# Human thymic T cell repertoire is imprinted with strong convergence to shared sequences

Heikkilä, Nelli

1    **Human thymic T cell repertoire is imprinted with strong convergence to shared sequences**

2

3    Heikkilä Nelli[a], Vanhanen Reetta[b], Yohannes Dawit A.[c], Kleino Iivari[d], Mattila Ilkka P.[e], Saramäki

4    Jari[f], Arstila T. Petteri[g]

5

6    a)  Research Programs Unit, Translational Immunology and Medicum, Department of

7        Bacteriology and Immunology, University of Helsinki. Haartmaninkatu 3, 00290 Helsinki,

8        Finland. nelli.heikkila@helsinki.fi

9    b)  Research Programs Unit, Translational Immunology and Medicum, Department of

10       Bacteriology and Immunology, University of Helsinki. Haartmaninkatu 3, 00290 Helsinki,

11       Finland. reetta.vanhanen@helsinki.fi

12   c)  Research Programs Unit, Translational Immunology and Medicum, Department of Medical

13       and Clinical Genetics, University of Helsinki. Haartmaninkatu 8, 00290 Helsinki, Finland.

14       dawit.yohannes@helsinki.fi

15   d)  Research Programs Unit, Translational Immunology, University of Helsinki.

16       Haartmaninkatu 3, 00290 Helsinki, Finland. iivari.kleino@helsinki.fi

17   e)  Department of Pediatric Cardiac and Transplantation Surgery, Hospital for Children and

18       Adolescents, Helsinki University Central Hospital. Stenbäckinkatu 9, 00290 Helsinki,

19       Finland. ilkka.mattila@hus.fi

20   f)  Department of Computer Science, Aalto University. Konemiehentie 2, 02150 Espoo,

21       Finland. jari.saramaki@aalto.fi

22   g)  Research Programs Unit, Translational Immunology and Medicum, Department of

23       Bacteriology and Immunology, University of Helsinki. Haartmaninkatu 3, 00290 Helsinki,

24       Finland. petteri.arstila@helsinki.fi

25

26

27     Corresponding author:

28     Nelli Heikkilä

29     Research Programs Unit, Translational Immunology, PB21, 00014 University of Helsinki, Finland

30     e-mail: nelli.heikkila@helsinki.fi

31     telephone: +358 29 41 26373

32     fax: +358 29 41 26382

33

**Abstract**

A highly diverse repertoire of T cell antigen receptors (TCR) is created in the thymus by recombination of gene segments and the insertion or deletion of nucleotides at the junctions. Using next-generation TCR sequencing we define here the features of recombination and selection in the human TCRα and TCRβ locus, and show that a strikingly high proportion of the repertoire is shared by unrelated individuals. The thymic TCRα nucleotide repertoire was more diverse than TCRβ, with $4.1 \times 10^6$ vs. $0.81 \times 10^6$ unique clonotypes, and contained nonproductive clonotypes at a higher frequency (69.2% vs. 21.2%). The convergence of distinct nucleotide clonotypes to the same amino acid sequences was higher in TCRα than in TCRβ repertoire (1.45 vs. 1.06 nucleotide sequences per amino acid sequence in thymus). The gene segment usage was biased, and generally all individuals favored the same genes in both TCRα and TCRβ loci. Despite the high diversity, a large fraction of the repertoire was found in more than one donor. The shared fraction was bigger in TCRα than TCRβ repertoire, and more common in in-frame sequences than in nonproductive sequences. Thus, both biases in rearrangement and thymic selection are likely to contribute to the generation of shared repertoire in humans.

**Keywords:** T cell antigen receptor, TCR repertoire, TCR recombination, thymus, next-generation sequencing

Abbreviations: T cell antigen receptor (TCR), V (variable), D (diversity), J (joining), CDR3 (complementarity-determining region 3), Pgen (generative probability)

## 1. Introduction

T cell antigen receptor (TCR) is a heterodimeric surface protein, consisting in most cells of α and β chains, while a small minority of cells use γ and δ chains. Both chains are encoded by genes assembled from incomplete segments via somatic recombination during development in the thymus. The TCRβ locus contains 47 variable (V), 2 diversity (D) and 13 joining (J) gene segments whereas the TCRα locus contains 42 V and 61 J segments but lacks the D segment. Further diversity is achieved at the gene segment junctions where a number of nucleotides may be removed and palindromic P-nucleotides and non-templated N-nucleotides inserted. Thus, most of the variability in the TCR concentrates in the junctional regions, called complementary determining region 3 (CDR3), which also form the main site of antigen recognition (Davis and Bjorkman, 1988).

The recombination process is capable of creating a high level of diversity. Direct sequencing of TCRβ repertoire has measured a lower limit of $1\text{-}3\text{x}10^6$ clonotypes, whereas a mathematical estimator suggested a total repertoire of about $100\text{x}10^6$ unique clonotypes (Qi et al., 2014). We have recently measured the lower limit of thymic TCR diversity in pediatric samples to be $10.3\text{x}10^6$ for TCRβ and $3.7\text{x}10^6$ for TCRα clonotypes, and statistical modelling suggested the total repertoire to consist of $40\text{-}70\text{x}10^6$ and $60\text{-}100\text{x}10^6$ clonotypes for TCRβ and TCRα respectively (Vanhanen et al., 2016). The pairing of TCRα to TCRβ has been studied but little. A sequencing of a limited TCR subset showed that on the average each TCRβ chain can bind to at least 24 different TCRα chains (Arstila et al., 1999), while a recent large scale single-cell analysis suggested that the pairing is more limited than would be compatible with a fully stochastic process (Grigaityte et al., 2017). The full TCRαβ repertoire thus consists of at least tens of millions of different receptors.

To date the human thymic TCR repertoire has been studied very little. The TCRβ locus is rearranged first and is subject to relatively stringent allelic exclusion. However, TCRβ locus may

83    also be rearranged in cells destined to the γδ T cell lineage, which may account for a part of the

84    nonfunctional TCRβ repertoire. Since most recombination events will result in an out-of-frame

85    sequence, the functionality of the rearranged TCRβ chain is ensured by pairing with a surrogate

86    TCRα chain, the preTα. The cells capable of signaling through pre-TCR then proliferate before

87    recombination begins in the TCRα locus (von Boehmer et al., 1998). Unlike the TCRβ locus, in

88    TCRα recombination both alleles are rearranged simultaneously, until a functional TCRαβ is

89    expressed, stopping the recombination. Thus, in a large proportion of cells both TCRα loci are

90    rearranged, although only one is likely to produce a functional protein chain (Casanova et al.,

91    1991). The newly generated TCRαβ+ cells are then subjected to positive and negative selections,

92    which remove cells incapable of interacting with HLA molecules or displaying too strong affinity to

93    self-antigens (Klein et al., 2014). Overall, only an estimated 3-5% of the developing thymocytes

94    survive the selection process to form the mature peripheral repertoire (Egerton et al., 1990; Yates,

95    2014).

96

97    In the present study, we characterize the composition of the thymic TCRβ and TCRα repertoire,

98    identifying differences in the two chains related to their biology. Our data also show a strikingly

99    strong convergence to shared repertoire in unrelated individuals.

## 2. Materials and Methods

The study was approved by the Pediatric Ethical Committee of the Helsinki University Hospital and parents gave a written informed consent. Thymus samples were obtained from eight immunologically healthy children undergoing a corrective operation for congenital cardiac defects (donors A-D and donors 1-4). Additionally, a peripheral blood sample was drawn from donors 1-4 during the operation. The donors were 7–244 days old and 2/8 were female (Table 1). Two of the subjects (donors A and B) were monozygotic twins. The impact of genetics on the repertoire has been analyzed in detail elsewhere (Heikkila et al., 2020). All thymus samples appeared macroscopically normal. Thymocyte populations from donors B-D were analyzed with flow cytometry for expression of CD4, CD8, TCRαβ, TCRγδ, CD3 and CD69.

From each subject, an aliquot of 10–30 million thymocytes and from donors 1-4 an aliquot of 0.5 mL peripheral blood was used for sequencing both TCRAD and TCRB repertoire. Thymocytes were extracted mechanically from the tissue. To remove red blood cells blood samples were treated with Gibco™ACK Lysing Buffer (Thermo Fisher Scientific, Massachusetts, USA), according to the manufacturer's orders. Genomic DNA was extracted from pelleted cell samples with QIAsymphony (Qiagen, Germany) according to the manufacturer's instructions. The TCRα and TCRβ CDR3 regions were amplified and sequenced from a standardized quantity of quality-controlled DNA using ImmunoSEQ assay (Adaptive Biotechnologies, Seattle, USA). In summary, the sequencing assay consists of a multiplex PCR system to amplify the rearranged CDR3 regions from the DNA samples at a length that is sufficient to subsequently identify the VDJ and VJ regions spanning each unique CDR3α and CDR3β regions, respectively. Amplicon sequencing was performed with Illumina platform. TCRα and TCRβ gene segment definitions were obtained from IMGT database (www.imgt.org). Primer bias was corrected as previously described (Vanhanen et al., 2016) and the resulting data filtered and clustered using both the relative frequency ratio between similar clones

125   and a modified nearest-neighbor algorithm to remove both PCR and sequencing errors. All

126   sequences are available at immuneACCESS database provided by Adaptive Biotechnologies

127   (clients.adaptivebiotech.com/immuneaccess).

