

1 **CD4+ T cell subsets present stable relationships in their T cell receptor repertoires**

2 Shiyu Wang^{1,2,3}, Longlong Wang^{1,2,3}, Ya Liu^{2,3,4}

3

4 1. BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, China

5 2. BGI-Shenzhen, Shenzhen 518083, China

6 3. China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China

7 4. Shenzhen Key Laboratory of Single-Cell Omics, BGI-Shenzhen, Shenzhen, 518100, China

8

9 Correspondence: liuyal@genomics.cn (Y. Liu)

10

11 **Abstract**

12 CD4⁺ T cells are key components of adaptive immunity. The cell differentiation equips CD4⁺
13 T cells with new functions. However, the effect of cell differentiation on T cell receptor (TCR)
14 repertoire is not investigated. Here, we examined the features of TCR beta (TCRB) repertoire
15 of the top clones within naïve, memory and regular T cell (Treg) subsets: repertoire structure,
16 gene usage, length distribution and sequence composition. First, we found that memory subsets
17 and Treg would be discriminated from naïve by the features of TCRB repertoire. Second, we
18 found that the correlations between the features of memory subsets and naïve were positively
19 related to differentiation levels of memory subsets. Third, we found that public clones presented
20 a reduced proportion and a skewed sequence composition in differentiated subsets. Furthermore,
21 we found that public clones led naïve to recognize a broader spectrum of antigens than other
22 subsets. Our findings suggest that TCRB repertoire of CD4⁺ T cell subsets is skewed in a
23 differentiation-dependent manner. Our findings show that the variations of public clones
24 contribute to these changes. Our findings indicate that the reduce of public clones in
25 differentiation trim the antigen specificity of CD4⁺ T cells. The study unveils the physiological
26 effect of memory formation and facilitates the selection of proper CD4⁺ subset for cellular
27 therapy.

28 **Key words:** CD4⁺ T cell, T cell receptor beta chain repertoire, differentiation, public clones

29

30 **Introduction**

31 CD4⁺ T cells play critical roles in mediating adaptive immunity. Via T cell receptors
32 (TCR), CD4⁺ T cells recognize the complex of epitopes and major histocompatibility complex
33 II and then induce the activation of other cells in infections^{1,2}, cancer³ and autoimmune diseases.

34 To acquire mature functions, CD4⁺ T cells undergo differentiation. NT is the prototype of
35 CD4⁺ T cell and has the greatest potential among CD4⁺ T subsets to differentiate to other
36 subsets. NT usually keep a serenity and can refresh themselves by proliferation. When NT
37 encounters pathogens, it will home to lymphatic organs and receive the help from dendritic
38 cells to initiate the polarization. The study on TCR repertoire suggests that NT has the most
39 large scale of evenness of TCR repertoire among all CD4⁺ subsets⁴, which indicates the
40 greatest potential to recognize antigens. In a classical differentiation model^{5,6}, naive (NT)
41 senses stimulations via TCR, polarizes and then differentiates to effector T (ET). ET plays the
42 key role to mediate adaptive immune response, although the amount of ET in peripheral blood
43 is limited. After few weeks' activation, a small part of ET differentiates to memory. Memory is
44 sensitive to antigens, while always keeps silent. In peripheral blood, central memory (CMT),
45 effector memory (EMT) and stem-cell like memory T cell (Tscm) are main subsets of memory.
46 CMT and Tscm have a potential to be self-renewal and is found to affect the infections. EMT,
47 compared to CMT⁷, is long-lived and has a lower threshold to reactivate to pathogens.
48 Compared to NT, memory subsets have a lower diversity of TCRB repertoire⁸. However, it is
49 unclear that large differences between NT and memory subsets exist in the functions of TCR
50 repertoire or not. Treg is a special subset of CD4⁺ cell and usually plays a role to regulate
51 functions, proliferation, and differentiation of conventional T cells⁹.

52 TCR plays the key role to determine T cell functions¹⁰. For NT, signals via diverse TCRs
53 reform the TCR repertoire, and skew their differentiation potential. For instance, strong TCR
54 signals during viral infection correspond to helper T cell differentiation, and comparative lower
55 signals facilitate the differentiation of memory and follicular T cell¹¹. Memory's TCR repertoire
56 composition determines the possibility to provide a rapid protection for individuals to against
57 former and, sometimes, novel antigens. For example, the architecture of the TCR repertoire
58 contributes to the performance of the adaptive immune response against pathogens, such as

59 SARS-CoV-2¹². Cross-reactivation from memory can provide a rapid protection to a novel
60 pathogen in some individuals, such as the case reports of COVID-19¹³. The importance of TCR
61 repertoire for memory cell functions was found in tissues, where the differential composition
62 of TCR repertoire of CD4⁺ memory among tissues equipped them with distinct functions¹⁴.
63 The function of Treg was restricted by TCR repertoire. The optimal diversity of TCR was
64 essential for the suppressive ability¹⁵, and limitations on TCR diversity disturbed the self-
65 tolerance of immune system¹⁶. Although evidences show that the features of TCR repertoire are
66 distinct among CD4⁺ T subsets, the effect of differentiation on T cell receptor (TCR) repertoire
67 of CD4⁺ T cells are not investigated.

68 To unveil the influence of differentiation on TCR repertoire, we analyzed the sequencing
69 data of TCR beta (TCRB) chain of NT, ET, EMT, CMT, Tscm and Treg. We detected repertoire
70 structure, germline gene usage, sequence composition and public clones of TCRB repertoire of
71 each subset. We found that NT, CMT, Tscm and Treg were discriminated from each other by
72 repertoire structure, gene usage and sequence composition, independently. The TCRB
73 repertoires of NT are similar to the TCRB repertoires of less-differentiated memory subset (for
74 example, Tscm and CMT). The TCRB repertoires of ET and the TCRB repertoires of EMT are
75 sensitive to the healthy state and have flexible relationships with other subsets. Public clones
76 account for the main part of top clones of NT and reduce along the CD4⁺ cell differentiation.
77 The enrichment of public clones shortens the length distribution of top clones in NT. Public
78 clones are polyfunctional and broaden the antigen spectrum recognized by NT. Our findings
79 disclose the differential functions of CD4⁺ cell subsets and the influence of differentiation on
80 TCR repertoire. Our findings facilitate the selection of CD4⁺ subsets for cellular therapy.

81

82 **Material and Methods**

83 **Datasets**

84 In this study, we conducted analyses on high-throughput TCR repertoire datasets for
85 CD4+ T cell subsets from publication described as follows.

86 Dataset1 included Naïve (NT), central memory (CMT), stem-cell like memory (Tscm)
87 and Treg from eight healthy individuals and eight type-one-diabetes (T1D) patients¹⁷. CD4+
88 T cells were sorted into subsets by fluorescence-activated cell sorting (FACS), and then RNA
89 was extracted and sequenced in parallel. For validation, we employed dataset2 of five CD4+
90 T cell subsets: NT, effector (ET), CMT, effector memory (EMT) and Treg from other ten
91 rheumatoid arthritis patients (RA)¹⁸. Top1000 clones are referred as the 1000 clones with the
92 highest frequency. The V-/J-segments used by top1000 clones were extracted for statistic.

