

# 1 Confirmation of HLA-II associations with TB susceptibility in admixed 2 African samples

3 Dayna Croock<sup>1</sup>, Yolandi Swart<sup>1</sup>, Haiko Schurz<sup>1</sup>, Desiree C. Petersen<sup>1</sup>, Marlo Möller<sup>1,2</sup>, Caitlin Uren<sup>1,2\*</sup>

4  
5 <sup>1</sup>DSI-NRF Centre of Excellence for Biomedical Tuberculosis Research, South African Medical Research Council Centre  
6 for Tuberculosis Research, Division of Molecular Biology and Human Genetics, Faculty of Medicine and Health  
7 Sciences, Stellenbosch University

8 <sup>2</sup>Centre for Bioinformatics and Computational Biology, Stellenbosch University

9 \*Corresponding author: [caitlinu@sun.ac.za](mailto:caitlinu@sun.ac.za)

10

## 11 Abstract

12 The International Tuberculosis Host Genetics Consortium (ITHGC) demonstrated the  
13 power of large-scale GWAS analysis across diverse ancestries in identifying tuberculosis  
14 (TB) susceptibility loci. Despite identifying a significant genetic correlate in the human  
15 leukocyte antigen (HLA)-II region, this association did not replicate in the African  
16 ancestry-specific analysis, due to small sample size and the inclusion of admixed samples.  
17 Our study aimed to build upon the findings from the ITHGC and identify TB susceptibility  
18 loci in an admixed South African cohort using the local ancestry allelic adjusted  
19 association (LAAA) model. We identified a near-genome-wide significant association  
20 (*rs3117230*,  $p$ -value =  $5.292 \times 10^{-6}$ , OR = 0.437, SE = 0.182) in the *HLA-DPB1* gene  
21 originating from KhoeSan ancestry. These findings extend the work of the ITHGC,  
22 underscore the need for innovative strategies in studying complex admixed populations,  
23 and confirm the role of the HLA-II region in TB susceptibility in admixed South African  
24 samples.

25

## 26 Keywords

27 Human leukocyte antigen (HLA)-II, tuberculosis (TB), local ancestry, admixture, KhoeSan  
28 ancestry

29

## 30 Introduction

31 Tuberculosis (TB) is a communicable disease caused by *Mycobacterium tuberculosis* (*M.tb*)  
32 (World Health Organization, 2023). *M.tb* infection has a wide range of clinical manifestations

33 from asymptomatic, non-transmissible, or so-called “latent” infections to active TB (Zaidi et

34 al., 2023). Approximately 1/4 of the global population is infected with *M.tb*, but only 5-15% of  
35 infected individuals will develop active TB (Menzies et al., 2021). Several factors increase the  
36 risk of progressing to active TB, including co-infection with human immunodeficiency virus  
37 (HIV) and comorbidities, such as diabetes mellitus, asthma and other airway and lung  
38 diseases (Glaziou et al., 2018). Socio-economic factors including smoking, malnutrition,  
39 alcohol abuse, intravenous drug use, prolonged residence in a high burdened community,  
40 overcrowding, informal housing and poor sanitation also influence *M.tb* transmission and  
41 infection (Cudahy et al., 2020; Escombe et al., 2019; Laghari et al., 2019; Matose et al.,  
42 2019; Smith et al., 2023). Additionally, individual variability in infection and disease  
43 progression has been attributed to variation in the host genome (Schurz et al., 2024; Caitlin  
44 Uren et al., 2021; Verhein et al., 2018). Numerous genome-wide association studies  
45 (GWASs) investigating TB susceptibility have been conducted across different population  
46 groups. However, findings from these studies often do not replicate across population groups  
47 (Möller & Kinnear, 2020; Möller et al., 2018; Caitlin Uren et al., 2017). This lack of  
48 replication could be caused by small sample sizes, variation in phenotype definitions among  
49 studies, variation in linkage disequilibrium (LD) patterns across different population groups  
50 and the presence of population-specific effects (Möller & Kinnear, 2020). Additionally,  
51 complex LD patterns within population groups, produced by admixture, impede the  
52 detection of statistically significant loci when using traditional GWAS methods (Swart et al.,  
53 2020).

54  
55 The International Tuberculosis Host Genetics Consortium (ITHGC) performed a meta-  
56 analysis of TB GWAS results including 14 153 TB cases and 19 536 controls of African, Asian  
57 and European ancestries (Schurz et al., 2024). The multi-ancestry meta-analysis identified  
58 one genome-wide significant variant (*rs28383206*) in the human leukocyte antigen (HLA)-II  
59 region ( $p = 5.2 \times 10^{-9}$ , OR = 0.89, 95% CI = 0.84-0.95). The association peak at the *HLA-II* locus  
60 encompassed several genes encoding crucial antigen presentation proteins (including *HLA-*  
61 *DR* and *HLA-DQ*). While ancestry-specific association analyses in the European and Asian  
62 cohorts also produced suggestive peaks in the HLA-II region, the African ancestry-specific  
63 association test did not yield any associations or suggestive peaks. The authors described  
64 possible reasons for the lack of associations, including the smaller sample size compared to  
65 the other ancestry-specific meta-analyses, increased genetic diversity within African  
66 individuals and population stratification produced by two admixed cohorts from the South

67 African Coloured (SAC) population (Schurz et al., 2024). The SAC population (as termed in  
68 the South African census (Lehohla, 2012)) form part of a multi-way (up to five-way) admixed  
69 population with ancestral contributions from Bantu-speaking African (~30%), KhoeSan  
70 (~30%), European (~20%), and East (~10%) and Southeast Asian (~10%) populations  
71 (Chimusa et al., 2013). The diverse genetic background of admixed individuals can lead to  
72 population stratification, potentially introducing confounding variables. However, the power  
73 to detect statistically significant loci in admixed populations can be improved by leveraging  
74 admixture-induced local ancestry (Swart et al., 2021; Swart, van Eeden, et al., 2022). Since  
75 previous computational algorithms were not able to include local ancestry as a covariate for  
76 GWASs, the local ancestry allelic adjusted association model (LAAA) was developed to  
77 overcome this limitation (Duan et al., 2018). The LAAA model identifies ancestry-specific  
78 alleles associated with the phenotype by including the minor alleles and the corresponding  
79 ancestry of the minor alleles (obtained by local ancestry inference) as covariates. The LAAA  
80 model has been successfully applied in a cohort of multi-way admixed SAC individuals to  
81 identify novel variants associated with TB susceptibility (Swart et al., 2021; Swart, van  
82 Eeden, et al., 2022).

83

84 Our study builds upon the findings from the ITHGC (Schurz et al., 2024) and aim to resolve  
85 the challenges faced in African ancestry-specific association analysis. Here, we explore  
86 host genetic correlates of TB in a complex admixed SAC population using the LAAA  
87 model.

88

## 89 **Methods**

### 90 *Data*

91 This study included the two SAC admixed datasets from the ITHGC analysis [RSA(A) and  
92 RSA(M)] as well as four additional TB case-control datasets obtained from admixed South  
93 African population groups (Table 1). Like the SAC population, the Xhosa population are  
94 admixed with rain-forest forager and KhoeSan ancestral contributions (Choudhury et al.,  
95 2021). All datasets were collected over the past 30 years under different research projects  
96 (Daya et al., 2013; Kroon et al., 2020; Schurz et al., 2018; Smith et al., 2023; Ugarte-Gil et  
97 al., 2020) and individuals that were included in the analyses consented to the use of their  
98 data in future research regarding TB host genetics. Across all datasets, TB cases were  
99 bacteriologically confirmed (culture positive) or diagnosed by GeneXpert. Controls were

100 healthy individuals with no previous or current history of TB disease or treatment. However,  
101 given the high prevalence of TB in South Africa [852 cases (95% CI 679-1026) per 100 00  
102 individuals 15 years and older (Cudahy et al., 2020)], most controls have likely been exposed  
103 to *M.tb* at some point (Gallant et al., 2010). For all datasets, cases and controls were  
104 obtained from the same community and thus share similar socio-economic status and  
105 health care access.

106

107 **Table 1.** Summary of the datasets included in analysis.

Dataset	Genotyping platform	Self-reported ethnicity	Cases/controls	Reference
RSA(A)	Affymetrix 500k	SAC	642/91	(Daya et al., 2013)
RSA(M)	MEGA array 1.1M	SAC	555/440	(Schurz et al., 2018; Swart et al., 2021)
RSA(TANDEM)	H3Africa array	SAC and Bantu-speaking African	161/133	(Swart, Uren, et al., 2022)
RSA(NCTB)	H3Africa array	SAC	49/111	(Oyageshio et al., 2023)
RSA(Worcester)	H3Africa array	SAC	61 cases	Unpublished
RSA(Xhosa)	Whole genome sequencing	IsiXhosa	44/120	Unpublished

108

109 A list of sites genotyped on the Infinium™ H3Africa array  
110 (<https://chipinfo.h3abionet.org/browse>) were extracted from the whole-genome sequenced  
111 [RSA(Xhosa)] dataset and treated as genotype data in subsequent analyses. Quality control  
112 (QC) of raw genotype data was performed using PLINK v1.9 (Purcell et al., 2007). In all  
113 datasets, individuals were screened for sex concordance and discordant sex information  
114 was corrected based on X chromosome homozygosity estimates ( $F_{\text{estimate}} < 0.2$  for females  
115 and  $F_{\text{estimate}} > 0.8$  for males). In the event that sex information could not be corrected based  
116 on homozygosity estimates, individuals with missing or discordant sex information were  
117 removed. Individuals with genotype call rates less than 90% and SNPs with more than 5%  
118 missingness were removed as described previously (Swart et al., 2021). Monomorphic sites  
119 were removed. Individuals were screened for deviations in Hardy-Weinberg Equilibrium  
120 (HWE) for each SNP and sites deviating from the HWE threshold of  $10^{-5}$  were removed. Sex  
121 chromosomes were excluded from the analysis. The genome coordinates across all datasets  
122 were checked for consistency and, if necessary, converted to GRCh37 using the UCSC  
123 liftOver tool (Kuhn et al., 2013).

