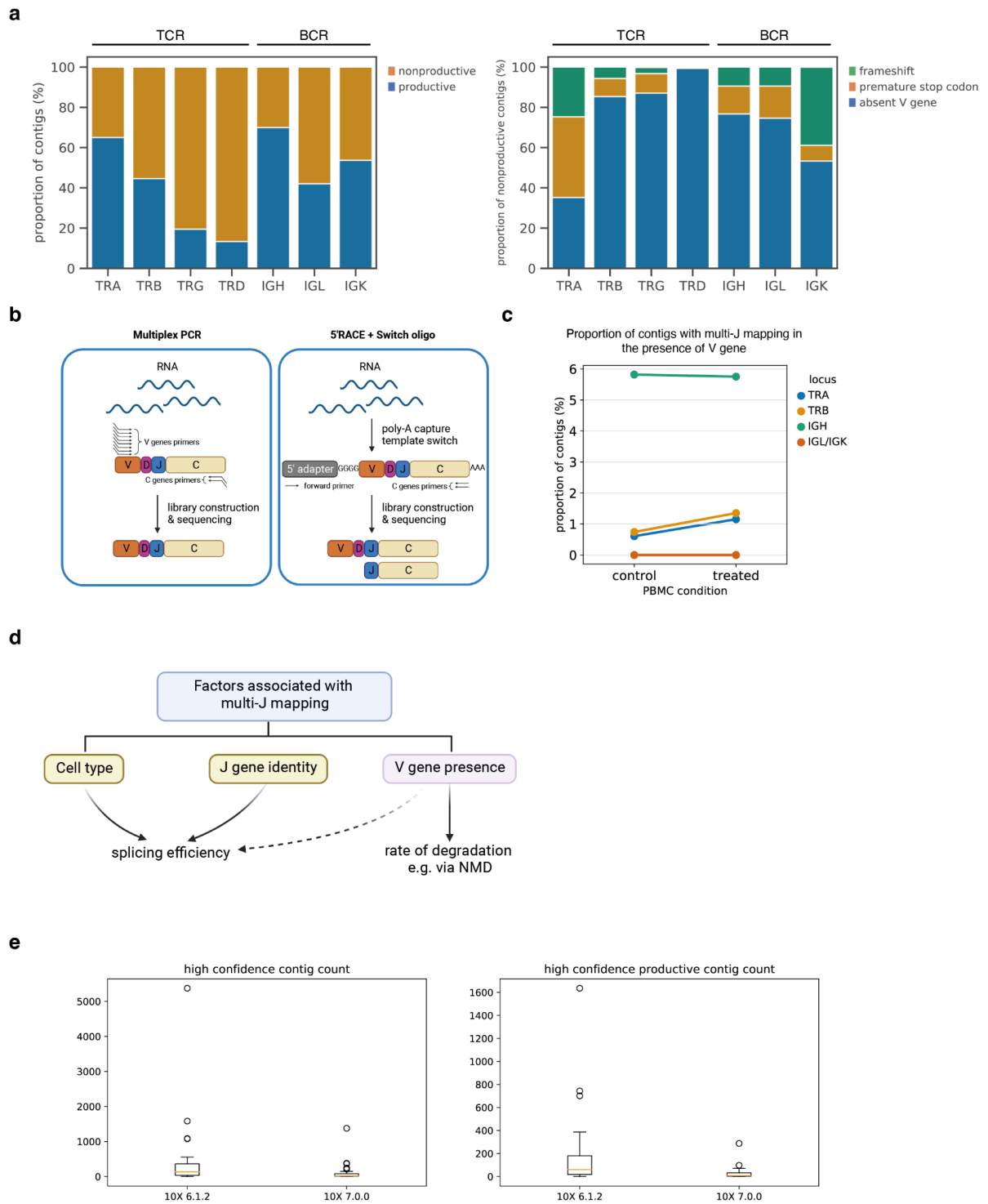


1 Supplementary Figures

	 Immccantation Framework	 enclone	 scRepertoire	VDJView	Immunarch	 scirpy	 dandelion
Programming Language	R & Python	Rust	R	R	R	Python	Python
Version Control	Bitbucket	Github	Github	Bitbucket	Github	Github	Github
AIRR Software Certified	✓	✗	✗	✗	✗	✓	✓
TCR/BCR Centric	BCR	BCR	Both	Both	Both	Both	Both
V(D)J Re-annotation	✓	✗	✗	✗	✗	✗	✓
Clone Definition	✓	✓	✓	✗	✗	✓	✓
BCR Mutation Quantification	✓	✗	✗	✗	✗	✗	Through Immccantation
Diversity Estimation	✓	✗	✓	✗	✓	✓	✓
Visualization	Minimal	✓	✓	✓	✓	✓	✓
Single-cell Integration	Minimal	✓	✓	✓	✗	✓	✓
Phylogenetic Lineage Inference	✓	✓	✗	✗	✗	✗	✗
Trajectory Inference	✗	✗	✗	Through Monocle2 (GEX only)	✗	✗	✓

