

Progression of recent *Mycobacterium tuberculosis* exposure to active tuberculosis is a highly heritable complex trait driven by 3q23 in Peruvians

Yang Luo^{1,2,3,4,5}, Sara Suliman¹, Samira Asgari^{1,2,3,4,5}, Tiffany Amariuta^{1,2,3,4,5}, Roger Calderon⁶, Leonid Lecca⁶, Segundo R. León⁶, Judith Jimenez⁶, Rosa Yataco⁶, Carmen Contreras⁶, Jerome T. Galea⁷, Mercedes Becerra⁸, Sergey Nejentsev⁹, Marta Martínez-Bonet¹, Peter A. Nigrovic^{1,10}, D. Branch Moody¹, Megan B Murray⁸, Soumya Raychaudhuri^{1,2,3,4,5,11}

Affiliations

¹Division of Rheumatology, Immunology, and Allergy, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

²Division of Genetics, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

³Broad Institute of MIT and Harvard, Cambridge, MA, USA

⁴Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA

⁵Center for Data Sciences, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115

⁶Socios En Salud, Lima, Peru

⁷School of Social Work, University of South Florida, FL, USA

⁸Department of Global Health and Social Medicine, and Division of Global Health Equity, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

⁹Department of Medicine, University of Cambridge, Cambridge, UK

¹⁰Division of Immunology, Boston Children's Hospital, Boston MA 02115, USA

¹¹Arthritis Research UK Centre for Genetics and Genomics, Manchester Academic Health Science Centre, University of Manchester, Manchester, UK

Corresponding authors

Soumya Raychaudhuri

77 Avenue Louis Pasteur, Harvard New Research Building, Suite 250D
Boston, MA 02115, USA.

Megan Murray

Department of Global Health and Social Medicine, Harvard Medical School, Boston, Massachusetts, USA.

Department of Epidemiology, Harvard School of Public Health, Boston, Massachusetts, USA.

1 **Among 1.8 billion people worldwide infected with *Mycobacterium tuberculosis*, 5-15% are**
2 **expected to develop active tuberculosis (TB). Approximately half of these will progress**
3 **to active TB within the first 18 months after infection, presumably because they fail to**
4 **mount the initial immune response that contains the local bacterial spread. The other half**
5 **will reactivate their latent infection later in life, likely triggered by a loss of immune**
6 **competence due to factors such as HIV-associated immunosuppression or ageing. This**
7 **natural history suggests that undiscovered host genetic factors may control early**
8 **progression to active TB. Here, we report results from a large genome-wide genetic study**
9 **of early TB progression. We genotyped a total of 4,002 active TB cases and their**
10 **household contacts in Peru and quantified genetic heritability (h_g^2) of early TB**
11 **progression to be 21.2% under the liability scale. Compared to the reported h_g^2 of**
12 **genome-wide TB susceptibility (15.5%), this result indicates early TB progression has a**
13 **stronger genetic basis than population-wide TB susceptibility. We identified a novel**
14 **association between early TB progression and variants located in an enhancer region on**
15 **chromosome 3q23 (rs73226617, OR=1.19; $P < 5 \times 10^{-8}$). We used *in silico* and *in vitro***
16 **analyses to identify likely functional variants and target genes, highlighting new**
17 **candidate mechanisms of host response in early TB progression.**

18
19 The infectious pathogen *Mycobacterium tuberculosis* (*M.tb*) infects about one quarter of the
20 world's population¹. Approximately 5-15% of infected individuals progress to active TB while the
21 vast majority remain infected with viable latent *M.tb* (**Figure 1a**). From the approximately 10.4
22 million patients with active TB, an estimated ~1.3 million people died in 2016². Active TB can
23 develop immediately (within the first 18 months) after recent *M.tb* infection or after many years
24 of latency, presumably caused via distinct disease mechanisms. Late progression or TB
25 reactivation is more likely the consequence of acquired immune compromise due to other

26 diseases or ageing, whereas early progression is presumably due to failure in mounting the
27 initial immune response that contains the bacterial spread. Previous studies have indicated a
28 strong heritable component of population-wide TB susceptibility, that includes early disease
29 progression, reactivation and infection³⁻⁵. But whether early progression has a different genetic
30 architecture compared to population-wide susceptibility has yet to be defined.

31
32 Reported associations for TB, and other infectious diseases, has to be considered in the context
33 of TB diagnostic criteria and selected control groups^{6,7}. To date, genome-wide association
34 studies (GWAS) of TB have compared mixed pools of TB patients with early progression or
35 reactivation, to population controls, who may not have been exposed to *M.tb* at all⁸⁻¹². Hence,
36 known human genetic loci associations with clinical outcomes might represent risk factors for
37 *M.tb* infection, progression from recent *M.tb* exposure to active TB, or reactivation of TB after a
38 period of latency. Infection, progression and reactivation represent pathophysiologically distinct
39 disease transitions likely involving distinct mechanisms of transmission, early innate immune
40 response and control by adaptive immunity. Thus, the study of mixed TB populations using
41 controls of unknown exposure status may underestimate or miss genetic associations for these
42 separate stages of disease.

43
44 To identify host factors that drive pulmonary early TB progression, we conducted a large,
45 longitudinal genetic study in Lima, Peru (**Fig. 1b**), where the TB incidence rate is one of the
46 highest in the region². We enrolled patients with microbiologically confirmed pulmonary TB.
47 Within two weeks of enrolling an index patient, we identified their household contacts (HHCs)
48 and screened for infection as measured by a tuberculin skin test (TST) and for signs and
49 symptoms of pulmonary and extra-pulmonary TB. HHCs were re-evaluated at two, six and
50 twelve months. We considered individuals to be early progressors if they are (1) index patients
51 whose *M.tb* isolates shared a molecular fingerprint with isolates from other enrolled patients; (2)

52 HHCs who developed TB disease within one year after exposure to an index patient and (3)
53 index patients who were 40 years old or younger at time of diagnosis. We considered HHCs
54 who were TST positive at baseline or any time during the 12 month follow up period, but who
55 had no previous history of TB disease and remained disease free, as non-progressing controls
56 **(Methods, Figure 1b)**. In total, we genotyped 2,175 recently exposed pulmonary TB cases
57 (early progressors) versus 1,827 HHCs with latent tuberculosis infection, who had not
58 progressed to active TB during one year of follow-up (non-progressors), as controls **(Methods,**
59 **Supplementary Table 1)**.

60
61 To our knowledge, this represents the most extensive genetic study conducted in Peru to date.
62 Peru is a country with a complex demographic history and underexplored genomic variation.
63 When Spanish conquistadors arrived in the region in the 16th century, Peru was the center of
64 the vast Inca Empire and was inhabited by a large Native American population^{13,14}. During the
65 colonial period, Europeans and Africans (brought in as slaves) arrived in large numbers to Peru.
66 After Peru gained its independence in 1821, there was a flow of immigrants from southern
67 China to all regions of Peru as a replacement for slaves^{15,16}. As a result, the genetic background
68 of the current Peruvian population is shaped by different levels of admixture between Native
69 Americans, Europeans, African and Asian immigrants that arrived in waves with specific and
70 dated historical antecedents. When compared to individuals from other South American
71 countries^{17,18}, Peruvians tend to share a greater genetic similarity with Andean indigenous
72 people such as Quechua and Aymara **(Figure 2, Supplementary Figure 1, Methods)**.

