

Supplemental Information

Supplemental Methods:

Generation of humanized mice

NOD-scid common cytokine gamma chain knockout (NOD.Cg-Prkdc^{scid} Il2rg^{tm1Wjl}/SzJ) (NSG) mice were obtained from the Jackson Laboratory and housed in a specific pathogen-free microisolator environment. Discarded human fetal thymus and liver tissues (gestational age 17 to 20 weeks) were obtained from Advanced Biosciences Resource. Fetal thymus fragments were cryopreserved in 10% dimethyl sulfoxide and 90% human AB serum (Atlanta Biologicals). In Experiment 1, three NSG mice were sublethally irradiated (100cGy) and injected i.v with 2×10^5 human fetal liver (FL)-derived CD34⁺ cells (referred to as hematopoietic stem cells, HSCs) (Figure 1A). Autologous human fetal thymus fragments measuring about 1 mm³ were cryopreserved, thawed and transplanted under the kidney capsule of these mice, as described¹⁶. In Experiment 2, six mice received i.v injection of 2×10^5 human FL-derived CD34⁺ HSCs from another donor. Three of these mice (mice 2autoA, 2autoB and 2autoC) received autologous human fetal thymus and the other three (mice 2alloA, 2alloB and 2alloC) received an allogeneic human fetal thymus transplant (Figure 1B). In Experiment 3, two NSG mice were thymectomized, sublethally irradiated (100cGy) and injected i.v with 2×10^5 human fetal liver (FL)-derived CD34⁺ HSCs from a different donor than those used in Experiments 1 and 2 (Figure 1C). To ensure that the transplanted donor thymus T cells were not able to persist, we froze and thawed the thymus tissues and also physically removed residual cells by repeated pipetting up and down before transplantation. To further deplete passenger thymocytes that might migrate to the periphery and limit allogeneic HSC engraftment, an anti-human CD2 antibody was injected to the mice in 2 weekly doses (400µg/mouse, i.p) as we have described¹⁶. For analysis of human reconstitution, mice were bled at regular intervals for FCM analysis of human T cells, B cells and monocytes and their naïve/memory state.

FACS sorting of different subsets of grafted thymus and peripheral cells

At weeks 14, 20 and 22 after thymus transplantation, mice from Experiments 1, 2 and 3, respectively, were euthanized. Grafted thymi (for mice in all experiments) and spleen and lymph nodes (LNs) (only mice in Experiment 3) were harvested and the thymocytes and pooled spleen and LN cells were isolated by physical force (crushing the thymus tissue between two slides and crushing the spleen and LNs through a 70µm cell strainer using a syringe plunger). After counting the total number of cells, they were stained with the following antibodies for FACS sorting: anti-human CD3 (PerCP-Cy5.5), anti-human CD5 (FITC), anti-human CD4 (PE-Cy7), anti-human CD8 (APC-Cy7), anti-human CD25 (PE) and anti-human CD127 (BV421). In Experiment 2, besides staining with these antibodies, a portion of cells were stained in a separate tube with the following markers: anti-human CD3 (PerCP-Cy5.5), anti-human CD4 (PE-Cy7), anti-human CD8 (APC-Cy7), anti-human CD69 (BV650) and anti-human TCRα/β (PE). In both experiments, after gating out the dead cells and doublets, Tregs, single positive (SP) CD8 cells and Treg-depleted SP CD4 cells were sorted within a CD3⁺ CD5⁺ gate. Tregs were sorted as CD8⁻ CD25^{high} CD127⁻ CD4⁺ cells (Figure S2A). In Experiment 2, the cells in the second tube were first gated out for doublets and dead cells. Within the population of CD4 and CD8 double positive (DP) cells, CD69⁺ TCRα/β^{high} cells were sorted as positively-selected DP cells. The remaining DP cells were sorted as non-selected DP cells (Figure S2B). Thymic SP cells in Experiment 3 were sorted with the same panel as in Experiment 2. To sort peripheral (pooled spleen and LN) CD4 and CD8 cells in Experiment 3, after gating out the dead cells and doublets, CD4 and CD8 cells were sorted within a CD3⁺ gate. Sorting was done using a BD Influx cell sorter. The purity of sorted cells was %90-%96 for different cell subsets (Figure S2C).

Single cell TCR sequencing

Single cell TCR sequencing was performed according to the manuals provided by the 10X Genomics company (Chromium Single Cell 5' Library & Gel Bead Kit, PN-1000006). Briefly, after sorting thymic SP-CD4 cells, 17,000 cells from each thymus graft were loaded into the chip along with partitioning oil,

the Gel beads and a master mix containing RT enzyme and poly-dt RT primers. The assembled chips were placed into the Chromium Controller, where the cells are mixed with the beads, master mix reagents and oil. Gel Beads-in-emulsion (GEMs) were generated, where all generated cDNAs shared a common 10x Barcode. After cDNA amplification and TCR locus target enrichment, enriched libraries were constructed. Each sample was indexed with a unique barcode for each well of the Chromium i7 Index Plate. After quantifying the amplified DNA using a Bioanalyser, the same amount of DNA from different samples was pooled and sequenced with an Illumina NextSeq machine. The output files were converted to FASTQ files using the Cell Ranger pipeline. The Loupe V(D)J Browser was used for preliminary analysis within each sample. Further analysis to compare different samples was done in R. The vloupe files of the single cell TCR sequences are available at <https://github.com/Aleksobrad/Humanized-Mouse-Data>.

Computational and statistical analysis

Adaptive ImmunoSeq performs PCR amplification, read sequencing, and mapping, with bias correction and internal controls. These analyses return tabulated read counts corresponding to unique clonal CDR3 DNA sequences across all samples, and including information on the CDR3 amino acid sequence and VJ usage of these clones. From this, we normalize read counts to frequency of clonal expression for each sample on the level of distinct CDR3 nucleotide sequence, distinct CDR3 amino acid sequence, and distinct V-J pair. This repertoire characterization process is done separately for read-count tables of productive clones and nonproductive clones, which are identified as being out of frame or including a stop codon.

For each sample, then, we generate clone frequency tables at the level of non-productive nucleotide sequence, productive nucleotide sequence, amino acid sequence, and VJ usage. Template counts, clonality scores, unique clone counts and entropies for each sample were calculated. Templates are cell count estimates for each clone, derived by Adaptive ImmunoSeq in their TCR-sequencing pipeline. Each unique TCR DNA-sequence in the repertoire (unique clone) may be represented by multiple

sequenced templates, with a greater number of templates indicating a higher-frequency clone. In every sample, clonality is calculated as an inverse measure of repertoire diversity, in order to ensure that repertoires are comparable. Clonality is entropy normalized for the number of clones N , where: $\forall i$ with frequency p_i , $H_{obs} = \sum p_i \log_2 p_i$, $H_{max} = \sum \frac{1}{N} \log_2 \frac{1}{N}$ and $clonality = 1 - H_{obs}/H_{max}$ such that clonality of 1 indicates a single dominant clone, and clonality of 0 indicates uniform distribution of clone frequencies. Our definition of clonality is based on CDR3 β sequences and not the entire TCR β chain. Entropy (H) is a measure of diversity in a system, such that high-frequency clones in the repertoire decrease entropy and entropy for a sample is maximized if all clones are present at the same frequency. It is not a normalized metric and has no upper bound. Entropy is expected to be larger for larger samples, so only samples with similar numbers of unique clones can be compared in terms of relative diversity using entropy

We compared repertoires for the same cell population across mice using shared clone fraction, a non-symmetric measure such that the shared clone fraction of repertoire p compared to repertoire q is equal to the number of clonotypes present in both repertoires divided by the total number of unique clones defined in repertoire p . We alternately compared repertoires defined by their CDR3 non-productive nucleotide sequences, CDR3 productive nucleotide sequences, and CDR3 amino acid sequences for each thymic sub-population in both experiments. We also performed systematic comparison of repertoires using the Jensen-Shannon Divergence (JSD), which accounts for clone frequencies and scales for repertoire sizes. JSD is an information theory-based measure of the divergence of TCR repertoires. This is a symmetric value defined for any two repertoires p and q as: $JSD(p, q) = H_{obs}(0.5 * (p + q)) - 0.5 * (H_{obs}(p) + H_{obs}(q))$. JSD values range between 0 and 1, where 0 indicates identical repertoires, and 1 indicates complete divergence. For both shared clone fraction and JSD, we established a statistical baseline to distinguish any observed repertoire divergences across samples from divergence due to under-sampling of rare clones. This was done by $\frac{1}{4}$ sub-sampling (with replacement) of each repertoire 100 times, and computing mean and standard deviation of divergence by JSD and clone fraction when comparing all subsamples drawn from the same sample, thus

approximating divergence due to repertoire under-sampling and capturing any potential biases towards lower divergence across thymic sub-populations due to the presence of dominant high-frequency clones. All repertoire comparisons were validated for robustness to sample size differences by sub-sampling repertoires to the same low template count (2000 templates) three times each and repeating comparisons made between whole samples across the sub-samples.

We further plot the V and J gene frequencies across samples per cell population. Mann-Whitney U-tests are performed comparing the V and J distributions of different samples, as well as the observed distributions of combined VJ frequencies to the frequencies expected by stochastic pairing of 60 possible V genes with 13 possible J genes according to the background frequency of each V and J.

To identify correlations between amino acid use at P6 or P7 and hydrophobicity, the amino acid sequence data provided by Adaptive Immunoseq were tabulated for each of the five cell populations (DPCD69⁻, DPCD69⁺, SP CD4, SP CD8, SP Treg) in each animal and the amino acid and the corresponding relative frequency at P6 and at P7 was recorded for each of the CDR3 β lengths. These frequencies were normalized such that the sum of all the amino acids within a given cell population and given CDR3 β length in each mouse is one. These frequencies were subsequently chain-length matched, and the fold-change value was obtained as the ratio of the amino acid's relative frequency in cell population 2 to its relative frequency in cell population 1. The average fold change of the amino acid was determined as the numerical average of the fold changes across the mice.

To identify motifs at the sequence level comparing sequences shared between any two mice for a given cell population and sequences unique to a single mouse, a length-matched unshared sequences dataset of the same size as the shared sequences dataset was generated for each population by randomly selecting a sequence of the same length from the unshared sequence set for each sequence in the shared sequence set. Methods from Greiff et al. ²⁷, which successfully distinguished between public and private antibody repertoires were applied to this dataset to identify subsequence level

features which can be used to distinguish between shared and unshared sequences. This method uses normalized gapped k-mer (two subsequences of length k, separated by a gap of up to m amino acids) count as an input to a support vector machine to predict shared/unshared status. SVM analysis was run using $k = 1$, $m = 1$ and cost = 100, and 10 fold cross-validation was performed to assess performance of the classifier, using balanced accuracy (mean of sensitivity and specificity) as a performance metric. This was repeated on 10 length-matched datasets generated as described above. To analyze differential usage of amino acids at each position as defined by IMGT, Fisher's exact test was performed for all sequences in one length matched dataset of shared and unshared sequences. Frequency differences of amino acid and position combinations were analyzed and plotted for all cases where $p < 0.05$ by Fisher's exact test.

