

1 **Unbiased definition of a shared T-cell receptor motif enables population-based studies of**
2 **tuberculosis.**

3

4 DeWitt, W.S.^{1*}, Quan, K.K.², Wilburn, D.³, Sherwood, A.¹, Vignali, M.¹, De Rosa, S.C.⁴, Day
5 C.L.⁵, Scriba T.J.⁶, Robins H.S.^{1,4}, Swanson W.³, Emerson R.O.¹, Seshadri, C.²

6

7 1. Adaptive Biotechnologies, Seattle, USA

8 2. Department of Medicine, University of Washington, Seattle, USA

9 3. Department of Genome Sciences, University of Washington, Seattle, USA

10 4. Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle,
11 USA

12 5. Department of Microbiology and Immunology, Emory University School of Medicine and
13 Emory Vaccine Center, Atlanta, USA

14 6. South African Tuberculosis Vaccine Initiative and Institute of Infectious Disease and
15 Molecular Medicine, Division of Immunology, Department of Pathology, University of Cape
16 Town, Cape Town, South Africa.

17 * Current affiliation: Computational Biology Program, Fred Hutchinson Cancer Research Center,
18 Seattle, USA

19

20 Corresponding author:

21 Chetan Seshadri, M.D.

22 University of Washington Medical Center

23 750 Republican Street, Suite E663

24 Seattle, WA 98109

25 Email: seshadri@u.washington.edu

26 Phone: 206-543-6709; Fax: 206-616-4898

27 Running Title: Shared TCR Motifs in Tuberculosis

28

29 **Brief Summary**

30 We used human genetics and immunosequencing to define a shared T-cell receptor motif that is
31 specific for a mycobacterial lipid antigen and associated with tuberculosis independently of host
32 genetic background.

33

34 **Conflict of Interest Statement**

35 H.S.R has full-time employment, equity ownership and patents & royalties at Adaptive
36 Biotechnologies Corporation. R.O.E, M.V, A.S, and W.S.D have full-time employment and
37 equity ownership at Adaptive Biotechnologies Corporation

38 **ABSTRACT**
39

40 Peptide-specific T cells that are restricted by highly polymorphic major histocompatibility
41 complex (MHC) proteins express diverse T-cell receptors (TCRs) that are rarely shared among
42 unrelated individuals. T-cells can also recognize bacterial lipid antigens that bind the relatively
43 non-polymorphic CD1 family of proteins. However, genetic variation in human CD1 genes and
44 TCR diversity expressed by CD1-restricted T-cells have not been quantitatively determined.
45 Here, we show that CD1B is nearly nucleotide-identical across all five continental ancestry
46 groups, providing evidence for purifying selection during human evolution. We used CD1B
47 tetramers loaded with a mycobacterial glycolipid antigen to isolate T-cells from four genetically
48 unrelated South African adults and cataloged thousands of TCRs from *in-vitro* expanded T-cells
49 using immunosequencing. We identified highly conserved motifs that were co-expressed as a
50 functional heterodimer and significantly enriched among tetramer-positive T-cells sorted directly
51 from peripheral blood. Finally, we show that frequencies of these TCR motifs are increased in
52 the blood of patients with active tuberculosis compared to uninfected controls, a finding that is
53 confirmed by ex-vivo frequencies of tetramer-positive T-cells determined by flow cytometry.
54 These data provide a framework for unbiased definition of TCRs targeting lipid antigens, which
55 can be tested for clinical associations independently of host genetic background.

56 INTRODUCTION

57 Classically, T cells recognize peptide antigens when bound to major histocompatibility
58 complex (MHC) molecules (1). The six MHC class I and class II genes are among the most
59 diverse in the human genome. MHC diversity provides a population advantage against
60 pathogens by reducing the probability that a mutation will evade detection by the adaptive
61 immune system. This leads to a selective preference for MHC heterozygosity (2). T cells
62 express an antigen-specific T cell receptor (TCR), which enables recognition of peptide antigen
63 only when it is bound to MHC. A consequence of MHC allelic diversity is that each host must
64 develop a personalized set of T cells able to recognize peptide-based antigens produced by any
65 pathogen. Even among identical twins, the stochastic nature of TCR formation results in largely
66 distinct naïve T cell repertoires, although increased sharing is observed among antigen-
67 experienced T cells (3). Thus, the collection of all T cells (repertoire) between two genetically
68 unrelated individuals rarely overlaps.

69 The TCR is a heterodimer consisting of an α and β chain that is generated by somatic
70 recombination of germline-encoded segments. Further junctional diversity is provided by
71 nucleotide additions and deletions, increasing the potential diversity to nearly one trillion unique
72 sequences (4, 5). However, the actual number of unique T cells (clonotypes) in peripheral
73 blood is estimated at 3-4 million, a substantially lower number because of the processes of
74 positive and negative selection that occur in the thymus and because of the finite number of T
75 cells present in a single individual at a specific point in time (6). Among genetically unrelated
76 individuals that share a dominant MHC allele, a minority of T cells recognizing a viral peptide
77 antigen share a common TCR- β sequence (7). These examples demonstrate the challenge for
78 generalizing information derived from TCRs without first conditioning on genetic similarity.

79 Notably, human T cells have evolved mechanisms independent of MHC to facilitate the
80 recognition of non-peptide antigens. Such antigens include bacterial lipids, which bind the CD1

81 family of antigen-presenting molecules (8). There are five CD1 proteins in humans (CD1A,
82 CD1B, CD1C, CD1D, and CD1E) capable of processing and presenting lipid antigens to T cells.
83 Rapid evolution is common among immunity-associated genes, and MHC class I molecules are
84 frequently under strong directional selection (9). By contrast, CD1 genes do not display the
85 level of polymorphism inherent to MHC, though the actual levels of human genetic variation
86 have not been quantitatively determined. CD1 gene conservation may result in CD1-restricted
87 T cells with limited TCR diversity. For example, a number of studies have documented the
88 presence of invariant NK T (iNKT) cells in the blood of unrelated subjects from different
89 populations (10-12). These cells are activated by glycolipid antigens bound to CD1D, and
90 canonically express a TCR- α consisting of a germline rearrangement including TRAV10-1 and
91 TRAJ18-1 gene segments in humans (13, 14). Similarly, germline-encoded mycolyl lipid-
92 reactive (GEM) T cells are activated by mycobacterial glycolipids bound to CD1B and express a
93 TCR- α consisting of TRAV1-2 and TRAJ9-1 gene segments (15). Both iNKT and GEM cells
94 express TCR- β with biased gene segment usage, thus permitting the recognition of multiple
95 antigens.

96 Another class of T cells that are shared at the population level are mucosal-associated
97 invariant T (MAIT) cells, which are activated by bacterial metabolites of vitamin B that bind the
98 MR1 antigen-presenting molecule (16). Like iNKT and GEM cells, MAIT cells express a semi-
99 invariant TCR- α consisting of TRAV-1 and TRAJ33 gene segments in humans, and are
100 detectable in the blood of unrelated individuals (14, 17). MAIT cells are activated in an MR1-
101 dependent manner by M.tb-infected cells, though the M.tb-specific ligand(s) remain to be
102 defined (18, 19). By contrast, a number of mycobacterial lipid antigens presented by CD1
103 proteins to T cells have already been discovered (20). However, it is not known whether these
104 lipid antigens are recognized by antigen-specific T cells bearing TCRs that are shared across a
105 population or whether these TCRs are associated with disease.

106 Our data reveal a remarkable conservation of CD1B across diverse populations, which is
107 in stark contrast to the polymorphism that is inherent in MHC. We hypothesized that this
108 structural constraint would enable a lipid-antigen specific T-cell response to be shared across
109 genetically unrelated members of a population. By analyzing the TCRs of T cells specific for a
110 mycobacterial lipid antigen, we identified a conserved T cell response that was detectable in the
111 blood of multiple unrelated individuals. Using TCR immunosequencing and flow cytometry, we
112 show that this lipid-specific T cell population is expanded during tuberculosis infection and
113 disease. Our data reveal one example of an antigen-specific shared T cell repertoire that could
114 be leveraged to develop molecular markers for diseases in which lipid antigens are targeted by
115 T cells, such as tuberculosis, psoriasis, or leukemia.

116 RESULTS

117 Comparing genetic variation in Human HLA-A and CD1B

118 To compare the variation between classical and non-classical antigen presenting
119 molecules, we examined DNA sequence diversity in a representative MHC Class I gene (HLA-
120 A) and CD1 gene (CD1B) using available data from Phase 3 of the 1000 Genomes Project (21).
121 This resource includes whole genome sequences from 2,504 individuals covering all five major
122 continental ancestry groups, and thus it serves as a comprehensive resource for studying
123 human genetic diversity. For each of these two loci, we quantified nucleotide diversity (π) and
124 evidence of selection (Tajima's D) at each position along the gene (22, 23). Within HLA-A, the
125 median value of π was 0.015, with hotspots noted in exons 2 and 3, which are known to code for
126 the peptide antigen-binding domains. We also observed a median value of 0.92 for Tajima's D,
127 suggesting balancing selection of multiple alleles at high and low frequencies at the population
128 level (Figure 1A). We then examined the consequences of the observed nucleotide diversity by
129 mapping these onto the crystal structure of HLA-A (24). This analysis highlights the distribution
130 of missense, nonsense, and frameshift mutations that are particularly enriched among residues
131 lining the peptide-binding groove (Figure 1B). These data are consistent with published studies
132 of population diversity within MHC genes, and with our current understanding of how sequence
133 diversity leads to diversity in ligand binding at the molecular level (25, 26). By contrast, CD1B
134 was nearly invariant with a median π of 0.00014, and a median Tajima's D value of -1.14,
135 suggesting neutrality or weak purifying selection at the population level (Figure 1C). The tertiary
136 structure of the lipid-binding domain of CD1B was completely conserved with only two
137 synonymous polymorphisms outside of the antigen-binding groove (Figure 1D) (27). These
138 data support purifying selection on CD1B during human evolution and stand in stark contrast to
139 the immense diversity inherent within HLA-A and other MHC Class I genes.

140

141 *Derivation of GMM-specific T cell lines from healthy South African donors*

142 The lack of allelic diversity within CD1B suggests that genetically unrelated individuals
143 might generate a shared T-cell response to a given lipid antigen. To test this hypothesis
144 directly, we used CD1B tetramers loaded with the immunodominant mycobacterial lipid glucose
145 monomycolate (GMM) to isolate antigen-specific T cells from four healthy South African adults.
146 Tetramer-positive cells were expanded *in-vitro* for four weeks to derive GMM-specific T cell lines
147 (Figure 2A). We confirmed that GMM-specific T cells were primarily CD4+, that they were
148 activated and produced IFN- γ in the presence of GMM and CD1b, and that they did not cross-
149 react even with closely related antigens, such as mycolic acid (Figure 2B and data not shown).
150 Finally, we confirmed that GMM-specific T cell lines were functional, expressing IFN- γ , TNF- α ,
151 IL-2, IL-17, and CD40L when stimulated through cognate interactions with antigen or
152 independently of the T-cell receptor (Figure 2C and data not shown). These data confirm
153 previously published data by our group and others showing that GMM-specific T cells with
154 similar functional profiles are present in the blood of unrelated donors (28, 29).