73
74 This unique genetic heritage provides both a challenge and an opportunity for biomedical
75 research. To optimally capture genetic variation, and particularly rare variations in Peruvians,
76 we designed a 712,000-SNP customized array (LIMAArray) with genome-wide coverage based
77 on whole-exome sequencing data from 116 active TB cases **(Methods, Supplementary Table**

78 **2, Supplementary Figure 2**). When compared to other more comprehensive genotyping
79 platforms available at the time, LIMAArray showed an approximately 5% increase in imputation
80 accuracy, particularly for population-specific and low-frequency variants (**Supplementary Table**
81 **3**). We derived estimated genotypes for ~8 million variants using the 1000 Genomes Project
82 Phase 3¹⁷ as the reference panel and tested single marker and rare-variant burden associations
83 with linear mixed models that account for both population stratification and relatedness in the
84 cohort (**Supplementary Figure 3-4, Methods**). Genome-wide association results of 2,160
85 cases and 1,820 controls after quality control (**Methods**) are summarized in **Supplementary**
86 **Figure 5**. We observed no inflation of test statistics ($\lambda_{GC} = 1.03$, $\lambda_{GC} = 1.00$ for common and
87 rare association analyses respectively), which suggests potential biases were strictly controlled
88 in our study.

89
90 To investigate the genetic basis of early TB progression, we first estimated its variant-based
91 heritability (h_g^2). Using GCTA¹⁹ we estimated h_g^2 of TB progression to be 21.2% (standard error
92 (s.e.)=0.08, $P = 2.64 \times 10^{-3}$) on the liability scale with assumed incidence rate of 0.05 in the
93 cohort (**Methods**). To avoid biases introduced from calculating genetic relatedness matrices
94 (GRMs) in admixed individuals, we calculated two different GRMs based on admixture-aware
95 relatedness estimation methods^{20,21} and removed related individuals. Both admixture-aware
96 methods reported similar h_g^2 estimates (**Supplementary Table 4**), indicating our reported
97 heritability estimation is robust under different model assumptions. We quantified h_g^2 of TB
98 progression and observed a surprisingly strong genetic basis. This degree of heritability is
99 comparable to traits with a well-established genetic basis (**Supplementary Table 5**). For
100 example, GWAS have identified ~200 risk loci for Crohn's disease^{22,23}, which has a reported h_g^2
101 of 28.4% (s.e.=0.02, $P = 8.62 \times 10^{-71}$)²². In contrast, using LD Score regression²⁴ on summary
102 statistics from a GWAS of population-wide TB susceptibility in Russia¹⁰, we estimated the h_g^2 of

103 population-wide TB susceptibility to be 15.5% (s.e.=0.04, $P = 5.33 \times 10^{-5}$) with assumed
104 prevalence of 0.04²⁵. These data suggest that recently exposed TB progression may have a
105 stronger host genetic basis compared to population-wide TB susceptibility, and larger
106 progression studies may be well-powered to discover additional variants.

107

108 We next identified a novel risk locus associated with TB progression on chromosome 3q23,
109 which is comprised of 11 variants in non-coding regions downstream of *RASA2* and upstream of
110 *RNF7* ($P < 1 \times 10^{-5}$) (**Figure 3, Supplementary Table 6**). The strongest association was at a
111 genotyped variant rs73226617 (OR=1.19; $P = 3.93 \times 10^{-8}$). To test for artifacts and to identify
112 stronger associations that might have been missed due to genotyping and imputation, we first
113 checked the genotype intensity cluster plot of the top associated variant which showed clear
114 separation between genotypes AA, AG and GG (**Supplementary Figure 6**). We then designed
115 individual TaqMan genotyping assays for four top associated variants (**Methods,**
116 **Supplementary Table 7**). We genotyped these four SNPs in 4,002 initial subjects and
117 concluded that all four variants show a high concordance rate (>99%) with imputed genotypes
118 (**Supplementary Table 6**). Because all 11 variants in the risk locus are in high linkage
119 disequilibrium (LD) with each other (**Supplementary Figure 7**), the other imputed variants are
120 also likely to have high imputation quality.

121

122 To determine whether the reported association is specific to TB progression from recent *M.tb*
123 infection or instead derived from reactivation of latent TB, we conducted a case-only analysis
124 removing age from our case selection criteria. This approach is based on the premise that TB
125 cases that share a DNA fingerprint for *M.tb* and HHCs who developed active TB are
126 epidemiologically related while cases in which *M.tb* fingerprints are different might have resulted
127 from remote infection that reactivated during the study assessment²⁶. 1,472 out of 2,175

128 presumed early progressors shared molecular fingerprint of *M.tb* isolates with another case or
129 developed active TB during the one year of follow-up (**Supplementary Figure 8**). Other cases
130 did not have a shared the molecular fingerprint among *M.tb* isolates or did not come from the
131 same household as the index case, leading to a lower degree of certainty in the early
132 progression status of these cases. In this case-only analysis, the top associated signal
133 rs73226617 was nominally associated with early progression ($P = 0.016$, OR=1.09). A
134 heritability analysis restricted to those that shared the same molecular fingerprint or from the
135 same household estimated in a larger h_g^2 (22.1%, s.e.=0.06, $P = 1.32 \times 10^{-4}$) despite the smaller
136 number of samples. These results provide further evidence that the signals we reported at 3q23
137 are only associated with TB progression after recent exposure to *M.tb* and not from reactivation
138 of latent infection.

139
140 We examined our 11 most associated variants for early TB progression identified in the
141 Peruvian cohort in previously published GWAS datasets^{8-10,27} (**Supplementary Table 8**). These
142 SNPs were less frequent (<1%) in the African populations than in the European and Peruvian
143 populations, resulting in lower statistical power to detect association. We therefore examined
144 the SNPs in two previously published Russian¹⁰ (5,530 TB cases and 5,607 controls) and
145 Icelandic²⁷ (4,049 TB cases and 6,543 TST+ controls) GWAS datasets. We observed that the
146 effects in the Russian cohort were similar, as they shared comparable ORs of 1.10 (Peru) and
147 1.18 (Russia) for rs73226617 ($P_{\text{Russia}}=0.065$). In contrast, there was no signal observed in the
148 Icelandic cohort (OR=1.06, $P_{\text{Iceland}}=0.437$). Consistent with our previous case-only analysis, the
149 weaker signals observed in both European cohorts indicate that 3q23 is specifically associated
150 with early TB progression. The association signals were therefore most likely diluted due to the
151 inclusion of reactivation cases and non-infected controls in the cohort collection.