2
3

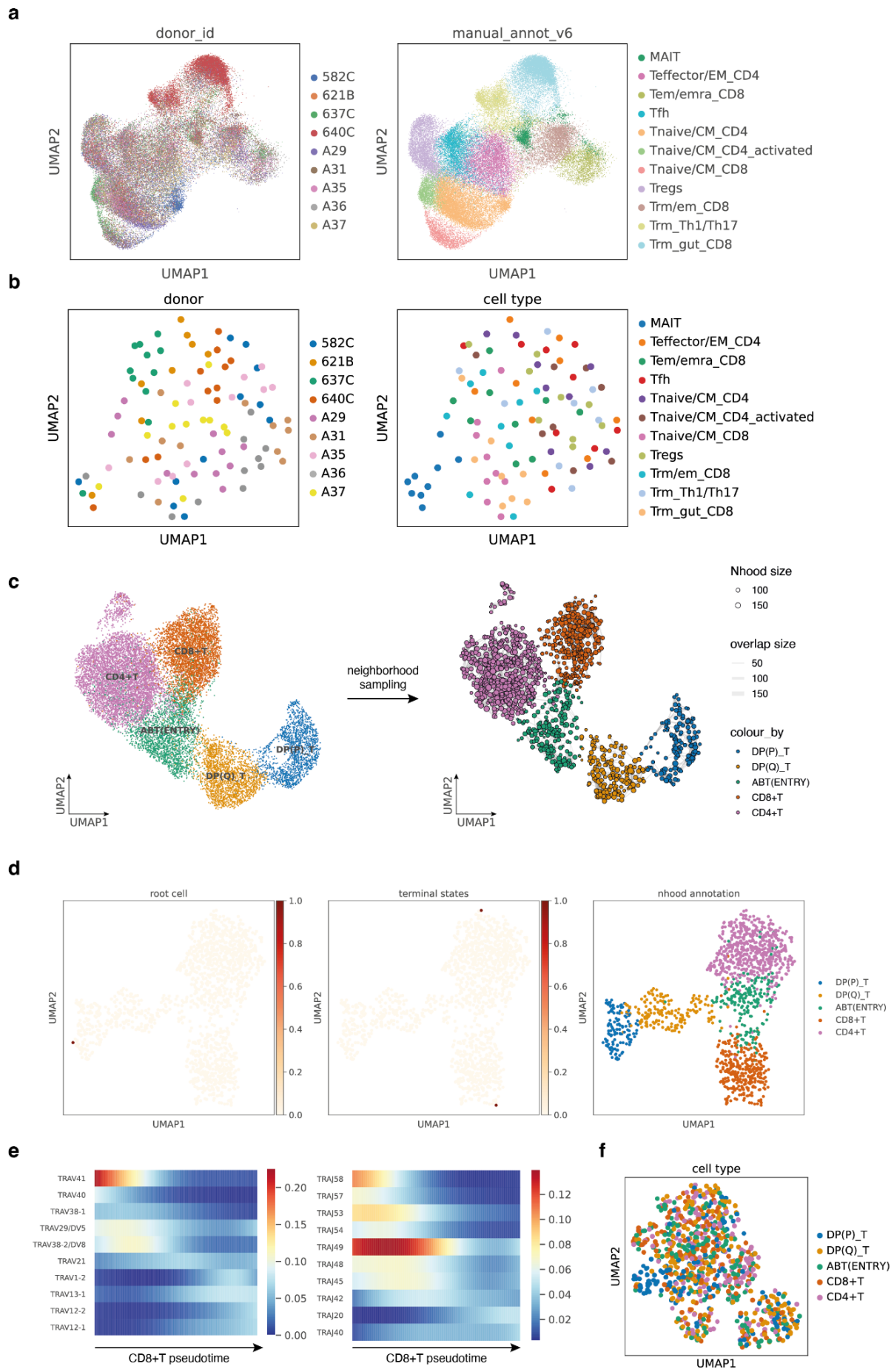
4 **Extended Data Fig. 1 | List of features included in AgR repertoire analysis pipelines.** A
5 table outlining the features of other methods compared to *Dandelion*. As the output from
6 *Dandelion* is compatible with any AIRR-compliant softwares e.g. *Dandelion* output can be
7 passed to *Immccantation* to perform phylogenetic lineage inference.