124

125 Genotype datasets were pre-phased using SHAPEIT v2 (Delaneau et al., 2013) and imputed  
126 using the Positional Burrows-Wheeler Transformation (PBWT) algorithm through the Sanger  
127 Imputation Server (SIS) (Durbin, 2014). The African Genome Resource (AGR) panel (n=4 956),  
128 accessed via the SIS, was used as the reference panel for imputation (Gurdasani et al., 2015)  
129 since it has been shown that the AGR is the best reference panel for imputation of missing  
130 genotypes for samples from the SAC population (Schurz et al., 2019). Imputed data were  
131 filtered to remove sites with imputation quality INFO scores less than 0.95. Individual  
132 datasets were screened for relatedness using KING software (Manichaikul et al., 2010) and  
133 individuals up to second degree relatedness were removed. A total of 7 544 769 markers  
134 overlapped across all six datasets. This list of intersecting markers was extracted from each  
135 dataset using PLINK --extract flag. The datasets were then merged using the PLINK v1.9. After  
136 merging, all individuals missing more than 10% genotypes were removed, markers with more  
137 than 5% missing data were excluded and a HWE filter was applied to controls (threshold <math>10^{-5}</math>  
138 <sup>5</sup>). The merged dataset was screened for relatedness using KING and individuals up to second  
139 degree relatedness were subsequently removed. The final merged dataset after QC and data  
140 filtering (including the removal of related individuals) consisted of 1 544 individuals (952 TB  
141 cases and 592 healthy controls). A total of 7 510 057 variants passed QC and filtering  
142 parameters.

143

#### 144 *Global ancestry inference*

145 ADMIXTURE was used to determine the correct number of contributing ancestral proportions  
146 in our multi-way admixed population cohort (Alexander & Lange, 2011). ADMIXTURE  
147 estimates the number of contributing ancestral populations (denoted by K) and population  
148 allele frequencies through cross-validation (CV). All 1 544 individuals were grouped into  
149 running groups of equal size together with 191 reference populations (Table 2). Running  
150 groups were created to ensure approximately equal numbers of reference populations and  
151 admixed populations. Xhosa and SAC samples were divided into separate running groups.

152

153 **Table 2.** Ancestral populations included for global ancestry deconvolution.

Population	n	Source
European (British – GBR)	40	1000 Genomes (1000G) phase 3 (1000 Genomes Project Consortium et al., 2015)
East Asian (Chinese – CHB)	40	1000G phase 3

Bantu-speaking African (Yoruba – YRI)	40	1000G phase 3
Southeast Asian (Malaysian)	38	Singapore Sequencing Malay Project (SSMP) (Wong et al., 2013)
KhoeSan (Nama)	33	African Genome Variation Project (AGVP/ADRP) (Gurdasani et al., 2015)

154

155 Redundant SNPs were removed by PLINK through LD pruning by removing each SNP with LD  
156  $r^2 > 0.1$  within a 50-SNP sliding window (advanced by 10 SNPs at a time). Ancestral  
157 proportions were inferred in an unsupervised manner for  $K = 3-6$  (1 iteration). The best value  
158 of  $K$  for the data was selected by choosing the  $K$  value with the lowest CV error across all  
159 running groups. Ten iterations of  $K = 3$  and  $K = 5$  was run for the Xhosa and SAC individuals  
160 respectively. Since it has been shown that RFMix (Maples et al., 2013) outperforms  
161 ADMIXTURE in determining global ancestry proportions (C Uren et al., 2020), RFMix was also  
162 used to refine inferred global ancestry proportions. Global ancestral proportions were  
163 visualised using PONG (Behr et al., 2016).

164

#### 165 *Local ancestry inference*

166 The merged dataset and the reference file (containing reference populations from Table 2)  
167 were phased separately using SHAPEIT2. The local ancestry for each position in the genome  
168 was inferred using RFMix (Maples et al., 2013). Default parameters were used, but the  
169 number of generations since admixture was set to 15 for the SAC individuals and 20 for the  
170 Xhosa individuals (as determined by previous studies) (Caitlin Uren et al., 2016). RFMix was  
171 run with three expectation maximisation iterations and the --reanalyse-reference flag.

172

#### 173 *Batch effect screening and correction*

174 Merging separate datasets generated at different timepoints and/or facilities, as we have  
175 done here, will undoubtedly introduce batch effects. Principal component analysis (PCA) is  
176 a common method used to visualise batch effects, where the first two principal components  
177 (PCs) are plotted with each sample coloured by batch, and a separation of colours is  
178 indicative of a batch effect (Nyamundanda et al., 2017). However, it is difficult to  
179 differentiate between separation caused by population structure and separation caused by  
180 batch effect using PCA alone. An alternative method to detect batch effects (Chen et al.,  
181 2022) involves coding case/control status by batch followed by running an association

182 analysis testing each batch against all other batches. If any single dataset has more positive  
183 signals compared to the other datasets, then batch effects may be responsible for producing  
184 spurious results. Batch effects can be resolved by removing those SNPs which pass the  
185 genome-wide significance threshold from the merged dataset. We have adapted this batch  
186 effect correction method for application in a highly admixed cohort with complex population  
187 structure (Croock et al., 2024). Our modified method was used to remove 36 627 SNPs  
188 affected by batch effects from our merged dataset.

189

#### 190 *Local ancestry allelic adjusted association analysis*

191 The LAAA association model was used to investigate if there are allelic, ancestry-specific or  
192 ancestry-specific allelic associations with TB susceptibility in our merged dataset. Global  
193 ancestral components inferred by RFMix, age and sex were included as covariates in the  
194 association tests. Variants with minor allele frequency (MAF) < 1% were removed to improve  
195 the stability of the association tests. Dosage files, which code the number of alleles of a  
196 specific ancestry at each locus across the genome, were compiled. Separate regression  
197 models for each ancestral contribution were fitted to investigate which ancestral  
198 contribution is associated with TB susceptibility. Details regarding the models have been  
199 described elsewhere (Swart, van Eeden, et al., 2022); but in summary, four regression  
200 models were tested to detect the source of the association signals observed:

201

#### 202 *(1) Null model or global ancestry (GA) model:*

203 The null model only includes global ancestry, sex and age covariates. This test investigates  
204 whether an additive allelic dose exerts an effect on the phenotype (without including local  
205 ancestry of the allele).

206

#### 207 *(2) Local ancestry (LA) model:*

208 This model is used in admixture mapping to identify ancestry-specific variants associated  
209 with a specific phenotype. The LA model evaluates the number of alleles of a specific  
210 ancestry at a locus and includes the corresponding marginal effect as a covariate in  
211 association analyses.

212

#### 213 *(3) Ancestry plus allelic (APA) model:*



214 The APA model simultaneously performs model (1) and (2). This model tests whether an  
215 additive allelic dose exerts an effect of the phenotype whilst adjusting for local ancestry.  
216

#### 217 (4) Local ancestry adjusted allelic (LAAA) model:

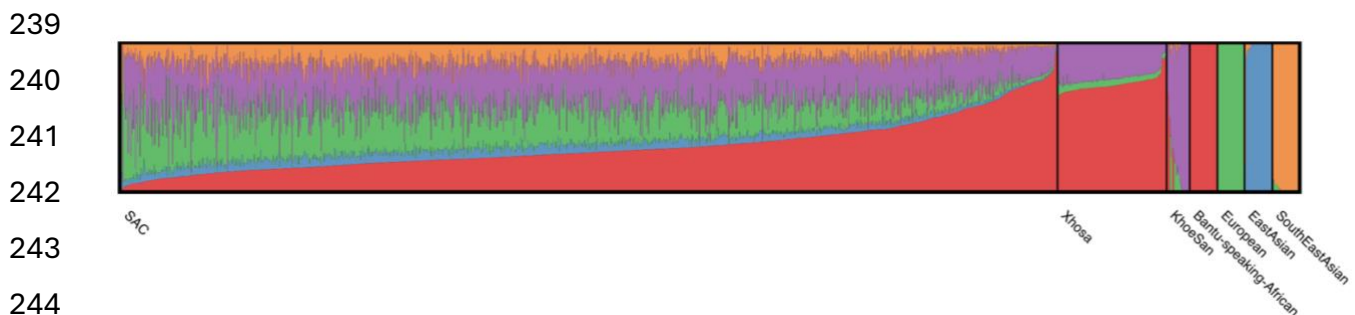
218 The LAAA model is an extension of the APA model, which models the combination of the  
219 minor allele and ancestry of the minor allele at a specific locus and the effect this interaction  
220 has on the phenotype.  
221

222 The R package *STEAM* (Significance Threshold Estimation for Admixture Mapping) (Grinde et  
223 al., 2019) was used to determine the genome-wide significance threshold given the global  
224 ancestral proportions of each individual and the number of generations since admixture ( $g =$   
225 15). *STEAM* permuted these factors 10 000 times to derive a threshold for significance.  
226 Results were visualised in RStudio. A genome-wide significance threshold of  $p$ -value  $< 2.5 \times$   
227  $10^{-6}$  was deemed significant by *STEAM*.  
228

## 229 Results

### 230 Global and local ancestry inference

231 After close inspection of global ancestry proportions generated using ADMIXTURE, the K  
232 number of contributing ancestries (the lowest k-value determined through cross-validation)  
233 was  $K = 3$  for the Xhosa individuals and  $K = 5$  for the SAC individuals (Figure 1). This is  
234 consistent with previous global ancestry deconvolution results (Chimusa et al., 2014;  
235 Choudhury et al., 2021). It is evident that our cohort is a complex, highly admixed group with  
236 ancestral contributions from the indigenous KhoeSan (~22 - 30%), Bantu-speaking African  
237 (~30 - 72%), European (~5 - 24%), Southeast Asian (~11%) and East Asian (~5%) population  
238 groups.  
239



245 **Figure 1.** Genome-wide ancestral proportions of all individuals in the merged dataset. Ancestral proportions for each  
246 individual are plotted vertically with different colours representing different contributing ancestries.



247

248 Local ancestry was estimated for all individuals. Admixture between geographically distinct  
249 populations creates complex ancestral and admixture-induced LD blocks, which can be  
250 visualised using local ancestry karyograms. Figure 2 shows karyograms for three individuals  
251 from the merged dataset. It is evident that, despite individuals being from the same  
252 population group, each possesses unique patterns of local ancestry arising from differing  
253 numbers and lengths of ancestral segments.