152

153 We next examined how previously published TB GWAS risk loci are associated in this study.
154 We detected evidence of association in a previously reported TB locus at rs9272785 in the HLA
155 region²⁷ (OR=1.04, $P = 4.49 \times 10^{-3}$), but did not detect signals at other reported risk loci
156 (**Supplementary Table 9**). Thus, previously reported loci may relate to infection or reactivation
157 phenotypes, rather than early TB progression whereas HLA association may affect both early
158 progression and reactivation. The strongest association observed in the HLA region after
159 imputation (**Method, Supplementary Figure 9**) was rs7739434 located upstream of HLA-A
160 (OR=1.10 $P = 4.59 \times 10^{-7}$), indicating a possible HLA class I association with TB progression.
161
162 To try to identify which of the variants in our reported risk locus is likely to be the functional
163 polymorphism affecting the risk of pulmonary TB progression, we employed the FINEMAP²⁸
164 software (**Methods**). The 90% credible set includes seven genomic variants, with rs73226617
165 having the highest posterior probability (0.54), followed by rs58538713 (0.16) and the indel
166 rs148722713 (0.05) among 713 variants in the region (**Supplementary Table 6**). To identify
167 likely functional variants and target genes, we employed a method called IMPACT (Inference
168 and Modeling of Phenotype-related ACtive Transcription)²⁹. Briefly, IMPACT identifies regions
169 predicted to be involved in transcriptional regulatory processes related to a cell-type-specific key
170 transcription factor (**Method, Figure 3**) by leveraging information from nearly 400 *in silico*
171 epigenomic and sequence annotations from public databases (**Supplementary Table 10**). We
172 trained IMPACT on the epigenetic chromatin signature of binding sites of the transcription factor
173 IRF1 to identify active regulatory regions specific to macrophages. Among 11 variants in the risk
174 locus, the leading associated variant rs73226617 had the highest predicted probability (0.704)
175 of lying in an active macrophage-specific regulatory region. Overexpression of *IRF1*, along with
176 other Type I interferon response genes, was detected early in tuberculosis contacts who
177 progressed to active disease^{30,31}. Overall, we saw an enrichment of the interferon response
178 factor in the 3q23 locus (**Figure 3d**). We then performed electrophoretic mobility shift assays

179 (EMSA) and luciferase assays to functionally identify the most likely causal variant among the
180 seven variants that constitute 90% credible set (**Method**). EMSA tests whether the variants
181 differentially bound nuclear complexes in an allele-specific manner. Four variants (rs73226617,
182 rs148722713, rs11710569 and rs73226608) showed differential EMSA signals in the risk
183 variants that were suppressed with unlabeled probes, consistent with allele-specific protein
184 binding in the Jurkat76 cell line. We performed luciferase assays on the four candidate variants
185 but failed to detect allele-specific enhancer activity on human embryonic kidney (HEK293T)
186 cells (**Methods, Supplementary Figure 10**). This negative result may be driven by the
187 sensitivity limits of the assay or the variants having cell-type-specific activities which might not
188 have been captured by the designed assays. Using GeneHancer³² version 4.7, the most
189 plausible target genes include *RNF7* (GH03I141681, gene-enhancer score =11.7, 56,393 base
190 pairs (bp) from the top associated variant rs73226617). *RNF7* is a highly conserved ring finger
191 protein. It is an essential subunit of SKP1-cullin/CDC53-F box protein ubiquitin ligases, which
192 are a part of the protein degradation machinery important for cell cycle progression and signal
193 transduction. Among its related pathways are innate immune system and class I MHC mediated
194 antigen processing and presentation. Another candidate target gene located near the risk locus
195 is *RASA2*³², which lies 66,469 bp upstream from our top associated variant and is expressed in
196 human lung cells³³. The *RASA2* protein is a member of the GAP1 family of GTPase-activating
197 proteins (GAPs) that is involved in cellular proliferation and differentiation. It has been indicated
198 that *RASA2* acts as a regulator of alveolar macrophage activation³⁴. Interestingly, previously
199 reported TB-associated gene *ASAP1*¹⁰ also encodes GAPs, indicating that the family of GAPs
200 may play an important role in TB pathogenesis.

201
202 Overall, our results argue that rapid TB progression is a highly heritable trait, comparable to
203 other human diseases with an established genetic origin. More generally, these results begin to
204 address general questions about genomic approaches to infectious diseases, which have

205 lagged behind other disease in terms of locus discovery in comparison to other complex traits
206 **(Supplementary Table 5)**. Infections, especially chronic infectious diseases, play out in highly
207 distinct phases that involve exposure, crossing epithelial boundaries, pathogen expansion,
208 locating a host niche, and in the case of TB, decades-long persistence, reactivation and re-
209 transmission. Each of these stages can be controlled by distinct host factors. Our analysis
210 indicates that progression from recent *M.tb* exposure to active TB has a different genetic basis
211 compared to TB reactivation. Specific analysis of clinical progression as a distinct phase allows
212 for a more powerful detection of risk factors for an equal number of samples, as compared to
213 case-control studies, which are an amalgamation of different phenotypes. Thus, this work
214 argues strongly that while detailed, stage-specific phenotypic profiling may be more costly, it
215 may offer key advantages for infectious disease genetic studies. Specifically, it allows for
216 precise phenotype definitions, greater heritability and identification of biological targets with
217 specific implications. Therefore detailed phenotypic profiling should become a general strategy
218 for future genetic studies of infectious diseases.

219

220 **Methods**

221 **Ethics statement**

222 We recruited 4,002 subjects from a large catchment area of Lima, Peru that included 20 urban
223 districts and approximately 3.3 million residents to donate a blood sample for use in our study.
224 We obtained written informed consent from all the participants. The study protocol was
225 approved by the Institutional Review Board of Harvard School of Public Health and by the
226 Research Ethics Committee of the National Institute of Health of Peru.

227 Preparation of genome-wide genetic data

228 We enrolled index cases as adults (aged 15 and older) who presented with clinically suspected
229 pulmonary TB at any of 106 participating health centers. We excluded patients who resided
230 outside the catchment area, who had received treatment for TB before and those who were
231 unable to give informed consent. Pulmonary TB patients have been diagnosed by the presence
232 of acid fast bacilli in sputum smear or a positive *M.tb* culture at any time from enrollment to the
233 end of treatment. All cultures of the index cases were genotyped using mycobacterial
234 interspersed repetitive units-variable number of tandem repeats (MIRU-VNTR). Within two
235 weeks of enrolling an index patient, we enrolled his or her household contacts (HHCs). The *M.tb*
236 status was determined using the Tuberculin Skin Test (TST). All HHCs were evaluated for signs
237 and symptoms of pulmonary and extra-pulmonary TB disease at two, six and 12 months after
238 enrollment. To select cases who were likely to have recently exposed TB, we chose HIV-
239 negative, culture-positive, drug-sensitive pulmonary TB cases from one of three groups: (1)
240 exposed HHCs who developed active TB during a 12 month follow up period; (2) index patients
241 whose *M.tb* isolates shared a molecular fingerprint with isolates from other enrolled patients and
242 (3) index patients who were 40 years old or younger at time of diagnosis. To maximize the
243 likelihood that controls were exposed to *M.tb* but did not develop active disease, we chose them
244 from among TST positive HHCs with no previous history of TB disease, and who remained
245 disease free at the time of recruitment both by directly re-contacting individuals to inquire
246 about their latest medical history and by checking their names against lists of notified TB
247 patients at all of the 106 health clinics. Where possible, we chose controls who are less than
248 second-degree related to the index cases.

249 Customized Axiom array for Peruvian populations

250 We developed a custom array (LIMAArray) based on whole-exome sequencing data from 116
251 active TB cases to optimize the capture of genome-wide genetic variation in Peruvians. Many
252 markers were included because of known associations with, or possible roles in, phenotypic
253 variation, particularly TB-related (**Supplementary Table 11**). The array also includes coding
254 variants across a range of minor allele frequencies (MAFs), including rare markers (<1% MAF),
255 and markers that provide good genome-wide coverage for imputation in Peruvian populations in
256 the common (>5%), low frequency (1-5%) and rare (0.5-1%) MAF ranges (**Supplementary**
257 **Table 3**). This approach allowed the detection of rare population specific coding variants and
258 those which predisposed individuals to TB risk.

