 Relationship. *Am. J. Surg. Pathol.* **40**, e94–e102 (2016).
- 585 5. Freedman, M. L. *et al.* Admixture mapping identifies 8q24 as a prostate cancer risk locus in African-
586 American men. *Proc. Natl. Acad. Sci. U. S. A.* **103**, 14068–14073 (2006).
- 587 6. Bateman, E. D. *et al.* Global strategy for asthma management and prevention: GINA executive
588 summary. *Eur. Respir. J.* **31**, 143–78 (2008).
- 589 7. Daya, M. & Barnes, K. C. African American ancestry contribution to asthma and atopic dermatitis. *Ann.*
590 *Allergy. Asthma Immunol.* **122**, 456–462 (2019).
- 591 8. Wyss, A. B. *et al.* Multiethnic meta-analysis identifies ancestry-specific and cross-ancestry loci for
592 pulmonary function. *Nat. Commun.* **9**, 2976 (2018).
- 593 9. Demenais, F. *et al.* Multiancestry association study identifies new asthma risk loci that colocalize with
594 immune-cell enhancer marks. *Nat. Genet.* **50**, 42–50 (2018).
- 595 10. Benetos, A. & Aviv, A. Ancestry, Telomere Length, and Atherosclerosis Risk. *Circ. Cardiovasc. Genet.*
596 **10**, (2017).
- 597 11. Mozaffarian, D. *et al.* Heart Disease and Stroke Statistics—2015 Update. *Circulation* **131**, (2015).
- 598 12. Sirugo, G., Williams, S. M. & Tishkoff, S. A. The Missing Diversity in Human Genetic Studies. *Cell* **177**,

- 599 26–31 (2019).
- 600 13. Popejoy, Alice B., Fullerton, S. M. Genomics is falling. *Nature* **538**, 161–164 (2016).
- 601 14. Sul, J. H., Martin, L. S. & Eskin, E. Population structure in genetic studies: Confounding factors and
602 mixed models. *PLOS Genet.* **14**, e1007309 (2018).
- 603 15. Huang, H. *et al.* Bootstrat: Population Informed Bootstrapping for Rare Variant Tests. *bioRxiv* 068999
604 (2016). doi:10.1101/068999
- 605 16. Sohail, M. *et al.* Polygenic adaptation on height is overestimated due to uncorrected stratification in
606 genome-wide association studies. *Elife* **8**, (2019).
- 607 17. Berg, J. J. *et al.* Reduced signal for polygenic adaptation of height in UK biobank. *Elife* **8**, (2019).
- 608 18. Lander, E. S. & Schork, N. J. Genetic dissection of complex traits. *Science* (80-.). **265**, 2037–2048
609 (1994).
- 610 19. Coram, M. A., Fang, H., Candille, S. I., Assimes, T. L. & Tang, H. Leveraging Multi-ethnic Evidence for
611 Risk Assessment of Quantitative Traits in Minority Populations. *Am. J. Hum. Genet.* **101**, 218–226
612 (2017).
- 613 20. Walters, R. K. *et al.* Transancestral GWAS of alcohol dependence reveals common genetic
614 underpinnings with psychiatric disorders. *Nat. Neurosci.* **21**, 1656–1669 (2018).
- 615 21. Martin, E. R. *et al.* Properties of global- and local-ancestry adjustments in genetic association tests in
616 admixed populations. *Genet. Epidemiol.* **42**, 214–229 (2018).
- 617 22. Stevenson, A. *et al.* Neuropsychiatric Genetics of African Populations-Psychosis (NeuroGAP-
618 Psychosis): a case-control study protocol and GWAS in Ethiopia, Kenya, South Africa and Uganda. *BMJ*
619 *Open* **9**, e025469 (2019).
- 620 23. Consortium, T. H. Enabling the genomic revolution in Africa. *Science* **344**, 1346–8 (2014).
- 621 24. TOPMed Whole Genome Sequencing Project. Freeze 5b, Phases 1 and 2. (2018).
622 doi:10.1155/2013/865181
- 623 25. Precision Medicine Initiative (PMI) Working Group. The precision medicine initiative cohort program –
624 building a research foundation for 21st century medicine. *Precis. Med. Initiat. Work. Gr. Rep. to Advis.*

- 625 *Comm. to Dir. NIH Sept 17*, 1–108 (2015).
- 626 26. Logue, M. W. *et al.* The Psychiatric Genomics Consortium Posttraumatic Stress Disorder Workgroup:
627 Posttraumatic Stress Disorder Enters the Age of Large-Scale Genomic Collaboration.
628 *Neuropsychopharmacology* **40**, 2287–97 (2015).
- 629 27. Bien, S. A. *et al.* The Future of Genomic Studies Must Be Globally Representative: Perspectives from
630 PAGE. (2019). doi:10.1146/annurev-genom-091416
- 631 28. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat.*
632 *Genet.* **51**, 584–591 (2019).
- 633 29. Peterson, R. E. *et al.* Genome-wide Association Studies in Ancestrally Diverse Populations:
634 Opportunities, Methods, Pitfalls, and Recommendations. *Cell* **179**, 589–603 (2019).
- 635 30. Hero, J. O., Zaslavsky, A. M. & Blendon, R. J. The United States leads other nations in differences by
636 income in perceptions of health and health care. *Health Aff.* **36**, 1032–1040 (2017).
- 637 31. Williams, D. R., Priest, N. & Anderson, N. B. Understanding associations among race, socioeconomic
638 status, and health: Patterns and prospects. *Heal. Psychol.* **35**, 407–411 (2016).
- 639 32. Agency for Healthcare Research & Quality. *2016 National Healthcare Quality and Disparities Report.*
640 (2017).
- 641 33. Li, Y. R. & Keating, B. J. *Trans-ethnic genome-wide association studies: advantages and challenges of*
642 *mapping in diverse populations.* (2014). doi:10.1186/s13073-014-0091-5
- 643 34. Spain, S. L. & Barrett, J. C. Strategies for fine-mapping complex traits. *Hum. Mol. Genet.* **24**, R111-9
644 (2015).
- 645 35. Schaid, D. J., Chen, W. & Larson, N. B. *From genome-wide associations to candidate causal variants*
646 *by statistical fine-mapping.* *Nature Reviews Genetics* **19**, 491–504 (2018).
- 647 36. Wu, Y. *et al.* Trans-Ethnic Fine-Mapping of Lipid Loci Identifies Population-Specific Signals and Allelic
648 Heterogeneity That Increases the Trait Variance Explained. *PLoS Genet.* **9**, e1003379 (2013).
- 649 37. van de Bunt, M. *et al.* Evaluating the Performance of Fine-Mapping Strategies at Common Variant
650 GWAS Loci. *PLOS Genet.* **11**, e1005535 (2015).