93 **Statistical analysis and plots**

94 Statistical analyses were performed with R. Significance was examined by Willcox ranked
95 test. The correlation coefficients and significance was tested using `cor.test()` in R with default
96 parameters. Graphics were generated with R package `ggplot2`. Principle component analysis
97 (PCA) was conducted with R package `forcats`. Data was treated with R package `readr`, `dplyr`
98 and `tidyr`.

99 **Definition of a clone**

100 For all analyses, clones were defined as the amino acid sequence identity of CDR3 (TCRB)
101 regions. CDR3s from dataset1 were defined and annotated by IMonitor¹⁹, and CDR3s from
102 dataset2 were reported by authors.

103 **Determination of diversity**

104 Renyi entropy was used in our study to evaluate the diversity with alpha value from 0 to
105 20. When alpha is 1, the index equals to the Shannon index. Renyi entropy formula is $H =$
106 $\frac{1}{1-q} \ln(\sum_{i=1}^R p_i^q)$, and Shannon diversity index formula is $H = -\sum_{i=1}^R p_i \ln p_i$, where H is
107 the diversity index, q is alpha value, R is the total number of clones for analysis, and p_i is the
108 frequency of the i th clone.

109

110 **Determination of Jensen-Shannon distance**

111 Jensen-Shannon distance (JSD) is used to evaluate the similarity in repertoire architecture
112 among subsets²⁰. We calculated JSD with `JSD()`, a function included in `philentropy` (²¹), a R
113 package. A low JSD indicates that TCRB repertoire structures are similar.

114 **Identify the contribution of k-mer to PCA classification**

115 We used `prcomp()` to calculate the principle components (PC) for data, and estimated the
116 contribution of each k-mer to PC1 and PC2 with `cos2()`.

117 **KeBABS SVM analysis**

118 SVM analysis was performed using kernel-based analysis of biological sequences with a
119 R package `KeBABS`²². Amino acid sequence of clones was split into features with length $k =$
120 3, and cost parameter $C = 100$ was used for the misclassification of a sequence. For all SVM
121 analyses, data was split into training (80%) and test (20%) set. SVM training was performed on
122 the training set, and class prediction was performed on the test set. Prediction accuracy of
123 classification was qualified by calculating $BACC = \frac{1}{2} \times (spec + sens)$, where specificity was
124 calculated as $spec = \frac{TN}{TN + FP}$, and sensitivity was defined as $sens = \frac{TP}{TP + FN}$, (where TN =
125 true negative, FP = false negative, TP = true positive and FN = false negative). The area under
126 the receiver operating characteristic curve (AUC) was calculated, where the AUC = 0.5 means
127 a random classification (BACC = 50%), and AUC = 1 means a perfect classification (BACC =
128 100%).

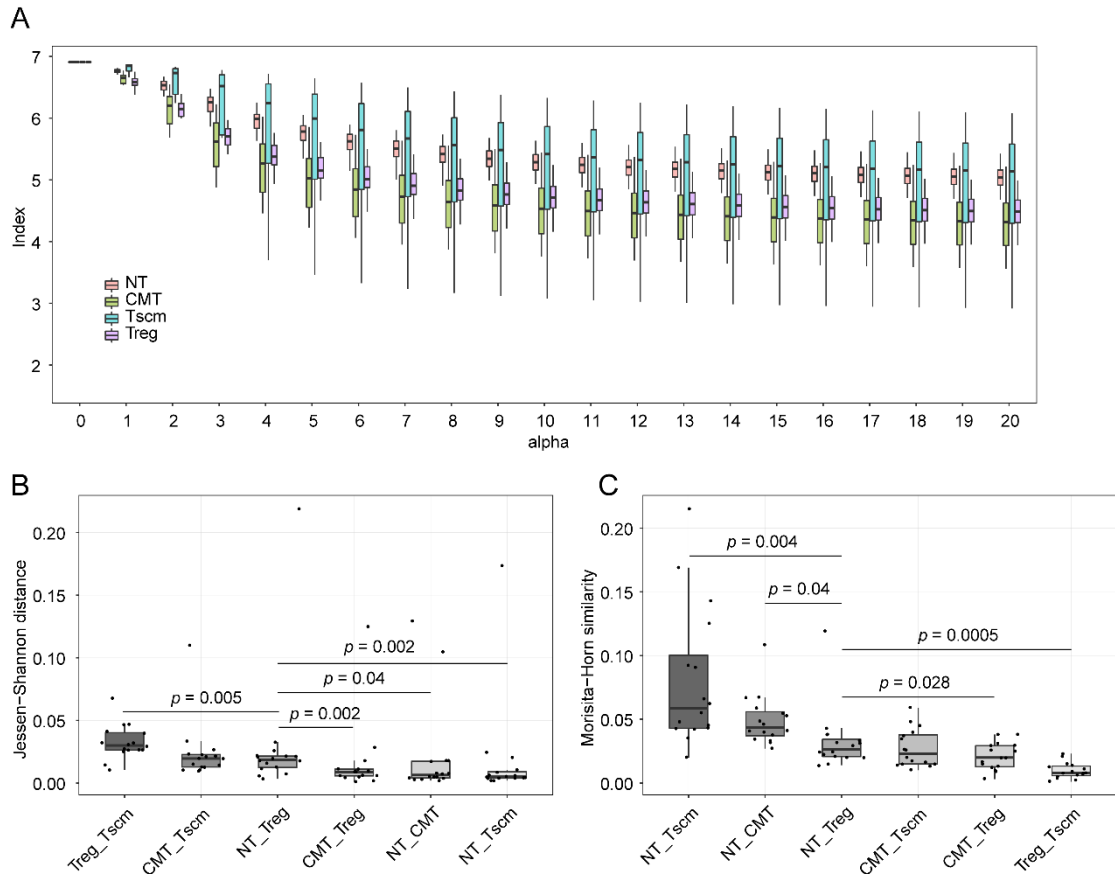
129 **Determination of epitope specificity of clones by GLIPH2**

130 GLIPH2²³ is a robust tool to predict the cluster of clones targeting the same epitope. Here, we
131 use this method to unveil the diversity of potential epitopes targeted by top1000 clones in each
132 subset. The reference of CD4+ T clones and their gene usage, the distribution of gene usage of
133 CD4+ T and length distribution of CDR3 were included in `ref_CD4.txt`, `ref_V_CD4.txt` and
134 `ref_L_CD4.txt`. These reference files were downloaded from the official website of GLIPH2
135 (<http://50.255.35.37:8080/>). A filter with a high stringency (`Fisher_score < 0.0001`,
136 `number_subject >= 3` and `number_unique_cdr3 >= 3`) was used to improve the prediction
137 accuracy.

