

# Quantifying changes in the T cell receptor repertoire during thymic development

Francesco Camaglia<sup>1§</sup>, Arie Ryvkin<sup>2§</sup>, Erez Greenstein<sup>2</sup>, Shlomit Reich-Zeliger<sup>2</sup>, Benny Chain<sup>3</sup>, Thierry Mora<sup>1\*</sup>, Aleksandra M. Walczak<sup>1\*</sup>, Nir Friedman<sup>2\*‡</sup>

**\*For correspondence:**

[thierry.mora@phys.ens.fr](mailto:thierry.mora@phys.ens.fr),  
[aleksandra.walczak@phys.ens.fr](mailto:aleksandra.walczak@phys.ens.fr)

<sup>§</sup>These authors contributed equally.

\*These authors contributed equally.

‡Deceased.

<sup>1</sup>Laboratoire de physique de l'École normale supérieure, CNRS, PSL University, Sorbonne Université, and Université de Paris, 75005 Paris, France; <sup>2</sup>Department of Immunology, Weizmann Institute of Science, Rehovot 76100, Israel; <sup>3</sup>Division of Infection and Immunity, University College London, London, United Kingdom

**Abstract** One of the feats of adaptive immunity is its ability to recognize foreign pathogens while sparing the self. During maturation in the thymus, T cells are selected through the binding properties of their antigen-specific T-cell receptor (TCR), through the elimination of both weakly (positive selection) and strongly (negative selection) self-reactive receptors. However, the impact of thymic selection on the TCR repertoire is poorly understood. Here, we use transgenic Nur77-mice expressing a T-cell activation reporter to study the repertoires of thymic T cells at various stages of their development, including cells that do not pass selection. We combine high-throughput repertoire sequencing with statistical inference techniques to characterize the selection of the TCR in these distinct subsets. We find small but significant differences in the TCR repertoire parameters between the maturation stages, which recapitulate known differentiation pathways leading to the CD4+ and CD8+ subtypes. These differences can be simulated by simple models of selection acting linearly on the sequence features. We find no evidence of specific sequences or sequence motifs or features that are suppressed by negative selection. These results favour a collective or statistical model for T-cell self non-self discrimination, where negative selection biases the repertoire away from self recognition, rather than ensuring lack of self-reactivity at the single-cell level.

## Introduction

In order to protect themselves against infection, jawed vertebrates have evolved an adaptive immune system. T lymphocytes play a leading role in this system. Each T lymphocyte expresses a unique T-cell receptor (TCR) capable of binding short protein fragments presented by the host's Major Histocompatibility Complexes (MHC), subsequently triggering clonal expansion and differentiation of immune effector function. The T cell system discriminates pathogen derived "foreign" proteins from the body's own "self" proteins, in such a way that an immune response is usually triggered only by peptides from exposure to a potentially harmful threat. We ask if we can identify specific TCR features which allow the system to discriminate foreign and self-peptides.

TCRs are generated in a stochastic assembly process based on random recombinations of genomic templates and additional non-templated insertions and deletions *Hozumi and Tonegawa (1976)*. The ability to discriminate between self and non-self targets cannot therefore be exclusively inherited, but must at least in part be learned afresh in each individual. This process is widely believed to occur during the development of haemopoietic precursors into mature T cells, which

41 occurs in a specialized microenvironment within the thymus. This process has been studied in  
42 considerable detail. T cells precursors first produce a  $\beta$  chain and if the generated chain is functional,  
43 the cell proliferates and an  $\alpha$  chain is generated. While the TCR chains are being assembled, CD4  
44 and CD8 surface markers are expressed as precursor cells transit to the Double Positive state (DP).  
45 DP TCR are subject to thymic selection, a process that tests receptor binding by presenting them  
46 with the organism's own proteins, and eliminates very weak binders (positive selection), but also too  
47 strongly self-reactive receptors (negative selection) *Yates (2014)*. During thymic selection, DP cells  
48 differentiate into CD4<sup>+</sup> or CD8<sup>+</sup> cells by keeping expression of only one of these molecules, which  
49 determines their function. While this picture is well-established and the maturation trajectory has a  
50 well established gene expression signature *Park et al. (2020)*, the TCR sequences removed during  
51 thymic selection, which should be manifested as "holes" in the repertoire, have never been directly  
52 observed. The lack of quantifiable signatures of thymic selection, differentiation and proliferation  
53 hinders a dynamic description of TCR maturation *Robert et al. (2021)*.

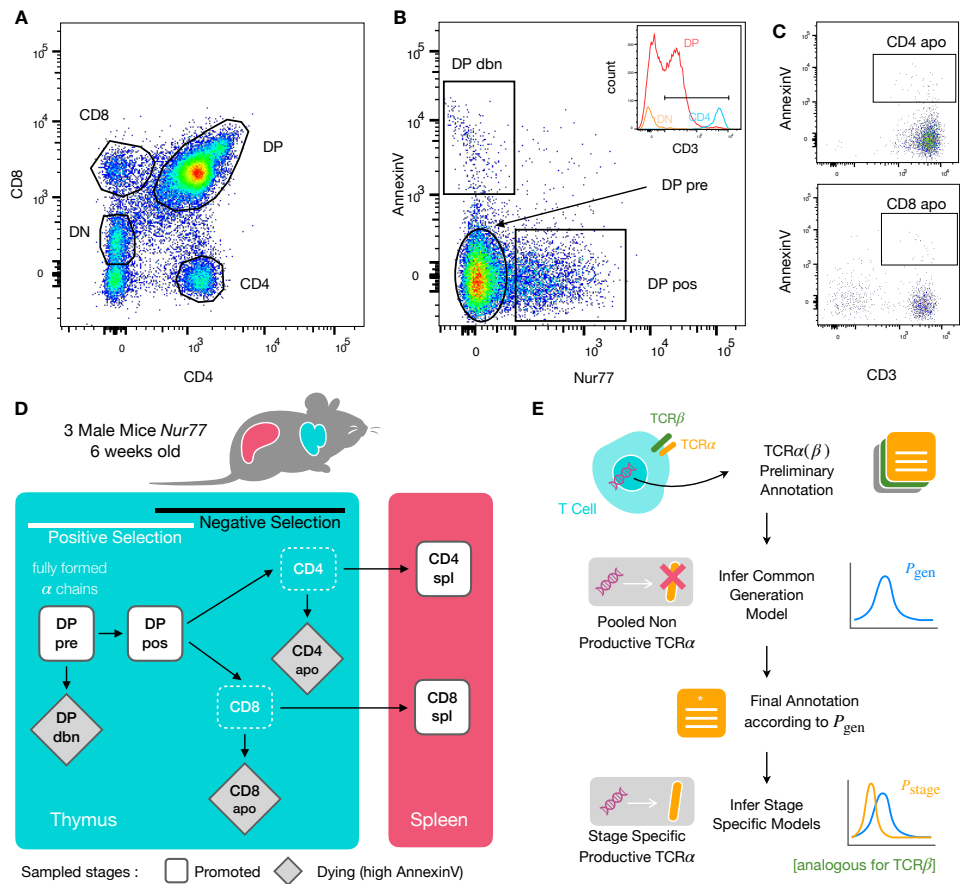
54 Positive and negative selection imposes upper and lower boundaries on the binding energy  
55 of the interaction between TCR and self peptide-MHC complexes *Košmrlj et al. (2009)*. However,  
56 it remains unclear whether every thymocyte is exposed to every self-antigen, or how efficient  
57 the process of selection is. Negative selection is known to be leaky *Yu et al. (2015)*, letting auto-  
58 reactive cells differentiate into regulatory cells *Bains et al. (2013)*; *Wing and Sakaguchi (2010)*.  
59 The efficiency of negative selection for the naive conventional (non-regulatory) effector T cell  
60 compartment remains unclear *Yu et al. (2015)*; *Gallegos and Bevan (2006)*. Partial or incomplete  
61 negative selection may limit its impact on the repertoire.

62 The difficulty of characterizing selection is partly due to survivor bias when sampling functional  
63 immune repertoires in the periphery *Madi et al. (2014, 2017)*; *Izraelson et al. (2018)*; *Sethna et al.*  
64 *(2017)*. To overcome this limitation, we sequenced the TCR repertoire of thymocyte subpopulations  
65 isolated from mice carrying a reporter transgene linked to Nur77, a marker of T cell activation  
66 both within the thymus and in the periphery. Nur77 expression, in combination with Annexin V, a  
67 marker of cell death, allows us to identify cells that are more likely to pass thymic selection, and  
68 those that are most likely not to pass selection. Although the CD4<sup>+</sup>CD8<sup>+</sup> Annexin V population may  
69 still contain some cells which will be negatively selected, but have not yet expressed Annexin V,  
70 the overall strategy provides us with a window into the repertoire at various stages of selection.  
71 By comparing the sequenced repertoires to statistical models of mouse TCR generation *Sethna*  
72 *et al. (2017)*, and subset-specific models of thymic selection, we searched for specific TCR sequence  
73 features that correlate with the different stages of intra-thymic T-cell development.

## 74 Results

### 75 Tracking T cell development stages by flow cytometry

76 To identify specific sequence features of TCR during each step of thymic selection, we performed  
77 high-throughput sequencing of TCR repertoires from different subpopulations of thymocytes from  
78 transgenic Nur77 reporter expressing mice. These mice carry a fluorescent reporter gene which  
79 is co-expressed with Nur77, a marker of T cell activation *Liebmann et al. (2018)*. Three genetically  
80 identical Nur77 reporter mice were sacrificed at the age of 6 weeks, when thymus development is  
81 completed and its cell population is stable *Gray et al. (2006)*. All animals were handled according  
82 to Weizmann Institute's Animal Care guidelines, in compliance with national and international  
83 regulations. Thymus and spleen were removed, and stained for fluorescence-activated cell sorting  
84 (see Materials and Methods). The cells were sorted based on Nur77 reporter expression (to detect  
85 activation), Annexin V (to detect early apoptosis) in combination with CD3, CD4, and CD8 cell surface  
86 markers. We used the gating strategy illustrated in *Figure 1A, B, C* to isolate double positive DP  
87 cells preceding selection (CD4<sup>+</sup>CD8<sup>+</sup>, Nur77<sup>-</sup>, Annexin V<sup>-</sup>: DP pre), DP cells in the process of being  
88 positively selected (CD4<sup>+</sup>CD8<sup>+</sup>, Nur77<sup>+</sup> Annexin V<sup>-</sup>: DP pos), DP cells dying by neglect or possibly by  
89 damage during the preparation (CD4<sup>+</sup>CD8<sup>+</sup>, Nur77<sup>-</sup> Annexin V<sup>+</sup>: DP dbn); and single positive (SP)



**Figure 1.** Experiment outline and repertoire sampling. **(A)** Flow cytometry scatterplots of T cell population from the thymus according to the markers CD4 and CD8. **(B)** The DP population is separated from DN according to CD3 expression (insert). Cells are then FACS sorted according to the expression of Nur77 and AnnexinV. **(C)** CD4 cells in the spleen (above) and CD8 (below) are FACS sorted according to the expression of CD3 and AnnexinV. **(D)** Schematic evolution of the sampled cell types during thymic maturation. **(E)** Analysis workflow: annotated reads in sampled repertoires are input for model inference (see Materials and Methods). Out-of-frame TCR sequences are pooled from all mice and stages to learn a generation model. In-frame sequences are used to learn maturation stage specific selection models with the generation model as background.

**Figure 1—figure supplement 1.** Summary of the RepSeq datasets.

90 cells: CD4<sup>+</sup>CD8<sup>-</sup>, Annexin V<sup>+</sup> (CD4 apo), and CD4<sup>-</sup>CD8<sup>+</sup>, Annexin V<sup>+</sup> cells (CD8 apo). The Annexin V  
91 staining was not very strong and did not give a very clear separation between positive and negative  
92 populations. In addition, Annexin V<sup>+</sup> subsets may be contaminated by cells that are dying for other  
93 reasons than negative selection. Nevertheless, we may still assume that the two apo subsets  
94 are enriched in negatively selected cells. In addition, we sequenced the repertoires of mature  
95 (post-selection) single positive SP CD4<sup>+</sup> and CD8<sup>+</sup> cells from the spleen (CD4 spl and CD8 spl). The  
96 proposed differentiation pathway between these populations at different maturation stages are  
97 schematically represented in **Figure 1D**. Together, these seven repertoires should contain both the  
98 selected thymocytes and the pre-selection repertoires, as well as the thymocytes that fail either  
99 positive or negative selection and die in the thymus.

### 100 **TCR repertoire sequencing**

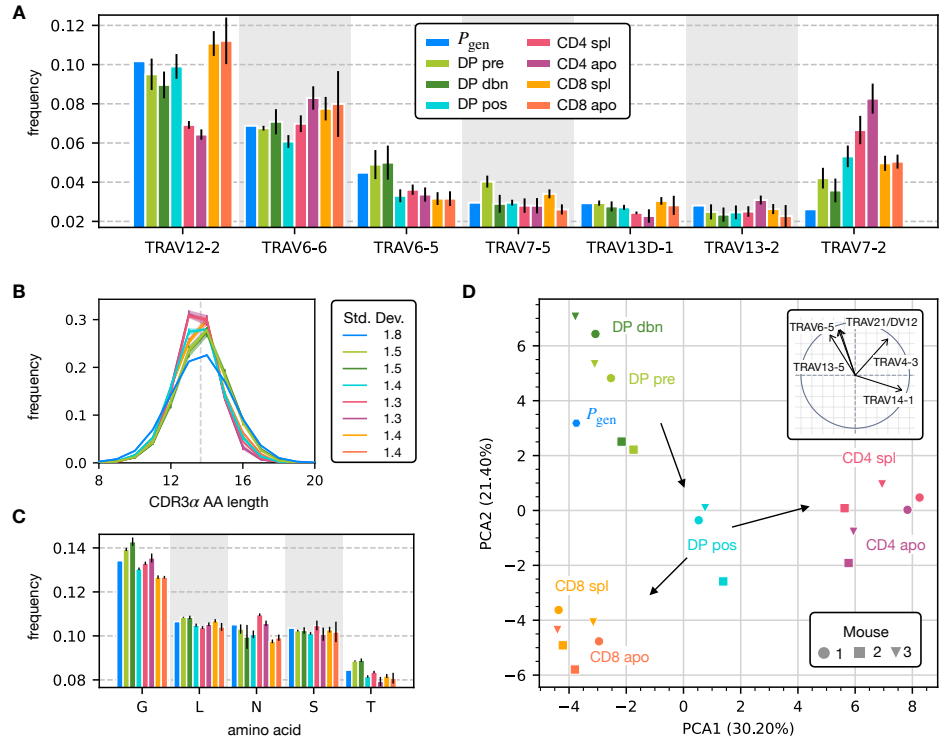
101 We sequenced and annotated the TCR repertoires of each subset as described in Materials and  
102 Methods. The cDNA of individual  $\alpha$  and  $\beta$  genes (TRA and TRB) were barcoded with unique molecular  
103 identifiers (UMI) in order to allow for correction of sequencing errors and PCR bias. However, in this  
104 analysis we focused on unique sequences (discarding count information) to avoid expression and  
105 amplification biases. As a quality control of the whole procedure, we showed that the number of  $\alpha$   
106 and  $\beta$  sequences within each population was highly correlated (**Figure 1—figure Supplement 1A**).  
107 We further verified that the relative fraction of TCR $\alpha$  sequences associated with iNKT cells (identified  
108 by TRAV11 and TRAJ18 genes *Garner et al. (2018)*) is higher in CD4 than in CD8 cells (see **Figure 1—**  
109 **figure Supplement 1B**).

110 We obtained seven datasets for both chains and for each of the 3 mice. A small fraction of  
111 sequences contain stop codons, usually because of a frameshift in the CDR3. These sequences  
112 likely come from transcription from a chromosome carrying a nonproductive chain, which is known  
113 to persist despite allelic exclusion acting on the TRB locus. The rest of the sequences are assumed  
114 to be productive. Since nonproductive TCR owe their survival to the productive gene on the other  
115 chromosome, they are not affected by selection. We combined all nonproductive sequences from  
116 all subsets to infer a generative mechanistic model of the V(D)J recombination process using IGoR  
117 *Marcou et al. (2018)*. Once trained, the model can be used to assign a generation probability  $P_{\text{gen}}$  to  
118 any TCR sequence observed *Murugan et al. (2012)*; *Marcou et al. (2018)* (see Materials and Methods  
119 and **Figure 1E**).

120 The datasets contain  $\sim 1,000$ -50,000 unique productive sequences per subset (**Figure 1—figure**  
121 **Supplement 1C** for the  $\alpha$  chain and **Figure 1—figure Supplement 1D** for the  $\beta$  chain). Since the 3  
122 mice were isogenic and shared the same MHC haplotype, we expect their repertoires to be subject  
123 to the same processes of recombination and selection *Madi et al. (2014)*. Unless specified otherwise,  
124 all downstream analyses were therefore carried out on pooled productive TCR sequences from  
125 each population from the three individuals to increase statistical power.

### 126 **Repertoires from different T cell populations have different statistical parameters.**

127  
128 To assess how selection acts at the different maturation stages, we studied the distribution of  
129 sequence features in TCR $\alpha$  repertoires. We compared TRAV and TRAJ gene usage at the different  
130 maturation stages with each other and with their expected frequency from the generation model  
131 learned from nonproductive sequences, which we will refer to as the pre-selection model or  $P_{\text{gen}}$ .  
132 TRAV usage broadly follows the pattern of the pre-selection model (**Figure 2A**), although SP CD4<sup>+</sup>  
133 repertoires have a lower proportion of TRAV12-2, and most populations have an increased pro-  
134 portion of TRAV7-2. TRAJ gene usage also broadly agrees with the pre-selection model predictions  
135 (**Figure 2—figure Supplement 2A**), although SP CD8<sup>+</sup> repertoires have a lower proportion of TRAJ31,  
136 SP CD4<sup>+</sup> repertoires have an increased proportion of TRAJ27 and TRAJ32 which is underrepresented  
137 in all cell types. For both V and J genes, we see little difference between the repertoires of spleen  
138 CD4 and CD8 cells, and their discarded counterparts in the thymus (apo). We also observe strong



**Figure 2.** Properties of the  $\alpha$  chain sequence (the analogous plot for the  $\beta$  chain is shown in *Figure 2—figure Supplement 1*). The color code is common to all subplots. **(A)** TRAV gene distribution at different maturation stages compared to the pre-selection model distribution  $P_{gen}$  (see *Figure 2—figure Supplement 2A* for TRA). Only the most frequent according to the  $P_{gen}$  model are reported. Errorbars correspond to the empirical standard deviation across the three different mice. **(B)** CDR3 length distribution of TCR $\alpha$  sequences. The errors associated with mouse variability are minor and illustrated via the shaded curves. See *Figure 2—figure Supplement 3A* for individual curves. The dashed line is the average CDR3 length from the  $P_{gen}$  model. Standard deviations of the average length distributions are shown at right. **(C)** Distribution of the most frequent amino acids at different maturation stages. The counts correspond to the number of observations within the CDR3 (i.e. excluding the first two and the last positions), summed for all the sequences in the subpopulation. Error bars represent the empirical standard deviation across mice. **(D)** Principal component analysis of the TRAV gene distribution at each maturation stage. Insert: projection on the principal axis of the five most abundant TRAV genes (see Materials and Methods). Analogous results for TRAJ are shown in *Figure 2—figure Supplement 2C*.