259 Genotyping and quality control

260 We extracted genomic DNA from whole blood of the participating subjects. Genotyping of all
261 samples was performed using our customized Affymetrix LIMAArray. Genotypes were called in
262 a total of 4,002 samples using the apt-genotype-axiom³⁵. Individuals were excluded if they were
263 missing more than 5% of the genotype data, had an excess of heterozygous genotypes (± 3.5
264 standard deviations, **Supplementary Table 12**), duplicated with identity-by-state >0.9 or index
265 cases with age at diagnosis greater than 40 years old. After excluding these individuals, we
266 excluded variants with a call rate less than 95%, with duplicated position markers, those with a
267 batch effect ($P < 1 \times 10^{-5}$), Hardy-Weinberg (HWE) P-value below 10^{-5} in controls, and a
268 missing rate per SNP difference in cases and controls greater than 10^{-5} (**Supplementary**
269 **Table 13**). In total, there were 3,980 samples and 677,232 SNPs left for imputation and
270 association analyses after quality control.

271 Imputation and association analyses

272 The genotyped data were pre-phased using SHAPEIT³⁶. IMPUTE2³⁷ was then used to impute
273 genotypes at untyped genetic variants using the 1000 Genomes Project Phase 3 dataset¹⁷ as a
274 reference panel. For chromosome X, males are coded as diploid. That is male genotypes are
275 coded as 0/2 and females genotypes are coded as 0/1/2. HLA imputation was performed using
276 SNP2HLA³⁸ and a multi-ethnic HLA imputation reference pane³⁹. Imputed SNPs were excluded
277 if the imputation quality score r^2 was less than 0.4, HWE P-value $< 10^{-5}$ in controls or a
278 missing rate per SNP greater than 5%. After filtering, 7,756,401 SNPs were left for further
279 association analyses.

280

281 Common single variant associations were tested with a linear mixed model (LMM) implemented
282 in GEMMA⁴⁰ version 0.94.1 on genotype likelihood from imputation assuming an additive
283 genetic model. We use the genetic relatedness matrix (GRM) as random effects to correct for
284 cryptic relatedness between collected individuals. Sex and age were included as fixed effects to
285 correct for population stratification (**Supplementary Figure 2**). The GRM was obtained from an
286 LD-pruned ($r^2 < 0.2$), with MAF $\geq 1\%$ after removing large high-LD regions⁴¹ (**Supplementary**
287 **Table 14**) dataset of 154,660 SNPs using GEMMA⁴² version 0.7-1.

288

289 Gene-based rare variant (MAF $<1\%$) burden test was performed using GMMAT⁴² version 0.7-1,
290 a generalized linear mixed model framework. For each gene j , we aggregated the information
291 for multiple rare variants into a single burden score ($C_i = \sum_{j=1}^M G_{ij}$) for each subject i . Where G_{ij}
292 denotes the allele counts {0,1,2} for m variants in the gene. The genomic control inflation factor
293 (λ_{GC}) for variants after imputation was 1.03 and 1.00 for common and rare association study
294 respectively (**Supplementary Figure 4**), indicating that we have successfully controlled for any
295 residual population structure or cryptic relatedness between genotyped samples.

296

297 To avoid false-positive signals due to population stratification and heterogeneity of effects due
298 to differential LD in admixed populations, we also computed GRMs based on methods^{20,21} that
299 account for inflation of identity-by-state statistics due to admixture LD. LMM with admixture-
300 aware GRMs resulted in numerically similar association statistics to those from unadjusted
301 analyses (**Supplementary Table 15**).

302

303 To identify likely causal variants in the identified risk locus, we used FINEMAP²⁸ method to
304 calculate marginal likelihoods and Bayes factor for each variant assuming that there is one true
305 causal variant in the region, and it has been included in the analysis and has been well imputed
306 (--n-causal-max 1). We used the in-sample LD scores calculated using LDstore⁴³ to further
307 increase the accuracy of the fine-mapping analysis.

308

309 TaqMan SNPs and Genotyping

310 Selection of SNPs in the 3q23 locus was conducted based on information from the dbSNP
311 database (<http://www.ncbi.nlm.nih.gov/projects/SNP/>). Two polymorphisms rs73226617,
312 rs73226619, rs73239724 and rs73226608 were included for the genotyping tests. Real-time
313 PCR using the following calculations: 2.5uL Genotyping Master Mix, 0.25uL SNP Assay-probes,
314 and 2.25uL DNA template (at 5ng/uL= 11.25ng total). Thermal cycling conditions were as
315 follows: 60C 30secs Pre-read, 95 °C for 10 min, followed by 40 cycles at 95°C for 15 s and at 60
316 °C for 1 min, then 60C 30secs Post-read. Genotyping of the polymorphisms was carried out
317 using the 5' exonuclease TaqMan Allelic Discrimination assay, which was performed utilizing
318 minor groove binder probes fluorescently labeled with VIC or FAM and the protocol
319 recommended by the supplier (Applied Biosystems, Foster City, CA, USA). Analysis for
320 interpretation was performed with Via7 software and Taqman Genotyper software calls. Per

321 variant concordance rate was obtained by comparing genotypes obtained from imputation and
322 from TaqMan assays (**Supplementary Table 6**)

323 Heritability estimation

324 The genetic heritability based on genome-wide markers (h_g^2) was first estimated from the
325 genetic relatedness matrix (GRM) after removing related individuals (--grm-cutoff 0.125) and
326 corrected for population stratifications using the top 10 principal component (--qcovar), as
327 implemented in GCTA^{19,44}. Among a total of 14,044 enrolled HHCs, 692 progressed to active
328 TB. Based on these numbers, we estimated the incidence rate in the Lima cohort for recent TB
329 progression is 5%. Using this rate, we report h_g^2 on the liability scale to be 0.21 (s.e. = 0.07). If
330 the true prevalence was in fact half as high, our estimate would instead be 0.17 (s.e. = 0.02); if
331 twice as high, 0.26 (s.e. = 0.09). h_g^2 on the observed scale is 0.24 (s.e. = 0.09).

332 *In silico* functional annotation of candidate causal variants

333 We combined multiple sources of *in silico* genome-wide functional annotations from publicly
334 available databases to help identify potential functional variants and target genes in the 3q23
335 novel risk locus. To investigate functional elements enriched across the region encompassing
336 the strongest candidate causal variants, We aggregated approximately 400 genomic and
337 sequence annotations including cell-type-specific annotation types such as ATAC-seq, DNase-
338 seq, FAIRE-seq, HiChIP-H3K27ac, HiChIP-CTCF, polymerase and elongation factor ChIP-seq,
339 and histone modification ChIP-seq, as well as cell-type-nonspecific annotations such as
340 conservation, coding annotation, and distance to TSS. A list of all included resources is
341 summarized in **Supplementary Table 10**.

342

343 The influence of candidate causal variants on transcription factor binding sites was identified
344 using HaploReg⁴⁵ version 4.1. Among possible motif changes, IRF1 is a key transcriptional
345 regulator (TF) that plays critical roles in activation of macrophages by proinflammatory signals
346 such as interferon- γ and highly relevant to tuberculosis pathogenesis^{46,47}. We subsequently
347 determined genome-wide TF occupancy from publicly available ChIP-seq of IRF1
348 (**Supplementary Table 16**), a key regulator of monocyte-derived macrophages (GSE100381)⁴⁸.
349 Briefly, CD14+ human monocytes were purified from PBMCs and then treated with a
350 macrophage colony-stimulating factor (M-CSF). ChIP-seq peaks were called by macs⁴⁹ [v1.4.2
351 20120305] (FDR<0.05). Using IMPACT (Inference and Modeling of Phenotype-related ACtive
352 Transcription)²⁹, we built a model that predicts TF binding on a motif by learning the epigenomic
353 profiles of the TF binding sites. We train IMPACT on gold standard regulatory and non-
354 regulatory elements of IRF1. To build the regulatory class, we scanned the IRF1 ChIP-seq
355 peaks, mentioned above, for matches to the IRF1 binding motif, using HOMER⁵⁰ [v4.8.3] and
356 retained the highest scoring match for each ChIP-seq peak. To build the non-regulatory class,
357 we scanned the entire genome for IRF1 motif matches, again using HOMER, and selected motif
358 matches with no overlap with IRF1 ChIP-seq peaks. IMPACT learns the feature in 10-fold cross
359 validation (CV) of the complete sets of regulatory and non-regulatory elements. We scored
360 regions of interest with the learned from this CV.