Supplementary Figures:

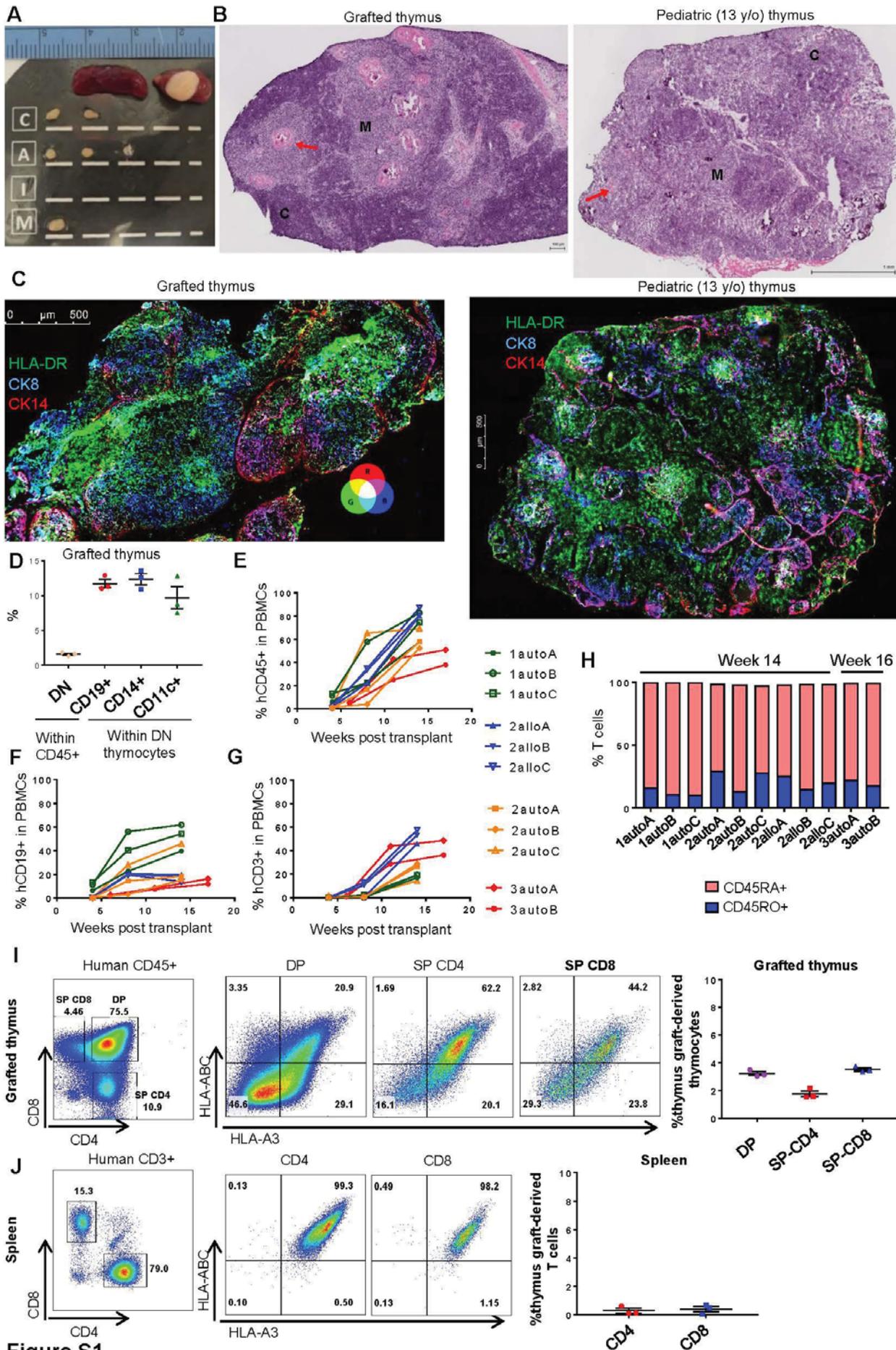


Figure S1

Figure S1. Grafted human thymus structure, the kinetics of development of human immune cells and the level of thymus graft-derived T cells. A shows a grafted human thymus under the kidney capsule, in addition to the spleen and LNs (cervical (C), axillary (A), Inguinal (I) and mesenteric (M)) of the same humanized mouse at the time of harvest (24 weeks post transplantation). B shows H&E staining of a grafted human thymus. Cortical (hypercellular areas (C)) and medullary (hypocellular areas (M)) areas as well as Hassall's corpuscles (shown by red arrows) are evident. C shows immunofluorescent staining of a grafted human thymus stained with antibodies to CK8, CK14 and HLA-DR. D shows the percentage of CD4⁻ CD8⁻ (double negative, DN) cells among human CD45⁺ cells of each grafted human thymus and the percentage of human B cells(CD19⁺), monocytes(CD14⁺) and dendritic cells(CD11c⁺) among DN cells. E-H show the kinetics of development of human immune cells (hCD45⁺), B cells (CD19⁺) and T cells (CD3⁺) in peripheral blood as well as the T cell naïve/memory phenotype. I and J show the level of thymus graft-derived T cells (HLA-ABC⁺ HLA-A3⁻) in the grafted thymi and spleens at 24 weeks of humanized mice generated with allogeneic fetal HSCs (HLA-A3⁺) and thymus tissue (HLA-A3⁻) (n=3).

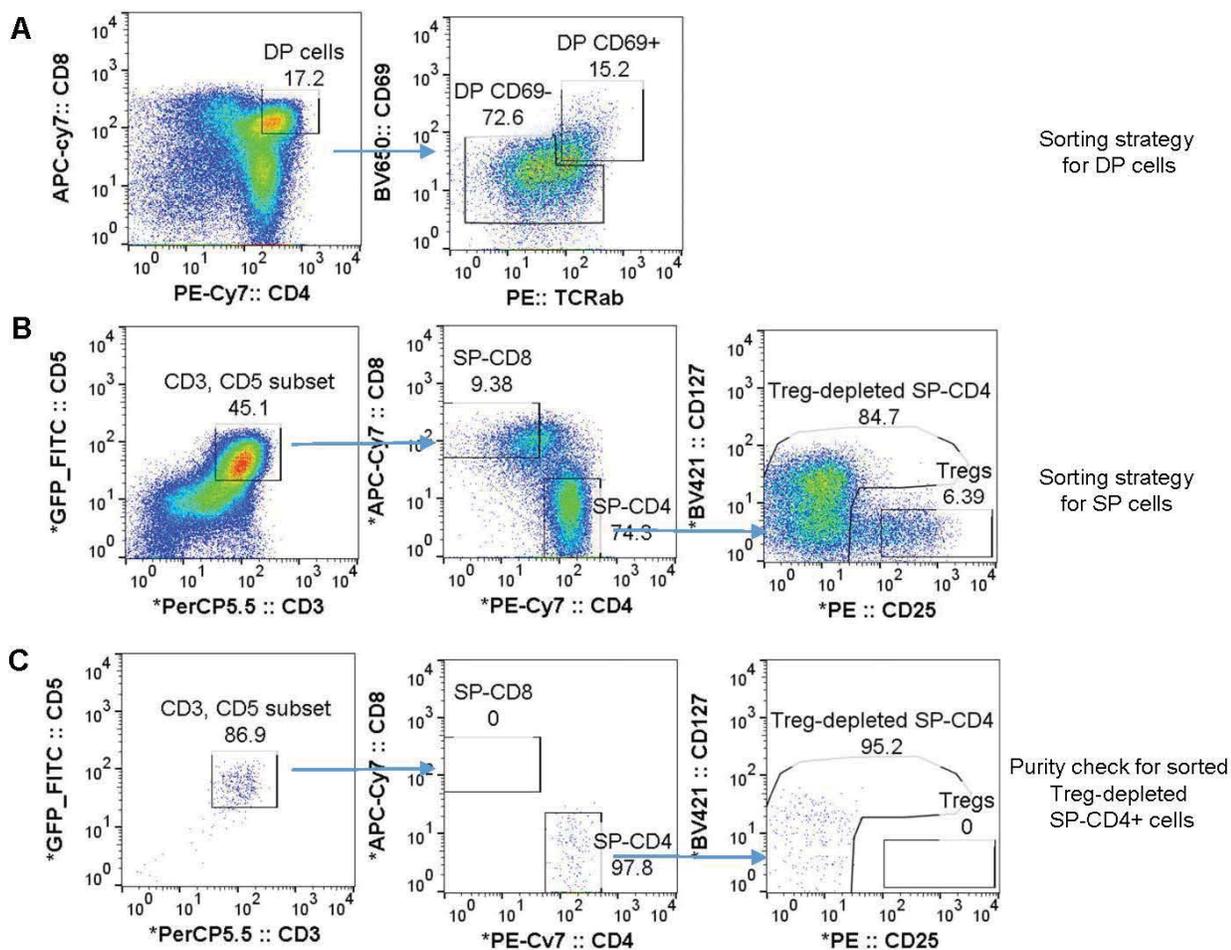


Figure S2. Sorting strategy and purity check. A shows the sorting strategy for DP cells. After gating out the dead cells and doublets, within the population of double positive (DP) CD8 and CD4 cells, CD69⁺ TCRα/β^{high} cells were sorted as positively selected DP cells. The remaining DP cells were sorted as non-selected-DP cells as shown. B shows the sorting strategy for SP cells. First, dead cells and doublets were gated out. Tregs, SP CD8 cells and Treg-depleted SP CD4 cells were sorted within a CD3^{high} CD5^{high} gate. Tregs were sorted as CD25^{high} CD127⁻ CD4⁺ cells. C shows the purity of sorted Treg-depleted SP-CD4 cells.