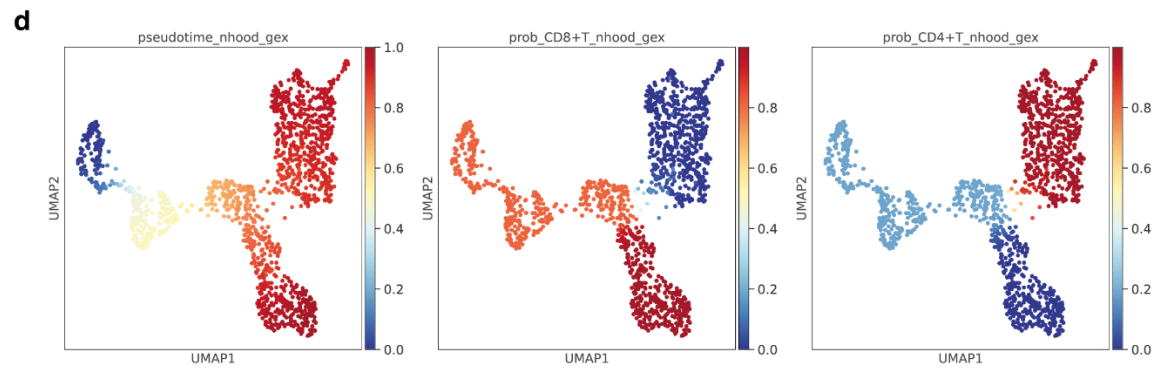
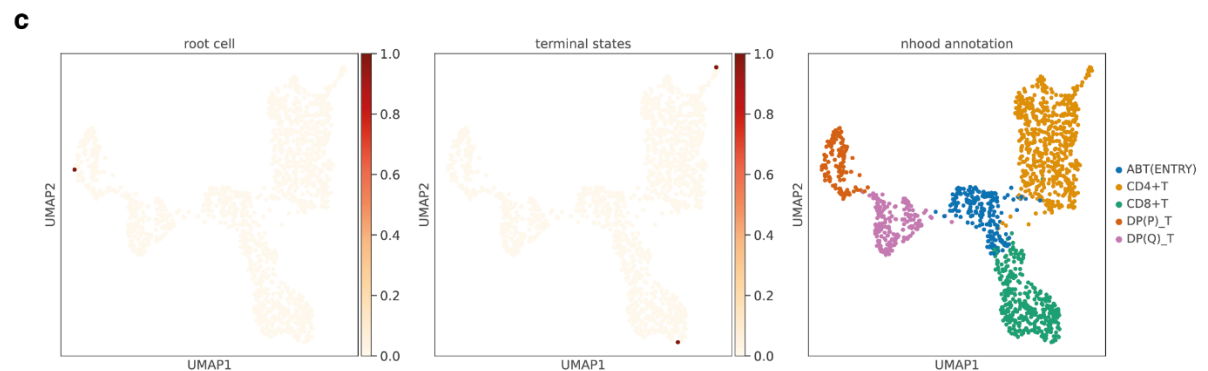
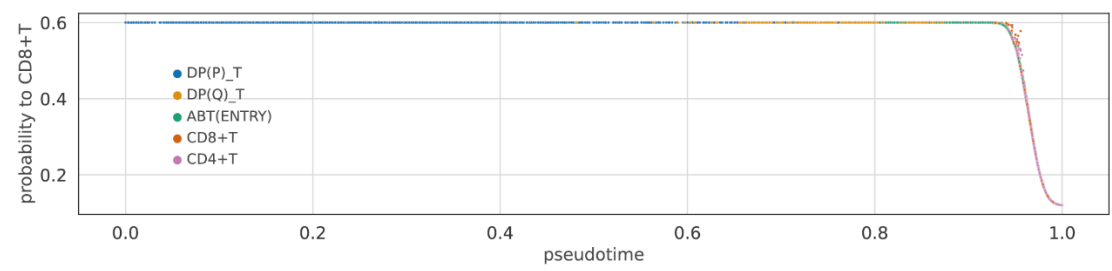
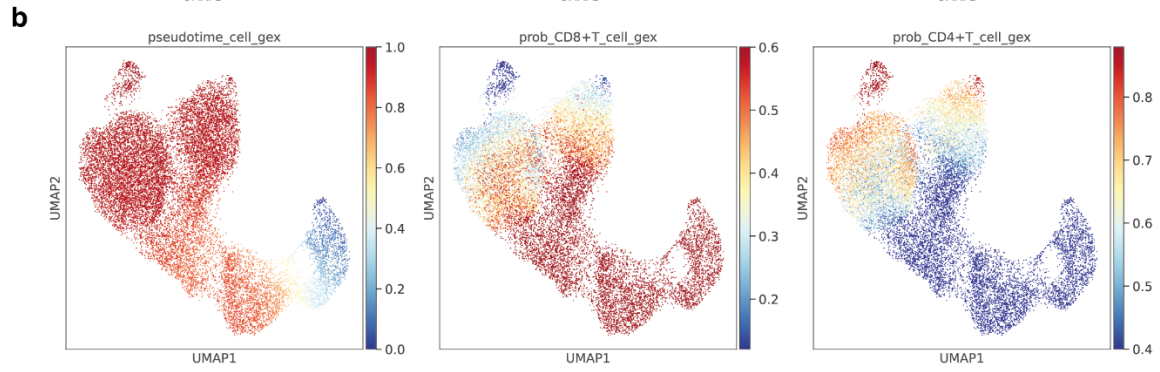
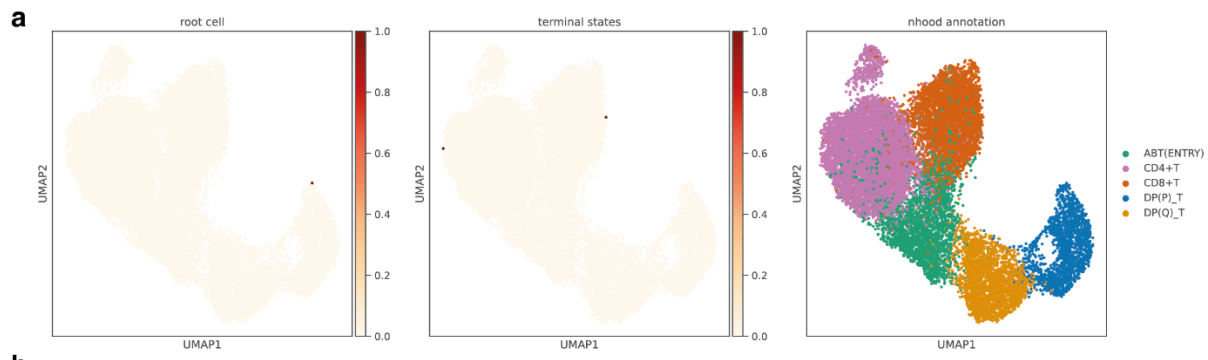
1

2 **Extended Data Fig. 2 | *Dandelion* offers improved contig annotations.** **a**, Left: barplot of
 3 proportion of contigs that are productive or non-productive in each locus. Right: barplot
 4 showing the causes of non-productive contigs in each locus. For both plots, $\text{sc-}\gamma\delta\text{TCR}$, -
 5 $\alpha\beta\text{TCR}$ and -BCR data were taken from Suo et al. 2022³ excluding thymus samples. **b**,
 6 Schematic illustration showing that mRNA without V genes would be captured by 5'RACE +
 7 Switch oligo technique but not by multiplex PCR strategy. **c**, Pointplot of proportion of
 8 contigs with multi-J mapping in the presence of V gene in control and cycloheximide-treated
 9 PBMC samples. Points are colored by locus of TCR/BCR. **d**, Schematic illustration showing