155

156 *The diversity of T cells that bind CD1B-GMM tetramer*

157 We then re-sorted those T cells that bound GMM-loaded CD1b tetramer with high avidity
158 and exhaustively profiled the TCRs in both the tetramer-positive and tetramer-negative cell
159 populations using high-throughput immunosequencing (Figure 2A)(6). The surprising diversity
160 of V and J gene segments present in both TCR- α and TCR- β sequences from tetramer-positive
161 cells suggest that a number of clonotypes (i.e., unique T-cells expressing a specific TCR- α and
162 TCR- β) are able to bind CD1B tetramer loaded with GMM (Figure 3A, 3B and Supplementary
163 Figure 1). The clonotype diversity observed after *in vitro* expansion was qualitatively similar to
164 that observed in *ex vivo* samples (Figure 3A and Supplementary Figure 1). For example, there
165 was marked enrichment of TRAJ26 in the *ex vivo* tetramer-positive and expanded and resorted

166 populations compared to tetramer-negative cells (Figure 3A). A similar pattern was observed
167 within the TCR- β chain sequences, in which tetramer-positive cells showed a clear bias towards
168 the use of certain variable and joining segments as compared to those observed in the tetramer-
169 negative cells (Figure 3B and Supplementary Figure 1). We then compiled lists of the most
170 highly enriched V and J gene segments from TCR- α (TRAV1, TRAJ9) and TCR- β (TRBV6,
171 TRBJ2) by comparing expanded and resorted tetramer-positive cells to tetramer-negative cells
172 and analyzing the length of the CDR3 region. We found that CDR3 length was severely
173 restricted in these cell populations (Figure 3C and Supplementary Figure 2). Collectively, these
174 data provide further evidence of substantial restriction in the TCR- α and TCR- β chains that bind
175 CD1B-GMM tetramers.

176

177 Unbiased definition of a shared GMM-specific TCR motif

178 To determine whether there is a GMM-specific T-cell repertoire that is shared among
179 unrelated donors, we standardized the analysis described above by constructing a simple test to
180 identify TCR motifs. We defined motifs as sequences with common V and J family usage and
181 CDR3 length that were significantly enriched in the expanded and re-sorted samples, as
182 compared to their corresponding tetramer-negative samples. This unbiased approach resulted
183 in one TCR- α motif and one TCR- β motif that were significantly enriched independently in each
184 of the four donors. The TCR- α motif contained a 13-amino acid CDR3 sequence that is
185 consistent with the recently published TRAV01/TRAJ09 rearrangement present in GEM T cells
186 (Figure 4A and Supplementary Table 1) (15). The TCR- β motif we discovered in this manner
187 included a 14-amino acid CDR3 sequence that results from a TRBV06/TRBJ02 rearrangement
188 (Figure 4B and Supplementary Table 1). Furthermore, several lines of evidence support co-
189 expression of these motifs as a heterodimeric TCR in unrelated individuals. First, the dominant
190 TCR- α and TCR- β sequences for one of our T cell lines analyzed in this study matches these

191 motifs, suggesting they are co-expressed by a dominant T-cell clone (data not shown).
192 Second, when we sorted GMM-CD1B tetramer avid cells from a fifth unrelated South African
193 donor and performed limiting dilution to establish T-cell clones, we identified two T-cell clones
194 whose TCRs match the TCR- α and TCR- β motifs described above (Supplementary Table 2).
195 Finally, Van Rhijn et al. reported a GMM-specific T-cell clone derived from an M.tb-infected
196 subject whose TCR also matches these TCR- α and TCR- β motifs (15). Taken together, these
197 data suggest that specific TCR- α and TCR- β chains that conform to the motif described above
198 are co-expressed as a heterodimer, are specific for the mycobacterial antigen GMM, and are
199 shared among multiple unrelated donors.

200

201 *Predicted binding interactions and ontogeny of the GMM-specific TCR motifs*

202 We next addressed the question of how, given the stochastic nature of T-cell
203 development, these motifs might arise independently in unrelated donors. To accomplish this,
204 we first assessed the enrichment of particular amino acid residues at each position along the
205 CDR3 of expanded and re-sorted samples. In this analysis, we used tetramer-negative samples
206 as a natural control and looked for differences in the CDR3 sequence of the motifs described
207 above. In the TCR- α motif, we observed that the arginine at position 4 and leucine at position 5
208 were enriched, suggesting that these residues are particularly important for mediating binding to
209 the GMM-CD1B tetramer (Figure 4C and Supplementary Tables 3 and 4). This inference is
210 supported by TCR reconstitution experiments, in which a TCR containing a TCR- α chain that
211 lacks arginine at position failed to respond functionally to GMM (15). We also identified several
212 enriched residues between positions 4 and 13 of the TCR- β motif (Figure 4D and
213 Supplementary Table 4). Figure 4E shows the location of these enriched residue positions
214 identified through sequence analysis mapped to the crystal structure of a GMM-specific TCR
215 whose TCR- α and TCR- β chains match our motifs (15). Interestingly, both chains of the

216 heterodimer appear to contribute equally to binding the CD1B-GMM tetramer, which is
217 supported by the recently published co-crystal of a GEM TCR with CD1B-GMM (30). These
218 data highlight how immunosequencing can be used to infer structural requirements for antigen
219 binding by TCRs.

220 We next explored the distribution of nucleotide additions and deletions leading to
221 junctional diversity in the TCR- α motif described above. We observed a lower level of 3'
222 deletion in the V segment of TCR- α sequences from expanded and resorted T-cells than in the
223 those from tetramer-negative cells (Figure 4F). We note that the germline definition for TRAV01
224 ends with CAVR (Figure 4F), which includes the conserved arginine at position 4 required for
225 binding CD1B-GMM, suggesting that this CDR3 is germline-encoded, as previously proposed
226 (15). In contrast, the germline sequence of TRAJ09 begins with GNTGGFKTIF, implying that
227 position 5 of the selected motif results from selection of non-templated insertions that result in
228 the encoding of a leucine residue in CD1B-GMM specific TCRs. Alternatively, it is possible that
229 modulation of the residue distribution at position 5 could be a consequence of the restriction
230 observed at position 4, since non-templated insertion has been shown to be context dependent
231 (31). When we performed the same analysis for TCR- β , we observed an increased level of 3' V
232 deletion in CD1B-GMM specific TCRs as compared to those observed in TCRs from tetramer-
233 negative cells, suggesting that enriched residues at positions 4, 5, and 6 are likely non-
234 templated. On the other hand, the enriched residues at positions 10 to 13 are located entirely
235 within the germline-encoded TRBJ02 gene segment (Figure 4G). These data show that a
236 GMM-specific shared TCR could arise through a combination of germline rearrangements and
237 high-probability events occurring at junctions.

238

239 *Ex-vivo analysis of GMM-specific TCR motifs in healthy donors*

240 To extend the results described above, which were obtained from *in-vitro* expanded and
241 resorted T-cell lines, we next examined sequences obtained from tetramer-positive cells sorted
242 directly *ex vivo* (Figure 2A). We pooled sequences from all four subjects studied initially and
243 considered sequences concordant with the previously defined TCR- α and TCR- β motifs in terms
244 of V and J gene family and CDR3 length as matches (Figure 4A and 4B). We observed a
245 significant enrichment of both the TCR- α and TCR- β motifs among tetramer-positive cells as
246 compared to tetramer-negative cells (Table 1). These data support the utility of the TCR motifs
247 identified in this study as markers of antigen-specific T-cell responses *ex vivo*. We next
248 examined the prevalence of the TCR- β motif in a previously published dataset that included
249 TCR- β sequences obtained from the PBMC of 587 bone marrow donors at low risk for
250 *Mycobacterium tuberculosis* infection (32). We found that the motif was present at a frequency
251 very similar to that observed in tetramer-negative T cells in the pooled data (Table 1). The
252 relative absence of the GMM-specific TCR- β motif in T cells from healthy donors suggests that it
253 could be used as a specific marker of mycobacterial infection in population-based studies.

254

255 GMM-specific TCR motifs during tuberculosis

256 Finally, we investigated the association between GMM-specific TCR motifs and
257 mycobacterial infection by analyzing a cohort of South Africans with known M.tb infection status.
258 This cohort included three clinical groups: adolescents without latent tuberculosis infection
259 (IGRA-negative), adolescents with latent tuberculosis (IGRA-positive), and adults with a new
260 diagnosis of active tuberculosis disease (Active TB, n=10). We comprehensively profiled the
261 TCR- α and TCR- β chains from an average of approximately 100,000 T cells per donor by
262 immunosequencing cryopreserved PBMCs, and we pooled the data into each of the three
263 clinical groups. Next, we calculated and compared the 'motif burden' of the TCR- α and TCR- β
264 motifs described above, defined as the fraction of unique sequences that matched the V and J

265 gene family, and CDR3 length of these motifs (Table 2). For the TCR- α motif, we noted an
266 increased motif burden in IGRA-positive subjects ($p=0.0014$) and active TB patients ($p=8.3 \times 10^{-7}$)
267 compared to IGRA-negative control subjects. For the TCR- β motif, we also observed an
268 increased motif burden in IGRA-positive subjects ($p=4.9 \times 10^{-5}$), and active TB patients ($p=8.3 \times$
269 10^{-7}) when compared to IGRA-negative subjects (Table 2). When we calculated subject-specific
270 motif burdens, we confirmed enrichment in the TCR- α when comparing active TB patients to
271 IGRA-negative subjects ($p=0.004$), although in this case we did not see an enrichment in the
272 TCR- β motif burden (Figure 5A and Supplementary Figure 3). Notably, we did not see
273 enrichment of canonical MAIT cell or iNKT cell TCR- α rearrangements in active TB, highlighting
274 the specificity of the association between a GMM-specific TCR- α motif and tuberculosis
275 (Supplementary Figure 3).

276 We further validated these results using flow cytometry to quantify GMM-specific T cells
277 ex vivo (Supplementary Figure 4). We found that the frequency of CD3⁺ T cells that stained
278 with CD1B-GMM tetramer was higher in patients with active TB as compared to IGRA-negative
279 control subjects ($p=0.001$, Figure 5B). Finally, in one patient with active tuberculosis, we
280 analyzed blood samples at diagnosis, during treatment, and at the completion of therapy, and
281 we observed a decrease in the frequency of sequences that matched both the TCR- α and TCR-
282 β motifs over time (Supplementary Figure 3). Together, these data support clonal expansion
283 and contraction of GMM-specific T cells during mycobacterial infection and clearance,
284 respectively, and demonstrate the association between GMM-specific TCR motifs and M.tb
285 infection and disease.

286 **DISCUSSION**

287 The phenomenon of MHC restriction has been a cardinal feature of T-cell immunology
288 for more than forty years. In practical terms, this means that the ability to activate T cells
289 depends not only on the foreign antigens they recognize, but also on the MHC molecules that
290 bind the antigen. By contrast, we show here that CD1B is nearly invariant among humans, and
291 we report T-cell receptor motifs that are specific for a mycobacterial lipid antigen presented by
292 CD1B, and consistently shared among multiple unrelated donors. We also show that T cells
293 bearing these motifs are more frequent in South Africans with latent M.tb infection or
294 tuberculosis disease than in M.tb-uninfected subjects. These data support the existence of a
295 lipid antigen-specific shared T-cell repertoire that is independent of genetic background and
296 could be leveraged to develop molecular diagnostics based on TCR sequence motifs.

