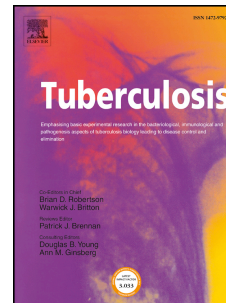


Journal Pre-proof

Host transcriptional response to TB preventive therapy differentiates two sub-groups of IGRA-positive individuals

Claire Broderick, Jacqueline M. Cliff, Ji-Sook Lee, Myrsini Kaforou, David AJ. Moore



PII: S1472-9792(20)30200-6

DOI: <https://doi.org/10.1016/j.tube.2020.102033>

Reference: YTUBE 102033

To appear in: *Tuberculosis*

Received Date: 2 August 2020

Revised Date: 24 November 2020

Accepted Date: 24 November 2020

Please cite this article as: Broderick C, Cliff JM, Lee J-S, Kaforou M, Moore DA, Host transcriptional response to TB preventive therapy differentiates two sub-groups of IGRA-positive individuals, *Tuberculosis* (2020), doi: <https://doi.org/10.1016/j.tube.2020.102033>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Ltd.

AUTHOR CONTRIBUTIONS

Claire Broderick: Conceptualisation, Methodology, Investigation, Formal analysis, Data curation, Writing- original draft, Visualisation, Project administration, Funding acquisition.

Jackie Cliff: Methodology, Investigation, Formal analysis, Data curation, Resources, Writing- original draft, Funding acquisition.

Ji-Sook Lee: Investigation, Data curation, Resources.

Myrsini Kaforou: Methodology, Formal analysis, Writing- review and editing, Visualisation.

David Moore: Conceptualisation, Methodology, Formal analysis, Writing- original draft, Supervision, Funding acquisition.

1 **Host transcriptional response to TB preventive therapy differentiates two sub-**
2 **groups of IGRA-positive individuals**

3
4 Claire Broderick ^{a,b,c*}, Jacqueline M Cliff ^{b,d}, Ji-Sook Lee ^{b,d}, Myrsini Kaforou ^c, David AJ Moore ^{a,b}

5

6

7 a. Clinical Research Department, Faculty of Infectious and Tropical Diseases, London School of
8 Hygiene and Tropical Medicine, Keppel St, London WC1E 7HT, UK

9 b. TB Centre, London School of Hygiene and Tropical Medicine, Keppel St, London WC1E 7HT,
10 UK

11 c. Section for Paediatric Infectious Disease, Department of Infectious Disease, Faculty of
12 Medicine, Imperial College London, London W2 1PG, UK

13 d. Department of Infection Biology, Faculty of Infectious and Tropical Diseases, London
14 School of Hygiene and Tropical Medicine, Keppel St, London WC1E 7HT, UK

15

16

17 Claire Broderick (Corresponding author): c.broderick@imperial.ac.uk, Tel +447850162757 ¹

18 Jacqueline M Cliff: jackie.cliff@lshtm.ac.uk

19 Ji-Sook Lee: JiSook.Lee@lshtm.ac.uk

20 Myrsini Kaforou: m.kaforou@imperial.ac.uk

21 David AJ Moore: david.moore@lshtm.ac.uk

¹ Present address: Section for Paediatric Infectious Disease, Department of Infectious Disease, Faculty of Medicine, Imperial College London, London W2 1PG, UK

22 **ABSTRACT**

23

24 We hypothesised that individuals with immunological sensitisation to *Mycobacterium tuberculosis*
25 (*Mtb*), conventionally regarded as evidence of latent tuberculosis infection (LTBI), would
26 demonstrate binary responses to preventive therapy (PT), reflecting the differential
27 immunological consequences of the sterilisation of viable infection in those with active *Mtb*
28 infection versus no *Mtb* killing in those who did not harbour viable bacilli.

29

30 We investigated longitudinal whole blood transcriptional profile responses to PT of Interferon
31 gamma release assay (IGRA)-positive tuberculosis contacts and IGRA-negative, tuberculosis-
32 unexposed controls. Longitudinal unsupervised clustering analysis with a subset of 474 most
33 variable genes in antigen-stimulated blood separated the IGRA-positive participants into two
34 distinct subgroups, one of which clustered with the IGRA-negative controls. 117 probes were
35 differentially expressed over time between the two cluster groups, many of them associated with
36 immunological pathways important in mycobacterial control.

37

38 We contend that the differential host RNA response reflects lack of *Mtb* viability in the group that
39 clustered with the IGRA-negative unexposed controls, and *Mtb* viability in the group (1/3 of IGRA-
40 positives) that clustered away.

41

42 Gene expression patterns in the blood of IGRA-positive individuals emerging during the course of
43 PT, which reflect *Mtb* viability, could have major implications in the identification of risk of
44 progression, treatment stratification and biomarker development.

45 **KEYWORDS**

46

47 Latent tuberculosis infection

48 Preventive therapy

49 Transcriptome

Journal Pre-proof

50 **1 INTRODUCTION**

51
52 The term latent tuberculosis infection (LTBI) is loaded with the inference that viable
53 *Mycobacterium tuberculosis* (*Mtb*) organisms are present in the affected individual which, under
54 the right circumstances, have the capacity to cause reactivation and TB disease. Tests of
55 immunological reactivity, whether delayed type hypersensitivity reactions measured in the
56 tuberculin skin test (TST) or T lymphocyte stimulation through antigen recognition in the interferon
57 gamma release assays (IGRAs) are widely referred to as tests for LTBI [1].

58
59 However, neither approach demonstrates presence of viable *Mtb* bacilli and there is no
60 histopathological hallmark of LTBI. The lifetime risk of reactivation disease from an *Mtb* infection
61 acquired remotely in time is around 10% [2]. In the interval between acquisition of infection and
62 development of disease, *Mtb* maintains viability and is assumed to be slowly replicating, either
63 under close immunological control or in a relatively immunologically privileged location. Thus, LTBI
64 induces immunological sensitisation as reflected in the TST and IGRA, tests that demonstrate
65 immunological memory for prior exposure to mycobacterial antigens.

66
67 Nevertheless, 90% of individuals demonstrating immunological recognition of *Mtb* antigens by
68 positive IGRA or TST never develop active TB disease. Taking the inherent assumption that TST and
69 IGRA are indicators of LTBI to its logical conclusion, the 90% who escape development of TB do so
70 because the immune control-pathogen balance remains in favour of the human host. An
71 alternative explanation might be that a large proportion of those with positive TST and IGRA
72 testing do not harbour viable organisms and are thus incapable of progressing to reactivation TB.

73
74 Preventive therapy (PT), in which a limited course of anti-TB antibiotics is used to sterilise
75 presumed viable infection in individuals with positive TST and/or IGRA tests, has been shown to be
76 highly effective in reducing the risk of future TB disease [3].

77
78 We hypothesised that differentiation of LTBI with viable bacilli from immunological sensitisation
79 without viable infection could be achieved by investigating the whole blood transcriptomic
80 response to effective PT. We hypothesised that mycobacterial killing from effective LTBI PT would
81 lead to a detectable alteration in the transcriptome that would not be seen in those individuals in

82 whom there were no *Mtb* to be killed, whether these were IGRA/TST positive or healthy IGRA/TST
83 negative controls with no known prior TB exposure.

84

85 **2 MATERIALS AND METHODS**

86

87 **2.1 Ethics statement**

88 The study procedures and protocol were approved by City & East NHS Research Ethics Committee,
89 London (reference 16/LO/1206) and the London School of Hygiene and Tropical Medicine
90 Research Ethics Committee (reference 11603). Written informed consent was given by all
91 participants before inclusion in the study.

92

93 **2.2 Participants**

94 Study participants were recruited from National Health Service (NHS) tuberculosis (TB) outpatient
95 clinics in London (Whittington Health NHS Trust, Royal Free London NHS Foundation Trust, Barts
96 Health NHS Trust, Homerton University Hospital NHS Foundation Trust). Healthy controls were
97 recruited from the London School of Hygiene and Tropical Medicine.

98

99 Participants were recruited who were aged 18 years and above, had positive Interferon Gamma
100 Release Assay (IGRA) (performed by the local hospital laboratories, using the QuantiFERON-TB
101 Gold In-tube assay [Qiagen, Manchester, UK]), with known exposure to an index person with
102 isoniazid- and rifampicin- susceptible pulmonary TB (contact history unconfirmed for three
103 individuals) and who planned to initiate a 12-week course of combined rifampicin/ isoniazid (RH)
104 as preventive therapy (once daily rifampicin 600 mg/ isoniazid 300 mg as Rifinah) plus once daily
105 pyridoxine 10 mg. Adult volunteers aged 18 years and above were recruited as healthy control
106 participants.

107

108 Once consented, demographic information, TB exposure history, and medical history were
109 recorded on a data capture sheet and testing for human immunodeficiency virus (HIV) was
110 performed. Healthy volunteers additionally underwent IGRA testing (performed using the
111 QuantiFERON-TB Gold In-tube assay according to the manufacturer's recommendations) and were
112 excluded if they were found to be IGRA+. Individuals were excluded if they had a prior history of
113 TB infection, of having taken anti-TB treatment or exposure to drug-resistant TB. Participants who
114 were pregnant, breastfeeding or trying to conceive, those with immunosuppressive disorders

115 including HIV and those who had taken immunosuppressant medication in the preceding six
116 months were also excluded. Healthy control participants reporting prior exposure to TB were also
117 excluded.

118
119 Healthy controls were given a two-week course of RH (once daily rifampicin 600 mg/ isoniazid 300
120 mg as Rifinah) plus once daily pyridoxine 10 mg.

121
122 Blood samples were collected from all participants at baseline (V1) and 2 weeks after initiating RH
123 (V2), with an additional sample point in IGRA+ participants within 6 weeks of completion of the
124 12-week course of treatment (V3). At all sampling timepoints, all participants were asked about
125 their adherence to treatment, and whole blood was collected in a PAXgene blood RNA tube
126 (PreAnalytiX GmbH, Hombrechtikon, Switzerland) for RNA expression analysis and a lithium
127 heparin tube (Becton Dickinson, Berkshire, UK) for subsequent stimulation assays. The PAXgene
128 tubes were frozen within 4 hours of collection.

129

130 **2.3 Stimulation of whole blood**

131 Stimulation was performed using QuantiFERON-TB Gold Plus In-tube Assay (QFT-TB Plus) (Qiagen).
132 Within four hours of collection, 1 ml of blood was transferred from the lithium heparin tube to
133 each of the four QFT-TB Plus tubes: TB1 antigen, TB2 antigen (both containing peptides from
134 ESAT-6 and CFP-10 antigens), mitogen positive control and (unstimulated) negative control. The
135 tubes were gently shaken to dissolve the lyophilized peptides in the blood. The QFT-TB Plus tubes
136 were immediately incubated upright at 37°C for 22 -24 hours. After incubation, the blood was
137 transferred into a 1.5 ml microcentrifuge tube and centrifuged for 15 minutes at 3000 RCF(g).
138 Supernatants were removed and the remaining cell pellet (500 µl) was transferred into a 15 ml
139 tube containing 2.5 ml RNeasy Protect[®] Cell Reagent (Qiagen). The cells were resuspended by
140 vortexing, and incubated for 2 hours for complete cell lysis before freezing at -80°C.