- 651 38. Mahajan, A. *et al.* Genome-wide trans-ancestry meta-analysis provides insight into the genetic
652 architecture of type 2 diabetes susceptibility. *Nat. Genet.* **46**, 234–244 (2014).
- 653 39. Project, T. T. G. C. An integrated map of genetic variation from 1,092 human genomes. *Nature* **135**,
654 (2012).
- 655 40. Skotte, L., Korneliussen, T. S. S., Moltke, I. & Albrechtsen, A. Ancestry specific association mapping in
656 admixed populations. *bioRxiv* 1–33 (2018). doi:10.1101/014001
- 657 41. Zhang, J. & Stram, D. O. The Role of Local Ancestry Adjustment in Association Studies Using Admixed
658 Populations. *Genet. Epidemiol.* **38**, 502–515 (2014).
- 659 42. Smith, E. N. *et al.* Genome-wide association study of bipolar disorder in European American and African
660 American individuals. *Mol. Psychiatry* **14**, 755–763 (2009).
- 661 43. Szulc, P., Bogdan, M., Frommlet, F. & Tang, H. Joint genotype- and ancestry-based genome-wide
662 association studies in admixed populations. *Genet. Epidemiol.* **41**, 555–566 (2017).
- 663 44. Tang, H., Siegmund, D. O., Johnson, N. A., Romieu, I. & London, S. J. Joint testing of genotype and
664 ancestry association in admixed families. *Genet. Epidemiol.* **34**, 783–791 (2010).
- 665 45. Coram, M. A. *et al.* Genome-wide Characterization of Shared and Distinct Genetic Components that
666 Influence Blood Lipid Levels in Ethnically Diverse Human Populations. *Am. J. Hum. Genet.* **92**, 904–916
667 (2013).
- 668 46. Aschard, H., Gusev, A., Brown, R. & Pasaniuc, B. Leveraging local ancestry to detect gene-gene
669 interactions in genome-wide data. *BMC Genet.* **16**, 1–9 (2015).
- 670 47. Zaitlen, N., Pas, B., Gur, T., Ziv, E. & Halperin, E. ARTICLE Leveraging Genetic Variability across
671 Populations for the Identification of Causal Variants. *Am. J. Hum. Genet.* **86**, 23–33
- 672 48. Pasaniuc, B. *et al.* Enhanced statistical tests for GWAS in admixed populations: assessment using
673 African Americans from CARE and a Breast Cancer Consortium. *PLoS Genet.* **7**, e1001371 (2011).
- 674 49. Pasaniuc, B. *et al.* Analysis of Latino populations from GALA and MEC studies reveals genomic loci with
675 biased local ancestry estimation. *Bioinformatics* **29**, 1407–1415 (2013).
- 676 50. Chimusa, E. R. *et al.* Genome-wide association study of ancestry-specific TB risk in the South African

- 677 coloured population. *Hum. Mol. Genet.* **23**, 796–809 (2014).
- 678 51. Shriner, D. Overview of admixture mapping. in *Current Protocols in Human Genetics* 1.23.1-1.23.8
679 (2013). doi:10.1002/cphg.44
- 680 52. Chen, M. *et al.* Admixture mapping analysis in the context of GWAS with GAW18 data. in *BMC*
681 *Proceedings* **8**, S3 (BioMed Central Ltd., 2014).
- 682 53. Chen, W. *et al.* A Generalized Sequential Bonferroni Procedure for GWAS in Admixed Populations
683 Incorporating Admixture Mapping Information into Association Tests. *Hum. Hered.* **79**, 80–92 (2015).
- 684 54. Hoggart, C. J., Shriver, M. D., Kittles, R. A., Clayton, D. G. & McKeigue, P. M. Design and Analysis of
685 Admixture Mapping Studies. *Am. J. Hum. Genet.* **74**, 965–978 (2004).
- 686 55. Patterson, N. *et al.* Methods for High-Density Admixture Mapping of Disease Genes. *Am. J. Hum.*
687 *Genet.* **74**, 979–1000 (2004).
- 688 56. Spear, M. L. *et al.* A genome-wide association and admixture mapping study of bronchodilator drug
689 response in African Americans with asthma. *Pharmacogenomics J.* **19**, 249–259 (2019).
- 690 57. Gignoux, C. R. *et al.* An admixture mapping meta-analysis implicates genetic variation at 18q21 with
691 asthma susceptibility in Latinos. *J. Allergy Clin. Immunol.* **143**, 957–969 (2019).
- 692 58. Shetty, P. B. *et al.* Variants for HDL-C, LDL-C, and triglycerides identified from admixture mapping and
693 fine-mapping analysis in African American families. *Circ. Cardiovasc. Genet.* **8**, 106–113 (2015).
- 694 59. Shetty, P. B. *et al.* Variants in CXADR and F2RL1 are associated with blood pressure and obesity in
695 African-Americans in regions identified through admixture mapping. *J. Hypertens.* **30**, 1970–1976
696 (2012).
- 697 60. Reiner, A. P. *et al.* Genome-wide association and population genetic analysis of c-reactive protein in
698 african american and hispanic american women. *Am. J. Hum. Genet.* **91**, 502–512 (2012).
- 699 61. Florez, J. C. *et al.* Strong Association of Socioeconomic Status and Genetic Ancestry in Latinos:
700 Implications for Admixture Studies of Type 2 Diabetes. in *Racial Identities, Genetic Ancestry, and Health*
701 *in South America: Argentina, Brazil, Colombia, and Uruguay* (eds. Gibbon, S., Santos, R. V. & Sans, M.)
702 137–153 (Palgrave Macmillan US, 2011). doi:10.1057/9781137001702_7