297 There are numerous examples of ‘public’ TCR sequences that are shared across
298 unrelated donors. Most of these have been reported as reactive to viral peptide antigens
299 (reviewed in (33) and (34)). However, these reports also make it clear that most T-cell
300 responses are ‘private’ to a single individual, and that public T-cell responses are the exception
301 rather than the rule (7). Importantly, since most of these TCRs likely bind peptide-based
302 epitopes, they are by definition MHC-restricted. On the other hand, some T cells that share an
303 invariant TCR- α chain, such as iNKT and MAIT cells, are widely prevalent but not clearly
304 antigen-specific or related to pathogen exposure. In this study, we show that for CD1B
305 presentation of GMM, a shared antigen-specific TCR response is seen reliably in more than 30
306 unrelated individuals. Our data support a model that requires a conserved, but not necessarily
307 invariant, rearrangement in both α and β chains of the TCR heterodimer to recognize a foreign
308 lipid antigen. Additionally, by comparing the TCR sequences of GMM tetramer-positive and
309 tetramer-negative T cells, we provide a structural rationale for this notion that is confirmed by
310 the recent publication of GEM T cells and associated crystal structures (15, 30). Our data

311 support the emerging paradigm that TCR recognition of foreign lipids follows a 'parallel docking
312 mode' that is more reminiscent of peptide-MHC than of CD1d- α GalCer TCR interactions, which
313 is relatively biased toward the TCR- α chain (35).

314 Our unbiased approach to motif discovery is thus validated in the setting of GMM,
315 suggesting it can be applied more broadly to lipid antigens outside of infectious diseases. For
316 example, T cells have been described that recognize methyl-lyso phosphatidic acid when bound
317 to CD1C expressed by acute myeloid and B cell leukemias (36). T cells recognizing
318 lysophospholipids bound to CD1A contribute to pathogenic inflammation in atopic dermatitis and
319 psoriasis (37, 38). Human CD1A and CD1C tetramers as well as synthetic lipid antigens are
320 available as validated reagents, so that an approach similar to the one described here could
321 potentially be used to define allergic or cancer-specific shared TCRs and to test their clinical
322 associations.

323 We detected GMM-specific TCR motifs in subjects without *M.tb* infection, suggesting
324 either a high frequency of these motifs in the naïve T cell repertoire as a consequence of near-
325 germline rearrangements, T-cell priming through mechanisms other than infection with *M.*
326 *tuberculosis*, or both. Because GMM is produced by mycobacteria other than *M.tb*, this
327 phenotype could result from exposure to non-tuberculous mycobacteria in the environment, or
328 from vaccination with the *M. bovis* derived-BCG vaccine at birth, which is standard practice in
329 South Africa. Consistent with this possibility, we found that the TCR- β motif was not enriched in
330 a cohort of U.S. bone marrow donors who are at low risk for *M.tb* exposure and did not receive
331 BCG. We can apply our unbiased approach to identify shared TCRs for lipids that are
332 preferentially expressed by virulent *M.tb*. Unlike GMM, sulfoglycolipids are not expressed by
333 BCG or most environmental bacteria, and T-cell responses to sulfoglycolipids have been shown
334 to be greater in *M.tb*-infected subjects when compared to BCG-vaccinated subjects (39). The

335 discovery of shared TCRs covering a variety of mycobacterial lipids could then be tested
336 independently or in combination for utility as a molecular marker for tuberculosis.

337 We used both immunosequencing and flow cytometry to demonstrate that GMM-specific
338 T cells in peripheral blood are expanded during active tuberculosis. The increased abundance
339 of the GMM-specific TCR motifs we identified among active TB patients suggests expansion of
340 high-avidity clones in the presence of antigen. Moreover, we show a decrease in GMM-specific
341 TCRs during treatment for active tuberculosis. These dynamic changes are consistent with the
342 development of immunologic memory to lipid antigens. By contrast, we show that the frequency
343 of TCRs expressed by innate-like T cells, such as iNKT and MAIT cells, are not associated with
344 active tuberculosis. In this respect, GMM-specific T cells may more closely resemble
345 conventional MHC-restricted T cells rather than innate-like T cells. Future studies could
346 address whether these TCRs specific for mycobacterial lipids could be used as molecular
347 markers of response to drug treatment or immunogenicity of whole cell vaccines, such as BCG.

348 The advantage of maintaining a highly polymorphic pool of MHC genes in an outbred
349 population have been extensively explored (40). In contrast, the advantages of maintaining a
350 nearly invariant CD1B gene are less obvious. While MHC finds its origins in jawed vertebrates,
351 CD1 genes have only been reported in mammals and birds, suggesting CD1 evolved from MHC
352 genes to perform a non-redundant function. A number of pathogens have evolved mechanisms
353 to avoid detection by the adaptive immune system, by varying the sequence of peptides
354 available to bind host MHC molecules (41). Thus, CD1 may have evolved to facilitate T-cell
355 detection of non-peptide antigens that are less likely to undergo mutational escape because of
356 their importance in maintaining cell wall integrity. Supporting this notion is the success of
357 mycolic acid synthesis inhibitors, such as isoniazid, for the treatment of tuberculosis. Isoniazid
358 is predicted to also inhibit the production of GMM (42). Whether the findings we report here can
359 be generalized to other lipids or non-peptide antigens remains to be determined, as do the true
360 evolutionary origins of non-MHC systems of antigen presentation.

361 In summary, we provide evidence of purifying selection in human CD1B and identify a
362 set of conserved TCRs specific for CD1B and a mycobacterial lipid that is enriched during
363 tuberculosis infection and disease. These data demonstrate a general framework for the rapid
364 and unbiased identification of T cells specific for CD1-presented antigens as well as their
365 associations with disease. These T cells are part of a potentially much larger shared T-cell
366 repertoire that might have a number of potential clinical applications. One possibility is
367 molecular diagnosis of M.tb or any other disease characterized by lipid antigen-specific T-cells,
368 such as atopic dermatitis, psoriasis, or leukemia. The other possibility is adoptive cell therapies
369 that do not depend upon the genetic background of the patient. This is in contrast to current
370 clinical trials of TCR-modified T cells, for example, that are typically limited to the dominant
371 MHC haplotypes present in a population at risk for cancer (43). The realization of such novel
372 cell therapies will require the identification of disease-specific lipid antigens and shared TCRs,
373 but much of the technology for translating this information into a clinical therapy is already
374 available.

375 **MATERIALS AND METHODS**

376 Generation of GMM-loaded CD1B tetramers.

377 Soluble biotinylated CD1B monomers were provided by the National Institutes of Health
378 Tetramer Core Facility (Emory University, Atlanta, GA). Glucose monomycolate (GMM)-loaded
379 tetramers were generated as previously described (29). In brief, C32-GMM was dried down in a
380 glass tube using a nitrogen evaporator and sonicated into 0.25% CHAPS/sodium citrate at pH 4
381 (preparation of CHAPS in sodium citrate; CHAPS Hydrate, Sigma; Sodium Citrate Dihdrate,
382 Fisher) for two minutes at 37°C. The lipid solution was transferred to a microfuge tube, and 9
383 μ L of CD1B monomer was added. The CD1B-GMM preparation was then incubated in a 37°C
384 water bath for 2 hours with vortexing every 30 minutes. At the end of the incubation, the
385 solution was neutralized to pH 7.4 with 6 μ L of 1M Tris pH 9. Finally, 10 μ L of Streptavidin
386 conjugated to allophycocyanin or phycoerythrin (Life Technologies) was added in ten aliquots
387 of 1 μ L every 10 minutes to facilitate tetramerization. The final product was filtered through a
388 SpinX column (Sigma) to remove aggregates and stored at 4°C until use.

389

390 Antigens and Media

391 C32-GMM purified from *Rhodococcus equi* was generously provided by the laboratory of D.
392 Branch Moody (44). Our base T cell media consists of RPMI 1640 (Gibco) supplemented with
393 100 U/mL Penicillin and 100 μ g/mL Streptomycin (Gibco), 55 μ M 2-mercaptoethanol (Gibco),
394 0.3X Essential Amino Acids (Gibco), 60 μ M Non-essential Amino Acids (Gibco), 11mM HEPES
395 (Gibco), and 800 μ M L-Glutamine (Gibco). We additionally supplemented this media with either
396 fetal calf serum (HyClone) or human serum that was derived from healthy donors.

397

398 Generation of T cell lines using GMM-loaded CD1B tetramers

399 Peripheral blood mononuclear cells (PBMC) were isolated from South African adults with latent
400 tuberculosis infection as defined by a positive QuantiFERON-TB Gold blood test. PBMC were
401 depleted of CD14-expressing monocytes for a separate study and cryopreserved. At the time of
402 this study, PBMC were thawed and washed in warm RPMI 1640 (Gibco) supplemented with
403 10% fetal calf serum (Hyclone) and 10 $\mu\text{L}/\text{mL}$ Benzoylase (Millipore) and enumerated using
404 Trypan blue exclusion. PBMC were then plated in a 24-well plate at a density of three million
405 cells per well in T cell media and allowed to rest overnight at 37°C in humidified incubators
406 supplemented with 5% CO₂. The following day, PBMC were washed and blocked with 50%
407 human serum (Valley Biomedical) in PBS supplemented with 0.2% BSA (Sigma) (FACS buffer)
408 for 10 min at 4°C. The samples were washed twice with PBS and stained with Aqua Live/Dead
409 stain (Life Technologies) according to the manufacturer's instructions. Following two additional
410 PBS washes, cells were resuspended in 50 μL FACS buffer and 1 μL of either unloaded CD1B
411 tetramer or GMM-loaded CD1B tetramer and incubated at room temperature for 40 minutes in
412 the dark. Finally, cells were stained with anti-CD3 ECD (Beckman Coulter) as well as surface
413 markers for activation and memory phenotype, washed twice in T cell media, and screened
414 through a cell strainer tube (Falcon) prior to sorting. Tetramer-positive T cells were sorted at the
415 UW Department of Immunology Flow Cytometry Core using a FACS Aria II (BD) cell sorter
416 equipped with 407nm, 488nm, and 641nm lasers.