254

255

256

257

258

259

260

261

262

263

264

265

266

267

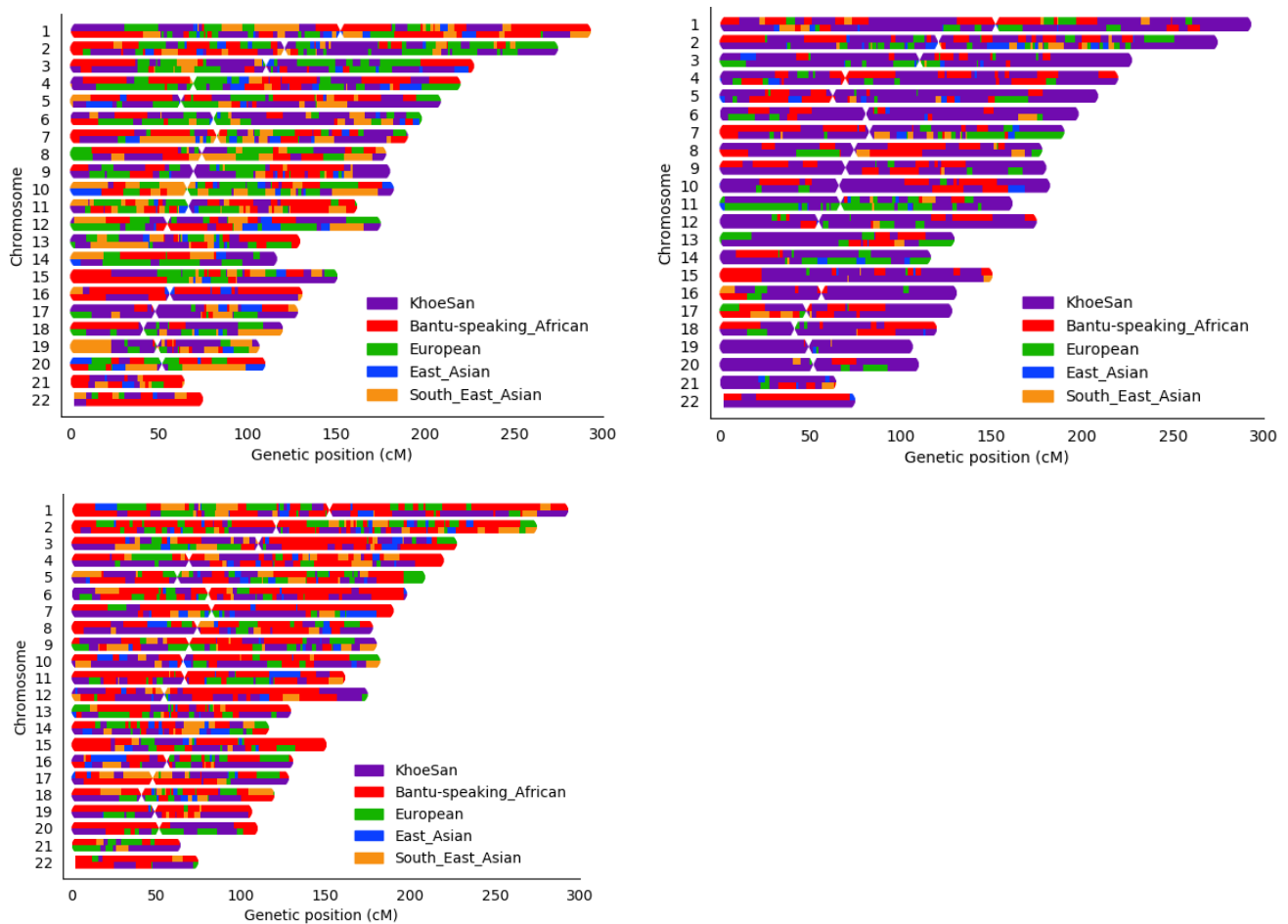
268

269

270

271

272



273 **Figure 2.** Local ancestry karyograms of three admixed individuals from the SAC population. Each admixed individual  
274 has unique local ancestry patterns generated by admixture among geographically distinct ancestral population groups.

275

### 276 *Local ancestry-allelic adjusted analysis*

277 A total of 784 557 autosomal markers (with MAF > 1%) and 1 544 unrelated individuals (952  
278 TB cases and 592 healthy controls) were included in logistic regression models to assess  
279 whether any loci and/or ancestries were significantly associated with TB status (whilst

280 adjusting for sex, age, and global ancestry proportions). LAAA models were successfully  
281 applied for all five contributing ancestries (Khoesan, Bantu-speaking African, European, East  
282 Asian and Southeast Asian). Only one variant (*rs74828248*) was significantly associated with  
283 TB status ( $p$ -value  $< 2.5 \times 10^{-6}$ ) whilst utilising the LAAA model and whilst adjusting for Bantu-  
284 speaking African ancestry on chromosome 20 ( $p$ -value =  $2.272 \times 10^{-6}$ , OR = 0.316, SE = 0.244)  
285 (Supplementary Figure 1). No genomic inflation was detected in the QQ-plot for this region  
286 (Supplementary Figure 2). However, this variant is located in an intergenic region and the link  
287 to TB susceptibility is unclear (Supplementary Figure 3).

288

289 Although no other variants passed the genome-wide significance threshold, multiple lead  
290 SNPs were identified. Notably, an appreciable peak was identified in the HLA-II region of  
291 chromosome 6 when using the LAAA model and adjusting for Khoesan ancestry (Figure 3).  
292 The QQ-plot suggested minimal genomic inflation, which was verified by calculating the  
293 genomic inflation factor ( $\lambda = 1.05289$ ) (Supplementary Figure 4). The lead variants identified  
294 using the LAAA model whilst adjusting for Khoesan ancestry in this region on chromosome 6  
295 are summarised in Table 3. The association peak encompasses the *HLA-DPA1/B1* (major  
296 histocompatibility complex, class II, DP alpha 1/beta 1) genes (Figure 4). It is noteworthy that  
297 without the LAAA model, this association peak would not have been observed for this cohort.  
298 This highlights the importance of utilising the LAAA model in future association studies when  
299 investigating disease susceptibility loci in admixed individuals, such as the SAC population.