**Figure 2—source code 1.** [https://github.com/statbiophys/thymic\\_development\\_2022/blob/main/fig2.ipynb](https://github.com/statbiophys/thymic_development_2022/blob/main/fig2.ipynb)

**Figure 2—figure supplement 1.** Analysis of the annotated productive  $\beta$  clonotypes for the different maturation stages.

**Figure 2—figure supplement 2.** Statistics of the J gene usage and the  $P_{gen}$  distributions.

**Figure 2—figure supplement 3.** Separate amino acid CDR3 length distributions across all stages.

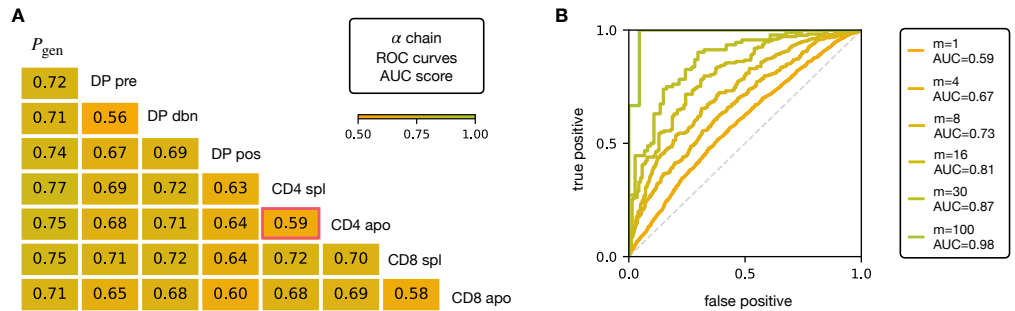
139 similarities between all the DP subsets. TRB gene usage follows similar trends, although there are  
140 some differences in J gene usage between selected and unselected SP CD4<sup>+</sup> and CD8<sup>+</sup> cells. Overall  
141 the biases of the recombination process dominate any effects of selection on V and J region usage  
142 (**Figure 2—figure Supplement 1A**, **Figure 2—figure Supplement 2B**).

143 For both chains, CDR3 amino acid length of SP CD4<sup>+</sup> and CD8<sup>+</sup> has a sharper distribution  
144 compared to earlier maturation stages (DP) (see **Figure 2B**, **Figure 2—figure Supplement 1B** and  
145 **Figure 2—figure Supplement 3**). This has previously been interpreted as a signature of selection  
146 due to structural constraints on the pMHC-TCR complex *Madi et al. (2017)*; *Lu et al. (2019)*; *Carter*  
147 *et al. (2019)*. We also compared the single amino acid usage (excluding the constant regions) across  
148 the different repertoires (**Figure 2C** for  $\alpha$  chain, **Figure 2—figure Supplement 1C** for  $\beta$  chain). We  
149 observe similarities between the DP stages, the CD4 stages and the CD8 stages, as observed for  
150 the gene usage. The repertoires from different maturation stages cannot be distinguished by any  
151 one individual feature discussed above. However, Principal Component Analysis (PCA) on the TRAV  
152 gene usage distributions in individual mice at different stages identified clusters of related cell  
153 types (**Figure 2D**). The DP Nur77<sup>-</sup> populations cluster with the pre-selection model, the SP CD4<sup>+</sup>  
154 and CD8<sup>+</sup> populations form distinct clusters, and the DP pos Nur77<sup>+</sup> cells, which we hypothesise  
155 are cells in the process of positive selection, occupy an intermediate position between these three  
156 clusters. This pattern is consistent with the known developmental trajectory as illustrated by the  
157 arrows in **Figure 2D**. PCA of TRAJ usage also shows similar clustering patterns (**Figure 2—figure**  
158 **Supplement 2C**). The PCA of TRBV and TRBJ usage also discriminates between SP CD4<sup>+</sup> and CD8<sup>+</sup>  
159 populations, and from the pre-selection populations, although the overall pattern is less clear  
160 (**Figure 2—figure Supplement 1D** and **Figure 2—figure Supplement 2D**). The overall distribution  
161 of TCR generation probabilities,  $P_{\text{gen}}$ , does not change from the pre-selection and post-selection  
162 thymic stages to the mature peripheral SP repertoires (**Figure 2—figure Supplement 2E** and **F**),  
163 consistent with previous reports comparing thymic and peripheral repertoires *Sethna et al. (2017)*.  
164 In summary, the effects of selection impose subtle changes on the pattern of TCR variable gene  
165 usage, which cannot be adequately captured by looking at any single V or J gene, but only by a  
166 combination of features.

167 V and J gene usage, and CDR3 length are coarse grained measures of a TCR repertoire. We  
168 therefore explored whether the repertoires of different maturation stages could be linked to more  
169 precise features of the TCR sequence, in particular incorporating the sequence of the CDR3. We  
170 encoded each TCR as a sparse  $\{0,1\}$  binary vector  $\vec{\sigma}$  which captures V gene, J gene and CDR amino  
171 acid sequence (for details see Materials and Methods). We then trained a logistic regression  
172 model on the set of  $\vec{\sigma}$  from repertoires of different subsets. We trained and tested the classifier  
173 to distinguish pairs of repertoires from different subsets. The classifier achieved only moderate  
174 Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) scores (**Figure 3A** for  
175 TCR $\alpha$ , and **Figure 3—figure Supplement 1A** for TCR $\beta$ ), in agreement with previous studies *Emerson*  
176 *et al. (2013)*; *Isacchini et al. (2021)*. We verify that this result is not an artifact introduced by  
177 pooling repertoires of different mice, by testing the same techniques on the individual with the  
178 largest datasets (mouse 3). The AUC scores for the  $\alpha$  and  $\beta$  repertoires are shown respectively in  
179 **Figure 3—figure Supplement 1C** and **D**.

180 Controls in which population labels were shuffled, resulted in AUC close to 0.5 (**Figure 3—**  
181 **figure Supplement 2A** and **B** for the  $\alpha$  chain, **Figure 3—figure Supplement 2C** and **D** for  $\beta$ ). The  
182 results shown in **Figure 3** indicate that the TCR populations differ at a statistical level (i.e. have  
183 different distributions of sequence features), but that each individual TCR is only a weak predictor  
184 of repertoire class. However, better classification efficiencies can be achieved by combining the  
185 predictions from sets of TCRs. For example, multiplying the predictions from 30 TCR sequences from  
186 the same repertoire (**Figure 3B**), we can distinguish CD4 spl and CD4 apo TCR $\alpha$  with an AUC score  
187 of  $>0.85$ ; see **Figure 3—figure Supplement 1B** for TCR $\beta$ . Thus statistical properties of a repertoire  
188 can distinguish it from another repertoire, even when the feature distributions of individual TCRs  
189 are largely overlapping.





**Figure 3. (A)** Area under the curves (AUC) values computed from Receiver Operating Characteristic (ROC) curves of linear classifiers of TCR $\alpha$  between two subsets. The training/testing set is a random subsample containing 70%/30% of the full dataset at a given maturation stage. **(B)** ROC curves for classifying a group of  $m$  sequences from the same maturation stage, between CD4 spl and CD4 apo (red frame in **Figure 3A**), illustrating the improvement with increasing number of TCRs. See **Figure 3—figure Supplement 1A** and **B** for the analogous analysis on TCR $\beta$ .

**Figure 3—source code 1.** [https://github.com/statbiophys/thymic\\_development\\_2022/blob/main/fig3.ipynb](https://github.com/statbiophys/thymic_development_2022/blob/main/fig3.ipynb)

**Figure 3—figure supplement 1.** AUC scores for the pooled  $\beta$  datasets and for an individual mouse.

**Figure 3—figure supplement 2.** Validation of the stages discrimination.

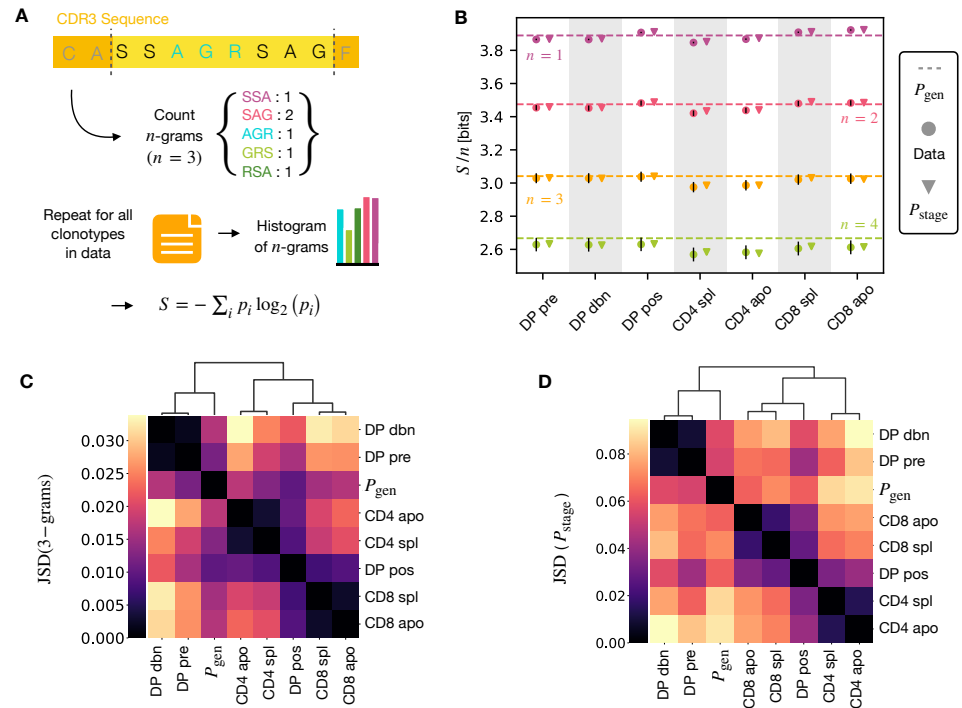
## 190 Selection models and $n$ -grams capture the relations between the stages of thymic 191 development

192 A number of studies have highlighted the importance of short amino acid motifs ( $k$ -mers or  $n$ -grams)  
193 within the CDR3 sequence in determining TCR specificity *Thomas et al. (2014)*; *Sun et al. (2017)*;  
194 *Cinelli et al. (2017)* (see **Figure 4A**). Specifically,  $n$ -grams can be used to reduce the dimensionality  
195 of the TCR space, while capturing amino acid correlations or patterns which might play a role  
196 in antigen recognition. We therefore counted the frequency of  $n$ -grams in each repertoire. We  
197 excluded from the analysis the most conserved regions (the first two positions in the CDR3 that  
198 are usually a cysteine and alanine, and the last one, typically a phenylalanine). We then used  
199 these  $n$ -gram frequency distributions to calculate the diversity of the repertoire as quantified  
200 by the Shannon entropy  $S$  (see Materials and Methods). In practice, the Shannon entropy is  
201 computationally too expensive to calculate exactly for very large data sets, and we therefore  
202 restricted our analysis to  $n$ -grams of length 4 or less, using the approximate Nemenman-Shafee-  
203 Bialek (NSB) entropy estimator *Nemenman et al. (2002)* to correct for finite sampling bias (see  
204 Materials and Methods). This estimator outcompetes alternative entropy estimators on synthetic  
205 data (**Figure 4—figure Supplement 1A** and **B**). Our analysis combines together CDR3 of different  
206 amino acid lengths which may influence the entropy measurements. However, detailed analysis  
207 of the entropy of DP repertoires, using different CDR3 lengths separately, demonstrated that the  
208 differences observed due to lengths effects were small compared to error due to sequencing  
209 (**Figure 4—figure Supplement 1C** and **D**). Another advantage of the Nemenman-Shafee-Bialek  
210 estimator is that it was shown to converge at the sizes of the smallest datasets ( $\sim 10^3 - 10^4$   
211 clonotypes), as reported in **Figure 4—figure Supplement 2**. Once computed the set of entropy  
212 measurements based on  $n$ -gram frequencies for each different repertoire, we compared the data-  
213 derived entropy measurements with the prediction of a simple generative model of each repertoire  
214 which treated each feature of each TCR (V gene, J gene and each CDR3 amino acid) as independent.  
215 Taking the set of TCR vectors  $\vec{\sigma}$  we fitted a set of parameters  $E_{\text{stage}}$  by maximising the posterior  
216 probability over all of the TCRs for each repertoire separately  $P_{\text{stage}}(\vec{\sigma}) = (1/Z)e^{-E_{\text{stage}}(\vec{\sigma})} P_{\text{gen}}(\vec{\sigma})$ , where  
217  $P_{\text{gen}}(\vec{\sigma})$  are the pre-selection generative probabilities for all the TCRs,  $E_{\text{stage}}(\vec{\sigma})$  is a linear function of the  
218 features *Elhanati et al. (2014)*; *Sethna et al. (2020)*, and  $Z$  is a normalization factor (**Figure 1E** and

219 Materials and Methods). The enrichment factors  $E_{\text{stage}}(\bar{\sigma})$  encode the intuition that due to selection,  
220 a given TCR in a given repertoire is seen with higher or lower frequency than expected by the  
221 pre-selection generation model. Once we had learnt the enrichment factors for each repertoire, we  
222 used the resulting model to generate *in silico* synthetic repertoires of  $3 \times 10^6$  TCRs, and recalculated  
223  $n$ -gram frequency distributions and entropy estimates for each synthetic repertoire.

224 The comparison of the estimated entropy for each  $n$ -gram length, and each subpopulation  
225 of T cells, using both directly data-derived and model-derived repertoires is illustrated for TCR $\alpha$   
226 (Figure 4B) and TCR $\beta$  (Figure 4—figure Supplement 3A) chains. An upper bound for the entropy is  
227 given by uniformly distributed amino acids,  $S_{\text{max}}/n = \log_2 20 \sim 4.3$  bits, while amino acids distributed  
228 according to their frequency in the overall vertebrate proteome gives a slightly smaller value of  
229  $\sim 4.2$  bits per position King and Jukes (1969). Both the observed and model-derived entropies are  
230 less than this maximum even for single amino acids ( $n$ -grams of length  $n = 1$ ), and decrease further  
231 with  $n$ -gram length (see Fig Figure 4—figure Supplement 3B and C). This reflects strong bias on  
232 the abundance of individual amino acids, and strong correlations between amino acids within the  
233 CDR3 which are observed in all CDR3 repertoires, and are captured by the frequency distribution  
234 of the longer  $n$ -grams. Two additional important points can be observed. First, the entropy of  
235 the repertoires after selection and lineage commitment (in the single positive populations) is  
236 less than the earlier pre-selection DP repertoires, which match closely the entropy of the pre-  
237 selection generative model (shown by the dotted line for each  $n$ -gram length). This decrease  
238 becomes more evident with longer  $n$ -gram length (the circles lie below the dotted lines). Thus, as  
239 predicted, selection does impose some decrease in repertoire diversity, although this is a much  
240 smaller effect than the decrease in diversity imposed by the generation process itself. The second  
241 key observation is that the entropy calculated directly from  $n$ -gram frequency in the data is very  
242 similar to that of the synthetic repertoires generated using the linear generative models in which  
243 individual TCR amino acids are treated as independent variables. Thus, at least at the level of  
244 diversity of  $n$ -grams, there is no evidence that selection at any step involves complex sequence  
245 motifs, or amino acid interactions, which would not be captured by the linear model. We looked  
246 in more detail at the  $n$ -gram ( $n = 3$ ) distributions derived by the linear selection models for the  
247 different maturation stages. A plot of the Jensen-Shannon divergences (JSD) between all pairwise  
248 comparisons largely recapitulated the expected relationships between the subsets, with DP pre  
249 and DP dbn clustering with the pre-selection generative model, while the single positive CD4 and  
250 CD8 populations clustered separately, and DP pos have an intermediate position (Figure 4C). A  
251 comparison for both TCR $\alpha$  and  $\beta$  for different  $n$  is shown in Figure 4—figure Supplement 4). Since  
252 some differences between populations were seen even for amino acid usage ( $n = 1$ ), we compared  
253 the discriminatory power of models based on  $n$ -grams with  $n = 1$  and  $n = 3$  (Figure 4—figure  
254 Supplement 5). The 3-gram model outperformed the 1-gram model in almost all cases. We can go  
255 beyond  $n$ -grams and use the subset-specific  $P_{\text{stage}}$  models to predict the entropy of the full sequence  
256 (Materials and Methods), shown in Figure 4—figure Supplement 6A for  $\alpha$  and Figure 4—figure  
257 Supplement 6C for  $\beta$ . This entropy is substantially reduced from generation to the DP stages, and  
258 further reduced in the single positive stages, especially in CD4<sup>+</sup> subsets. We also computed the JSD  
259 of the distributions  $P_{\text{stage}}$  between subtypes (Figure 4C for TCR $\alpha$  and Figure 4—figure Supplement 6D  
260 for TCR $\beta$ ). These JSD showed similar patterns as with  $n$ -grams, except for CD8<sup>+</sup> spleen cells showing  
261 more similarity to the  $P_{\text{gen}}$  distribution in TCR $\beta$ . Note that the absolute values of the entropies and  
262 JSD are larger, since they include information about longer sequences, with additional V and J gene  
263 usage information. In summary, we fitted the data with a set of stage-specific generative models  
264 based on linear weighted combinations of TCR sequence features. The repertoires generated by  
265 this model accurately estimate the sequence and  $n$ -gram entropy derived directly from the data,  
266 and generate repertoires which differ in a small but reproducible manner from each other. The  
267 magnitude of these differences reflect the expected developmental relationships between the  
268 different populations.





**Figure 4.** *n*-gram frequency discriminates between repertoires. **(A)** *n*-gram definition. We count how many times *n*-gram amino acid subsequences are seen in the CDR3 across a repertoire. **(B)** Shannon entropy *S* of the *n*-gram distributions normalized by *n* for the maturation stages. The entropy is estimated with the Nemenmann-Shafee-Bialek *Nemenman et al. (2002)* estimator and it is expressed in bits. The error on the estimated Shannon entropy from data is estimated from the sequencing error (see Materials and Methods). **(C)** Clustering according to Jensen-Shannon divergence between the 3-gram distributions computed from the selection model  $P_{\text{stage}}$  on synthetic repertoires. Dendrogram are computed with the Ward method (see Materials and Methods). **(D)** Clustering based on Jensen-Shannon divergence for the full  $P_{\text{stage}}$  selection model using  $P_{\text{stage}}$ .

**Figure 4—source code 1.** [https://github.com/statbiophys/thymic\\_development\\_2022/blob/main/fig4.ipynb](https://github.com/statbiophys/thymic_development_2022/blob/main/fig4.ipynb)

**Figure 4—figure supplement 1.** Comparison of different entropy estimators and of the dependence on the CDR3aa length choice.

**Figure 4—figure supplement 2.** Convergence of the *n*-gram entropy estimations.

**Figure 4—figure supplement 3.** Shannon entropy on  $\beta$  *n*-grams and entropy dependency on *n*.

**Figure 4—figure supplement 4.** Jensen-Shannon divergence between *n*-gram distributions.

**Figure 4—figure supplement 5.** AUC values computed from the ROC curves of the linear classifiers learnt over *n*-grams features.

**Figure 4—figure supplement 6.** Measure of the Shannon entropy using the full stage models.

**Figure 4—figure supplement 7.** Logo plots for the relative enrichment of positional amino acid usage.

**Figure 4—figure supplement 8.** Hydrophobic score at different stages and AUC scores of classifiers on hydrophobic features.