361

362 Electrophoretic Mobility Shift Assay (EMSA)

363 Frozen cell pellets from the Jurkat76 cell line (ATCC) were used for preparation of nuclear
364 extracts using NE-PER Nuclear and Cytoplasmic Extraction reagent (ThermoFisher) according
365 to the manufacturer's instructions, then dialyzed overnight at 4°C with gentle stirring in 1L of
366 pre-cooled dialysis buffer (10% glycerol, 10mM Tris pH 7.5, 50mM KCl, 200mM NaCl, 1mM di-

367 thiothreitol, 1mM phenylmethanesulfonyl fluoride). Samples were quantified using BCA Protein
368 Assay Kit (ThermoFisher) and stored in 1X Halt protease inhibitor cocktail (ThermoFisher) at -
369 80°C until use. We designed single stranded oligonucleotides (30-34bp) corresponding to each
370 set of alleles (Integrated DNA Technologies, **Supplementary Table 17**), and biotinylated the
371 forward and reverse sequences separately using the Biotin 3'End DNA Labeling Kit
372 (ThermoFisher Scientific) following the manufacturer's instructions. Single stranded probes were
373 annealed by incubation for 5 minutes at 95°C followed by 1 hour at room temperature. EMSA
374 reactions were performed using the LightShift Chemiluminescent EMSA kit (ThermoFisher).
375 Binding reactions were performed in a volume of 20µL: 2µL of 10x binding buffer, 16µg nuclear
376 extract, 2.5% glycerol, 5mM MgCl₂, 0.05% Nonidet P-40 and 50ng Poly dl:dC as a non-specific
377 DNA competitor, and 20fmol of biotinylated probes with or without unlabeled competitor probes
378 at 200 fold molar excess. The assay was performed as previously described⁵¹

379 Luciferase reporter assays

380 We designed double stranded oligonucleotides matching the probes used for EMSA and flanked
381 by either BglIII or BamHI restriction sites (**Supplementary Table 18**) for cloning either upstream
382 or downstream of the firefly Luciferase (*Luc*) gene in the pGL3 promoter reporter vector
383 (Promega), respectively. Double stranded inserts were cloned into the pGL3 vector following
384 standard cloning protocols and verified by colony PCR reactions (**Supplementary Table 19**) as
385 well as plasmid Sanger sequencing (GENEWIZ). Plasmids with inserts cloned in-sense with the
386 luciferase promoter sequence were expanded and purified using PureLink HiPure Plasmid
387 Miniprep Kits (Thermofisher Scientific). For transfection, 10⁴ human embryonic kidney
388 (HEK293T) cells were plated per well in 60µL Dulbecco's Modified Eagle Medium-10 media
389 (DMEM, Gibco, 10% fetal bovine serum, 1x penicillin-streptomycin) in flat-bottom 96-well plates
390 and transfected with 500ng of plasmids (4:1 of pGL3:pRL-TK), using lipofectamine LTX Reagent

391 with PLUS (Thermofisher) according to the manufacturer's instructions. Transfected cells were
392 incubated for 18-20 hours at 37°C, then analyzed using two-step Dual-Glo® Luciferase Assay
393 System (Promega) and read on Synergy H1 Hybrid Multi-Mode Reader (BioTek). Luminescence
394 is reported as the ratio of firefly (pGL3) to renilla (pRL-TK) luciferase luminescence, normalized
395 to pGL3. We compared the pooled averages of triplicates per plate, paired by transfection plate
396 from 3-10 independent experiments for each variant using a Wilcoxon signed-rank test.

397 Acknowledgements

398 The study was supported by the National Institutes of Health (NIH) TB Research Unit Network,
399 Grant U19 AI11224-01. The content is solely the responsibility of the authors and does not
400 necessarily represent the official views of the NIH. S.N. was supported by MRC
401 (MR/M012328/1), the ERC Starting grant (260477) and the National Institute for Health
402 Research (NIHR) Cambridge Biomedical Research Centre. The authors thank Garðar
403 Sveinbjornsson, Patrick Sulem, Ingileif Jonsdottir and Kari Stefansson at deCODE genetics,
404 Reykjavik, Iceland, for validating the association of rs73226617 with TB progression in the
405 Icelandic population.

406 Author contributions

407 Y.L. designed the genotyping array, performed statistical analysis of the GWAS data and wrote
408 the first draft of the manuscript. S.S. performed the EMSA and luciferase assays experiments.
409 S.A. carried out the rare association studies of the GWAS data. T.A. implemented the IMPACT
410 model. R.C., L.L., S.R. L., J. J., R. Y., C.C., J T. G., M.B. and M.B.M. participated in study
411 design, protocol development and sample collection. S.N. contributed the Russian data for the
412 meta and heritability analysis. M. MB. and P.A.N. helped to develop the protocols for EMSA and

413 luciferase assays experiments. D.B.M supervised the EMSA and luciferase assay experiments.

414 M.B.M. participated in study design, protocol development, and study conception. S.R.

415 conceived and supervised the study. All authors contributed to the writing of the manuscript.

416

417 Competing financial interests

418 The authors declare no competing financial interests.

419

420 References

- 421 1. Houben, R. M. G. J. & Dodd, P. J. The Global Burden of Latent Tuberculosis Infection: A
422 Re-estimation Using Mathematical Modelling. *PLoS Med.* **13**, e1002152 (2016).
- 423 2. WHO | Global tuberculosis report 2017. (2017).
- 424 3. van der Eijk, E. A., van de Vosse, E., Vandenbroucke, J. P. & van Dissel, J. T. Heredity
425 versus environment in tuberculosis in twins: the 1950s United Kingdom Prophit Survey
426 Simonds and Comstock revisited. *Am. J. Respir. Crit. Care Med.* **176**, 1281–1288 (2007).
- 427 4. Kallmann, F. J., Reisner, D. & Others. Twin Studies on Genetic Variations in Resistance to
428 Tuberculosis. *J. Hered.* **34**, 269–276 (1943).
- 429 5. Cobat, A. *et al.* High heritability of antimycobacterial immunity in an area of
430 hyperendemicity for tuberculosis disease. *J. Infect. Dis.* **201**, 15–19 (2010).
- 431 6. Stein, C. M. Genetic epidemiology of tuberculosis susceptibility: impact of study design.
432 *PLoS Pathog.* **7**, e1001189 (2011).
- 433 7. Abel, L., El-Baghdadi, J., Bousfiha, A. A., Casanova, J.-L. & Schurr, E. Human genetics of
434 tuberculosis: a long and winding road. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **369**,
435 20130428 (2014).