300

301

### Khoesan Ancestry

302

303

304

305

306

307

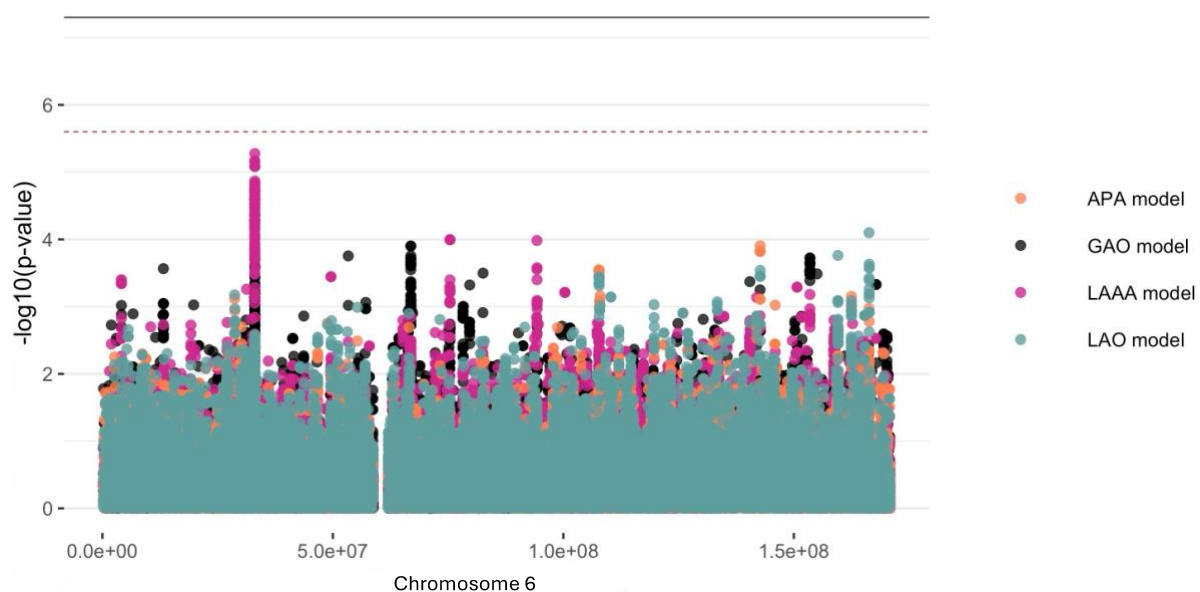
308

309

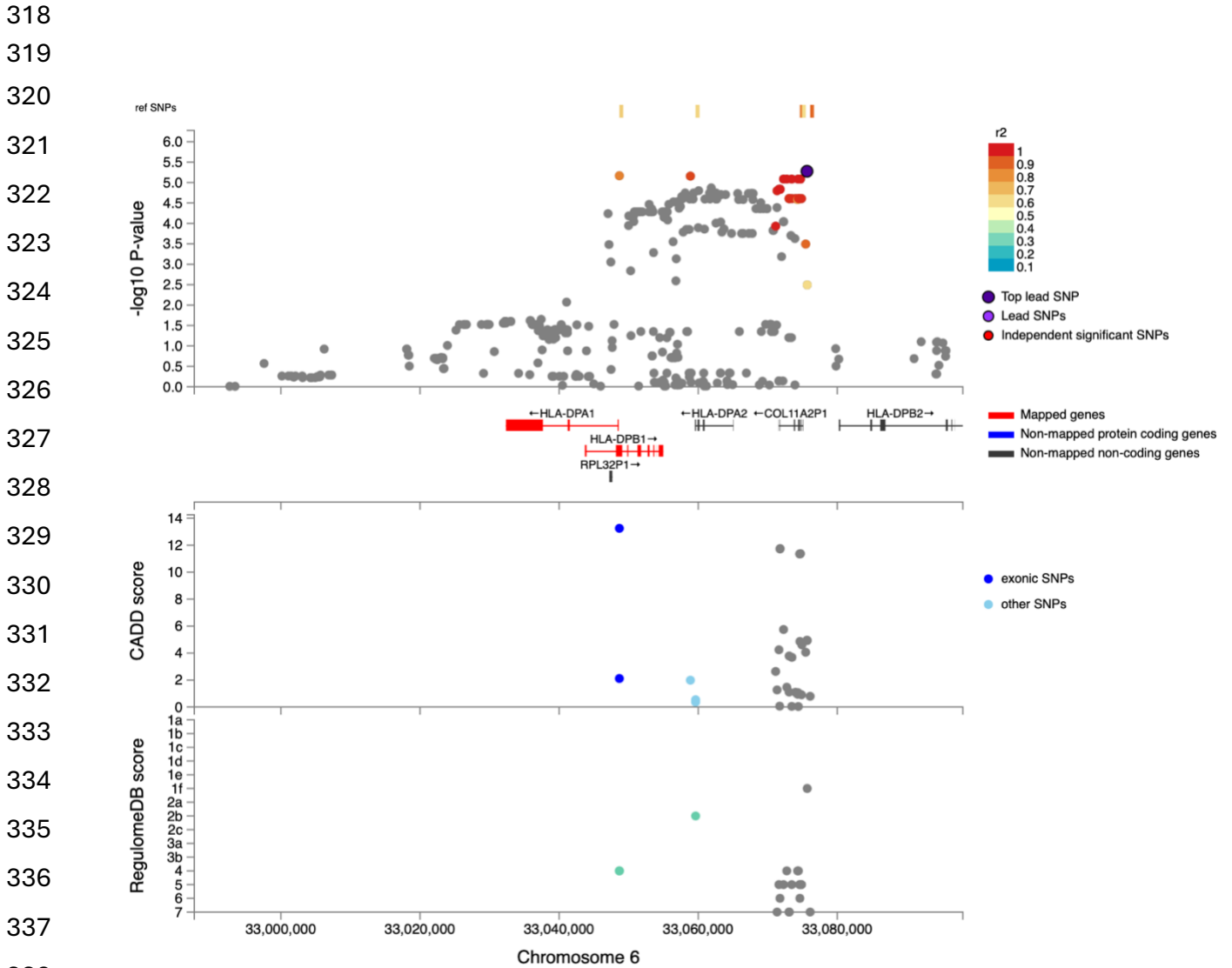
310

311

312



313 **Figure 3.** Log transformation of association signals obtained for KhoeSan ancestry whilst using the LAAA model on  
 314 chromosome 6. The dashed red line represents the significant threshold for admixture mapping calculated with the  
 315 software STEAM ( $p$ -value =  $2.5 \times 10^{-6}$ ) and the black solid line represents the genome wide significant threshold ( $p$ -  
 316 value =  $5 \times 10^{-8}$ ). The four different models are represented in black (global ancestry only - GAO), blue (local ancestry  
 317 effect - LAO), orange (ancestry plus allelic effect - APA) and pink (local ancestry adjusted allelic effect - LAAA).



339 **Figure 4.** Regional plot indicating the nearest genes in the region of the lead variant (*rs3117230*) observed on  
 340 chromosome 6. SNPs in linkage disequilibrium (LD) with the lead variant are coloured red/orange. The lead variant is  
 341 indicated in purple. Functional protein-coding genes are coded in red and non-functional (pseudo-genes) are indicated  
 342 in black.

343

344

345 The lead variant lies within *COL11A2P1* (collagen type X1 alpha 2 pseudogene 1).  
 346 *COL11A2P1* is an unprocessed pseudogene ([ENSG00000228688](https://www.ncbi.nlm.nih.gov/nuccore/ENSG00000228688)). Unprocessed  
 347 pseudogenes are seldomly transcribed and translated into functional proteins (Witek &

348 Mohiuddin, 2024). *HLA-DPB1* and *HLA-DPA1* are the closest functional protein-coding genes  
349 to our lead variants.

350

351 **Table 3.** Suggestive associations ( $p$ -value  $< 1e^{-5}$ ) for the LAAA analysis adjusting for KhoeSan local ancestry on  
352 chromosome 6.

Position	Marker name	Ref	Alt	AltFreq	OR (95% CI)	SE	$p$ -value ( $\times 10^{-6}$ )	Gene	Location	Imputed/typed	INFO score
33075635	<i>rs3117230</i>	A	G	0.370	0.437 (0.306; 0.624)	0.182	5.292	<i>HLA-DPB1</i>	Intergenic	Genotyped	NA
33048661	<i>rs1042151</i>	A	G	0.325	0.437 (0.305; 0.627)	0.184	6.806	<i>HLA-DPB1</i>	Exonic	Imputed	0.992
33058874	<i>rs2179920</i>	C	T	0.369	0.445 (0.313; 0.633)	0.180	6.960	<i>HLA-DPB1</i>	Intergenic	Genotyped	NA
33072266	<i>rs2064478</i>	C	T	0.371	0.447 (0.313; 0.637)	0.181	8.222	<i>HLA-DPB1</i>	Intergenic	Imputed	1
33072729	<i>rs3130210</i>	G	T	0.371	0.447 (0.313; 0.637)	0.181	8.222	<i>HLA-DPB1</i>	Intergenic	Imputed	0.999
33073440	<i>rs2064475</i>	G	A	0.371	0.447 (0.313; 0.637)	0.181	8.222	<i>HLA-DPB1</i>	Intergenic	Imputed	1
33074348	<i>rs3117233</i>	T	C	0.371	0.447 (0.313; 0.637)	0.181	8.222	<i>HLA-DPB1</i>	Intergenic	Imputed	1
33074707	<i>rs3130213</i>	G	A	0.371	0.447 (0.313; 0.637)	0.181	8.222	<i>HLA-DPB1</i>	Intergenic	Imputed	0.970

353 Ref, reference allele; Alt, alternate allele; AltFreq, alternate allele frequency; OR, odds ratio; SE, standard error

354

355 The lead variant identified in the ITHGC meta-analysis, *rs28383206*, was not present in our  
356 genotype or imputed datasets. The ITHGC imputed genotypes using the 1000 Genomes  
357 (1000G) reference panel (Schurz et al., 2024). Variant *rs28383206* has an alternate allele  
358 frequency of 11.26% in the African population subgroup within the 1000G dataset  
359 (<https://www.ncbi.nlm.nih.gov/snp/rs28383206>). However, *rs28383206* is absent from our  
360 in-house whole-genome sequencing (WGS) datasets, which include Bantu-speaking African  
361 and KhoeSan individuals. This absence suggests that *rs28383206* might not have been  
362 imputed in our datasets using the AGR reference panel, potentially due to its low alternate  
363 allele frequency in southern African populations. Our merged dataset contained two variants  
364 located within 800 base pairs of *rs28383206*: *rs482205* (6:32576009) and *rs482162*  
365 (6:32576019). However, these variants were not significantly associated with TB status in our  
366 cohort (Supplementary Table 1).

## 367 Discussion

368 The LAAA analysis of host genetic susceptibility to TB, involving 942 TB cases and 592  
369 controls, identified one suggestive association peak adjusting for KhoeSan local ancestry.  
370 The association peak identified in this study encompasses the *HLA-DPB1* gene, a highly  
371 polymorphic locus, with over 2 000 documented allelic variants (Robinson et al., 2020). This  
372 association is noteworthy given that *HLA-DPB1* alleles have been associated with TB  
373 resistance (Dawkins et al., 2022; Ravikumar et al., 1999; Selvaraj et al., 2008). The  
374 direction of effect the lead variants in our study (Table 3) similarly suggest a protective effect  
375 against developing active TB. However, variants in *HLA-DPB1* were not identified in the ITHGC  
376 meta-analysis.

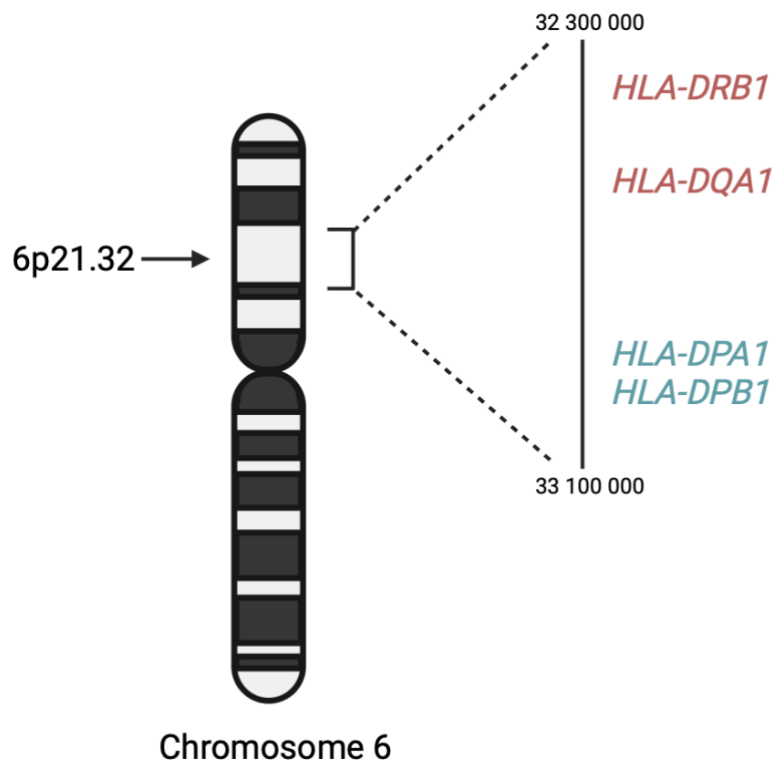
377

378 Population stratification arising from the highly heterogeneous admixed cohorts might have  
379 masked this association signal in the African ancestry-specific association analysis. The  
380 association peak in the HLA-II region was only captured using the LAAA model whilst  
381 adjusting for KhoeSan local ancestry. This underscores the importance of incorporating  
382 global and local ancestry in association studies investigating complex multi-way admixed  
383 individuals, as the genetic heterogeneity present in admixed individuals (produced as a result  
384 of admixture-induced and ancestral LD patterns) may cause association signals to be missed  
385 when using traditional association models (Duan et al., 2018; Swart, van Eeden, et al.,  
386 2022).

387

388 We did not replicate the significant association signal in *HLA-DRB1* identified by the ITHGC.  
389 However, the ITHGC also did not replicate this association in their own African ancestry-  
390 specific analysis. The significant association, *rs28383206*, identified by the ITHGC appears  
391 to be tagging the *HLA-DQA1\*02:1* allele, which is associated with TB in Icelandic and Asian  
392 populations (Li et al., 2021; Sveinbjornsson et al., 2016; Zheng et al., 2018). It is possible  
393 that this association signal is specific to non-African populations, but additional research is  
394 required to verify this hypothesis. Both our study and the ITHGC independently pinpointed  
395 variants associated with TB susceptibility in different genes within the HLA-II locus (Figure 5).  
396 The HLA-II region spans ~0.8Mb on chromosome 6p21.32 and encompasses the *HLA-DP*, -  
397 *DR* and *-DQ* alpha and beta chain genes. The HLA-II complex is the human form of the major  
398 histocompatibility complex class II (MHC-II) proteins on the surface of antigen presenting  
399 cells, such as monocytes, dendritic cells and macrophages. The innate immune response

400 against *M.tb* involves phagocytosis by alveolar macrophages. In the phagosome,  
401 mycobacterial antigens are processed for presentation on MHC-II on the surface of the  
402 antigen presenting cell. Previous studies have suggested that *M.tb* interferes with the MHC-II  
403 pathway to enhance intracellular persistence and delay activation of the adaptive immune  
404 response (Oliveira-Cortez et al., 2016). For example, *M.tb* can inhibit phagosome maturation  
405 and acidification, thereby limiting antigen processing and presentation on MHC-II molecules  
406 (Chang et al., 2005). Given that MHC-II plays an essential role in the adaptive immune  
407 response to TB and numerous studies have identified HLA-II variants associated with TB (Cai  
408 et al., 2019; Chihab et al., 2023; de Sá et al., 2020; Harishankar et al., 2018; Schurz et al.,  
409 2024; Selvaraj et al., 2008), additional research is required to elucidate the effects of HLA-II  
410 variation on TB risk status.