128

129   The TCR sequence analysis was performed using the immunoSEQ ANALYZER 3.0 (Adaptive

130   Biotechnologies, Seattle, USA), VDJTools software (Shugay et al., 2015) and in-house scripts for

131   computing languages R (www.r-project.org) and python 2.7 (www.python.org). The in-house

132   scripts generated for this study are published in Supplements 1&2. The similarity of two sets of

133   unique or total sequences was assessed calculating the Jaccard index, which is defined as the size of

134   the intersection of two data sets (A and B) divided by the size of their union: $J(A, B) = \frac{|A \cup B|}{|A \cap B|}$. The

135   abundance based Jaccard index was defined as $J_{abund} = UV/(U+V-UV)$, where U is the total relative

136   abundance of shared sequences in sample A and V the total relative abundance of shared sequences

137   in sample B (Chao et al., 2006).  The CDR3 nucleotide sequences were extracted separately for in-

138   frame and nonproductive sequences and subsequently the generative probabilities were calculated

139   using the OLGA software (Sethna et al., 2019).

140

141   **3.    Results**

142       **3.1.    TCRα and TCRβ repertoires differ in diversity and productivity**

143   Thymus samples were collected from eight pediatric patients (donors A-D and donors 1-4), two of

144   whom were monozygotic twins (donors A and B; Table 1). Flow cytometric analysis was performed

145   for donors B-D and showed a normal distribution of CD4 and CD8 double-negative (DN), double-

146   positive (DP), and single-positive (SP) thymocytes as well as normal pattern of TCRαβ and TCRγδ

147   expression (Figure 1A). Postselection thymocytes were defined as DPCD3highCD69+, CD4SP or

148   CD8SP (Swat et al., 1993; Yamashita et al., 1993). On the average, 23.1±3.7% of total thymocytes

149   represented postselection and 76.9±3.7% preselection population (Figure 1B).

150

151   Sequencing of thymic TCRs yielded $1.2 \times 10^5$-$1.6 \times 10^6$ (mean 810 000) unique TCRβ clonotypes of

152   which 78.8±2.7 % were in-frame, 19.3±2.4 % were out-of-frame and 2.1±0.5 % contained a

153   premature stop-codon (Fig. 1C). Consistent with our previous estimation on thymic TCR diversity,

154   the TCRα diversity was higher than TCRβ diversity, with $1.3$-$7.6 \times 10^6$ (mean $4.1 \times 10^6$) unique

155   clonotypes per sample (Vanhanen et al., 2016). However, the productivity in TCRα was much

156   lower, as only 30.8±0.8 % of the unique clonotypes were in-frame. Of the unique TCRα clonotypes

157   66.0±0.6 % were out-of-frame and 7.0±4.5 % contained a premature stop-codon (Fig. 1C). As the

158   sequencing assay is based on genomic DNA, it also provides a quantitative estimate of the number

159   of total genomes with rearranged TCR segments in the sample.

160

161   A small blood sample from donors 1-4 was sequenced simultaneously with the thymus samples,

162   producing an average of 84 000 unique TCRβ clonotypes and 150 000 unique TCRα clonotypes. In

163   the TCRβ repertoire, the fractions of in-frame and nonproductive clonotypes remained essentially

164   similar to that in the thymus (Figure 1D). In the TCRα repertoire, the fraction of in-frame

clonotypes was higher in the blood samples than in the thymus (38.5±1.4% vs. 30.8±0.8%; Figure 1D).

To estimate the convergence of distinct nucleotide clonotypes to identical amino acid chains we calculated the nucleotide-to-amino acid-ratio for each sample. The majority of amino acid chains in the TCRβ repertoire were encoded by a single nucleotide clonotype, the nucleotide-to-amino acid-ratio being for unique in-frame clonotypes 1.06±0.03 in the thymus and 1.05±0.02 in the periphery. In the TCRα repertoire the number of unique nucleotide clonotypes converging to the same amino acid chain was higher than in the TCRβ repertoire, particularly in the thymus (ratio 1.45±0.13) but also to some degree in the periphery (ratio 1.18±0.01).

### 3.2. The V and J segment usage is biased before thymic selections

Previous studies of peripheral repertoire have shown a biased usage of V and J genes in healthy subjects. Similarly, in thymus the use of V gene elements was uneven, and the same segments were favored in each individual both in thymus and in blood (e.g. TRBV5-1, TRBV27-01 and TRAV21-1, TRAV29-1; Supplement 3). Similar findings were also obtained for J gene usage (Supplement 4). The biased V and J gene usage pattern was largely observed both in the in-frame and nonproductive repertoire, indicating that it is due the recombination process rather than selection (Figure 2). Consistent with our previous study (Heikkila et al., 2020), the samples from the monozygotic twins A and B clustered together, indicating a genetic component in V and J gene usage. Interestingly, in the TCRα repertoire, the gene segment usage clustered thymic and peripheral blood samples mainly according to the sample type and not the identity of the donor (Figure 2A). In the TCRβ repertoire, in contrast, the gene segment usage clustered together blood and thymus samples taken from the same donor (Figure 2B).

190  Some of the gene segment bias might be caused by thymic generation of semi-invariant T cell

191  subsets, such as natural killer T cells (NKTs) or mucosal-associated T cells (MAITs). Human NKTs

192  prefer TRAV10/TRAJ18 combination and MAITs use invariable TRAV01-02/TRAJ33-01

193  combination. The β chain usage is less restricted, but with a preference of TRBV25 for NKTs and

194  TRBV6 and TRBV20 for MAITs. In our data none of the semi-invariant α chains was dominant

195  whereas some MAIT-associated TRBV6 genes were found at an elevated frequency. However,

196  these TRBV segments are also ubiquitously used by conventional variable T cells (Tickotsky et al.,

197  2017).

198

199  TCRδ gene segments are embedded within the TCRα locus and αβ and γδ lymphocytes may use

200  both TCRα and TCRδ gene segments in an overlapping manner (Verschuren et al., 1998). Since the

201  thymocytes we analyzed were not sorted, and the sequencing protocol included primers specific for

202  the entire TCRAD locus, we obtained a mixture of TCRα and TCRδ sequences. In the thymus, the

203  frequency of γδ TCR+ thymocytes, as measured by flow cytometry, was 0.80±0.20%. However, the

204  frequency of unique clonotypes using a combination of TRDV and TRDJ was 1.1±0.18%. In the

205  peripheral blood the frequency of TRDV-TRDJ combinations was slightly higher (1.7±0.96 % of

206  the unique clonotypes). We also identified relatively frequent combinations of TRDV to TRAJ

207  (2.4±0.20 % of the unique clonotypes), whereas sequences using a combination of TRAV and

208  TRDJ were rare both in thymus and in periphery (Table 2).

209

210  **3.3.  CDR3 region length reflects recombination and selection events**

211  The TCRβ chain comprises V, D, and J segments, whereas the TCRα chain lacks D segments and

212  thus contains only one junctional site. This difference was reflected in the higher number of non-

213  templated nucleotide insertions in the TCRβ than in the TCRα sequences with an average of 9.3 vs.

214  3.9 nucleotides in thymic and 7.0 vs. 3.7 in peripheral in-frame repertoires (Figure 3A). The

215 nonproductive sequences cannot be subject to TCR-mediated selection, and thus represent the non-

216 selected product of the recombination process. Consistent with the previously reported shortening

217 of CDR3 during thymic selection (Matsutani et al., 2011; Niemi et al., 2015; Yassai and Gorski,

218 2000), the mean CDR3 length was shorter in the in-frame rearrangements (41.6 base pairs (bp) for

219 TCRα and 45.7 bp for TCRβ) than in the nonproductive rearrangements (41.9 bp for TCRα and

220 46.3 bp for TCRβ) in the thymus (Figure 3B). In the peripheral in-frame repertoire, the CDR3

221 regions were still shorter (41.4 bp for TCRα and 43.4 bp for TCRβ).

222

223    **3.4.    Pgen distributions differ in TCRα and TCRβ repertoires**

224 In the process of V(D)J gene segment recombination and insertion of random nucleotides between

225 gene segments some sequences are generated more readily while the generation of others is more

226 unlikely. We used OLGA software to calculate the generative probabilities (Pgen) in the TCRα and

227 TCRβ repertoires (Sethna et al., 2019). For a large majority of nonproductive sequences we

228 obtained Pgen values 0, probably because the OLGA calculations are based on amino acid rather

229 than nucleotide sequences and CDR3 amino acid definition remains ambivalent for nonproductive

230 sequences. For the thymic in-frame sequences the Pgen was higher for TCRα (average Pgen 1.57e-

231 7) than for TCRβ (average Pgen 1.34e-9) repertoire, a finding likely due to the lower junctional

232 complexity in TCRα chains. The same was observed in the peripheral repertoires (average Pgen for

233 TCRα 1.56e-7 and for TCRβ 3.62e-9). In the TCRα repertoire, the thymic and peripheral Pgen

234 averages and distributions were largely identical, while for the TCRβ the thymic repertoires had

235 lower Pgen values than the peripheral repertoires (1.34e-9 vs. 3.62e-9; Figure 4).