- 703 62. Maples, B. K., Gravel, S., Kenny, E. E. & Bustamante, C. D. RFMix: a discriminative modeling approach
704 for rapid and robust local-ancestry inference. *Am. J. Hum. Genet.* **93**, 278–88 (2013).
- 705 63. Geza, E. *et al.* A comprehensive survey of models for dissecting local ancestry deconvolution in human
706 genome. *Brief. Bioinform.* **20**, 1709–1724 (2019).
- 707 64. Schubert, R., Andaleon, A. & Wheeler, H. E. *Comparing local ancestry inference models in populations*
708 *of two- and three-way admixture. Research Square* (2018).
- 709 65. Choi, Y., Chan, A. P., Kirkness, E., Telenti, A. & Schork, N. J. Comparison of phasing strategies for
710 whole human genomes. *PLOS Genet.* **14**, e1007308 (2018).
- 711 66. Andrés, A. M. *et al.* Understanding the accuracy of statistical haplotype inference with sequence data of
712 known phase. *Genet. Epidemiol.* **31**, 659–71 (2007).
- 713 67. O'Connell, J. *et al.* A general approach for haplotype phasing across the full spectrum of relatedness.
714 *PLoS Genet.* **10**, e1004234 (2014).
- 715 68. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 716 69. Natarajan, P. *et al.* Deep-coverage whole genome sequences and blood lipids among 16,324
717 individuals. *Nat. Commun.* **9**, 3391 (2018).
- 718 70. Musunuru, K. & Kathiresan, S. Genetics of Common, Complex Coronary Artery Disease. *Cell* **177**, 132–
719 145 (2019).
- 720 71. Rotimi, C. N. *et al.* The genomic landscape of African populations in health and disease. *Hum. Mol.*
721 *Genet.* **26**, R225–R236 (2017).
- 722 72. Superko, H. R., Momary, K. M. & Li, Y. Statins Personalized. *Medical Clinics of North America* **96**, 123–
723 139 (2012).
- 724 73. Fu, J. *et al.* Unraveling the regulatory mechanisms underlying tissue-dependent genetic variation of
725 gene expression. *PLoS Genet.* **8**, (2012).
- 726 74. Avery, C. L. *et al.* A phenomics-based strategy identifies loci on APOC1, BRAP, and PLCG1 associated
727 with metabolic syndrome phenotype domains. *PLoS Genet.* **7**, (2011).
- 728 75. Lettre, G. *et al.* Genome-Wide association study of coronary heart disease and its risk factors in 8,090

- 729 african americans: The nhlbi CARE project. *PLoS Genet.* **7**, (2011).
- 730 76. Talmud, P. J. *et al.* Gene-centric Association Signals for Lipids and Apolipoproteins Identified via the
731 HumanCVD BeadChip. *Am. J. Hum. Genet.* **85**, 628–642 (2009).
- 732 77. Sandhu, M. S. *et al.* LDL-cholesterol concentrations: a genome-wide association study. *Lancet* **371**,
733 483–491 (2008).
- 734 78. Sanna, S. *et al.* Fine mapping of five loci associated with low-density lipoprotein cholesterol detects
735 variants that double the explained heritability. *PLoS Genet.* **7**, (2011).
- 736 79. Fox, C. S. *et al.* Genome-wide association to body mass index and waist circumference: the
737 Framingham Heart Study 100K project. *BMC Med. Genet.* **8**, S18 (2007).
- 738 80. Kathiresan, S. *et al.* A genome-wide association study for blood lipid phenotypes in the Framingham
739 Heart Study. *BMC Med. Genet.* **8**, S17 (2007).
- 740 81. Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese population links cell types to
741 complex human diseases. *Nat. Genet.* **50**, 390–400 (2018).
- 742 82. Karczewski, K. J. *et al.* Variation across 141,456 human exomes and genomes reveals the spectrum of
743 loss-of-function intolerance across human protein-coding genes. *bioRxiv* 531210 (2019).
744 doi:10.1101/531210
- 745 83. Lin, M. *et al.* Population specific reference panels are crucial for the genetic analyses of Native
746 Hawaiians: an example of the CREBRF locus. *bioRxiv* 789073 (2019). doi:10.1101/789073
- 747 84. Ntzani, E. E., Liberopoulos, G., Manolio, T. A. & Ioannidis, J. P. A. Consistency of genome-wide
748 associations across major ancestral groups. *Hum. Genet.* **131**, 1057–1071 (2012).
- 749 85. Marigorta, U. M. & Navarro, A. High trans-ethnic replicability of GWAS results implies common causal
750 variants. *PLoS Genet.* **9**, e1003566 (2013).
- 751 86. Waters, K., Stram, D., ... M. H.-PI. & 2010, undefined. Consistent association of type 2 diabetes risk
752 variants found in europeans in diverse racial and ethnic groups. *ncbi.nlm.nih.govPaperpile*
- 753 87. Lam, M. *et al.* Comparative genetic architectures of schizophrenia in East Asian and European
754 populations. *Nat. Genet.* (2019). doi:10.1038/s41588-019-0512-x