## 269 **Models capture modulations of hydrophobic residues in different subpopulations**

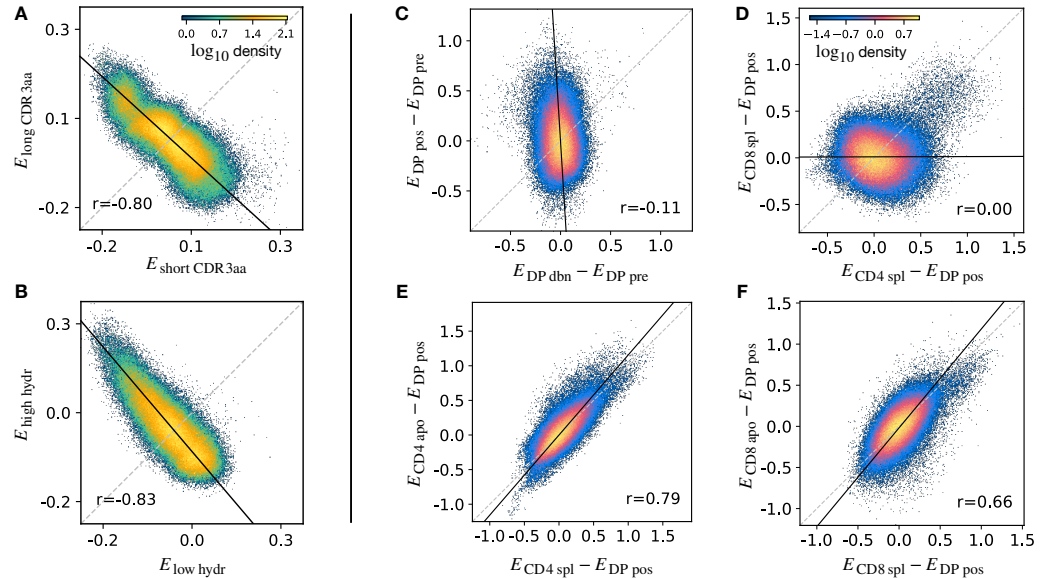
270 We inspected single amino acid usage in terms of the model marginals to check for relative positional  
271 enrichments between pairs of repertoires (Eq. (3)), but we did not observe any striking signal for  
272 amino acid charge properties. The logo plots with a visualization of this results are shown in  
273 **Figure 4—figure Supplement 7**. Hydrophobic residues in the central positions of the CDR3 have  
274 been reported to be enriched in the TCRs of regulatory versus conventional T cells *Lagattuta et al.*  
275 **(2022)**. This suggests hydrophobicity may function as a proxy for auto-reactivity, and might be  
276 enriched in cells selected for negative selection *Stadinski et al. (2016)*; *Daley et al. (2019)*. To test  
277 this idea, we defined a stage-specific hydrophobicity score  $U$ , obtained by summing the enrichment  
278 factors of hydrophobic residues CFILMWY at central positions of the CDR3 as learnt by our model  
279 at each stage (see Eq. (4) in Materials and Methods).

280 We observe a clear increase of this score from DP pre to DP pos, suggesting that positive  
281 selection introduces a bias towards more hydrophobic TCRs (**Figure 4—figure Supplement 8A** and  
282 B). The score also decreases in the single positive sets (CD4 and CD8), in agreement with the known  
283 role of negative selection to prune too strongly self-reactive T cells *Butler et al. (2013)*. Finally,  
284 AnnexinV+ single positive sets ('apo') show a slightly higher score than their respective spleen ('spl')  
285 sets (with the exception of the CD8  $\alpha$  chain scores). Overall, these changes in hydrophobicity are  
286 consistent with the hypothesised position of the different populations defined in our study in the  
287 stages of TCR selection.

288 Note that this score (like other scores found in the literature *Isacchini et al. (2021)*; *Lagattuta*  
289 *et al. (2022)*) is statistical and can not be used to classify individual sequences. To assess how much  
290 of single-sequence discriminability is explained by the presence of hydrophobic residues, we then  
291 introduced an empirical "hydrophobicity index"  $u$ , here defined as the number of hydrophobic  
292 residues (again CFILMWY) contained in the CDR3, normalized by its amino acid length (see Materials  
293 and Methods, Eq. 5). The classifiers using this feature yielded poor performance (**Figure 4—figure**  
294 **Supplement 8C** and D), worse than the 1-gram models (**Figure 4—figure Supplement 5A** and B).

## 295 **Discriminatory power of thymic selection**

296 The stage-specific enrichment factors in the generative models described above can be considered  
297 as capturing the combination of features which drives a particular selection step. A prediction of  
298 this idea is that, at each selection point, the TCRs which are selected and those which are not would  
299 have a distribution of model probabilities ( $P_{\text{stage}}$ ) which are anti-correlated. For example, a TCR that  
300 is present in the DP pos repertoire but "forbidden" from the CD4 repertoire (e.g. because of cross-  
301 reactivity to a Class II self pMHC) would be expected to have a large positive  $P_{\text{DP pos}}$  and a  $P_{\text{CD4 spl}} \approx 0$ ,  
302 reflecting the large enrichment factor between these two populations. A toy example illustrating this  
303 idea is illustrated in (**Figure 5A**). We consider a simple model in which TCRs are selected according  
304 to their CDR3 length into a "long" population with probability  $P(\text{long}|L) = L^h / (L^h + L_0^h)$  and into a  
305 "short" population otherwise. We apply this selection process *in silico* to  $P_{\text{gen}}$ -generated TCRs, and  
306 fit a separate  $P_{\text{population}}$  model on the synthetic sequences found in each subset. We then calculate  
307  $E_{\text{population}}$  for each TCR according to both subset models, and plot these values against each other.  
308 The distribution of enrichment strengths according to the two models are clearly anti-correlated  
309 (**Figure 5A**). In other words, if a TCR is more likely to be classified as a "long" sequence, it is in  
310 general less likely to be classified as a "short" one. Interestingly, however, the enrichment strengths  
311 distributions from the two models are significantly overlapping. As a result, attempts to classify  
312 individual TCRs according to their enrichment strengths is poor,  $\text{AUC} \sim 0.57$ . We then consider a  
313 different toy model where TCRs are selected according to the empirical hydrophobicity index  $u$  (Eq. 5).  
314 Similarly, we choose to generate a synthetic "high hydrophobicity" (HH) population filtering  $P_{\text{gen}}$ -  
315 generated TCRs with a probability  $P(\text{HH}|u) = u^h / (u^h + u_0^h)$ , otherwise "low hydrophobicity". Repeating  
316 the analysis performed with the length example, we observe in the corresponding scatterplot that  
317 enrichment strengths are again anti-correlated (**Figure 5B**).



**Figure 5.** Density scatter-plots of TCRa sequences comparing the selection energies learnt at two different stages. **(A)** Synthetic example of soft discrimination between “short” and “long” CDR3, where sequences are randomly assigned into either of the two populations with a bias that depends on their CDR3 length. The density scatter plot shows a clear anti-correlation between the selection energies learnt from these two populations. Yet, sequence classification is imprecise, as quantified by the low AUC=0.57. The parameters chosen for the filter in this example are  $L_0 = 13$  and  $h = 2$ . **(B)** Synthetic example of soft discrimination between “low” and “high hydrophobic” CDR3 showing clear anti-correlation between these two populations. Sequence classification is again poor AUC=0.60. The parameters chosen for the filter on the “hydrophobic index”  $u$  in this example are  $u_0 \simeq 0.2$  (the median value over a set of  $P_{\text{gen}}$ -distributed sequences) and  $h = 1$ . **(C)** The differential enrichment parameter of each TCR calculated according to  $P_{\text{DPdbn}}$  model is plotted against the energy calculated against the  $P_{\text{DPpos}}$  model. To correct for bias imposed by the TCRa generation process, the DP pre energy, which encodes background selection common to both stages, is subtracted. The black line is the direction of the major eigenvector of the dots moments matrix. The value  $r$  reported in each plot is the Pearson’s correlation coefficient (see Materials and Methods). **(D)** Differential enrichment parameter according to CD4 spl and CD8 spl models, relative to DP pos. **(E)** Differential enrichment parameter according to CD4 spl and CDd apo models, relative to DP pos. **(F)** Differential enrichment parameter according to CD8 spl and CD8 apo models, relative to DP pos.

**Figure 5—source code 1.** [https://github.com/statbiophys/thymic\\_development\\_2022/blob/main/fig5.ipynb](https://github.com/statbiophys/thymic_development_2022/blob/main/fig5.ipynb)

**Figure 5—figure supplement 1.** Differential increments scatterplots for all pairs of stages.

318 We extended this approach to look for relationships between enrichment strengths for TCRs at  
319 different developmental stages. Since all cells pass through a preceding selection stage, we must  
320 consider it as a common background distribution for all the successive thymic stages. We therefore  
321 considered the differential enrichment parameter  $E_{\text{stage}} - E_{\text{pre-stage}}$ , a linear operator which predicts  
322 whether a sequence is more or less likely to be present in a particular developmental stage as  
323 compared to the previous stage. We generated a set of sequences using the generation model  $P_{\text{gen}}$   
324 (thus with no selection bias), and then computed differential enrichment parameters for each TCR  
325 according to all the stage specific models. The full set of pairwise correlations between enrichment  
326 values for the different populations relative to  $P_{\text{DP pre}}$  are shown in (**Figure 5—figure Supplement 1A**  
327 for TCR $\alpha$  and **Figure 5—figure Supplement 1C** for TCR $\beta$ ). The DP dbn repertoire showed a narrow  
328 distribution of values, which was uncorrelated to any other subset, in particular to DP pos (**Figure 5C**).  
329 This would be consistent with the DP dbn repertoire containing a random sample of the DP pre  
330 repertoire, unrelated to its TCR sequence. To check if the signal coming from DP pos stage is the  
331 principal cause of the high correlation between the single positive stages, we repeated the analysis  
332 for CD4 and CD8 using  $P_{\text{DP pos}}$  as the common background distribution (the full set of scatterplots  
333 for TCR $\alpha$  is shown in **Figure 5—figure Supplement 1B**, in **Figure 5—figure Supplement 1D** for TCR $\beta$ ).  
334 There was therefore no evidence of selection pressure operating on TCR sequence to distinguish  
335 these two populations. The correlation between the CD4<sup>+</sup> and CD8<sup>+</sup> subsets was negligible ( $r \sim 0$ ),  
336 suggesting that the selection pressures operating on the two populations are distinct (**Figure 5D**).  
337 The spleen SP and the thymic apo populations were also highly correlated for both CD4<sup>+</sup> and CD8<sup>+</sup>  
338 cells ( $r = 0.79$  for CD4 spl vs CD4 apo, in **Figure 5E**;  $r = 0.66$  for CD8 spl vs CD8 apo, in **Figure 5F**).  
339 Similar results are obtained for the sequences of the  $\beta$  chain (**Figure 5—figure Supplement 1C,D**).  
340 In contrast to the examples illustrated above, most plots showed a positive correlation between  
341 enrichment values for two models. Thus a common dominant selection process is driving the  
342 repertoire shift between the DP pos and all subsequent stages, which dominates the impact of  
343 individual stage specific selection processes. In summary, the TCR enrichment value distributions  
344 differ between different thymic populations, but do not show evidence of dominant exclusive  
345 sequence-based selection operating at any step of the selection process.

## 346 Discussion

347 Thymic selection is often portrayed as a simple discrimination process that eliminates TCRs capable  
348 of strongly binding any self-peptide, while promoting TCRs that bind them weakly. However, this  
349 simple picture has been challenged and the fidelity of the negative selection process and the  
350 proportion of the self-repertoire which can effectively be scanned by each individual thymocyte  
351 during the window of negative selection remains incompletely understood *Yu et al. (2015); Gallegos*  
352 *and Bevan (2006)*. If significant number of T cells escape negative selection and enter the peripheral  
353 repertoire, no sequence feature will unambiguously distinguish TCRs from pre and post-selection  
354 repertoires. Many efforts have been made to connect TCR sequences to peptide recognition *Weber*  
355 *et al. (2021); Montemurro et al. (2021); Isacchini et al. (2021)*. However, these approaches cannot  
356 yet be used to define the target peptidome of entire repertoires. Here we take the complementary  
357 approach, by looking for TCR sequence features that are linked to thymic selection.

358 Although there has been a lot of work on understanding and modeling thymic development *Yates*  
359 *(2014); Robert et al. (2021)* our study presents the first comprehensive analysis of TCR repertoire of  
360 different developmental stages of thymic maturation. By incorporating a reporter for the activation  
361 marker Nur77, which is activated during thymic selection, and an early marker of apoptosis, Annexin  
362 V, we were able to enrich for identifying subpopulations during the process of positive or negative  
363 selection. Although this more sophisticated strategy in principal allows the unbiased isolation of the  
364 major stages of thymic selection, some limitations remain. For example, the time interval during  
365 which negatively selected cells survive after they received their instruction to go into apoptosis may  
366 depend on signal strength. If strong TCR signal strength translates into short subsequent lifetime,  
367 then the AnnexinV<sup>+</sup> cells sorted may be enriched for cells receiving a rather weaker negative signal.

368 We examined the repertoires from two perspectives. In the first part of the paper, we compare  
369 statistical properties of the sequences of the repertoires using features of different dimensionalities,  
370 which include V gene, J gene and CDR3 length frequency distributions, and individual CDR3 se-  
371 quences represented as sparse {0,1} binary vectors. The analysis incorporated both coarse-grained  
372 (V, J and CDR3 length) and fine-grained (individual CDR3 sequence) features, and the results were  
373 remarkably consistent. No single feature adequately discriminated between any pair of repertoires.  
374 Combination of features when averaged across a repertoire did show subtle but reproducible  
375 differences between repertoires, which could be used to discriminate between subpopulations  
376 using both unsupervised (PCA) and supervised (logistic regression) analysis. Furthermore, the  
377 difference between these statistical parameters captured the known developmental trajectory of  
378 thymic development, illustrated schematically in *Figure 1D*. Interestingly, the smallest distances  
379 observed were between mature CD4 or CD8 cells, and their thymic SP negatively selected (apo)  
380 counterparts. This suggests that negative selection of single positives is only weakly associated with  
381 the sequence properties of single TCRs, or at least single chains. It is in principle possible that larger  
382 differences exist in the paired  $\alpha$ - $\beta$  repertoires, which would not be detectable in either the alpha or  
383 beta repertoires alone, but previous work on the functional alpha-beta repertoire has suggested  
384 that pairing was largely random, with weak associations between some germline genes *Grigaityte*  
385 *et al. (2017)*; *Dupic et al. (2019)*.

386 An additional possibility which must be considered is that Annexin V staining does not exclusively  
387 capture the negatively selected population, but also identifies cells which were damaged during  
388 the preparation. Contamination of the AnnexinV+ population by these damaged or dying cells will  
389 weaken the selection signature observed, although the fact we do manage to discriminate between  
390 the apo and spleen subpopulations (*Figure 3A*) indicates that these differences, however small, do  
391 exist. Conversely, cells marked for deletion may not have the time to express Annexin V, so that the  
392 DPpos subset may contain cells that are being negatively selected against, in addition to cells that  
393 are being positively selected *Stritesky et al. (2013)*.

394 A limitation of our sorting strategy is that we do not identify Treg from conventional CD4+ T cells.  
395 It has been suggested that regulatory T cells (Tregs), which are more auto-reactive and should thus  
396 bear the same marks as the cells that fail negative selection, have distinctive TCR features, notably  
397 the presence of hydrophobic residues at key positions *Stadinski et al. (2016)*; *Daley et al. (2019)*.  
398 TCR scores based on more detailed features than hydrophobicity have been proposed *Isacchini*  
399 *et al. (2021)*; *Lagattuta et al. (2022)*. We note that these scores are statistical and do not classify  
400 individual sequences. Consistent with these previous results, we can project our model parameters  
401 to build a single hydrophobicity index, which we observe to be significantly increased in positively  
402 selected cells (DP pos) versus DP pre, and decreased in single positive sets (*figure Supplement 8*).  
403 Beyond hydrophobicity, it remains an open question whether the features that drive Treg fate are  
404 the same that drive negative selection.

405 Although the statistical properties of the repertoires differed between subpopulations, it was not  
406 possible to classify individual TCRs at high accuracy. As discussed above, this may in part be due to  
407 the fact that the populations we define only imperfectly correlate with their fate and self-reactivity.  
408 However, the differences between CD4+ and CD8+ repertoires, which are much less likely to be  
409 affected by issues of functional or physical cross-contamination, are also seen only at a statistical  
410 population level, and not an individual TCR sequence level. Learning the collective properties of at  
411 least a few dozen TCRs was required in order to achieve good discrimination between repertoires.

412 The statistical population-level differences between populations of thymocytes and mature T  
413 cells which we observe is reminiscent of previous models emphasising the importance of collective,  
414 rather than individual T cell behaviour. *Butler et al. (2013)* proposed that a minimum number of  
415 T-cells must collectively recognize a peptide to trigger a response, proposing quorum sensing as a  
416 mechanistic explanation of this collective decision making. Recent experiments have confirmed that  
417 quorum sensing between TCRs can occur, mediated via cytokine signaling *Polonsky et al. (2018)*, and  
418 estimating a minimum quorum size of activated T cells to be  $\geq 30$ . Our results suggest that thymic

419 selection imposes only a rather weak selective pressure on the repertoire, which is consistent with  
420 *Butler et al. (2013)*'s hypothesis that most self-peptides are not screened by TCR during negative  
421 selection. Our results are consistent with their model, in which even a subtle depletion, rather than  
422 complete elimination of non-self TCRs, may still translate into robust self/non-self discrimination  
423 in populations of reactive TCRs. Self versus foreign peptide discrimination by TCR is somewhat  
424 the conjugate task of self-reactive versus a non self-reactive T cell discrimination during negative  
425 selection. While the performance of the two tasks cannot be directly compared at first sight, they are  
426 related in that both are impaired by a factor  $(1 - f)$  due to partial screening of self peptides, where  
427  $f$  is the fraction of self-peptides that are presented during thymic development. The common  
428 point is that even when  $f$  is small, the law of large numbers can rescue the discrimination task  
429 when there are multiple observations. In Materials in Methods, we argue using the model of *Butler*  
430 *et al. (2013)* how the idealized performance of repertoire discrimination using multiple ( $m$ ) TCR (akin  
431 to the task of Fig. 3B) may be compared to the task of telling self from foreign peptides in the  
432 periphery, when the number of T cells specific to one particular peptide and recruited to the site of  
433 infection is  $m \times \bar{n}$ , where  $\bar{n}$  is the average number of self-peptides recognized by a random TCR. While  
434 those numbers cannot be applied directly to the results of Fig. 3B, which are based on an imperfect  
435 classifier from a single chain, they give a sense of how the same principle of discrimination apply to  
436 both cases.

437 In the second part of the study we explore in more detail whether we can discover any evidence  
438 that thymic selection depends on specific sequence motifs (i.e. a strong correlative structure be-  
439 tween CDR3 amino acids). For this purpose, we build on our previous work which have established  
440 a framework for the development of generative statistical models of repertoire generation, based  
441 firmly on a mechanistic understanding of TCR generation and selection. Specifically, we construct  
442 models which incorporate only linear combinations of CDR3 sequences to capture the selective  
443 process which can transform one repertoire into another. These models produce an "enrichment  
444 factor" for each TCR which estimates its relative likelihood of being in a particular stage-specific  
445 population. Intuitively, one can consider these factors as capturing the probable enrichment or  
446 depletion of a TCR with a particular sequence when comparing two repertoires. We demonstrate  
447 that these linear models effectively capture the progressive decrease in repertoire diversity which  
448 we observe in the preselected DP to the SP transition. They also effectively capture the known  
449 developmental relationships between the thymic subpopulations. Thus we find no evidence that  
450 complex non-linear amino acid sequence interactions are required to explain the observed changes  
451 in repertoire observed in our data. We also compared the distributions of enrichment factors  
452 between populations. We demonstrate that, contrary to the predictions of a strong binary selection  
453 model, we do not observe any negative correlation between enrichment factor distributions be-  
454 tween selected and non-selected repertoires. Instead, we observe a set of positive correlations,  
455 revealing a dominant conserved selection process spanning the developmental stages between  
456 pre-selection DP and mature SP. Consistent with the clustering data discussed above, we find  
457 strong correlation between the enrichment factor distributions of mature SP and thymic negatively  
458 selected population, and no evidence of binary selection between these two populations.