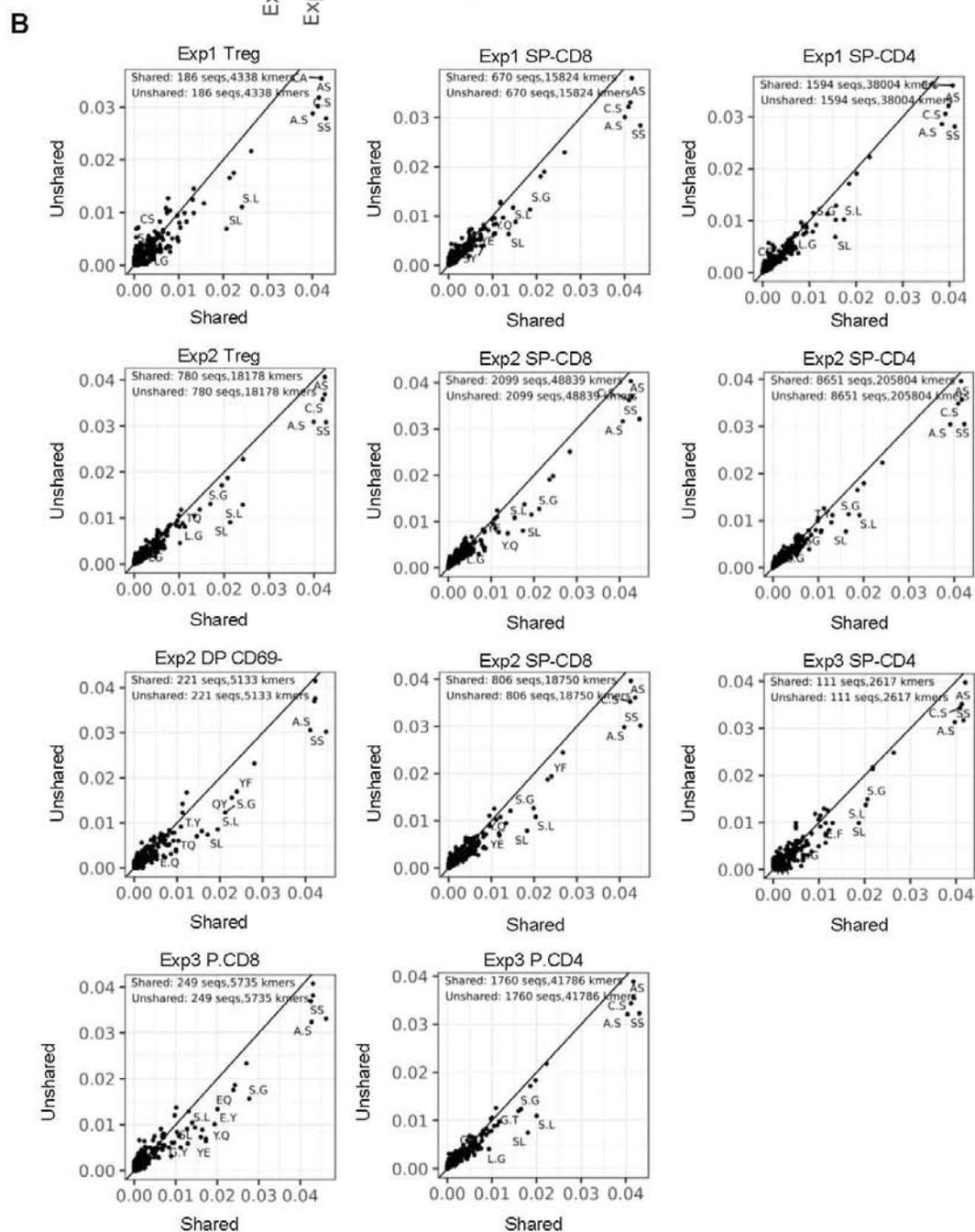
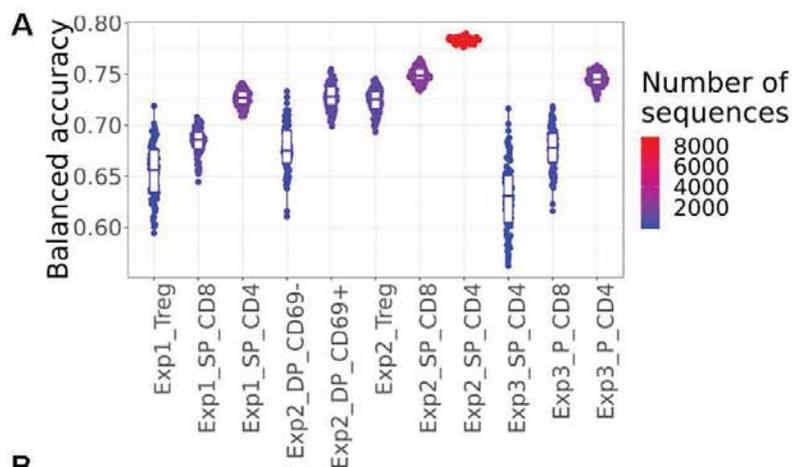


Figure S4. Subsequence features in shared vs unshared CDR3 β s. A, Support vector machine (SVM) analysis using normalized count of gapped k-mers showing that these features can be used to predict shared or unshared status of sequences with a median balanced accuracy of ~62-78% for all cell subsets and developmental stages. B, Frequency of gapped k-mers in shared sequences plotted against the frequency in unshared sequences.

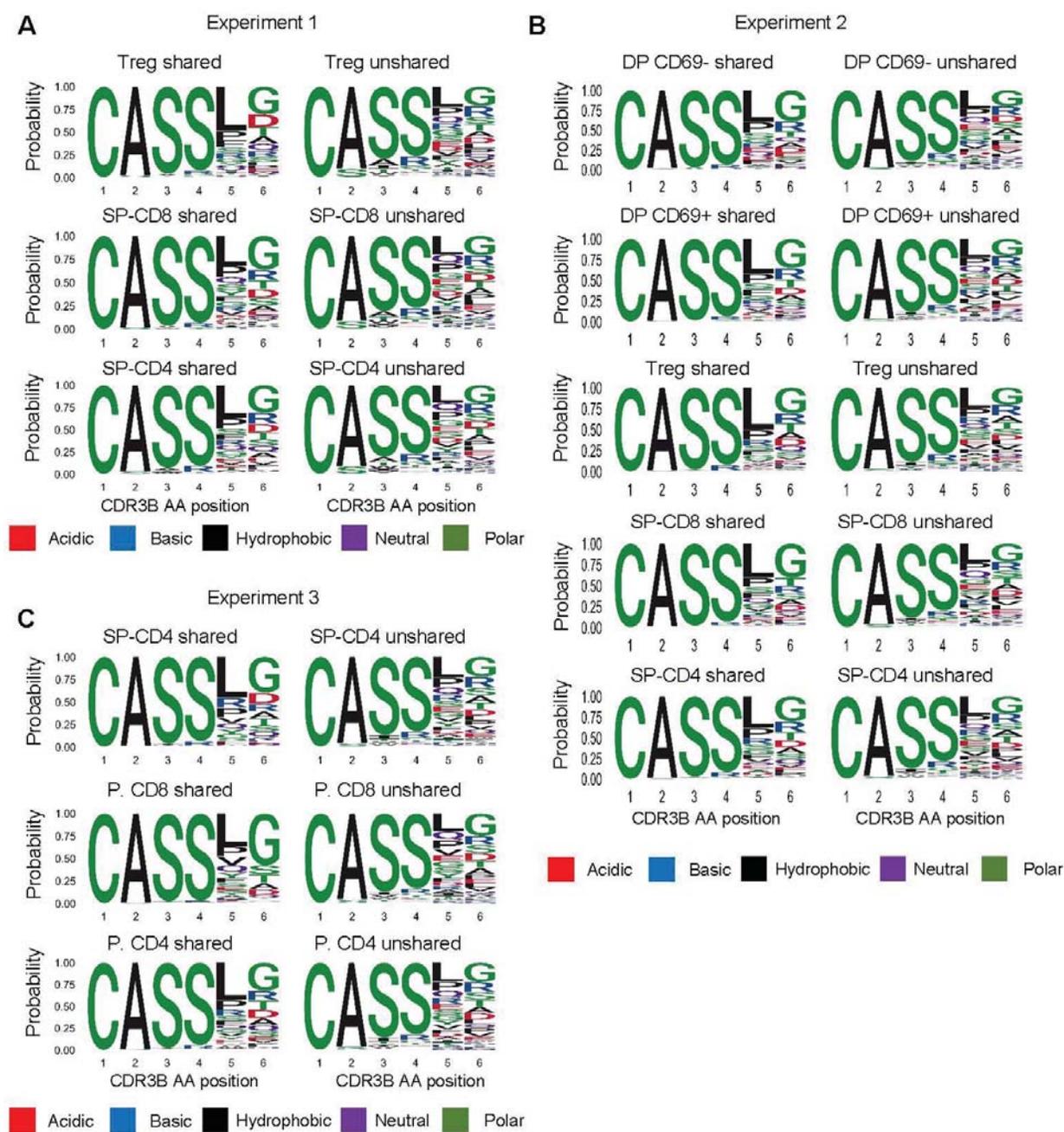


Figure S5. Amino acid usage in shared and unshared CDR3 β s. Stacked amino acid usage plots showing frequency of each amino acid in the first six positions of CDR3 β , for the set of all shared and all unshared sequences within each cell population in Experiments 1, 2 and 3. Amino acids are stacked such that the highest-usage amino acid is on top, the height of each amino acid represents its frequency, and the total frequencies of all amino acids at a given position sum to 1.