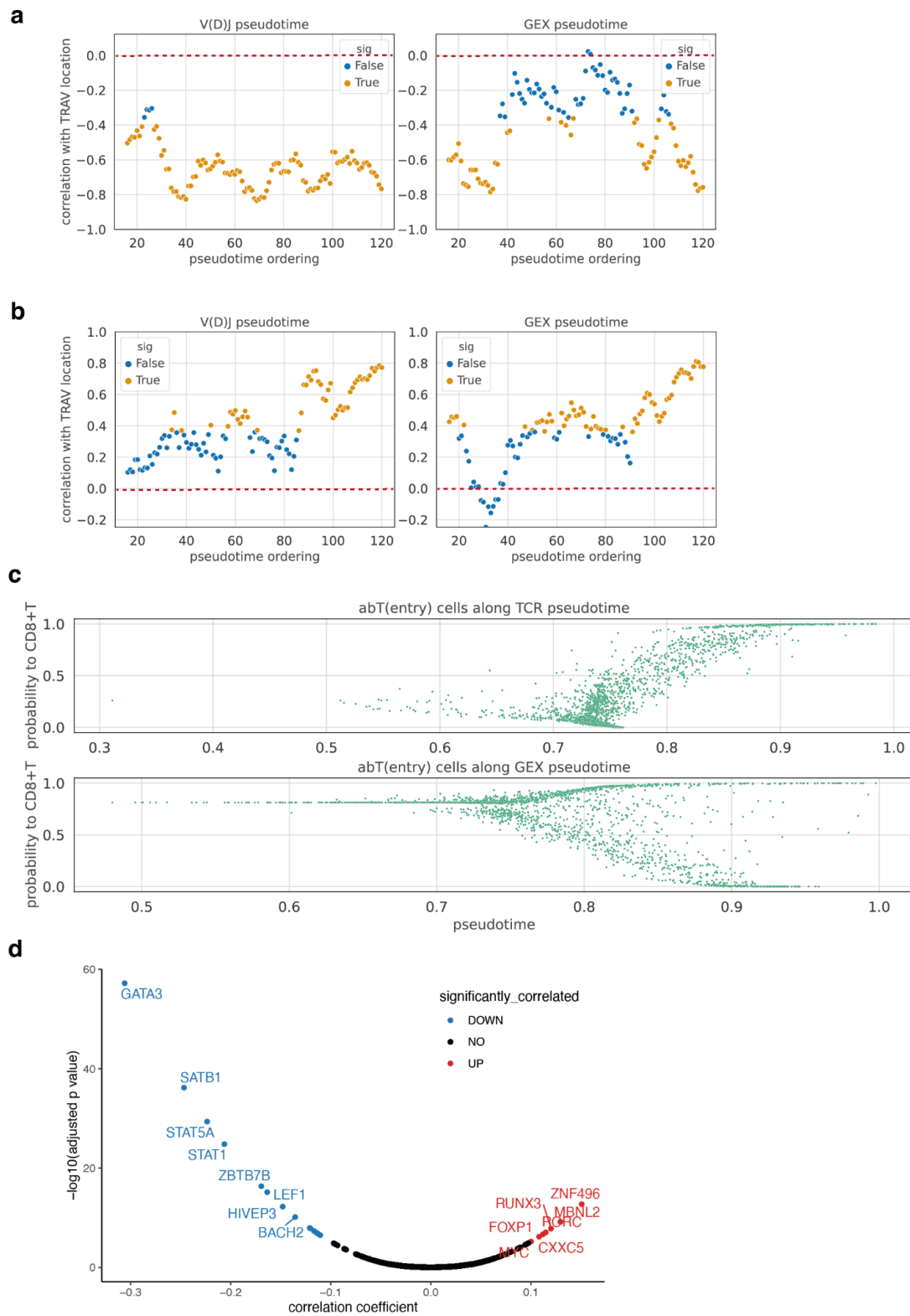
1 the factors associated with multi-J mapping and the proposed mechanisms. **e**, Boxplots of $\gamma\delta$ TCR
2 contig counts annotated by 10X *cellranger vdj* v6.1.2 versus v7.0.0 using data from
3 Suo et al. 2022³. Left: all high confidence contigs (P -value 5.43e-6, r 0.91 in the Wilcoxon
4 signed-rank test). Right: high confidence productive contigs (P -value 1.69e-6, r 0.96 in the
5 Wilcoxon signed-rank test).



1 **Extended Data Fig. 3 | V(D)J feature space. a**, Gene expression UMAP of all T cells from
2 Conde et al. 2022⁵, colored by donor ID (left) or high-level cell type annotations (right). Each
3 point represents a cell. **b**, UMAP of the pseudo-bulk V(D)J feature space of the same cells as
4 in **a**, colored by donor ID (left) or high-level cell type annotations (right). Each point
5 represents a cell pseudo-bulk. **c**, Left: UMAP of DP to mature T cells with paired productive
6 $\alpha\beta$ TCR in data from Suo et al. 2022³. Each point represents a cell, colored by cell types.
7 Right: cell neighborhood graph on the same UMAP embedding. Each point represents a cell
8 neighborhood, colored by cell types. The point size represents neighborhood size, with
9 connecting edges representing overlapping cell numbers between any two neighborhoods.
10 Only edges with more than 30 overlapping cells are shown. The layout of nodes is
11 determined by the position of the neighborhood index cell in the UMAP on the left. **d**, The
12 root cell and terminal states selected for pseudotime inference in **Fig. 3c**. **e**, Gene expression
13 trends over CD8+T pseudotime imputed with *palantir*²³. Only the top 10 most frequently
14 used TRAV or TRAJ genes are shown. **f**, UMAP representation of tcrdist-derived PCA
15 coordinates of VDJ data computed by CoNGA²⁴, with the same dataset as used in **c**, colored
16 by cell types.