141

142 **2.4 Peripheral blood RNA expression by microarray**

143 Total RNA was extracted from the PAXgene tubes using the PAXgene Blood miRNA Kit (Qiagen),
144 and from the QFT-TB Plus stimulated samples, which had been lysed in RNeasy Protect, using the
145 RNeasy mini kit (Qiagen), according to the manufacturer's instructions, incorporating on-column
146 DNase digestion. Globin depletion was performed using the GLOBINclear Kit (ThermoFisher),

147 quantified by Nanodrop and the quality was assessed using an Agilent Bioanalyzer (Agilent,
148 Cheshire, UK. The two-color low input Quick Amp Labelling kit (Agilent) was used to Cy3- or Cy5-
149 fluorescently label cRNA samples, which were then hybridized to SurePrint G3 Human Gene
150 Expression 60K GeneChip microarrays (Agilent) according to the manufacturer's instructions.
151 Hybridization intensity was quantified via a SureScan Microarray Scanner (Agilent). Microarray
152 data are deposited at Gene Expression Omnibus, Series **GSE153342**.

153

154 Individual channel intensities from the GeneChip data were extracted independently and analysed
155 as separate observations [4].

156

157 **2.5 Statistical analyses**

158 Clinical data were analysed using 'R' Language and Environment for Statistical Computing 3.5.2.
159 Fishers, Chi-squared and Kruskal Wallis tests of significance were used for categorical data. Mann-
160 Whitney U tests of significance were used for continuous data.

161

162 Expression data were analysed using 'R' Language and Environment for Statistical Computing
163 3.5.2. Pre-processing, log-2 transformation and normalisation were performed using the Agilp
164 package [5]. Microarrays were run using two batches of microarray slides and Principal
165 Component Analysis identified an associated batch effect. Batch correction was performed using
166 the COmBat function in the Surrogate Variable Analysis (sva) package in R [6, 7]. To minimise the
167 potential influence of batch correction on subsequent clustering analyses, no reference batch was
168 used and independent COmBat-corrections were performed for each dataset of interest
169 (individual PAXgene, TB1 and TB2 tube datasets and a combined TB1/ TB2/ negative tube dataset).
170 Post-Combat correction PCA plots were undertaken to confirm the removal of the batch effect and
171 identify outliers.

172

173 Differential gene expression analysis was performed using the limma package in R [8] which uses
174 linear models. Where paired samples were available and analysis was relevant, paired t-tests were
175 performed, with this being stated in the results. Adjustment for false discovery rate was
176 performed using Benjamini-Hochberg (BH) correction with a significance level of adjusted p-value
177 <0.05.

178

179 Prior to longitudinal analyses, the gene expression set was filtered to remove noise. Lowly
180 expressed transcripts for which expression values did not exceed a value of 6 for any of the
181 samples, were removed. Transcripts with extreme outlying values were removed, which were
182 defined as values $< (\text{Quartile1} - [3 * \text{Inter-Quartile Range}])$ or $> (\text{Quartile3} + [3 * \text{Inter-Quartile}$
183 $\text{Range}])$. Transcripts with the greatest temporal and interpersonal variability were then selected
184 based on their variance, with those transcripts with variance >0.1 taken forwards to the
185 longitudinal analysis. X-chromosome transcripts which were significantly differentially expressed
186 with gender at V1, V2 and/ or V3 were identified using linear models in limma (BH corrected p
187 value < 0.05) and were excluded, as were Y-chromosome transcripts.

188
189 Unsupervised longitudinal clustering analyses were performed using the BClustLong package in 'R'
190 [9], which uses a Dirichlet process mixture model for clustering longitudinal gene expression data.
191 A linear mixed-effects framework is used to model the trajectory of genes over time and it bases
192 clustering on the regression coefficients obtained from all genes. 500 iterations were run (thinning
193 by 2, so 1000 iterations in total).

194
195 Longitudinal differential gene expression analyses were performed using the MaSigPro package in
196 R [10]. MaSigPro follows a two-step regression strategy to find genes with significant temporal
197 expression changes and significant differences between groups. Coefficients obtained in the
198 second regression model are then used to cluster together significant genes with similar
199 expression patterns. Adjustment for false discovery rate was performed using BH correction with a
200 significance level of adjusted p-value <0.05 . Given the three timepoints from the IGRA+ individuals
201 and the two timepoints from the healthy control groups, we employed both quadratic and linear
202 approaches to account for all the potential curve shapes in the gene expression data.

203
204 Estimations of relative cellular abundances were calculated from the normalised full gene
205 expression matrix (58,201 gene probes) using CibersortX [11], which uses gene expression data to
206 deconvolve mixed cell populations. We used the LM22 [12] leukocyte gene signature matrix as
207 reference, that comprises 22 different immune cell types, and ran 1,000 permutations. Total
208 monocyte fraction was calculated as the sum of the fractions of monocytes, macrophages and
209 dendritic cells. Total lymphocyte fraction was calculated as the sum of B cells, Plasma cells, CD8+ T
210 cells, CD4+ T cells, Helper follicular T cells, Regulatory T cells, Gamma delta T cells, and NK cells. A

211 polynomial model (degrees of freedom = 2) was fitted in R to estimate relationships between the
212 monocyte: lymphocyte ratio and time, in IGRA+ subgroups A and B.

213

214 **3 RESULTS**

215

216 **3.1 Recruitment of participants**

217 Thirty adult IGRA-positive (IGRA+) participants were recruited to the study in the period October
218 2016 to January 2018, of whom 20 took a 12-week course of daily combined rifampicin/ isoniazid
219 (RH) as preventive therapy (PT) and completed study follow-up. Adult IGRA-negative (IGRA-)
220 healthy volunteers were recruited to the study and completed a two-week course of daily RH.
221 After quality control and pre-processing, 18 IGRA+ individuals and 4 IGRA- healthy controls were
222 taken forward for comparator analyses (Figure 1 and Supplementary figures 1 and 2). Recent
223 exposure to drug-susceptible pulmonary TB was confirmed for 15/18 IGRA+s. There were no
224 significant differences in age, gender, ethnicity or BCG status between the 18 IGRA+s and 4 IGRA-
225 healthy controls (Table 1).

226

227 **3.2 Comparing gene expression profiles for IGRA+ versus IGRA- participants**

228 First, we evaluated whether there were discernable differences in gene expression between the
229 IGRA+ participants and IGRA- healthy controls, using linear models[8]. In the unstimulated
230 PAXgene blood samples, no transcripts were found to be significantly differentially expressed
231 (SDE) between the IGRA+ and IGRA- participants at baseline (V1) or visit 2 (V2) (Benjamini-
232 Hochberg [BH] corrected p value < 0.05).

233

234 In this study, QuantIFERON-TB Gold Plus TB1 and TB2 tubes were used to stimulate whole blood.
235 While both tubes contain peptides from ESAT-6 and CFP-10 *Mycobacterium tuberculosis* (*Mtb*)
236 antigens, the TB1 tube peptides are designed to stimulate CD4+ T cells, and the TB2 peptides to
237 stimulate both CD4+ and CD8+ T cells [13]. In contrast to the PAXgene tube whole blood samples,
238 in the TB1-stimulated samples, 123 transcripts were SDE between IGRA+ and IGRA- individuals in
239 the baseline (V1) samples and 93 were SDE between IGRA+ and IGRA- individuals in the V2
240 samples (BH corrected p value < 0.05) (Figure 2A and 2B and listed in Supplementary File 1). In the
241 TB2-stimulated blood samples, when IGRA+ individuals were compared to IGRA-, 43 transcripts
242 were found to be SDE in the V1 samples and 86 in the V2 samples. (BH corrected p value < 0.05)

243 (Figure 2C and 2D and listed in Supplementary File 1). In summary, in vitro stimulation was
244 necessary to distinguish the IGRA+ group from the IGRA- group.

245

246 **3.3 Effects of stimulation on whole blood gene expression**

247 In addition to the TB1 and TB2 *Mtb*-peptide-containing tubes, the QuantiFERON-TB Gold Plus kit
248 also includes a “negative” tube which contains no mycobacterial antigen peptides We assessed
249 the effects of stimulation by comparing gene expression in the TB1- and TB2- stimulated tubes
250 versus the negative tube at visit 1, using paired t-tests. In the IGRA+ group, when TB1 tube
251 samples were compared to the negative tube, 3578 transcripts were SDE, while 3217 transcripts
252 were SDE in the TB2 tube samples versus the negative tube samples (BH corrected p value < 0.05),
253 2495 of which overlapped with the TB1 comparison (Supplementary figure 3A and 3B; SDE
254 transcripts listed in Supplementary File 2). No genes were found to be SDE for the TB1- versus
255 TB2-stimulated samples comparison.

256

257 In the IGRA- healthy controls, 37 transcripts were SDE in the TB1-stimulated samples compared to
258 the negative tubes at visit 1 whereas just four transcripts were SDE in the TB2-stimulated samples
259 (BH corrected p value < 0.05) (Supplementary figure 3C and 3D; SDE transcripts listed in
260 Supplementary File 3).

261

262 **3.4 Filtering the gene expression dataset**

263 Analyses were focused on the stimulated samples, as there had been no detectable differences
264 between the IGRA+ and IGRA- participants in the unstimulated PAXgene samples. As described
265 above, stimulation induced changes in gene expression in the IGRA- healthy controls, with a higher
266 number of SDE genes observed with TB1-stimulation than TB2- stimulation, suggesting a greater
267 non-specific effect independent of *Mtb* infection in the TB1 stimulation. We were concerned these
268 non-specific effects could provide interference, so focused on the TB2-stimulated samples for the
269 next stage of the analysis.

270

271 The gene set was filtered to eliminate noise. Expression values of the 58,201 transcripts ranged
272 from 4.4 to 18.7, so a conservative noise threshold of 6 was chosen. Of the remaining 34,110
273 transcripts, those with the greatest variability between participants and over time were selected

274 for the analysis as described in 2.5. Through this process, a dataset with the “most variable genes”
275 was generated for the TB2-stimulated samples (474 transcripts, listed in Supplementary File 4).
276

277 **3.5 Clustering analysis of longitudinal gene expression**

278 We hypothesised that the IGRA+ group is heterogeneous, containing individuals with viable
279 mycobacteria who would demonstrate a transcriptomic response to PT, and IGRA+ individuals
280 without viable mycobacteria, who would not demonstrate a transcriptomic response to PT and
281 would more closely resemble the healthy control IGRA- group. To unmask the PT-specific
282 transcriptomic responses, we sought to stratify the IGRA+ group of individuals in an agnostic way.
283 We employed unsupervised clustering analysis of longitudinal gene expression in the 18 IGRA+
284 patients and the 4 IGRA- controls, aiming to identify IGRA+ subgroups, using the most variable 474
285 transcripts in the TB2-stimulated dataset. The BClustLong package in ‘R’ [14] was utilised, which
286 uses a linear mixed-effects framework to model the trajectory of genes over time and bases
287 clustering on the regression coefficients obtained from all genes.