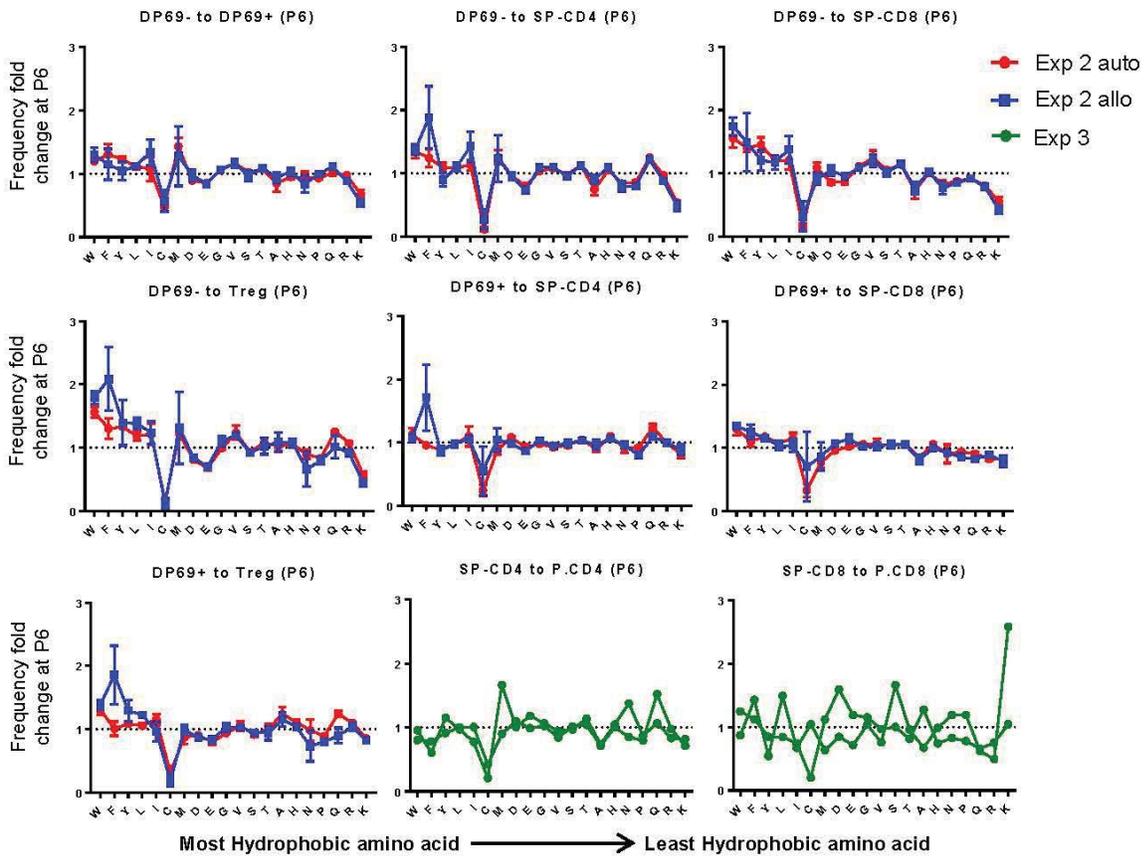
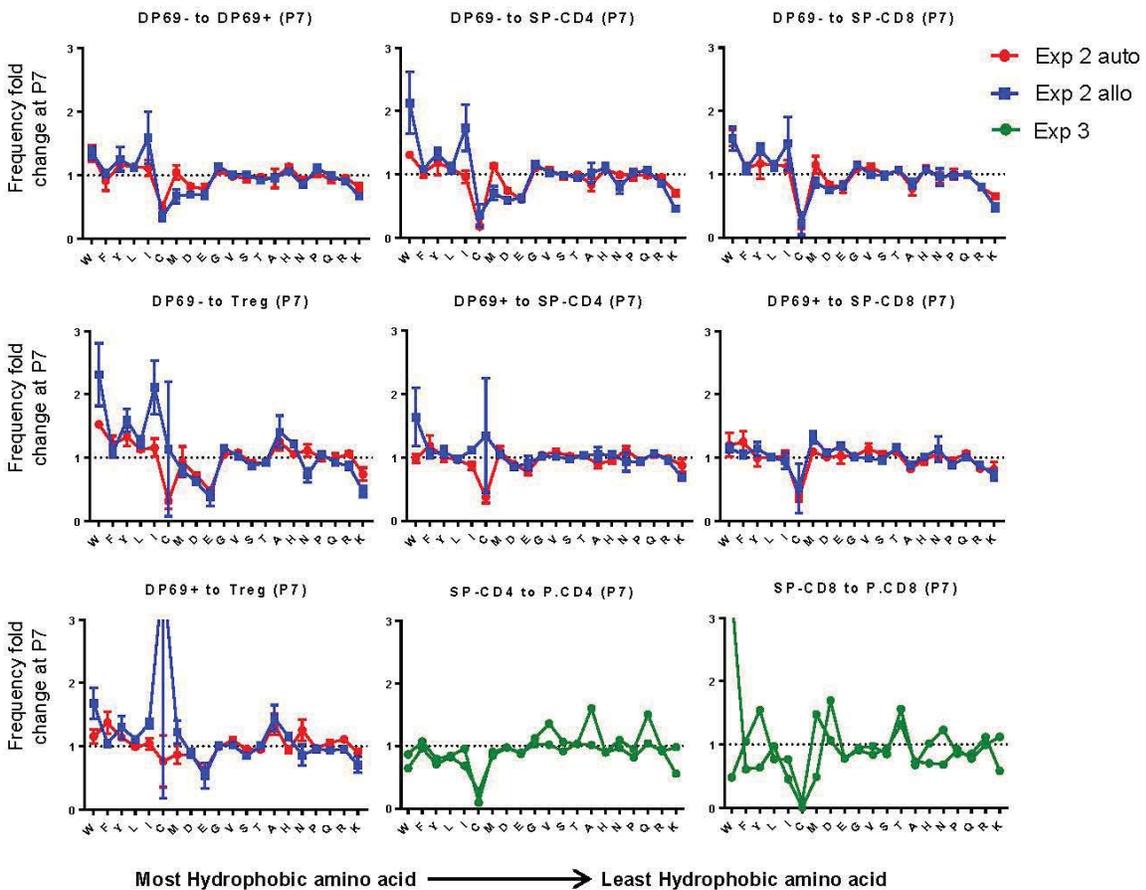
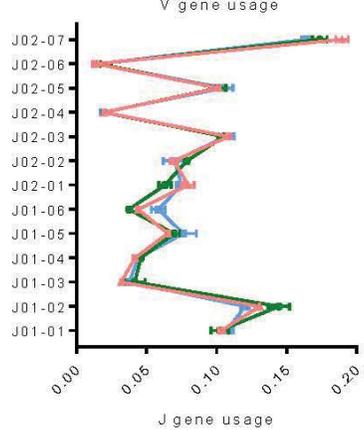
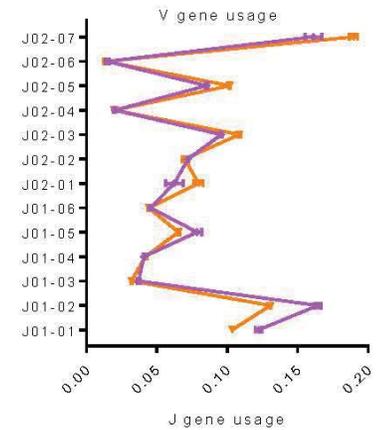
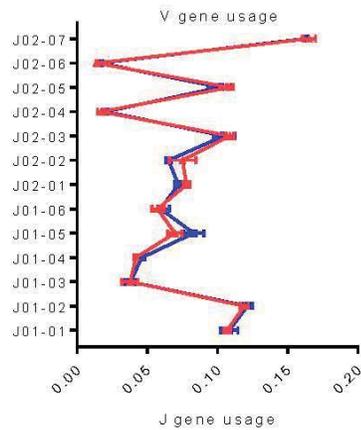
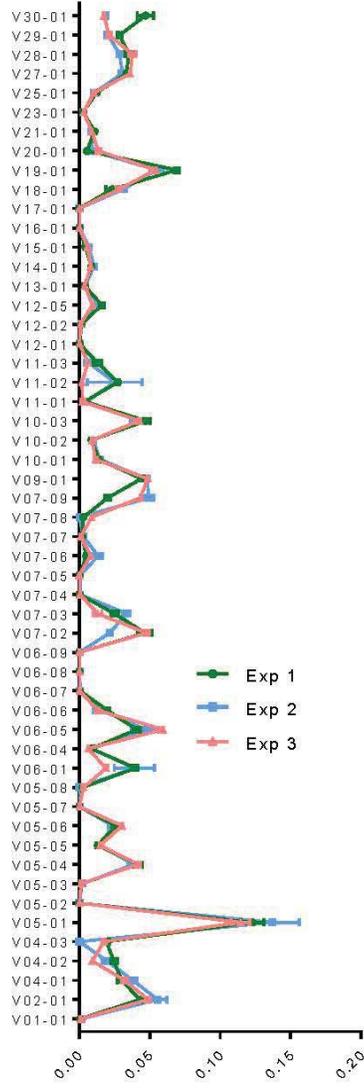
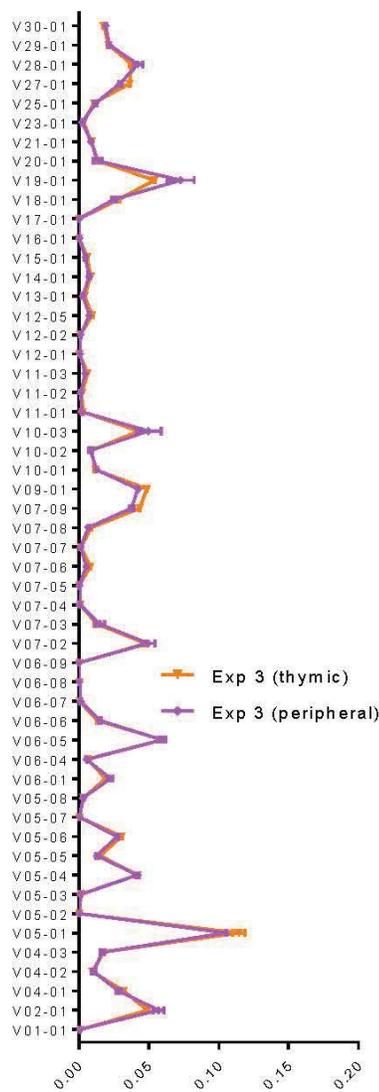
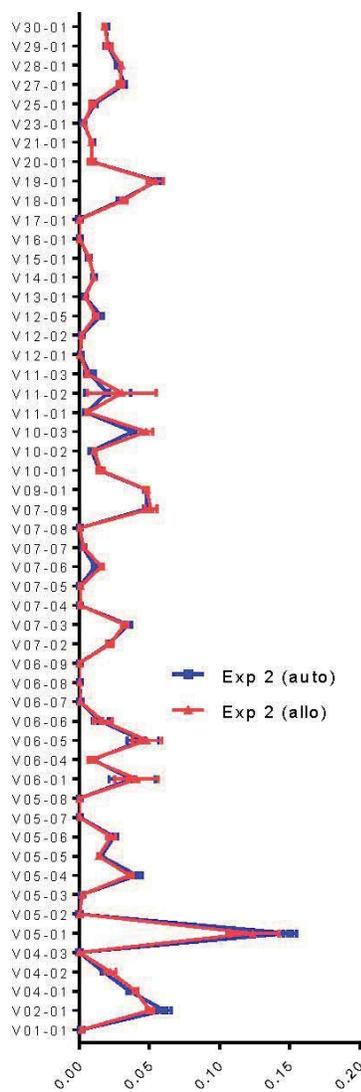
A**B**

Figure S6. Fold changes in amino acid frequency at Position 7 of the CDR3 β s. Fold changes in relative amino acid frequencies at Positions 6 (A) and 7 (B) of CDR3 β s are shown in transition from DP CD69⁻ to DP CD69⁺ cells and from there to SP cell subsets for mice with allogeneic vs autologous thymus in Experiment 2 (shown as mean \pm SEM). Also the fold changes in transition from SP-CD4 to peripheral CD4 cells and from SP-CD8 to peripheral CD8 cells is shown for mice in Experiment 3. Amino acids are listed in decreasing order of hydrophobicity from left to right.



| Exp 2 auto vs allo | | p-value |
|--------------------|------|---------|
| J genes | n.s. | n.s. |
| V genes | n.s. | n.s. |