- 755 88. Liu, J. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel disease and
756 highlight shared genetic risk across populations. *nature.comPaperpile*
- 757 89. Carlson, C. S. *et al.* Generalization and Dilution of Association Results from European GWAS in
758 Populations of Non-European Ancestry: The PAGE Study. *PLoS Biol.* **11**, (2013).
- 759 90. Easton, D., Pooley, K., Dunning, A., Nature, P. P.- & 2007, undefined. Genome-wide association study
760 identifies novel breast cancer susceptibility loci. *nature.comPaperpile*
- 761 91. Mägi, R. *et al.* Trans-ethnic meta-regression of genome-wide association studies accounting for
762 ancestry increases power for discovery and improves fine-mapping resolution. *Hum. Mol. Genet.* **26**,
763 3639–3650 (2017).
- 764 92. Wegmann, D. *et al.* Recombination rates in admixed individuals identified by ancestry-based inference.
765 *Nat. Genet.* **43**, 847–853 (2011).
- 766 93. Zhang, R. The ANGPTL3-4-8 model, a molecular mechanism for triglyceride trafficking. *Open Biol.* **6**,
767 150272 (2016).
- 768 94. Fu, Z., Abou-Samra, A. B. & Zhang, R. A lipasin/Angptl8 monoclonal antibody lowers mouse serum
769 triglycerides involving increased postprandial activity of the cardiac lipoprotein lipase. *Sci. Rep.* **5**, 18502
770 (2015).
- 771 95. Zhang, R. Lipasin, a novel nutritionally-regulated liver-enriched factor that regulates serum triglyceride
772 levels. *Biochem. Biophys. Res. Commun.* **424**, 786–792 (2012).
- 773 96. Siddiqa, A. *et al.* Visualizing the regulatory role of Angiotensin-like protein 8 (ANGPTL8) in glucose and
774 lipid metabolic pathways. *Genomics* **109**, 408–418 (2017).
- 775 97. Yamada, H. *et al.* Circulating betatrophin is elevated in patients with type 1 and type 2 diabetes. *Endocr.*
776 *J.* **62**, 417–421 (2015).
- 777 98. Espes, D., Martinell, M. & Carlsson, P.-O. Increased circulating betatrophin concentrations in patients
778 with type 2 diabetes. *Int. J. Endocrinol.* **2014**, 323407 (2014).
- 779 99. Hu, H. *et al.* Increased circulating levels of betatrophin in newly diagnosed type 2 diabetic patients.
780 *Diabetes Care* **37**, 2718–22 (2014).

- 781 100. Fu, Z. *et al.* Elevated circulating lipasin/betatrophin in human type 2 diabetes and obesity. *Sci. Rep.* **4**,
782 5013 (2015).
- 783 101. Cannon, M. E. *et al.* Trans-ancestry Fine Mapping and Molecular Assays Identify Regulatory Variants at
784 the ANGPTL8 HDL-C GWAS Locus. (2017). doi:10.1534/g3.117.300088
- 785 102. Adzhubei, I., Jordan, D. M. & Sunyaev, S. R. Predicting functional effect of human missense mutations
786 using PolyPhen-2. *Curr. Protoc. Hum. Genet.* **Chapter 7**, Unit7.20 (2013).
- 787 103. Ng, P. C. & Henikoff, S. SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids*
788 *Res.* **31**, 3812–4 (2003).
- 789 104. Atkinson, E. G. *et al.* No Evidence for Recent Selection at FOXP2 among Diverse Human Populations.
790 *Cell* **174**, 1424-1435.e15 (2018).
- 791 105. Deng, L., Ruiz-Linares, A., Xu, S. & Wang, S. Ancestry variation and footprints of natural selection along
792 the genome in Latin American populations. *Sci. Rep.* **6**, 1–7 (2016).
- 793 106. Jin, W. *et al.* Genome-wide detection of natural selection in African Americans pre- and post-admixture.
794 *Genome Res.* **22**, 519–527 (2012).
- 795 107. Bomba, L., Walter, K. & Soranzo, N. The impact of rare and low-frequency genetic variants in common
796 disease. *Genome Biol.* **18**, 77 (2017).
- 797 108. Mathieson, I. & McVean, G. Differential confounding of rare and common variants in spatially structured
798 populations. **44**, 243–246 (2012).
- 799 109. The Hail team. Hail. (2018). Available at: <https://github.com/hail-is/hail>. (Accessed: 16th January 2019)
- 800 110. Kluyver, T. *et al.* *Jupyter Notebooks—a publishing format for reproducible computational workflows.*
801 *Positioning and Power in Academic Publishing: Players, Agents and Agendas* (2016). doi:10.3233/978-
802 1-61499-649-1-87
- 803 111. International Hapmap Consortium, T. The International HapMap Project. *Nature* **426**, 789–796 (2003).
- 804 112. Tishkoff, S. a *et al.* The genetic structure and history of Africans and African Americans. *Science* **324**,
805 1035–44 (2009).
- 806 113. Gravel, S. *et al.* Demographic history and rare allele sharing among human populations. *Proc. Natl.*

- 807 *Acad. Sci. U. S. A.* **108**, 11983–8 (2011).
- 808 114. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction across Diverse
809 Populations. (2017). doi:10.1016/j.ajhg.2017.03.004
- 810 115. Chen, C. Y. *et al.* Improved ancestry inference using weights from external reference panels.
811 *Bioinformatics* **29**, 1399–1406 (2013).
- 812 116. Williams, A. admix-simu: program to simulate admixture between multiple populations. (2016).
813 doi:10.5281/ZENODO.45517
- 814 117. Skotte, L., Jørsboe, E., Korneliussen, T. S., Moltke, I. & Albrechtsen, A. Ancestry-specific association
815 mapping in admixed populations. *Genet. Epidemiol.* **43**, 506–521 (2019).
- 816 118. Cann, H. M. *et al.* A human genome diversity cell line panel. *Science* **296**, 261–2 (2002).
- 817 119. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
818 individuals. *Genome Res.* **19**, 1655–64 (2009).
- 819 120. Google Cloud Platform Blog. Google Compute Engine launches, expanding Google’s cloud offerings.
820 Available at: <https://cloudplatform.googleblog.com/2012/06/google-compute-engine-launches.html>.
821 (Accessed: 16th January 2019)
- 822 121. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage
823 analyses. *Am. J. Hum. Genet.* **81**, 559–75 (2007).
- 824 122. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association scan results.
825 *Bioinformatics* **26**, 2336–2337 (2010).
- 826 123. Bokeh Development Team. Bokeh: Python library for interactive visualization. (2019). Available at:
827 <https://bokeh.org/citation/>. (Accessed: 31st March 2020)
- 828 124. Loh, P. R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts.
829 *Nat. Genet.* **47**, (2015).
- 830 125. Benner, C. *et al.* FINEMAP: Efficient variable selection using summary data from genome-wide
831 association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 832 126. Benner, C., Havulinna, A., Salomaa, V., Ripatti, S. & Pirinen, M. Refining fine-mapping: effect sizes and