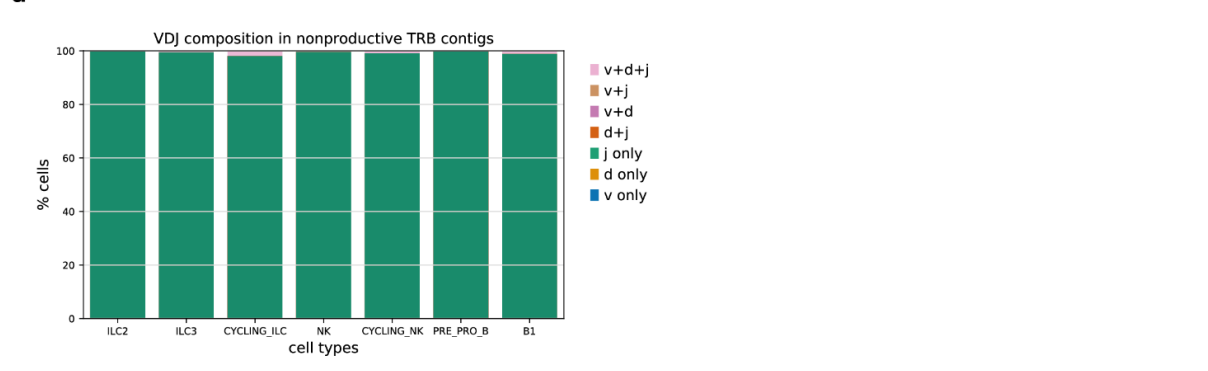
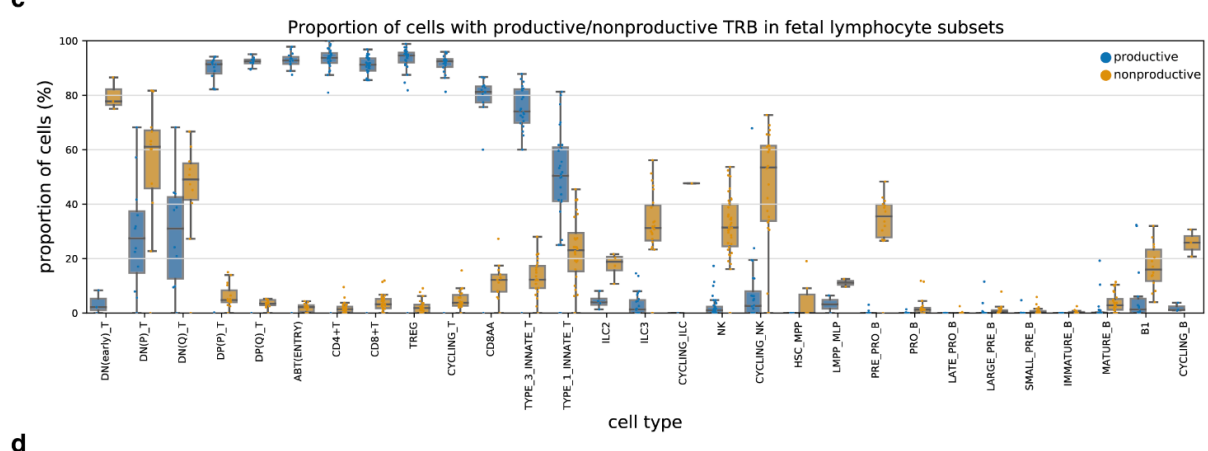
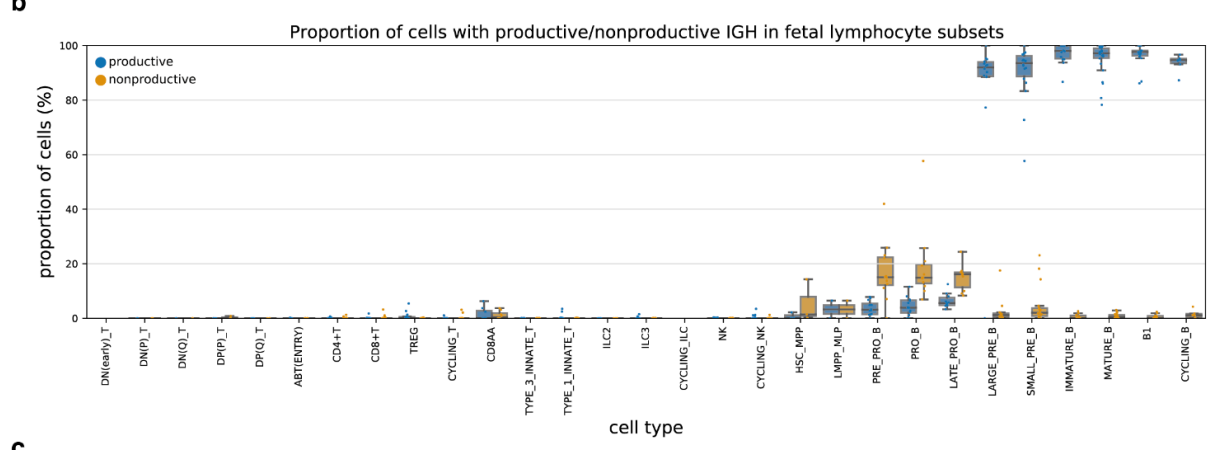
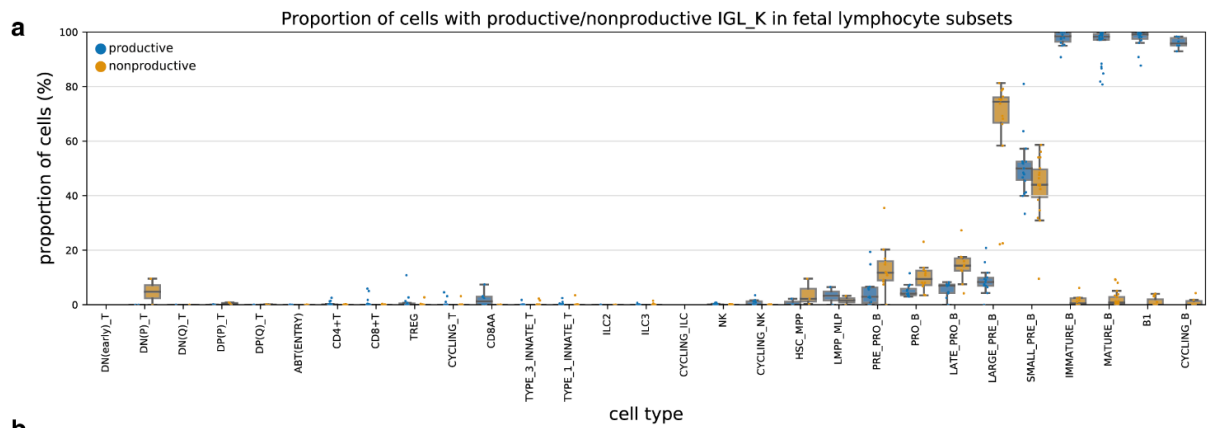


1 **Extended Data Fig. 4 | T cell development pseudotime inference comparison. a**, DP to
2 mature T cells with paired productive $\alpha\beta$ TCR in data from Suo et al. 2022³, on the same
3 UMAP embedding as in **Fig. 4a** and **Supplementary Fig. 3c**. The first two panels show the
4 root cell and terminal states selected for pseudotime inferred directly from single-cell gene
5 expression. The last panel shows the cell types. **b**, Top: pseudotime and branch probabilities
6 inferred directly from single-cell gene expression on the same UMAP embedding as in **a**.
7 Bottom: scatterplot of branch probability to CD8+T against pseudotime. Each point
8 represents a cell. **c**, UMAP of neighborhood GEX space, with the same neighborhoods as
9 sampled in **Supplementary Fig. 3c** and UMAP embedding computed on gene expression
10 pseudo-bulked by neighborhoods. Each point represents a cell neighborhood. The first two
11 panels show the root cell and terminal states selected for pseudotime inferred from
12 neighborhood GEX space. The last panel shows the cell types. **d**, Inferred pseudotime, and
13 branch probabilities to CD8+T and to CD4+T respectively overlaid onto the same UMAP
14 embedding in **c**.