459 In conclusion, we report a comprehensive analysis of the TCR repertoire at various stages of  
460 thymic development. We then combine data-driven and model-based analysis of these repertoires.  
461 Our conclusions are incompatible with a model of thymic developments which involves a sequence  
462 of clear-cut binary selection processes, based on TCR sequence features. Rather, our data suggest  
463 a probabilistic fuzzy decision making process at each selection step. We propose that this model is  
464 compatible with robust self/non-self discrimination, if T cell responses to antigen are governed by  
465 collective quorum based decision making. Further experimental and theoretical work is required to  
466 test these hypotheses, which have fundamental implications for strategies to modulate the immune  
467 response for prophylaxis or therapy of human disease.



## 468 **Methods**

### 469 **Animals**

470 The experiment was carried out using three 6-weeks old male inbred Nur77-GFP/Foxp3-mCherry  
471 (C57BL/6 background) *Moran et al. (2011)*. The cross was bred and maintained at the Weizmann  
472 institute. This study was performed in strict accordance with the recommendations in the Guide  
473 for the Care and Use of Laboratory Animals of the National Institutes of Health. All of the animals  
474 were handled according to approved institutional animal care and use committee (IACUC) protocols  
475 (#21661115-2) of the Weizmann Institute of Science. The protocol was approved by the Committee  
476 on the Ethics of Animal Experiments of the Weizmann Institute of Science. Every effort was made to  
477 minimize suffering.

### 478 **Sample preparation and T cell isolation**

479 Thymocytes and splenocytes were isolated from Nur77-GFP/Foxp3-mCherry 6-weeks old mice.  
480 Erythrocytes were removed by hypotonic lysis in ammonium chloride. Thymocytes were stained  
481 with fluorescent antibodies, and sorted using a flow cytometer as described below. Splenic CD4 and  
482 CD8 cells were purified in two steps: (1) CD4+ positive selection (CD4 (L3T4) MicroBeads, mouse, #  
483 130-117-043, Miltenyi) to generate the "CD4 spl" samples (2) the negative cells fraction were further  
484 selected for the CD8+ positive cells (CD8a (Ly-2) MicroBeads, mouse, # 130-117-044, Miltenyi Biotec)  
485 to generate "CD8 spl" samples.

### 486 **Flow cytometry analysis and cells sorting**

487 The following fluorochrome-labeled mouse antibodies were used according to the manufacturers'  
488 protocols: PerCP/Cy5.5 anti-CD4, PB anti-CD8, PE/cy7 anti-CD3, APC annexinV (Biolegend). UV  
489 LIVE/DEAD™ (ThermoFisher Scientific, # L23105). Labelled cells were sorted on a SORP-FACS-Ariall  
490 using a 70 µm nozzle to 5 populations (see Table 1). Cell counts are reported in Table S2. Cells were  
analyzed using *FlowJo* (Tree Star) software.

Sample\Marker	CD4	CD8	CD3	AnnexinV	Nur77
DP pre	+	+	+	-	-
DP pos	+	+	+	-	+
DP dbn	+	+	+	+	-
CD4 apo	+	-	+	+	
CD8 apo	-	+	+	+	

**Table 1.** Cell sorting based on fluorochrome-labeled mouse antibodies.

491

### 492 **Library preparation for TCR-seq**

493 All libraries in this work were prepared according to the published method *Oakes et al. (2017)*, with  
494 minor adaptations as described below. Briefly, total RNA was extracted from each of the seven  
495 populations using RNeasy Micro Kit (# 74004, Qiagen) and cleaned from excess DNA with DNase 1  
496 enzyme (# M6101, Promega). RNA samples were reverse transcribed to cDNA (SuperScript™ III, #  
497 12574026, Invitrogen) using primers for the mouse α chain (mAlpha\_RC2) and for the mouse beta  
498 chain (mBeta\_RC2) (see Table S1). Following reverse transcription the samples were purified on  
499 minielute spin columns (# 28004, QIAGEN). The cDNA was ligated to an oligonucleotide containing  
500 a unique 12 basepair molecular identifier (UMI) (6N\_I8.1\_6N\_I8.1\_SP2, see Table S1) using T4 RNA  
501 ligase (M0204S, NEB). Ligation products were purified using Agencourt AMPure XP beads (# A63881,  
502 BeckmanCoulter). Next, three rounds of extension PCR were executed (using KAPA HiFi DNA  
503 Polymerase, KAPA Biosystems) to add illumina sequencing adaptors and Illumina sample indices  
504 for multiplex sequencing (see Table S2). The thermal cycler parameters are an initial denaturation  
505 step (3 minutes at 95°C) followed by cycles of denaturation (98°C for 20 seconds), annealing (61°C

506 for 15 seconds), and extension 72°C for 30 seconds. The final extension step was at 72°C for  
507 five minutes. The lid was maintained at 105°C. After the first round PCR (5 cycles), PCR products  
508 were purified using Agencourt AMPure XP beads and split in two, and  $\alpha$  and  $\beta$  TCR genes were  
509 processed separately in subsequent steps. After the second PCR (8 cycles), PCR products were again  
510 purified using Agencourt AMPure XP beads. The final amplification using the adapter sequences P5  
511 and P7 were carried out on a real-time qPCR machine, and the amplification was tracked by the  
512 incorporation of SYBR green. The cyclor was stopped manually when the fluorescent signal reached  
513 a predetermined threshold, thus preventing overamplification.

514 The final library concentration was measured using Qubit Fluorometric Quantification (Ther-  
515 moFisher Scientific) and the presence of the correct 600-700 bp product confirmed by electrophore-  
516 sis on a High Sensitivity D1000 ScreenTape cassette using a 4200 TapeStation System (Agilent).  
517 Multiple samples were pooled in equal molarity, and then sequenced using NextSeq 550 (200 bp  
518 forward read, 100 bp reverse) (Illumina).

### 519 **Pre-Processing and Error Correction for Raw Reads**

520 Data were processed using an in-house pipeline, coded in R. First, UMI sequences were transferred  
521 from read 2 to read 1. Trimmomatic was used to filter out the raw reads containing bases with  
522 Q-value  $\leq 20$  and trim reads containing adaptors sequences *Bolger et al. (2014)*. The remaining  
523 reads were separated according to their barcodes and reads containing the constant region for  
524  $\alpha$  or  $\beta$  chain primers sequences were filtered (CAGCAGTTCTGGTTCTGGATG / TGGGTGGAGT-  
525 CACATTTCTCAGATCCT  $\alpha$  and  $\beta$  chain, respectively), allowing up to three mismatches. To correct  
526 for possible sequence errors, we cluster the sequences UMIs' in two steps; (1) The UMIs with the  
527 highest frequency are grouped within a Levenshtein distance of 1 *Levenshtein et al. (1966)*. (2) Out  
528 of these sequences, CDR3AA sequences (starting from the most frequent sequence in a group)  
529 were clustered using a Hamming distance threshold of 4 *Hamming (1950)*. Finally, the UMI of each  
530 CDR3 sequence was counted.

### 531 **Annotation and Generation Model**

532 From the raw nucleotide reads, we performed a preliminary annotation using the python module  
533 *PyIR* (version 1.3.0) *Soto et al. (2020)*, which provides a wrapper and parser of the open source  
534 software *IgBlast* *Ye et al. (2013)*. We then separated the productive clonotypes from the out of  
535 frame reads and/or reads containing stop codons. We define a clonotype as TCRs sharing V genes,  
536 J genes, and the same CDR3 nucleotide sequence. If different reads are annotated as the same  
537 clonotype in the same dataset, only the read with highest UMI counts is considered.

538 For our models, we use a reduced set of genes from the IMGT free online repository *Lefranc*  
539 *et al. (2015)* in order to have a single allele per gene, preferring functional alleles to open reading  
540 frame or pseudo genes. A further reduction is done for the V genes of the  $\alpha$  chain, clustering to  
541 a single representative all of the those genes that result indistinguishable in the region from the  
542 maximum observed V offset for the annotation to the conserved cysteine. Two genes are said to  
543 be indistinguishable if the Hamming distance *Hamming (1950)* between the considered regions is  
544 equal to 0. For each TRAV cluster, we choose as the representative the most frequent gene in the  
545 preliminary annotation. In this way we obtain 76 V genes and 51 J genes for the  $\alpha$  chain, 26 V genes,  
546 2 D genes and 14 J genes for the  $\beta$  chain.

547 In order to infer a generation model we use the open source software *IGoR* *Marcou et al. (2018)*  
548 on all out-of-frame clonotypes pooled from all maturation stages of all mice. The generation  
549 model associates to each  $\alpha$  ( $\beta$ ) read a probability  $P_{\text{gen}}$  of being generated through the VJ (VDJ)  
550 recombination process. After learning a generation model, we annotate the reads using the most  
551 probable alignment scenario using the *IGoR* software, as the clonotype (V, J gene choice, CDR3  
552 nucleotide sequence) with the highest  $P_{\text{gen}}$  among all possible recombination scenarios.

553 The PCA was computed in *R* (version 3.6.0) using the function "PCA" from the *FactoMineR* package  
554 (version 2.4).

## 555 **Statistical Classification**

556 The features are assigned to each  $\alpha$  chain as a binary vector  $\vec{\sigma}$ , where each entry is equal to 1  
557 if the feature is observed, 0 otherwise. In this study the set of features is encoded using the  
558 “SoniaLeftposRightpos” class (from the Python package *Sonia* version 0.0.45) which provides 5033  
559 features: 30 for the CDR3 amino acid lengths, 25 left to right positions for each of the 20 amino  
560 acids (500 features), 25 right to left positions (500), the joint V/J gene usage ( $76 \times 51=3876$ ) and  
561 the independent usage ( $76+51=127$ ). Analogously for the  $\beta$  we obtain 1434 features (without  
562 considering D genes).

563 To learn the models for the statistical classification of two stages, we first remove all sequences  
564 that share the same features between the two sets (i.e. same amino acid CDR3, V and J gene). Then,  
565 we balance the size of the sets sub-sampling the larger one so that its size does not exceed 25%  
566 of the size of the smaller. Each of the resulting sets is divided into a train and a test set by a ratio  
567 70%/30% (“StratifiedShuffleSplit”, module “model\_selection” from the Python package *scikit-learn*,  
568 version 0.24.2). The classifiers are learned with linear models, defined by a single layer with binary  
569 cross entropy as a loss function, binary accuracy as metrics, a sigmoid as activation function, coded  
570 using the “keras” module from the Python package *tensorflow* (version 2.4.1). We obtained similar  
571 performance for the classification task by learning with a random forest algorithm as provided by the  
572 function “RandomForestModel” in the module “keras” from the package *tensorflow\_decision\_forests*  
573 (version 0.2.4).

## 574 **Selection Model**

575 To learn a  $P_{\text{stage}}$  selection model for each maturation stage, we pooled together the annotated  
576 sequences from all mice for the given maturation stage, discarding all clonotypes annotated with  
577 non-functional and pseudo genes. We learn a selection model using the open source software  
578 *Sonia* for each maturation stage. *Sonia* performs a linear regression over the features of the  
579 sequences in the dataset to infer the enrichment ratio between the maturation specific dataset and  
580 the generation model. The feature choice for the enrichment model is similar, except for the fact  
581 that only independent gene usage is considered, reducing features to 1157 for  $\alpha$  chain (1070 for  $\beta$   
582 chain). The probability of observing a sequence in a stage is modeled as

$$P_{\text{stage}}(\vec{\sigma}) = \frac{1}{Z} e^{-E_{\text{stage}}(\vec{\sigma})} P_{\text{gen}}(\vec{\sigma}) \quad (1)$$

583 where  $Z$  is a normalization factor and the energy  $E_{\text{stage}}(\vec{\sigma})$  for a sequence showing a set of features  
584  $\mathcal{F}(\vec{\sigma})$  is defined as

$$E_{\text{stage}}(\vec{\sigma}) = \sum_{f \in \mathcal{F}(\vec{\sigma})} \epsilon_{\text{stage}}^{(f)}(\vec{\sigma}) \quad (2)$$

585 Here  $\epsilon^{(f)}$  is a weight associated to the feature  $f$  and is learnt from data. To look at specific enhanced  
586 features between stages  $a$  and  $b$  one can obtain the average weights difference from the respective  
587  $P_{\text{stage}}$  models as

$$\langle \epsilon_a^{(f)} - \epsilon_b^{(f)} \rangle = \frac{p_a^{(f)} + p_b^{(f)}}{2} \cdot (\epsilon_a^{(f)} - \epsilon_b^{(f)}) \quad (3)$$

588 where  $p_{\text{stage}}^{(f)}$  is the marginal associated by the model to the feature.

589 The limited amount of clonotypes for certain maturation stages precludes using deep neural  
590 network based selection models, although we do not expect the conclusions to change with the  
591 DNN SoNNia model *Isacchini et al. (2021)*.

## 592 **Hydrophobicity Score**

593 To study the hydrophobicity increase with respect to the generation, we define a stage-wide score  
594 as

$$U = \sum_{\substack{a \in \text{hydro} \\ x \in \text{CDR3cr}}} \epsilon^{(a|x)} \cdot p_{\text{gen}}^{(a|x)} \quad (4)$$

595 where  $\epsilon^{(a|x)}$  is the weight associated by the model to the amino acid  $a$  at position  $x$ ; the marginal  
 596  $p_{\text{gen}}^{(a|x)}$  is obtained by the generation model on the same feature (see Materials and Methods). The  
 597 sum runs over the hydrophobic amino acids CFILMWY, following the definition from *Lagattuta*  
 598 *et al. (2022)*, considering just the positions of our model which correspond to the central region  
 599 p108-p114 of the CDR3 in IMGT convention (model positions (4:10) from the left, and, (-11:-5) from  
 600 the right). We also define an index  $u$  for hydrophobicity which can now be associated to each  
 601 sequence as follow

$$u = \frac{1}{L} \sum_{\substack{a \in \text{hydro} \\ x \in \text{CDR3cr}}} 1 \quad (5)$$

602 i.e. the number of hydrophobic residues found in the central region (same choices as above),  
 603 normalized by the CDR3 length  $L$ .

#### 604 ***n*-gram Shannon Entropy Estimation**

605 As a diversity measure we consider the Shannon entropy defined as :

$$S = - \sum_i p(i) \log_2 p(i) \quad (6)$$

606 where  $p(i)$  is the probability of finding a clonotype in the data. Since  $n$ -grams are sampled from  $20^n$   
 607 possible motifs, undersampling could bias a naive estimation of the entropy. We overcome this bias  
 608 by estimating the Shannon entropy using the Nemenman-Shafee-Bialek (NSB) estimator *Nemenman*  
 609 *et al. (2002)*. The NSB estimator is computationally tractable and calculates an estimation error.  
 610 We implement the entropy and variance estimators as given in *Archer et al. (2014)*. To check for  
 611 convergence we subsample the clonotypes in the dataset at increasing sizes and estimate the  
 612 entropy for each sub-sample (*figure Supplement 2*). Convergence sets a limit of  $n = 4$  due to  
 613 sample size constraints of the smallest dataset. We repeat the same computation for synthetic  
 614 repertoires. We verified the NSB estimators better performance for our datasets compared to other  
 615 non-parametric estimators (*figure Supplement 1A and B*), consistently with previous reports *Archer*  
 616 *et al. (2014)*.

#### 617 **Full Model Shannon Entropy Estimation**

618 The Shannon entropy in Eq. 6 associated to the full  $p = P_{\text{stage}}(\vec{\sigma})$  model requires summing over all  
 619 possible clonotypes  $i = \vec{\sigma}$ . Practically we evaluate the entropy by producing synthetic sequences  
 620 according to the selection model  $P_{\text{stage}}$  and averaging the value of  $\log_2 P_{\text{stage}}$

$$S(P_{\text{stage}}) \simeq \frac{1}{N} \sum_{k=1}^N \log_2 P_{\text{stage}}(\vec{\sigma}_k^*) \quad (7)$$

621 with clonotypes  $\vec{\sigma}_k^*$  sampled from the  $P_{\text{stage}}$  distribution.

622 Because of sequencing errors, the entropy of  $n$ -grams is systematically overestimated in the  
 623 data. To estimate and correct for this bias, we measured the error rate from the data, provided  
 624 as a byproduct of the IGoR training procedure *Marcou et al. (2018)*. We used this rate to produce  
 625 synthetic sequences with simulated sequencing errors. The difference in  $n$ -gram entropy between  
 626 error-prone and error-free sequences was then applied as a subtractive correction factor to the  
 627 data.

#### 628 ***n*-gram Jensen-Shannon Divergence**

To quantify the distance between two distributions  $p_a$  and  $p_b$  defined on the same support, we use  
 the symmetric Jensen-Shannon divergence JSD:

$$\text{JSD}(p_a, p_b) = \frac{1}{2} \sum_i p_a(i) \log_2 \frac{2p_a(i)}{p_a(i) + p_b(i)} + \quad (8)$$

$$+ (a \leftrightarrow b)$$

629 where the sum runs over all possible observables  $i$  and the term  $(a \leftrightarrow b)$  corresponds to the same  
 630 expression in the first one with  $a$  and  $b$  inverted. The Jensen-Shannon divergence is bounded  
 631 between 0 and 1 bits, with  $\text{JSD} = 0$  bits if the distributions are identical and a maximal difference  
 632 of  $\text{JSD} = 1$  bit. We use JSD to assess the divergence between  $n$ -gram distributions and between  
 633 selection models.

### 634 Full Model Jensen-Shannon Divergence

To compare selection models of complete clonotypes at two maturation stages, the divergence between the  $P_{\text{stage}}$  distribution of model  $a$  and the model  $b$  is:

$$\text{JSD}(P_a, P_b) \simeq \frac{1}{2N} \sum_{k=1}^N \log_2 \frac{2e^{-E_a(\bar{\sigma}_k^a)}}{e^{-E_a(\bar{\sigma}_k^a)} + e^{-E_b(\bar{\sigma}_k^a)}} + (a \leftrightarrow b) \quad (9)$$

635 where the clonotypes  $\bar{\sigma}_k^a$  are sampled from the  $P_a$  distribution. In Eq. 9, we used the fact that a given  
 636 sequence has the same background generation probability  $P_{\text{gen}}$  in both selection models.

### 637 Discrimination in the thymus vs discrimination in the periphery

638 Here we show a formal link between discrimination of negatively selected vs non-negatively selected  
 639 TCR on the one hand, and foreign vs self-peptide recognition on the other.

640 We start by considering (negative) thymic selection. Following *Butler et al. (2013)*, we assume  
 641 that a random TCR will recognize any peptide with probability  $p$ . Then the number of recognized  
 642 self-peptides  $n$  by a random TCR is distributed according to a Poisson law of mean  $\bar{n} = pN$ , where  $N$   
 643 is the number of self-peptides,  $P(n) = \text{Poiss}[pN](n) \equiv e^{-pN}(pN)^n/n!$ .