- 833 regional heritability. *bioRxiv* 318618 (2018). doi:10.1101/318618
- 834 127. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable selection in
835 regression, with application to genetic fine-mapping. *bioRxiv* 501114 (2019). doi:10.1101/501114
- 836 128. Akiyama, M. *et al.* Characterizing rare and low-frequency height-associated variants in the Japanese
837 population. *Nat. Commun.* **10**, 1–11 (2019).
- 838 129. Kanai, M., Tanaka, T. & Okada, Y. Empirical estimation of genome-wide significance thresholds based
839 on the 1000 Genomes Project data set. *J. Hum. Genet.* **61**, 861–866 (2016).
- 840 130. HARRELL, F. E. & DAVIS, C. E. A new distribution-free quantile estimator. *Biometrika* **69**, 635–640
841 (1982).
- 842
- 843

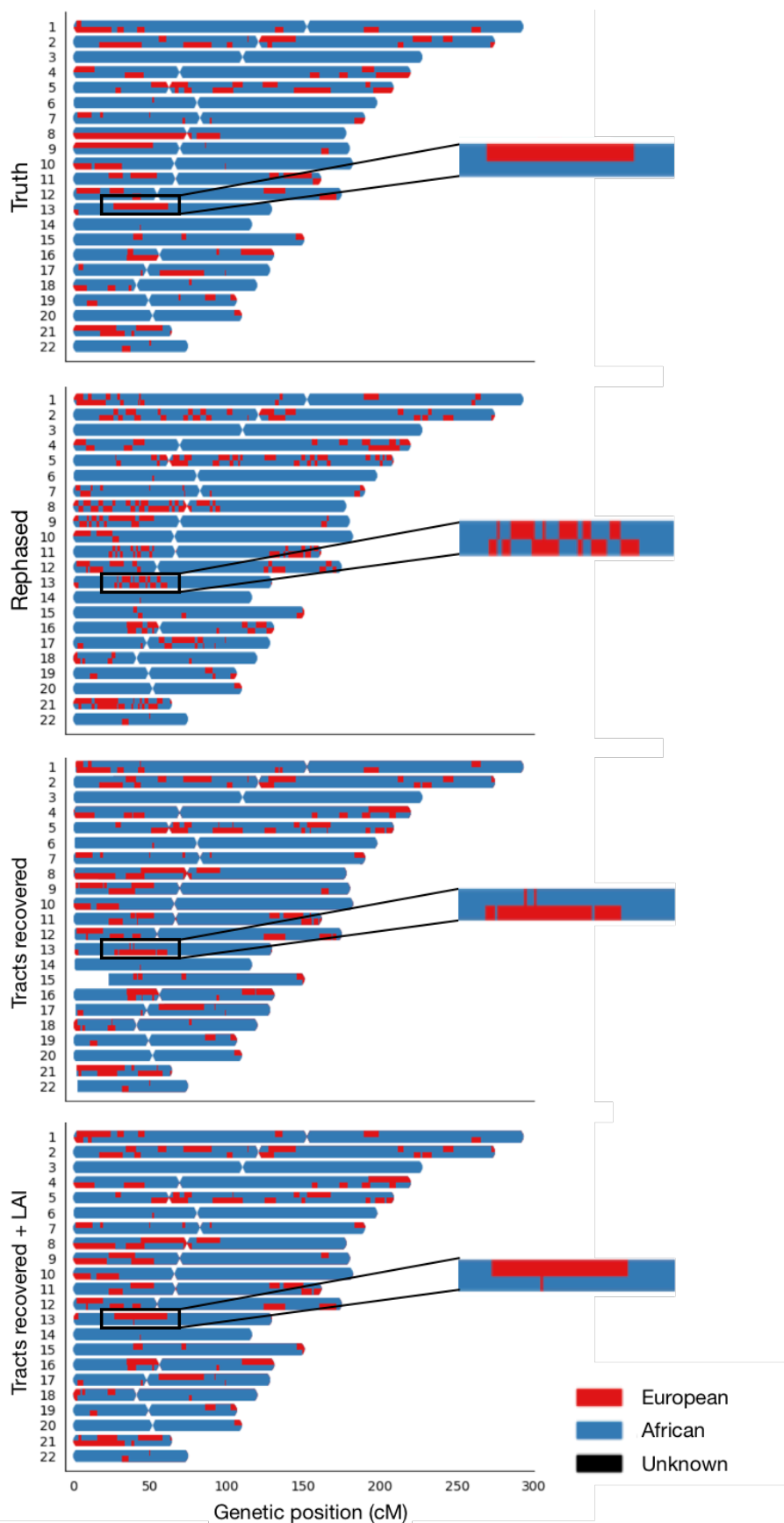


Figure 1. Painted karyograms of a simulated AA individual showing EUR (red) and AFR (blue) ancestral tracts across data treatments. The top panel shows the truth results for an example individual in our simulated AA cohort. A painted karyogram after statistical phasing is shown in the second row – note the disruption of long haplotypes. The third panel illustrates our recovery of tracts broken by switch errors in phasing. The bottom panel shows the smoothing and further improvement of tracts acquired through an additional round of LAI. The same section of chr13 showing an example tract at higher resolution is pictured on the right to highlight tract recovery.