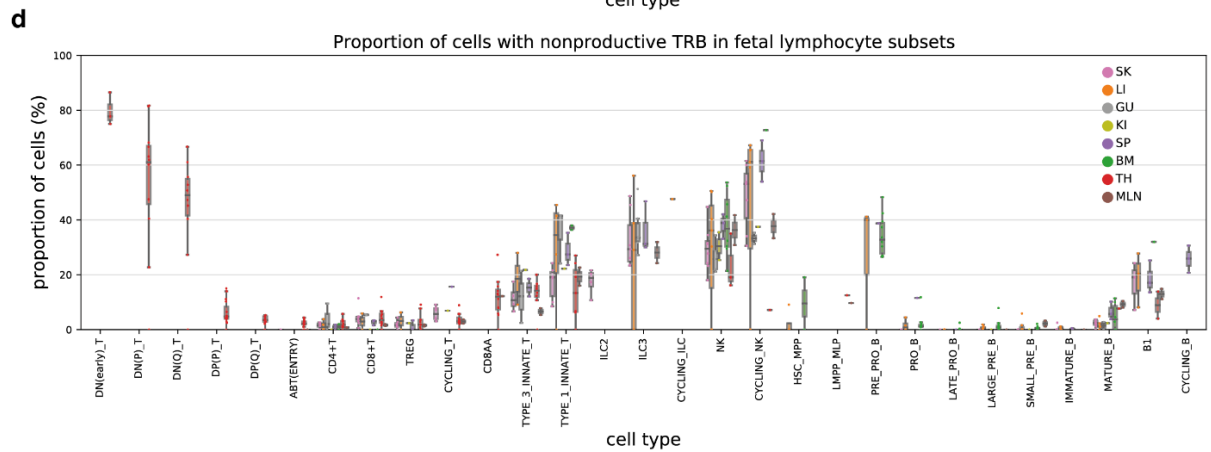
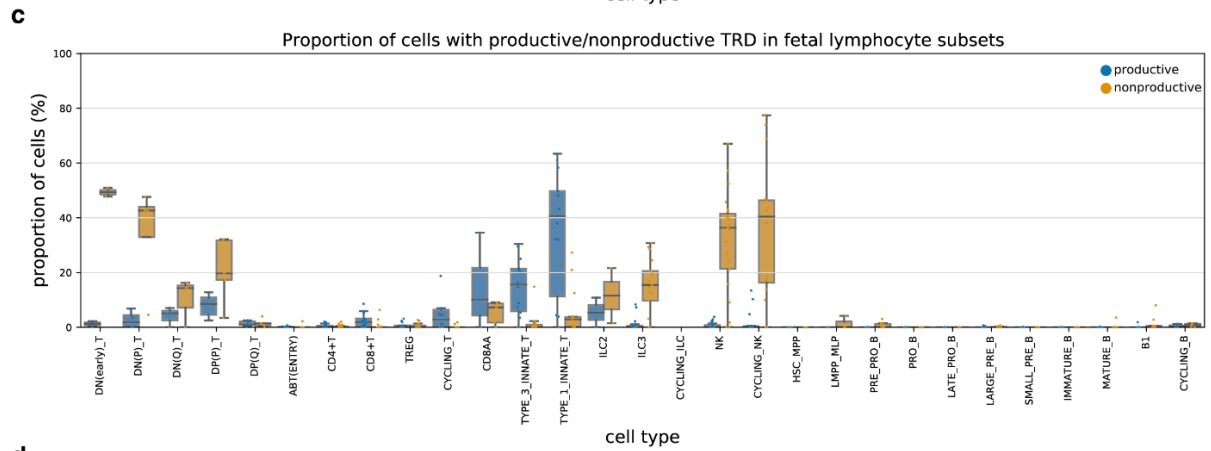
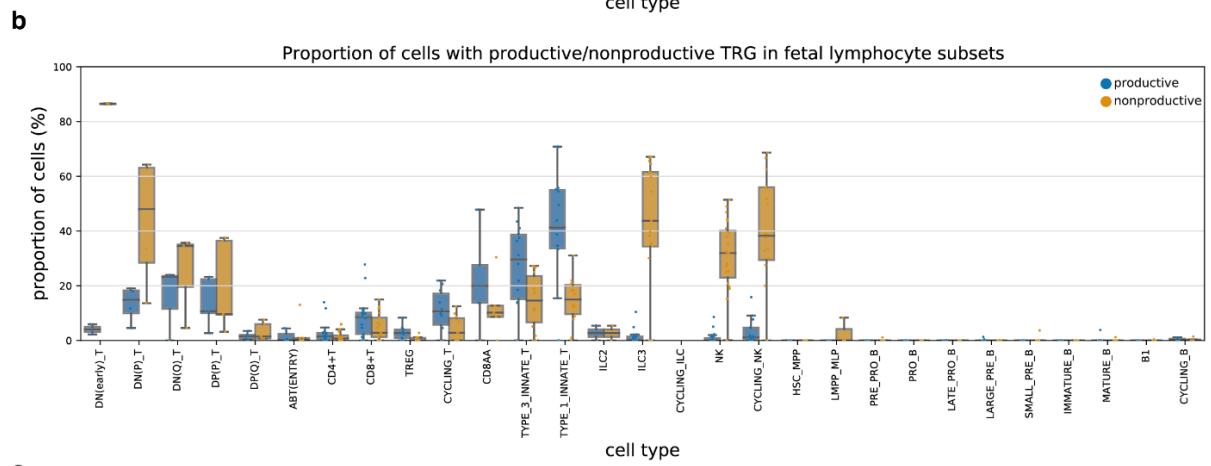
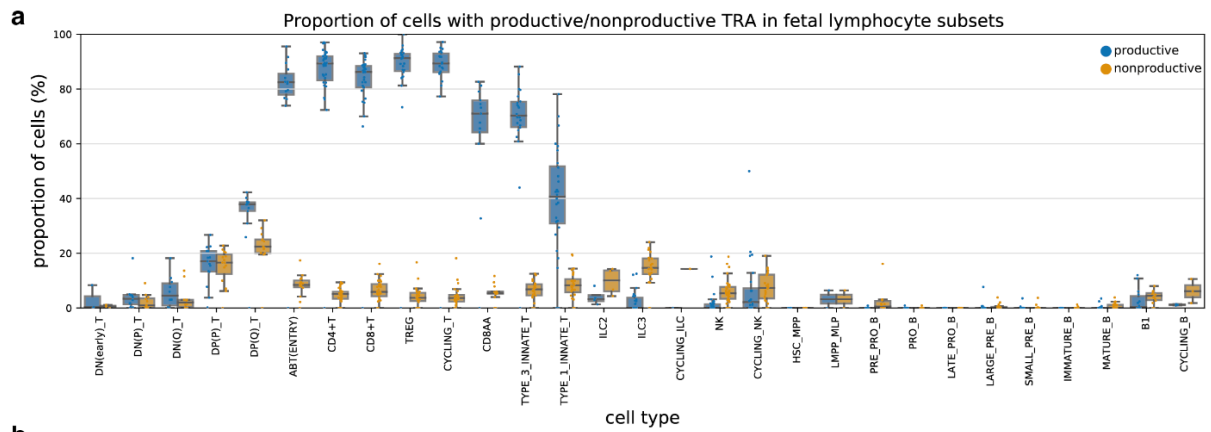