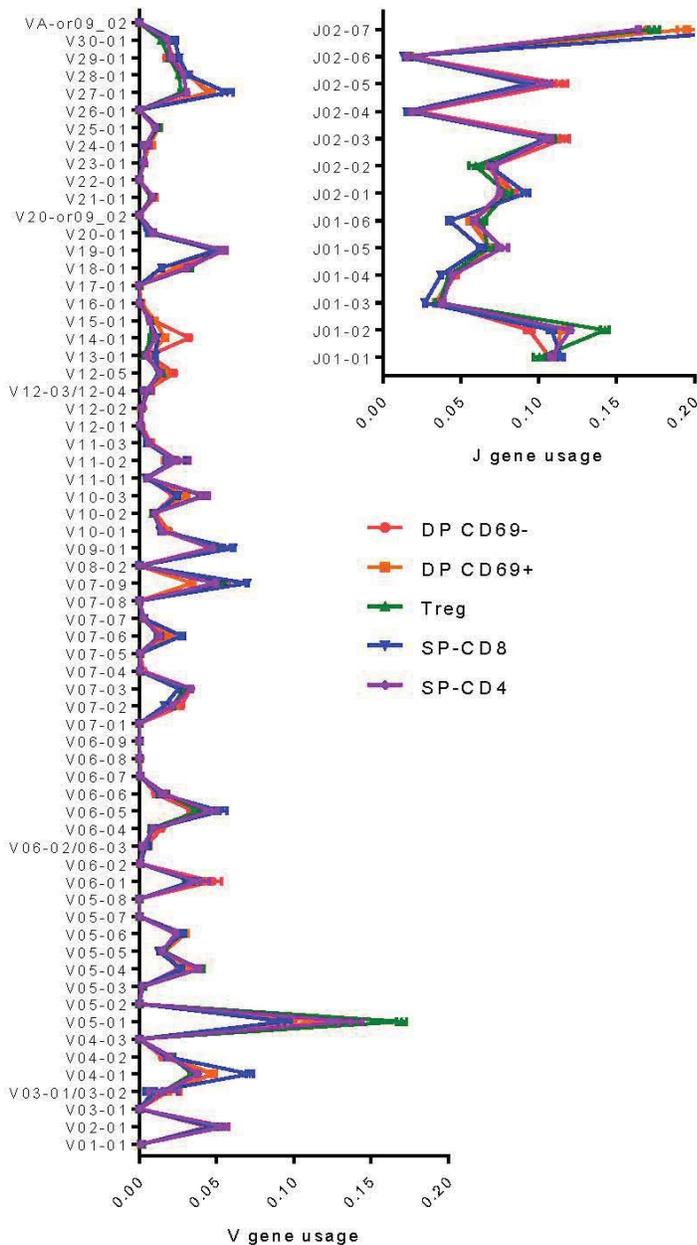
| Exp 3 thymic vs peripheral | | p-value |
|----------------------------|------|---------|
| J genes | n.s. | n.s. |
| V genes | n.s. | n.s. |

| Exp 1 vs Exp 2 | | p-value |
|----------------|--------|-------------|
| J genes | J01-06 | 0.01068645 |
| V genes | V04-03 | 0.0254639 |
| V genes | V07-09 | 0.000869265 |
| V genes | V17-01 | 0.009298685 |

| Exp 1 vs Exp 3 | | p-value |
|----------------|--------|------------|
| J genes | n.s. | n.s. |
| V genes | V06-01 | 0.02936416 |
| V genes | V07-09 | 0.04591294 |
| V genes | V11-02 | 0.01275489 |

| Exp 2 vs Exp 3 | | p-value |
|----------------|--------|------------|
| J genes | n.s. | n.s. |
| V genes | V07-03 | 0.0173452 |
| V genes | V07-08 | 0.01892448 |

Figure S7. V β and J β gene usage of thymic and peripheral CD4 cells. V β and J β gene usage distributions of SP-CD4 thymocytes of mice with allogeneic (n=3) vs autologous (n=3) thymus in Experiment 2 are shown in the left graphs. The V β and J β gene usage distributions of thymic (n=2) vs peripheral (n=2) CD4 T cells in Experiment 3 are shown in the middle graphs. The V β and J β gene usage distributions of thymic SP-CD4 T cells comparing mice in Experiments 1 (n=3), 2 (n=6) and 3 (n=2) are shown in the right graphs. Results are shown as mean \pm SEM. Unpaired t-test with Bonferroni correction for multiple testing was performed to compare the V β and J β gene usages for each gene. Genes with statistically significant differences in usage are shown in the tables below. p-value<0.05 was considered significant.



| DP69- vs DP69+ | | P-value |
|------------------|------------|-------------|
| J genes | n.s. | |
| V genes | TCRBV07-09 | 0.006875296 |
| | TCRBV12-05 | 0.020999554 |
| | TCRBV14-01 | 0.029324548 |
| DP69+ vs Treg | | P-value |
| J genes | TCRBJ02-07 | 0.03670862 |
| V genes | TCRBV14-01 | 0.016555176 |
| | TCRBV27-01 | 0.001426217 |
| DP69+ vs SP-CD8 | | P-value |
| J genes | TCRBJ02-07 | 7.81E-03 |
| V genes | TCRBV04-01 | 0.002171162 |
| | TCRBV07-09 | 0.000318317 |
| DP69+ vs SP-CD4 | | P-value |
| J genes | TCRBJ02-07 | 1.00E-02 |
| V genes | TCRBV07-09 | 0.01939807 |
| | TCRBV27-01 | 0.02853274 |
| Treg vs SP-CD8 | | P-value |
| J genes | TCRBJ01-02 | 1.15E-02 |
| | TCRBJ01-06 | 9.56E-03 |
| | TCRBJ02-07 | 1.84E-03 |
| V genes | TCRBV04-01 | 2.88E-03 |
| | TCRBV05-01 | 5.49E-05 |
| Treg vs SP-CD4 | | P-value |
| J genes | n.s. | |
| V genes | n.s. | |
| SP-CD8 vs SP-CD4 | | P-value |
| J genes | TCRBJ01-03 | 1.08E-02 |
| | TCRBJ01-06 | 2.26E-02 |
| | TCRBJ02-01 | 2.61E-03 |
| | TCRBJ02-07 | 0.00111738 |
| V genes | TCRBV04-01 | 0.02541822 |
| | TCRBV07-06 | 0.03352391 |
| | TCRBV07-09 | 0.01988107 |

Figure S8. Vβ and Jβ gene usage distributions during thymic selection. The left and right plots show the Vβ and Jβ gene usage distributions of different thymic cell populations (DP CD69⁻, DP CD69⁺, SP CD4, SP CD8 and Treg cell populations) of the 6 mice in Experiment 2. Paired t-test with Bonferroni correction for multiple testing was performed to compare the V and J gene usage between DP CD69⁻ and DP CD69⁺ cells and also between DP CD69⁺ cells and SP cell populations for each gene. Genes with statistically significant differences in usage are shown in the table to the right. p-value<0.05 was considered significant. All results are shown as mean±SEM.

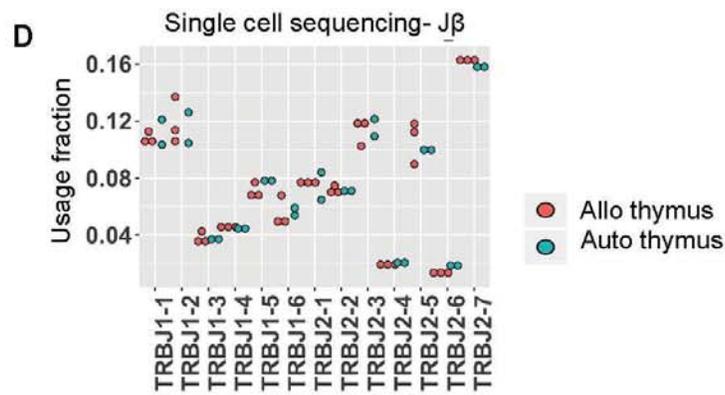
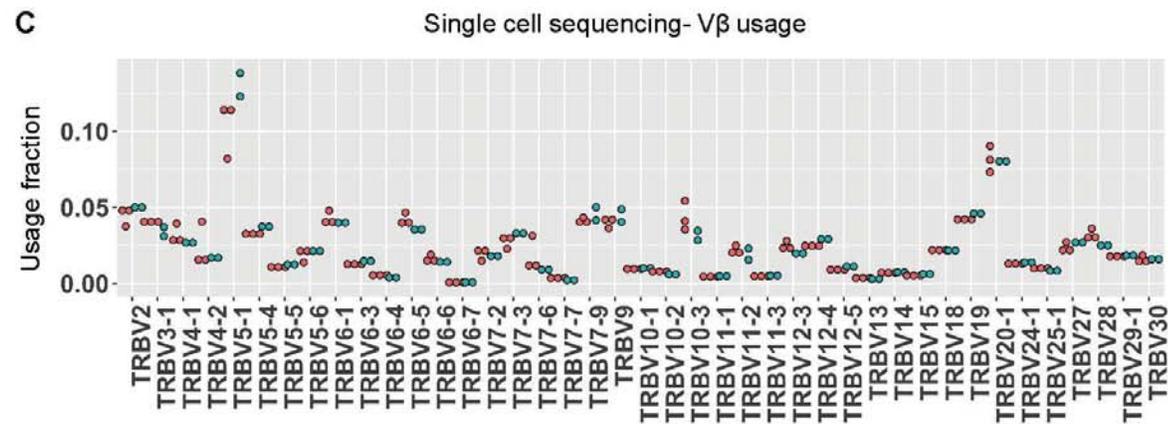
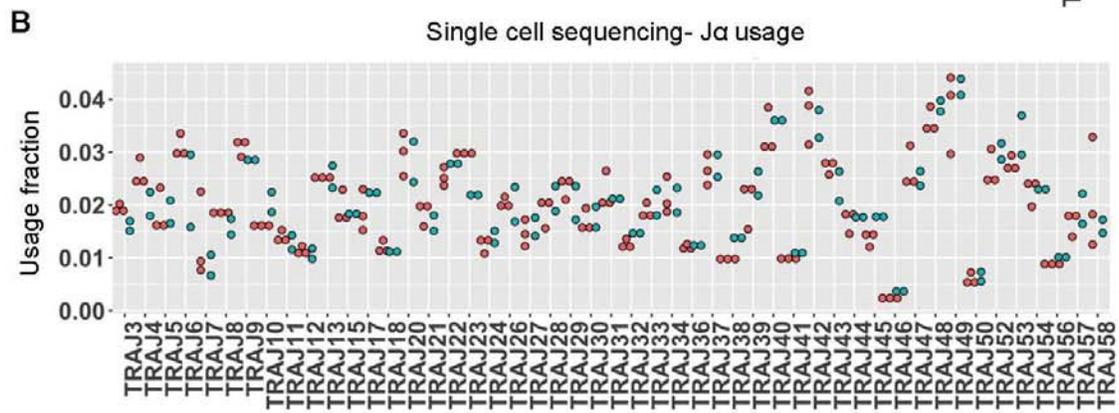
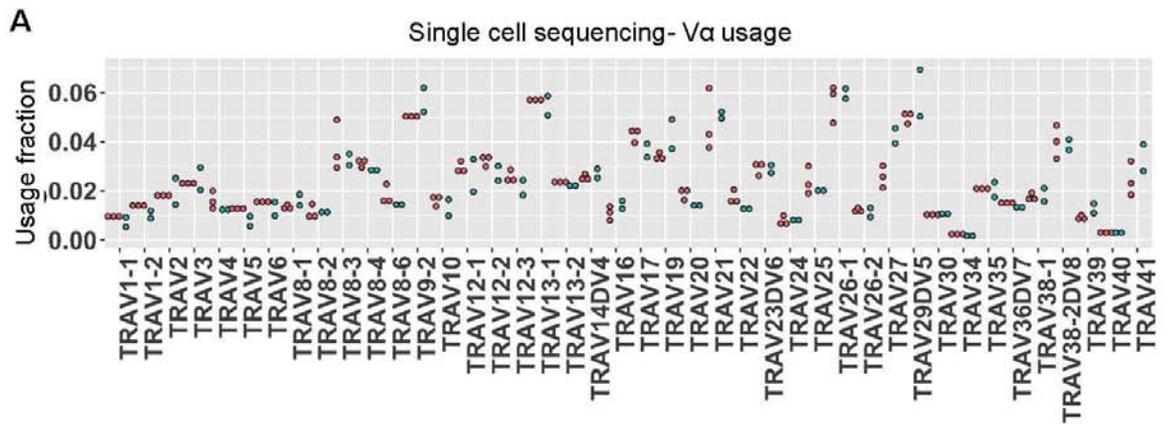


Figure S9. V and J gene usage for both α and β chains obtained in single cell TCR sequencing.