417 Sorted T cells were washed and resuspended in T cell media supplemented with 10%
418 human serum. T cells were then divided among eight wells of a 96-well U-bottom tissue culture
419 plate into which irradiated PBMC (150,000 cells per well) were added as feeder cells along with
420 phytohaemagglutinin (Remel) at a final concentration of 1.6 $\mu\text{g}/\text{ml}$. After two days in culture at
421 37°C, 5% CO₂, 10 μL natural IL-2 (Hemagen) was added to each well. Half the media was
422 replaced every two days with T cell media supplemented with 10% human serum and natural IL-
423 2. When the cell clusters were large and round (approximately after eight days of growth), they

424 were pooled into a 24-well plate. After 10 days in culture, cell lines were screened by tetramer
425 staining or functional response to GMM. We then further expanded T cell lines by modifying a
426 previously published rapid expansion protocol (45). Briefly, 200,000 T cells were mixed with 5
427 million irradiated EBV-transformed B cells and 25 million irradiated PBMC as feeder cells in T
428 cell media. Anti-CD3 (clone OKT3) was added a final concentration of 30 ng/mL, and the
429 mixture was incubated overnight at 37°C, 5% CO₂. The following day, recombinant IL-2 (rIL-2)
430 (UWMC Clinical Pharmacy) was added to a final concentration of 50 U/mL. On day 4, the cells
431 were washed twice in T cell media to remove OKT3, and fresh media supplemented with rIL-2
432 at 50 U/mL was added. Half the media was replaced every three days or split into new T25
433 tissue culture flasks (Costar) as determined by cell confluency. After 13 days in culture, the
434 lines were screened by tetramer staining and then cryopreserved on day 14.

435

436 Clinical Cohorts

437 As recently published, 6363 adolescents were enrolled into a study that aimed to
438 determine the incidence and prevalence of tuberculosis infection and disease (46). Twelve to
439 18 year-old adolescents were enrolled at eleven high schools in the Worcester region of the
440 Western Cape of South Africa. Subjects were screened for the presence of latent tuberculosis
441 by a tuberculin skin test and IFN- γ release assay (IGRA) QuantiFERON-TB GOLD In-Tube
442 (Cellestis Inc.) at study entry. Peripheral blood mononuclear cells (PBMC) were isolated from
443 freshly collected heparinized blood via density centrifugation and cryopreserved. For this work,
444 a subset of samples from 10 M.tb-infected and 10 M.tb-uninfected adolescents were selected
445 based on matching for age and sex and availability of cryopreserved specimens after
446 completion of the primary objectives of the parent study. We also accessed a recently
447 published cohort of South African adults with a new diagnosis of active tuberculosis (47). Only
448 participants were included when \geq 18 years of age and seronegative for HIV were included. All
449 included patients had either positive sputum smear microscopy and/or positive culture for *M.tb*.

450 Blood was obtained and PBMC archived prior to or within 7 days of starting standard course
451 anti-TB treatment, which was provided according to South African national health guidelines.
452 For this work, a convenience sample from 10 adults with active tuberculosis was selected based
453 on availability of cryopreserved specimens after completion of the primary study. For T-cell
454 cloning studies, we used cryopreserved PBMC from a cohort of healthy South African adults
455 (Supplemental Table 5).

456 Human peripheral blood samples were also obtained from a cohort of healthy bone
457 marrow donors from the Fred Hutchinson Cancer Research Center Research Cell Bank
458 biorepository, under a protocol approved and supervised by the Fred Hutchinson Cancer
459 Research Center Institutional Review Board, following written informed consent. This cohort has
460 been previously described (32) (Emerson et al., manuscript submitted).

461 *Flow Cytometry*

462 PBMC from South African adolescents and adults were thawed and rested overnight as
463 described above. The following day, one million PBMC per subject were counted and plated
464 into each of two wells of a 96-well U-bottom plate. The cells were spun down and blocked in
465 human serum and FACS buffer at 1:1 for ten minutes at 4°C. Following a wash with FACS
466 buffer, the cells were resuspended in 50 μ L FACS buffer, and 1 μ L of either unloaded CD1B
467 tetramer conjugated to PE or 1 μ L GMM-loaded CD1B tetramer conjugated to APC. The
468 samples were incubated at room temperature for 40 minutes and then washed twice with PBS.
469 Aqua Live/Dead stain was then prepared at a 1:100 dilution in PBS and added to each sample.
470 All samples were incubated for 15 min at room temperature in the dark. The samples were then
471 washed twice and resuspended in 50 μ L of staining cocktail containing anti-CD3 ECD in FACS
472 buffer for 30 minutes at 4°C in the dark. Cells were then washed twice in FACS buffer and
473 centrifuged at 1800 rpm for 3 min. All samples were then resuspended in 200 μ L of 1%

474 paraformaldehyde prior to acquisition on LSRII (BD Biosciences) equipped with 488nm
475 (100mW), 532nm (150mW), 628nm (200mW), 405nm (100mW), and 355nm (20mW) lasers.

476

477 Functional Assays

478 We used IFN- γ ELISPOT to examine antigen-specificity of T cell lines. Briefly, K562
479 cells stably transfected with empty vector (K562-EV) or CD1B (K562-CD1B)(48) were
480 maintained in RPMI 1640 (GIBCO) supplemented with 10% fetal calf serum and G418 (Sigma)
481 at a concentration of 200 μ g/ml and periodically assessed for CD1B expression by flow
482 cytometry. Multiscreen-IP filter plates (Millipore) were coated with 1-D1K antibody (Mabtech)
483 overnight at 4°C. Lipids were evaporated to dryness from chloroform-based solvents under a
484 sterile nitrogen stream and then sonicated into media. This lipid suspension was added to co-
485 cultures of 50,000 mock transfected or CD1B-transfected K562 antigen-presenting cells and
486 2000 T cells with a final concentration of 1 μ g/ml for MA and GMM (8, 49). The co-cultures
487 were incubated for 16 hours at 37°C. The following day, the cells were lysed with water and
488 then incubated with 7-B6-1 biotin conjugate (Mabtech) for two hours at room temperature. The
489 plate was washed with PBS (Life Technologies) and then incubated with ExtraAvidin-Alkaline
490 Phosphatase (Sigma) for 1 hour at room temperature. Lastly, the wells were washed and then
491 developed using BCIP/NBT substrate (Sigma) for 5 min at room temperature in the dark. The
492 wells were imaged and the IFN- γ spots were counted using an ImmunoSpot S6 Core Analyzer
493 (Cellular Technology Limited).

494 To observe intracellular cytokine staining following GMM antigen stimulation, 3.3 million
495 K562 cells/ml were incubated with 5 μ g/mL GMM at a final volume of 100 μ l for 18 hours at
496 37°C, 5% CO₂ to facilitate lipid loading. The following day, T cell lines were plated at 1 million
497 cells/mL and 80 μ L of pre-loaded K562 cells was added to the T cells. The cell mixture was
498 allowed to incubate for six hours in the presence of anti-CD28/49d antibodies (BD Biosciences),

499 Brefeldin A at a final concentration of 10 μ g/ml (Sigma), and GolgiStop containing Monensin
500 (BD), and anti-CD107a (BD) after which EDTA, at a final concentration of 2mM, was added to
501 disaggregate cells. Plates were stored at 4°C until the following day when they were stained
502 and acquired by FACS. We used a previously published optimized and validated 12-color panel
503 (50, 51). Briefly, cells were first stained with Avid Live/Dead (Life Technologies) viability dye
504 and anti-CD14. After washing, the cells were permeabilized with FACS Perm II (BD) and
505 stained for the remaining markers (CD3, CD4, CD8, CD154, IFN- γ , TNF- α , IL-2, IL-4, IL-17a,
506 and CD40L). Fully stained cells were washed and resuspended in 1% paraformaldehyde. Data
507 were acquired on an LSR II (BD) equipped with a high-throughput sampler and configured with
508 405nm (50mW), 488nm (20mW), 561nm (50mW), and 641nm (150mW) lasers.

509

510 Ethics

511 The study was approved by the IRB of the University of Washington and the University
512 of Cape Town. Written informed consent was obtained from all adult participants as well as
513 from the parents and/or legal guardians of the adolescents who participated. In addition, written
514 informed assent was obtained from the adolescents.

515

516 Comparative Genomics

517 Data for CD1B (Chr1:158,327,951-158,331,531) and HLA-A (Chr6:29,941,260-
518 29,945,884) were extracted from the 1000 Genomes Project (Phase 3; aligned to GRCh38), and
519 measures of divergence calculated using VCFtools. Nucleotide diversity (π)
520 (<http://www.pnas.org/content/76/10/5269.abstract>) was estimated with sliding windows of 50
521 base pair with two base pair steps, and Tajima's D
522 (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1203831/>) was estimated over static 50 base
523 pair windows. Coding SNP variation with >1% minor allele frequency was extracted, and

524 manually mapped onto the corresponding protein structures (CD1B: PDB #1UQS, HLA-A: PDB
525 #3MRG) according to variant type (synonymous or nonsynonymous) (24, 27).

526

527 Immunosequencing

528 High-throughput sequencing of TCRs was performed using the ImmunoSEQ assay
529 (Adaptive Biotechnologies) with TCR- β and/or TCR- α/δ assays for each sample using a
530 multiplex PCR approach following by Illumina high-throughput sequencing (52).

531

532 Putative CD1B-GMM specific TCR motifs defined by V/J gene usage and CDR3 length

533 Sequence motifs for CD1B-GMM specific TCR- α and TCR- β were defined for each
534 subject by assessing enrichment of specific V and J gene family usage and CDR3 lengths within
535 the expanded and resorted sample as compared to the respective tetramer negative sample.
536 For each possible combination of V family, J family, and CDR3 length, we enumerated unique
537 sequences matching this combination, as well as the mismatching sequences for each sample.
538 From these four counts (expanded/resorted, tetramer-negative, matching, and mismatching), we
539 generated a 2x2 contingency table and used Fisher's Exact test to assign a p-value for the
540 significance of the enrichment of matches in the expanded and resorted sample. Such a test
541 was performed for all combinations for each subject and for both TCR- α and TCR- β data, and
542 we selected those with p-values below a pre-specified significance threshold ($p < 0.001$). This
543 resulted in a set of CD1B-GMM specific TCR motifs for each subject. To define the 'public'
544 CD1B-GMM specific TCR motifs within each locus, we identified motifs that were significantly
545 enriched as described above among all subjects.

546

547 Sequence conservation among CD1B-GMM specific TCRs matching putative motifs

548 We determined the sequence conservation among putative motif matches of expanded
549 and resorted TCRs (i.e. considered truly CD1B-GMM specific),s using putative motif matches
550 among sequences derived from tetramer-negative cells as a natural control. For each possible
551 amino acid at a given position in the CDR3, sequences matching the putative motif that
552 contained this amino acid at the given position were enumerated for expanded and resorted
553 data and tetramer-negative data, as were the sequences not matching that specific amino acid
554 at the given position. From these four counts, we generated a 2x2 contingency table and used
555 Fisher's Exact test to assign a p-value for significance of enrichment of this amino acid at the
556 given position among the expanded and resorted putative motif matches. Such a test was
557 performed for amino acids and CDR3 positions and for both TCR- α and TCR- β data, and we
558 selected those with p-values below a pre-specified significance threshold (p-value < 0.001).
559 Using the annotations of V, N, and J region boundaries (for TCR- α sequences) and V, N1, D,
560 N2, and J region boundaries (for TCR- β sequences), we were also able to assess differences in
561 recombination structure among CD1B-GMM specific TCRs with respect to nonspecific TCRs
562 that also matched the putative motifs (53). By integrating these observations with the positions
563 and identities of conserved residues, we were able to assess the conservation of germline-
564 encoded sequence versus specific non-templated insertions.