236

237    **3.5.    Overlap of thymic clonotypes between two individuals**

238 Despite the high diversity of the junctional CDR3 sequences, a considerable overlap of peripheral

239 TCR repertoires between different individuals has been reported (Shugay et al., 2013). In our

240  thymic samples, a substantial fraction of TCR sequences were shared between two individuals, and

241  some of the TCRα and TCRβ clonotypes were shared even between multiple individuals (Figure

242  5A). This phenomenon was more marked in the TCRα than TCRβ repertoire. Indeed, in the samples

243  1-4, in which the sequencing depth was shallower, there were no TCRβ clonotypes shared by all

244  four donors.

245

246  To estimate the fraction of thymic repertoire shared by two individuals, we used the Jaccard index

247  (JI), calculated as the intersection of two samples divided by the union of the samples, with a

248  maximum index of 1 for fully overlapping repertoires. In the nonproductive TCRβ clonotypes the JI

249  was low (mean JI 6.3e-5), but increased clearly in the in-frame repertoire (mean JI 4.6e-4). When

250  unique amino acid CDR3 regions were analyzed, the shared fraction was higher still (mean JI

251  0.013; Figure 5B). In the TCRα repertoire, the shared fraction was generally higher than in the

252  TCRβ repertoire, and in the nonproductive clonotypes the mean JI was 0.029. A small but

253  consistent increase to mean JI of 0.032 was found in the in-frame repertoire. In the unique amino

254  acid CDR3 regions the shared fraction was again clearly higher (mean JI 0.10; Figure 5B). As

255  previously reported (Heikkila et al. 2020), comparison of the twins A and B produced slightly

256  higher JIs than the other pairs. In general, samples 1-4 were sequenced to a lesser depth than

257  samples A-D, affecting the observed number of shared clonotypes, and the JI values were

258  consequently smaller. However, the increasing trend in JI from nonproductive to in-frame and

259  amino acid sequences was clear in all samples.

260

261  The shared sequences contained fewer non-templated insertions than the individual private

262  repertoires. The average number of non-templated insertions in was 1.4 and 2.6 respectively for

263  shared in-frame and nonproductive TCRβ clonotypes. In TCRα the shared in-frame clonotypes

264  contained on the average 1.4 and nonproductive 1.6 insertions. Also, the Pgen was higher in the

265 shared repertoire compared to the full repertoires. In the in-frame repertoire the average Pgen for

266 unique shared in-frame TCRβ clonotypes was 4.71e-8 and in the full repertoire 1.34e-9. For in-

267 frame TCRα clonotypes the difference in Pgen between shared (2.38e-7) and full (1.57e-7)

268 repertoires was smaller than for TCRβ but still distinct.

269

270 Since our sequencing method uses genomic DNA instead of messenger-RNA as starting material, it

271 has been optimized for quantitative analysis and provides us with a reasonable estimate of the

272 clonal abundance (Robins et al., 2009; Vanhanen et al., 2016). Thus, the analysis of the shared

273 fraction of total genomes reflects the actual size of repertoire common to different individuals. For

274 total genomes, a similar increasing trend in JIs from nonproductive to in-frame and to amino acid

275 repertoires was observed as seen for unique sequences. In total in-frame nucleotide genomes the

276 mean JI for TCRα repertoire was 0.083 and for TCRβ repertoire 0.00063. In total amino acids the

277 shared part of the repertoire was extremely large (mean JI 0.30 for TCRα and 0.026 for TCRβ;

278 Figure 5C). In percentages, on the average, of the total TCRβ amino acid repertoire of any given

279 individual 6.1% was also found in the repertoire of another individual (range 1.55-11.4%). In the

280 TCRα repertoire the overlap in percentages was strikingly high (mean 46.7%, range 32.6-62.7%;

281 Supplement 5).

282

283    **3.6.    Sharing of high abundance clones**

284 To analyze the relationship between clone size and the likelihood of sharing, we calculated the

285 Jaccard indexes for the most abundant 1%, 2%, 5%, 10%, 20% and 50% of clones. For this analysis

286 samples 1-4 were excluded, because the relatively shallow sequencing produced very little overlap

287 among the top 1-5% clonotypes. In TCRα repertoire we observed a clear correlation between the

288 sharing and the clonotype abundance. JI values were clearly highest in the top 1-2% most abundant

289 clonotypes and decreased gradually when less abundant clonotypes were included (Figure 6A).  In

13

290      contrast, there was no similar correlation in the TCRβ repertoire and the interindividual variation in

291      JIs among the top 1% most abundant clones was very wide (Figure 6A). The number of non-

292      templated nucleotide inserts also showed a correlation with the sharing among highly abundant

293      clones. Non-templated inserts were rare among the most abundant shared clones. In the TCRα

294      repertoire the average number of inserts in the shared repertoire increased steadily with the analysis

295      of less abundant clonotypes (Figure 6B). In the TCRβ repertoire the number of inserts was typically

296      zero among the top 2% most abundant shared repertoire and increased abruptly for the top 5-50%

297      most abundant clonotypes (Figure 6B).

298

299      **3.7.     Sequence overlap in the peripheral samples**

300 Despite the clearly smaller number of cells analyzed, clonotype sharing was also observed in the

301 peripheral blood. Similarly to the thymus, sharing was higher in the TCRα than in the TCRβ

302 repertoire and some clonotypes were shared between all four samples (Figure 7A). Also in the

303 peripheral samples, sharing was lowest in the nonproductive nucleotide repertoire, increased in the

304 in-frame nucleotide and even more so in the amino acid CDR3 repertoires (Figure 7B-C).

## 4. Discussion

Until recently, our understanding of the human thymus has been largely based on extrapolation from circulating repertoire and from murine studies. However, studies on organ donors combined with high-throughput techniques and next-generation sequencing have begun to provide information on the various types of cells in the human thymus (Park et al., 2020; Thome et al., 2016). A single-cell sequencing study coupled with TCRαβ profiling identified approximately 200 000 individual lymphoid cells among 24 fetal and mature thymi and showed a biased V(D)J usage originating from recombination and modified by selection (Park et al., 2020). We have previously estimated the total thymic TCR diversity to be 60-100x10$^6$ for TCRα and 40-70x10$^6$ for TCRβ repertoire and thus currently beyond the coverage of single-cell experiments (Vanhanen et al., 2016). Our current data from eight pediatric thymi comprises a total of 161 million TCRα reads and 55 million TCRβ reads, representing the most extensive characterization of the thymic TCR repertoire so far. Although our analysis was performed on unsorted cells and thus allows little conclusions on the developmental stage and functionality of the TCRs, the large scale provides an opportunity to compare specific features of TCRα and TCRβ repertoires and, particularly, to measure thymic repertoire overlap across individuals.

As previously reported for peripheral blood samples and recently for thymus as well (Park et al., 2020; Quiros Roldan et al., 1995; Zvyagin et al., 2014) the usage of V and J gene segments is clearly biased in the thymus. The same gene segments were dominant in every individual, in both the TCRα and TCRβ chains. Some of this bias has been ascribed to selection by HLA molecules, which interact with protein loops encoded by the germ-line parts of TCR V genes (Huseby et al., 2005; Rudolph et al., 2006; Wu et al., 2002). However, the same biased usage was also observed in the nonproductive repertoire, which cannot be subjected to selection by antigen-HLA complexes. This suggests that the bias is partly generated in the recombination itself. We have previously

15

330 reported that genetic factors influence the gene segment usage in the thymus, a finding confirmed

331 here with an increased number of samples. Our data also show that the use of TCRD elements in αβ

332 T cells is common, with ca. 6% of thymic sequences containing TCRD gene segments, while the

333 frequency of γδ TCR+ thymocytes was less than 1%. However, combining TCRAV to TCRDJ

334 seems to be largely prevented.