288
289 This longitudinal clustering analysis revealed two subgroups of IGRA+ participants. One subgroup
290 of IGRA+s (IGRA+ subgroup A, N=12) clustered with the four healthy controls (Cluster 1),
291 suggesting their gene expression over time was more similar to this *Mtb*-unexposed IGRA-
292 population than it was to the remaining IGRA+s (IGRA+ subgroup B, N=6) who formed Cluster 2.
293 There were no significant differences in age, gender, ethnicity, BCG vaccination status or the
294 IGRA+ participants’ TB contact history between clusters 1 and 2 (Table 2).

295

296 **3.6 Longitudinal differential gene expression analysis**

297 In order to unravel the underlying blood transcriptomic differences between the two cluster
298 groups generated by the unsupervised clustering, we performed longitudinal differential gene
299 expression analysis using MaSigPro package in R [10]. MaSigPro identifies genes with significant
300 temporal expression changes and genes which are significantly differentially expressed between
301 groups.

302

303 Of the 474 transcripts in the dataset, 117 transcripts corresponding to 109 genes, were SDE over
304 time between the two patient groups (with degrees of freedom=1 capturing linear trends, BH
305 corrected p value < 0.05, listed in Supplementary File 5), while 2 of these genes had significant

306 linear terms associated with time (*P2RY6*, *SLC2A3*). Setting the degrees of freedom to 2, 69 out of
307 the 117 genes were SDE over time between the two cluster groups (BH corrected p value < 0.05,
308 listed in Supplementary File 5), while 4 of these genes (*MSR1*, *MT1CP*, *IGHG3*, *IGHG1*) also had
309 significant linear and quadratic terms associated with time. In comparing Cluster 1 versus Cluster
310 2, when one of the clusters is heterogeneous (IGRA+ subgroup A plus IGRA- healthy controls), it is
311 expected that some of the differences will be due to the IGRA+ subgroup B versus IGRA-
312 comparison and not the IGRA+ subgroup B vs IGRA+ subgroup A comparison.

313

314 **3.7 Biological relevance of the significantly differentially expressed genes**

315 The biological relevance of the 117 transcripts significantly differentially expressed over time
316 between the two patient cluster groups was investigated. Around one quarter of these SDE genes
317 have been previously reported in transcriptomics studies comparing blood from TB patients with
318 healthy controls (31 transcripts, 25 genes) or with other diseases (9 transcripts, 7 genes) [15-21];
319 (Supplementary File 5). Functional classification of these genes using PANTHER [22, 23] revealed
320 that 44/84 of the coding genes encode proteins with specific immunological functions, including
321 cytokines, cytokine receptors and cytokine signaling (12), chemokines and chemokine-like proteins
322 (11), immunoglobulins (9), immune cell receptors (4), antimicrobial peptides (3), complement (1)
323 and antigen presentation (1) (Supplementary File 6).

324

325 Coefficients obtained using MaSigPro were used to cluster significant genes with similar
326 longitudinal expression patterns (Figure 3). Often the proteins contained within a gene set had
327 similar function, such as the CXC chemokines CXCL9, 10 and 11 in gene set 2 which were more
328 highly expressed in patient cluster 2 and increased at V2, and the pro-inflammatory NF- κ B
329 transcription factor-inducing proteins IFN γ , IL-1R associated kinase 2 (IRAK2) and TNF superfamily
330 member 15 (TNFSF15) in gene set 4, which were more highly expressed in patient cluster 2 and
331 decreased through PT. BATF2, GCH1 and GBP3 all grouped in gene set 9, with consistently higher
332 expression in patient cluster 2. Gene expression was higher in patient cluster 1 in only one gene
333 set (gene set 3).

334

335 Biological pathways analysis was performed using Reactome pathway knowledgebase [24], with
336 80/117 transcripts successfully mapping to the database. Eleven pathways had significant over-
337 representation of transcripts within our dataset (BH corrected p value < 0.05; listed in
338 Supplementary file 7): these were all related to the immune system and encompassed pathways

339 related to chemokine receptor binding, cytokine signaling – including IL10, TNF and regulatory T
340 cells, metal ion binding and Complement cascade activation. There were a further 39 pathways
341 with borderline over-representation: these largely encompassed biological functions related to
342 innate immunity, antimicrobial peptides, phagocytosis, intracellular infection, and further cytokine
343 signaling and Complement activation pathways.

344

345 **3.8 Differing cellular responses to preventive therapy**

346 Relative cellular abundances were estimated from the gene expression data using CibersortX [11].
347 The estimated abundances of monocytes and lymphocytes were used to calculate the monocyte:
348 lymphocyte ratio (MLR) for the two cluster groups at all three visits. At visits 1 and 3, the MLRs
349 were similar between Clusters 1 and 2. However, at Visit 2, they were higher in Cluster 2 (median=
350 0.52) compared to Cluster 1 (median= 0.29, $p=0.03$). This difference at Visit 2 remained when the
351 IGRA- healthy controls were removed from the analysis, with the MLR higher in IGRA+ subgroup B
352 (median=0.52) compared to subgroup A (median=0.35, $p=0.04$) (Figure 4A).

353

354 Using a second-degree polynomial model, the MLR was found to change over the time-course of
355 the study period in IGRA+ subgroup B, and was close to the threshold of significance (linear term
356 $p=0.07$, quadratic term $p=0.06$). This was not observed in IGRA+ subgroup A (linear term $p=0.6$,
357 quadratic term $p=0.8$) (Figure 4B and C).

358

359 The relative abundances of other cell types including total monocytes, total lymphocytes, total
360 CD4+ T cells and neutrophils were also observed to change with time in IGRA+ subgroup B and not
361 subgroup A (Supplementary figure 4).

362

363 **4 DISCUSSION**

364

365 This analysis has demonstrated that IGRA+ participants could be stratified according to their whole
366 blood transcriptome into two distinct populations, one of which clustered with IGRA-, tuberculosis
367 (TB)-unexposed controls. This separation was not clearly discernible when the transcriptomes of
368 participants were evaluated at baseline in unstimulated whole blood, but rather was unmasked by
369 TB-specific peptide stimulation after 14 days of TB preventive therapy (PT).

370

371 We hypothesised that PT would mediate mycobacterial death in participants for whom IGRA
372 positivity was attributable to ongoing viable *Mycobacterium tuberculosis* (*Mtb*) infection and that
373 the resulting immunological response, detected as a whole blood transcriptomic readout, would
374 differentiate such individuals from a group of IGRA+ participants in whom PT would have no anti-
375 mycobacterial effect due to the absence of viable *Mtb*. Our agnostic clustering approach clustered
376 all four IGRA- healthy controls with a subgroup of IGRA+s (IGRA+ A), which is strongly suggestive
377 that if indeed these clusters do define *Mtb* viability status then the true latent tuberculosis
378 infection (LTBI) participants lie within the other subgroup (IGRA+ B). The genes differentially
379 expressed between the two clusters through PT were predominantly involved in the immune
380 system, particularly related to intracellular infection, inflammation, chemotaxis and cytokine
381 signalling, indicating a biologically plausible specific response in the IGRA+ B subgroup.

382

383 Alternative explanations for the clear separation of these two groups were considered. Rifampicin
384 has important antimicrobial effects against gram-positive organisms and can eliminate upper
385 respiratory tract carriage of gram-negative organisms such as *Neisseria meningitidis* and
386 *Haemophilus influenzae* within 2-4 days. The inclusion of rifampicin/isoniazid treated, IGRA-
387 negative control participants was an attempt to capture and isolate any such non-
388 mycobactericidal effect. In the absence of microbiological sampling and/or microbiome analysis
389 we cannot entirely exclude the possibility that the separation of the groups is attributable to an
390 effect completely unrelated to *Mtb* infection; however two factors which weigh against this
391 alternative explanation are the low prevalence of *N. meningitidis* and *H. influenzae* carriage in this
392 population (<10% combined) and the identification amongst the differentially expressed genes of
393 several genes known to be associated with *Mtb* response pathways. The changes through PT
394 overlapped with reported changes in blood transcriptome during treatment of active TB cases [25,
395 26] and during Isoniazid PT [27]. The monocyte-to-lymphocyte ratio transiently increased only in
396 the IGRA+ B subgroup: this ratio has been linked with TB disease susceptibility and blood
397 transcriptomes [28]. The prevalence of carriage of non-tuberculous mycobacteria in this London-
398 resident population would also be expected to be very low. We considered the possibility that our
399 observations could reflect differences in drug metabolism. Rifampicin induces gene expression
400 changes in hepatocytes [29, 30], but after reviewing this literature and publicly available RNASeq
401 data (Gene Expression Omnibus, Series **GEO139896**) [29] we found no evidence for this (data not
402 shown). This could also be attributed to the fact that our study focused on peripheral blood
403 associated gene expression changes as compared to the liver-derived hepatocytes described in

404 these previously reported studies To further investigate any non-specific effects of Rifampicin and
405 Isoniazid, we also compared gene expression at Visit 2 versus Visit 1 in the healthy controls, using
406 a paired t- test in limma, and found no significantly differentially expressed genes. Finally, we
407 were concerned to exclude all possible artefactual explanations related to sample handling and
408 found no effect association with study site, time to sample processing, study personnel or date of
409 enrolment.

410

411 We contend that interferon gamma release assays (IGRA) and tuberculin skin tests (TST) are mis-
412 represented as tests for LTBI, a term which infers viability of *Mtb* with potential to cause future
413 reactivation disease. We believe that the observation that 90% of individuals with positive testing
414 by IGRA/TST do not develop TB disease is more likely to reflect low frequency of persistent viable
415 (“reactivate-able”) infection than low frequency of breakout of *Mtb* replication from long-term
416 immunological control. The empirical evidence that we present in support of this contention is
417 consistent with recent re-evaluations of epidemiological data which suggest that (1) duration of
418 *Mtb* infection viability is likely to be much shorter than previously believed [31] and that (2)
419 reactivation rates in IGRA or TST positive individuals unprotected by PT undergoing
420 immunosuppressive therapy are much lower than would be expected if such testing represented
421 infection truly capable of reactivation [32]. Emerging mathematical modelling outputs add weight
422 to this paradigm shift, suggesting that a significant proportion of *Mtb*-infected individuals achieve
423 self-clearance, leaving a much smaller population with persisting viable *Mtb* infection than
424 previously assumed [33]. Finally, a precedent for lasting anti-mycobacterial immunological
425 reactivity in the absence of bacterial viability already exists in the form of erythema nodosum
426 leprosum, type II reactions to persistent *M. leprae* antigens which are known to occur years after
427 mycobacterial cure.

428

429 These blood transcriptional responses to PT suggest that around one third of our IGRA+ study
430 participants had true (viable) LTBI. This study was performed in TB contacts with recent exposure,
431 who are an IGRA+ population at high risk of progression. The proportion with viable infection is
432 predicted to be lower with increasing remoteness in time since exposure, for example in migrants
433 now resident in low-incidence countries [31]. The implications for national and global estimates of
434 LTBI prevalence that rely upon IGRA/ TST data are clear and suggest a large overestimation of the
435 size of the global reservoir of potentially reactivatable latent infection; we contend that such data
436 should in future be presented as prevalence of tuberculin sensitivity and that the term LTBI should

437 be used more judiciously. Since all incident reactivation arises from the true LTBI pool, the
438 incidence rate in this subgroup of all IGRA positives will be considerably higher than, for example,
439 the 0.6 per 100 person-years seen in the placebo arm of a recent vaccine trial [34]. The
440 development of tools and strategies to readily identify this true LTBI subgroup would facilitate
441 more efficient targeting of interventions to interrupt reactivation and would accelerate evaluation
442 of novel interventions because the sample size required for future vaccine trials and trials of
443 preventive therapy would be considerably reduced. Evaluations of risk factors associated with
444 infection, premised on the use of IGRA/TST to define infection, have likely been using a very
445 imperfect endpoint with the associated high likelihood of misclassification error.