Usage of $V\alpha$, $J\alpha$, $V\beta$ and $J\beta$ genes for SP-CD4 cells of the mice in Experiment 2 with autologous (n=2) vs allogeneic (n=3) thymus obtained in single cell TCR sequencing are plotted in panels A-D.

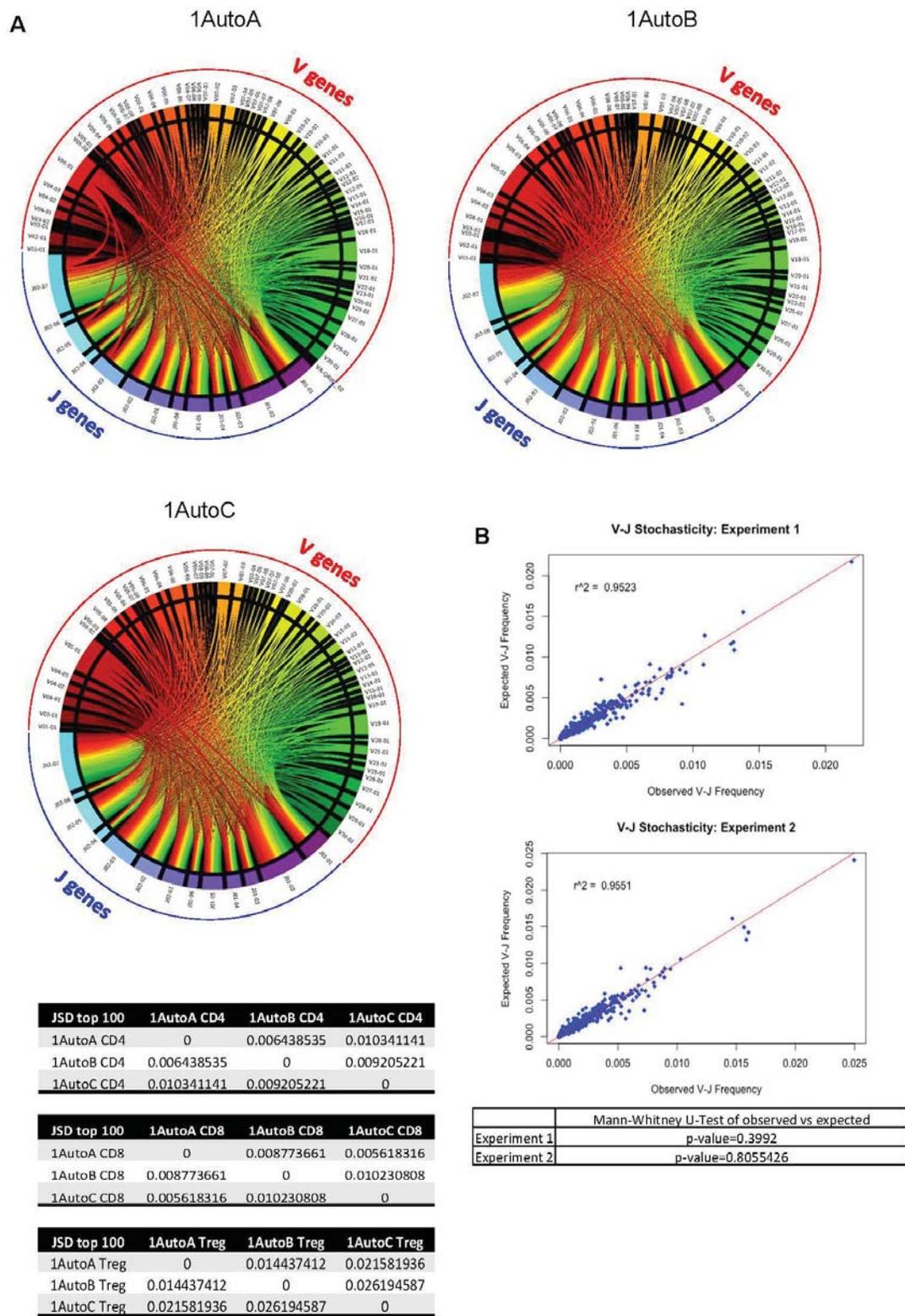


Figure S10. Pairing of V and J genes and comparison of observed vs expected VJ frequency distributions. A, Circos plots of VJ usage for SP-CD4 repertoires of the 3 mice in Experiment 1 are

shown. V genes are shown at the top and J genes are shown at the bottom. Width of a line connecting the top and bottom represents the relative frequency of a given VJ pair, and the widths of bars on the circumference represent the relative frequencies of each V and J. The color of each connecting line corresponds to a particular V gene, such that range of colors along the bottom of the plot shows diversity of Vs connecting to Js. The tables show JSD values comparing the V-J distribution between each pair of mice in Experiment 1 for different cell populations. B, Linear regressions of observed vs. expected frequency of each V-J pair are shown for the cumulative total of all samples in Experiments 1 and 2. Expected frequency is calculated as the product of V and J frequency. The correlation between the observed frequencies and the expected frequencies were plotted. Mann-Whitney U-Tests were performed with the null hypothesis that the VJ combination is stochastic.

Supplementary Tables:

Table S1. Grafted thymus cell counts, FACS-sorted cell counts, template and unique cell counts and clonality scores for each mouse in Experiments 1, 2 and 3 and the donor thymus for Experiment 2.

| Mouse/Tissue | Thymus cell count | Cell population | No. of FACS-sorted cells | nucleotide nonproductive | | | nucleotide productive | | | amino acid | |
|--------------|-------------------|-----------------|--------------------------|--------------------------|-----------|--------------------|-----------------------|-----------|--------------------|------------|--------------------|
| | | | | Template count | Clonality | Unique clone count | Template count | Clonality | Unique clone count | Clonality | Unique clone count |
| 1AutoA | 33x10e6 | Treg | 87,355 | 4,209 | 0.01 | 3,416 | 17,085 | 0.01 | 13,019 | 0.013 | 12,707 |
| | | SP_CD4 | 583,226 | 19,578 | 0.0099 | 15,375 | 74,649 | 0.014 | 53,952 | 0.023 | 49,836 |
| | | SP_CD8 | 491,310 | 13,921 | 0.0074 | 11,907 | 49,776 | 0.0085 | 41,270 | 0.012 | 38,209 |
| 1AutoB | 58x10e6 | Treg | 53,789 | 2,136 | 0.0066 | 1,875 | 7,722 | 0.006 | 6,531 | 0.0077 | 6,382 |
| | | SP_CD4 | 427,728 | 13,671 | 0.0048 | 12,288 | 49,949 | 0.0058 | 43,227 | 0.0096 | 41,134 |
| | | SP_CD8 | 314,351 | 4,535 | 0.0022 | 4,352 | 16,098 | 0.0044 | 15,195 | 0.0049 | 14,352 |
| 1AutoC | 25x10e6 | Treg | 34,048 | 1,801 | 0.0068 | 1,575 | 6,431 | 0.0076 | 5,456 | 0.0097 | 5,340 |
| | | SP_CD4 | 201,086 | 6,046 | 0.0076 | 5,298 | 21,577 | 0.011 | 17,937 | 0.016 | 17,121 |
| | | SP_CD8 | 228,890 | 8,090 | 0.0065 | 7,070 | 30,152 | 0.0074 | 25,695 | 0.012 | 24,624 |
| 2AutoA | 150x10e6 | DP_CD69- | 792,889 | 6,534 | 0.0018 | 6,266 | 18,536 | 0.0021 | 17,576 | 0.0031 | 17,080 |
| | | DP_CD69+ | 204,108 | 5,520 | 0.0012 | 5,382 | 17,783 | 0.0011 | 17,286 | 0.0025 | 16,652 |
| | | Treg | 222,288 | 7,547 | 0.0022 | 7,228 | 25,318 | 0.0031 | 24,001 | 0.005 | 22,904 |
| | | SP_CD4 | 820,703 | 29,597 | 0.0035 | 27,539 | 104,688 | 0.0041 | 95,847 | 0.0066 | 89,788 |
| | | SP_CD8 | 337,912 | 8,237 | 0.0019 | 7,879 | 29,185 | 0.0021 | 27,693 | 0.0034 | 26,853 |
| 2AutoB | 400x10e6 | DP_CD69- | 433,613 | 4,465 | 0.0015 | 4,329 | 11,661 | 0.0013 | 11,295 | 0.0022 | 11,046 |
| | | DP_CD69+ | 111,822 | 6,518 | 0.0011 | 6,355 | 18,638 | 0.0011 | 18,127 | 0.0025 | 17,473 |
| | | Treg | 121,528 | 4,344 | 0.0055 | 3,865 | 14,422 | 0.0066 | 12,457 | 0.0071 | 12,290 |
| | | SP_CD4 | 821,448 | 10,798 | 0.0015 | 10,462 | 33,300 | 0.0015 | 32,034 | 0.0036 | 30,382 |
| | | SP_CD8 | 238,420 | 10,195 | 0.0017 | 9,797 | 33,243 | 0.0018 | 31,748 | 0.0033 | 30,661 |
| 2AutoC | 100x10e6 | DP_CD69- | 200,000 | 1,691 | 0.0025 | 1,606 | 4,915 | 0.0024 | 4,644 | 0.0028 | 3,913 |
| | | DP_CD69+ | 300,000 | 2,052 | 0.0029 | 1,928 | 6,460 | 0.0026 | 6,079 | 0.0028 | 5,194 |
| | | Treg | 70,000 | 656 | 0.0032 | 618 | 2,038 | 0.0035 | 1,915 | 0.0037 | 1,647 |
| | | SP_CD4 | 300,000 | 3,397 | 0.0022 | 3,236 | 11,193 | 0.0026 | 10,531 | 0.0032 | 9,012 |
| | | SP_CD8 | 248,000 | 2,202 | 0.0029 | 2,070 | 7,471 | 0.0029 | 6,932 | 0.0035 | 6,127 |
| 2AlloA | 220x10e6 | DP_CD69- | 489,849 | 3,764 | 0.0019 | 3,616 | 9,754 | 0.0016 | 9,409 | 0.0025 | 9,183 |
| | | DP_CD69+ | 103,440 | 4,812 | 0.0014 | 4,666 | 15,408 | 0.0013 | 14,909 | 0.0021 | 14,587 |
| | | Treg | 195,514 | 6,611 | 0.003 | 6,336 | 19,930 | 0.0049 | 18,792 | 0.0067 | 17,987 |
| | | SP_CD4 | 779,000 | 22,054 | 0.002 | 21,011 | 69,169 | 0.011 | 64,132 | 0.014 | 60,352 |
| | | SP_CD8 | 302,336 | 7,012 | 0.0014 | 6,803 | 21,439 | 0.0023 | 20,551 | 0.0045 | 19,545 |
| 2AlloB | 150x10e6 | DP_CD69- | 400,317 | 628 | 0.00094 | 618 | 1,597 | 0.00099 | 1,566 | 0.0019 | 1,540 |
| | | DP_CD69+ | 101,851 | 2,913 | 0.0044 | 2,629 | 9,496 | 0.0046 | 8,357 | 0.0053 | 8,214 |
| | | Treg | 178,315 | 2,300 | 0.002 | 2,212 | 7,178 | 0.0042 | 6,812 | 0.0051 | 6,661 |
| | | SP_CD4 | 800,000 | 30,088 | 0.0034 | 27,698 | 97,477 | 0.0061 | 87,590 | 0.0085 | 82,340 |
| | | SP_CD8 | 330,315 | 2,385 | 0.0013 | 2,318 | 7,358 | 0.0013 | 7,161 | 0.0029 | 6,892 |
| 2AlloC | 168x10e6 | DP_CD69- | 200,000 | 1,047 | 0.0033 | 982 | 2,957 | 0.0026 | 2,796 | 0.0028 | 2,358 |
| | | DP_CD69+ | 300,000 | 1,940 | 0.0026 | 1,834 | 5,972 | 0.0026 | 5,619 | 0.0031 | 4,766 |
| | | Treg | 42,000 | 150 | 0.005 | 141 | 397 | 0.0017 | 386 | 0.0017 | 332 |
| | | SP_CD4 | 300,000 | 1,744 | 0.0044 | 1,619 | 5,596 | 0.0056 | 5,110 | 0.0065 | 4,347 |
| | | SP_CD8 | 156,000 | 2,399 | 0.0034 | 2,268 | 7,495 | 0.0039 | 7,004 | 0.0045 | 6,147 |
| Fetal Thymus | 15x10e6 | Treg | 23,904 | 278 | 0.0071 | 248 | 725 | 0.0069 | 633 | 0.0089 | 609 |
| | | SP_CD4 | 414,877 | 18,082 | 0.0052 | 16,030 | 44,541 | 0.017 | 33,829 | 0.049 | 23,388 |
| | | SP_CD8 | 111,449 | 8,357 | 0.0067 | 6,945 | 21,015 | 0.013 | 15,901 | 0.04 | 11,816 |
| 3AutoA | 107x10e6 | SP_CD4 | 748,000 | 6,269 | 0.0023 | 5,930 | 23,178 | 0.0027 | 21,558 | 0.0036 | 18,056 |
| | | SP_CD8 | 245,000 | 1,066 | 0.0029 | 1,010 | 4,176 | 0.0028 | 3,942 | 0.0031 | 3,436 |
| | | P_CD4 | 800,000 | 11,396 | 0.05 | 8,082 | 48,326 | 0.067 | 31,767 | 0.07 | 26,553 |
| | | P_CD8 | 800,000 | 7,423 | 0.097 | 5,052 | 30,741 | 0.13 | 18,866 | 0.14 | 16,420 |
| 3AutoB | 175x10e6 | SP_CD4 | 437,000 | 2,124 | 0.0021 | 2,033 | 7,752 | 0.0029 | 7,326 | 0.0036 | 6,133 |
| | | SP_CD8 | 144,000 | 765 | 0.0037 | 714 | 2,477 | 0.0032 | 2,307 | 0.0039 | 1,954 |
| | | P_CD4 | 800,000 | 33,554 | 0.06 | 22,077 | 124,806 | 0.06 | 80,609 | 0.06 | 65,157 |