335

336 Despite the structural and functional similarity of the two TCR chains, the generation of TCRα and

337 TCRβ repertoire has several differences, which are also reflected in our data. First, the number of

338 non-templated nucleotide insertions was much higher in the TCRβ locus, most likely explainable by

339 the fact that, unlike TCRα chain, TCRβ chain undergoes two recombination events (D to J followed

340 by V to DJ). This is also displayed by the slightly longer CDR3 region length and lower calculated

341 Pgen in TCRβ than in TCRα repertoire. Second, the number of non-templated inserts and CDR3

342 length were lower and respectively Pgen was higher in the peripheral than in thymic samples in

343 TCRβ repertoire while in TCRα these features remained relatively similar in thymus and periphery.

344 Third, the fraction of in-frame rearrangements was higher in TCRβ than in TCRα locus (78.8% vs.

345 30.8% in the thymus). This reflects the difference in allelic exclusion in TCRβ and TCRα locus. In

346 TCRβ locus the exclusion is strict, whereas both TCRα loci are rearranged simultaneously and a

347 large fraction of cells will end up with a nonfunctional rearrangement in the other TCRα locus

348 (Borgulya et al., 1992; Casanova et al., 1991). The frequency of nonproductive sequences in our

349 samples is also increased by the presence of immature thymocytes not yet subjected to TCR-

350 mediated selection, and the ongoing TCRα locus rearrangement in some of the cells. Furthermore,

351 the repertoire overlap was much higher in TCRα than TCRβ repertoire, consistent with previous

352 analyses (Khosravi-Maharlooei et al., 2019; Zvyagin et al., 2014), but here shown in a large-scale

353 analysis of thymic repertoire. Although we obtained fewer TCRβ than TCRα sequences, it is clear

354 <u>that the higher sequence overlap in TCRα compared with TCRβ is mostly biological and not due to</u>

355 <u>differences in sequencing depth, a finding also confirmed by others.</u>

356

357 Indeed, the remarkably high degree of clonal sharing between individuals is the most interesting

358 observation in our current data. Here, it must be noted that two of our donors were monozygotic

359 twins, which introduces a bias to the analysis. However, these two samples only shared a

360 marginally higher fraction of sequences than the unrelated samples (Heikkila et al., 2020). This is

361 largely consistent with a recent analysis of the peripheral repertoire in three pairs of identical twins,

362 which concluded that there was no difference between the twins and unrelated donors in the sharing

363 of CDR3 sequences (Zvyagin et al., 2014).

364

365 Given the enormous diversity of possible TCRs, the expected likelihood to detect identical

366 receptors in two individuals is practically nonexistent. Still, previous studies of inbred mouse lines

367 have reported that roughly 30% of the peripheral TCRβ repertoire is shared (Bousso et al., 1998;

368 Furmanski et al., 2008). More recent studies have used next-generation sequencing methods,

369 analyzing much larger numbers of sequences. Sequencing of TRBV12-4/TRBJ1-2 expressing

370 peripheral blood CD8+ T cells in four unrelated healthy donors yielded in average 29 000 unique

371 clonotypes per individual and the overlap of unique amino acid CDR3 sequences was 3.8–9.8%

372 (Venturi et al., 2011). Zvyagin et al. measured the overlap of both TCRβ and TCRα repertoires in

373 three pairs of monozygotic twins reaching an overlap of 3–10% and 10–26.5% of unique amino

374 acid CDR3 clonotypes in TCRβ and TCRα repertoires, respectively, without higher similarity

375 between the twins than unrelated pairs (Zvyagin et al., 2014). It was also estimated that if the

376 predicted peripheral TCRβ diversity of $5 \times 10^6$ unique sequences was entirely sequenced, the CDR3

377 overlap between two individuals would reach 44.1% in the amino acid and 3.6% in the nucleotide

378 repertoire (Shugay et al., 2013).

379

380 In our thymus data, taking into account the clonal abundances of the clonotypes, the average

381 fraction of sequences found in any other donor for the total TCRβ repertoire was 0.2% for in-frame

382 nucleotide chains and 6.1% for amino acid CDR3 chains. In the TCRα repertoire, the sharing was

383 much higher: an average of 15.7% total in-frame nucleotide and 46.7% total amino acid CDR3

384 chains were shared between any two donors. Furthermore, in the TCRα repertoire the overlap

385 showed a strong correlation with clone abundance, whereas the same was not true of the TCRβ

386 repertoire. The average number of non-templated inserts was also lower among the most abundant

387 clonotypes both for TCRα and TCRβ repertoire. Together, these data suggest that some TCRα

388 sequences are generated easily and preferred across different unrelated individuals.

389

390 Notably, the surprisingly high repertoire overlap in thymus was directly measured from samples of

391 10 million thymocytes, taken from an organ with an estimated 50 billion cells (Ganusov and De

392 Boer, 2007; Rodewald, 2008). Recent analyses have shown that the degree of sharing is correlated

393 with the size of the sample sequenced (Campregher et al., 2010; Putintseva et al., 2013; Shugay et

394 al., 2013; Venturi et al., 2011). It is thus possible that exhaustive sequencing of the thymic

395 repertoire would reveal an even higher proportion of shared sequences. Indeed, it may be that TCR

396 repertoire is really individualized only by α-to-β pairing, although even this may be less stochastic

397 than previously assumed (Grigaityte et al., 2017).

398

399 It is clear that some of this sharing reflects convergent recombination, i.e., the recombination

400 machinery favoring certain gene segments and particular types of CDR3 sequences. Previous

401 analysis of peripheral repertoire has shown that shared sequences have relatively few nucleotide

402 additions and are generally closer to germline sequences (Pogorelyy et al., 2017; Quigley et al.,

403 2010; Venturi et al., 2008a; Venturi et al., 2006). This is also seen in our thymus samples, where the

404 shared sequences had on average fewer nucleotide insertions and higher calculated Pgen than the

405 repertoire in general. This implies that some junctional sequences are easier to generate and

406 therefore appear repeatedly, and their high frequency may therefore not require peripheral

407 expansion (Venturi et al., 2008b).

408

409 However, our quantitation showed a strong enrichment of the shared repertoire the further the

410 sequences receded from the recombination process. In every donor pair the shared fraction was

411 higher in the in-frame than in the nonproductive repertoire and higher still in amino acid sequences

412 and total number of genomes. This was particularly striking in the TCRβ repertoire, in which the

413 average JI increased from $6.3 \times 10^{-5}$ in the nonproductive repertoire to 0.026 in total amino acid

414 genomes, or by a factor of ~400. In the TCRα chain the increase was by a factor of ~10, from 0.026

415 to 0.30. Since the nonproductive nucleotide sequences are not subject to any form of TCR-mediated

416 selection, this enrichment indicates that a substantial fraction of the clonal sharing is due to antigen-

417 driven selection in the thymus.

418

419 In the periphery, although shared clones specific to defined antigens have been described, the

420 antigen-dependent selection seems in general to lead to divergence in the repertoire. Analysis of

421 naive and memory CD8+ T cells found fewer shared clones in the latter, antigen-experienced

422 population, while a comparison of preterm neonates with adults showed that the shared fraction of

423 TCRβ (CDR3 amino acid chains) decreased from 8% to 1% (Carey et al., 2017). Similarly, donors

424 in younger age groups shared a larger fraction of TCRβ repertoire than older individuals and while

425 TCRβ repertoires in young are similarly high in diversity, with age clonal expansions accumulate

426 and the individual repertoires develop to divergent directions (Britanova et al., 2014; Britanova et

427 al., 2016). In our data the shared fraction of CDR3 amino acid sequences in the peripheral blood

428  was 5.3% in the TCRβ and 17.1% in the TCRα repertoire, in donors ranging from 7 days to 5

429  months of age.

430

431  A further point relates to the transitory nature of thymic function. Since thymus is a primary

432  lymphoid organ constantly producing new T cells, any given clone will spend only a limited time in

433  the thymus before either failing selection and dying or maturing and emigrating to periphery. The

434  repertoire might thus also be expected to be transitory, with a different snapshot of the repertoire

435  obtained at different points in time. In contrast, the high degree of interindividual clonal sharing

436  suggests by extension that at different time points a given thymus is producing similar clones.

437  Indirectly, our results imply that although the thymic T cell population and TCR repertoire is

438  transitory, the clonal composition of human thymus is surprisingly stable.

439

440  In conclusion, our study provides the first detailed characterization of the human thymic TCRα and

441  TCRβ repertoire, showing similarities and differences in the features of these two TCR chains. We

442  also show an unexpectedly high overlap of thymic TCR repertoire between unrelated donors,

443  especially in the TCRα chain. Moreover, our data indicate that this convergence is substantially

444  driven by thymic selection. Finally, it must be noted that the specificity of any TCR is determined

445  by α-to-β pairing, which our data do not address. As shown by Grigaityte et al., novel technology is

446  finally allowing this part of the repertoire to be analyzed, as well (Grigaityte et al., 2017).