138 **Results**

139 **The TCRB repertoire structure of NT is similar to The TCRB repertoire structure of CMT**
140 **and Tscm**

141 Frequent clones affect the immune repertoire structure²⁴. We thus performed the analyses
142 on top1000 clones within each subset. Renyi entropy with alpha values from zero to twenty was
143 used to evaluate the diversity. In dataset1, The TCRB repertoire of NT and Tscm present similar
144 diversities at all alpha values, and are more diverse than the TCRB repertoire of CMT and the
145 TCRB repertoire of Treg (Figure 1A). In dataset2, NT has the most diverse TCRB repertoire
146 among all subset whereas ET has the lowest. The TCRB repertoire of CMT is more diverse
147 than that of ETM and Treg. (Supplementary figure 1A). The similarity of TCRB repertoire
148 structure of subsets was estimated by Jensen-Shannon distance. In dataset1, the TCRB
149 repertoire structure of NT is similar to the TCRB repertoire structure of less-differentiated
150 subsets (CMT and Tscm), but the TCR repertoire structures of Tscm and CMT are different
151 from each other; the TCRB repertoire of Treg is different to the TCRB repertoire of NT and
152 CMT with high JSDs. It indicates that Treg has a structure of TCRB repertoire like that of more-
153 differentiated memory subsets (Figure 1B). In dataset2, NT and CMT have similar TCRB
154 repertoire structures, and the TCRB repertoire structure of Treg is similar to the TCRB
155 repertoire structure of EMT rather than that of CMT (Supplemental Figure 1B). These findings
156 fit with the trend found in dataset1. To consider the overlapping usage of CDR3 clones, we
157 further evaluated the similarity of TCRB repertoire among subsets with the Morisita-Horn
158 similarity index. In this analysis, NT keeps a similar TCRB repertoire like Tscm and CMT,
159 while the TCRB repertoire of NT is different from the TCRB repertoire of EMT and Treg
160 (Figure 1C; Supplemental Figure 1C). In conclusion, the TCRB repertoire structure of CD4+ T
161 cells is skewed along with cell differentiation levels.



162

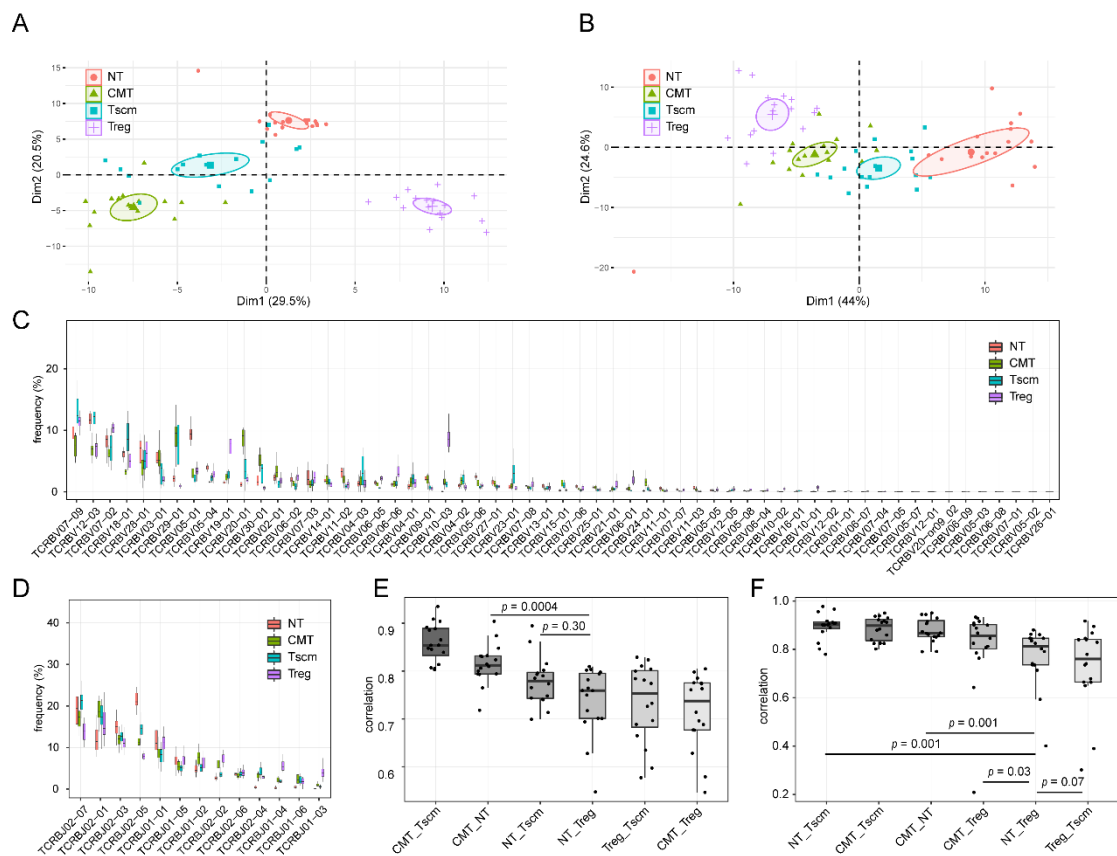
163 **Figure 1. The diversity of TCRB repertoire of four subsets and the relationship of their**
164 **TCRB repertoire structures. (A)** The Renyi entropy index of all subsets with alpha value from
165 0 to 20. When alpha is equal to 1, the index was calculated as Shannon-index. **(B)** The Jessen-
166 Shannon distance between subsets. The difference between NT and other subsets, and Treg and
167 other subsets were tested. **(C)** The Morisita-Horn similarity between subsets. Paired Wilcox-
168 ranked test was used in **B** and **C**.

169 **Gene usage is distinct among subsets and a part of genes are skewed along subsets**

170 We used principle component analysis (PCA) to examine the discrimination of gene usage
171 among subsets. NT, CMT, Tscm and Treg could be distinguished from each other by V- and J-
172 gene respectively (Figure 2A and B; Supplemental Figure 2A). However, EMT and ET could
173 not be separated from other subsets in dataset2 (Supplemental Figure 2A). We examined the
174 genes usage of each subset, and found that 24 of 72 genes were differently used by NT, Tscm
175 and CMT in dataset1, 23 of 72 genes were differently used by NT and CMT in dataset2. There
176 are 11 increased genes and 13 decreased genes in CMT, compared with NT in dataset1; and 3

177 increased genes and 20 decreased genes in dataset2, compared with NT in dataset2. Notably,
 178 these genes continuously changed according to NT, CMT, ET, EMT and Treg (Supplemental
 179 Figure 2C and D).

180 With the correlations of V-gene usage among subsets, we found that NT, CMT and Tscm
 181 show high correlations with each other, while their correlations with Treg are low (Figure 2E;
 182 Supplementary Figure 2E). For J-gene, NT also shows a high correlation with CMT and Tscm,
 183 whereas a low correlation with Treg. In dataset2, the V-gene usage of NT highly correlates to
 184 the V-gene usage of CMT rather than EMT and Treg (Supplementary Figure 2F). Notably, the
 185 V-gene usage of Treg is similar to the V-gene usage of NT rather than to the V-gene usage of
 186 CMT in the T1D donors ($p=0.0005$, Supplementary Figure 4) but not in health donors ($p=0.05$).
 187 This phenomenon suggests Treg has a flexible relationship with NT depending on the healthy
 188 states of donors. In conclusion, NT, memory subsets and Treg have distinct gene usages, and
 189 the gene usage of NT is more similar to less-differentiated memory subsets (CMT and Tscm)
 190 rather than EMT.