| | | | | | | | | | | | |
|--|--|---------------|---------|-------|------|-------|--------|------|-------|------|-------|
| | | P. CD8 | 600,000 | 4,417 | 0.14 | 2,587 | 15,888 | 0.14 | 9,011 | 0.15 | 7,823 |
|--|--|---------------|---------|-------|------|-------|--------|------|-------|------|-------|

Table S2. HLA typing of fetal tissues used to generate humanized mice in all three experiments.

Typing for some MHCs are not done (shown as ND).

| HLA | Experiment 1 thymus and HSCs | Experiment 2 HSCs and auto thymus | Experiment 2 allo thymus | Experiment 3 thymus and HSCs |
|----------------|---|--|-------------------------------------|---|
| A1 | 02 | 02:01:01 | 31:01:02 | 68:01:01:01 |
| A2 | 24 | 23:01:01 | 24:02:01 | 34:02:01 |
| B1 | ND | 49:01 | 35:12:02 | 49:01 |
| B2 | ND | 35:08 | 39:06:02 | 81:01 |
| C1 | ND | 07:01:01 | 04:01:01 | 07:01:01 |
| C2 | ND | 04:01:01 | 07:02:01 | 08:04:01 |
| DRB1-1 | 4:01 | 11:04:01 | 08:02:01 | 08:02:01 |
| DRB1-2 | 4:02 | 04:05:01 | 14:06:01 | 14:06:01 |
| DQB1-24 | 3:02 | 03:01:01 | 04:02:01 | 03:01:01:01 |
| DQB1-24 | ND | 03:02:01 | 03:01:01:01 | 06:09:01 |
| DQA1 | 3:01 | ND | 04:01:01 | 04:01:01 |
| DQA2 | ND | ND | 05:01:01:01 | 01:02:01:01 |

Table S3. Jensen-Shannon Divergence (JSD) scores comparing each pair of mice in Experiment 2 for different cell subsets (amino acid level).

| JSD | Experiment 2: amino acid level | | | | | |
|---------------|---------------------------------------|---------------|---------------|---------------|---------------|---------------|
| DP69- | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 0 | 0.99431849 | 0.99554971 | 0.99515414 | 0.99841448 | 0.99731111 |
| 2autoB | 0.99431849 | 0 | 0.99830573 | 0.99534656 | 0.99876902 | 0.99712877 |
| 2autoC | 0.99554971 | 0.99830573 | 0 | 0.99692997 | 0.99919302 | 1 |
| 2alloA | 0.99515414 | 0.99534656 | 0.99692997 | 0 | 0.99930653 | 0.99734988 |
| 2alloB | 0.99841448 | 0.99876902 | 0.99919302 | 0.99930653 | 0 | 0.99890367 |
| 2alloC | 0.99731111 | 0.99712877 | 1 | 0.99734988 | 0.99890367 | 0 |
| DP69+ | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 0 | 0.98814838 | 0.99251902 | 0.98885698 | 0.99198406 | 0.99382505 |
| 2autoB | 0.98814838 | 0 | 0.99179387 | 0.98878889 | 0.99057187 | 0.99351998 |
| 2autoC | 0.99251902 | 0.99179387 | 0 | 0.9945415 | 0.99315126 | 0.99330556 |
| 2alloA | 0.98885698 | 0.98878889 | 0.9945415 | 0 | 0.99239498 | 0.99356011 |
| 2alloB | 0.99198406 | 0.99057187 | 0.99315126 | 0.99239498 | 0 | 0.99536621 |
| 2alloC | 0.99382505 | 0.99351998 | 0.99330556 | 0.99356011 | 0.99536621 | 0 |
| Treg | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 0 | 0.98747171 | 0.99497389 | 0.97855264 | 0.98574743 | 0.9986055 |
| 2autoB | 0.98747171 | 0 | 0.99559425 | 0.98666388 | 0.9930674 | 0.99762974 |
| 2autoC | 0.99497389 | 0.99559425 | 0 | 0.99612371 | 0.99612923 | 0.99888098 |
| 2alloA | 0.97855264 | 0.98666388 | 0.99612371 | 0 | 0.98481821 | 0.99807573 |
| 2alloB | 0.98574743 | 0.9930674 | 0.99612923 | 0.98481821 | 0 | 0.9990455 |
| 2alloC | 0.9986055 | 0.99762974 | 0.99888098 | 0.99807573 | 0.9990455 | 0 |
| SP-CD8 | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 0 | 0.96974267 | 0.98418995 | 0.97236814 | 0.98337473 | 0.98485022 |
| 2autoB | 0.96974267 | 0 | 0.98420004 | 0.96871132 | 0.9833744 | 0.98465465 |
| 2autoC | 0.98418995 | 0.98420004 | 0 | 0.98618989 | 0.99020655 | 0.99030095 |
| 2alloA | 0.97236814 | 0.96871132 | 0.98618989 | 0 | 0.9822219 | 0.98472348 |
| 2alloB | 0.98337473 | 0.9833744 | 0.99020655 | 0.9822219 | 0 | 0.99007559 |
| 2alloC | 0.98485022 | 0.98465465 | 0.99030095 | 0.98472348 | 0.99007559 | 0 |
| SP-CD4 | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 0 | 0.96057632 | 0.97953517 | 0.95198027 | 0.94699361 | 0.98754906 |
| 2autoB | 0.96057632 | 0 | 0.98143646 | 0.96208864 | 0.96339328 | 0.98903941 |
| 2autoC | 0.97953517 | 0.98143646 | 0 | 0.97866014 | 0.97861591 | 0.99206847 |
| 2alloA | 0.95198027 | 0.96208864 | 0.97866014 | 0 | 0.94102448 | 0.98845932 |
| 2alloB | 0.94699361 | 0.96339328 | 0.97861591 | 0.94102448 | 0 | 0.98718573 |
| 2alloC | 0.98754906 | 0.98903941 | 0.99206847 | 0.98845932 | 0.98718573 | 0 |

Table S4. Fraction of shared CDR3 β s comparing each pair of mice in Experiment 2 for different cell subsets (amino acid level).