425

426

427 **Figure 5.** A schematic diagram the location of HLA-II genes associated with TB susceptibility. Genes in red were  
428 identified by the ITHGC. Genes in blue were identified by this study.

429

430 This analysis has a few limitations. First, unlike the ITHGC manuscript, we did not validate  
431 our SNP peak in the HLA-II region through fine mapping. Although we initially considered  
432 performing HLA imputation and fine-mapping using the HIBAG R package, as described in

433 the ITHGC article ([https://hibag.s3.amazonaws.com/hlares\\_index.html#estimates](https://hibag.s3.amazonaws.com/hlares_index.html#estimates)), the  
434 African HIBAG model was trained on genotype data from African American and HapMap YRI  
435 populations, which have minimal to no KhoeSan ancestry. Since our association peak likely  
436 originates from KhoeSan ancestral haplotype blocks, using an imputation reference panel  
437 that includes individuals with KhoeSan ancestry is essential to this analysis. We  
438 acknowledge that HLA typing could validate the importance of our lead SNPs in the HLA-II  
439 region and support the LAAA model, but this was not feasible due to the absence of a suitable  
440 reference panel that includes KhoeSan ancestry. Second, our analysis has a notable case-  
441 control imbalance (cases/controls = 1.610). While many studies discuss methods for  
442 addressing case-control imbalances with more controls than cases (which can inflate type 1  
443 error rates (Dai et al., 2021; Öztornaci et al., 2023; Zhou et al., 2018), few address the  
444 implications of a large case-to-control ratio like ours (952 cases to 592 controls). To assess  
445 the impact of this imbalance, we used the Michigan genetic association study (GAS) power  
446 calculator (Skol et al., 2006). Under an additive disease model with an estimated prevalence  
447 of 0.15, a disease allele frequency of 0.3, a genotype relative risk of 1.5, and a default  
448 significance level of  $7 \times 10^{-6}$ , we achieved an expected power of approximately 75%. With a  
449 balanced sample size of 950 cases and 950 controls, power would exceed 90%, but it would  
450 drop significantly with a smaller balanced cohort of 590 cases and 590 controls. Given these  
451 results, we proceeded with our analysis to maximize statistical power despite the case-  
452 control imbalance.

453

454 In conclusion, application of the LAAA to a highly admixed SAC cohort revealed a suggestive  
455 association signal in the HLA-II region associated with protection against TB. Our study builds  
456 on the results of the ITHGC by demonstrating an alternative method to identify association  
457 signals in cohorts with complex genetic ancestry. This analysis shows the value of including  
458 individual global and local ancestry in genetic association analyses. Furthermore, we  
459 confirm HLA-II loci associations with TB susceptibility in an admixed South African  
460 population and hope that this publication will encourage greater appreciation for the role of  
461 the adaptive immune system in TB susceptibility and resistance.

462

### 463 **Acknowledgements**

464 We acknowledge the support of the DSI-NRF Centre of Excellence for Biomedical  
465 Tuberculosis Research, South African Medical Research Council Centre for Tuberculosis



466 Research (SAMRC CTR), Division of Molecular Biology and Human Genetics, Faculty of  
467 Medicine and Health Sciences, Stellenbosch University, Cape Town, South Africa. We also  
468 acknowledge the Centre for High Performance Computing (CHPC), South Africa, for  
469 providing computational resources. This research was partially funded by the South African  
470 government through the SAMRC and the Harry Crossley Research Foundation.

471

#### 472 **Author ORCIDs**

473 Dayna Croock: 0000-0002-5107-8006

474 Yolandi Swart: 0000-0002-9840-3646

475 Haiko Schurz: 0000-0002-0009-3409

476 Desiree C. Petersen: 0000-0002-0817-2574

477 Marlo Möller: 0000-0002-0805-6741

478 Caitlin Uren: 0000-0003-2358-0135

479

#### 480 **Ethics**

481 Ethics approval was granted by the Health Research Ethics Committee (HREC) of  
482 Stellenbosch University, South Africa (project number S22/02/031).

483

#### 484 **Competing interests**

485 None declared.

486

#### 487 **References**

488 1000 Genomes Project Consortium, Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P.,

489 Kang, H. M., Korbel, J. O., Marchini, J. L., McCarthy, S., McVean, G. A., & Abecasis, G. R.

490 (2015). A global reference for human genetic variation. *Nature*, 526(7571), 68–74.

491 <https://doi.org/10.1038/nature15393>

492 Alexander, D. H., & Lange, K. (2011). Enhancements to the ADMIXTURE algorithm for  
493 individual ancestry estimation. *BMC Bioinformatics*, 12, 246.

494 <https://doi.org/10.1186/1471-2105-12-246>

- 495 Behr, A. A., Liu, K. Z., Liu-Fang, G., Nakka, P., & Ramachandran, S. (2016). pong: fast  
496 analysis and visualization of latent clusters in population genetic data. *Bioinformatics*,  
497 32(18), 2817–2823. <https://doi.org/10.1093/bioinformatics/btw327>
- 498 Cai, L., Li, Z., Guan, X., Cai, K., Wang, L., Liu, J., & Tong, Y. (2019). The research progress of  
499 host genes and tuberculosis susceptibility. *Oxidative Medicine and Cellular Longevity*,  
500 2019, 9273056. <https://doi.org/10.1155/2019/9273056>
- 501 Chang, S. T., Linderman, J. J., & Kirschner, D. E. (2005). Multiple mechanisms allow  
502 Mycobacterium tuberculosis to continuously inhibit MHC class II-mediated antigen  
503 presentation by macrophages. *Proceedings of the National Academy of Sciences of the*  
504 *United States of America*, 102(12), 4530–4535.  
505 <https://doi.org/10.1073/pnas.0500362102>
- 506 Chen, D., Tashman, K., Palmer, D. S., Neale, B., Roeder, K., Bloemendal, A., Churchhouse,  
507 C., & Ke, Z. T. (2022). A data harmonization pipeline to leverage external controls and  
508 boost power in GWAS. *Human Molecular Genetics*, 31(3), 481–489.  
509 <https://doi.org/10.1093/hmg/ddab261>
- 510 Chihab, L. Y., Kuan, R., Phillips, E. J., Mallal, S. A., Rozot, V., Davis, M. M., Scriba, T. J.,  
511 Sette, A., Peters, B., Lindestam Arlehamn, C. S., & SATVI Study Group. (2023). Expression  
512 of specific HLA class II alleles is associated with an increased risk for active tuberculosis  
513 and a distinct gene expression profile. *HLA : Immune Response Genetics*, 101(2), 124–  
514 137. <https://doi.org/10.1111/tan.14880>
- 515 Chimusa, E. R., Daya, M., Möller, M., Ramesar, R., Henn, B. M., van Helden, P. D., Mulder,  
516 N. J., & Hoal, E. G. (2013). Determining ancestry proportions in complex admixture

- 517 scenarios in South Africa using a novel proxy ancestry selection method. *Plos One*, 8(9),  
518 e73971. <https://doi.org/10.1371/journal.pone.0073971>
- 519 Chimusa, E. R., Zaitlen, N., Daya, M., Möller, M., van Helden, P. D., Mulder, N. J., Price, A.  
520 L., & Hoal, E. G. (2014). Genome-wide association study of ancestry-specific TB risk in  
521 the South African Coloured population. *Human Molecular Genetics*, 23(3), 796–809.  
522 <https://doi.org/10.1093/hmg/ddt462>
- 523 Choudhury, A., Sengupta, D., Ramsay, M., & Schlebusch, C. (2021). Bantu-speaker  
524 migration and admixture in southern Africa. *Human Molecular Genetics*, 30(R1), R56–  
525 R63. <https://doi.org/10.1093/hmg/ddaa274>
- 526 Cudahy, P. G. T., Wilson, D., & Cohen, T. (2020). Risk factors for recurrent tuberculosis after  
527 successful treatment in a high burden setting: a cohort study. *BMC Infectious Diseases*,  
528 20(1), 789. <https://doi.org/10.1186/s12879-020-05515-4>
- 529 Dai, X., Fu, G., Zhao, S., & Zeng, Y. (2021). Statistical Learning Methods Applicable to  
530 Genome-Wide Association Studies on Unbalanced Case-Control Disease Data. *Genes*,  
531 12(5). <https://doi.org/10.3390/genes12050736>
- 532 Dawkins, B. A., Garman, L., Cejda, N., Pezant, N., Rasmussen, A., Rybicki, B. A., Levin, A.  
533 M., Benchek, P., Seshadri, C., Mayanja-Kizza, H., Iannuzzi, M. C., Stein, C. M., &  
534 Montgomery, C. G. (2022). Novel HLA associations with outcomes of Mycobacterium  
535 tuberculosis exposure and sarcoidosis in individuals of African ancestry using nearest-  
536 neighbor feature selection. *Genetic Epidemiology*, 46(7), 463–474.  
537 <https://doi.org/10.1002/gepi.22490>
- 538 Daya, M., van der Merwe, L., Galal, U., Möller, M., Salie, M., Chimusa, E. R., Galanter, J. M.,  
539 van Helden, P. D., Henn, B. M., Gignoux, C. R., & Hoal, E. (2013). A panel of ancestry

540 informative markers for the complex five-way admixed South African coloured  
541 population. *Plos One*, 8(12), e82224. <https://doi.org/10.1371/journal.pone.0082224>

542 de Sá, N. B. R., Ribeiro-Alves, M., da Silva, T. P., Pilotto, J. H., Rolla, V. C., Giacoia-Gripp, C.  
543 B. W., Scott-Algara, D., Morgado, M. G., & Teixeira, S. L. M. (2020). Clinical and genetic  
544 markers associated with tuberculosis, HIV-1 infection, and TB/HIV-immune  
545 reconstitution inflammatory syndrome outcomes. *BMC Infectious Diseases*, 20(1), 59.  
546 <https://doi.org/10.1186/s12879-020-4786-5>

547 Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F., & Marchini, J. (2013). Haplotype estimation  
548 using sequencing reads. *American Journal of Human Genetics*, 93(4), 687–696.  
549 <https://doi.org/10.1016/j.ajhg.2013.09.002>