Figure 2. Tractor recovers disrupted tracts, improving tract distributions. The top row (**A-C**) shows the improvements to the distributions of the number of discrete EUR tracts observed in simulated AA individuals under demographic models of 1 pulse of admixture at 3, 9 (realistic for AA population history) and 20 generations ago. The bottom row (**D,E**) shows the results from different initial admixture fractions, of 70% and 50% AFR, respectively, at the realistic 9 generations since admixture. These can be compared to the inferred demographic model in AA with ~80% AFR ancestry shown in **B**. In all panels, the simulated truth dataset is shown in black, after statistical phasing in purple, immediately after tract recovery procedures is in orange, and after one additional round of LAI after tract recovery in yellow.

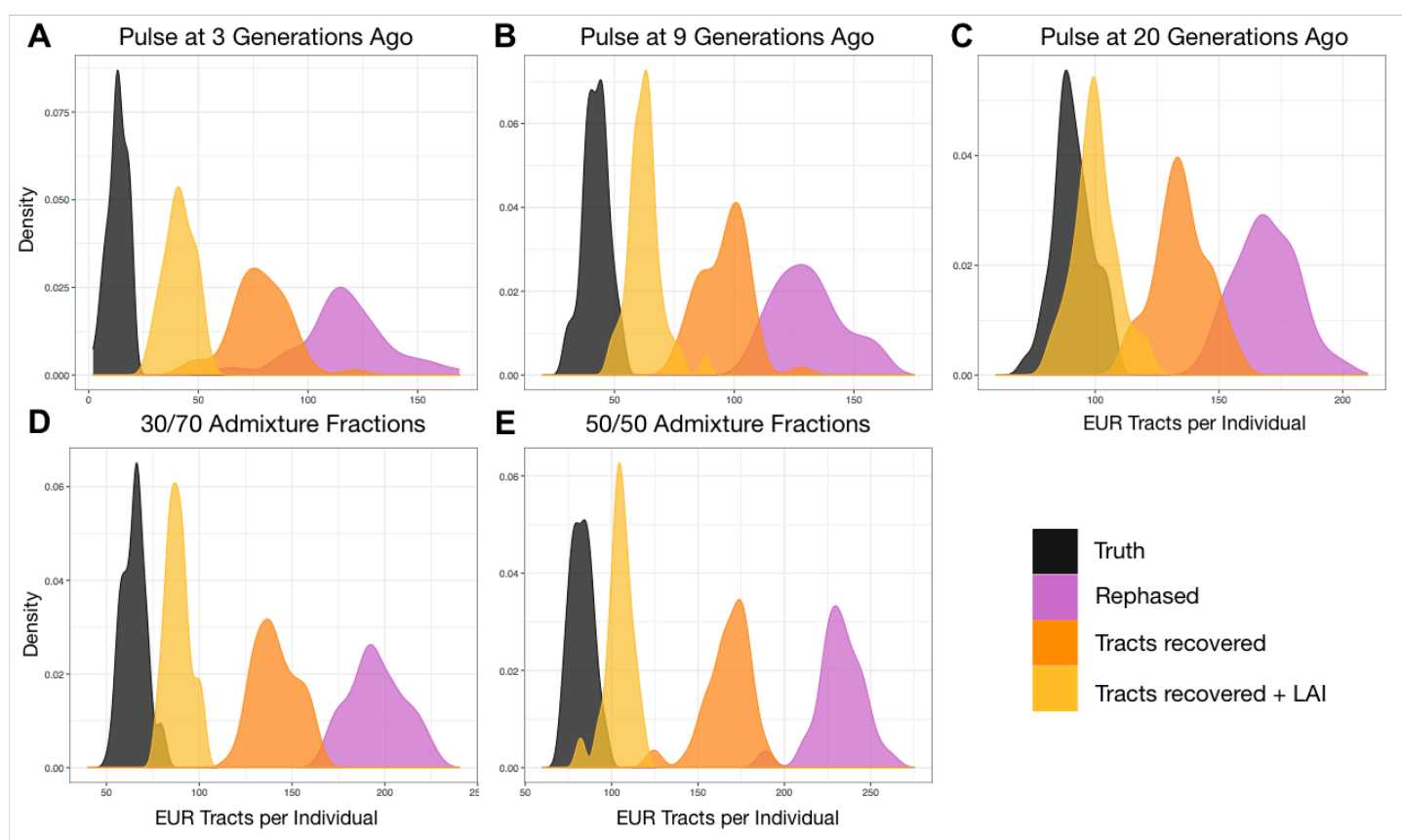


Figure 3. GWAS power gains across sample sizes, ancestral MAF differences, admixture proportions, and effect size differences. In all scenarios shown, dashed lines correspond to the power from the *Tractor* model incorporating local ancestry, solid lines are for a traditional GWAS model. In all panels we modeled a 10% disease prevalence. Unless otherwise noted, we used the parameters for a realistic demographic scenario for AA individuals: 80% AFR ancestry, an effect present only in the AFR genetic background, 12k cases and 30k controls, and 20% MAF. **(A)** There are similar gains in GWAS power when using the *Tractor* LAI-aware model across samples sizes of 4,000 (grey) and 12,000 (black) cases with 2x controls. **(B)** When there is a MAF difference between ancestries, the gains in power are even more pronounced. Gains vary across the allele frequency spectrum: black=MAF 10% AFR, 30% EUR; grey=MAF 20% AFR, 40% EUR. **(C)** Gains become more pronounced when the admixture fractions are modified to 50/50. **(D)** Dramatic gains are obtained when the effect is switched to instead only be present on the EUR background.