- 436 8. Thyre, T. *et al.* Genome-wide association analyses identifies a susceptibility locus for
437 tuberculosis on chromosome 18q11.2. *Nat. Genet.* **42**, 739–741 (2010).
- 438 9. Thyre, T. *et al.* Common variants at 11p13 are associated with susceptibility to tuberculosis.
439 *Nat. Genet.* **44**, 257–259 (2012).
- 440 10. Curtis, J. *et al.* Susceptibility to tuberculosis is associated with variants in the ASAP1 gene
441 encoding a regulator of dendritic cell migration. *Nat. Genet.* **47**, 523–527 (2015).
- 442 11. Mahasirimongkol, S. *et al.* Genome-wide association studies of tuberculosis in Asians
443 identify distinct at-risk locus for young tuberculosis. *J. Hum. Genet.* **57**, 363–367 (2012).
- 444 12. Chimusa, E. R. *et al.* Genome-wide association study of ancestry-specific TB risk in the
445 South African Coloured population. *Hum. Mol. Genet.* **23**, 796–809 (2014).
- 446 13. Sandoval, J. R. *et al.* Tracing the genomic ancestry of Peruvians reveals a major legacy of
447 pre-Columbian ancestors. *J. Hum. Genet.* **58**, 627–634 (2013).
- 448 14. Wang, S. *et al.* Geographic patterns of genome admixture in Latin American Mestizos.
449 *PLoS Genet.* **4**, e1000037 (2008).
- 450 15. Evelyn Hu-DeHart. From Slavery to Freedom: Chinese Coolies on the Sugar Plantations of
451 Nineteenth Century Cuba. *Labour Hist.* 31–51 (2017). doi:10.5263/labourhistory.113.0031
- 452 16. Gonzales, M. J. Chinese plantation workers and social conflict in Peru in the late nineteenth
453 century. *J. Lat. Am. Stud.* **21**, 385–424 (1989).
- 454 17. Sudmant, P. H. *et al.* An integrated map of structural variation in 2,504 human genomes.
455 *Nature* **526**, 75–81 (2015).
- 456 18. Reich, D. *et al.* Reconstructing Native American population history. *Nature* **488**, 370–374
457 (2012).
- 458 19. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide
459 complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 460 20. Conomos, M. P. *et al.* Genetic Diversity and Association Studies in US Hispanic/Latino
461 Populations: Applications in the Hispanic Community Health Study/Study of Latinos. *Am. J.*

- 462 *Hum. Genet.* **98**, 165–184 (2016).
- 463 21. Thornton, T. *et al.* Estimating Kinship in Admixed Populations. *Am. J. Hum. Genet.* **91**,
464 122–138 (2012).
- 465 22. Luo, Y. *et al.* Exploring the genetic architecture of inflammatory bowel disease by whole-
466 genome sequencing identifies association at ADCY7. *Nat. Genet.* **49**, 186–192 (2017).
- 467 23. Liu, J. Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel
468 disease and highlight shared genetic risk across populations. *Nat. Genet.* **47**, 979–986
469 (2015).
- 470 24. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity
471 in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
- 472 25. Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.* **49**,
473 986–992 (2017).
- 474 26. Murray, M. Determinants of cluster distribution in the molecular epidemiology of
475 tuberculosis. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 1538–1543 (2002).
- 476 27. Sveinbjornsson, G. *et al.* HLA class II sequence variants influence tuberculosis risk in
477 populations of European ancestry. *Nat. Genet.* **48**, 318–322 (2016).
- 478 28. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-
479 wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 480 29. Amariuta, T. *et al.* Rheumatoid arthritis heritability is concentrated in regulatory elements
481 with CD4+ T cell-state-specific transcription factor binding profiles. *bioRxiv* 366864 (2018).
482 doi:10.1101/366864
- 483 30. Scriba, T. J. *et al.* Sequential inflammatory processes define human progression from M.
484 tuberculosis infection to tuberculosis disease. *PLoS Pathog.* **13**, e1006687 (2017).
- 485 31. Zak, D. E. *et al.* A blood RNA signature for tuberculosis disease risk: a prospective cohort
486 study. *Lancet* **387**, 2312–2322 (2016).
- 487 32. Fishilevich, S. *et al.* GeneHancer: genome-wide integration of enhancers and target genes

- 488 in GeneCards. *Database* **2017**, (2017).
- 489 33. The GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue
490 gene regulation in humans. *Science* **348**, 648–660 (2015).
- 491 34. Hussell, T. & Bell, T. J. Alveolar macrophages: plasticity in a tissue-specific context. *Nat.*
492 *Rev. Immunol.* **14**, 81–93 (2014).
- 493 35. Schillert, A. & Ziegler, A. Genotype calling for the Affymetrix platform. *Methods Mol. Biol.*
494 **850**, 513–523 (2012).
- 495 36. O’Connell, J. *et al.* A general approach for haplotype phasing across the full spectrum of
496 relatedness. *PLoS Genet.* **10**, e1004234 (2014).
- 497 37. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and
498 accurate genotype imputation in genome-wide association studies through pre-phasing.
499 *Nat. Genet.* **44**, 955–959 (2012).
- 500 38. Jia, X. *et al.* Imputing amino acid polymorphisms in human leukocyte antigens. *PLoS One*
501 **8**, e64683 (2013).
- 502 39. Luo, Y. *et al.* Novel high-resolution multi-ethnic HLA imputation reference panels
503 constructed based on high-coverage whole-genome sequencing data. (2017).
- 504 40. Zhou, X. & Stephens, M. Efficient multivariate linear mixed model algorithms for genome-
505 wide association studies. *Nat. Methods* **11**, 407–409 (2014).
- 506 41. Price, A. L. *et al.* Long-range LD can confound genome scans in admixed populations. *Am.*
507 *J. Hum. Genet.* **83**, 132–5; author reply 135–9 (2008).
- 508 42. Chen, H. *et al.* Control for Population Structure and Relatedness for Binary Traits in
509 Genetic Association Studies via Logistic Mixed Models. *Am. J. Hum. Genet.* **98**, 653–666
510 (2016).
- 511 43. Benner, C. *et al.* Prospects of Fine-Mapping Trait-Associated Genomic Regions by Using
512 Summary Statistics from Genome-wide Association Studies. *Am. J. Hum. Genet.* **101**, 539–
513 551 (2017).

- 514 44. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height.
515 *Nat. Genet.* **42**, 565–569 (2010).
- 516 45. Ward, L. D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation,
517 and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.*
518 **40**, D930–4 (2012).
- 519 46. Langlais, D., Barreiro, L. B. & Gros, P. The macrophage IRF8/IRF1 regulome is required for
520 protection against infections and is associated with chronic inflammation. *J. Exp. Med.* **213**,
521 585–603 (2016).
- 522 47. Pine, R. Review: IRF and Tuberculosis. *J. Interferon Cytokine Res.* **22**, 15–25 (2002).
- 523 48. Park, S. H. *et al.* Type I interferons and the cytokine TNF cooperatively reprogram the
524 macrophage epigenome to promote inflammatory activation. *Nat. Immunol.* **18**, 1104–1116
525 (2017).
- 526 49. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
- 527 50. Heinz, S. *et al.* Simple combinations of lineage-determining transcription factors prime cis-
528 regulatory elements required for macrophage and B cell identities. *Mol. Cell* **38**, 576–589
529 (2010).
- 530 51. Westra, H.-J. *et al.* Fine-mapping identifies causal variants for RA and T1D in DNASE1L3,
531 SIRPG, MEG3, TNFAIP3 and CD28/CTLA4 loci. *bioRxiv* 151423 (2017).
532 doi:10.1101/151423