191

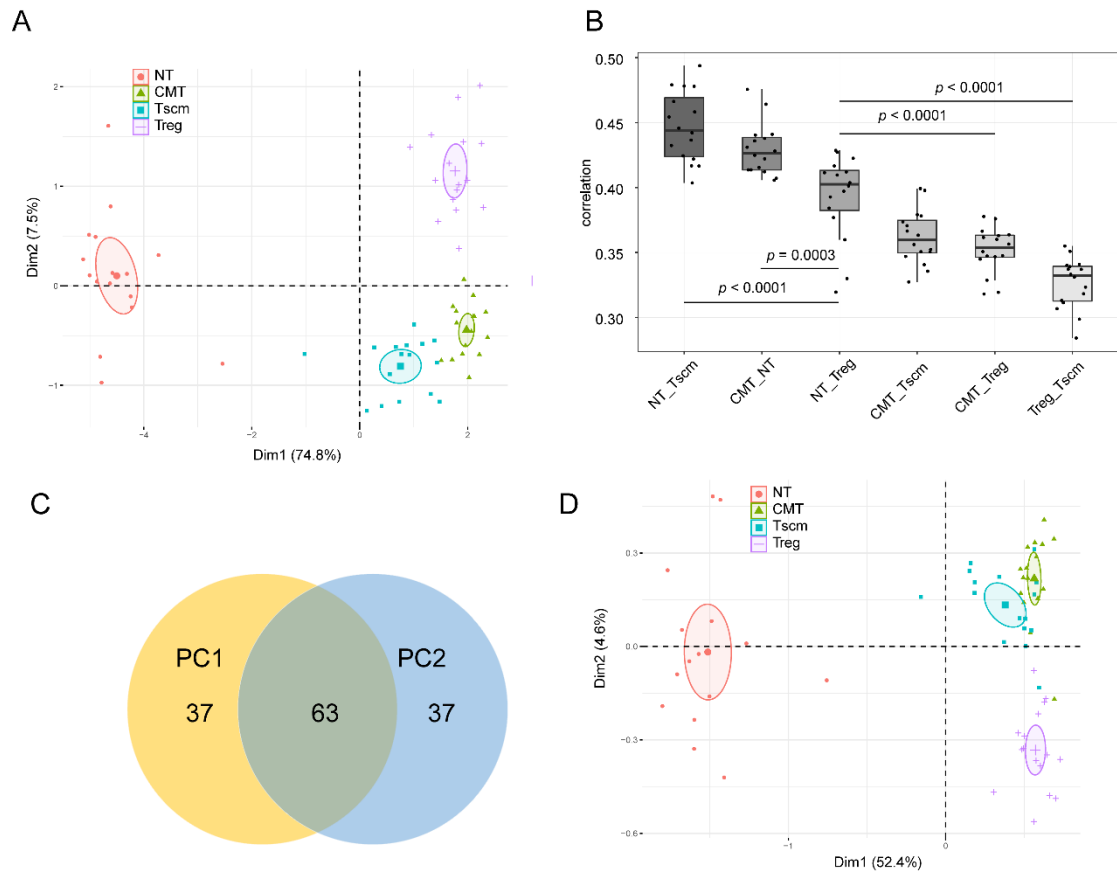
192 **Figure 2. The V- and J-gene usage of CD4+ T subsets and the relationships of the gene**

193 **usage of subsets.** (A) PCA for the frequency of V-genes of each subset. (B) PCA for the
194 frequency of J-genes of each subset. (C) The frequency of V-genes of each subset. (D) the
195 frequency of J-gene within each subset. (E) The Spearman correlation of V-gene usage between
196 subsets. (F) The Spearman correlation of J-gene usage between subsets. Paired Wilcox-ranked
197 test was used in E and F.

198 **CDR3 sequence composition are different among TCRB repertoire of subsets and indels**
199 **contribute to these differences highly**

200 We examined the CDR3 sequence composition by decomposing kernels containing three
201 amino acids. To identify the difference in the composition among subsets, we used PCA to
202 discriminate subsets. In PCA plot, subsets are distinguished from each other. Treg is close to
203 EMT, CMT mix with ET. NT has a clear boundary to others (Figure 3A; Supplemental Figure
204 5A). To evaluate the correlation of k-mer among subsets, we used Spearman correlation method.
205 The k-mer usage of NT exhibits weakened correlations with others according to Tscm, CMT
206 and Treg in dataset1, and CMT, ET, EMT and Treg in dataset2 (Figure 3B; Supplemental Figure
207 5B).

208 To identify the subregions where the k-mers contributes to PCA, we extracted the top100
209 k-mers by ranking their contributions to principle component 1 (PC1) and principle component
210 2 (PC2). 6371 unique k-mers were used by both of PC1 and PC2. We ranked k-mers by their
211 contributions to PC1 and PC2 respectively, and found that PC1 and PC2 shared 63 top100 k-
212 mers (Figure 3C; Supplemental Figure 5C). Then we aligned k-mer to references, and showed
213 that 29 top100 k-mers located in V-region and 47 in J-region. However, after we removed the
214 k-mer enriching in V-/J-segments, we found that NT, Tscm, CMT and Treg still kept distinctions
215 to each other (Figure 3D). we further perform a same operation on dataset2 to remove the k-
216 mer in V- and J-segments, and showed a similar result that NT, CMT and Treg are able to be
217 distinguished from each other (Supplemental Figure 5D). Since the gene usage, insertion and
218 deletion (indel) are factors to skew the CDR3 sequence compositions²⁵, these results suggests
219 that indel in N1-D-N2 region rather than gene usage contribute to the difference of sequence
220 composition among subsets.



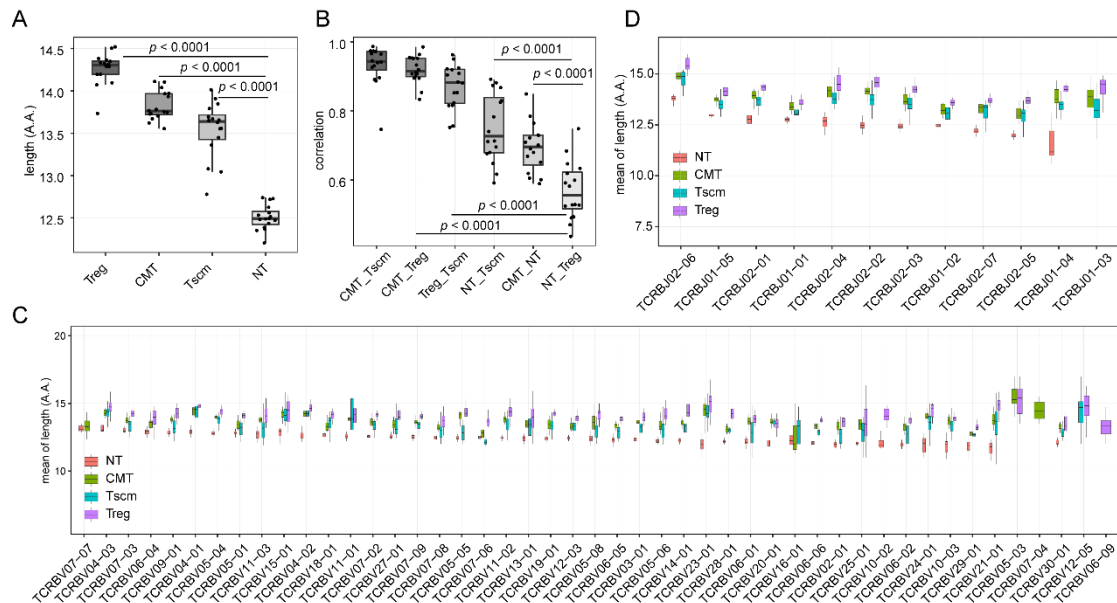
221

222 **Figure 3. The relationship of sequence composition among subsets.** (A) PCA based on the
 223 sequence composition for samples. (B) Spearman correlations of k-mer usage between subsets.
 224 (C) The overlap of top100 k-mers which mostly contribute to principle component 1 (PC1) and
 225 PC2. (D) the PCA based on sequence composition without those located on V- and J-segments
 226 for samples. Paired Wilcox-ranked test was used in B.