446
447 The temporal dynamics of the transcriptomic changes are such that evidence of a response can be
448 detected as early as 2 weeks into PT. This raises the possibility of a ‘treat and test’ approach to PT
449 wherein the absence of a specific change in a biomarker (or biomarker profile) at an early time
450 point, say 2 weeks into treatment, could be interpreted as an indication that further treatment will
451 have no effect and can then be discontinued. Recent TB host gene expression studies have shown
452 that biomarker signatures can be shrunk to small sets with the potential to be implemented as
453 diagnostic or prognostic tests in the field [35-37].

454
455 This is the first study to look at longitudinal transcriptomic responses in the blood of IGRA+
456 individuals post-stimulation during the course of PT. Despite its novelty and strengths, it has a
457 relatively modest number of participants. Sequential transcriptomic and cell count differential
458 testing on a larger study population in which defined secondary cases are identified, with a variety
459 of exposure histories and diverse PT regimens (including those under investigation for multidrug-
460 resistant LTBI) will help to elucidate the array of responses encountered. The hunt for predictors
461 of future disease amongst TB- exposed individuals has previously been directed towards
462 identification of biomarkers indicating increased risk, an approach that risks dismissal of future
463 changes in the host environment which it might not be possible to anticipate (e.g. transplant
464 immunosuppression). By removing from the pool of *Mtb*-sensitised participants (IGRA+ or TST+) a
465 significant proportion for whom reactivation is biologically impossible (because no viable *Mtb*
466 infection remains), the scale of the prevention challenge is drastically reduced and a more efficient
467 targeted and nuanced approach can be considered.

468

469 Important implications of a test that can distinguish IGRA+ or TST+ *Mtb* sensitised individuals at
470 zero risk of progression/reactivation include drastic reevaluation of the global burden of LTBI,
471 stratification of preventive therapy and post-exposure vaccine efficacy, higher resolution targeting
472 of LTBI preventive therapy, potential use as a biomarker for efficacy evaluation of novel PT
473 regimens for drug-susceptible and drug-resistant-TB, and PT test of cure.

474

475 **5 CONCLUSION**

476

477 Individuals with immunological memory of a prior encounter with *Mtb* (commonly referred to as
478 LTBI) who are treated with PT demonstrate two different phenotypes of transcriptomic response.
479 We propose that the clear responders are those who had truly viable latent *Mtb* infection, and
480 that the minimal responders, in common with the IGRA-negative, previously unexposed healthy
481 controls, had no viable *Mtb* organisms and were therefore not truly latently TB infected.

482

483 **ACKNOWLEDGEMENTS**

484

485 The authors wish to thank the patients and volunteers who participated in the study. We also
486 thank the clinical staff at Barts Health NHS Trust, Homerton University Hospital Foundation Trust
487 and TB Service North Central London, in particular Dr Heinke Kunst (Barts Health NHS Trust), Prof
488 Graham Bothamley (Homerton University Hospital Foundation Trust) and Prof Marc Lipman (TB
489 Service North Central London) for facilitating recruitment. The authors also wish to thank the
490 research nurses who assisted with this study, including Victoria Dean, Michelle Berin (University
491 College London) and Nirmala Ghimire (Barts Health), as well as Ortensia Vito and Dominic
492 Habgood-Coote (Imperial College London) for help with data analysis.

493

494 This work was supported by a British Infection Association Small Project Research Grant (2016)
495 and a Rosetrees Trust Seed Corn Award (# JS15 / M660). C.B. was funded by an Academic Clinical
496 Fellowship from the National Institute for Health Research (NIHR) (ACF-2012-18-008) and
497 currently receives support from an Imperial 4i Wellcome Trust/ NIHR Imperial BRC Clinical PhD
498 Fellowship. M.K. receives support from the NIHR Imperial College BRC and the Wellcome Trust
499 (Sir Henry Wellcome Fellowship grant no. 206508/Z/17/Z). J.C. receives support from the Medical
500 Research Council Newton Fund (#MR/P017568/1). The funders were not involved in study design,
501 data collection and analysis, decision to publish, nor in preparation of the manuscript.

502

503 **REFERENCES**

504

- 505 [1.] World Health Organisation. Latent tuberculosis infection: updated and consolidated
506 guidelines for programmatic management. Geneva; 2018.
- 507 [2.] Esmail H, Barry CE, 3rd, Young DB, Wilkinson RJ. The ongoing challenge of latent
508 tuberculosis. *Philos Trans R Soc Lond B Biol Sci.* 2014;369(1645):20130437.
- 509 [3.] Whalen CC, Johnson JL, Okwera A, Hom DL, Huebner R, Mugenyi P, et al. A trial of three
510 regimens to prevent tuberculosis in Ugandan adults infected with the human immunodeficiency
511 virus. Uganda-Case Western Reserve University Research Collaboration. *The New England journal
512 of medicine.* 1997;337(12):801-8.
- 513 [4.] Chain B, Bowen H, Hammond J, Posch W, Rasaiyaah J, Tsang J, et al. Error, reproducibility
514 and sensitivity: a pipeline for data processing of Agilent oligonucleotide expression arrays. *BMC
515 Bioinformatics.* 2010;11:344.
- 516 [5.] Chain B. *agilp: Agilent expression array processing package.* R package version 3.14.0. 2018
517 [Available from: <http://bioconductor.org/packages/release/bioc/html/agilp.html>].
- 518 [6.] Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
519 empirical Bayes methods. *Biostatistics.* 2007;8(1):118-27.
- 520 [7.] Leek JT, Johnson WE, Parker HS, Fertig EJ, Jaffe AE, Storey JD, et al. *sva: Surrogate Variable
521 Analysis.* R package version 3.34.0. 2019 [Available from:
522 <https://bioconductor.org/packages/release/bioc/html/sva.html>].
- 523 [8.] Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, et al. *limma* powers differential
524 expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.*
525 2015;43(7):e47.
- 526 [9.] Sun J, Herazo-Maya JD, Kaminski N, Zhao H, Warren JL. *BClustLonG: A Dirichlet Process
527 Mixture Model for Clustering Longitudinal Gene
528 Expression Data.* R package version 0.1.2. 2017 [Available from: [https://CRAN.R-
529 project.org/package=BClustLonG](https://CRAN.R-project.org/package=BClustLonG)].
- 530 [10.] Conesa A, Nueda MJ. *maSigPro: Significant Gene Expression Profile Differences in Time
531 Course Gene Expression Data.* R package version 1.54.0 2018 [Available from:
532 <http://bioinfo.cipf.es/>].
- 533 [11.] Newman AM, Steen CB, Liu CL, Gentles AJ, Chaudhuri AA, Scherer F, et al. Determining cell
534 type abundance and expression from bulk tissues with digital cytometry. *Nat Biotechnol.*
535 2019;37(7):773-82.
- 536 [12.] Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell
537 subsets from tissue expression profiles. *Nat Methods.* 2015;12(5):453-7.
- 538 [13.] Petruccioli E, Chiacchio T, Pepponi I, Vanini V, Urso R, Cuzzi G, et al. First characterization of
539 the CD4 and CD8 T-cell responses to QuantiFERON-TB Plus. *J Infect.* 2016;73(6):588-97.
- 540 [14.] Sun J, Herazo-Maya JD, Kaminski N, Zhao H, Warren JL. A Dirichlet process mixture model
541 for clustering longitudinal gene expression data. *Stat Med.* 2017;36(22):3495-506.
- 542 [15.] Anderson ST, Kaforou M, Brent AJ, Wright VJ, Banwell CM, Chagaluka G, et al. Diagnosis of
543 childhood tuberculosis and host RNA expression in Africa. *The New England journal of medicine.*
544 2014;370(18):1712-23.
- 545 [16.] Berry MP, Graham CM, McNab FW, Xu Z, Bloch SA, Oni T, et al. An interferon-inducible
546 neutrophil-driven blood transcriptional signature in human tuberculosis. *Nature.*
547 2010;466(7309):973-7.
- 548 [17.] Blankley S, Graham CM, Turner J, Berry MP, Bloom CI, Xu Z, et al. The Transcriptional
549 Signature of Active Tuberculosis Reflects Symptom Status in Extra-Pulmonary and Pulmonary
550 Tuberculosis. *PLoS One.* 2016;11(10):e0162220.