447

## References:

458  **References:**

459  **Arstila, T.P., Casrouge, A., Baron, V., Even, J., Kanellopoulos, J., Kourilsky, P.**, 1999. A direct

460  estimate of the human alphabeta T cell receptor diversity. Science 286, 958-961

461  **Borgulya, P., Kishi, H., Uematsu, Y., von Boehmer, H.**, 1992. Exclusion and Inclusion of Alpha-

462  T-Cell and Beta-T-Cell Receptor Alleles. Cell 69, 529-537. 10.1016/0092-8674(92)90453-J

463  **Bousso, P., Casrouge, A., Altman, J.D., Haury, M., Kanellopoulos, J., Abastado, J.P.,**

464  **Kourilsky, P.**, 1998. Individual variations in the murine T cell response to a specific peptide reflect

465  variability in naive repertoires. Immunity 9, 169-178

466  **Britanova, O.V., Putintseva, E.V., Shugay, M., Merzlyak, E.M., Turchaninova, M.A.,**

467  **Staroverov, D.B., Bolotin, D.A., Lukyanov, S., Bogdanova, E.A., Mamedov, I.Z., Lebedev,**

468  **Y.B., Chudakov, D.M.**, 2014. Age-related decrease in TCR repertoire diversity measured with

469  deep and normalized sequence profiling. J Immunol 192, 2689-2698. 10.4049/jimmunol.1302064

470  **Britanova, O.V., Shugay, M., Merzlyak, E.M., Staroverov, D.B., Putintseva, E.V.,**

471  **Turchaninova, M.A., Mamedov, I.Z., Pogorelyy, M.V., Bolotin, D.A., Izraelson, M., Davydov,**

472  **A.N., Egorov, E.S., Kasatskaya, S.A., Rebrikov, D.V., Lukyanov, S., Chudakov, D.M.**, 2016.

473  Dynamics of Individual T Cell Repertoires: From Cord Blood to Centenarians. J Immunol 196,

474  5005-5013. 10.4049/jimmunol.1600005

475  **Campregher, P.V., Srivastava, S.K., Deeg, H.J., Robins, H.S., Warren, E.H.**, 2010.

476  Abnormalities of the alphabeta T-cell receptor repertoire in advanced myelodysplastic syndrome.

477  Exp Hematol 38, 202-212. 10.1016/j.exphem.2009.12.004

478  **Carey, A.J., Hope, J.L., Mueller, Y.M., Fike, A.J., Kumova, O.K., van Zessen, D.B.H.,**

479  **Steegers, E.A.P., van der Burg, M., Katsikis, P.D.**, 2017. Public Clonotypes and Convergent

480  Recombination Characterize the Naive CD8(+) T-Cell Receptor Repertoire of Extremely Preterm

481  Neonates. Front Immunol 8, 1859. 10.3389/fimmu.2017.01859

482 **Casanova, J.L., Romero, P., Widmann, C., Kourilsky, P., Maryanski, J.L.**, 1991. T cell

483 receptor genes in a series of class I major histocompatibility complex-restricted cytotoxic T

484 lymphocyte clones specific for a Plasmodium berghei nonapeptide: implications for T cell allelic

485 exclusion and antigen-specific repertoire. J Exp Med 174, 1371-1383

486 **Chao, A., Chazdon, R.L., Colwell, R.K., Shen, T.J.**, 2006. Abundance-based similarity indices

487 and their estimation when there are unseen species in samples. Biometrics 62, 361-371.

488 10.1111/j.1541-0420.2005.00489.x

489 **Davis, M.M., Bjorkman, P.J.**, 1988. T-cell antigen receptor genes and T-cell recognition. Nature

490 334, 395-402. 10.1038/334395a0

491 **Egerton, M., Scollay, R., Shortman, K.**, 1990. Kinetics of mature T-cell development in the

492 thymus. Proc Natl Acad Sci U S A 87, 2579-2582

493 **Furmanski, A.L., Ferreira, C., Bartok, I., Dimakou, S., Rice, J., Stevenson, F.K., Millrain,**

494 **M.M., Simpson, E., Dyson, J.**, 2008. Public T cell receptor beta-chains are not advantaged during

495 positive selection. J Immunol 180, 1029-1039

496 **Ganusov, V.V., De Boer, R.J.**, 2007. Do most lymphocytes in humans really reside in the gut?

497 Trends in Immunology 28, 514-518. 10.1016/j.it.2007.08.009

498 **Grigaityte, K., Carter, J.A., Goldfless, S.J., Jeffery, E.W., Hause, R.J., Jiang, Y., Koppstein,**

499 **D., Briggs, A.W., Church, G.M., Vigneault, F., Atwal, G.S.**, 2017. Single-cell sequencing reveals

500 αβ chain pairing shapes the T cell repertoire. bioRxiv.

501 **Heikkila, N., Vanhanen, R., Yohannes, D.A., Saavalainen, P., Meri, S., Jokiranta, T.S., Jarva,**

502 **H., Mattila, I.P., Hamm, D., Sormunen, S., Saramaki, J., Arstila, T.P.**, 2020. Identifying the

503 inheritable component of human thymic T cell repertoire generation in monozygous twins. Eur J

504 Immunol 50, 748-751. 10.1002/eji.201948404

505  **Huseby, E.S., White, J., Crawford, F., Vass, T., Becker, D., Pinilla, C., Marrack, P., Kappler,**

506  **J.W.**, 2005. How the T cell repertoire becomes peptide and MHC specific. Cell 122, 247-260.

507  10.1016/j.cell.2005.05.013

508  **Khosravi-Maharlooei, M., Obradovic, A., Misra, A., Motwani, K., Holzl, M., Seay, H.R.,**

509  **DeWolf, S., Nauman, G., Danzl, N., Li, H., Ho, S.H., Winchester, R., Shen, Y., Brusko, T.M.,**

510  **Sykes, M.**, 2019. Crossreactive public TCR sequences undergo positive selection in the human

511  thymic repertoire. J Clin Invest 129, 2446-2462. 10.1172/JCI124358

512  **Klein, L., Kyewski, B., Allen, P.M., Hogquist, K.A.**, 2014. Positive and negative selection of the

513  T cell repertoire: what thymocytes see (and don't see). Nat Rev Immunol 14, 377-391.

514  10.1038/nri3667

515  **Matsutani, T., Ogata, M., Fujii, Y., Kitaura, K., Nishimoto, N., Suzuki, R., Itoh, T.**, 2011.

516  Shortening of complementarity determining region 3 of the T cell receptor alpha chain during

517  thymocyte development. Mol Immunol 48, 623-629. 10.1016/j.molimm.2010.11.003

518  **Niemi, H.J., Laakso, S., Salminen, J.T., Arstila, T.P., Tuulasvaara, A.**, 2015. A normal T cell

519  receptor beta CDR3 length distribution in patients with APECED. Cell Immunol 295, 99-104.

520  10.1016/j.cellimm.2015.03.005

521  **Park, J.E., Botting, R.A., Dominguez Conde, C., Popescu, D.M., Lavaert, M., Kunz, D.J., Goh,**

522  **I., Stephenson, E., Ragazzini, R., Tuck, E., Wilbrey-Clark, A., Roberts, K., Kedlian, V.R.,**

523  **Ferdinand, J.R., He, X., Webb, S., Maunder, D., Vandamme, N., Mahbubani, K.T., Polanski,**

524  **K., Mamanova, L., Bolt, L., Crossland, D., de Rita, F., Fuller, A., Filby, A., Reynolds, G.,**

525  **Dixon, D., Saeb-Parsy, K., Lisgo, S., Henderson, D., Vento-Tormo, R., Bayraktar, O.A.,**

526  **Barker, R.A., Meyer, K.B., Saeys, Y., Bonfanti, P., Behjati, S., Clatworthy, M.R., Taghon, T.,**

527  **Haniffa, M., Teichmann, S.A.**, 2020. A cell atlas of human thymic development defines T cell

528  repertoire formation. Science 367. 10.1126/science.aay3224

529 **Pogorelyy, M.V., Elhanati, Y., Marcou, Q., Sycheva, A.L., Komech, E.A., Nazarov, V.I.,**

530 **Britanova, O.V., Chudakov, D.M., Mamedov, I.Z., Lebedev, Y.B., Mora, T., Walczak, A.M.**,

531 2017. Persisting fetal clonotypes influence the structure and overlap of adult human T cell receptor

532 repertoires. PLoS Comput Biol 13. 10.1371/journal.pcbi.1005572

533 **Putintseva, E.V., Britanova, O.V., Staroverov, D.B., Merzlyak, E.M., Turchaninova, M.A.,**

534 **Shugay, M., Bolotin, D.A., Pogorelyy, M.V., Mamedov, I.Z., Bobrynina, V., Maschan, M.,**

535 **Lebedev, Y.B., Chudakov, D.M.**, 2013. Mother and child T cell receptor repertoires: deep

536 profiling study. Front Immunol 4, 463. 10.3389/fimmu.2013.00463

537 **Qi, Q., Liu, Y., Cheng, Y., Glanville, J., Zhang, D., Lee, J.Y., Olshen, R.A., Weyand, C.M.,**

538 **Boyd, S.D., Goronzy, J.J.**, 2014. Diversity and clonal selection in the human T-cell repertoire.

539 Proc Natl Acad Sci U S A 111, 13139-13144. 10.1073/pnas.1409155111

540 **Quigley, M.F., Greenaway, H.Y., Venturi, V., Lindsay, R., Quinn, K.M., Seder, R.A., Douek,**

541 **D.C., Davenport, M.P., Price, D.A.**, 2010. Convergent recombination shapes the clonotypic

542 landscape of the naive T-cell repertoire. Proc Natl Acad Sci U S A 107, 19414-19419.