227 **Top1000 clones within NT are shorter than those in other subsets and the shortness is little**
 228 **affected by V-/J-gene usage**

229 The entire repertoire of NT was reported to be longer than the repertoire of memory²⁶.
 230 However, we found that the top clones in NT were shorter than clones in other subsets in all
 231 datasets (Figure 4A; Supplemental Figure 6). Via calculation of the Pearson correlations, the
 232 length distribution of NT is different to the length distribution of CMT and Tscm. It suggests
 233 that the length distribution of TCRB repertoire of NT is highly skewed in these less-
 234 differentiated memory cells (Figure 4B). Since the naïve cells are sorted without antibody
 235 against CD27 in dataset2, the length distribution of NT can be affected by the contamination
 236 from cell sorting. We examined the length distribution in dataset 1.

237 To identify whether the gene usage affects the CDR3 length distribution, we calculated the
 238 mean length of clones for each gene. The mean length is different among clones by varied V-
 239 and J-genes (Figure 4C and D), however, clones using all of genes are shorter in NT than -
 240 clones in CMT and Tscm. Therefore, for top clones, the clones of NT are shorter than the clones
 241 of other subsets, and the gene usage contributes less to the distinct length distributions among
 242 subsets.



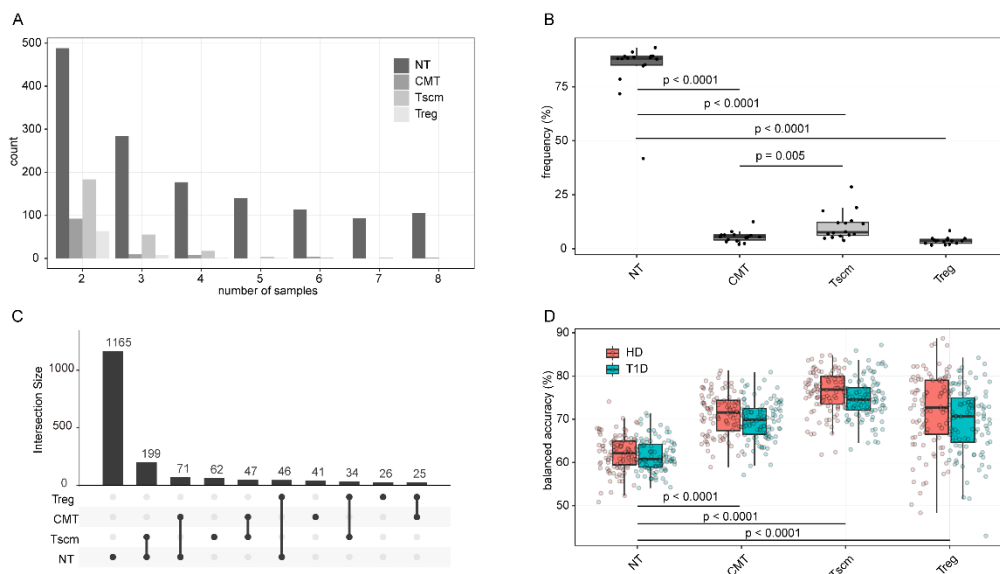
243

244 **Figure 4. The correlations of length distribution among subsets and the influence of gene**
 245 **usage on length distribution.** (A) The length distribution of CDR3s of each subset. (B) The
 246 correlation of length distribution between subsets. (C) The mean length of clones using each V-
 247 gene. (D) The mean length of clones using each J-gene. Paired Wilcox-ranked test was used in
 248 A and B.

249 **Public clones enrich in NT and have more differences with private clones in memory**
 250 **subsets than in NT**

251 Public clones that are shared by individuals were shown to be different from private clones
 252 in sequence composition²⁵. Our analyses showed that public clones were shorter than private
 253 ones within top1000 clones (Supplemental Figure 7A). It suggests that public clones may affect
 254 the features of top clones. We referred clones found in no less than two individuals as public
 255 clones. We found more public clones in NT than in other subsets: 1,400 in NT, 262 in Tscm,
 256 113 in CMT, and 72 in Treg from HD; 1544 in NT, 146 in Tscm, 128 in CMT and 92 in Treg

257 from T1D (Figure 5A; Supplemental Figure 7B). Via calculating abundance, we found that
 258 public clones were composed of ~70% of top1000 in NT, and ~5% in other subsets. It suggests
 259 that the public clones in NT have a larger effect on the repertoire of top clones than the public
 260 clones in others (Figure 5B). Most of public clones in NT were a little presented in other subsets
 261 (Figure 5C), and about 50% public clones in each subset could be found in NT. This result
 262 indicates that public clones in NT are less maintained than the top clones in other subsets. To
 263 detect the differences between public clones and private clones within each subset, support
 264 vector machine (SVM) was used. To avoid the influence of differential sample sizes, we
 265 randomly down-sampled 400 public clones and private clones for each subset, respectively. The
 266 prediction was repeated for 100 times. The prediction accuracy (BACC) in NT was found to be
 267 lower than BACC in other subsets (Figure 5D). It suggests that the differences between private
 268 and public clones in differentiated subsets are larger than the differences in naïve. Further
 269 analyses showed that the gene usage of public clones is similar to the gene usage of all top1000
 270 clones (Supplemental Figure 7C). It suggests that gene usage is not skewed in public clones.



271

272 **Figure 5. The differential usage of public clones among subsets and the discrimination of**
 273 **sequence composition between public clones and private clones. (A)** The number of public
 274 clones shared by from two to eight HDs. **(B)** The percentage of unique public clones within top
 275 clones in each subset. **(C)** The overlap of public clones from HDs among subsets. **(D)** The
 276 prediction accuracy (BACC) of SVM ($k = 3$) for public clones and private clones based on

277 sequence composition within each subset. Paired Wilcoxon-ranked test was used in **B**, Wilcoxon-
278 ranked test was used in **D**.