550 Duan, Q., Xu, Z., Raffield, L. M., Chang, S., Wu, D., Lange, E. M., Reiner, A. P., & Li, Y. (2018).  
551 A robust and powerful two-step testing procedure for local ancestry adjusted allelic  
552 association analysis in admixed populations. *Genetic Epidemiology*, 42(3), 288–302.  
553 <https://doi.org/10.1002/gepi.22104>

554 Durbin, R. (2014). Efficient haplotype matching and storage using the positional Burrows-  
555 Wheeler transform (PBWT). *Bioinformatics*, 30(9), 1266–1272.  
556 <https://doi.org/10.1093/bioinformatics/btu014>

557 Escombe, A. R., Ticona, E., Chávez-Pérez, V., Espinoza, M., & Moore, D. A. J. (2019).  
558 Improving natural ventilation in hospital waiting and consulting rooms to reduce  
559 nosocomial tuberculosis transmission risk in a low resource setting. *BMC Infectious  
560 Diseases*, 19(1), 88. <https://doi.org/10.1186/s12879-019-3717-9>

561 Gallant, C. J., Cobat, A., Simkin, L., Black, G. F., Stanley, K., Hughes, J., Doherty, T. M.,  
562 Hanekom, W. A., Eley, B., Beyers, N., Jaïs, J. P., van Helden, P., Abel, L., Alcaïs, A., Hoal,

- 563 E. G., & Schurr, E. (2010). Impact of age and sex on mycobacterial immunity in an area of  
564 high tuberculosis incidence. *The International Journal of Tuberculosis and Lung Disease*,  
565 14(8), 952–959.
- 566 Glaziou, P., Floyd, K., & Raviglione, M. C. (2018). Global epidemiology of tuberculosis.  
567 *Seminars in Respiratory and Critical Care Medicine*, 39(3), 271–285.  
568 <https://doi.org/10.1055/s-0038-1651492>
- 569 Grinde, K. E., Brown, L. A., Reiner, A. P., Thornton, T. A., & Browning, S. R. (2019). Genome-  
570 wide Significance Thresholds for Admixture Mapping Studies. *American Journal of*  
571 *Human Genetics*, 104(3), 454–465. <https://doi.org/10.1016/j.ajhg.2019.01.008>
- 572 Gurdasani, D., Carstensen, T., Tekola-Ayele, F., Pagani, L., Tachmazidou, I., Hatzikotoulas,  
573 K., Karthikeyan, S., Iles, L., Pollard, M. O., Choudhury, A., Ritchie, G. R. S., Xue, Y., Asimit,  
574 J., Nsubuga, R. N., Young, E. H., Pomilla, C., Kivinen, K., Rockett, K., Kamali, A., ...  
575 Sandhu, M. S. (2015). The African Genome Variation Project shapes medical genetics in  
576 Africa. *Nature*, 517(7534), 327–332. <https://doi.org/10.1038/nature13997>
- 577 Harishankar, M., Selvaraj, P., & Bethunaickan, R. (2018). Influence of genetic polymorphism  
578 towards pulmonary tuberculosis susceptibility. *Frontiers in Medicine*, 5, 213.  
579 <https://doi.org/10.3389/fmed.2018.00213>
- 580 Kroon, E. E., Kinnear, C. J., Orlova, M., Fischinger, S., Shin, S., Boolay, S., Walzl, G., Jacobs,  
581 A., Wilkinson, R. J., Alter, G., Schurr, E., Hoal, E. G., & Möller, M. (2020). An observational  
582 study identifying highly tuberculosis-exposed, HIV-1-positive but persistently TB,  
583 tuberculin and IGRA negative persons with M. tuberculosis specific antibodies in Cape  
584 Town, South Africa. *EBioMedicine*, 61, 103053.  
585 <https://doi.org/10.1016/j.ebiom.2020.103053>

- 586 Kuhn, R. M., Haussler, D., & Kent, W. J. (2013). The UCSC genome browser and associated  
587 tools. *Briefings in Bioinformatics*, 14(2), 144–161. <https://doi.org/10.1093/bib/bbs038>
- 588 Laghari, M., Sulaiman, S. A. S., Khan, A. H., Talpur, B. A., Bhatti, Z., & Memon, N. (2019).  
589 Contact screening and risk factors for TB among the household contact of children with  
590 active TB: a way to find source case and new TB cases. *BMC Public Health*, 19(1), 1274.  
591 <https://doi.org/10.1186/s12889-019-7597-0>
- 592 Lehohla, P. (2012). *South African Census 2011 Meta-data* (Report No. 03-01-47; p. 130).  
593 South African Census.
- 594 Li, M., Hu, Y., Zhao, B., Chen, L., Huang, H., Huai, C., Zhang, X., Zhang, J., Zhou, W., Shen,  
595 L., Zhen, Q., Li, B., Wang, W., He, L., & Qin, S. (2021). A next generation sequencing  
596 combined genome-wide association study identifies novel tuberculosis susceptibility  
597 loci in Chinese population. *Genomics*, 113(4), 2377–2384.  
598 <https://doi.org/10.1016/j.ygeno.2021.05.035>
- 599 Manichaikul, A., Mychaleckyj, J. C., Rich, S. S., Daly, K., Sale, M., & Chen, W.-M. (2010).  
600 Robust relationship inference in genome-wide association studies. *Bioinformatics*,  
601 26(22), 2867–2873. <https://doi.org/10.1093/bioinformatics/btq559>
- 602 Maples, B. K., Gravel, S., Kenny, E. E., & Bustamante, C. D. (2013). RFMix: a discriminative  
603 modeling approach for rapid and robust local-ancestry inference. *American Journal of*  
604 *Human Genetics*, 93(2), 278–288. <https://doi.org/10.1016/j.ajhg.2013.06.020>
- 605 Matose, M., Poluta, M., & Douglas, T. S. (2019). Natural ventilation as a means of airborne  
606 tuberculosis infection control in minibus taxis. *South African Journal of Science*,  
607 115(9/10). <https://doi.org/10.17159/sajs.2019/5737>

- 608 Menzies, N. A., Swartwood, N., Testa, C., Malyuta, Y., Hill, A. N., Marks, S. M., Cohen, T., &  
609 Salomon, J. A. (2021). Time Since Infection and Risks of Future Disease for Individuals  
610 with Mycobacterium tuberculosis Infection in the United States. *Epidemiology*, 32(1), 70–  
611 78. <https://doi.org/10.1097/EDE.0000000000001271>
- 612 Möller, M., Kinnear, C. J., Orlova, M., Kroon, E. E., van Helden, P. D., Schurr, E., & Hoal, E. G.  
613 (2018). Genetic Resistance to Mycobacterium tuberculosis Infection and Disease.  
614 *Frontiers in Immunology*, 9, 2219. <https://doi.org/10.3389/fimmu.2018.02219>
- 615 Möller, M., & Kinnear, C. J. (2020). Human global and population-specific genetic  
616 susceptibility to Mycobacterium tuberculosis infection and disease. *Current Opinion in*  
617 *Pulmonary Medicine*, 26(3), 302–310. <https://doi.org/10.1097/MCP.0000000000000672>
- 618 Nyamundanda, G., Poudel, P., Patil, Y., & Sadanandam, A. (2017). A novel statistical  
619 method to diagnose, quantify and correct batch effects in genomic studies. *Scientific*  
620 *Reports*, 7(1), 10849. <https://doi.org/10.1038/s41598-017-11110-6>
- 621 Oliveira-Cortez, A., Melo, A. C., Chaves, V. E., Condino-Neto, A., & Camargos, P. (2016). Do  
622 HLA class II genes protect against pulmonary tuberculosis? A systematic review and  
623 meta-analysis. *European Journal of Clinical Microbiology & Infectious Diseases*, 35(10),  
624 1567–1580. <https://doi.org/10.1007/s10096-016-2713-x>
- 625 Oyageshio, O. P., Myrick, J. W., Saayman, J., van der Westhuizen, L., Al-Hindi, D., Reynolds,  
626 A. W., Zaitlen, N., Uren, C., Möller, M., & Henn, B. M. (2023). Strong effect of demographic  
627 changes on tuberculosis susceptibility in south africa. *MedRxiv*.  
628 <https://doi.org/10.1101/2023.11.02.23297990>
- 629 Öztornaci, R. O., Syed, H., Morris, A. P., & Taşdelen, B. (2023). The use of class imbalanced  
630 learning methods on ULSAM data to predict the case–control status in genome-wide



- 631 association studies. *Journal of Big Data*, 10(1), 174. [https://doi.org/10.1186/s40537-023-](https://doi.org/10.1186/s40537-023-00853-x)  
632 00853-x
- 633 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A. R., Bender, D., Maller, J.,  
634 Sklar, P., de Bakker, P. I. W., Daly, M. J., & Sham, P. C. (2007). PLINK: a tool set for whole-  
635 genome association and population-based linkage analyses. *American Journal of Human*  
636 *Genetics*, 81(3), 559–575. <https://doi.org/10.1086/519795>
- 637 Ravikumar, M., Dheenadhayalan, V., Rajaram, K., Lakshmi, S. S., Kumaran, P. P.,  
638 Paramasivan, C. N., Balakrishnan, K., & Pitchappan, R. M. (1999). Associations of HLA-  
639 DRB1, DQB1 and DPB1 alleles with pulmonary tuberculosis in south India. *Tubercle and*  
640 *Lung Disease : The Official Journal of the International Union against Tuberculosis and*  
641 *Lung Disease*, 79(5), 309–317. <https://doi.org/10.1054/tuld.1999.0213>
- 642 Robinson, J., Barker, D. J., Georgiou, X., Cooper, M. A., Flicek, P., & Marsh, S. G. E. (2020).  
643 IPD-IMGT/HLA Database. *Nucleic Acids Research*, 48(D1), D948–D955.  
644 <https://doi.org/10.1093/nar/gkz950>
- 645 Schurz, H., Kinnear, C. J., Gignoux, C., Wojcik, G., van Helden, P. D., Tromp, G., Henn, B.,  
646 Hoal, E. G., & Möller, M. (2018). A Sex-Stratified Genome-Wide Association Study of  
647 Tuberculosis Using a Multi-Ethnic Genotyping Array. *Frontiers in Genetics*, 9, 678.  
648 <https://doi.org/10.3389/fgene.2018.00678>
- 649 Schurz, H., Müller, S. J., van Helden, P. D., Tromp, G., Hoal, E. G., Kinnear, C. J., & Möller, M.  
650 (2019). Evaluating the Accuracy of Imputation Methods in a Five-Way Admixed  
651 Population. *Frontiers in Genetics*, 10, 34. <https://doi.org/10.3389/fgene.2019.00034>
- 652 Schurz, H., Naranbhai, V., Yates, T. A., Gilchrist, J. J., Parks, T., Dodd, P. J., Möller, M., Hoal,  
653 E. G., Morris, A. P., Hill, A. V. S., & International Tuberculosis Host Genetics Consortium.