543 10.1073/pnas.1010586107

544 **Quiros Roldan, E., Sottini, A., Bettinardi, A., Albertini, A., Imberti, L., Primi, D.**, 1995.

545 Different TCRBV genes generate biased patterns of V-D-J diversity in human T cells.

546 Immunogenetics 41, 91-100

547 **Robins, H.S., Campregher, P.V., Srivastava, S.K., Wacher, A., Turtle, C.J., Kahsai, O.,**

548 **Riddell, S.R., Warren, E.H., Carlson, C.S.**, 2009. Comprehensive assessment of T-cell receptor

549 beta-chain diversity in alphabeta T cells. Blood 114, 4099-4107. 10.1182/blood-2009-04-217604

550 **Rodewald, H.-R.**, 2008. Thymus Organogenesis. 10.1146/annurev.immunol.26.021607.090408

551 **Rudolph, M.G., Stanfield, R.L., Wilson, I.A.**, 2006. How TCRs bind MHCs, peptides, and

552 coreceptors. Annu Rev Immunol 24, 419-466. 10.1146/annurev.immunol.23.021704.115658

553    **Sethna, Z., Elhanati, Y., Callan, C.G., Walczak, A.M., Mora, T.**, 2019. OLGA: fast computation

554    of generation probabilities of B- and T-cell receptor amino acid sequences and motifs.

555    Bioinformatics 35, 2974-2981. 10.1093/bioinformatics/btz035

556    **Shugay, M., Bagaev, D.V., Turchaninova, M.A., Bolotin, D.A., Britanova, O.V., Putintseva,**

557    **E.V., Pogorelyy, M.V., Nazarov, V.I., Zvyagin, I.V., Kirgizova, V.I., Kirgizov, K.I.,**

558    **Skorobogatova, E.V., Chudakov, D.M.**, 2015. VDJtools: Unifying Post-analysis of T Cell

559    Receptor Repertoires. PLoS Comput Biol 11, e1004503. 10.1371/journal.pcbi.1004503

560    **Shugay, M., Bolotin, D.A., Putintseva, E.V., Pogorelyy, M.V., Mamedov, I.Z., Chudakov,**

561    **D.M.**, 2013. Huge Overlap of Individual TCR Beta Repertoires. Front Immunol 4, 466.

562    10.3389/fimmu.2013.00466

563    **Swat, W., Dessing, M., von Boehmer, H., Kisielow, P.**, 1993. CD69 expression during selection

564    and maturation of CD4+8+ thymocytes. Eur J Immunol 23, 739-746. 10.1002/eji.1830230326

565    **Thome, J.J., Bickham, K.L., Ohmura, Y., Kubota, M., Matsuoka, N., Gordon, C., Granot, T.,**

566    **Griesemer, A., Lerner, H., Kato, T., Farber, D.L.**, 2016. Early-life compartmentalization of

567    human T cell differentiation and regulatory function in mucosal and lymphoid tissues. Nat Med 22,

568    72-77. 10.1038/nm.4008

569    **Tickotsky, N., Sagiv, T., Prilusky, J., Shifrut, E., Friedman, N.**, 2017. McPAS-TCR: a manually

570    curated catalogue of pathology-associated T cell receptor sequences. Bioinformatics 33, 2924-2929.

571    10.1093/bioinformatics/btx286

572    **Vanhanen, R., Heikkila, N., Aggarwal, K., Hamm, D., Tarkkila, H., Patila, T., Jokiranta, T.S.,**

573    **Saramaki, J., Arstila, T.P.**, 2016. T cell receptor diversity in the human thymus. Mol Immunol 76,

574    116-122. 10.1016/j.molimm.2016.07.002

575    **Venturi, V., Chin, H.Y., Asher, T.E., Ladell, K., Scheinberg, P., Bornstein, E., van Bockel, D.,**

576    **Kelleher, A.D., Douek, D.C., Price, D.A., Davenport, M.P.**, 2008a. TCR beta-chain sharing in

577    human CD8+ T cell responses to cytomegalovirus and EBV. J Immunol 181, 7853-7862

578    **Venturi, V., Chin, H.Y., Price, D.A., Douek, D.C., Davenport, M.P.**, 2008b. The role of

579    production frequency in the sharing of simian immunodeficiency virus-specific CD8+ TCRs

580    between macaques. J Immunol 181, 2597-2609

581    **Venturi, V., Kedzierska, K., Price, D.A., Doherty, P.C., Douek, D.C., Turner, S.J., Davenport,**

582    **M.P.**, 2006. Sharing of T cell receptors in antigen-specific responses is driven by convergent

583    recombination. Proc Natl Acad Sci U S A 103, 18691-18696. 10.1073/pnas.0608907103

584    **Venturi, V., Quigley, M.F., Greenaway, H.Y., Ng, P.C., Ende, Z.S., McIntosh, T., Asher, T.E.,**

585    **Almeida, J.R., Levy, S., Price, D.A., Davenport, M.P., Douek, D.C.**, 2011. A mechanism for

586    TCR sharing between T cell subsets and individuals revealed by pyrosequencing. J Immunol 186,

587    4285-4294. 10.4049/jimmunol.1003898

588    **Verschuren, M.C., Wolvers-Tettero, I.L., Breit, T.M., van Dongen, J.J.**, 1998. T-cell receptor V

589    delta-J alpha rearrangements in human thymocytes: the role of V delta-J alpha rearrangements in T-

590    cell receptor-delta gene deletion. Immunology 93, 208-212

591    **von Boehmer, H., Aifantis, I., Azogui, O., Feinberg, J., Saint-Ruf, C., Zober, C., Garcia, C.,**

592    **Buer, J.**, 1998. Crucial function of the pre-T-cell receptor (TCR) in TCR beta selection, TCR beta

593    allelic exclusion and alpha beta versus gamma delta lineage commitment. Immunol Rev 165, 111-

594    119

595    **Wu, L.C., Tuot, D.S., Lyons, D.S., Garcia, K.C., Davis, M.M.**, 2002. Two-step binding

596    mechanism for T-cell receptor recognition of peptide MHC. Nature 418, 552-556.

597    10.1038/nature00920

598    **Yamashita, I., Nagata, T., Tada, T., Nakayama, T.**, 1993. CD69 cell surface expression identifies

599    developing thymocytes which audition for T cell antigen receptor-mediated positive selection. Int

600    Immunol 5, 1139-1150. 10.1093/intimm/5.9.1139

601 **Yassai, M., Gorski, J.**, 2000. Thymocyte maturation: selection for in-frame TCR alpha-chain

602 rearrangement is followed by selection for shorter TCR beta-chain complementarity-determining

603 region 3. J Immunol 165, 3706-3712

604 **Yates, A.J.**, 2014. Theories and quantification of thymic selection. Front Immunol 5, 13.

605 10.3389/fimmu.2014.00013

606 **Zvyagin, I.V., Pogorelyy, M.V., Ivanova, M.E., Komech, E.A., Shugay, M., Bolotin, D.A.,**

607 **Shelenkov, A.A., Kurnosov, A.A., Staroverov, D.B., Chudakov, D.M., Lebedev, Y.B.,**

608 **Mamedov, I.Z.**, 2014. Distinctive properties of identical twins' TCR repertoires revealed by high-

609 throughput sequencing. Proc Natl Acad Sci U S A 111, 5980-5985. 10.1073/pnas.1319389111

610

611

**Tables**

**Table 1.** Description of the samples. The details of each sequenced sample and the numbers of

obtained unique clonotypes and total reads per sample for TCRα and TCRβ repertoires.

| Sample | Age (days) | Sex | TCRα | | TCRβ | |
|---|---|---|---|---|---|---|
| | | | Unique | Total | Unique | Total |
| Thymus A | 243 | M | 6 907 422 | 39 865 283 | 1 254 760 | 8 431 833 |
| Thymus B | 244 | M | 7 578 104 | 45 335 572 | 1 540 161 | 11 558 445 |
| Thymus C | 225 | F | 5 347 824 | 30 309 225 | 1 568 528 | 23 581 729 |
| Thymus D | 126 | M | 6 743 495 | 36 762 724 | 1 462 150 | 11 159 872 |
| Thymus 1 | 7 | M | 2 089 557 | 3 179 774 | 223 725 | 237 063 |
| Thymus 2 | 52 | M | 1 262 845 | 1 747 487 | 173 368 | 182 356 |
| Thymus 3 | 107 | M | 1 289 728 | 2 158 043 | 138 544 | 142 903 |
| Thymus 4 | 156 | F | 1 419 013 | 1 848 851 | 122 195 | 128 228 |
| **Average** | | | **4 079 749** | **20 150 870** | **810 429** | **6 927 804** |
| | | | | | | |
| Blood 1 | 7 | M | 138 159 | 154 682 | 77 868 | 82 418 |
| Blood 2 | 52 | M | 109 171 | 123 523 | 69 875 | 73 945 |
| Blood 3 | 107 | M | 180 100 | 245 126 | 104 236 | 134 110 |
| Blood 4 | 156 | F | 167 266 | 199 326 | 82 550 | 88 901 |
| **Average** | | | **148 674** | **180 664** | **83 632** | **94 844** |