565

566 *CD1B-GMM specific TCR motif burden*

567 For population-level comparisons between T-cell populations or between patient groups,
568 we defined the 'motif burden' of a sample simply as the fraction of unique TCRs from that
569 sample that matched the CD1B-GMM specific TCR motif defined for the corresponding TCR- α
570 or TCR- β . We first assessed significance by pooling all the sequences within each clinical group
571 and counting the number that matched and the number that mismatched each motif using
572 Fisher's exact test on the resulting contingency table. Pooling sequences from multiple donors

573 is sensitive to one donor with an excess of sequences biasing the result. To address this, we
574 examined the sampling depth across donors in in each clinical group, and found that healthy
575 subjects were not systematically undersampled (data not shown). We also assessed
576 differences between clinical groups using subject-specific fractional motif burdens using the
577 non-parametric Kruskal-Wallis test with Dunn post test (a generalization of the Mann-Whitney
578 U test to more than two groups).

579 **AUTHOR CONTRIBUTIONS**

580 C.S. and R.O.E. designed the study. W.S.D. and K.K.Q performed the experiments and
581 analyzed the data. D.W. and W.S. contributed human genetic analyses. C.L.D. and T.J.S.
582 established the clinical cohorts. S.C.D facilitated flow cytometry studies. A.S. and H.S.R.
583 facilitated immunosequencing. W.S.D, M.V., R.O.E., and C.S. wrote the manuscript with
584 contributions from all authors.

585 **ACKNOWLEDGEMENTS**

- 586 • This work was supported by the National Institutes of Health (K08-AI089938 to CS)
587 University of Washington Department of Medicine and UW Royalty Research Fund (CS).
- 588 • The authors would like to thank Branch Moody and Martine Gilleron for providing lipid
589 antigens and John Hansen and the Research Cell Bank at the Fred Hutchinson Cancer
590 Research Center for assembling the cohort of bone marrow donors. We acknowledge
591 the NIH Tetramer Core Facility (contract HHSN272201300006C) for provision of
592 biotinylated CD1B monomers.

593 **REFERENCES**

- 594 1. Garcia KC, Degano M, Stanfield RL, Brunmark A, Jackson MR, Peterson PA, Teyton L,
595 and Wilson IA. An alphabeta T cell receptor structure at 2.5 Å and its orientation in the
596 TCR-MHC complex. *Science*. 1996;274(5285):209-19.
- 597 2. Hughes AL, and Nei M. Pattern of nucleotide substitution at major histocompatibility
598 complex class I loci reveals overdominant selection. *Nature*. 1988;335(6186):167-70.
- 599 3. Zvyagin IV, Pogorelyy MV, Ivanova ME, Komech EA, Shugay M, Bolotin DA, Shelenkov
600 AA, Kurnosov AA, Staroverov DB, Chudakov DM, et al. Distinctive properties of identical
601 twins' TCR repertoires revealed by high-throughput sequencing. *Proc Natl Acad Sci U S*
602 *A*. 2014;111(16):5980-5.
- 603 4. Davis MM, and Bjorkman PJ. T-cell antigen receptor genes and T-cell recognition.
604 *Nature*. 1988;334(6181):395-402.
- 605 5. Robins HS, Srivastava SK, Campregher PV, Turtle CJ, Andriesen J, Riddell SR, Carlson
606 CS, and Warren EH. Overlap and effective size of the human CD8+ T cell receptor
607 repertoire. *Science translational medicine*. 2010;2(47):47ra64.
- 608 6. Robins HS, Campregher PV, Srivastava SK, Wacher A, Turtle CJ, Kahsai O, Riddell SR,
609 Warren EH, and Carlson CS. Comprehensive assessment of T-cell receptor beta-chain
610 diversity in alphabeta T cells. *Blood*. 2009;114(19):4099-107.
- 611 7. Venturi V, Chin HY, Asher TE, Ladell K, Scheinberg P, Bornstein E, van Bockel D,
612 Kelleher AD, Douek DC, Price DA, et al. TCR beta-chain sharing in human CD8+ T cell
613 responses to cytomegalovirus and EBV. *J Immunol*. 2008;181(11):7853-62.
- 614 8. Beckman EM, Porcelli SA, Morita CT, Behar SM, Furlong ST, and Brenner MB.
615 Recognition of a lipid antigen by CD1-restricted alpha beta+ T cells. *Nature*.
616 1994;372(6507):691-4.

- 617 9. Yang Z, and Swanson WJ. Codon-substitution models to detect adaptive evolution that
618 account for heterogeneous selective pressures among site classes. *Molecular biology*
619 *and evolution*. 2002;19(1):49-57.
- 620 10. Im JS, Kang TJ, Lee SB, Kim CH, Lee SH, Venkataswamy MM, Serfass ER, Chen B,
621 Illarionov PA, Besra GS, et al. Alteration of the relative levels of iNKT cell subsets is
622 associated with chronic mycobacterial infections. *Clinical immunology*. 2008;127(2):214-
623 24.
- 624 11. Chan AC, Leeansyah E, Cochrane A, d'Udekem d'Acoz Y, Mittag D, Harrison LC,
625 Godfrey DI, and Berzins SP. Ex-vivo analysis of human natural killer T cells
626 demonstrates heterogeneity between tissues and within established CD4(+) and CD4(-)
627 subsets. *Clin Exp Immunol*. 2013;172(1):129-37.
- 628 12. Montoya CJ, Pollard D, Martinson J, Kumari K, Wasserfall C, Mulder CB, Rugeles MT,
629 Atkinson MA, Landay AL, and Wilson SB. Characterization of human invariant natural
630 killer T subsets in health and disease using a novel invariant natural killer T cell-
631 clonotypic monoclonal antibody, 6B11. *Immunology*. 2007;122(1):1-14.
- 632 13. Lantz O, and Bendelac A. An invariant T cell receptor alpha chain is used by a unique
633 subset of major histocompatibility complex class I-specific CD4+ and CD4-8- T cells in
634 mice and humans. *J Exp Med*. 1994;180(3):1097-106.
- 635 14. Porcelli S, Yockey CE, Brenner MB, and Balk SP. Analysis of T cell antigen receptor
636 (TCR) expression by human peripheral blood CD4-8- alpha/beta T cells demonstrates
637 preferential use of several V beta genes and an invariant TCR alpha chain. *J Exp Med*.
638 1993;178(1):1-16.
- 639 15. Van Rhijn I, Kasmar A, de Jong A, Gras S, Bhati M, Doorenspleet ME, de Vries N,
640 Godfrey DI, Altman JD, de Jager W, et al. A conserved human T cell population targets
641 mycobacterial antigens presented by CD1b. *Nat Immunol*. 2013;14(7):706-13.

- 642 16. Kjer-Nielsen L, Patel O, Corbett AJ, Le Nours J, Meehan B, Liu L, Bhati M, Chen Z,
643 Kostenko L, Reantragoon R, et al. MR1 presents microbial vitamin B metabolites to
644 MAIT cells. *Nature*. 2012;491(7426):717-23.
- 645 17. Treiner E, Duban L, Bahram S, Radosavljevic M, Wanner V, Tilloy F, Affaticati P, Gilfillan
646 S, and Lantz O. Selection of evolutionarily conserved mucosal-associated invariant T
647 cells by MR1. *Nature*. 2003;422(6928):164-9.
- 648 18. Gold MC, Cerri S, Smyk-Pearson S, Cansler ME, Vogt TM, Delepine J, Winata E,
649 Swarbrick GM, Chua WJ, Yu YY, et al. Human mucosal associated invariant T cells
650 detect bacterially infected cells. *PLoS Biol*. 2010;8(6):e1000407.
- 651 19. Le Bourhis L, Martin E, Peguillet I, Guihot A, Froux N, Core M, Levy E, Dusseaux M,
652 Meyssonier V, Premel V, et al. Antimicrobial activity of mucosal-associated invariant T
653 cells. *Nat Immunol*. 2010;11(8):701-8.
- 654 20. Van Rhijn I, Ly D, and Moody DB. CD1a, CD1b, and CD1c in immunity against
655 mycobacteria. *Adv Exp Med Biol*. 2013;783(181-97).
- 656 21. Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel
657 JO, Marchini JL, McCarthy S, McVean GA, et al. A global reference for human genetic
658 variation. *Nature*. 2015;526(7571):68-74.
- 659 22. Nei M, and Li WH. Mathematical model for studying genetic variation in terms of
660 restriction endonucleases. *Proc Natl Acad Sci U S A*. 1979;76(10):5269-73.
- 661 23. Tajima F. Statistical method for testing the neutral mutation hypothesis by DNA
662 polymorphism. *Genetics*. 1989;123(3):585-95.
- 663 24. Reiser JB, Legoux F, Gras S, Trudel E, Chouquet A, Leger A, Le Gorrec M, Machillot P,
664 Bonneville M, Saulquin X, et al. Analysis of relationships between peptide/MHC
665 structural features and naive T cell frequency in humans. *J Immunol*.
666 2014;193(12):5816-26.

- 667 25. Parham P, and Ohta T. Population biology of antigen presentation by MHC class I
668 molecules. *Science*. 1996;272(5258):67-74.
- 669 26. Bjorkman PJ, and Parham P. Structure, function, and diversity of class I major
670 histocompatibility complex molecules. *Annual review of biochemistry*. 1990;59(253-88).
- 671 27. Batuwangala T, Shepherd D, Gadola SD, Gibson KJ, Zaccari NR, Fersht AR, Besra GS,
672 Cerundolo V, and Jones EY. The crystal structure of human CD1b with a bound bacterial
673 glycolipid. *J Immunol*. 2004;172(4):2382-8.
- 674 28. Seshadri C, Lin L, Scriba TJ, Peterson G, Freidrich D, Frahm N, DeRosa SC, Moody
675 DB, Prandi J, Gilleron M, et al. T Cell Responses against Mycobacterial Lipids and
676 Proteins Are Poorly Correlated in South African Adolescents. *J Immunol*.
677 2015;195(10):4595-603.
- 678 29. Kasmar AG, van Rhijn I, Cheng TY, Turner M, Seshadri C, Schiefner A, Kalathur RC,
679 Annand JW, de Jong A, Shires J, et al. CD1b tetramers bind alphabeta T cell receptors
680 to identify a mycobacterial glycolipid-reactive T cell repertoire in humans. *J Exp Med*.
681 2011;208(9):1741-7.
- 682 30. Gras S, Van Rhijn I, Shahine A, Cheng TY, Bhati M, Tan LL, Halim H, Tuttle KD, Gapin
683 L, Le Nours J, et al. T cell receptor recognition of CD1b presenting a mycobacterial
684 glycolipid. *Nature communications*. 2016;7(13257).
- 685 31. Murugan A, Mora T, Walczak AM, and Callan CG, Jr. Statistical inference of the
686 generation probability of T-cell receptors from sequence repertoires. *Proc Natl Acad Sci*
687 *U S A*. 2012;109(40):16161-6.
- 688 32. Dean J, Emerson RO, Vignali M, Sherwood AM, Rieder MJ, Carlson CS, and Robins
689 HS. Annotation of pseudogenic gene segments by massively parallel sequencing of
690 rearranged lymphocyte receptor loci. *Genome medicine*. 2015;7(123).
- 691 33. Li H, Ye C, Ji G, and Han J. Determinants of public T cell responses. *Cell research*.
692 2012;22(1):33-42.