- 551 [18.] Bloom CI, Graham CM, Berry MP, Rozakeas F, Redford PS, Wang Y, et al. Transcriptional
552 blood signatures distinguish pulmonary tuberculosis, pulmonary sarcoidosis, pneumonias and lung
553 cancers. PLoS One. 2013;8(8):e70630.
- 554 [19.] Kaforou M, Wright VJ, Oni T, French N, Anderson ST, Bangani N, et al. Detection of
555 tuberculosis in HIV-infected and -uninfected African adults using whole blood RNA expression
556 signatures: a case-control study. PLoS medicine. 2013;10(10):e1001538.
- 557 [20.] Maertzdorf J, Weiner J, 3rd, Mollenkopf HJ, Network TB, Bauer T, Prasse A, et al. Common
558 patterns and disease-related signatures in tuberculosis and sarcoidosis. Proceedings of the
559 National Academy of Sciences of the United States of America. 2012;109(20):7853-8.
- 560 [21.] Ottenhoff TH, Dass RH, Yang N, Zhang MM, Wong HE, Sahiratmadja E, et al. Genome-wide
561 expression profiling identifies type 1 interferon response pathways in active tuberculosis. PLoS
562 One. 2012;7(9):e45839.
- 563 [22.] Mi H, Muruganujan A, Ebert D, Huang X, Thomas PD. PANTHER version 14: more genomes,
564 a new PANTHER GO-slim and improvements in enrichment analysis tools. Nucleic Acids Research.
565 2018;47(D1):D419-D26.
- 566 [23.] Mi H, Muruganujan A, Huang X, Ebert D, Mills C, Guo X, et al. Protocol Update for large-
567 scale genome and gene function analysis with the PANTHER classification system (v.14.0). Nature
568 Protocols. 2019;14(3):703-21.
- 569 [24.] Jassal B, Matthews L, Viteri G, Gong C, Lorente P, Fabregat A, et al. The reactome pathway
570 knowledgebase. Nucleic Acids Res. 2020;48(D1):D498-D503.
- 571 [25.] Bloom CI, Graham CM, Berry MP, Wilkinson KA, Oni T, Rozakeas F, et al. Detectable
572 changes in the blood transcriptome are present after two weeks of antituberculosis therapy. PLoS
573 One. 2012;7(10):e46191.
- 574 [26.] Cliff JM, Lee JS, Constantinou N, Cho JE, Clark TG, Ronacher K, et al. Distinct phases of
575 blood gene expression pattern through tuberculosis treatment reflect modulation of the humoral
576 immune response. J Infect Dis. 2013;207(1):18-29.
- 577 [27.] de Oyarzabal E, García-García L, Rangel-Escareño C, Ferreyra-Reyes L, Orozco L, Herrera
578 MT, et al. Expression of USP18 and IL2RA Is Increased in Individuals Receiving Latent Tuberculosis
579 Treatment with Isoniazid. Journal of Immunology Research. 2019;2019:1297131.
- 580 [28.] Naranbhai V, Fletcher HA, Tanner R, O'Shea MK, McShane H, Fairfax BP, et al. Distinct
581 Transcriptional and Anti-Mycobacterial Profiles of Peripheral Blood Monocytes Dependent on the
582 Ratio of Monocytes: Lymphocytes. EBioMedicine. 2015;2(11):1619-26.
- 583 [29.] Dyavar SR, Mykris TM, Winchester LC, Scarsi KK, Fletcher CV, Podany AT. Hepatocytic
584 transcriptional signatures predict comparative drug interaction potential of rifamycin antibiotics.
585 Scientific reports. 2020;10(1):12565.
- 586 [30.] Williamson B, Dooley KE, Zhang Y, Back DJ, Owen A. Induction of influx and efflux
587 transporters and cytochrome P450 3A4 in primary human hepatocytes by rifampin, rifabutin, and
588 rifapentine. Antimicrobial agents and chemotherapy. 2013;57(12):6366-9.
- 589 [31.] Behr MA, Edelstein PH, Ramakrishnan L. Revisiting the timetable of tuberculosis. BMJ.
590 2018;362:k2738.
- 591 [32.] Behr MA, Edelstein PH, Ramakrishnan L. Is Mycobacterium tuberculosis infection life long?
592 BMJ. 2019;367:l5770.
- 593 [33.] Emery JC, Richards AS, Dale KD, McQuaid CF, White RG, J.T. D, et al., editors. Self-clearance
594 of *Mycobacterium tuberculosis* infection: implications for lifetime risk and population at-risk of
595 tuberculosis disease. . 50th Union World Conference on Lung Health; 2019; Hyderabad, India.
- 596 [34.] Tait DR, Hatherill M, Van Der Meeren O, Ginsberg AM, Van Brakel E, Salaun B, et al. Final
597 Analysis of a Trial of M72/AS01E Vaccine to Prevent Tuberculosis. The New England journal of
598 medicine. 2019;381(25):2429-39.

- 599 [35.] Roe JK, Thomas N, Gil E, Best K, Tsaliki E, Morris-Jones S, et al. Blood transcriptomic
600 diagnosis of pulmonary and extrapulmonary tuberculosis. JCI Insight. 2016;1(16):e87238.
- 601 [36.] Suliman S, Thompson E, Sutherland J, Weiner Rd J, Ota MOC, Shankar S, et al. Four-gene
602 Pan-African Blood Signature Predicts Progression to Tuberculosis. Am J Respir Crit Care Med. 2018.
- 603 [37.] Sweeney TE, Braviak L, Tato CM, Khatri P. Genome-wide expression for diagnosis of
604 pulmonary tuberculosis: a multicohort analysis. The Lancet Respiratory medicine. 2016;4(3):213-
605 24.
- 606
- 607
- 608

Journal Pre-proof

609 **AUTHOR CONTRIBUTIONS**

610

611 Claire Broderick: Conceptualisation, Methodology, Investigation, Formal analysis, Data curation,
612 Writing- original draft, Visualisation, Project administration, Funding acquisition.

613

614 Jackie Cliff: Methodology, Investigation, Formal analysis, Data curation, Resources, Writing-
615 original draft, Funding acquisition.

616

617 Ji-Sook Lee: Investigation, Data curation, Resources.

618

619 Myrsini Kaforou: Methodology, Formal analysis, Writing- review and editing, Visualisation.

620

621 David Moore: Conceptualisation, Methodology, Formal analysis, Writing- original draft,
622 Supervision, Funding acquisition.

Journal Pre-proof

623 **TABLES**

624

625 **Table 1**626 **Subject Characteristics.**

		IGRA+ group	IGRA- Healthy control group
Number		18	4
Age in years: Median (IQR)		34 (28-38)	28 (27-29)
Gender	Male	10 (56%)	3 (75%)
	Female	8 (44%)	1 (25%)
Confirmed recent drug-susceptible TB exposure	Yes	15 (83%)	0 (0%)
	No	3 (17%)	4 (100%)
BCG	Yes	14 (78%)	2 (50%)
	No	2 (11%)	2 (50%)
	Unknown	2 (11%)	0 (0%)
Continent of Birth	Africa	4 (22%)	0 (0%)
	Asia	4 (22%)	0 (0%)
	Australasia	0 (0%)	1 (25%)
	Europe	9 (50%)	2 (50%)
	North America	0 (0%)	1 (25%)
	South America	1 (6%)	0 (0%)
	Unknown	0 (0%)	0 (0%)
Ethnicity	Asian ¹	5 (28%)	2 (50%)
	Black ²	4 (22%)	0 (0%)
	White ³	8 (44%)	2 (50%)
	Other ⁴	1 (6%)	0 (0%)

627 ¹Includes Bengali, Hong Kong, Kurdish, Sri Lankan, Turkish; ²Includes Black African; ³Includes White British,628 Polish, Romanian, White other; ⁴Includes Latin American, Unknown

629

630 **Table 2**631 **Characteristics of Cluster groups 1 and 2.**

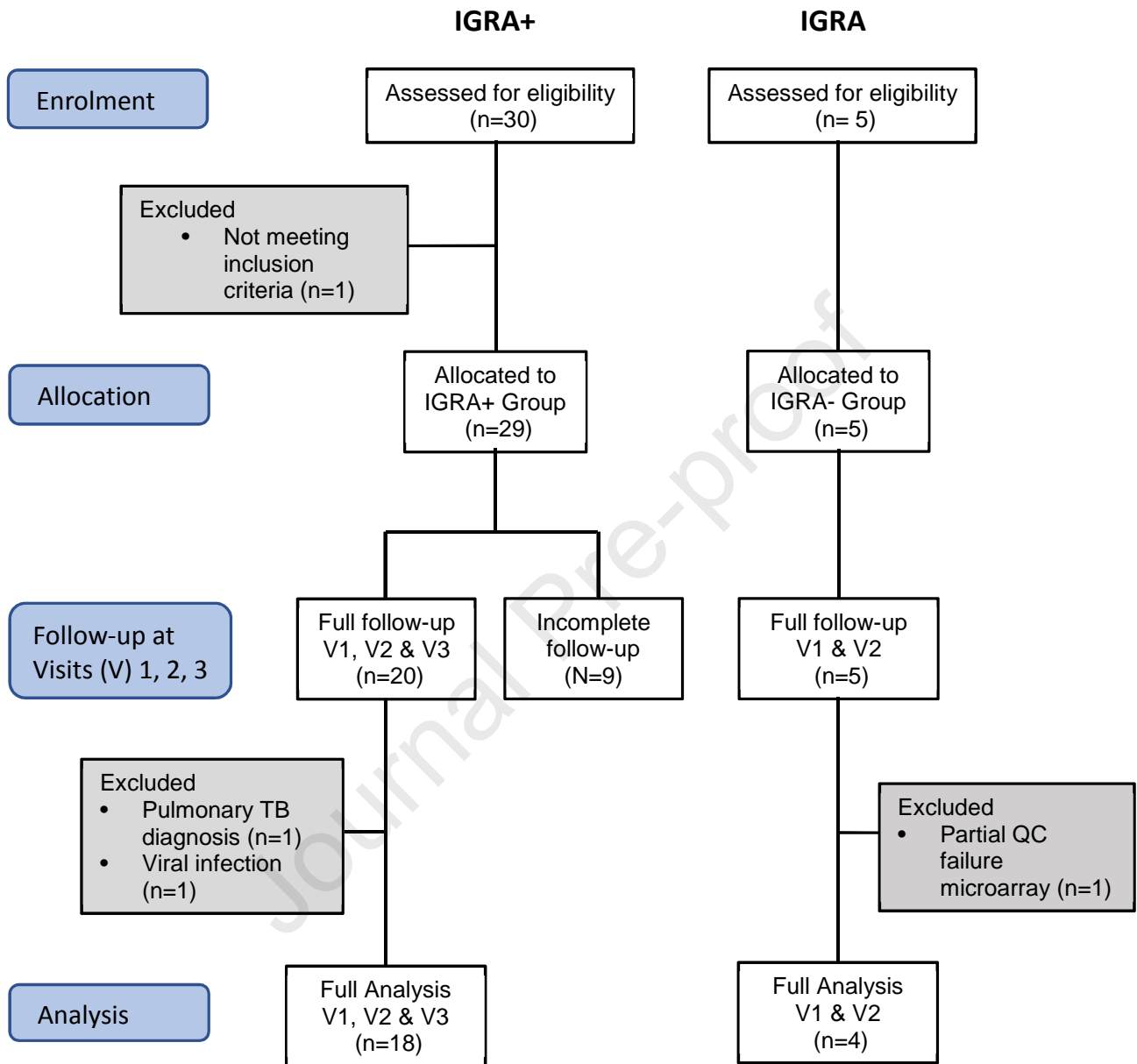
		BClustLong clustering group		
		Cluster 1	Cluster 2	p value
Number of participants		16	6	N/A
Patient IDs		HC51 HC53 HC54 HC55 LTBI1 LTBI2 LTBI3 LTBI5 LTBI7 LTBI9 LTBI12 LTBI15 LTBI16 LTBI27 LTBI28	LTBI6 LTBI10 LTBI14 LTBI22 LTBI23 LTBI30	N/A
Age in years: Median (IQR)		32.5 (24-41)	33.5 (29-38)	0.6
Gender	Male	9 (56%)	4 (66%)	1
	Female	7 (44%)	2 (33%)	
Confirmed recent exposure to DS-TB ¹	Yes	10 (83%)	5 (83%)	1
	No	2 (17%)	1 (17%)	
BCG	Yes	10 (62%)	6 (100%)	0.2
	No	4(25%)	0 (0%)	
	Unknown	2 (13%)	0 (0%)	
Continent of Birth	Africa	3 (19%)	1 (17%)	0.2
	Asia	1 (6%)	3 (50%)	
	Australasia	1 (6%)	0 (0%)	
	Europe	9 (56%)	2 (33%)	
	North America	1(6%)	0 (0%)	
	South America	1 (6%)	0 (0%)	
Ethnicity	Asian ²	4 (25%)	3 (50%)	0.7
	Black ³	3 (19%)	1 (17%)	
	White ⁴	8 (50%)	2 (33%)	
	Other ⁵	1 (6%)	0 (0%)	

632 ¹ for IGRA+ participants only633 ²Includes Bengali, Hong Kong, Kurdish, Sri Lankan, Turkish; ³Includes Black African; ⁴Includes White British,634 Polish, Romanian, White other; ⁵Includes Latin American, Unknown