644 If each TCR only screens  $M = fN$  self-peptides, with  $f < 1$ , then the probability of passing  
 645 selection (and ending up in spleen) is  $P(\text{spleen}|n) = (1 - n/N)^M \approx e^{-nf}$ , and  $P(\text{apo}|n) = 1 - e^{-nf}$  for the  
 646 probability of ending up in apo (as apoptosis, i.e. single-positive cells expressing Annexin V as in our  
 647 experiments).

We assume that the discriminator of apo vs spleen single positives is perfect, in the sense it can perfectly deduce  $n$  from the TCR sequence. In this idealized case, discrimination errors are entirely attributable to the partial screening of self-peptides. Using Bayes' law, one can show that the distributions of  $n$  in spleen and apo read:

$$P(n|\text{spleen}) = \frac{P(\text{spleen}|n)P(n)}{P(\text{spleen})} = \text{Poiss}[pN(1-f)](n), \quad (10)$$

$$P(n|\text{apo}) = \frac{P(\text{apo}|n)P(n)}{P(\text{apo})} = \frac{(1 - e^{-nM/N})(pN)^n e^{-pN}}{n!(1 - e^{-pM})} \approx \frac{nf}{pNf} \frac{(pN)^n e^{-pN}}{n!} = \text{Poiss}[pN](n-1), \quad (11)$$

648 where the first equation results from direct algebra, and the second is obtained in the limit of  
 649 small  $f$ . The AUC of the discrimination task is then given by the probability that drawing a random  
 650 number from a Poisson of mean  $\bar{n}(1-f)$  yields a smaller number than drawing a random number  
 651 from a Poisson of mean  $\bar{n}$ , and adding 1 to it. If we now use the observation of  $m$  TCRs from the  
 652 same subset (apo or spleen) instead of a single one, the task becomes easier: We can form a  
 653 collective score by adding up the  $n$ 's of each TCR (since they are independent draws from either  
 654 the apo or spleen ensembles) so that the two Poisson distributions, of respective means  $m\bar{n}(1-f)$   
 655 and  $m\bar{n}$ , become better separated. This is qualitatively the result of Fig. 3B, which is based on the  
 656 learned score, rather than on an idealized one.

657 We now turn to the case of self vs foreign peptide discrimination by a group of  $R$  T cells recruited  
 658 to a site of infection. If the peptide is from the self, then the probability for a given circulating TCR to  
 659 recognize it is  $p(1-f)$  (*Butler et al., 2013*). Then the number of specific TCR is Poisson distributed  
 660 with mean  $p(1-f)R$ . If the peptide is foreign, that number is also Poisson distributed, but with  
 661 mean  $pR$ . Again, the AUC of the discrimination task is given by the probability of drawing a smaller  
 662 number from the former distribution than from the latter. This task is expected to be at least as

663 hard as that of apo vs spleen TCR discrimination when  $pR \approx m\bar{n}$ , where  $pR$  is the expected number  
664 of TCR specific to the foreign antigen.

## 665 **Other Software for Statistical Analysis**

666 The Jensen-Shannon dendrograms linkage is computed by the Ward method as provided by the  
667 function "linkage", reordered according to the function "optimal\_leaf\_ordering", both from the  
668 Python module "cluster.hierarchy" in *scipy* package (version 1.7.3). The Pearson correlation coefficient  
669 is computed with the Python function "pearsonr" as contained in the module "stats" in the  
670 *scipy* package. The coefficient of determination  $R^2$  is computed with the Python function "r2\_score"  
671 as contained in the module "metrics" in the *sklearn* package.

## 672 **Code availability**

673 All code for reproducing the figures of this paper can be found at [https://github.com/statbiophys/  
674 thymic\\_development\\_2022.git](https://github.com/statbiophys/thymic_development_2022.git).

## 675 **Acknowledgments**

676 The study was supported by a '80 prime' CNRS-Weizmann PhD scholarship, European Research  
677 Council COG 724208 and ANR-19-CE45-0018 'RESP- REP' from the Agence Nationale de la Recherche  
678 and DFG grant CRC 1310 'Predictability in Evolution'.

## 679 **References**

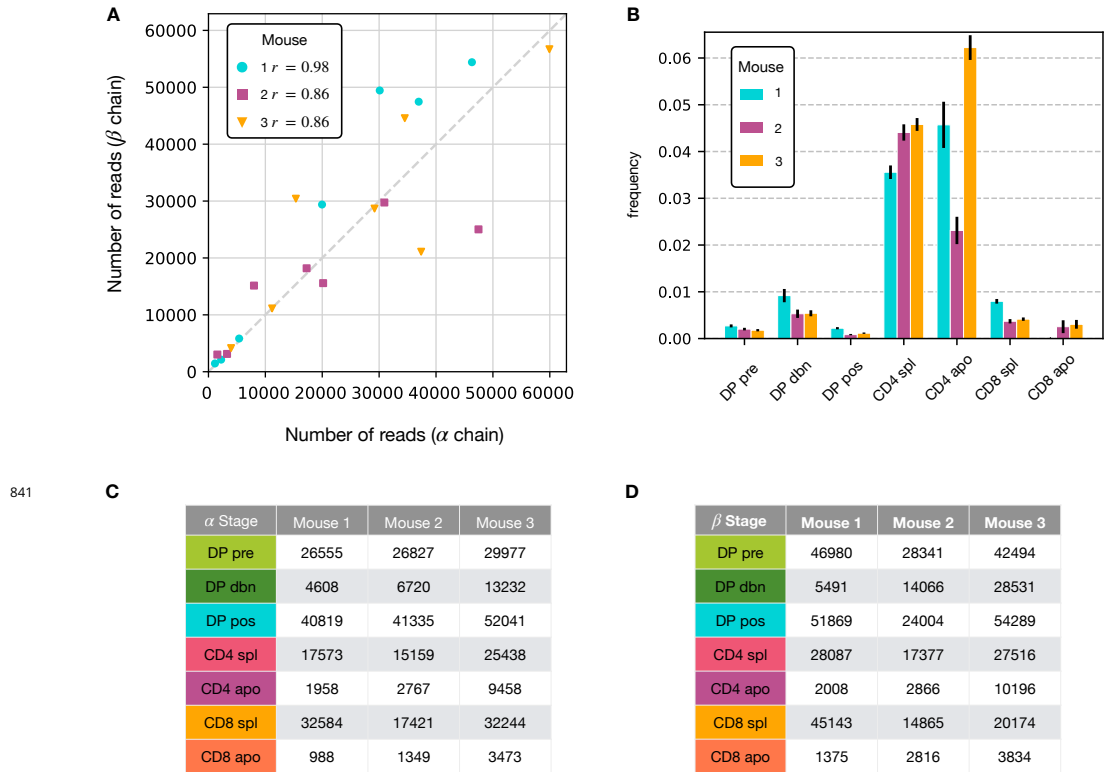
- 680 **Archer E**, Park IM, Pillow J. Bayesian Entropy Estimation for Countable Discrete Distributions. arXiv:13020328  
681 [cs, math]. 2014 Apr; <http://arxiv.org/abs/1302.0328>, arXiv: 1302.0328.
- 682 **Bains I**, Yates AJ, Callard RE. Heterogeneity in Thymic Emigrants: Implications for Thymectomy and Immunose-  
683 nescence. PLoS ONE. 2013 Feb; 8(2):e49554. doi: 10.1371/journal.pone.0049554.
- 684 **Bolger AM**, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics. 2014  
685 Aug; 30(15):2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>, doi: 10.1093/bioinformatics/btu170.
- 686 **Butler TC**, Kardar M, Chakraborty AK. Quorum sensing allows T cells to discriminate between self and nonself.  
687 Proceedings of the National Academy of Sciences. 2013 Jul; 110(29):11833–11838. [http://www.pnas.org/cgi/  
688 doi/10.1073/pnas.1222467110](http://www.pnas.org/cgi/doi/10.1073/pnas.1222467110), doi: 10.1073/pnas.1222467110.
- 689 **Carter JA**, Preall JB, Grigaityte K, Goldfless SJ, Jeffery E, Briggs AW, Vigneault F, Atwal GS. Single T Cell Sequencing  
690 Demonstrates the Functional Role of  $A\beta$  TCR Pairing in Cell Lineage and Antigen Specificity. Frontiers in  
691 Immunology. 2019 Jul; 10:1516. doi: 10.3389/fimmu.2019.01516.
- 692 **Chao A**, Shen TJ. Nonparametric estimation of Shannon's index of diversity when there are unseen species in  
693 sample. Environmental and Ecological Statistics. 2003; 10(4):429–443. [http://link.springer.com/10.1023/A:  
694 1026096204727](http://link.springer.com/10.1023/A:1026096204727), doi: 10.1023/A:1026096204727.
- 695 **Cinelli M**, Sun Y, Best K, Heather JM, Reich-Zeliger S, Shifrut E, Friedman N, Shawe-Taylor J, Chain B. Fea-  
696 ture selection using a one dimensional naïve Bayes' classifier increases the accuracy of support vector  
697 machine classification of CDR3 repertoires. Bioinformatics. 2017 Jan; p. btw771. [https://academic.oup.com/  
698 bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw771](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btw771), doi: 10.1093/bioinformatics/btw771.
- 699 **Daley SR**, Koay HF, Dobbs K, Bosticardo M, Wirasinha RC, Pala F, Castagnoli R, Rowe JH, Ott de Bruin LM, Keles S,  
700 Lee YN, Somech R, Holland SM, Delmonte OM, Draper D, Maxwell S, Niemela J, Stoddard J, Rosenzweig SD,  
701 Poliani PL, et al. Cysteine and Hydrophobic Residues in CDR3 Serve as Distinct T-cell Self-Reactivity Indices.  
702 Journal of Allergy and Clinical Immunology. 2019 Jul; 144(1):333–336. doi: 10.1016/j.jaci.2019.03.022.
- 703 **Dupic T**, Marcou Q, Walczak AM, Mora T. Genesis of the  $A\beta$  T-Cell Receptor. PLOS Computational Biology. 2019  
704 Mar; 15(3):e1006874. doi: 10.1371/journal.pcbi.1006874.
- 705 **Elhanati Y**, Murugan A, Callan CG, Mora T, Walczak AM. Quantifying selection in immune receptor repertoires.  
706 Proceedings of the National Academy of Sciences. 2014 Jul; 111(27):9875–9880. [http://www.pnas.org/lookup/  
707 doi/10.1073/pnas.1409572111](http://www.pnas.org/lookup/doi/10.1073/pnas.1409572111), doi: 10.1073/pnas.1409572111.



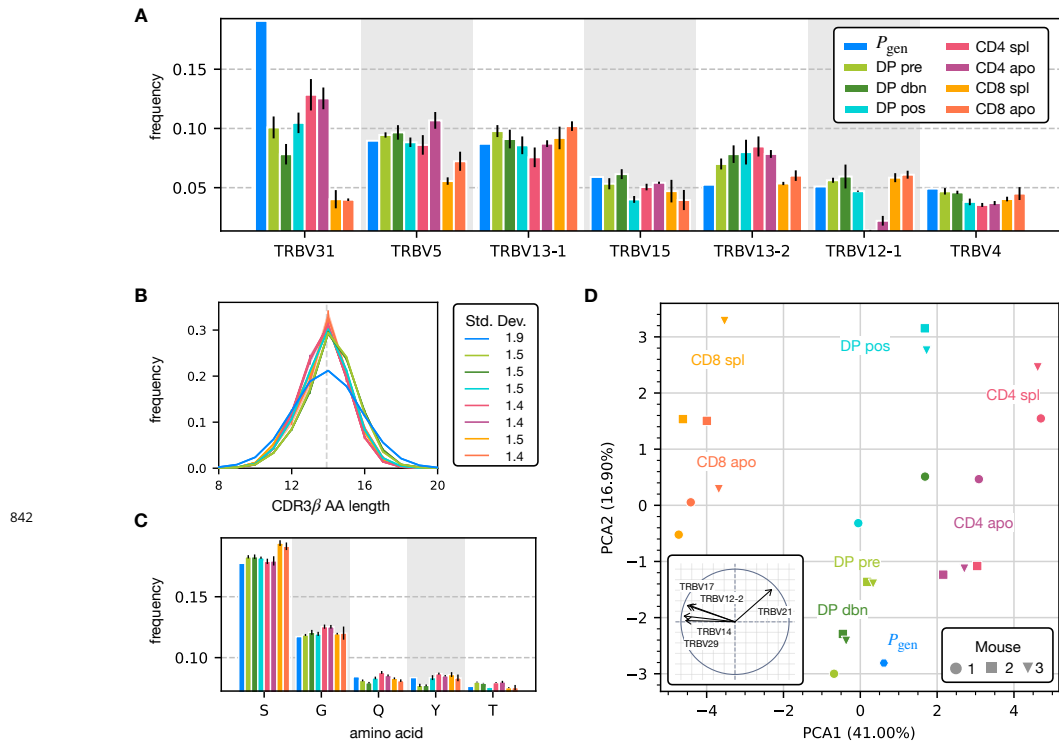
- 708 **Emerson R**, Sherwood A, Desmarais C, Malhotra S, Phippard D, Robins H. Estimating the Ratio of CD4+ to CD8+  
709 T Cells Using High-Throughput Sequence Data. *Journal of Immunological Methods*. 2013 May; 391(1-2):14-21.  
710 doi: [10.1016/j.jim.2013.02.002](https://doi.org/10.1016/j.jim.2013.02.002).
- 711 **Gallegos AM**, Bevan MJ. Central tolerance: good but imperfect. *Immunological Reviews*.  
712 2006; 209(1):290-296. <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.0105-2896.2006.00348.x>, doi:  
713 <https://doi.org/10.1111/j.0105-2896.2006.00348.x>.
- 714 **Garner LC**, Klenerman P, Provine NM. Insights Into Mucosal-Associated Invariant T Cell Biology From Studies of  
715 Invariant Natural Killer T Cells. *Frontiers in Immunology*. 2018; 9. <https://www.frontiersin.org/article/10.3389/fimmu.2018.01478>.  
716
- 717 **Gray DHD**, Seach N, Ueno T, Milton MK, Liston A, Lew AM, Goodnow CC, Boyd RL. Developmental kinetics,  
718 turnover, and stimulatory capacity of thymic epithelial cells. *Blood*. 2006 Dec; 108(12):3777-3785. <https://doi.org/10.1182/blood-2006-02-004531>, doi: 10.1182/blood-2006-02-004531.  
719
- 720 **Grigaityte K**, Carter JA, Goldfless SJ, Jeffery EW, Hause RJ, Jiang Y, Koppstein D, Briggs AW, Church GM, Vigneault  
721 F, Atwal GS. Single-Cell Sequencing Reveals  $\alpha\beta$  Chain Pairing Shapes the T Cell Repertoire. *Immunology*; 2017.
- 722 **Hamming RW**. Error detecting and error correcting codes. *The Bell System Technical Journal*. 1950 Apr;  
723 29(2):147-160. doi: [10.1002/j.1538-7305.1950.tb00463.x](https://doi.org/10.1002/j.1538-7305.1950.tb00463.x), conference Name: The Bell System Technical  
724 Journal.
- 725 **Hozumi N**, Tonegawa S. Evidence for somatic rearrangement of immunoglobulin genes coding for variable  
726 and constant regions. *Proceedings of the National Academy of Sciences*. 1976 Oct; 73(10):3628-3632.  
727 <https://www.pnas.org/content/73/10/3628>, doi: [10.1073/pnas.73.10.3628](https://doi.org/10.1073/pnas.73.10.3628), publisher: National Academy of  
728 Sciences Section: Research Article.
- 729 **Isacchini G**, Walczak AM, Mora T, Nourmohammad A. Deep generative selection models of T and B cell  
730 receptor repertoires with soNNia. *Proceedings of the National Academy of Sciences*. 2021 Apr; 118(14).  
731 <https://www.pnas.org/content/118/14/e2023141118>, doi: [10.1073/pnas.2023141118](https://doi.org/10.1073/pnas.2023141118), publisher: National  
732 Academy of Sciences Section: Physical Sciences.
- 733 **Izraelson M**, Nakonechnaya TO, Moltedo B, Egorov ES, Kasatskaya SA, Putintseva EV, Mamedov IZ, Staroverov  
734 DB, Shemiakina II, Zakharova MY, Davydov AN, Bolotin DA, Shugay M, Chudakov DM, Rudensky AY, Britanova  
735 OV. Comparative Analysis of Murine T-Cell Receptor Repertoires. *Immunology*. 2018 Feb; 153(2):133-144. doi:  
736 [10.1111/imm.12857](https://doi.org/10.1111/imm.12857).
- 737 **King JL**, Jukes TH. Non-Darwinian Evolution. *Science*. 1969; 164(3881):788-798. <https://www.science.org/doi/abs/10.1126/science.164.3881.788>, doi: [10.1126/science.164.3881.788](https://doi.org/10.1126/science.164.3881.788).
- 739 **Košmrlj A**, Chakraborty AK, Kardar M, Shakhnovich EI. Thymic Selection of T-Cell Receptors as an Extreme Value  
740 Problem. *Physical Review Letters*. 2009 Aug; 103(6):068103. <https://link.aps.org/doi/10.1103/PhysRevLett.103.068103>.  
741 [068103](https://doi.org/10.1103/PhysRevLett.103.068103), doi: [10.1103/PhysRevLett.103.068103](https://doi.org/10.1103/PhysRevLett.103.068103).
- 742 **Lagattuta KA**, Kang JB, Nathan A, Pauken KE, Jonsson AH, Rao DA, Sharpe AH, Ishigaki K, Raychaudhuri S.  
743 Repertoire Analyses Reveal T Cell Antigen Receptor Sequence Features That Influence T Cell Fate. *Nature*  
744 *Immunology*. 2022 Mar; 23(3):446-457. doi: [10.1038/s41590-022-01129-x](https://doi.org/10.1038/s41590-022-01129-x).
- 745 **Lefranc MP**, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, Carillon E, Duvergey H, Houles A,  
746 Paysan-Lafosse T, Hadi-Saljoqi S, Sasorith S, Lefranc G, Kossida S. IMGT®, the international ImMunoGeneTics  
747 information system® 25 years on. *Nucleic Acids Research*. 2015 Jan; 43(D1):D413-D422. <https://doi.org/10.1093/nar/gku1056>, doi: [10.1093/nar/gku1056](https://doi.org/10.1093/nar/gku1056).
- 749 **Levenshtein VI**, et al. Binary codes capable of correcting deletions, insertions, and reversals. In: *Soviet physics*  
750 *doklady*, vol. 10 Soviet Union; 1966. p. 707-710.
- 751 **Liebmann M**, Hücke S, Koch K, Eschborn M, Ghelman J, Chasan AI, Glander S, Schädlich M, Kuhlencord M,  
752 Daber NM, Eveslage M, Beyer M, Dietrich M, Albrecht P, Stoll M, Busch KB, Wiendl H, Roth J, Kuhlmann  
753 T, Klotz L. Nur77 serves as a molecular brake of the metabolic switch during T cell activation to restrict  
754 autoimmunity. *Proceedings of the National Academy of Sciences of the United States of America*. 2018 Aug;  
755 115(34):E8017-E8026. doi: [10.1073/pnas.1721049115](https://doi.org/10.1073/pnas.1721049115).
- 756 **Lu J**, Van Laethem F, Bhattacharya A, Craveiro M, Saba I, Chu J, Love NC, Tikhonova A, Radaev S, Sun X, Ko  
757 A, Arnon T, Shifrut E, Friedman N, Weng NP, Singer A, Sun PD. Molecular constraints on CDR3 for thymic  
758 selection of MHC-restricted TCRs from a random pre-selection repertoire. *Nature Communications*. 2019 Dec;  
759 10(1):1019. <http://www.nature.com/articles/s41467-019-08906-7>, doi: [10.1038/s41467-019-08906-7](https://doi.org/10.1038/s41467-019-08906-7).