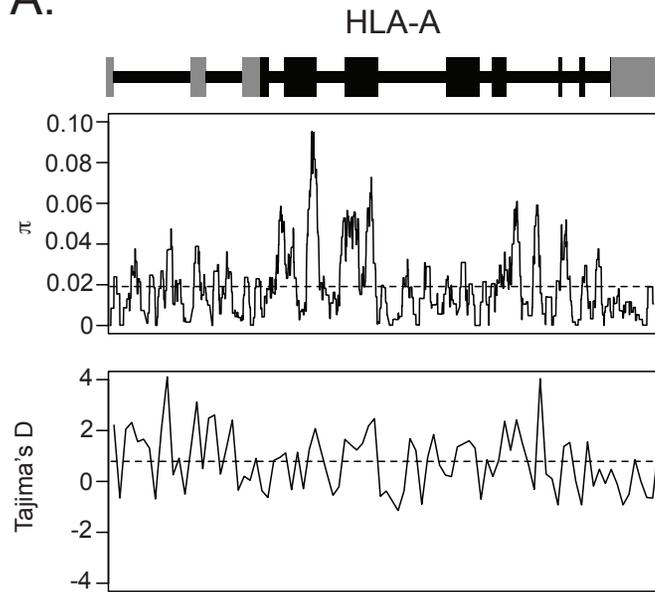
- 693 34. Venturi V, Price DA, Douek DC, and Davenport MP. The molecular basis for public T-cell
694 responses? *Nature reviews*. 2008;8(3):231-8.
- 695 35. Borg NA, Wun KS, Kjer-Nielsen L, Wilce MC, Pellicci DG, Koh R, Besra GS, Bharadwaj
696 M, Godfrey DI, McCluskey J, et al. CD1d-lipid-antigen recognition by the semi-invariant
697 NKT T-cell receptor. *Nature*. 2007;448(7149):44-9.
- 698 36. Lepore M, de Lalla C, Gundimeda SR, Gsellinger H, Consonni M, Garavaglia C,
699 Sansano S, Piccolo F, Scelfo A, Haussinger D, et al. A novel self-lipid antigen targets
700 human T cells against CD1c(+) leukemias. *J Exp Med*. 2014;211(7):1363-77.
- 701 37. Jarrett R, Salio M, Lloyd-Lavery A, Subramaniam S, Bourgeois E, Archer C, Cheung KL,
702 Hardman C, Chandler D, Salimi M, et al. Filaggrin inhibits generation of CD1a neolipid
703 antigens by house dust mite-derived phospholipase. *Science translational medicine*.
704 2016;8(325):325ra18.
- 705 38. Cheung KL, Jarrett R, Subramaniam S, Salimi M, Gutowska-Owsiak D, Chen YL,
706 Hardman C, Xue L, Cerundolo V, and Ogg G. Psoriatic T cells recognize neolipid
707 antigens generated by mast cell phospholipase delivered by exosomes and presented
708 by CD1a. *J Exp Med*. 2016;213(11):2399-412.
- 709 39. Gilleron M, Stenger S, Mazorra Z, Wittke F, Mariotti S, Bohmer G, Prandi J, Mori L, Puzo
710 G, and De Libero G. Diacylated sulfoglycolipids are novel mycobacterial antigens
711 stimulating CD1-restricted T cells during infection with *Mycobacterium tuberculosis*. *J*
712 *Exp Med*. 2004;199(5):649-59.
- 713 40. Spurgin LG, and Richardson DS. How pathogens drive genetic diversity: MHC,
714 mechanisms and misunderstandings. *Proceedings Biological sciences*.
715 2010;277(1684):979-88.
- 716 41. Deitsch KW, Lukehart SA, and Stringer JR. Common strategies for antigenic variation by
717 bacterial, fungal and protozoan pathogens. *Nat Rev Microbiol*. 2009;7(7):493-503.

- 718 42. Colijn C, Brandes A, Zucker J, Lun DS, Weiner B, Farhat MR, Cheng TY, Moody DB,
719 Murray M, and Galagan JE. Interpreting expression data with metabolic flux models:
720 predicting Mycobacterium tuberculosis mycolic acid production. *PLoS Comput Biol*.
721 2009;5(8):e1000489.
- 722 43. June CH, Riddell SR, and Schumacher TN. Adoptive cellular therapy: a race to the finish
723 line. *Science translational medicine*. 2015;7(280):280ps7.
- 724 44. Moody DB, Briken V, Cheng TY, Roura-Mir C, Guy MR, Geho DH, Tykocinski ML, Besra
725 GS, and Porcelli SA. Lipid length controls antigen entry into endosomal and
726 nonendosomal pathways for CD1b presentation. *Nat Immunol*. 2002;3(5):435-42.
- 727 45. Riddell SR, Watanabe KS, Goodrich JM, Li CR, Agha ME, and Greenberg PD.
728 Restoration of viral immunity in immunodeficient humans by the adoptive transfer of T
729 cell clones. *Science*. 1992;257(5067):238-41.
- 730 46. Mahomed H, Hawkrige T, Verver S, Geiter L, Hatherill M, Abrahams DA, Ehrlich R,
731 Hanekom WA, and Hussey GD. Predictive factors for latent tuberculosis infection among
732 adolescents in a high-burden area in South Africa. *Int J Tuberc Lung Dis*.
733 2011;15(3):331-6.
- 734 47. Day CL, Abrahams DA, Lerumo L, Janse van Rensburg E, Stone L, O'Rie T, Pienaar B,
735 de Kock M, Kaplan G, Mahomed H, et al. Functional capacity of Mycobacterium
736 tuberculosis-specific T cell responses in humans is associated with mycobacterial load. *J*
737 *Immunol*. 2011;187(5):2222-32.
- 738 48. de Jong A, Pena-Cruz V, Cheng TY, Clark RA, Van Rhijn I, and Moody DB. CD1a-
739 autoreactive T cells are a normal component of the human alphabeta T cell repertoire.
740 *Nat Immunol*. 2010;11(12):1102-9.
- 741 49. Moody DB, Reinhold BB, Guy MR, Beckman EM, Frederique DE, Furlong ST, Ye S,
742 Reinhold VN, Sieling PA, Modlin RL, et al. Structural requirements for glycolipid antigen
743 recognition by CD1b-restricted T cells. *Science*. 1997;278(5336):283-6.

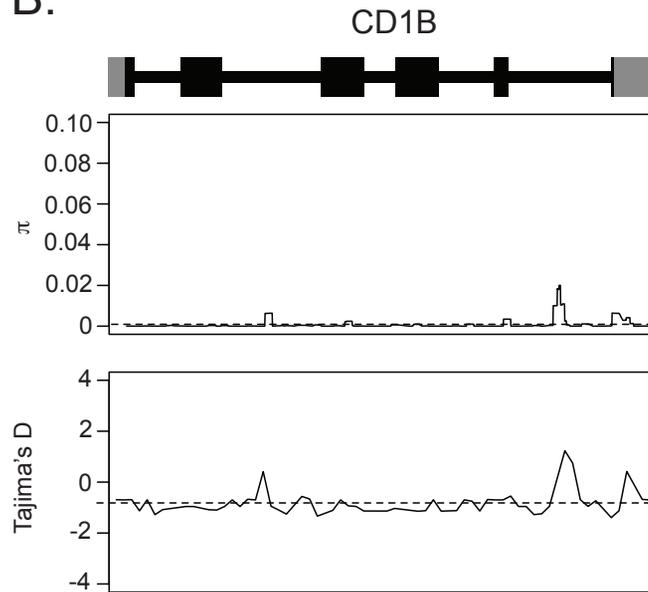
- 744 50. De Rosa SC, Carter DK, and McElrath MJ. OMIP-014: validated multifunctional
745 characterization of antigen-specific human T cells by intracellular cytokine staining.
746 *Cytometry Part A : the journal of the International Society for Analytical Cytology*.
747 2012;81(12):1019-21.
- 748 51. Horton H, Thomas EP, Stucky JA, Frank I, Moodie Z, Huang Y, Chiu YL, McElrath MJ,
749 and De Rosa SC. Optimization and validation of an 8-color intracellular cytokine staining
750 (ICS) assay to quantify antigen-specific T cells induced by vaccination. *J Immunol*
751 *Methods*. 2007;323(1):39-54.
- 752 52. Carlson CS, Emerson RO, Sherwood AM, Desmarais C, Chung MW, Parsons JM, Steen
753 MS, LaMadrid-Herrmannsfeldt MA, Williamson DW, Livingston RJ, et al. Using synthetic
754 templates to design an unbiased multiplex PCR assay. *Nature communications*.
755 2013;4(2680).
- 756 53. Yousfi Monod M, Giudicelli V, Chaume D, and Lefranc MP. IMGT/JunctionAnalysis:
757 the first tool for the analysis of the immunoglobulin and T cell receptor complex V-J
758 and V-D-J JUNCTIONS. *Bioinformatics*. 2004;20 Suppl 1(i379-85).

Dewitt et al. Figure 1

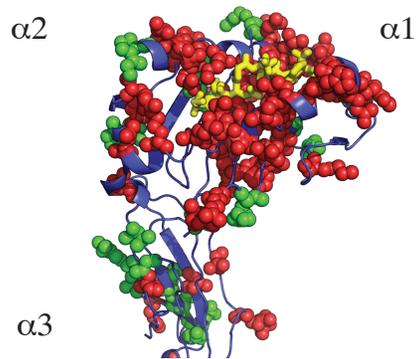
A.



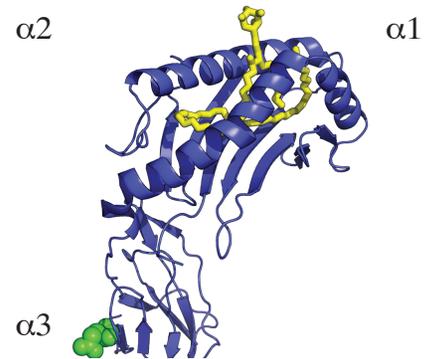
B.



C.



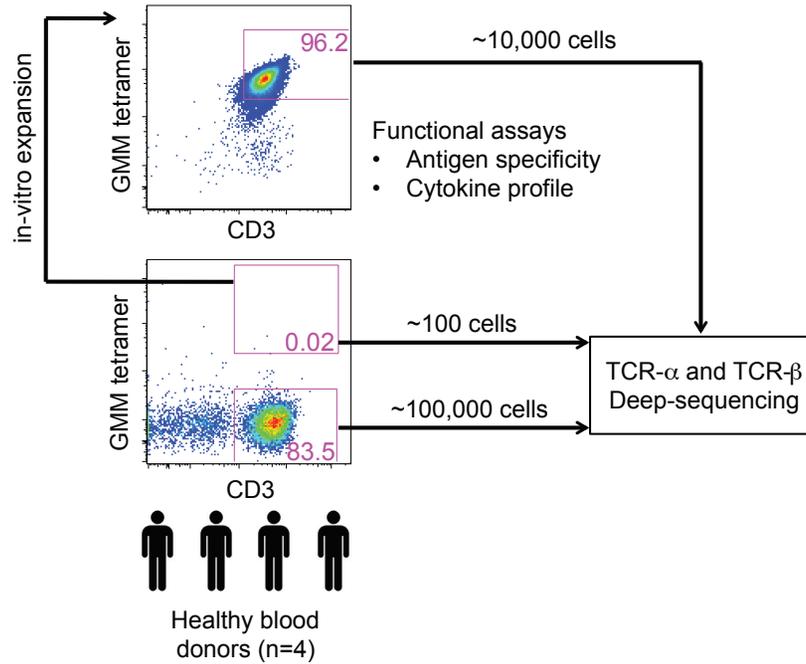
D.