635

636 **FIGURES:**

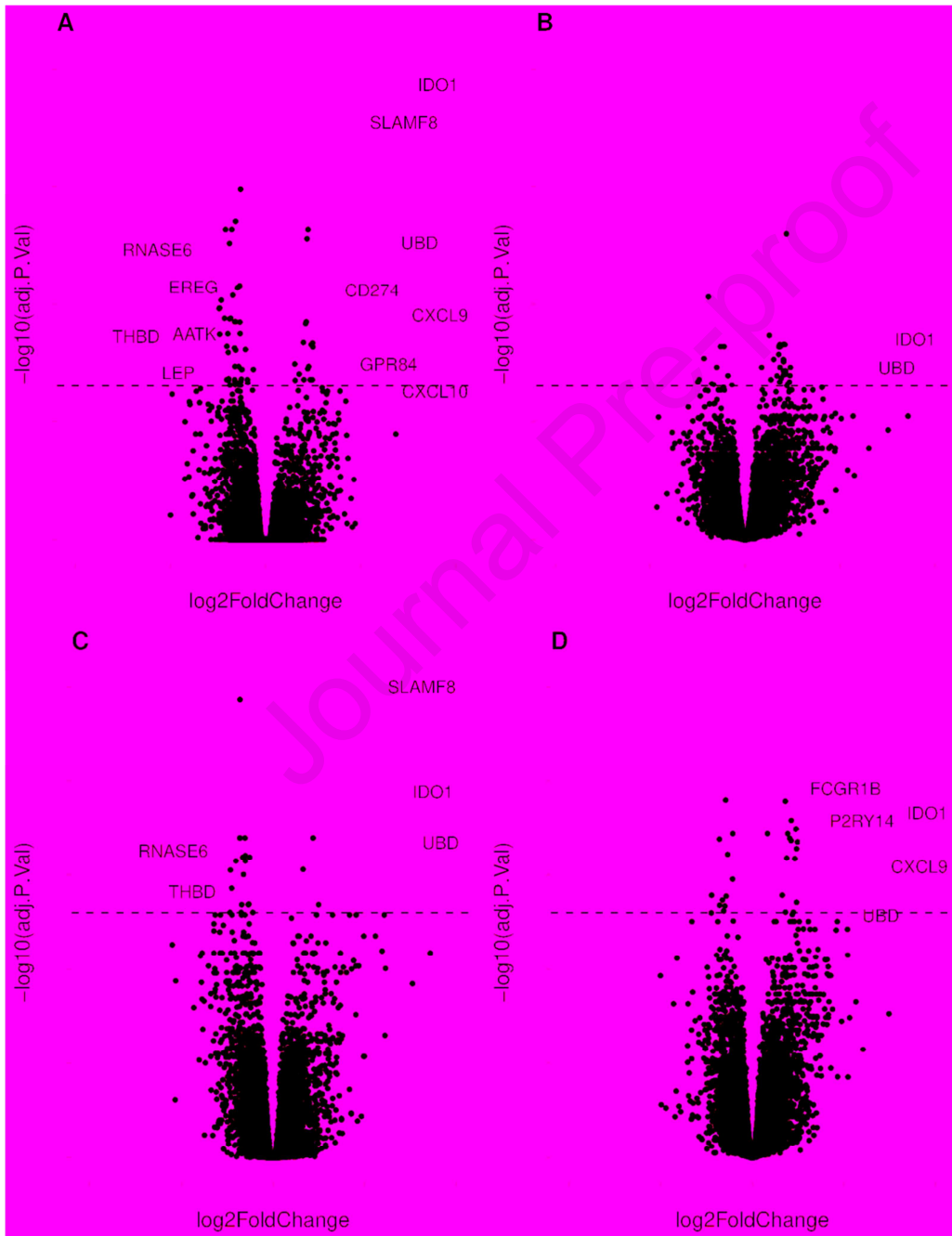
637

638 **Figure 1. Study overview, showing patient numbers and exclusions.**

639

640

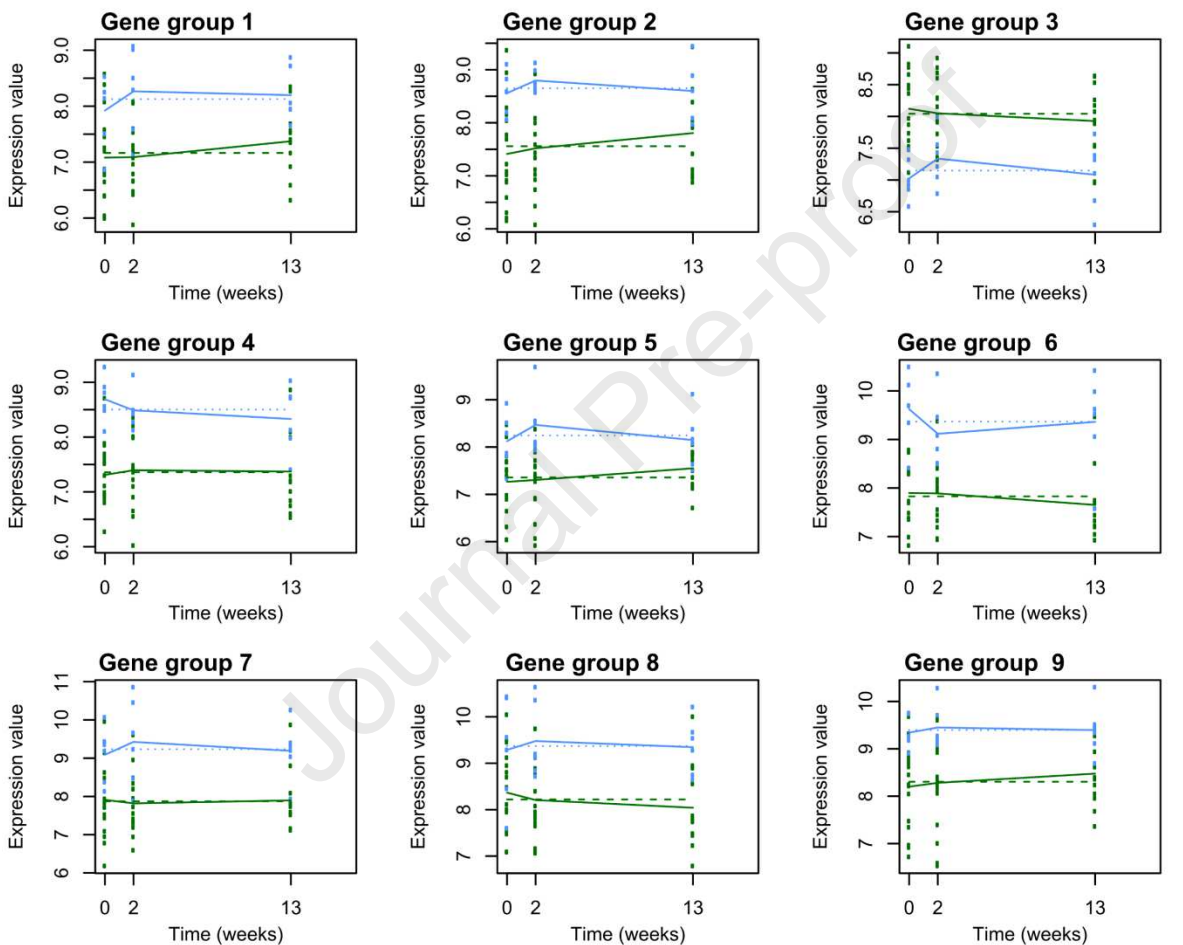
641 **Figure 2. Volcano plots showing genes significantly differentially expressed between IGRA+**
 642 **and IGRA- individuals.** Plots are shown for TB1-stimulated samples at Visit (V) 1 [A] and V2 [B] and
 643 TB2-stimulated samples at V1 [C] and V2 [D]. Genes overexpressed in IGRA+s with log2Foldchange
 644 (LFC) >1 and Benjamini-Hochberg adjusted p value <0.05 are shown in red. Genes underexpressed
 645 in IGRA+ individuals with LFC <-1 and BH adjusted p value <0.05 are shown in blue. Genes with LFC
 646 >2.7 and < -1.7 are annotated with their gene symbols. Dotted line denotes the significance cut-off
 647 (BH adjusted p value <0.05).



648

649

650 **Figure 3. Longitudinal differential gene expression analysis between patient cluster groups 1 and**
 651 **2 in TB2-stimulated whole blood samples.** With 1 degree of freedom, 117/474 transcripts were
 652 SDE over time and between cluster groups 1 and 2 (BH corrected p value < 0.05). The coefficients
 653 obtained were used to group together significant genes with similar longitudinal expression
 654 patterns. MaSigPro identified 9 gene groups. Plots of gene expression against time for these gene
 655 groups are shown for patient cluster groups 1 (green) and 2 (blue). Lines join the median
 656 expression values of the gene groups at each timepoint. The gene symbols are listed for each gene
 657 group.