- 760 **Madi A**, Poran A, Shifrut E, Reich-Zeliger S, Greenstein E, Zaretsky I, Arnon T, Laethem FV, Singer A, Lu J, Sun PD,  
761 Cohen IR, Friedman N. T cell receptor repertoires of mice and humans are clustered in similarity networks  
762 around conserved public CDR3 sequences. *eLife*. 2017 Jul; 6:e22057. <https://elifesciences.org/articles/22057>,  
763 doi: 10.7554/eLife.22057.
- 764 **Madi A**, Shifrut E, Reich-Zeliger S, Gal H, Best K, Ndifon W, Chain B, Cohen IR, Friedman N. T-cell receptor  
765 repertoires share a restricted set of public and abundant CDR3 sequences that are associated with self-  
766 related immunity. *Genome Research*. 2014 Oct; 24(10):1603–1612. <http://genome.cshlp.org/lookup/doi/10.1101/gr.170753.113>, doi: 10.1101/gr.170753.113.
- 768 **Marcou Q**, Mora T, Walczak AM. High-throughput immune repertoire analysis with IGoR. *Nature Communica-*  
769 *tions*. 2018 Dec; 9(1):561. <http://www.nature.com/articles/s41467-018-02832-w>, doi: 10.1038/s41467-018-  
770 02832-w.
- 771 **Miller G**. Note on the bias of information estimates. *Information Theory in Psychology: Problems and Methods*.  
772 1955; p. 95–100.
- 773 **Montemurro A**, Schuster V, Povlsen HR, Bentzen AK, Jurtz V, Chronister WD, Crinklaw A, Hadrup SR, Winther  
774 O, Peters B, Jessen LE, Nielsen M. NetTCR-2.0 enables accurate prediction of TCR-peptide binding by using  
775 paired TCR $\alpha$  and  $\beta$  sequence data. *Communications Biology*. 2021 Dec; 4(1):1060. [https://www.nature.com/](https://www.nature.com/articles/s42003-021-02610-3)  
776 [articles/s42003-021-02610-3](https://www.nature.com/articles/s42003-021-02610-3), doi: 10.1038/s42003-021-02610-3.
- 777 **Moran AE**, Holzapfel KL, Xing Y, Cunningham NR, Maltzman JS, Punt J, Hogquist KA. T cell receptor signal strength  
778 in Treg and iNKT cell development demonstrated by a novel fluorescent reporter mouse. *The Journal of*  
779 *Experimental Medicine*. 2011 Jun; 208(6):1279–1289. doi: 10.1084/jem.20110308.
- 780 **Murugan A**, Mora T, Walczak AM, Callan CG. Statistical inference of the generation probability of T-cell receptors  
781 from sequence repertoires. *Proceedings of the National Academy of Sciences*. 2012 Oct; 109(40):16161–16166.  
782 <http://www.pnas.org/cgi/doi/10.1073/pnas.1212755109>, doi: 10.1073/pnas.1212755109.
- 783 **Nemenman I**, Shafee F, Bialek W. Entropy and Inference, Revisited. *arXiv*. 2002; p. 8. [arXiv:physics/0108025](https://arxiv.org/abs/physics/0108025).
- 784 **Oakes T**, Heather JM, Best K, Byng-Maddick R, Husovsky C, Ismail M, Joshi K, Maxwell G, Noursadeghi M, Riddell  
785 N, Ruehl T, Turner CT, Uddin I, Chain B. Quantitative Characterization of the T Cell Receptor Repertoire of  
786 Naïve and Memory Subsets Using an Integrated Experimental and Computational Pipeline Which Is Robust,  
787 Economical, and Versatile. *Frontiers in Immunology*. 2017; 8. [https://www.frontiersin.org/article/10.3389/](https://www.frontiersin.org/article/10.3389/fimmu.2017.01267)  
788 [fimmu.2017.01267](https://www.frontiersin.org/article/10.3389/fimmu.2017.01267).
- 789 **Park JE**, Botting RA, Domínguez Conde C, Popescu DM, Lavaert M, Kunz DJ, Goh I, Stephenson E, Ragazzini  
790 R, Tuck E, Wilbrey-Clark A, Roberts K, Kedlian VR, Ferdinand JR, He X, Webb S, Maunder D, Vandamme N,  
791 Mahbubani KT, Polanski K, et al. A cell atlas of human thymic development defines T cell repertoire formation.  
792 *Science*. 2020 Feb; 367(6480):eaay3224. <https://www.sciencemag.org/lookup/doi/10.1126/science.aay3224>,  
793 doi: 10.1126/science.aay3224.
- 794 **Polonsky M**, Rimer J, Kern-Perets A, Zaretsky I, Miller S, Bornstein C, David E, Kopelman NM, Stelzer G, Porat  
795 Z, Chain B, Friedman N. Induction of CD4 T cell memory by local cellular collectivity. *Science*. 2018 Jun;  
796 360(6394):eaaj1853. <https://www.sciencemag.org/lookup/doi/10.1126/science.aaj1853>, doi: 10.1126/sci-  
797 [ence.aaj1853](https://www.sciencemag.org/lookup/doi/10.1126/science.aaj1853).
- 798 **Robert PA**, Kunze-Schumacher H, Greiff V, Krueger A. Modeling the Dynamics of T-Cell Development in the  
799 Thymus. *Entropy*. 2021 Apr; 23(4):437. <https://www.mdpi.com/1099-4300/23/4/437>, doi: 10.3390/e23040437.
- 800 **Sethna Z**, Elhanati Y, Dudgeon CR, Callan CG, Levine AJ, Mora T, Walczak AM. Insights into immune sys-  
801 tem development and function from mouse T-cell repertoires. *Proceedings of the National Academy of*  
802 *Sciences*. 2017 Feb; 114(9):2253–2258. <http://www.pnas.org/lookup/doi/10.1073/pnas.1700241114>, doi:  
803 [10.1073/pnas.1700241114](http://www.pnas.org/lookup/doi/10.1073/pnas.1700241114).
- 804 **Sethna Z**, Isacchini G, Dupic T, Mora T, Walczak AM, Elhanati Y. Population variability in the generation and  
805 thymic selection of T-cell repertoires. *PLoS Computational Biology*. 2020 Dec; 16(12):e1008394. [http:](http://arxiv.org/abs/2001.02843)  
806 [//arxiv.org/abs/2001.02843](http://arxiv.org/abs/2001.02843), doi: 10.1371/journal.pcbi.1008394, arXiv: 2001.02843.
- 807 **Soto C**, Finn JA, Willis JR, Day SB, Sinkovits RS, Jones T, Schmitz S, Meiler J, Branchizio A, Crowe JE. PyIR: a  
808 scalable wrapper for processing billions of immunoglobulin and T cell receptor sequences using IgBLAST.  
809 *BMC Bioinformatics*. 2020 Jul; 21(1):314. <https://doi.org/10.1186/s12859-020-03649-5>, doi: 10.1186/s12859-  
810 [020-03649-5](https://doi.org/10.1186/s12859-020-03649-5).

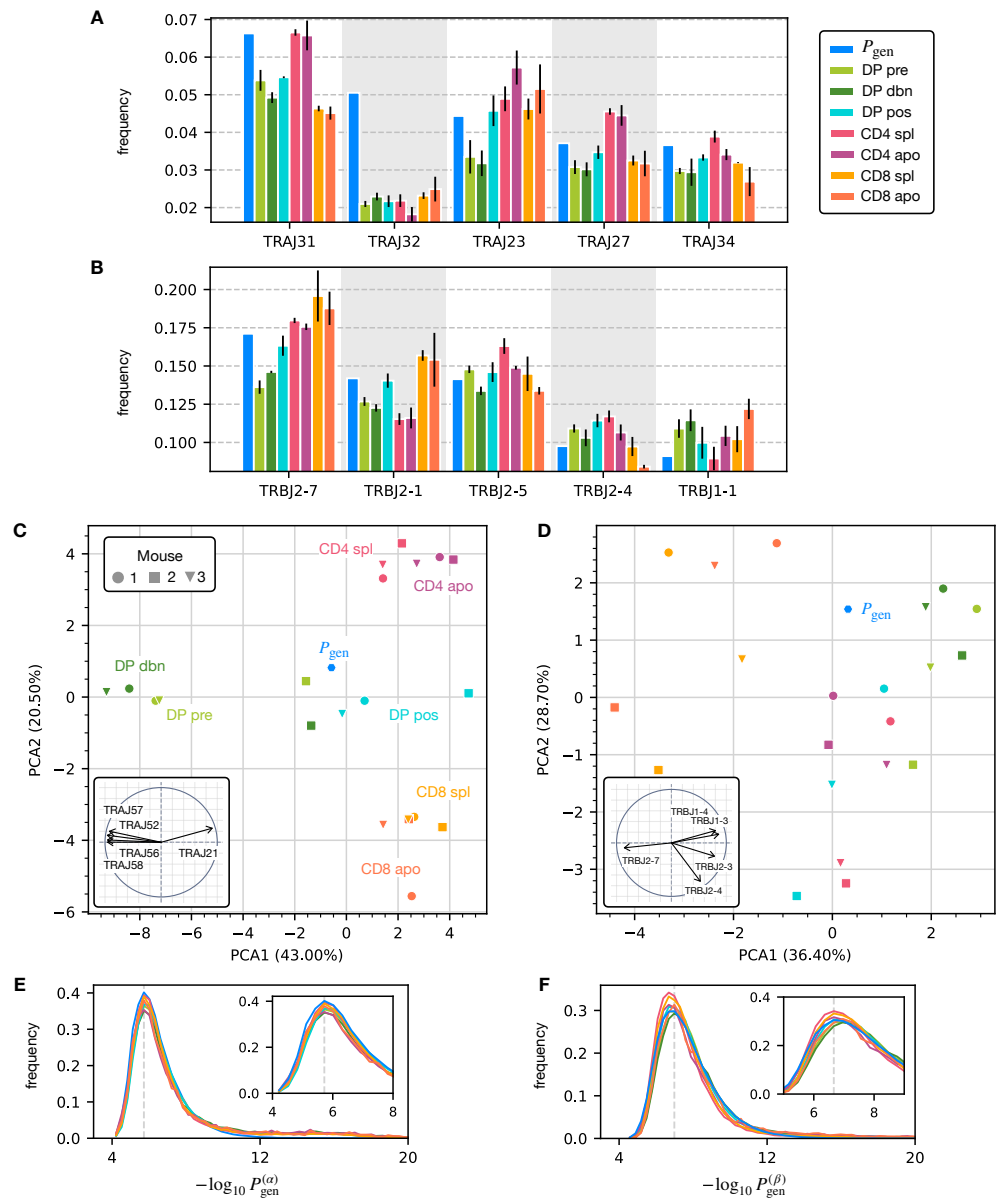
- 811 **Stadinski BD**, Shekhar K, Gómez-Touriño I, Jung J, Sasaki K, Sewell AK, Peakman M, Chakraborty AK, Huseby ES.  
812 Hydrophobic CDR3 Residues Promote the Development of Self-Reactive T Cells. *Nature Immunology*. 2016  
813 Aug; 17(8):946–955. doi: [10.1038/ni.3491](https://doi.org/10.1038/ni.3491).
- 814 **Stritesky GL**, Xing Y, Erickson JR, Kalekar LA, Wang X, Mueller DL, Jameson SC, Hogquist KA. Murine Thymic  
815 Selection Quantified Using a Unique Method to Capture Deleted T Cells. *Proceedings of the National Academy*  
816 *of Sciences*. 2013 Mar; 110(12):4679–4684. doi: [10.1073/pnas.1217532110](https://doi.org/10.1073/pnas.1217532110).
- 817 **Sun Y**, Best K, Cinelli M, Heather JM, Reich-Zeliger S, Shifrut E, Friedman N, Shawe-Taylor J, Chain B. Speci-  
818 ficity, Privacy, and Degeneracy in the CD4 T Cell Receptor Repertoire Following Immunization. *Frontiers*  
819 *in Immunology*. 2017 Apr; 8. <http://journal.frontiersin.org/article/10.3389/fimmu.2017.00430/full>, doi:  
820 [10.3389/fimmu.2017.00430](https://doi.org/10.3389/fimmu.2017.00430).
- 821 **Thomas N**, Best K, Cinelli M, Reich-Zeliger S, Gal H, Shifrut E, Madi A, Friedman N, Shawe-Taylor J, Chain  
822 B. Tracking global changes induced in the CD4 T-cell receptor repertoire by immunization with a com-  
823 plex antigen using short stretches of CDR3 protein sequence. *Bioinformatics*. 2014 Nov; 30(22):3181–  
824 3188. <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btu523>, doi:  
825 [10.1093/bioinformatics/btu523](https://doi.org/10.1093/bioinformatics/btu523).
- 826 **Tubiana J**, Cocco S, Monasson R. Learning Compositional Representations of Interacting Systems with Restricted  
827 Boltzmann Machines: Comparative Study of Lattice Proteins. *Neural Computation*. 2019 08; 31(8):1671–1717.  
828 [https://doi.org/10.1162/neco\\_a\\_01210](https://doi.org/10.1162/neco_a_01210), doi: [10.1162/neco\\_a\\_01210](https://doi.org/10.1162/neco_a_01210).
- 829 **Weber A**, Born J, Rodriguez Martínez M. TITAN: T-cell receptor specificity prediction with bimodal attention  
830 networks. *Bioinformatics*. 2021 Aug; 37(Supplement\_1):i237–i244. [https://academic.oup.com/bioinformatics/](https://academic.oup.com/bioinformatics/article/37/Supplement_1/i237/6319659)  
831 [article/37/Supplement\\_1/i237/6319659](https://academic.oup.com/bioinformatics/article/37/Supplement_1/i237/6319659), doi: [10.1093/bioinformatics/btab294](https://doi.org/10.1093/bioinformatics/btab294).
- 832 **Wing K**, Sakaguchi S. Regulatory T Cells Exert Checks and Balances on Self Tolerance and Autoimmunity. *Nature*  
833 *Immunology*. 2010 Jan; 11(1):7–13. doi: [10.1038/ni.1818](https://doi.org/10.1038/ni.1818).
- 834 **Yates AJ**. Theories and Quantification of Thymic Selection. *Frontiers in Immunology*. 2014; 5. [http://journal.](http://journal.frontiersin.org/article/10.3389/fimmu.2014.00013/abstract)  
835 [frontiersin.org/article/10.3389/fimmu.2014.00013/abstract](http://journal.frontiersin.org/article/10.3389/fimmu.2014.00013/abstract), doi: [10.3389/fimmu.2014.00013](https://doi.org/10.3389/fimmu.2014.00013).
- 836 **Ye J**, Ma N, Madden TL, Ostell JM. IgBLAST: an immunoglobulin variable domain sequence analysis tool. *Nucleic*  
837 *Acids Research*. 2013 Jul; 41(Web Server issue):W34–40. doi: [10.1093/nar/gkt382](https://doi.org/10.1093/nar/gkt382).
- 838 **Yu W**, Jiang N, Ebert PJR, Kidd BA, Müller S, Lund PJ, Juang J, Adachi K, Tse T, Birnbaum ME, Newell EW, Wilson DM,  
839 Grotenbreg GM, Valitutti S, Quake SR, Davis MM. Clonal Deletion Prunes but Does Not Eliminate Self-Specific  
840  $A\beta$  CD8+ T Lymphocytes. *Immunity*. 2015 May; 42(5):929–941. doi: [10.1016/j.immuni.2015.05.001](https://doi.org/10.1016/j.immuni.2015.05.001).



**Figure 1—figure supplement 1.** (A) Number of reads for the alpha chain vs the number for the beta chain within the same dataset. In the box is shown the Pearson correlation coefficient. Distribution of iNKT clonotypes for the  $\alpha$  chain. (B) The relative amount of (TRAV11, TRAJ18) clonotypes is significantly higher for all CD4 stages in all mice. (C) Numbers of unique productive (in-frame and with no stopping codons) single chain obtained for the maturation stages in each mouse after annotation for the  $\alpha$  chain. (D) Numbers of unique productive for the  $\beta$  chain.

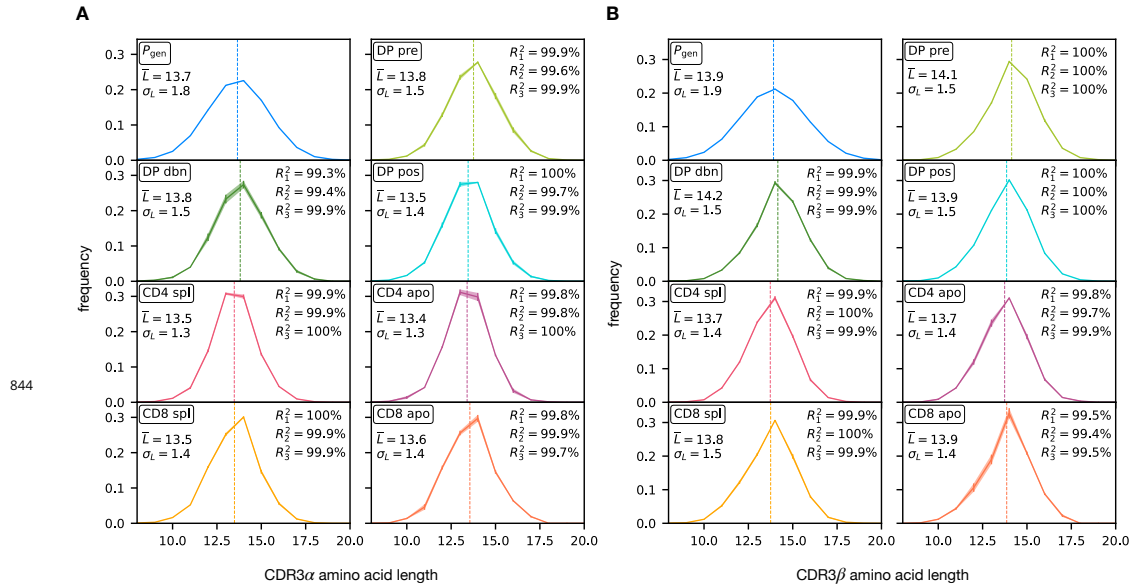


**Figure 2—figure supplement 1. (A)** The distribution of TRBV genes at different maturation stages compared to the generation distribution. **(B)** Distribution of  $\beta$  chain CDR3 amino acid sequence lengths. The CDR3 is defined between the typical cysteine and phenylalanine position. The dashed line represents the average length according to the  $P_{gen}$  model. **(C)** The distribution of the most frequent amino acid within the CDR3 region for the TRBV sequences. **(D)** Principal component analysis according to the TRBV gene distribution at each maturation stage. Insert: projection on the principal axis of the five most representative TRBV genes. Analogous results for TRBJ are shown in *Figure 2—figure Supplement 2D*.