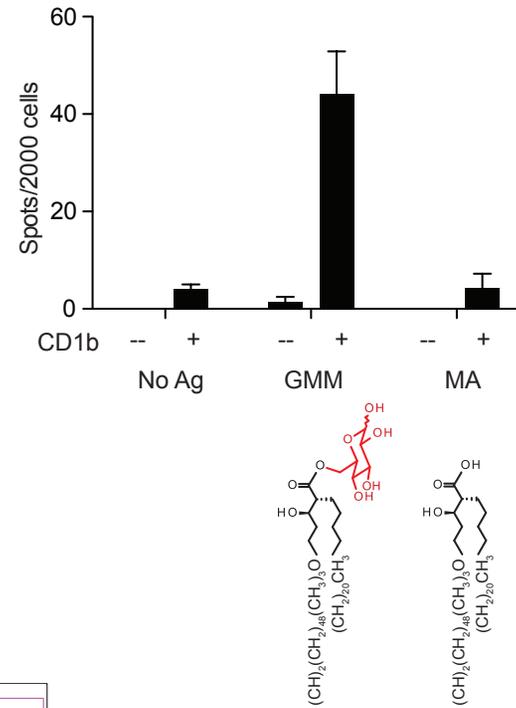
760 **Figure 1. Human genetic variation in HLA-A and CD1B.** The schematic of each gene structure, including introns (as lines) and
761 exons (as blocks), with coding sequence indicated in black and untranslated regions in grey. Nucleotide diversity (π) and Tajima's D
762 are reported for (A) HLA-A2 or (B) CD1B with the average value across each gene denoted by a horizontal dashed line. Variation in
763 protein coding sequence for (C) HLA-A2 or (D) CD1B is represented for SNPs with >1% minor allele frequency. Invariant residues
764 are denoted in blue as a ribbon structure. Spheres denote polymorphic residues, with green representing synonymous substitutions
765 and red representing non-synonymous substitutions. The ligand, either peptide or lipid, is shown in yellow.

Dewitt et al. Figure 2

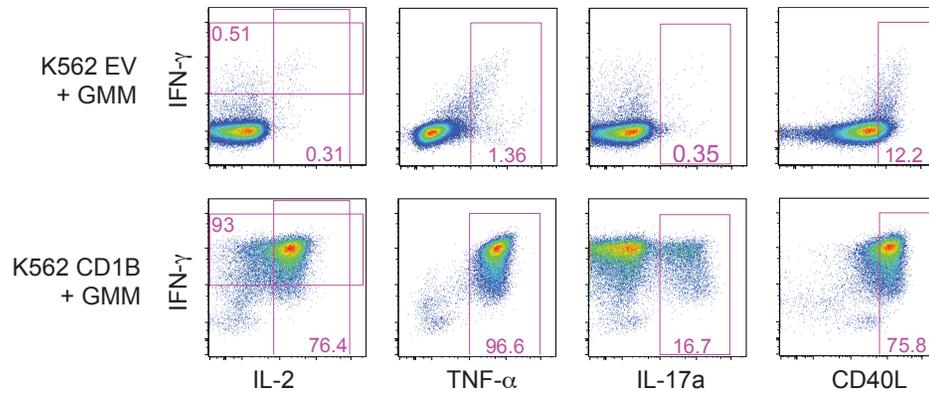
A.



B.

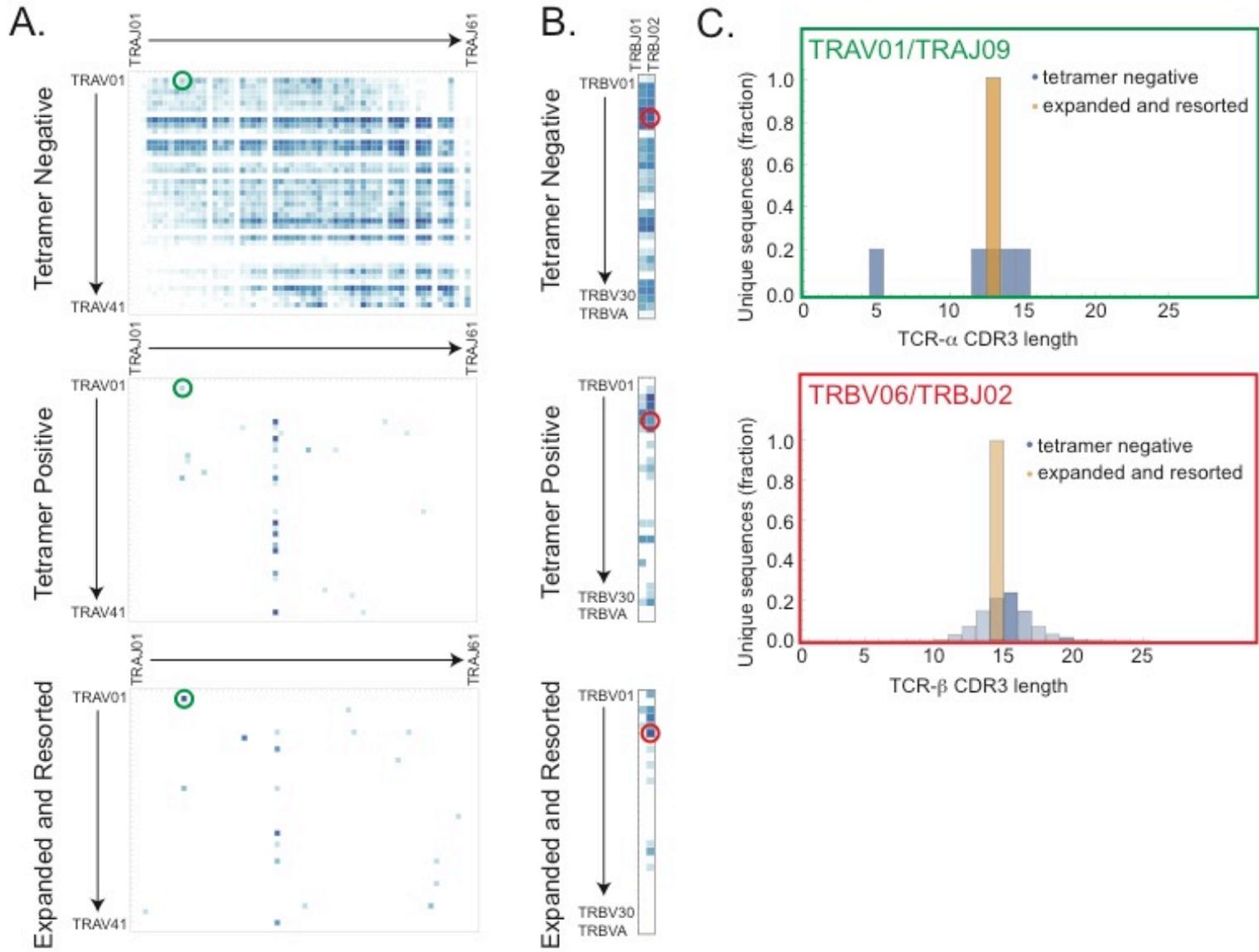


C.



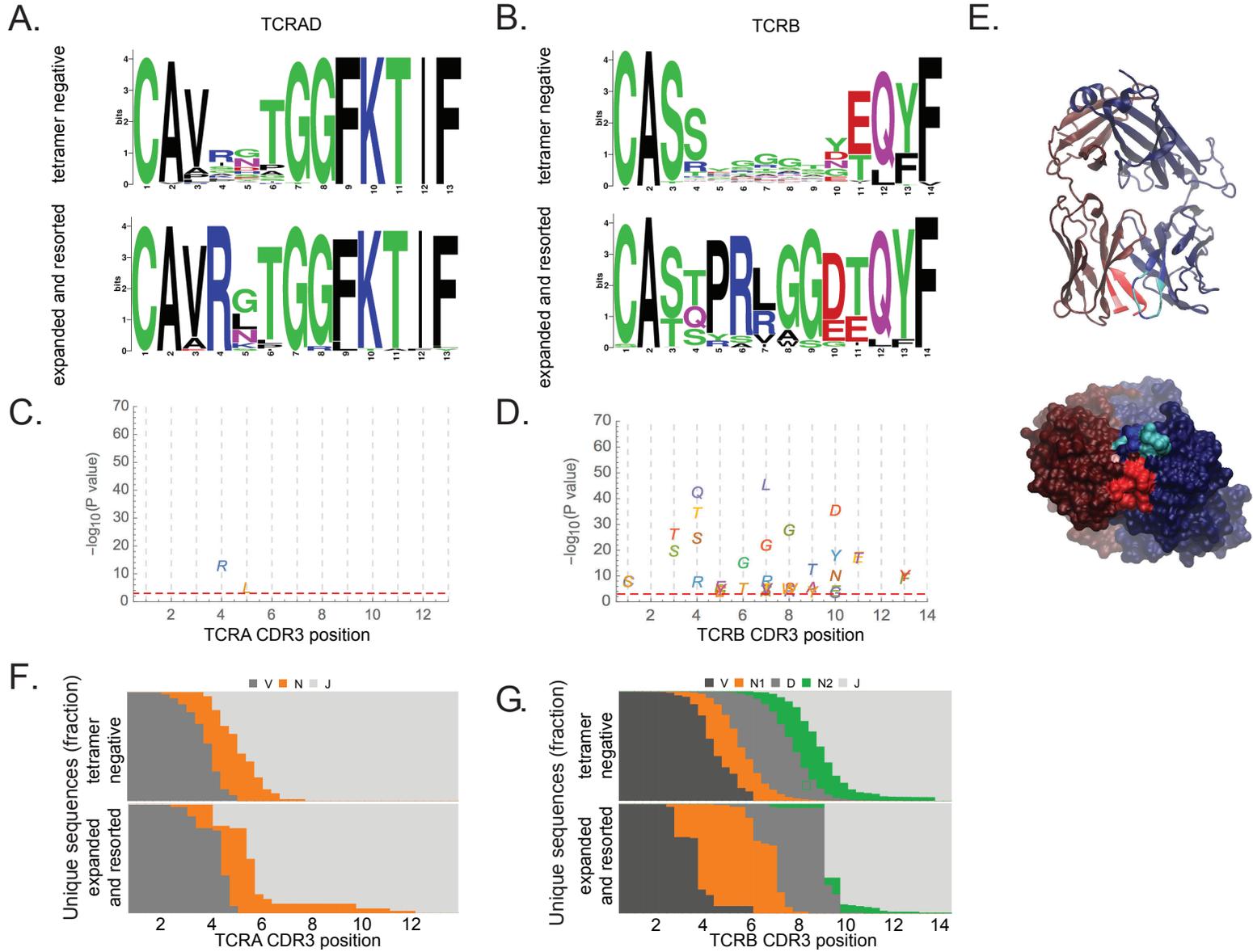
767 **Figure 2. Study schema and characterization of T-cell lines.** (A) We used GMM-loaded CD1B tetramers to stain and sort T cells
768 from cryopreserved PBMCs derived from four healthy South African blood donors. These cells were either used to generate T-cell
769 lines, or submitted directly for TCR immunosequencing. Tetramer-negative T cells as well as tetramer-positive T cells that were re-
770 sorted from T cell lines were also submitted for TCR immunosequencing. (B) T-cell lines were incubated with either mock
771 transfected or CD1B-transfected K562 cells as antigen-presenting cells in the presence of no antigen, GMM, or mycolic acid (MA).
772 Partial structures of the two lipid antigens are shown. IFN- γ production by ELISPOT was quantified after overnight incubation. (C)
773 Mock-transfected or CD1B-transfected K562 cells were loaded with GMM overnight and used to activate T cells for six hours prior to
774 intracellular cytokine staining. Data in B and C are representative of three independent experiments.
775