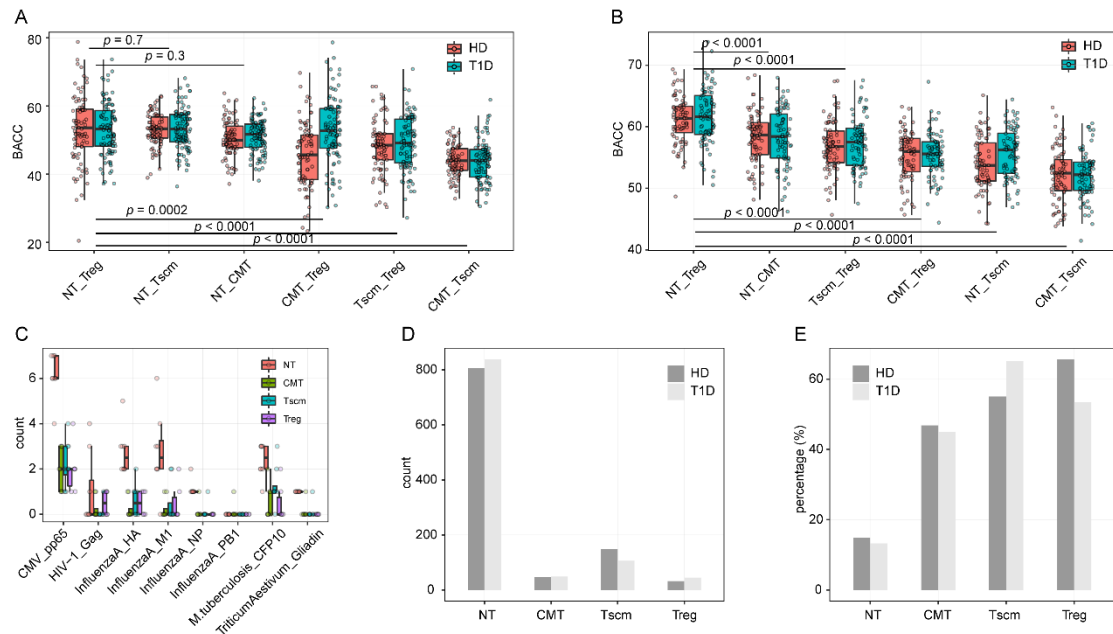
279 **The sequence compositions of public clones and private clones are both skewed along** 280 **differentiation**

281 To identify that public clones or private clones account for the increased difference in
282 memory and Treg, we performed SVM to discriminate public clones as well as private clones
283 from different subsets separately²⁵. For public clones, the BACC was from 50% to 60%; NT
284 was able to be discriminated from Treg, CMT and Tscm with ~ 55% BACC; whereas CMT was
285 incapable of separating from Tscm with ~ 50% BACC (Figure 6A). For private clones, the
286 BACC was from 50% to 70%. We were able to achieve a high prediction accuracy to
287 discriminate private clones of TN from private clones of Treg, but we failed to separate private
288 clones from CMT and Tscm (Figure 6B). When we increased the sample size of private clones
289 from 400 to 2500 for training SVM model, we found that the varied BACCs to discriminate
290 private clones from different subsets were still existed (Supplemental Figure 8A). These results
291 suggest that the sequence compositions of public clones and private clones are both skewed in
292 differentiation.

293 **The reduced number of public clones narrows the antigen spectrum recognized by** 294 **differentiated subsets**

295 To unveil the functions of clones among subsets, we annotated clones by VDJDdb²⁷. 1885
296 clones of CD4⁺ T cells targeting eight epitopes in total are recorded by this database.
297 Comparing with CMT, Tscm and Treg, NT has more clones recognizing antigens (HA, H1 and
298 NP) from influenza, pp65 from cytomegalovirus (CMV), CFP10 from *M. tuberculosis* and
299 gliadin from *Triticum Aestivum* (Figure 6C; Supplementary figure 8B). To estimate the
300 spectrum of antigens targeted by the top clones, we used GLIPH2²³ to predict the clusters
301 recognizing diverse antigens for each subset. With a stringency filter (see Methods), we found
302 out 806 clusters in NT from HC and 836 clusters in NT from T1D respectively; while less than
303 200 clusters in whole of Tscm, CMT and Treg (Figure 6D). When public clones were removed
304 from top clones, only 14.88% clusters remained in NT of HC and 13.25% clusters remained in
305 T1D, whereas over 45% clusters remained in other subsets (Figure 6E). It suggests that the

306 public clones enlarge the antigen spectrum recognized by top clones in NT. In conclusion, NT
307 recognizes a broader antigen profile contributed by public clones.



308

309 **Figure 6. public clones maintain stable composition across subsets and provide the large**
310 **part of the ability to recognize antigens. (A) SVM for discriminating the composition of**
311 **public clones among subsets. (B) SVM for discriminating the composition of private clones**
312 **among subsets. (C) The enrichment of antigen-related clones in each subset. (D) the number of**
313 **clusters targeting diverse antigens predicted by GLIPH2 within each subset. (E) The percentage**
314 **of clusters of private clones predicted by GLIPH2 within each subset. Wilcox-ranked test was**
315 **used in A and B.**

316

317 **Discussion**

318 It is essential for CD4+ T cells to recognize antigens with TCR, which is primarily
319 achieved by the CDR3 region. CD4+ T cells can acquire new functions via differentiation;
320 however, it is unclear how differentiation affects their TCR repertoire. We detected the
321 relationships among the TCRB repertoire of top1000 clones of naïve, memory and Treg subsets
322 (including NT, ET, Tcm, Tem, Tscm ET and Treg) by estimating the repertoire structure, the
323 germline gene usage, the sequence composition (K-mer) and public CDR3 clone usage.

324 We derive that the TRBV repertoire features of memory subsets are tightly regulated in
325 differentiation. We observed that 23 of 72 genes increased or decreased in an order of NT, Tscm,
326 CMT and EMT. It indicates that a mechanism exists to regulate the variations across subsets.
327 Furthermore, since Tscm is the least differentiated cell whereas EMT is the highest one among
328 the tree memory subsets²⁸, it indicates that the differentiation level is along with the mechanism.
329 CMT is formally considered as the primary memory subset which ET prefer to differentiate to,
330 and then part of CMT differentiates to EMT. In the past decade, Tscm has been found to mix
331 phenotypes of naïve and memory. Tscm is able to self-renew and replenish more differentiated
332 subsets of memory T cells, and therefore acts as the key intermediary of the generation of
333 memory^{29,30}. In together, differentiation levels of memory subsets reflect their differentiation
334 order. However, memory cells can be directly generated from naïve cells by asymmetric cell
335 division^{31,32}. It indicates that the differential order should not be the only factor skewing TCRB
336 repertoire. Shown by X. L. Hou *et al*³³, the early events in thymic T cell development are
337 different for CD4+ naïve and memory cells. It suggests that genetic factors affect TCRB
338 repertoire of CD4+ T cells in differentiation. In addition, events at the very early lifetime can
339 be involved in manipulating TCRB repertoire. Observations in newborns show that memory in
340 human develops at the very early period of lifetime³⁴. Newborns less encounter pathogenic
341 antigens. It implies that, for newborns, food, self-antigens and even cytokine driven clones
342 compose the large part of TCRB repertoire of memory. Since the highly frequent clones in NT
343 are self-antigen related, the features of frequent clones in NT will be delivered to memory at
344 this period³⁵. Furthermore, shown by Graeme *et al*, a half part of memory is maintained by self-
345 renewal influx during the lifetime³⁶. S. Jaafoura *et al* showed that less-differentiated memory
346 subset is more stable during pathogen infection²⁸. It suggests that Tscm and CMT rather than

347 EMT maintain the features of TCRB repertoire inherited from NT at the early lifetime. In
348 conclusion, it is reasonable to drive that events at the early lifetime, genetic factors and
349 differentiation order regulate the TCRB repertoire of CD4+ T subsets with differentiated levels.