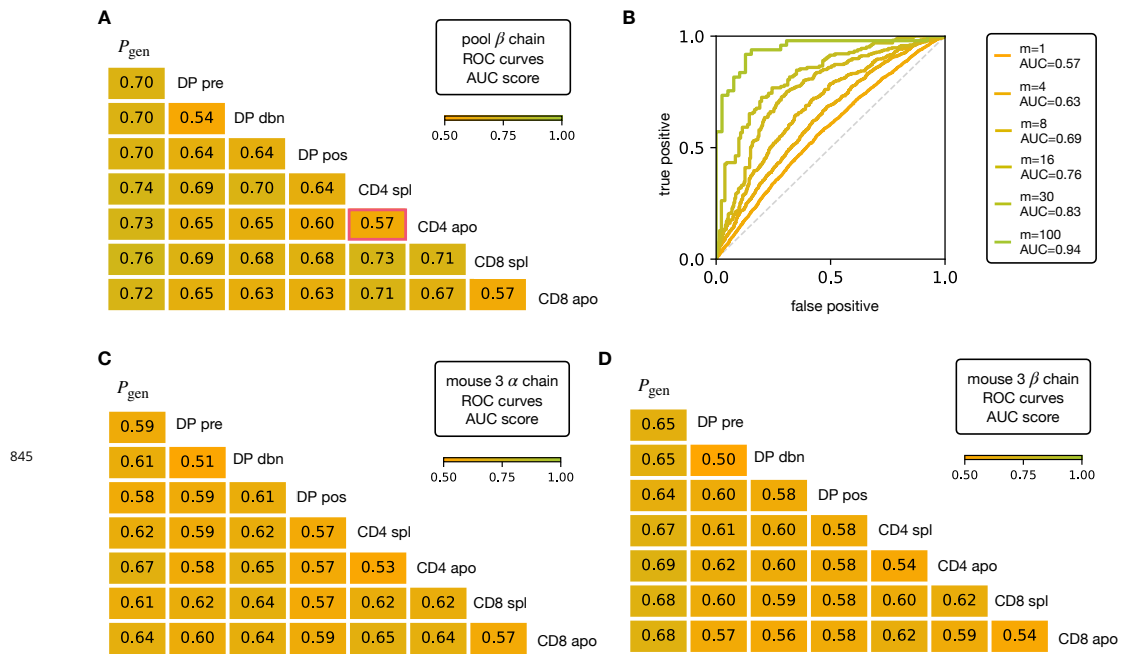


**Figure 2—figure supplement 2.** (A) The distribution of J genes at different maturation stages compared to the generation distribution. for the  $\alpha$  chain. (B) The distribution of J genes for the  $\beta$  chain. (C) Principal component analysis of  $\alpha$  chain J gene usage. Similarly as for the V genes (Figure 2C), DP, CD4 and CD8 maturation stages cluster by the cell types. (D) Principal component analysis of  $\beta$  chain J gene usage. (E)  $P_{gen}$  values distribution for the clonotypes at each maturation stage. Insert: blow up of the peak of the distribution (the dashed line) for the  $\alpha$  chain. (F)  $P_{gen}$  values distribution the  $\beta$  chain.

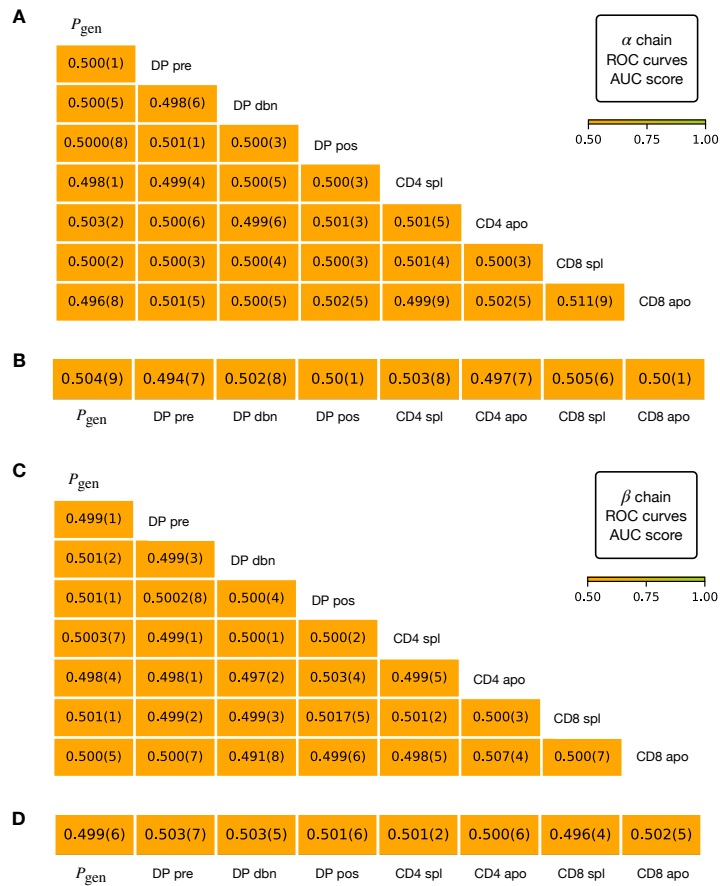




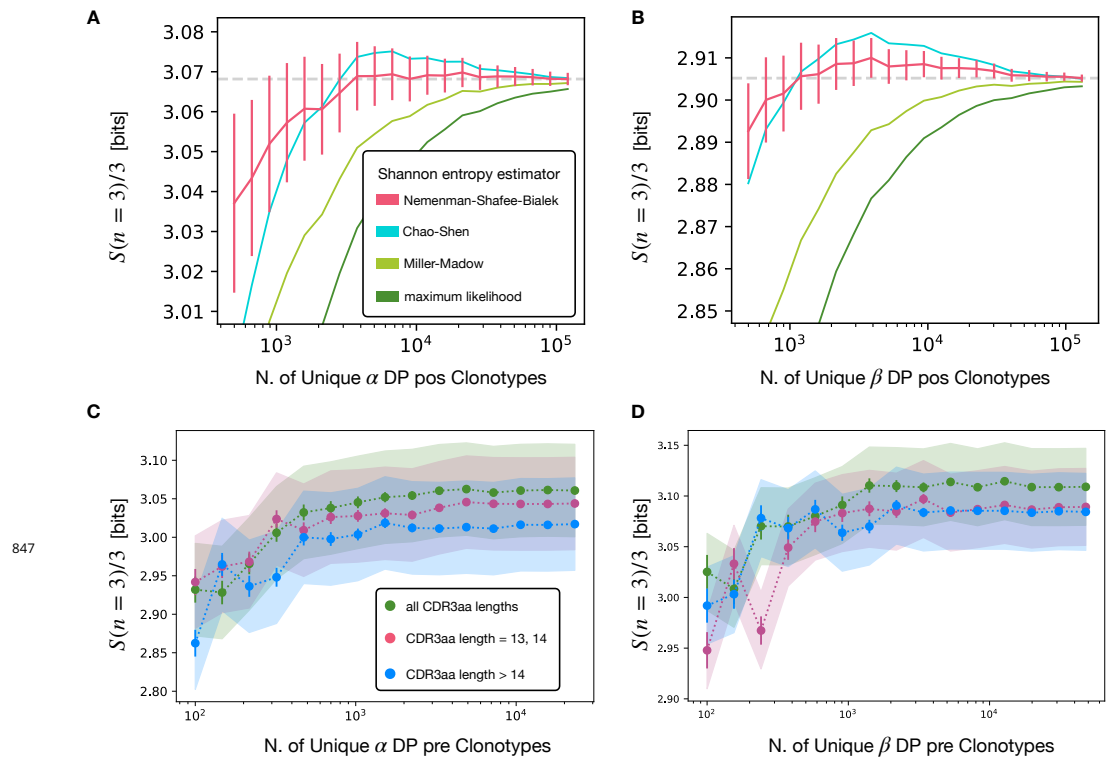
**Figure 2—figure supplement 3. (A)** Amino acid CDR3 length for TCR $\alpha$  sequences. The thick curve represent the average length across the three different mice, while the shaded part illustrates the mouse variability. On the left of each box we report the empirical average and standard deviation of the distribution; on the right the coefficient of determination  $R^2$  between each individual distribution and the average (see Materials and Methods). **(B)** Analogous for the TCR $\beta$  sequences.



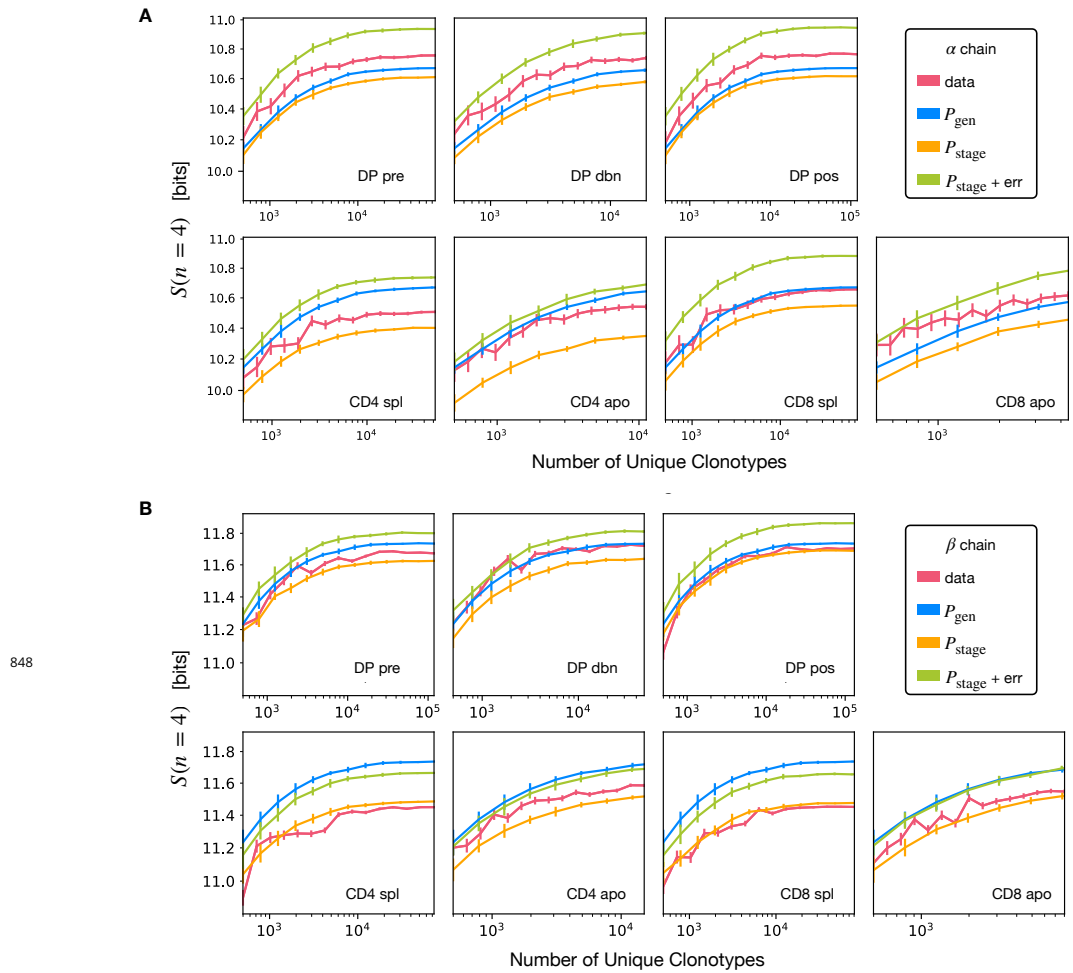
**Figure 3—figure supplement 1.** (A) AUC values computed from the ROC curves of the linear classifiers for TCR $\beta$  sequences between pairs of maturation stages. The training/testing set is a random subsample containing 70%/30% of the full dataset at a given maturation stage. (B) Illustration of the improvements of group discriminability between the stages CD4 spl and CD4 apo. (C) AUC values computed from the ROC curves of the linear classifiers for TCR sequences for the unpooled largest dataset for an individual (mouse 3). We observe that the score is never higher than for the pooled case and in fact it's typically worse for the  $\alpha$  chain. (D) Analogous for the  $\beta$  chain.



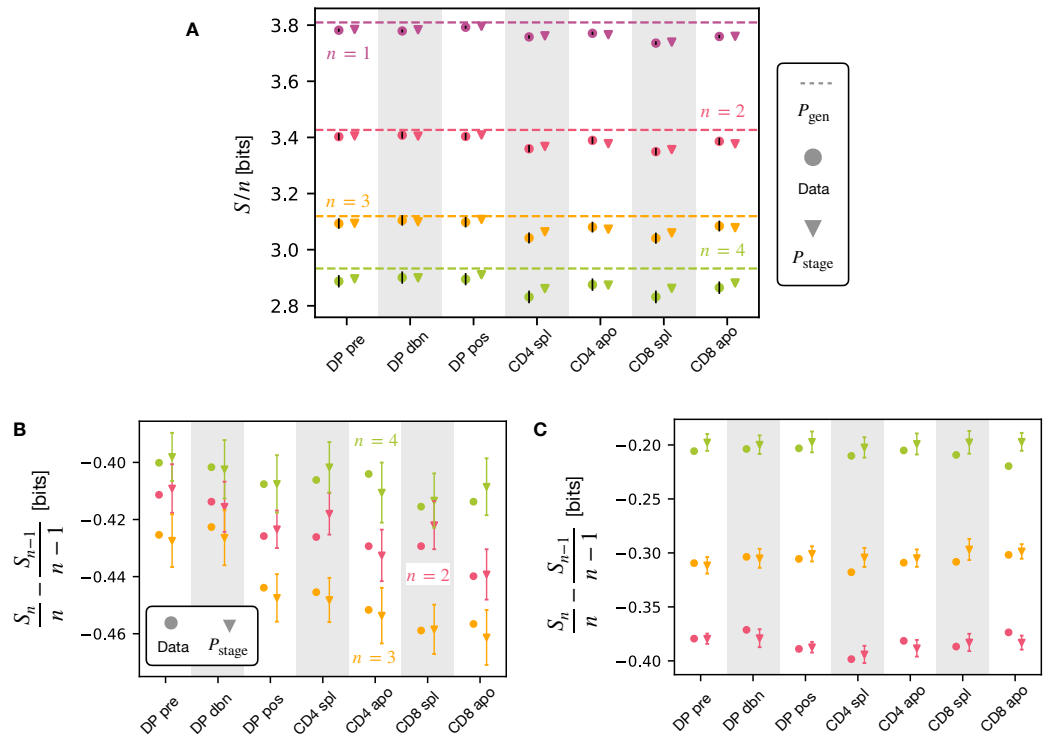
**Figure 3—figure supplement 2.** Here we organize the different train/test datasets as in the main text (**Figure 3**) and we learn logistic regression classifiers as in Materials and Methods, of which we report the AUC score. Analogous results are obtained by usage of a random decision forest. **(A)** We randomly shuffle the labels for each pair of a stages, observing that it is impossible to obtain two distinguishable repertoires through random mixing. The error for the last digit is expressed in brackets and is obtained from 20 realizations of the shuffling. **(B)** We randomly assign two labels to the TCR $\alpha$  sequences of a single repertoire and split the repertoire in a test and a train group. Again, we show that it's not possible to obtain the scores of the main text by randomly pick chains from a defined stage. In all these controls we sub-sampled to the size of the smallest dataset available in order to check for issue size. As in the main text (see Materials and Methods), we test linear and decision forest classifiers, imposing the size of the larger class to not exceed of more then 25% the size of the smaller, with the test set corresponding to 30% of data. **(C)** Classifiers learnt on pair of  $\beta$  stages with randomly shuffled labels. **(D)** Classifiers learnt on single  $\beta$ stages with randomly assigned labels.



**Figure 4—figure supplement 1.** (A) Averages over different subsamples of the same size for productive DP pos pooled  $\alpha$  clonotypes. The estimators compared are: Nemenman-Shafee-Bialek (see Materials and Methods) *Nemenman et al. (2002)*, Chao-Shen *Chao and Shen (2003)*, Miller-Madow *Miller (1955)* and maximum likelihood (naive estimator). (B) Analogous for the  $\beta$  chain. (C) Analysis of the length dependence of 3-gram entropy associated to different choices of the CDR3 amino acid lengths for the clonotypes considered in the  $\alpha$  DP pre stage. The errorbars are estimated with the NSB method, while the shaded curve represent the sequencing error. We notice how the difference between the different choices is greatly covered by the sequencing error. We prefer then to use all CDR3 lengths for the higher statistics. (D) Analogous for the  $\beta$  chain.

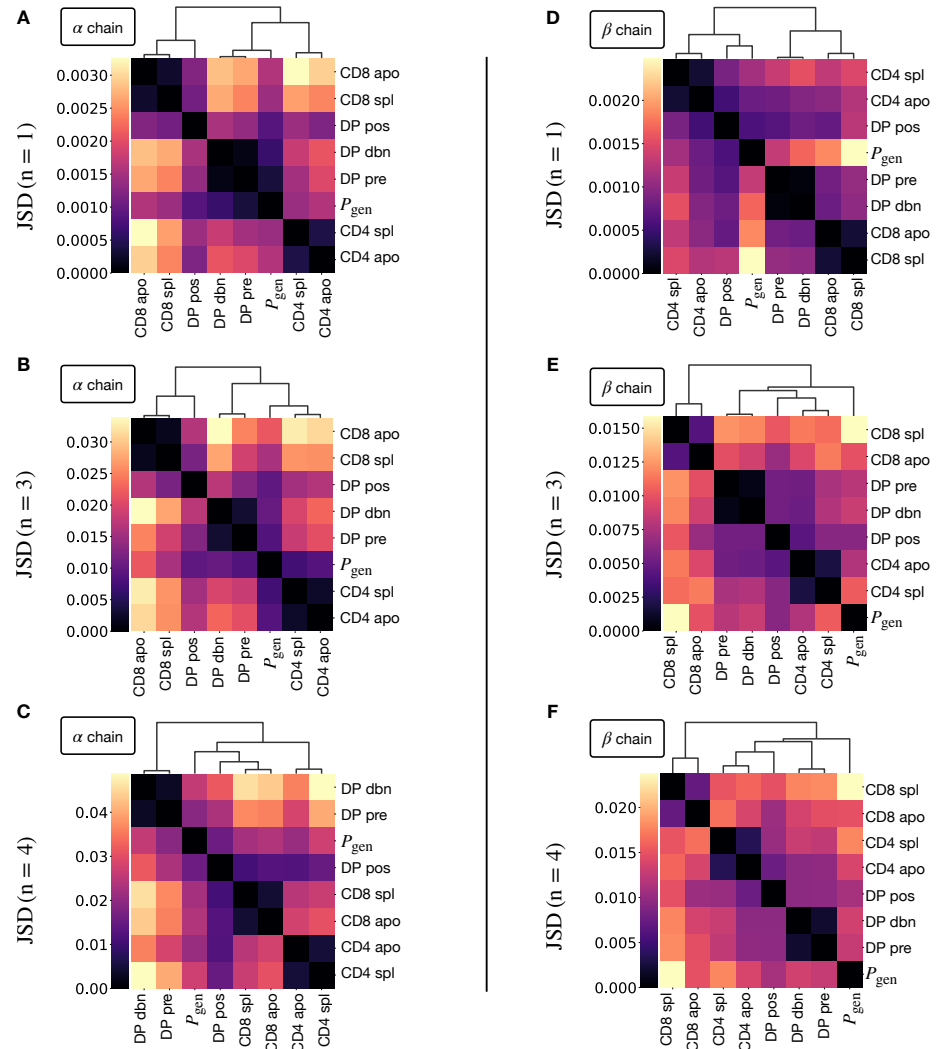


**Figure 4—figure supplement 2. (A)** We subsample the unique clonotypes for the TCR $\alpha$  sequences and we check that the 4-gram entropy estimations converge with increasing number of unique clonotypes. The synthetic sequences are produced with the generation model  $P_{\text{gen}}$  (same for all plots), the different selection models  $P_{\text{stage}}$  and a  $P_{\text{stage}}$  selection model with synthetic nucleotide sequencing error. The estimation is performed using the Nemenman-Shafee-Bialek (NSB) estimator. The error bars for data are obtained with the NSB method, while for synthetic sequences are estimated as the empirical standard deviation over different realizations of the simulation. Due to the increased statistics, the convergence is faster for  $n < 4$  (the number of possible  $n$ -grams grows as  $20^n$ , thus we decided to show this analysis only for the case  $n = 4$ ). **(B)** Analogous analysis for TCR $\beta$  sequences.