533

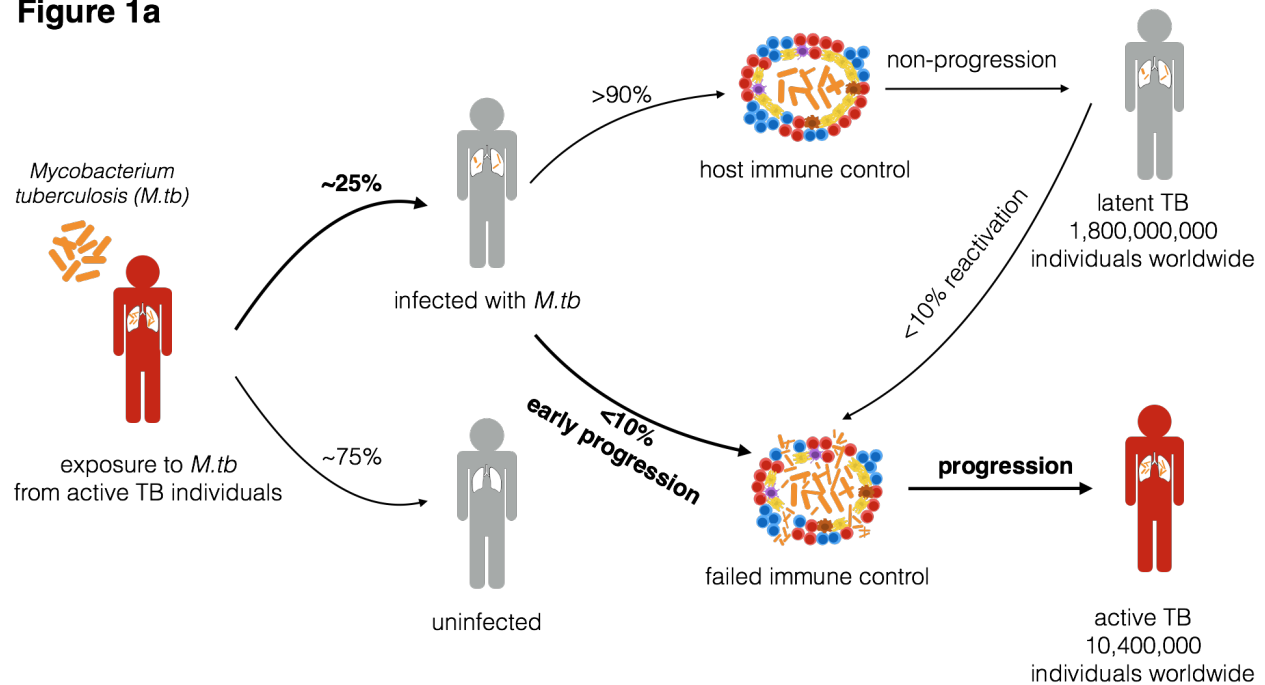
534

535

536

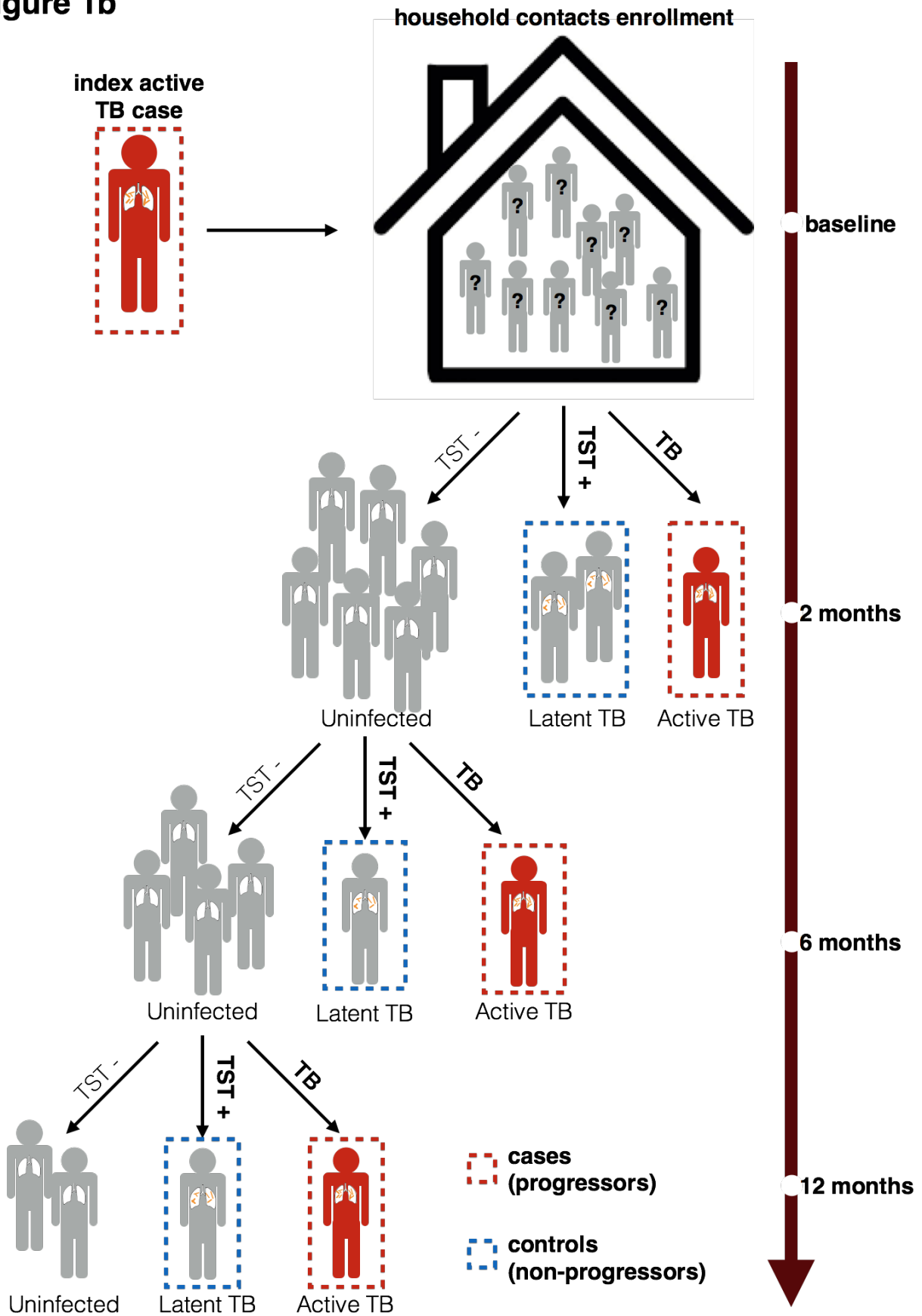
537 **Figure 1. Overview of phases of *Mycobacterium tuberculosis* (*M.tb*) infection and sample**
538 **collection.** (a) Pathophysiology of TB. Major steps following from the initial exposure to *M.tb*
539 are outlined, with the percentages of individuals progressing between steps taken from the
540 WHO TB report². (b) Schema of cohort collection. In this study, we focus on a genetic study
541 between recently exposed active pulmonary TB cases (progressors) and subjects with
542 tuberculin skin test (TST) positive results, who did not progress to active TB (non-progressors).
543 Index cases had sputum with confirmed TB. Controls were recruited in the same household as
544 index cases, with 12 month follow-up periods to confirm infection status using TST.