350 Public clones are key components that affect the features of TCRB repertoire in
351 differentiation. First, we found that public clone usage rather than gene usage shortens the
352 length distribution of top clones within NT. Second, the sequence composition of public clones
353 which is skewed in differentiated subsets contribute to the variations of TCRB repertoire in
354 differentiation. Third, decreased public clones induce a reduction in antigen spectrum
355 recognized by memory and Treg subsets. These results suggest that the skewed public clone
356 usage highly affect top clones in differentiated subsets. Furthermore, we showed that factors
357 affecting the generation of public clones in memory and Treg are different to that in NT. The
358 generation of public clones were largely attributed to genetic factors and thymic positive
359 selection in the previous study³⁷. In our study, public clones from NT are less maintained in
360 differentiated subsets, and SVM analyses indicate that sequence composition in memory
361 subsets is different from that of NT. These results suggest that other factors, such as antigen,
362 trim the sequence composition of public clones. However, there are uncertainties about public
363 clones in NT. First, it is unclear to the cause that the promiscuous ability of public clones does
364 not induce the clonality of the public clones in effector and memory subsets. Second, it is
365 unclear that the physiological function of public clones in NT. As a hypothesis, the promiscuous
366 public clones in NT, sensitive to many antigens, are important to initiate the primary immune
367 response, and this function is forbidden in differentiated subsets. Further studies with single-
368 cell RNA-seq and paired sequencing of TCR alpha and beta chains will unveil the biological
369 functions of the public clones in NT.

370 In addition to the sample size, we identified that the T cell subset also affect the SVM
371 prediction accuracy. Shown by Victor Grief *et al*²⁵, public clones and private clones can be
372 discriminated by SVM using k-mer distribution. SVM can obtain a better prediction accuracy
373 (BACC) by using a larger sample size. We extended the detection in NT, Tscm, CMT and Treg.
374 By normalizing the sample size for training, we found that the prediction accuracy for public
375 clones and private clones in Tscm, CMT and Treg is higher than the prediction accuracy in NT.

376 It suggests that the difference between public clones and private clones is enlarged in the
377 differentiated subsets. When we performed SVM on public clones and private clones among
378 subsets respectively, the sequence compositions of public clones and private clones were
379 skewed in differentiation.

380 A small part of peripheral Treg differentiated from conventional Treg. Shown by Golding
381 *A. et al*, the repertoire of Foxp3+ and Foxp3- cells did not overlap³⁸. Although peripheral Tregs
382 are differentiated from conventional T cells³⁹ and can introduce the features of NT into Treg,
383 the TCR repertoire of effector and memory subsets is similar to NT than to Treg. This
384 phenomenon suggests that the influx from naïve just composed a minor part of Treg in blood,
385 and comparing to Treg, the features of naïve are maintained in effector and memory subsets in
386 the differentiation. Our study includes samples of three healthy states (healthy, RA and T1D
387 individuals), and therefore highlights that our findings are consistent in healthy conditions and
388 datasets.

389 **Acknowledgments**

390 The authors would like to thank Z.Q. Ding from Singapore institute of technology for edition
391 of the manuscript; Y. Liu from BGI-Shenzhen and W. Zhang from department of computer
392 science, City University of Hong Kong for comments on this manuscript.

393 **Conflicts of interest**

394 The authors declare no conflicts of interest.

395 **Author contributions**

396 SY. W. designed the study, performed the analyses, and drafted the manuscript; Y. L. supervised
397 the study, and revised the manuscript.

398 **Financial support**

399 This work was supported by BGI-Shenzhen, China National GeneBank (CNGB), Science,
400 Technology and Innovation Commission of Shenzhen Municipality under grant No.
401 JCYJ20170817145845968, and Shenzhen Key Laboratory of Single-Cell Omics (NO.
402 ZDSYS20190902093613831).

403

404 **References**

- 405 1 Wang, H. *et al.* TNF-alpha/IFN-gamma profile of HBV-specific CD4 T cells is associated
406 with liver damage and viral clearance in chronic HBV infection. *J Hepatol* **72**, 45-56,
407 doi:10.1016/j.jhep.2019.08.024 (2020).
- 408 2 Gray, J. I., Westerhof, L. M. & MacLeod, M. K. L. The roles of resident, central and effector
409 memory CD4 T-cells in protective immunity following infection or vaccination.
410 *Immunology*, doi:10.1111/imm.12929 (2018).
- 411 3 Zander, R. *et al.* CD4(+) T Cell Help Is Required for the Formation of a Cytolytic CD8(+) T
412 Cell Subset that Protects against Chronic Infection and Cancer. *Immunity* **51**, 1028-1042
413 e1024, doi:10.1016/j.immuni.2019.10.009 (2019).
- 414 4 de Greef, P. C. *et al.* The naive T-cell receptor repertoire has an extremely broad
415 distribution of clone sizes. *Elife* **9**, doi:10.7554/eLife.49900 (2020).
- 416 5 Mueller, S. N., Gebhardt, T., Carbone, F. R. & Heath, W. R. Memory T cell subsets, migration
417 patterns, and tissue residence. *Annu Rev Immunol* **31**, 137-161, doi:10.1146/annurev-
418 immunol-032712-095954 (2013).
- 419 6 Kaech, S. M., Wherry, E. J. & Ahmed, R. Effector and memory T-cell differentiation:
420 implications for vaccine development. *Nat Rev Immunol* **2**, 251-262, doi:10.1038/nri778
421 (2002).
- 422 7 Jaigirdar, S. A. & MacLeod, M. K. Development and Function of Protective and Pathologic
423 Memory CD4 T Cells. *Front Immunol* **6**, 456, doi:10.3389/fimmu.2015.00456 (2015).
- 424 8 Qi, Q. *et al.* Diversity and clonal selection in the human T-cell repertoire. *Proc Natl Acad*
425 *Sci U S A* **111**, 13139-13144, doi:10.1073/pnas.1409155111 (2014).
- 426 9 Dowling, M. R. *et al.* Regulatory T Cells Suppress Effector T Cell Proliferation by Limiting
427 Division Destiny. *Front Immunol* **9**, 2461, doi:10.3389/fimmu.2018.02461 (2018).
- 428 10 Li, M. O. & Rudensky, A. Y. T cell receptor signalling in the control of regulatory T cell
429 differentiation and function. *Nat Rev Immunol* **16**, 220-233, doi:10.1038/nri.2016.26
430 (2016).
- 431 11 Snook, J. P., Kim, C. & Williams, M. A. TCR signal strength controls the differentiation of
432 CD4(+) effector and memory T cells. *Science immunology* **3**,
433 doi:10.1126/sciimmunol.aas9103 (2018).
- 434 12 Gutierrez, L., Beckford, J. & Alachkar, H. Deciphering the TCR Repertoire to Solve the
435 COVID-19 Mystery. *Trends Pharmacol Sci* **41**, 518-530, doi:10.1016/j.tips.2020.06.001
436 (2020).
- 437 13 Grifoni, A. *et al.* Targets of T Cell Responses to SARS-CoV-2 Coronavirus in Humans with
438 COVID-19 Disease and Unexposed Individuals. *Cell* **181**, 1489-1501 e1415,
439 doi:10.1016/j.cell.2020.05.015 (2020).
- 440 14 Schoettler, N., Hrusch, C. L., Blaine, K. M., Sperling, A. I. & Ober, C. Transcriptional
441 programming and T cell receptor repertoires distinguish human lung and lymph node
442 memory T cells. *Commun Biol* **2**, 411, doi:10.1038/s42003-019-0657-2 (2019).
- 443 15 Fohse, L. *et al.* High TCR diversity ensures optimal function and homeostasis of Foxp3+
444 regulatory T cells. *Eur J Immunol* **41**, 3101-3113, doi:10.1002/eji.201141986 (2011).
- 445 16 Adeegbe, D., Matsutani, T., Yang, J., Altman, N. H. & Malek, T. R. CD4(+) CD25(+) Foxp3(+)
446 T regulatory cells with limited TCR diversity in control of autoimmunity. *J Immunol* **184**,
447 56-66, doi:10.4049/jimmunol.0902379 (2010).