1
2 **Extended Data Fig. 5 | Comparing pseudotime inferred from neighborhood V(D)J space**
3 **or GEX space. a**, Pearson's correlation coefficients of pseudotime order and average relative
4 TRAV location over sliding windows of 30 adjacent neighborhoods on the pseudotime order
5 (left: pseudotime inferred from neighborhood V(D)J space; right: pseudotime inferred from
6 neighborhood GEX space). Y-axis is the correlation coefficient and the x-axis is the median
7 pseudotime order of the 30 adjacent neighborhoods. The color of the points represents

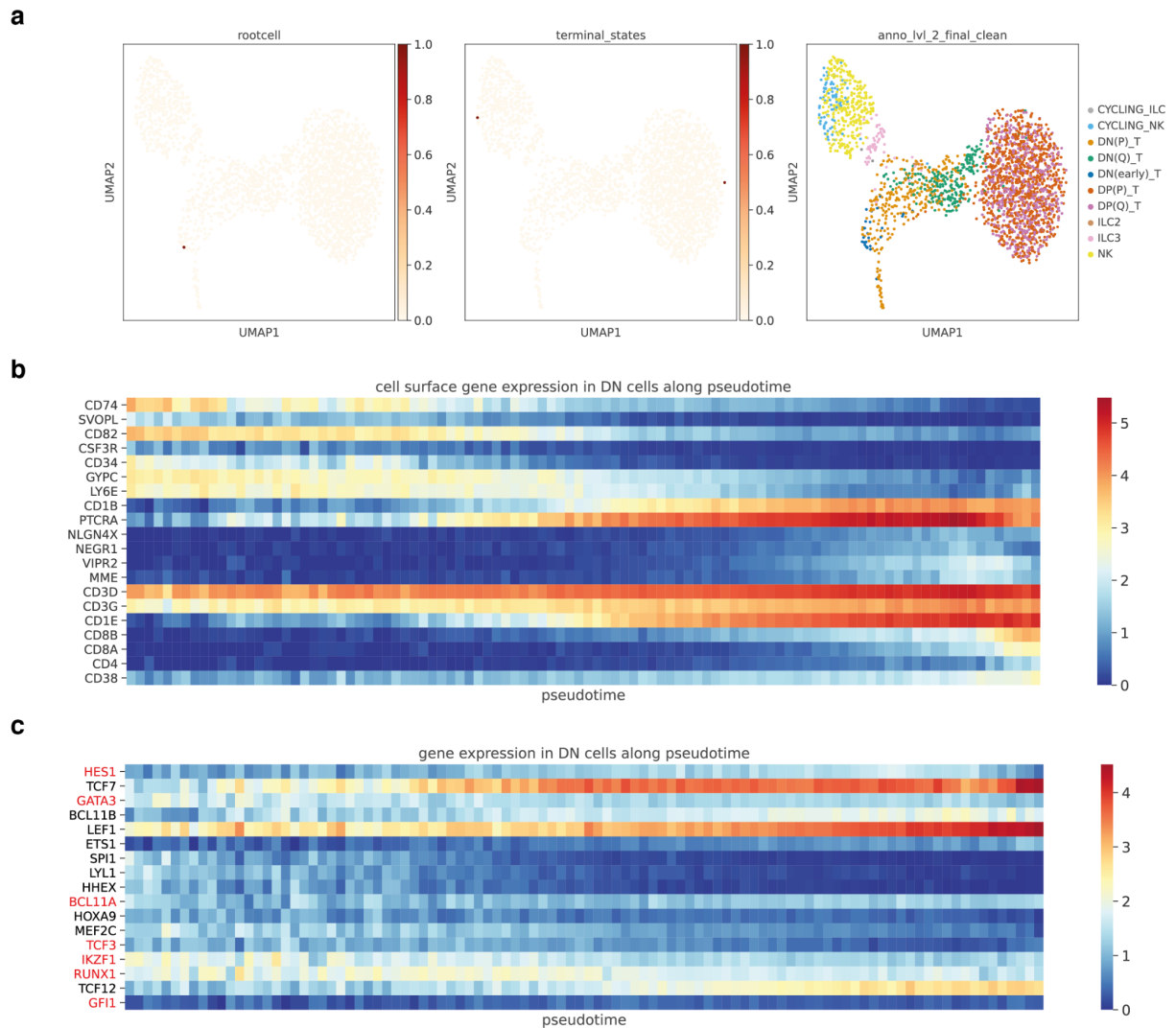
1 statistical significance (orange: P -value from the Pearson's correlation < 0.05 ; blue: P -value
2 ≥ 0.05). The red dashed lines mark the correlation coefficient of 0. **b**, The same plot as in **a**
3 but for TRAJ. **c**, Scatterplots of branch probability to CD8+T against pseudotime in
4 abT(entry) cells. Each point represents a cell. Top panel: pseudotime inferred from
5 neighborhood V(D)J space as in **Fig. 4a** top panel. Bottom panel: pseudotime inferred from
6 neighborhood GEX space as in **Fig. 4a** bottom right panel. **d**, Volcano plot summarizing
7 results of TFs that are correlated with branch probabilities to CD8+T lineage in V(D)J
8 pseudotime within abT(entry) cells. The y -axis is the $-\log_{10}$ (BH adjusted P -value) and the x -
9 axis is the correlation coefficient. Labeled TFs that had significant (BH adjusted P -value $<$
10 0.05) positive correlations (correlation coefficient > 0.1) were colored in red, the ones with
11 significant negative correlations (correlation coefficient < -0.1) were colored in blue, and the
12 rest were colored in black.