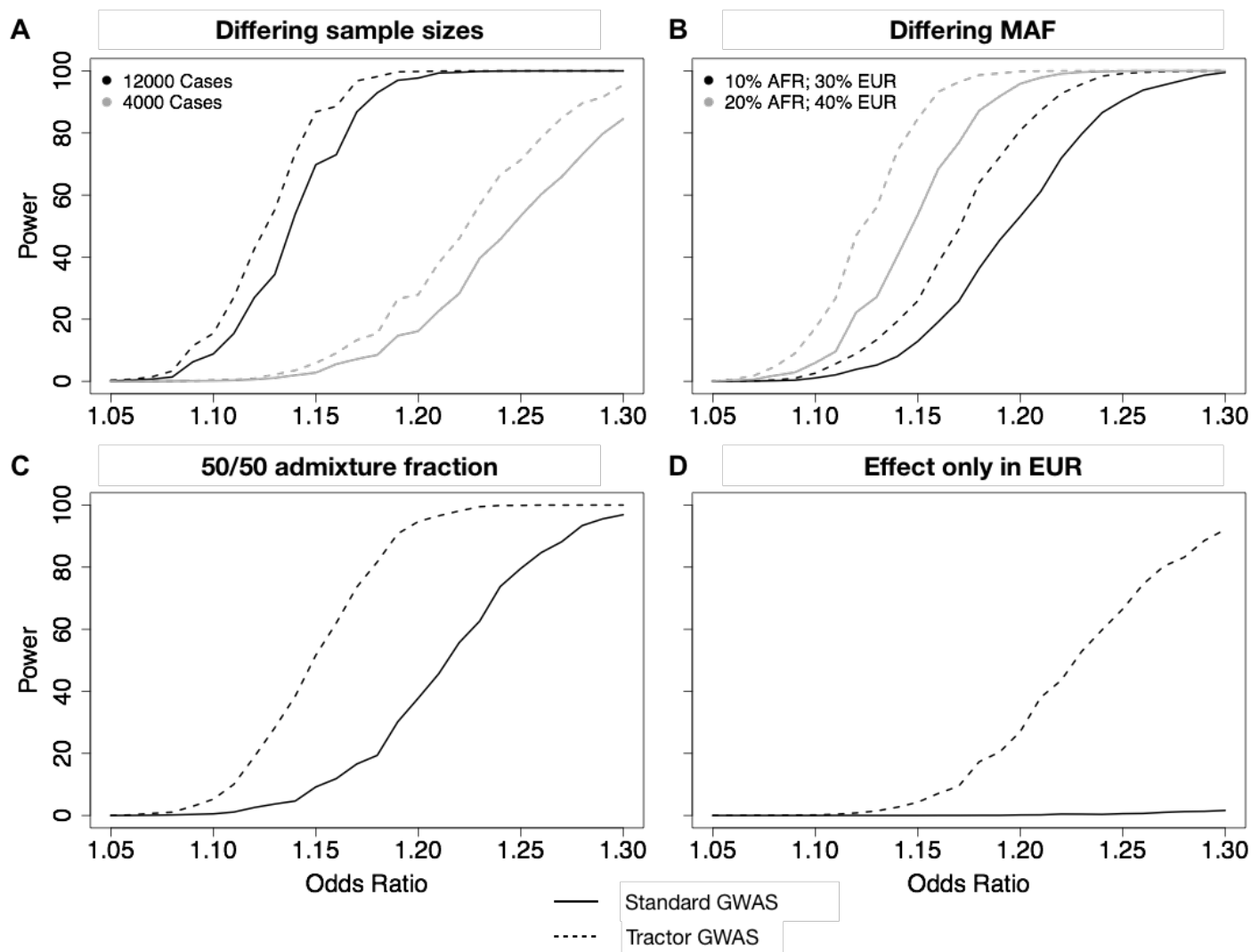
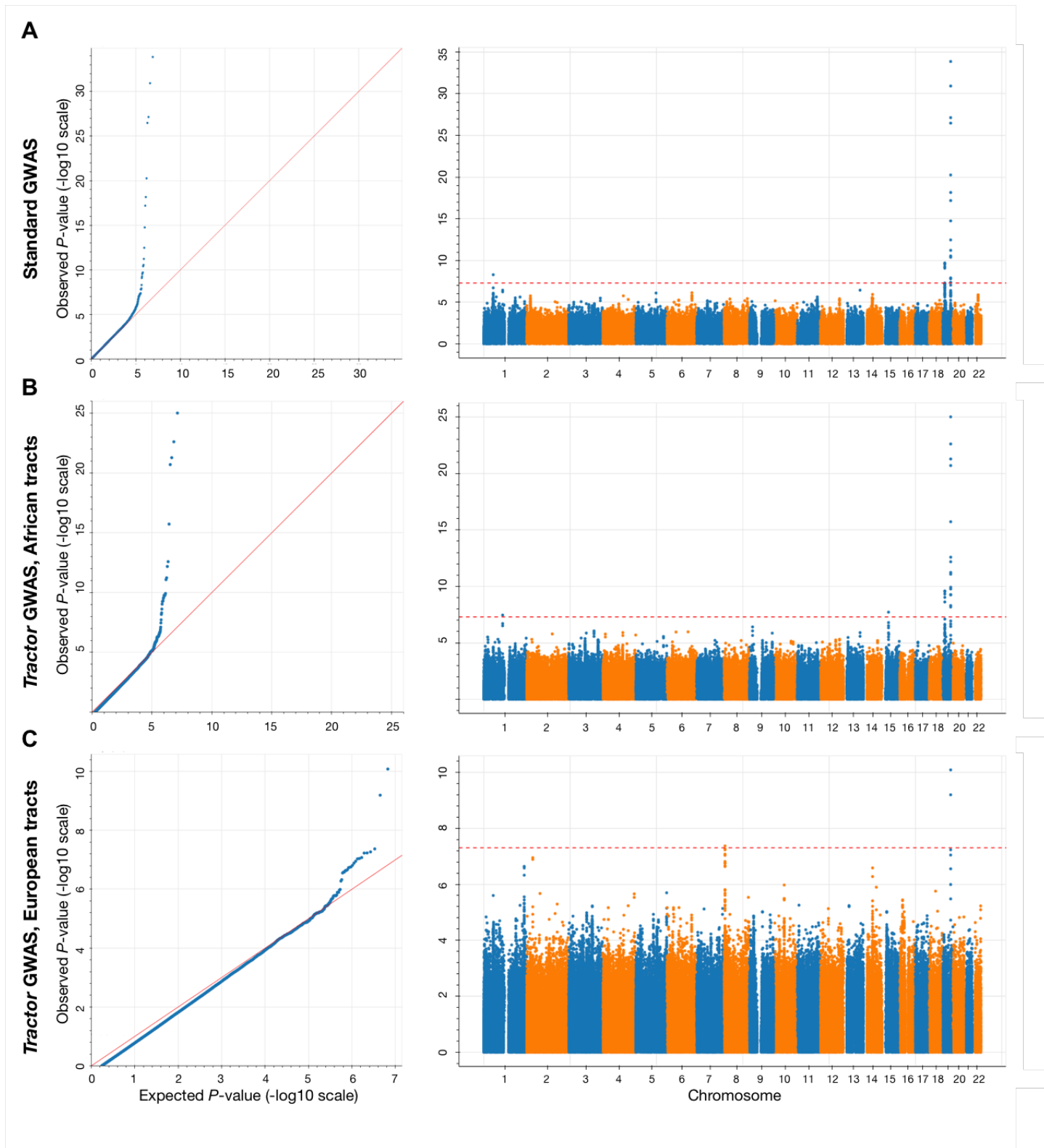