**Gene groups:**

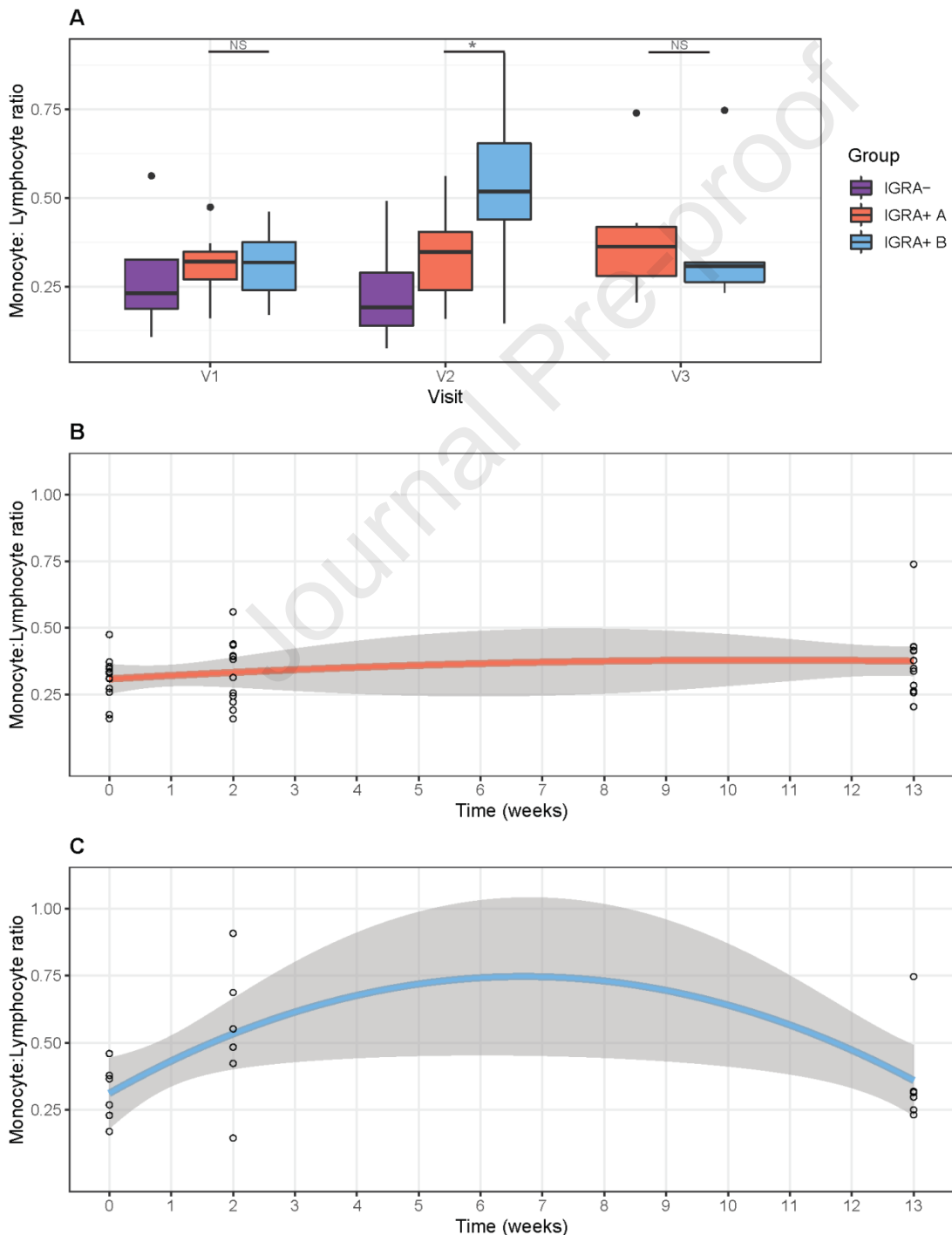
1: CCL8, CTSC, LILRB5, LINC00996, RETN, P2RY6, RGL1, SDSL
 2: BATF3, CD38, CH25H, CISH, CXCL9, CXCL10, CXCL11, EBI3, EXOC3L4, GGT5, HLA-DQA1, IL2RA, LINC02528, PIR, PROCR, SCIMP, TNFRSF4, UBD
 3: ANKRD20A9P, CMTM2, CXCR2, HIF1A-AS3, HIST2H2B, KRT72, LEP, LOC100507006, LOC254896, NFXL1, RNASE6, SLC2A3, TMEM45A, TNFRSF10C, TNFRSF8, VCAN, XLCO_I2_005490
 4: C3, DEFB1, EDN1, HLS4, IFNG, INHB4, IRAK2, GDF15, LINC02154, MFSD2A, MIR155HG, NEURL3, OLR1, PDE4DIP, PHLDA2, PI3, SLPI, TNFSF15
 5: CCL7, CTSL, CTSLP4, CTSLP8, MSR1
 6: CCL3, CCL3L3, CCL4L2, GPR84, HBEGF, HSD11B1, IL1A, LAMB3, Lnc-EGR-3, LOC388780, MDGA1
 7: LINC0109, MT1CP, MT1E, MT1G, MT1H, MT1M, MT1X
 8: CCL18, CPNE5, IGHA2, IGHG1, IGHG3, IGHG4, IGJ, IGKV1-16, IGLC6, IGLL5, SMIM1
 9: BATF2, C11orf85, FLJ27255, GCH, GBP3, Inc-PIK3CG-4, LOC731424, MIR3945HG, SLAMF8

Cluster group 1 ———
 Cluster group 2 ———

658

659

660 **Figure 4: Longitudinal changes in monocyte: lymphocyte ratio through preventive therapy in**
 661 **IGRA+ subgroups A and B.** Cibersortx was used to estimate the abundance of monocytes and
 662 lymphocytes in the TB2-stimulated whole blood samples at each visit, and the monocyte:
 663 lymphocyte ratio was calculated. (A) Boxplots showing the Monocyte: Lymphocyte ratios at Visits
 664 1, 2 and 3 for IGRA- healthy controls and IGRA+ groups A and B. NS denotes $p > 0.05$, * denotes
 665 $p \leq 0.05$. Scatterplots showing the change in Monocyte: lymphocyte ratio over the time-course of
 666 the study period for (B) IGRA+ subgroup A and (C) IGRA+ subgroup B, where Visit 1 is 0 weeks, Visit
 667 2 is 2 weeks and Visit 3 is 13 weeks, with 90% confidence intervals shown.



668

669 **SUPPLEMENTARY DATA**

670

671 **Supplementary Files**

672

673 **Supplementary File 1:** Significantly differentially expressed (SDE) transcripts IGRA+ vs IGRA- in TB1
674 tubes at Visit (V) 1 and V2 and in TB2 tubes at V1 and V2.

675 **Supplementary File 2:** SDE transcripts TB1 vs negative tube, TB2 vs negative tube at V1, in IGRA+.

676 **Supplementary File 3:** SDE transcripts TB1 vs negative tube, TB2 vs negative tube at V1, in IGRA-.

677 **Supplementary File 4:** 474 most variable transcripts (TB2-stimulated samples).

678 **Supplementary File 5:** MaSigPro results: transcripts SDE though time, Cluster 1 vs Cluster 2.

679 **Supplementary File 6:** Functional classification of transcripts differentially expressed between
680 Cluster 1 and Cluster 2.

681 **Supplementary File 7:** Results of biological pathways analysis using Reactome pathway
682 knowledgebase.

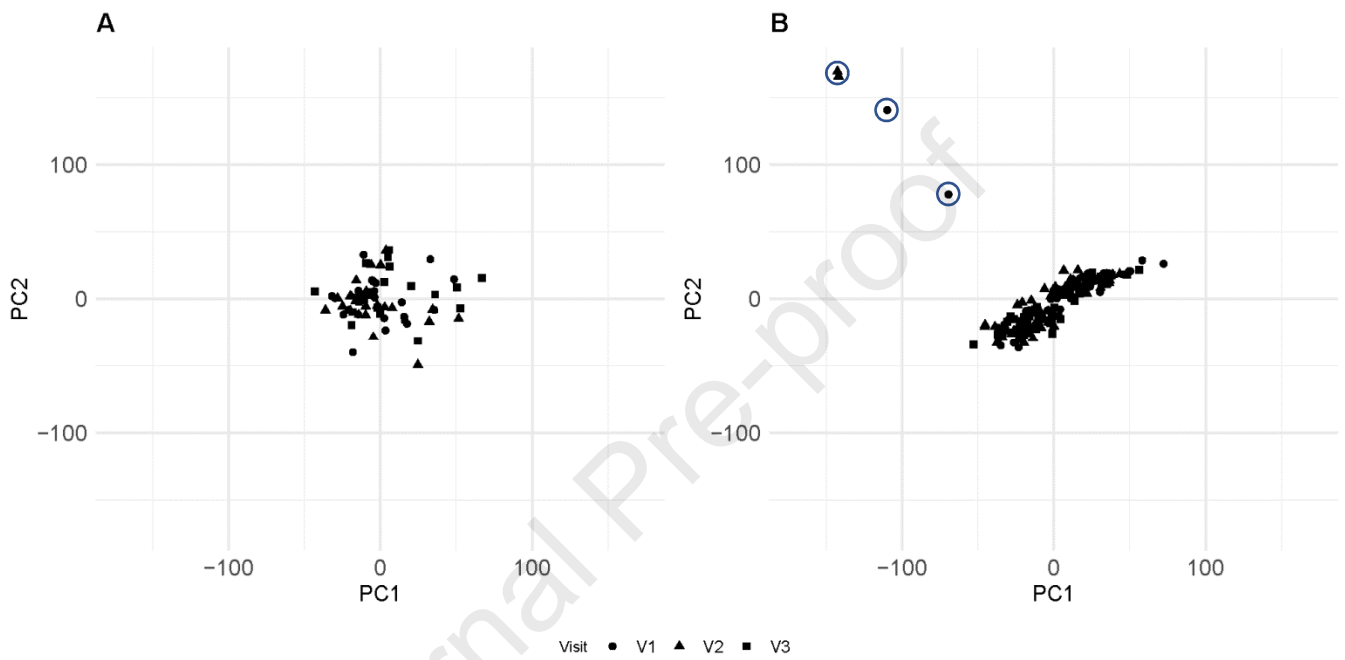
683

684

685 **Supplementary Figures**

686

687 **Supplementary Figure 1. Principle component analyses of the initial gene expression sets.** Plots
688 showing dimensions 1 and 2 of the principle component analyses of the PAXgene samples (A) and
689 the stimulated samples (B) before ComBat correction. In the stimulated samples, a healthy control
690 (HC52) was an outlier in dimensions 1 and 2 (circled) and this persisted after batch correction (not
691 shown), so HC52 was excluded from the subsequent analyses.