1 **Extended Data Fig. 6 | Non-productive BCR and TCR. a,b,c**, Boxplot of the proportion of
2 cells with productive (blue) or non-productive (orange) BCR light chain (**a**) and heavy chain
3 (**b**), and TRB (**c**) in different fetal lymphocyte subsets. Each point represents a sample and
4 data were taken from Suo et al. 2022³. Only samples with at least 20 cells are shown. Boxes
5 capture the first to third quartiles and whiskers span a further 1.5X interquartile range on each
6 side of the box. **d**, Barplot showing the VDJ composition of non-productive TRB contigs in
7 selected lymphocyte subsets from **Fig. 5a**.



1 **Extended Data Fig. 7 | Non-productive TCR. a,b,c**, Boxplot of the proportion of cells with
2 productive (blue) or non-productive (orange) TRA (**a**), TRG (**b**) and TRD (**c**) in different
3 fetal lymphocyte subsets. Each point represents a sample and data were taken from Suo et al.
4 2022³. Only samples with at least 20 cells are shown. Boxes capture the first to third quartiles
5 and whisks span a further 1.5X interquartile range on each side of the box. **d**, Boxplot of the
6 proportion of cells with non-productive TRB in different fetal lymphocyte subsets, colored by
7 organs. Each point represents a sample. Only samples with at least 20 cells are shown. Boxes
8 capture the first to third quartiles and whisks span a further 1.5X interquartile range on each
9 side of the box.



1

2 **Extended Data Fig. 8 | TRBJ-based trajectory for ILC/NK/T cell lineage. a,**
 3 Neighborhood V(D)J feature space covering ILC, NK and developing T cells with TRBJ on
 4 the same UMAP embedding as in **Fig. 5b**. The first two panels show the root cell and
 5 terminal states selected for pseudotime inference. The last panel shows the cell types. **b,**
 6 Heatmap of gene expression for genes encoding cell surface proteins across pseudotime in
 7 DN T cells. Pseudotime is equally divided into 100 bins, and the average gene expression is
 8 calculated for DN T cells with pseudotime that falls within each bin. Genes selected here had
 9 significantly high Chatterjee's correlation with pseudotime (BH adjusted P -value < 0.05 , and
 10 correlation coefficient > 0.1). **c,** Heatmap of gene expression for TFs known to be important
 11 in mouse DN T cell development⁴⁰, across pseudotime in human fetal DN T cells. TFs that
 12 showed discordant expression patterns between mouse and human are highlighted in red.

1 **Supplementary Tables**

2 **Supplementary Table 1: top_10_j_multimappers.csv (separate file)**

3 Top 10 J gene combinations with multi-J mapping for each locus in data from Suo et al.
4 2022³, with the number of contigs containing each combination shown next to it.

6 **Supplementary Table 2: LR_results.csv (separate file)**

7 Logistic regression results exploring factors associated with multi-J mapping presence in data
8 from Suo et al. 2022³.

10 **Supplementary Table 3: LR_results_combined.csv (separate file)**

11 Logistic regression results exploring factors associated with multi-J mapping presence in
12 control and cycloheximide-treated PBMC data.

14 **Supplementary Table 4: j_sequence_affect_j_multimapper.csv (separate file)**

15 List of leftmost (5' end) J genes that had significant association with increased or decreased
16 multi-J mapping, together with the sequences of their last 10 nucleotides at 3' ends and the
17 first 11 nucleotides of its 3' end intron.

19 **Supplementary Table 5: panimmune_differential_VDJ.csv (separate file)**

20 Differential V(D)J usage across CD4+T, CD8+T, and MAIT cells in data from Conde et al.
21 2022⁵.

23 **Supplementary Table 6: abtentry_cor_result.csv (separate file)**

24 Pearson's correlation coefficients and BH adjusted *P*-values of all genes with branch
25 probabilities to CD8+T lineage within abT(entry) cells.

26 [cor_tcr] Pearson's correlation coefficients for pseudotime inferred from neighborhood V(D)J
27 space

28 [pval_tcr] Pearson's correlation *P*-values for pseudotime inferred from neighborhood V(D)J
29 space

30 [adjp_tcr] *P*-values from pval_tcr adjusted by BH procedure

31 [cor_gex] Pearson's correlation coefficients for pseudotime inferred from neighborhood GEX
32 space

33 [pval_gex] Pearson's correlation *P*-values for pseudotime inferred from neighborhood GEX
34 space

35 [adjp_gex] *P*-values from pval_gex adjusted by BH procedure