- 654 (2024). Multi-ancestry meta-analysis of host genetic susceptibility to tuberculosis  
655 identifies shared genetic architecture. *ELife*, 13. <https://doi.org/10.7554/eLife.84394>
- 656 Selvaraj, P., Raghavan, S., Swaminathan, S., Alagarasu, K., Narendran, G., & Narayanan, P.  
657 R. (2008). HLA-DQB1 and -DPB1 allele profile in HIV infected patients with and without  
658 pulmonary tuberculosis of south India. *Infection, Genetics and Evolution*, 8(5), 664–671.  
659 <https://doi.org/10.1016/j.meegid.2008.06.005>
- 660 Skol, A. D., Scott, L. J., Abecasis, G. R., & Boehnke, M. (2006). Joint analysis is more  
661 efficient than replication-based analysis for two-stage genome-wide association studies.  
662 *Nature Genetics*, 38(2), 209–213. <https://doi.org/10.1038/ng1706>
- 663 Smith, M. H., Myrick, J. W., Oyageshio, O., Uren, C., Saayman, J., Boolay, S., van der  
664 Westhuizen, L., Werely, C., Möller, M., Henn, B. M., & Reynolds, A. W. (2023).  
665 Epidemiological correlates of overweight and obesity in the Northern Cape Province,  
666 South Africa. *PeerJ*, 11, e14723. <https://doi.org/10.7717/peerj.14723>
- 667 Sveinbjornsson, G., Gudbjartsson, D. F., Halldorsson, B. V., Kristinsson, K. G.,  
668 Gottfredsson, M., Barrett, J. C., Gudmundsson, L. J., Blondal, K., Gylfason, A.,  
669 Gudjonsson, S. A., Helgadóttir, H. T., Jonasdóttir, A., Jonasdóttir, A., Karason, A.,  
670 Kardum, L. B., Knežević, J., Kristjansson, H., Kristjansson, M., Love, A., ... Stefansson, K.  
671 (2016). HLA class II sequence variants influence tuberculosis risk in populations of  
672 European ancestry. *Nature Genetics*, 48(3), 318–322. <https://doi.org/10.1038/ng.3498>
- 673 Swart, Y., Uren, C., Eckold, C., Cliff, J. M., Malherbe, S. T., Ronacher, K., Kumar, V.,  
674 Wijmenga, C., Dockrell, H. M., van Crevel, R., Walzl, G., Kleynhans, L., & Möller, M.  
675 (2022). *cis* -eQTL mapping of TB-T2D comorbidity elucidates the involvement of African  
676 ancestry in TB susceptibility. *BioRxiv*. <https://doi.org/10.1101/2022.10.19.512814>

- 677 Swart, Y., Uren, C., van Helden, P. D., Hoal, E. G., & Möller, M. (2021). Local ancestry  
678 adjusted allelic association analysis robustly captures tuberculosis susceptibility loci.  
679 *Frontiers in Genetics*, 12, 716558. <https://doi.org/10.3389/fgene.2021.716558>
- 680 Swart, Y., van Eeden, G., Sparks, A., Uren, C., & Möller, M. (2020). Prospective avenues for  
681 human population genomics and disease mapping in southern Africa. *Molecular*  
682 *Genetics and Genomics*, 295(5), 1079–1089. [https://doi.org/10.1007/s00438-020-01684-](https://doi.org/10.1007/s00438-020-01684-8)  
683 8
- 684 Swart, Y., van Eeden, G., Uren, C., van der Spuy, G., Tromp, G., & Moller, M. (2022). GWAS  
685 in the southern African context. *Cold Spring Harbor Laboratory*.  
686 <https://doi.org/10.1101/2022.02.16.480704>
- 687 Ugarte-Gil, C., Alisjahbana, B., Ronacher, K., Riza, A. L., Koesoemadinata, R. C., Malherbe,  
688 S. T., Cioboata, R., Llontop, J. C., Kleynhans, L., Lopez, S., Santoso, P., Marius, C.,  
689 Villaizan, K., Ruslami, R., Walzl, G., Panduru, N. M., Dockrell, H. M., Hill, P. C., Mc  
690 Allister, S., ... van Crevel, R. (2020). Diabetes Mellitus Among Pulmonary Tuberculosis  
691 Patients From 4 Tuberculosis-endemic Countries: The TANDEM Study. *Clinical Infectious*  
692 *Diseases*, 70(5), 780–788. <https://doi.org/10.1093/cid/ciz284>
- 693 Uren, C, Hoal, E. G., & Möller, M. (2020). Putting RFMix and ADMIXTURE to the test in a  
694 complex admixed population. *BMC Genetics*, 21(1), 40. [https://doi.org/10.1186/s12863-](https://doi.org/10.1186/s12863-020-00845-3)  
695 020-00845-3
- 696 Uren, Caitlin, Henn, B. M., Franke, A., Wittig, M., van Helden, P. D., Hoal, E. G., & Möller, M.  
697 (2017). A post-GWAS analysis of predicted regulatory variants and tuberculosis  
698 susceptibility. *Plos One*, 12(4), e0174738. <https://doi.org/10.1371/journal.pone.0174738>

- 699 Uren, Caitlin, Hoal, E. G., & Möller, M. (2021). Mycobacterium tuberculosis complex and  
700 human coadaptation: a two-way street complicating host susceptibility to TB. *Human*  
701 *Molecular Genetics*, 30(R1), R146–R153. <https://doi.org/10.1093/hmg/ddaa254>
- 702 Uren, Caitlin, Kim, M., Martin, A. R., Bobo, D., Gignoux, C. R., van Helden, P. D., Möller, M.,  
703 Hoal, E. G., & Henn, B. M. (2016). Fine-Scale Human Population Structure in Southern  
704 Africa Reflects Ecogeographic Boundaries. *Genetics*, 204(1), 303–314.  
705 <https://doi.org/10.1534/genetics.116.187369>
- 706 Verhein, K. C., Vellers, H. L., & Kleeberger, S. R. (2018). Inter-individual variation in health  
707 and disease associated with pulmonary infectious agents. *Mammalian Genome*, 29(1–2),  
708 38–47. <https://doi.org/10.1007/s00335-018-9733-z>
- 709 Witek, J., & Mohiuddin, S. S. (2024). Biochemistry, Pseudogenes. In *StatPearls*. StatPearls  
710 Publishing.
- 711 Wong, L.-P., Ong, R. T.-H., Poh, W.-T., Liu, X., Chen, P., Li, R., Lam, K. K.-Y., Pillai, N. E., Sim,  
712 K.-S., Xu, H., Sim, N.-L., Teo, S.-M., Foo, J.-N., Tan, L. W.-L., Lim, Y., Koo, S.-H., Gan, L. S.-  
713 H., Cheng, C.-Y., Wee, S., ... Teo, Y.-Y. (2013). Deep whole-genome sequencing of 100  
714 southeast Asian Malays. *American Journal of Human Genetics*, 92(1), 52–66.  
715 <https://doi.org/10.1016/j.ajhg.2012.12.005>
- 716 World Health Organization. (2023). *Global Tuberculosis Report 2023* (World Health  
717 Organization, Ed.; p. 75). World Health Organization.
- 718 Zaidi, S. M. A., Coussens, A. K., Seddon, J. A., Kredo, T., Warner, D., Houben, R. M. G. J., &  
719 Esmail, H. (2023). Beyond latent and active tuberculosis: a scoping review of conceptual  
720 frameworks. *EClinicalMedicine*, 66, 102332.  
721 <https://doi.org/10.1016/j.eclinm.2023.102332>

722 Zheng, R., Li, Z., He, F., Liu, H., Chen, J., Chen, J., Xie, X., Zhou, J., Chen, H., Wu, X., Wu, J.,  
723 Chen, B., Liu, Y., Cui, H., Fan, L., Sha, W., Liu, Y., Wang, J., Huang, X., ... Ge, B. (2018).  
724 Genome-wide association study identifies two risk loci for tuberculosis in Han Chinese.  
725 *Nature Communications*, 9(1), 4072. <https://doi.org/10.1038/s41467-018-06539-w>

726 Zhou, W., Nielsen, J. B., Fritsche, L. G., Dey, R., Gabrielsen, M. E., Wolford, B. N., LeFaive, J.,  
727 VandeHaar, P., Gagliano, S. A., Gifford, A., Bastarache, L. A., Wei, W.-Q., Denny, J. C., Lin,  
728 M., Hveem, K., Kang, H. M., Abecasis, G. R., Willer, C. J., & Lee, S. (2018). Efficiently  
729 controlling for case-control imbalance and sample relatedness in large-scale genetic  
730 association studies. *Nature Genetics*, 50(9), 1335–1341. [https://doi.org/10.1038/s41588-](https://doi.org/10.1038/s41588-018-0184-y)  
731 [018-0184-y](https://doi.org/10.1038/s41588-018-0184-y)  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742

743

## Supplementary Material

744

### Bantu-speaking African Ancestry

745

746

747

748

749

750

751

752

753

754

755

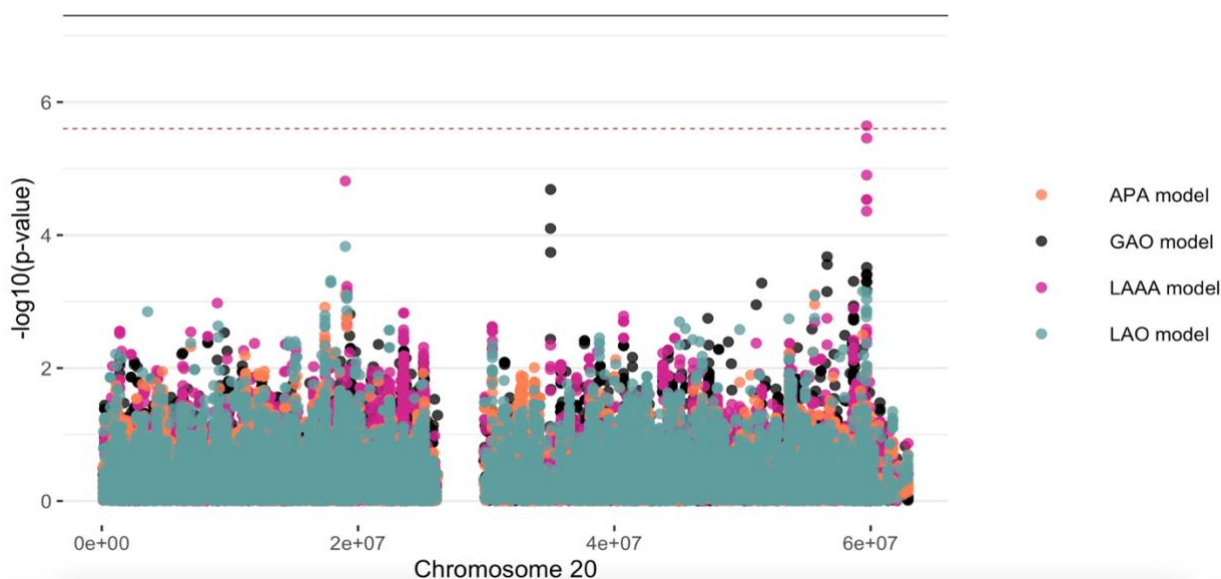
756

757

758

759

760



761

**Supplementary Figure 1.** Log transformation of association signals obtained for Bantu-speaking African ancestry whilst using the LAAA model on chromosome 20. The dashed red line represents the significant threshold for admixture mapping calculated with the software STEAM ( $p$ -value =  $2.5 \times 10^{-6}$ ) and the black solid line represents the genome-wide significant threshold ( $p$ -value =  $5 \times 10^{-8}$ ). The four different models are represented in black (global ancestry only - GAO), blue (local ancestry effect - LAO), orange (ancestry plus allelic effect - APA) and pink (local ancestry-adjusted allelic effect - LAAA).