Figure 1a



545

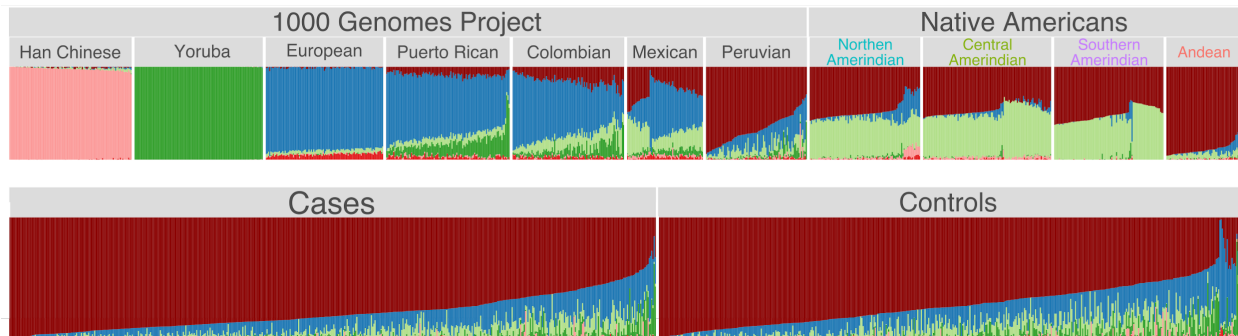
Figure 1b



547 **Figure 2. Global ancestry analysis of Peruvian populations.** (a) ADMIXTURE plot of
548 admixed individuals and continental reference panels. Each individual is represented as a thin
549 vertical bar. The colors can be interpreted as different ancestries. Reference panels are either
550 from the 1000 Genomes project¹⁷ (1000G) or Native American individuals collected from *Reich*
551 *et al. 2012 Nature*¹⁸. Han Chinese are from Beijing, China; Yoruba are from Ibadan, Nigeria;
552 European individuals are Utah Residents (CEPH) with Northern and Western European
553 Ancestry; Puerto Rican samples are from Puerto Rico; Colombian samples are from Medellin,
554 Colombia; Mexican individuals are from Los Angeles, California; Peruvian samples are from
555 Lima, Peru. Northern Amerindian includes individuals from Maya, Mixe and Kaqchikel. Central
556 Amerindian includes individuals from Pima, Zapotec, Mixtec, Yaqui, Chorotega, Tepehuano.
557 Southern Amerindian includes individuals from Piapoc, Karitiana, Surui, Wayuu, Jamaadi,
558 Parakana, Guarani, Kaingang, Ticuna, Palikur, Toba, Arara, Wichi, Chane and Guahibo.
559 Andean population includes Quechua and Aymara. $K = 6$ models are shown above, $K = 3$
560 through $K = 15$ models are available in **Supplementary Figure 1**. (b) Map of locations of
561 sampled Native American populations¹⁸.

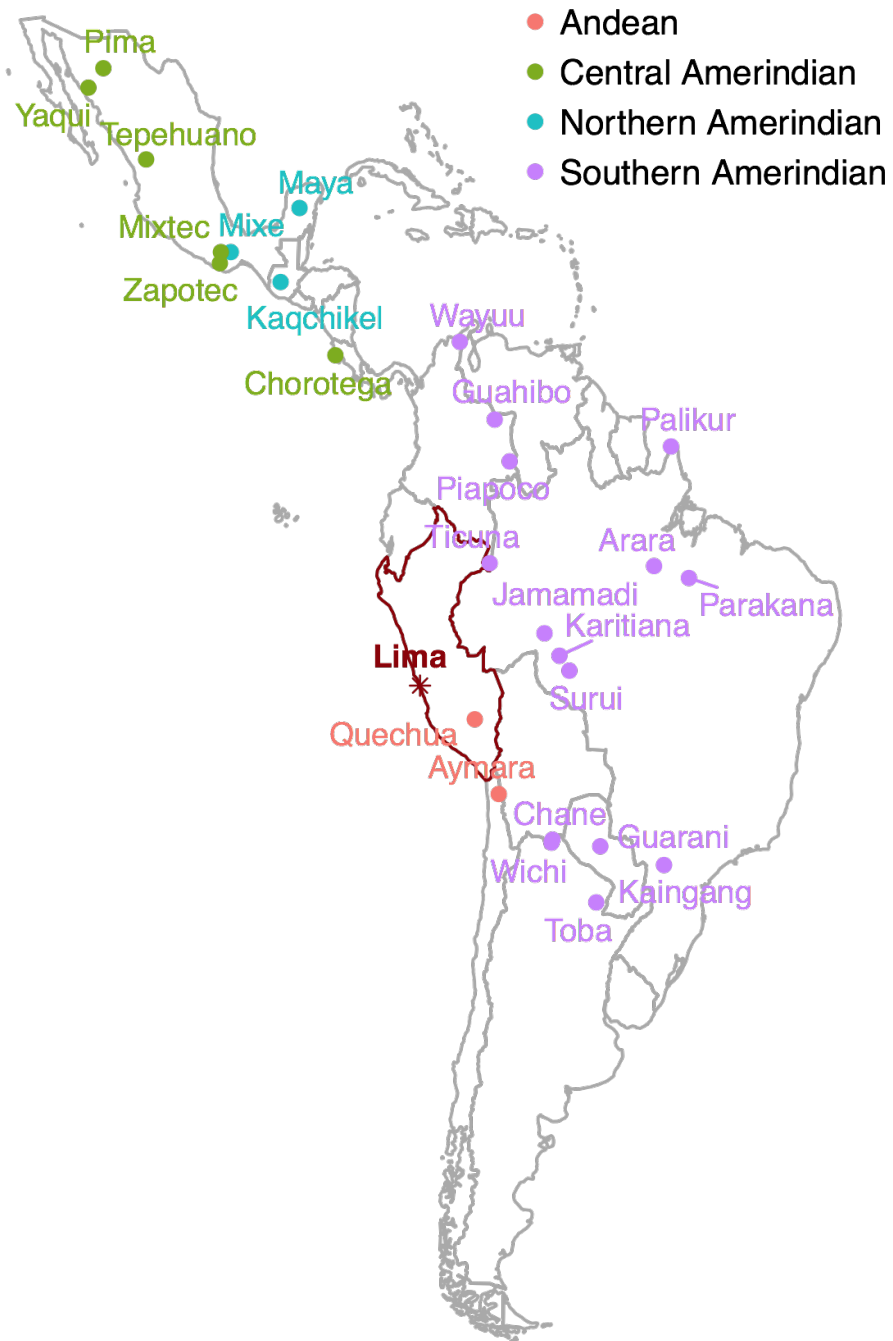
Figure 2a

ADMIXTURE $K=6$ ancestries



562

Figure 2b



563

564

565

566 **Figure 3. Genome-wide association details of the 3q23 locus.** (a) A regional association plot
567 of the 3q23 locus including all genotyped and imputed variants. (b) Fine-mapping posterior
568 probability of all variants in the chr3:140221602-145217859 region. (c) Number of overlaps
569 between all variants in the risk locus and ~400 epigenetic features. (d) Predicted posterior
570 probability of cell-type specific gene regulatory activity using Inference and Modeling of
571 Phenotype-related ACtive Transcription (IMPACT) based on the epigenetic chromatin signature
572 of binding sites of the transcription factor IRF1. Dashed lines highlights 11 top associated
573 variants. Genotyped variant rs73226617 is highlighted in red.