| Shared CDR3B Fraction | Experiment 2: amino acid level | | | | | |
|-----------------------|--------------------------------|---------------|---------------|---------------|---------------|---------------|
| DP69- | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 1 | 0.00431003 | 0.00241069 | 0.00343341 | 0.00043831 | 0.00124187 |
| 2autoB | 0.00671752 | 1 | 0.00113856 | 0.00375726 | 0.00045543 | 0.00136628 |
| 2autoC | 0.00843343 | 0.00255558 | 1 | 0.00383338 | 0.00051112 | 0 |
| 2alloA | 0.00645427 | 0.00453172 | 0.00205987 | 1 | 0.00027465 | 0.00151057 |
| 2alloB | 0.00494234 | 0.00329489 | 0.00164745 | 0.00164745 | 1 | 0.00164745 |
| 2alloC | 0.0072095 | 0.00508906 | 0 | 0.00466497 | 0.00084818 | 1 |
| DP69+ | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 1 | 0.01134866 | 0.00461268 | 0.00966467 | 0.00541807 | 0.00366086 |
| 2autoB | 0.01091396 | 1 | 0.00471765 | 0.00936488 | 0.00647796 | 0.00366146 |
| 2autoC | 0.01212938 | 0.0128995 | 1 | 0.00827878 | 0.00750866 | 0.00616095 |
| 2alloA | 0.01126857 | 0.01135394 | 0.00367082 | 1 | 0.00571965 | 0.00401229 |
| 2alloB | 0.01100045 | 0.01367623 | 0.00579753 | 0.00995986 | 1 | 0.00371637 |
| 2alloC | 0.01049098 | 0.01091062 | 0.00671423 | 0.00986152 | 0.00524549 | 1 |
| Treg | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 1 | 0.00883925 | 0.00169375 | 0.01503202 | 0.00529297 | 0.00037051 |
| 2autoB | 0.01668832 | 1 | 0.00179874 | 0.01438993 | 0.00479664 | 0.00059958 |
| 2autoC | 0.01942927 | 0.01092896 | 1 | 0.01275046 | 0.00667881 | 0.00060716 |
| 2alloA | 0.0194774 | 0.00987587 | 0.00144023 | 1 | 0.00596667 | 0.00054866 |
| 2alloB | 0.01842639 | 0.00884467 | 0.0020269 | 0.01603096 | 1 | 0.00036853 |
| 2alloC | 0.02108434 | 0.01807229 | 0.00301205 | 0.02409639 | 0.0060241 | 1 |
| CD8 | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 1 | 0.02638279 | 0.00763828 | 0.018964 | 0.00759438 | 0.00741879 |
| 2autoB | 0.02315635 | 1 | 0.00755182 | 0.02053633 | 0.00716653 | 0.00732064 |
| 2autoC | 0.02839889 | 0.03198955 | 1 | 0.02072793 | 0.00897666 | 0.00930308 |
| 2alloA | 0.02617547 | 0.0322952 | 0.0076951 | 1 | 0.00872516 | 0.0084222 |
| 2alloB | 0.03007127 | 0.03233096 | 0.00956023 | 0.02503042 | 1 | 0.0100817 |
| 2alloC | 0.02749309 | 0.03090939 | 0.00927282 | 0.02261266 | 0.0094355 | 1 |
| CD4 | 2autoA | 2autoB | 2autoC | 2alloA | 2alloB | 2alloC |
| 2autoA | 1 | 0.01967026 | 0.00701926 | 0.0316547 | 0.0401159 | 0.0033736 |
| 2autoB | 0.05815637 | 1 | 0.01033623 | 0.04464286 | 0.05003218 | 0.00494691 |
| 2autoC | 0.05725699 | 0.02851753 | 1 | 0.04727031 | 0.05614736 | 0.00554816 |
| 2alloA | 0.04780099 | 0.0228015 | 0.00875085 | 1 | 0.0457468 | 0.00382079 |
| 2alloB | 0.0444354 | 0.01874454 | 0.00762439 | 0.03355634 | 1 | 0.00361631 |
| 2alloC | 0.05705084 | 0.02829538 | 0.01150219 | 0.04278813 | 0.05521049 | 1 |

Table S5. Mean NT/AA ratio for shared vs unshared CDR3βs.

| mean NT/AA | shared | unshared |
|------------|----------|----------|
| DP69- | 2.171946 | 1.024848 |
| DP69+ | 2.380893 | 1.028829 |
| Treg | 2.421795 | 1.038576 |
| CD8 | 2.71512 | 1.03049 |
| CD4 | 2.954687 | 1.053213 |

Table S6. Comparison of NT/AA ratio of shared CDR3βs between different cell subsets.

| NT/AA ratio (shared sequences) | p-value |
|--------------------------------|-----------|
| DP69- vs DP69+ | 0.0002077 |
| DP69- vs CD4 | < 2.2e-16 |
| DP69- vs CD8 | < 2.2e-16 |
| DP69- vs Treg | 1.06E-05 |
| DP69+ vs CD4 | < 2.2e-16 |
| DP69+ vs CD8 | 3.73E-13 |
| DP69+ vs Treg | 0.3618 |

Table S7. Number of shared and unshared CDR3βs at the Nt/non-productive, Nt/productive and amino acid levels for different cell subsets of the mice in Experiment 1.

| Experiment 1 | Treg | | | SP-CD8 | | | SP-CD4 | | |
|-------------------|-------------------|---------------|------------|-------------------|---------------|------------|-------------------|---------------|------------|
| | Nt Non-productive | Nt productive | Amino acid | Nt Non-productive | Nt productive | Amino acid | Nt Non-productive | Nt productive | Amino acid |
| 1autoA | 3,415 | 13,007 | 12,544 | 11,906 | 41,238 | 37,674 | 15,372 | 53,583 | 48,543 |
| 1autoB | 1,874 | 6,520 | 6,254 | 4,352 | 15,168 | 13,997 | 12,285 | 43,140 | 39,877 |
| 1autoC | 1,575 | 5,451 | 5,249 | 7,069 | 25,276 | 24,157 | 5,296 | 17,895 | 16,408 |
| 1autoA and 1autoB | 0 | 9 | 95 | 0 | 18 | 203 | 2 | 70 | 881 |
| 1autoA and 1autoC | 0 | 3 | 58 | 0 | 14 | 315 | 1 | 25 | 337 |
| 1autoB and 1autoC | 0 | 2 | 23 | 0 | 9 | 135 | 1 | 13 | 301 |

| | | | | | | | | | |
|---|---|---|----|---|---|----|---|---|----|
| 1autoA and 1autoB and 1autoC | 0 | 0 | 10 | 0 | 0 | 17 | 0 | 4 | 75 |
|---|---|---|----|---|---|----|---|---|----|

Table S8. Number of shared and unshared CDR3βs at amino acid level for different cell subsets of the mice in Experiment 2.

| CDR3B unique sequences (amino acid level) | DP CD69- | DP CD69+ | SP CD8 | SP CD4 | Treg |
|--|---------------------|---------------------|---------------|---------------|-------------|
| 2autoA | 17,080 | 16,652 | 26,853 | 89,788 | 22,904 |
| 2autoB | 11,046 | 17,473 | 30,661 | 30,382 | 12,290 |
| 2autoC | 3,913 | 5,194 | 6,127 | 9,012 | 1,647 |
| 2alloA | 9,183 | 14,587 | 19,545 | 60,352 | 17,987 |
| 2alloB | 1,540 | 8,214 | 6,892 | 82,340 | 6,661 |
| 2alloC | 2,358 | 4,766 | 6,147 | 4,347 | 332 |
| shared between all three auto | 3 | 9 | 51 | 128 | 2 |
| shared between all three allo | 1 | 6 | 23 | 80 | 1 |
| shared between all six | 0 | 1 | 5 | 8 | 0 |

Table S9. Number of unshared CDR3βs and also CDR3βs that are shared between all mice in each experiment at the nucleotide and amino acid levels for each cell subset.

| # unique CDR3B sequences | | Nucleotide | | Amino acid | |
|---------------------------------|--------------------|-------------------|---------------|-------------------|---------------|
| experiment | cell subset | non-shared | shared | non-shared | shared |
| Exp1 | SP CD4 | 114,888 | 112 | 104,828 | 1,594 |
| | SP CD8 | 82,078 | 41 | 75,828 | 670 |
| | Treg | 24,978 | 14 | 24,047 | 186 |
| Exp2 | DP CD69- | 47,261 | 11 | 36,783 | 221 |
| | DP CD69+ | 70,310 | 33 | 54,529 | 806 |
| | SP CD4 | 294,009 | 588 | 207,188 | 8,651 |
| | SP CD8 | 100,893 | 95 | 78,583 | 2,099 |
| | Treg | 64,271 | 44 | 49,253 | 780 |
| Exp3 | SP CD4 | 28,872 | 6 | 23,967 | 111 |
| | SP CD8 | 6,245 | 2 | 5,362 | 14 |
| | P. CD4 | 112,178 | 99 | 88,190 | 1,760 |
| | p. CD8 | 27,845 | 16 | 23,745 | 249 |

Table S10. Comparison of CDR3 β length between different cell subsets for the top and bottom 1000 frequent sequences.

| CDR3B length | p-value (top1000) | p-value (bottom 1000) |
|----------------|-------------------|-----------------------|
| DP69- vs DP69+ | < 2.2e-16 | 0.874 |
| DP69- vs CD4 | < 2.2e-16 | 0.5433 |
| DP69- vs CD8 | < 2.2e-16 | 0.693 |
| DP69- vs Treg | < 2.2e-16 | 0.8581 |
| DP69+ vs CD4 | < 2.2e-16 | 0.4425 |
| DP69+ vs CD8 | 2.73E-15 | 0.8131 |
| DP69+ vs Treg | 7.43E-05 | 0.9837 |

Table S11. Mean NT/AA ratio for CDR3 β s unique to SP-CD4 and SP-CD8 cells and CDR3 β s shared between these two subsets for the mice in Experiment 2.

| Mean NT/AA | Shared in SP-CD4 and SP-CD8 | Unique to SP-CD4 | Unique to CD8 |
|------------|-----------------------------|------------------|---------------|
| 2autoA | 1.820886 | 1.058942 | 1.022433 |
| 2autoB | 2.086829 | 1.04837 | 1.029694 |
| 2autoC | 2.928571 | 1.04837 | 1.029694 |
| 2alloA | 2.008459 | 1.054977 | 1.041308 |
| 2alloB | 1.81445 | 1.060285 | 1.033697 |
| 2alloC | 2.291667 | 1.003684 | 1.006068 |

Table S12. Statistical significance of odds ratios of cross-reactivity and T1D-reactivity in shared vs unshared sequences for each cell subset in different experiments.

| P-value | Cross-reactivity in shared vs unshared | Sharing in cross-reactives vs allo-non-crossreactives | T1D-reactivity in shared vs unshared |
|-------------|--|---|--------------------------------------|
| Exp1 Treg | 0.000247 | 0.000173 | 0.149 |
| Exp1 SP-CD8 | 2.28E-03 | 0.00441 | 3.68E-09 |
| Exp1 SP-CD4 | 3.35E-21 | 2.86E-18 | 3.44E-10 |
| Exp2 Treg | 0.00379 | 0.00178 | 2.17E-05 |

| | | | |
|-------------|----------|----------|----------|
| Exp2 SP-CD8 | 5.33E-14 | 1.58E-11 | 1.97E-17 |
| Exp2 SP-CD4 | 4.97E-57 | 1.89E-42 | 2.76E-34 |
| Exp3 SP-CD4 | 0.0215 | 0.0144 | 0.00342 |
| Exp3 P. CD4 | 0.000711 | 4.24E-05 | 1.53E-08 |

Table S13. The number of unique CDR α s, CDR3 β s and paired CDR3 α -CDR3 β s, the fraction of cells with a β chain that have at least one paired α chain or two paired α chains and the fraction cells with an α chain that have a paired β chain for each of the five mice in Experiment 2, for which the single cell sequencing of SP-CD4 cells was done.

| Mouse | # unique CDR3α | # unique CDR3β | # unique CDR3α-CDR3β pairs | Fraction of cells with a β chain with at least a paired α chain | Fraction of cells with a β chain that have two paired α chains | Fraction of cells with an α chain that have a paired β chain |
|--------------|---|--|--|---|--|--|
| 2autoB | 6,312 | 7,063 | 6,236 | 0.8729003 | 0.12528 | 0.9608629 |
| 2autoC | 2,494 | 2,672 | 2,484 | 0.9166052 | 0.1136531 | 0.9829838 |
| 2alloA | 5,709 | 6,373 | 5,682 | 0.8792943 | 0.1219437 | 0.9630508 |
| 2alloB | 5,202 | 5,753 | 5,155 | 0.8842196 | 0.1221269 | 0.967167 |
| 2alloC | 4,087 | 4,446 | 4,156 | 0.9104249 | 0.1349102 | 0.9808872 |