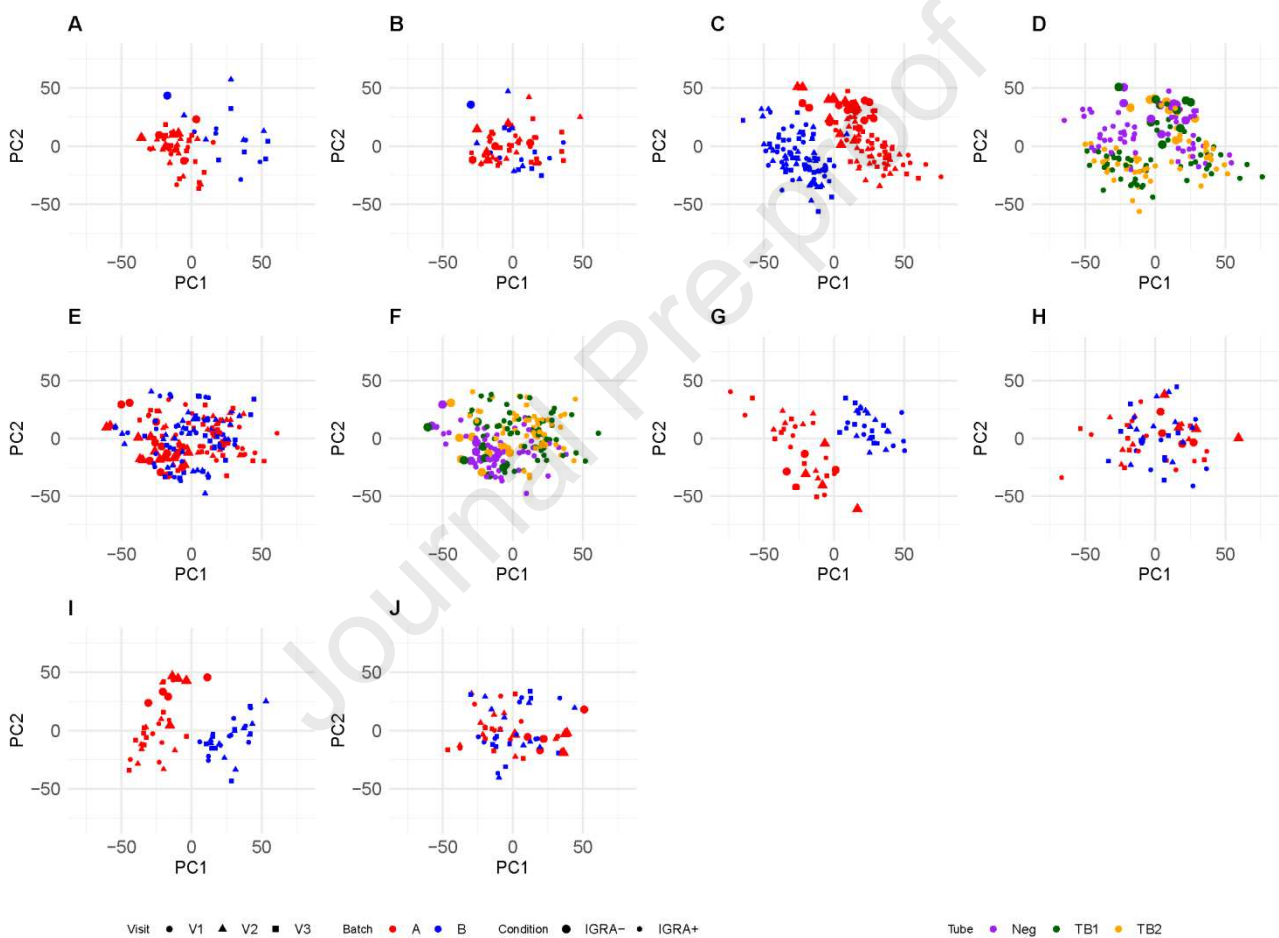


692

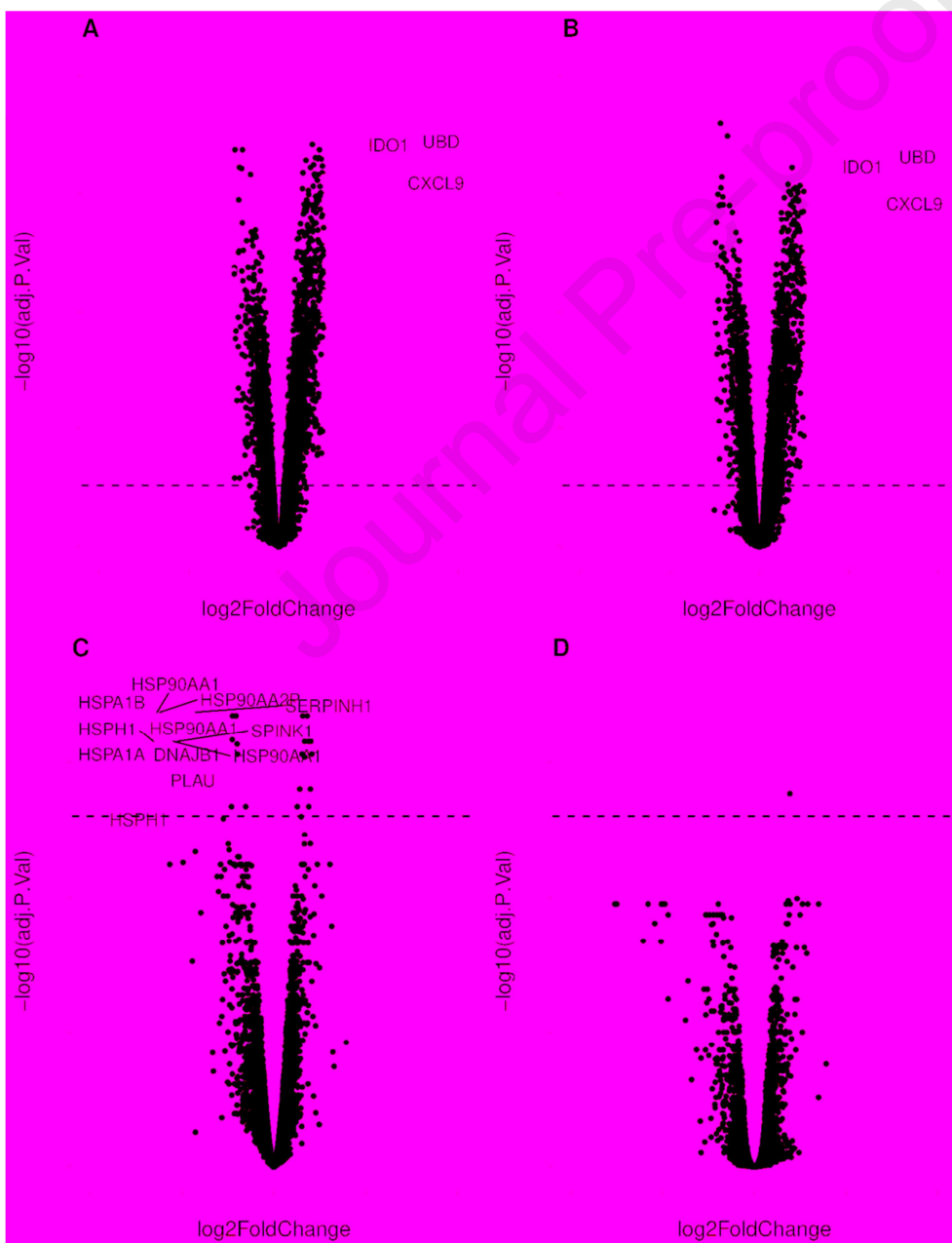
693

694 **Supplementary Figure 2. Principle component analyses of the gene expression sets before and**
 695 **after ComBat.** Gene expression data from 18 IGRA+ and 4 IGRA- participants were included in the
 696 final analyses. Batch correction was performed with ComBat. Plots showing dimensions 1 and 2 of
 697 the principle component analyses of the PAXgene tube samples before (A) and after ComBat (B);
 698 all stimulated samples (TB1, TB2 and Negative) before (C, D) and after ComBat (E, F) with C and E
 699 showing batch differentiation and D and F showing tube differentiation; TB1 samples before (G)
 700 and after Combat (H); TB2 samples before (I) and after Combat (J). Batch, visit, IGRA status and
 701 QuantiFERON TB Gold plus tube are provided.

702



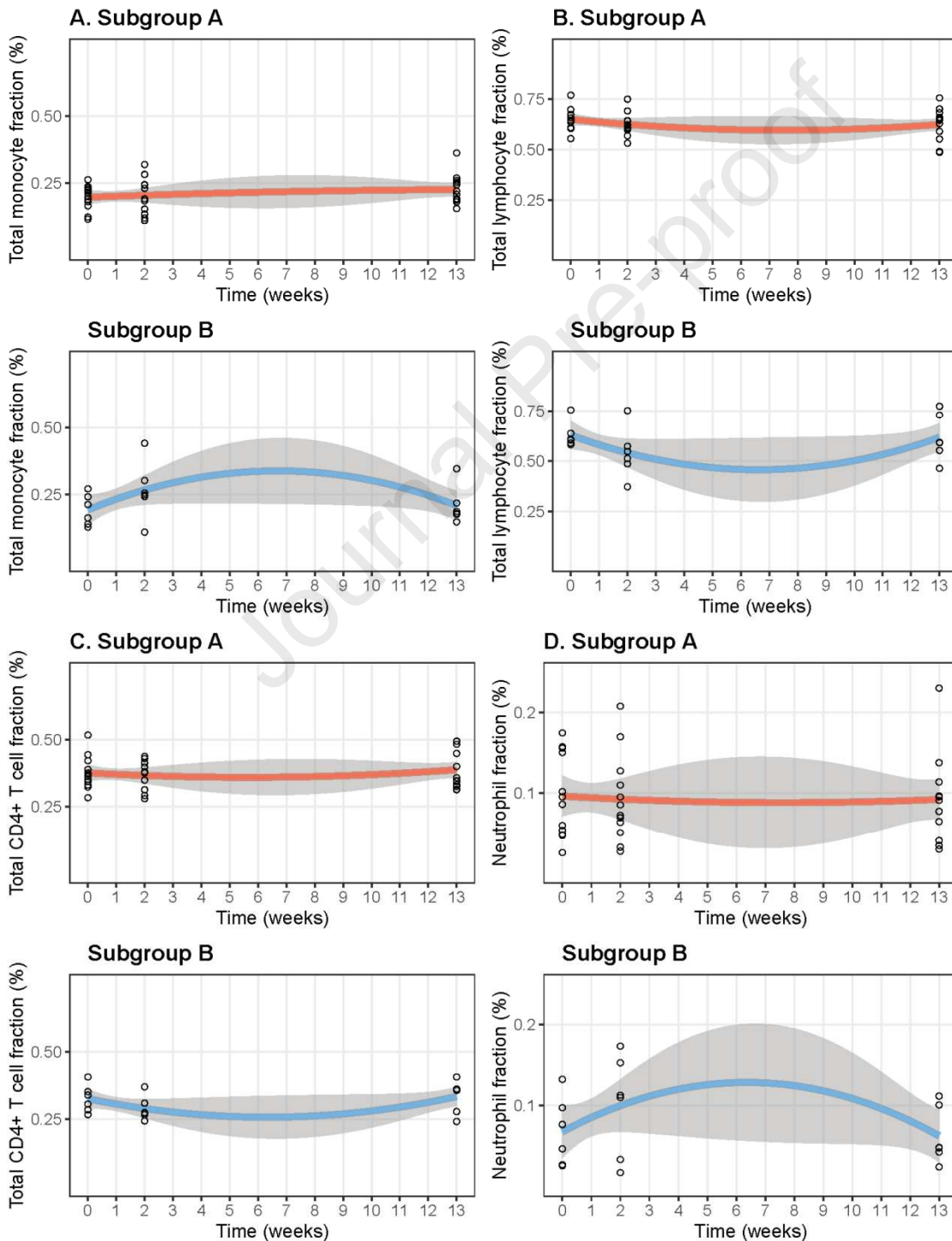
705 **Supplementary Figure 3. Volcano plots showing genes significantly differentially expressed**
 706 **between stimulated (QuantiFERON Gold Plus TB1 and TB2 tubes) and unstimulated**
 707 **(QuantiFERON Gold Plus negative tubes) blood samples.** Plots are shown for IGRA+ subjects,
 708 comparing TB1 vs. negative tube samples (A), and TB2 vs. negative tube samples (B) at visit 1. Also
 709 shown are plots for IGRA- subjects, comparing TB1 vs. negative tube samples (C), and TB2 vs.
 710 negative tube samples (D) at visit 1. Genes overexpressed in stimulated blood with
 711 $\log_2\text{Foldchange (LFC)} >1$ and BH adjusted p value <0.05 are shown in red. Genes underexpressed
 712 in stimulated blood with $\text{LFC} <-1$ and BH adjusted p value <0.05 are shown in blue. Genes with LFC
 713 >2.7 and <-1.7 are annotated with their gene symbols. Dotted line denotes the significance cut-off
 714 (BH adjusted p value <0.05).



715

716

717
 718 **Supplementary Figure 4. Longitudinal changes in cellular populations through preventive**
 719 **therapy in IGRA+ subgroups A and B.** Cibersortx was used to estimate the abundance of different
 720 cell types in the TB2-stimulated whole blood samples at each visit. Scatterplots showing the
 721 change cellular fractions over the time-course of the study period in IGRA+ subgroups A and B for
 722 Total monocyte fraction (A), Total lymphocyte fraction (B), Total CD4+ T cell fraction (C),
 723 Neutrophil fraction (D). Visit 1 is 0 weeks, Visit 2 is 2 weeks and Visit 3 is 13 weeks, with 90%
 724 confidence intervals shown.



725