Dewitt et al. Figure 3



777 **Figure 3. GMM-specific T-cell receptor diversity.** Heatmaps depicting the diversity of (A) TRAV/TRAJ or (B) TRBV/TRBJ
778 recombination events found among tetramer-negative, tetramer-positive, or expanded and resorted T cells obtained from one
779 representative blood donor with latent tuberculosis infection. Blue color scale indicates the relative abundance of each V/J
780 combination. The green circle represents TRAV1-2 and TRAJ9 rearrangement. The red circle represents the TRBV6 and TRBJ2
781 rearrangement (C) Histograms show the distribution of CDR3 length among the most highly enriched V and J gene segment
782 combinations for TCR- α (TRAV1/TRAJ9) and TCR- β (TRBV6/TRBJ2).

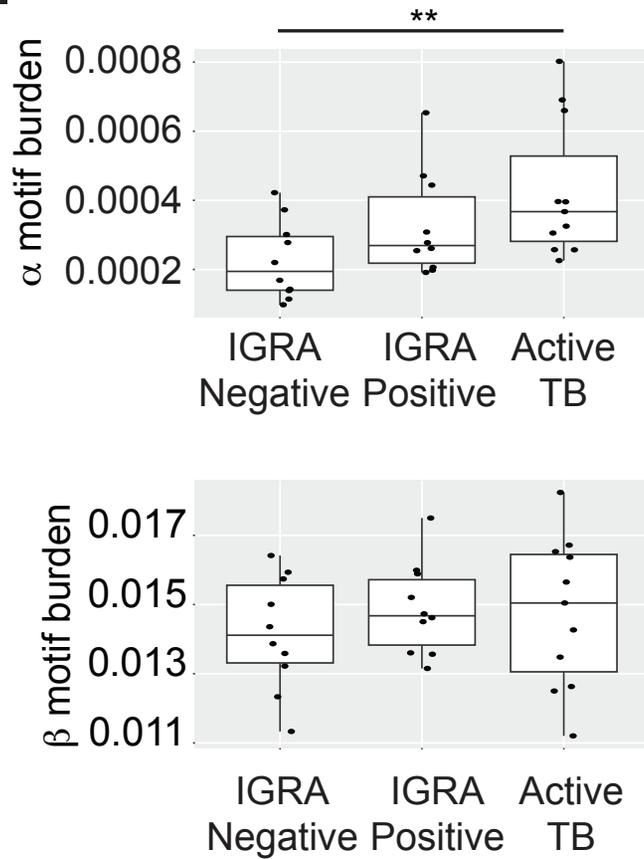
Dewitt et al. Figure 4



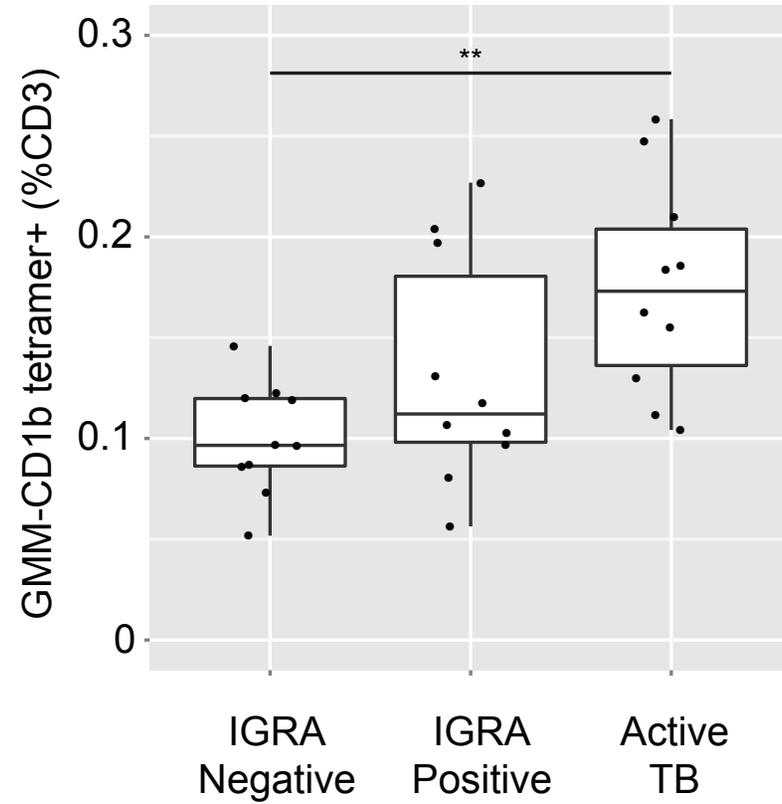
784 **Figure 4. Conservation among sequences matching the CD1B-GMM specific TCR motifs.** Logo plots indicating the positional
785 CDR3 amino acid usage among sequences within (A) TRAV01 and TRAJ09 rearrangement for TCR- α and (B) TRBV06 and TRBJ02
786 rearrangement for TCR- β among expanded and resorted samples as well as tetramer-negative samples. At each position along the
787 CDR3 region, the statistical significance of enrichment of each residue in the expanded and resorted samples as compared to
788 tetramer-negative samples is indicated for (C) TCR- α and (D) TCR- β . The selected significance threshold of p-value = 10^{-3} is
789 indicated with a red dashed line. (E) The crystal structure of a published CD1B-GMM specific TCR whose sequence matches both
790 the TCR- α and TCR- β motifs detected in this study. The TCR- α CDR3 is indicated in red, with position of significantly enriched
791 residues indicated in pink. The TCR- β CDR3 is indicated in blue, with position of significantly enriched residues indicated in cyan. A
792 portion of the TCR- α CDR3 was not resolved in the crystal structure, resulting in an apparent gap. At each position along the (F)
793 TCR- α or (G) TCR- β CDR3 region, the proportion of sequences annotated as originating in a germline V, J, or D gene, as well as
794 non-templated (N) bases at that position is depicted.

Dewitt et al. Figure 5

A.



B.



795

796

797 **Figure 5. GMM-specific TCR motifs and T cells are increased in frequency during tuberculosis infection and disease.** For
798 each subject, we calculated the fraction of TCRs matching the motifs defined in this study that were present in the peripheral blood of
799 IGRA-negative, IGRA-positive, and active TB patients. (A) Within TCR- α , there is an increased motif burden in active TB subjects
800 compared to IGRA-negative subjects (Kruskal-Wallis $p=0.03$, Dunn post-test $p=0.0038$). There was not a statistically significant
801 difference in subject-specific TCR- β motif burdens among the three groups (Kruskal-Wallis $p=0.67$). (B) The frequency of GMM-
802 CD1B tetramer positive T cells by flow cytometry is increased in patients with active TB as compared to IGRA-negative subjects
803 (Kruskal-Wallis $p=0.01$, Dunn post-test $p=0.0013$).
804

805

Table 1. Frequency of CD1B-GMM specific TCR motifs among sorted cells and PBMC from healthy bone marrow donors.

	Sequence Source	Motif Match	Motif Mismatch	Match/Mismatch	Fold Enrichment	p-value
TCR- α	Tetramer-negative	55	174409	0.0003152		
	Tetramer-positive	52	573	0.0832	263.92	6.21×10^{-98}
	Expanded & resorted	46	123	0.2721	863.40	3.34×10^{-113}
TCR- β	Tetramer-negative	7613	462440	0.01619		
	Tetramer-positive	126	566	0.1821	11.24	9.88×10^{-89}
	Expanded & resorted	119	255	0.3182	19.65	5.50×10^{-115}
	Bone marrow donors	2172898	130501755	0.0164	1.01	0.17

806

807 **Table 1. Frequency of CD1B-GMM specific TCR motifs among sorted cells and PBMC from bone marrow donors.** We pooled
808 unique sequences from all individuals for each sample type and each locus and counted the number matching (or mismatching) the
809 defined CD1B-GMM specific TCR motifs. Tetramer-positive and expanded and resorted samples exhibit significant enrichment of
810 motifs in both TCR- α and TCR- β as compared to tetramer-negative samples. PBMC from 587 healthy bone marrow donors have a
811 similar number of TCR- β motif matches as tetramer-negative samples. P-values were computed with Fisher's Exact test.

Table 2. Frequency of CD1B-GMM specific TCR motifs among PBMCs from South Africans with known M.tb infection status.

	Sequence Source	Motif Match	Motif Mismatch	Match/Mismatch	Fold Enrichment	p-value
TCR- α	IGRA-negative	180	817804	0.0002201		
	IGRA-positive	228	765928	0.0002976	1.35	1.40×10^{-3}
	Active TB	224	628088	0.0003565	1.62	8.29×10^{-7}
TCR- β	IGRA-negative	7565	532381	0.01401		
	IGRA-positive	8233	544217	0.01490	1.06	4.86×10^{-5}
	Active TB	8495	563097	0.01486	1.06	8.75×10^{-5}

812

813 **Table 2. Frequency of CD1B-GMM specific TCR motifs among M.tb-uninfected and M.tb-infected South African subjects.**

814 We performed immunosequencing for each locus from IGRA-negative (n=10), IGRA-positive (n=10), and active TB (n=10) subjects,
815 and then pooled the sequences for each locus from all subjects with a specific infection status. We tabulated the number of
816 sequences in the pool matching the defined CD1B-GMM specific TCR motif, as well as the complementary number of sequences not
817 matching the motif. Both IGRA-positive subjects and patients with active TB exhibit significant enrichment of both TCR- α and TCR- β
818 motifs as compared to IGRA-negative subjects. P-values were computed using Fisher's Exact test on the raw counts of matches and
819 mismatches. Because sampling depths varied widely among subjects, this pooling approach complements the subject-wise analysis
820 included in Fig. 5b, which eliminated sampling depth information by computing fractional motif burdens