Figure 4. *Tractor* GWAS replicates established hits for Total Cholesterol in admixed African-European individuals and identifies new ancestry-specific loci. QQ and Manhattan plots for Total Cholesterol using the standard GWAS model (A) compared to *Tractor* joint-analysis results for the AFR (B) and EUR (C) backgrounds.



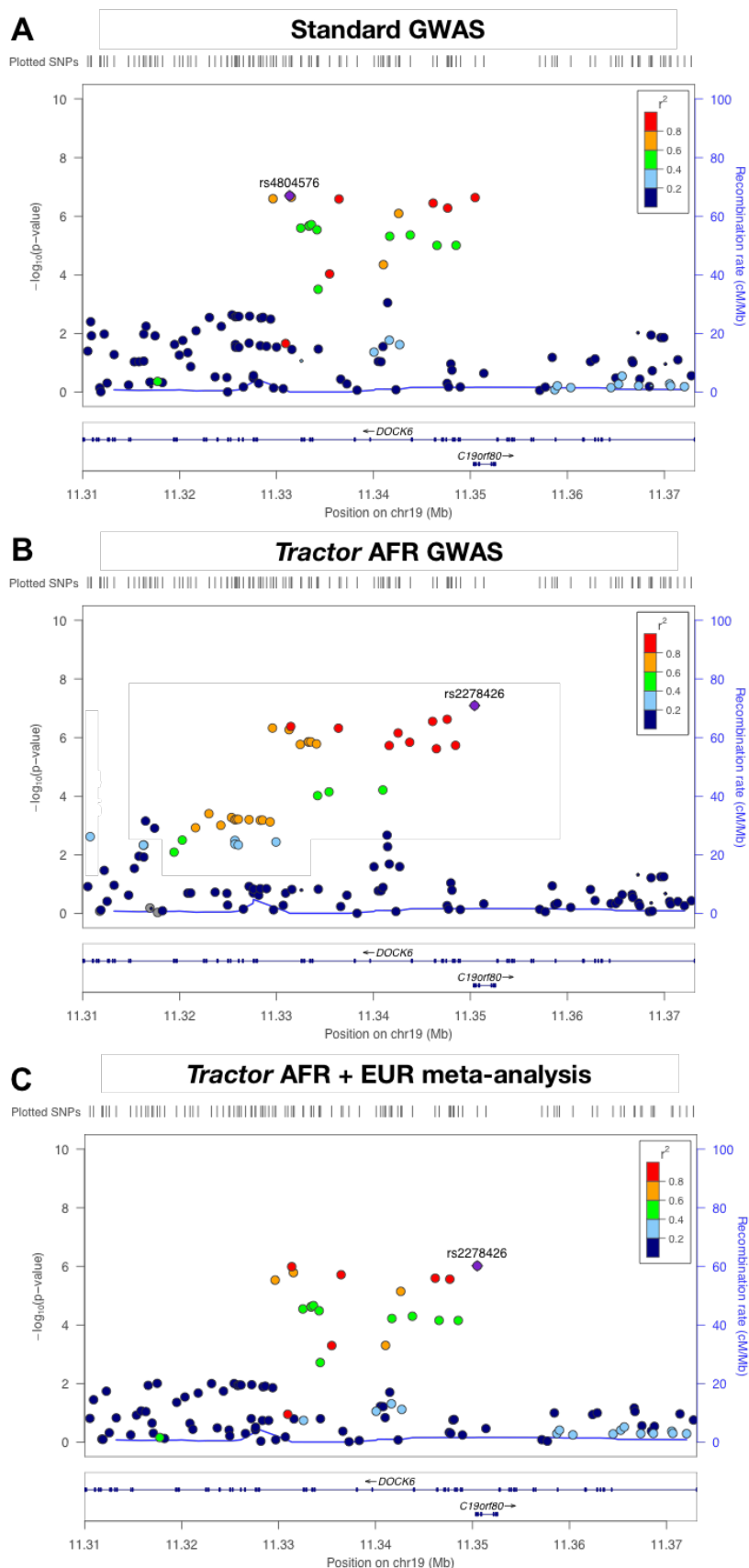


Figure 5. *Tractor* better localizes a top hit for Total Cholesterol. Previous hits for TC had pinpointed *DOCK6* as the gene of interest. Comparing runs on UKBB admixed individuals with a standard GWAS model (A), AFR-specific GWAS with *Tractor* (B), and a meta-analysis of GWAS runs on deconvolved EUR and AFR tracts (C), both ancestry-specific runs pinpoint a lead SNP ~20kb downstream in an intron of *DOCK6* spanning a better candidate gene *ANGPTL8* (also known as C19orf80) as the lead SNP. No significant signal was seen in the EUR segments. In all plots, point size is proportional to the number of samples included for that test, and color indicates r^2 to the named lead SNP. For B, the recombination rate line was generated from the AFR superpopulation of the 1000 Genomes Project, for other panels the EUR superpopulation rate is shown.