- 448 17 Gomez-Tourino, I., Kamra, Y., Baptista, R., Lorenc, A. & Peakman, M. T cell receptor beta-
449 chains display abnormal shortening and repertoire sharing in type 1 diabetes. *Nat*
450 *Commun* **8**, 1792, doi:10.1038/s41467-017-01925-2 (2017).
- 451 18 Jiang, X. *et al.* Comprehensive TCR repertoire analysis of CD4(+) T-cell subsets in
452 rheumatoid arthritis. *J Autoimmun*, 102432, doi:10.1016/j.jaut.2020.102432 (2020).
- 453 19 Zhang, W. *et al.* IMonitor: A Robust Pipeline for TCR and BCR Repertoire Analysis. *Genetics*
454 **201**, 459-472, doi:10.1534/genetics.115.176735 (2015).
- 455 20 Koch, H., Starenki, D., Cooper, S. J., Myers, R. M. & Li, Q. powerTCR: A model-based
456 approach to comparative analysis of the clone size distribution of the T cell receptor
457 repertoire. *PLoS Comput Biol* **14**, e1006571, doi:10.1371/journal.pcbi.1006571 (2018).
- 458 21 Drost, H.-G. Philentropy: Information Theory and Distance Quantification with R. *Journal*
459 *of Open Source Software* **3**, doi:10.21105/joss.00765 (2018).
- 460 22 Palme, J., Hochreiter, S. & Bodenhofer, U. KeBABS: an R package for kernel-based analysis
461 of biological sequences. *Bioinformatics* **31**, 2574-2576,
462 doi:10.1093/bioinformatics/btv176 (2015).
- 463 23 Huang, H., Wang, C., Rubelt, F., Scriba, T. J. & Davis, M. M. Analyzing the Mycobacterium
464 tuberculosis immune response by T-cell receptor clustering with GLIPH2 and genome-
465 wide antigen screening. *Nat Biotechnol*, doi:10.1038/s41587-020-0505-4 (2020).
- 466 24 Miho, E., Roskar, R., Greiff, V. & Reddy, S. T. Large-scale network analysis reveals the
467 sequence space architecture of antibody repertoires. *Nat Commun* **10**, 1321,
468 doi:10.1038/s41467-019-09278-8 (2019).
- 469 25 Greiff, V. *et al.* Learning the High-Dimensional Immunogenomic Features That Predict
470 Public and Private Antibody Repertoires. *J Immunol* **199**, 2985-2997,
471 doi:10.4049/jimmunol.1700594 (2017).
- 472 26 Hou, X. *et al.* Shorter TCR beta-Chains Are Highly Enriched During Thymic Selection and
473 Antigen-Driven Selection. *Front Immunol* **10**, 299, doi:10.3389/fimmu.2019.00299 (2019).
- 474 27 Bagaev, D. V. *et al.* VDJdb in 2019: database extension, new analysis infrastructure and a
475 T-cell receptor motif compendium. *Nucleic Acids Res* **48**, D1057-D1062,
476 doi:10.1093/nar/gkz874 (2020).
- 477 28 Jaafoura, S. *et al.* Progressive contraction of the latent HIV reservoir around a core of less-
478 differentiated CD4(+) memory T Cells. *Nat Commun* **5**, 5407, doi:10.1038/ncomms6407
479 (2014).
- 480 29 Ahmed, R. *et al.* Human Stem Cell-like Memory T Cells Are Maintained in a State of
481 Dynamic Flux. *Cell Rep* **17**, 2811-2818, doi:10.1016/j.celrep.2016.11.037 (2016).
- 482 30 Gattinoni, L. *et al.* A human memory T cell subset with stem cell-like properties. *Nat Med*
483 **17**, 1290-1297, doi:10.1038/nm.2446 (2011).
- 484 31 Chang, J. T. *et al.* Asymmetric T lymphocyte division in the initiation of adaptive immune
485 responses. *Science (New York, N.Y.)* **315**, 1687-1691, doi:10.1126/science.1139393 (2007).
- 486 32 Borsa, M. *et al.* Modulation of asymmetric cell division as a mechanism to boost CD8(+) T
487 cell memory. *Science immunology* **4**, doi:10.1126/sciimmunol.aav1730 (2019).
- 488 33 Hou, X. *et al.* Preselection TCR repertoire predicts CD4(+) and CD8(+) T-cell differentiation
489 state. *Immunology*, doi:10.1111/imm.13256 (2020).
- 490 34 Qazi, K. R. *et al.* Extremely Preterm Infants Have Significant Alterations in Their
491 Conventional T Cell Compartment during the First Weeks of Life. *J Immunol* **204**, 68-77,

- 492 doi:10.4049/jimmunol.1900941 (2020).
- 493 35 Madi, A. *et al.* T-cell receptor repertoires share a restricted set of public and abundant
494 CDR3 sequences that are associated with self-related immunity. *Genome research* **24**,
495 1603-1612, doi:10.1101/gr.170753.113 (2014).
- 496 36 Gossel, G., Hogan, T., Cownden, D., Seddon, B. & Yates, A. J. Memory CD4 T cell subsets
497 are kinetically heterogeneous and replenished from naive T cells at high levels. *Elife* **6**,
498 doi:10.7554/eLife.23013 (2017).
- 499 37 Khosravi-Maharlooei, M. *et al.* Crossreactive public TCR sequences undergo positive
500 selection in the human thymic repertoire. *J Clin Invest* **129**, 2446-2462,
501 doi:10.1172/JCI124358 (2019).
- 502 38 Golding, A., Darko, S., Wylie, W. H., Douek, D. C. & Shevach, E. M. Deep sequencing of the
503 TCR-beta repertoire of human forkhead box protein 3 (FoxP3)(+) and FoxP3(-) T cells
504 suggests that they are completely distinct and non-overlapping. *Clin Exp Immunol* **188**,
505 12-21, doi:10.1111/cei.12904 (2017).
- 506 39 Kraj, P. & Ignatowicz, L. The mechanisms shaping the repertoire of CD4(+) Foxp3(+)
507 regulatory T cells. *Immunology* **153**, 290-296, doi:10.1111/imm.12859 (2018).

508