617 **Table 2.** Mean frequency (%) of Vδ and Jδ segments in the TCRα repertoire

|  |  | Thymus | Peripheral blood |
|---|---|---|---|
| **Vδ-Jδ** | Unique | 1.09 | 1.71 |
|  | Total | 2.06 | 3.83 |
| **Vδ-Jα** | Unique | 2.39 | 2.71 |
|  | Total | 3.81 | 3.84 |
| **Vα-Jδ** | Unique | 0.36 | 0.42 |
|  | Total | 0.41 | 0.51 |

618

619

**Figure captions**

621

622   **Figure 1.** Analysis of the thymocyte subsets and repertoire productivity. The fraction of TCRαβ+

623   and TCRγδ+ in thymocytes and the distribution of CD4 and CD8 among TCRαβ+ thymocytes in a

624   representative thymus sample (donor C) with the applied backgating (A). The distribution of CD4

625   and CD8 expression in thymocytes and the fraction of CD3highCD69+ cells in CD4+CD8+ double

626   positive thymocytes (donor C) with the applied backgating (B). The fraction of sequences in-frame,

627   out-of-frame, or containing a premature stop codon among unique TCRα and TCRβ clonotypes for

628   thymic (C) and peripheral TCR repertoires (D).

629

630   **Figure 2.** The V gene usage in in-frame and nonproductive repertoires. The heatmaps display the

631   frequencies of different V gene segments and the attached dendrograms show the clustering of the

632   samples in in-frame and nonproductive TCRα (A) and TCRβ repertoires (B).

633

634   **Figure 3.** The number of non-templated insertions and the CDR3 lengths. The graphs show

635   the average and 95% confidence interval of the number of non-templated nucleotide insertions (A)

636   and of CDR3 lengths (B) in thymic and peripheral blood TCRα and TCRβ repertoires for in-frame

637   and nonproductive sequences.

638

639   **Figure 4.** The generation probability (Pgen) calculated with OLGA software. Thymic and

640   peripheral Pgen distribution plotted against probability density in the in-frame TCRα and TCRβ

641   repertoires for a representative thymus-blood pair (donor 1).

642

643   **Figure 5.** Sequence overlap between thymus samples. Venn diagrams show the overlap of unique

644   in-frame clonotypes separately for thymus samples A-D and 1-4 (A). Individual Jaccard indexes

645 (JI) between each thymus sample for nonproductive, in-frame and amino acid repertoires among

646 unique clonotypes (B) and total genomes (C). Monozygotic twins A and B are identified as open

647 circles, filled circles represent the JI between unrelated individuals. The average JI and the 95%

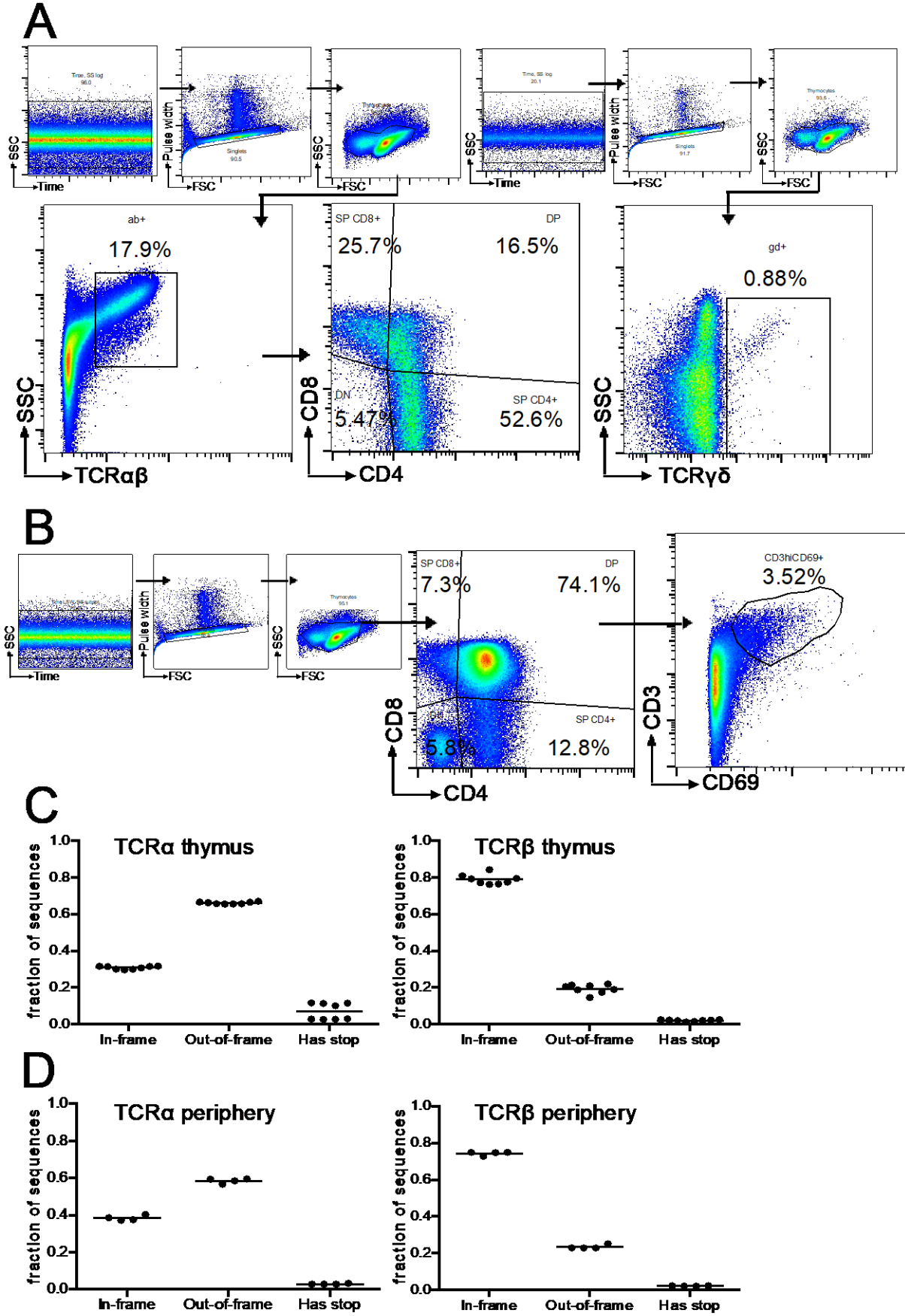648 confidence interval are shown.

649

650 **Figure 6.** Sequence overlap among the high abundance clonotypes. Jaccard indexes (A) and non-

651 templated insertions in the shared in-frame sequences (B) among the top 1%, 2%, 5%, 10%, 20%

652 and 50% most abundant clonotypes and full repertoire in thymus samples A-D. The horizontal bars

653 show the average and error bars indicate the 95% confidence interval.

654

655 **Figure 7.** Sequence overlap between peripheral blood samples. Venn diagrams show the overlap of

656 unique in-frame clonotypes (A). Jaccard indexes (JI) for nonproductive, in-frame and amino acid