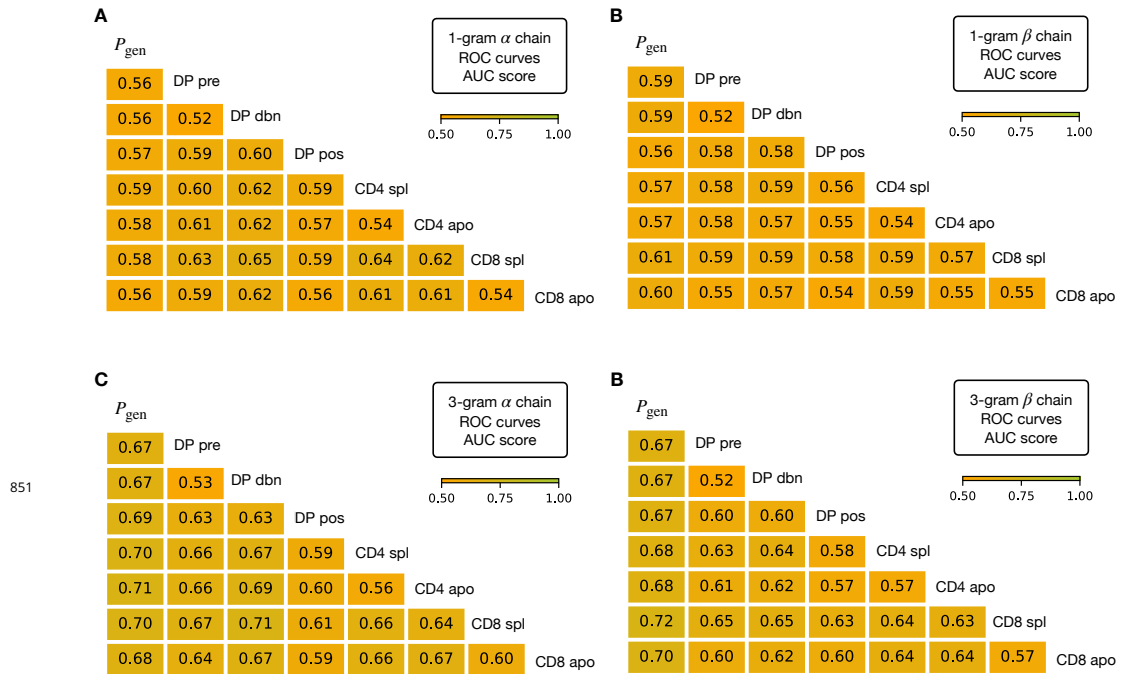


**Figure 4—figure supplement 3. (A)** NSB estimation of the Shannon entropy  $S$  normalized by  $n$ , associated to  $n$ -gram distributions within the CDR3 TCR $\beta$  chain of unique clonotypes from the different maturation stages. **(B)** Decrease in entropy per symbol between  $n$  and  $n - 1$  grams for the  $\alpha$  chain. The decreases are comparable. **(C)** For the  $\beta$  chain the decreases get smaller with  $n$ .

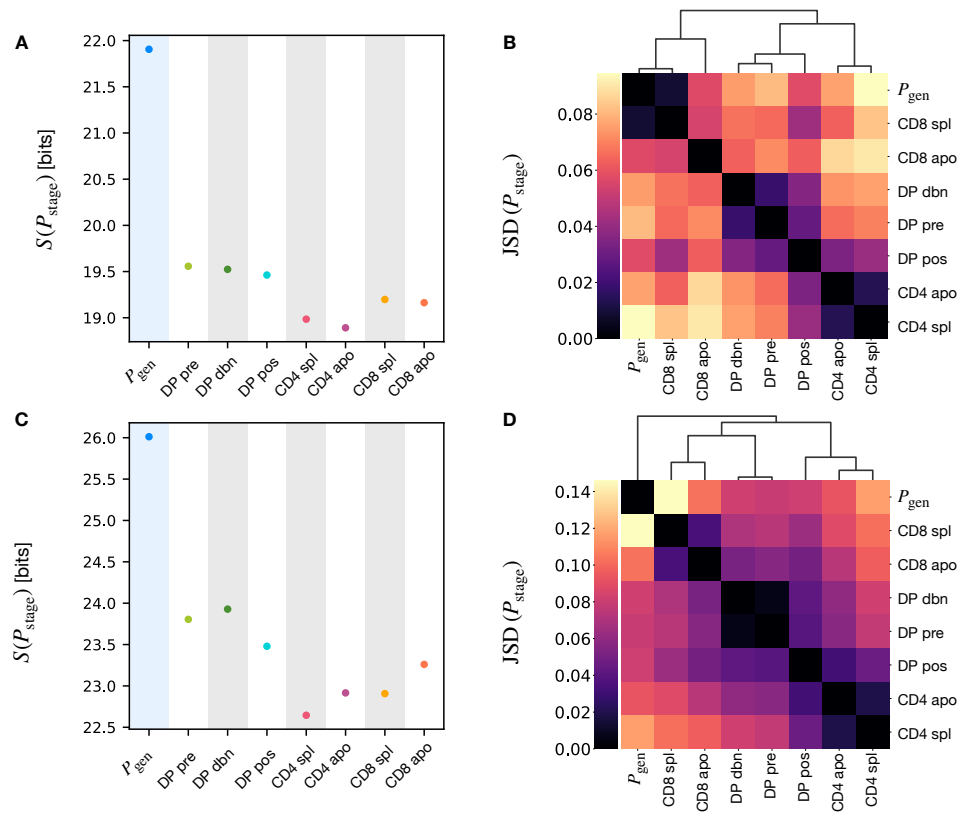




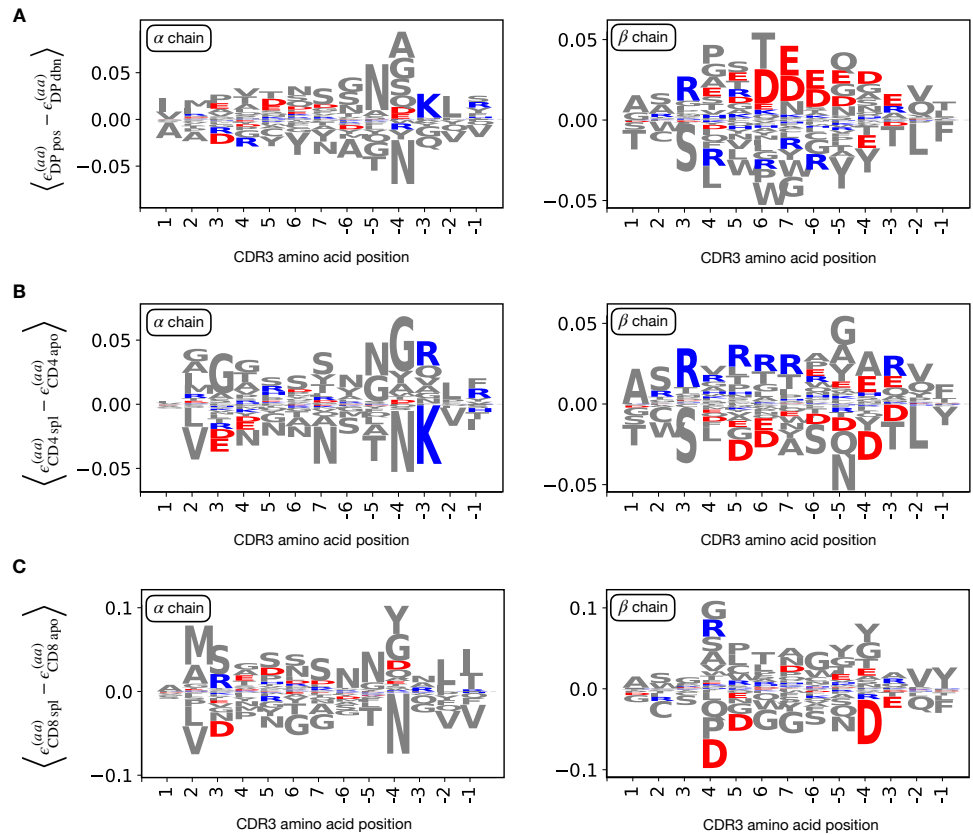
**Figure 4—figure supplement 4.** (A) Jensen-Shannon divergence (Eq. 8) for different  $n$ -gram distributions (here  $n = 1$ ) estimated on synthetic TCR $\alpha$  repertoires for different maturation stages. The dendrogram is computed with the Ward method (see Materials and Methods). (B) Divergence for TCR $\alpha$  in the case  $n = 3$ . We report here the same figure shown in the main text for the sake of comparison (Figure 4C). (C) Divergence for TCR $\alpha$  in the case  $n = 4$ . (D) Divergence for TCR $\beta$  in the case  $n = 1$ . (E) Divergence for TCR $\beta$  in the case  $n = 3$ . (F) Divergence for TCR $\beta$  in the case  $n = 4$ .



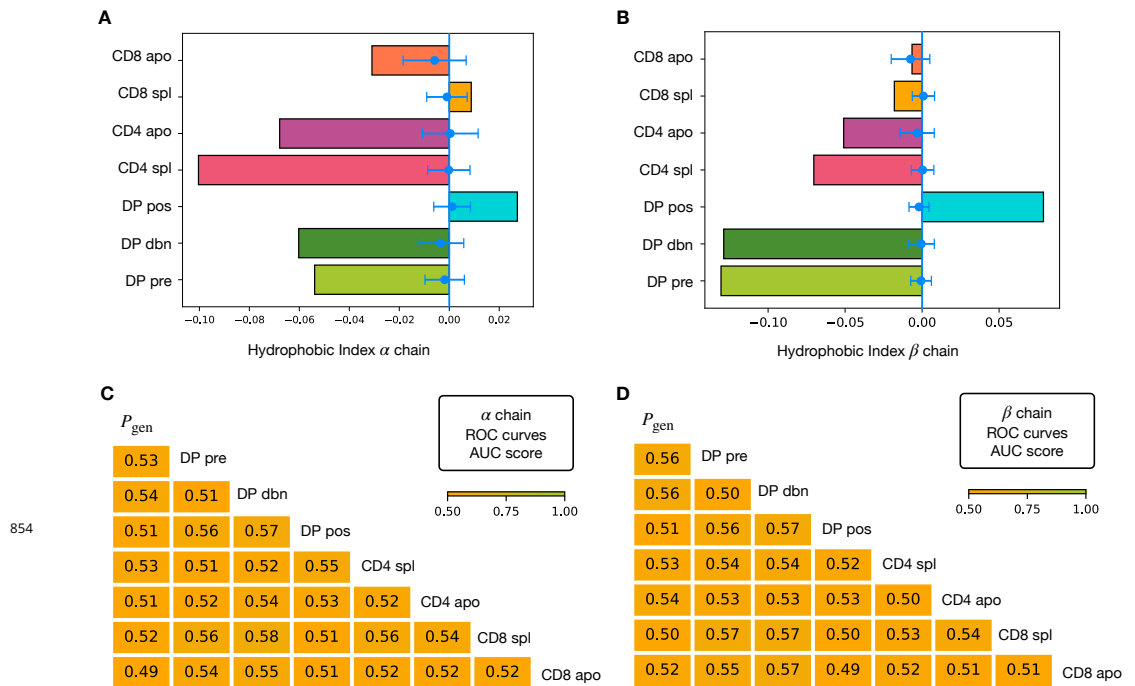
**Figure 4—figure supplement 5.** AUC values computed from the ROC curves of the linear classifiers learnt over n-grams features. **(A)** 1-gram classifiers for the  $\alpha$  chain. In the case of 1-grams, the features are assigned according to the counts of appearance of a certain amino acid within the CDR3 region (20 features). **(B)** 1-gram classifiers for the  $\beta$  chain. **(C)** 3-gram classifiers for the  $\alpha$  chain. In the case of 3-grams we choose a one-hot-encoding of the 8000 features. We observe the increased discrimination power of the 3-grams with respect to 1-gram as expected, the latter being generally worse than the models learnt on top of *Sonia* features. **(D)** 3-gram classifiers for the  $\beta$  chain.



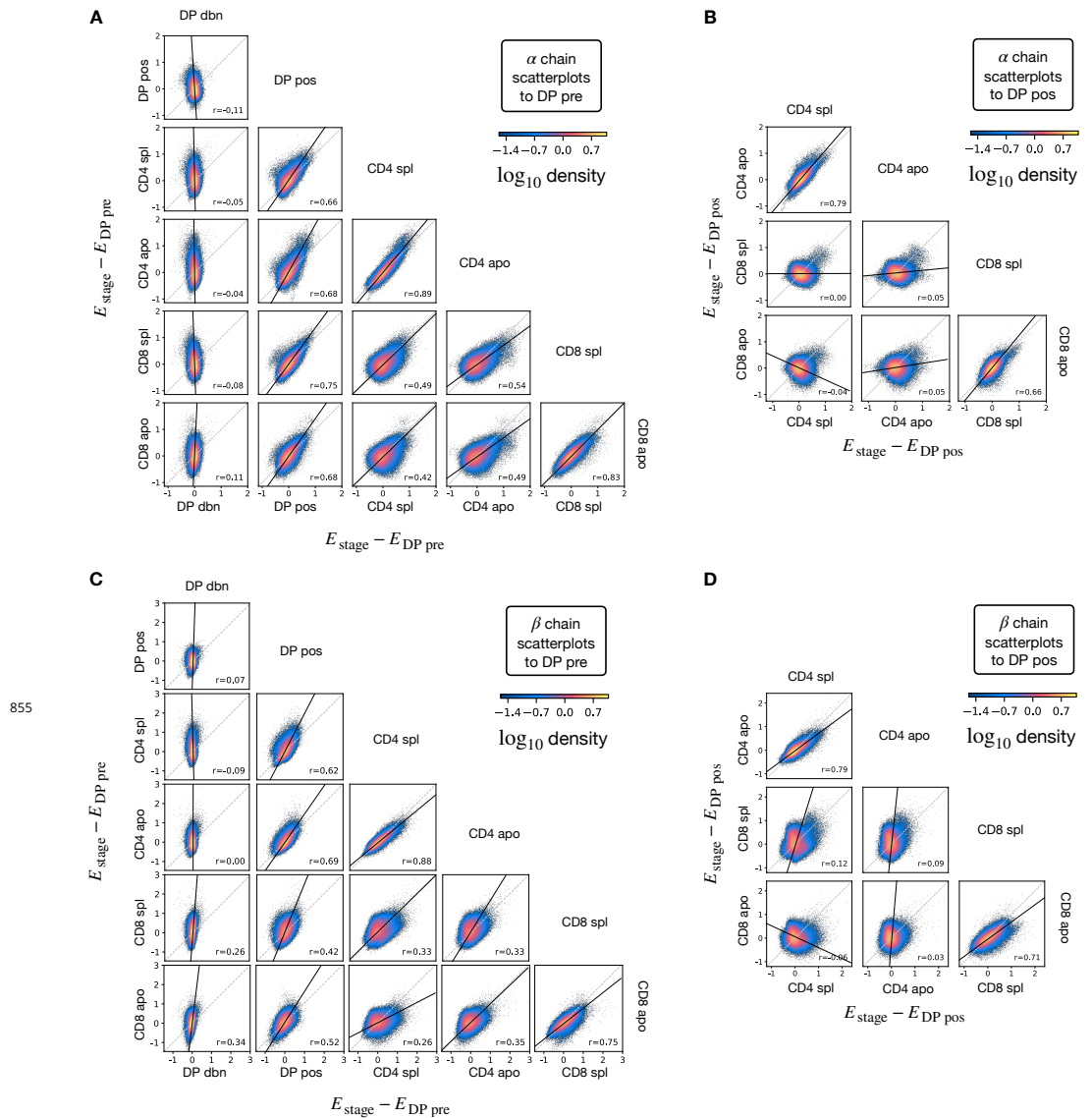
**Figure 4—figure supplement 6.** (A) Shannon entropy estimation associated to the full  $P_{\text{stage}}$  model for the  $\alpha$  chain (Eqn. 7) for the different selection models and the  $P_{\text{gen}}$  generation model. (B) We report here the same figure shown in the main text for the sake of comparison (Figure 4D). (C) Shannon entropy estimation associated to the full  $P_{\text{stage}}$  model for the  $\beta$  chain for the different selection models and the  $P_{\text{gen}}$  generation model. (D) Jensen-Shannon divergence for the  $\beta$  chain  $P_{\text{stage}}$  models.



**Figure 4—figure supplement 7. (A)** Logo plots for the CDR3 amino acid usage inferred by the model, from the left (positive position indexes) and from the right (negative), omitting the first and the last one. Here the quantity  $\langle \epsilon_1^{(aa)} - \epsilon_2^{(aa)} \rangle$  represents the average difference between weights associated to amino acid  $aa$  at the given position by  $P_{stage}$  models (see Materials and Methods). Analogously with the energy difference, a negative difference implies the feature is favoured in stage 1, vice versa for stage 2. We follow the color scheme from *Tubiana et al. (2019)* to highlight the charge properties (red for positive charge, blue for negative charge). On the left is shown the weights difference between stages DP pos and DP dbn for the  $\alpha$  chain. For the  $\beta$  chain (right) we see a reduction of positively charged amino acid in DP pos. **(B)** Stages CD4 spl and CD4 apo. Conversely, here the CD4 spl stage show enhancement in positively charged for the  $\beta$  chain (right). **(C)** Stages CD8 spl and CD8 apo. We observe just a slight enhancement of positively charged amino acid in CD8 spl.



**Figure 4—figure supplement 8. (A)** We measure the increase of hydrophobicity from the generation benchmark using a stage-wise score defined from the inferred  $P_{stage}$  models (see Materials and Methods, Eq. 4). As a control, we compute the same quantity over a set of models learnt on  $P_{gen}$ -generated repertoires of the same sizes (in blue). The score is showed at the various stages for the  $\alpha$  chain. We observe a clear increase from DP pre to DP pos and a subsequent decrease for the single positive sets, in agreement with the role of positive and negative selection. **(B)** Analogous analysis for the  $\beta$  chain. In this case we also observe AnnexinV+ sets with a higher score than the spleen sets. **(C)** AUC scores computed from the ROC curves of the logistic regression classifiers learnt over an empirical hydrophobic index of the  $\alpha$  repertoires (see Materials and Methods, Eq. 5). **(D)** AUC scores for the classifiers on on hydrophobic features for the  $\beta$  chain.



**Figure 5—figure supplement 1.** The differential enrichment parameters assigned by the stage specific selection models relative to the the preceding stage (i.e. energy differences from DP pre or DP pos). Each dot represents one of the  $3 \cdot 10^6$  synthetic sequences generated according to the generation model  $P_{\text{gen}}$ , here shown according to a dot density plot. Each figure uses the same set of synthetic sequences. **(A)** Density scatterplots of the energy differences between the energies of the TCR $\alpha$  models and the enregy of DP pre. **(B)** Density scatterplots for TCR $\alpha$  where DP pos energy is subtracted instead. **(C)** Density scatterplots for TCR $\beta$  where DP pre energy is subtracted. **(D)** Density scatterplots for TCR $\beta$  where DP pos energy is subtracted.