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

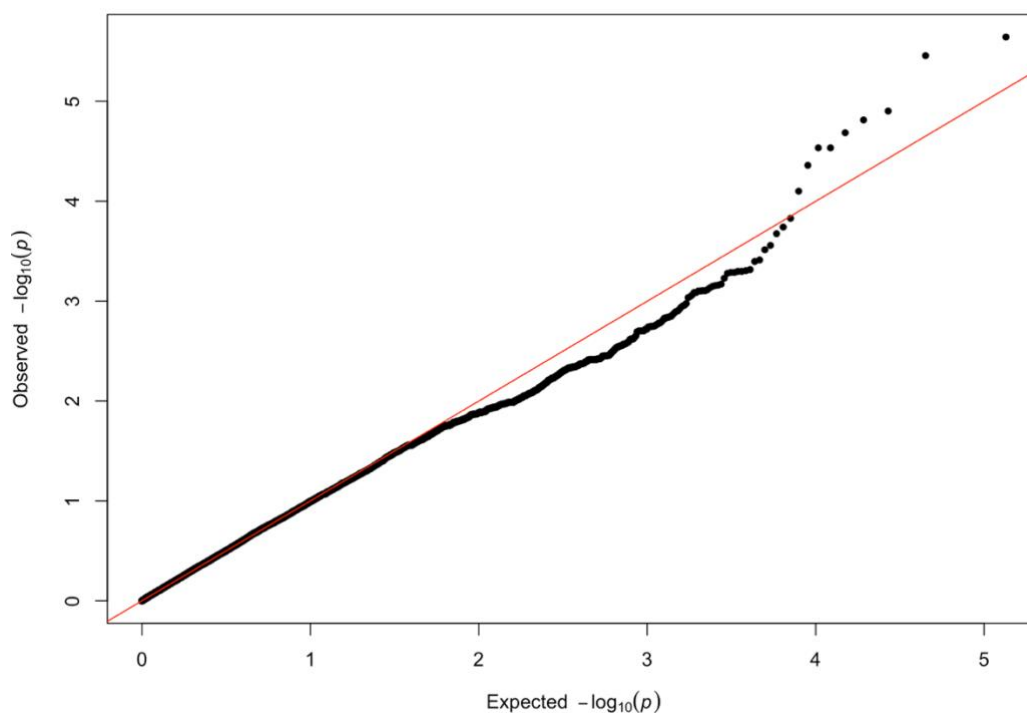
780

781

782

783

784



785

**Supplementary Figure 2.** QQ-plot of expected  $p$ -values and observed  $p$ -values for the association signals obtained for Bantu-speaking African ancestry located on chromosome 20.

786

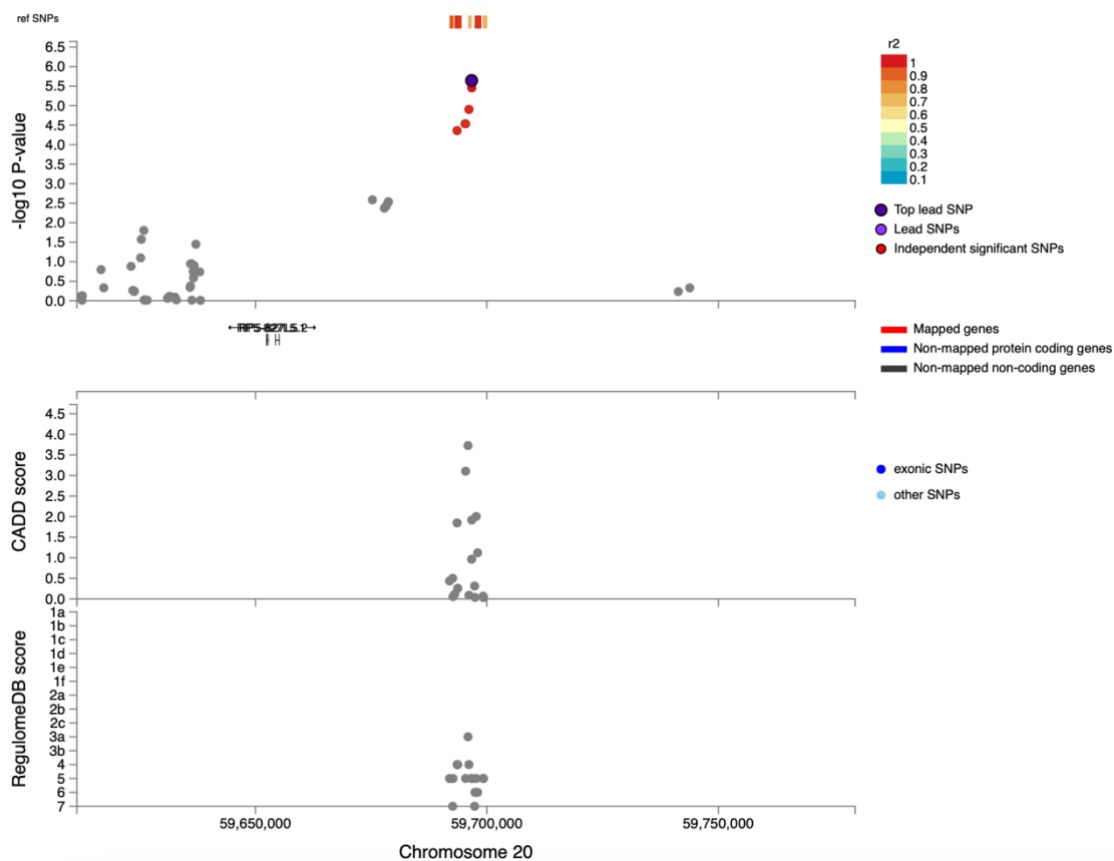
787

788

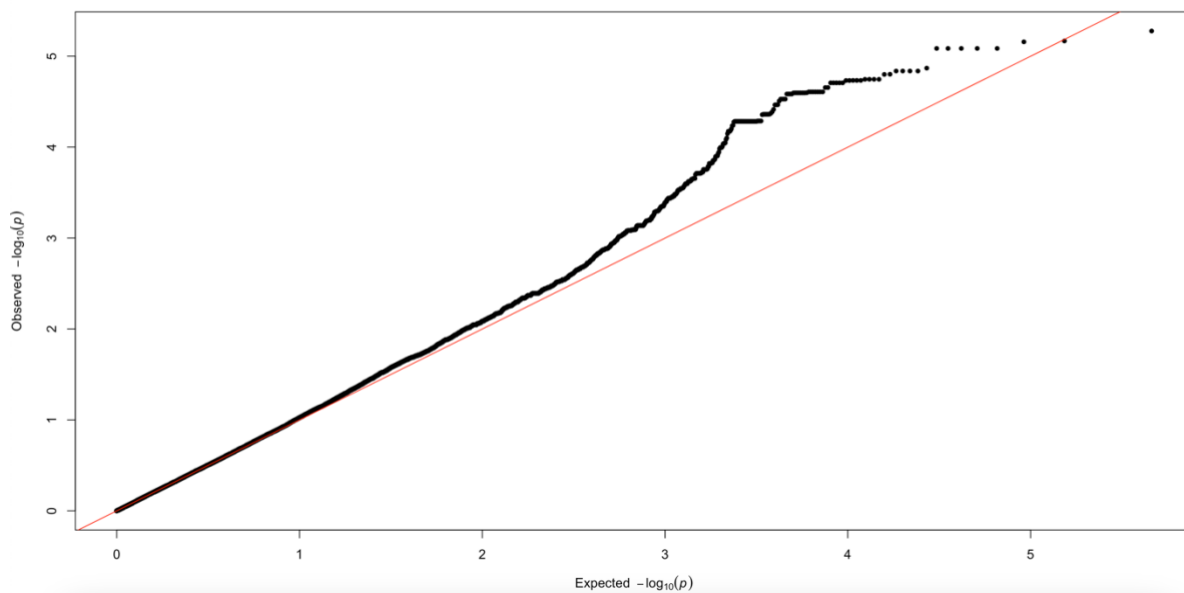
789

790

791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853



**Supplementary Figure 3.** Regional plot indicating the nearest genes in the region of the lead variant (*rs74828248*) observed on chromosome 20. SNPs in linkage disequilibrium (LD) with the lead variant are coloured red/orange. The lead variant is indicated in purple. Functional protein-coding genes are coded in red and non-functional (pseudo-genes) are indicated in black.



**Supplementary Figure 4.** QQ-plot of expected  $p$ -values and observed  $p$ -values for the association signals obtained for Khoisan ancestry located on chromosome 6.



854  
855  
856  
857  
858

**Supplementary Table 1.** Summary statistics for two variants within 800 base pairs of the ITHGC lead SNP (*rs28383206*) on chromosome 6 for the LAAA analysis adjusting for KhoeSan and Bantu-speaking African local ancestry.

Position	Marker name	Ref	Alt	AltFreq	<i>p</i> -value (KhoeSan local ancestry)	<i>p</i> -value (Bantu-speaking African local ancestry)
32576009	<i>rs482205</i>	T	G	0.322	0.032	0.116
32576019	<i>rs482162</i>	T	C	0.322	0.032	0.116

859  
860  
861  
862  
863  
864