657 repertoires among unique clonotypes (B) and total genomes (C). The average JI and the 95%
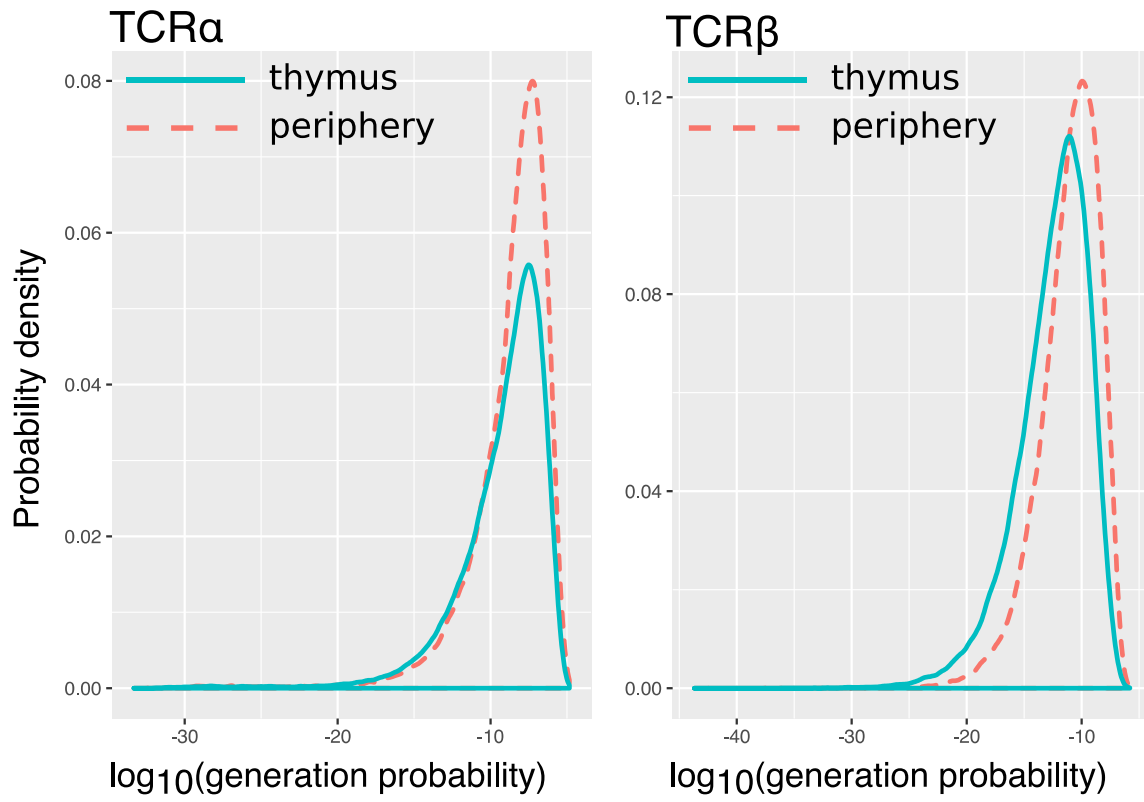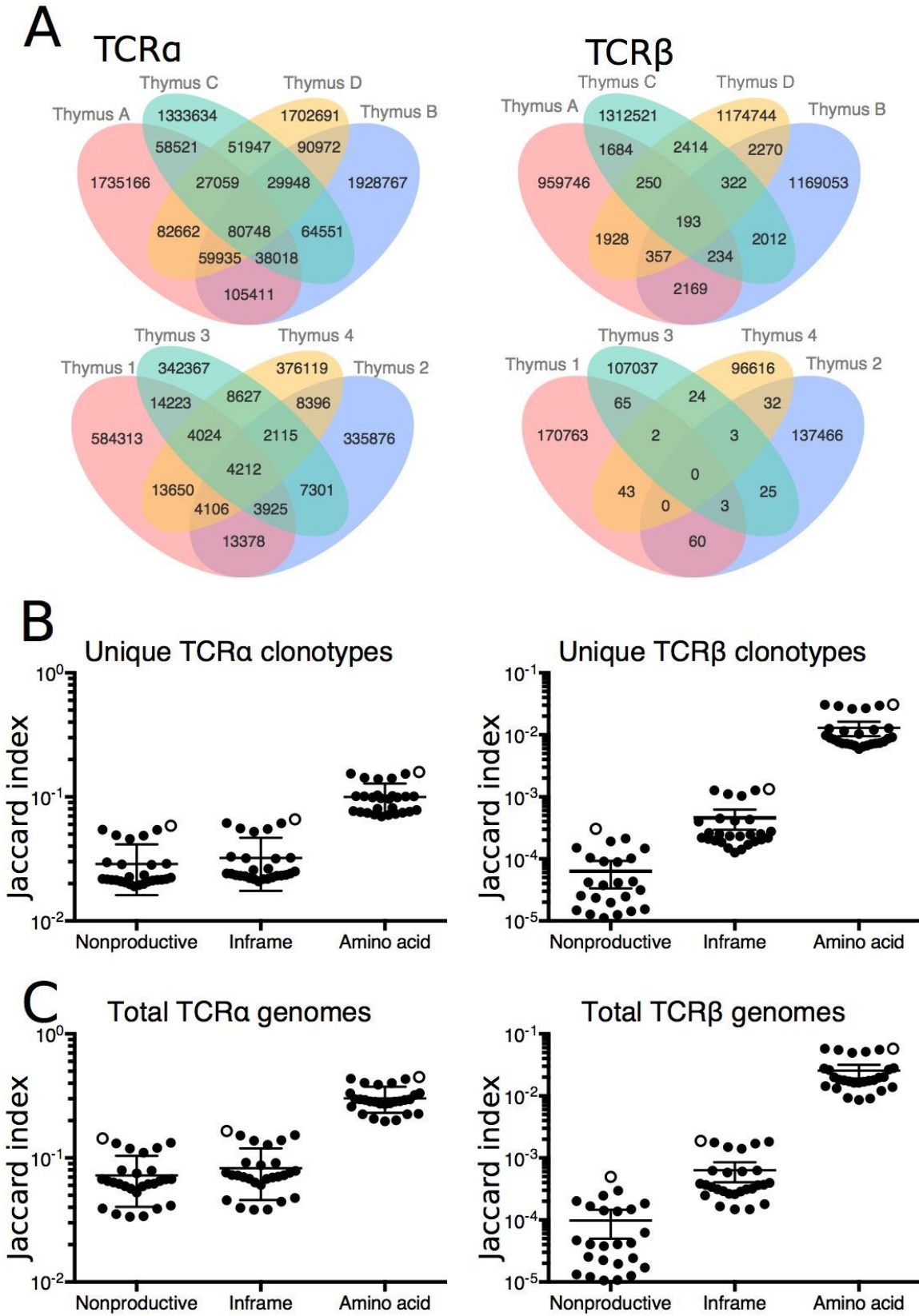
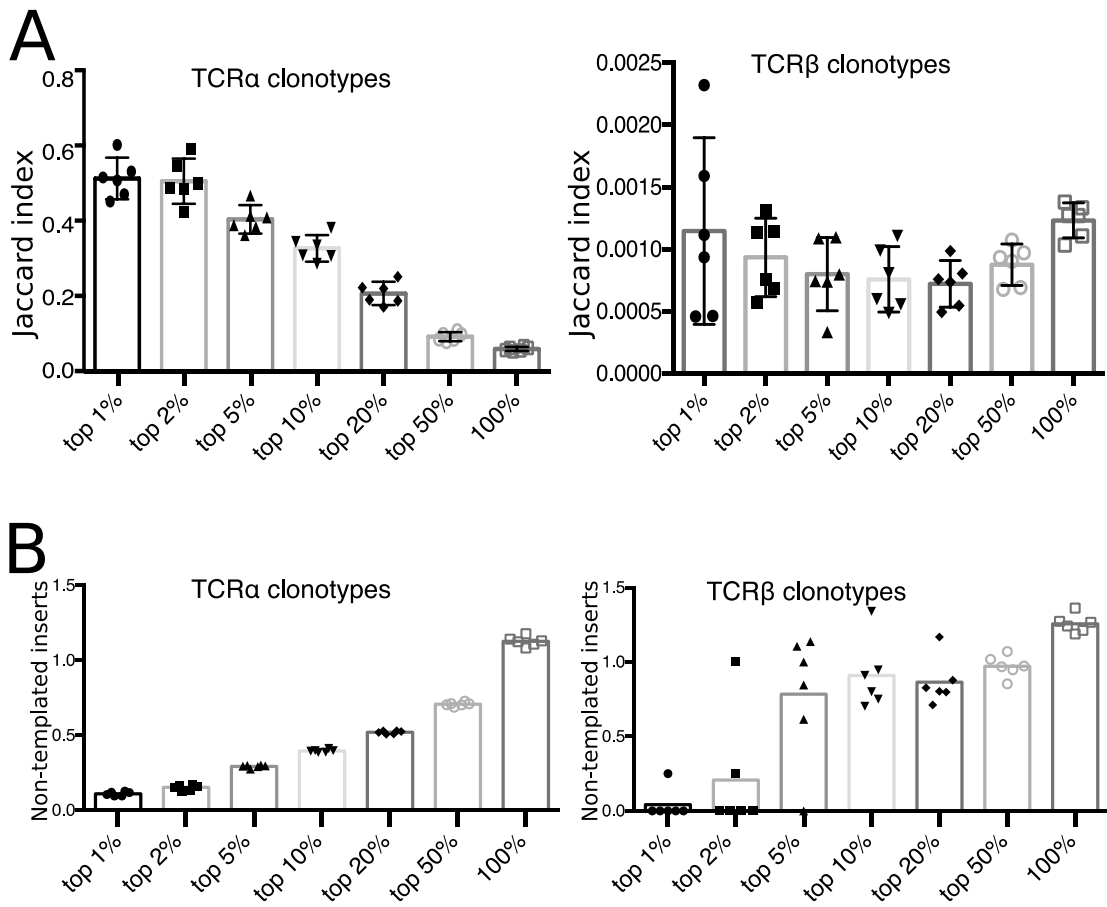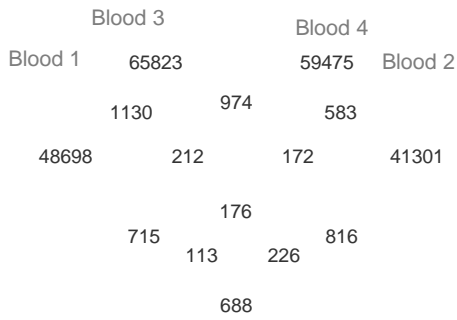658 confidence interval are displayed.

659

660

661    Figure 1

Figure 6