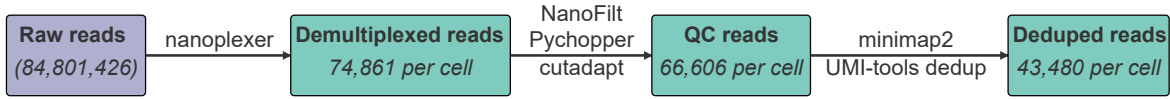
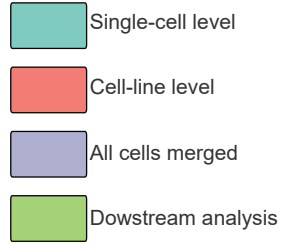


# Supplementary Fig. S1

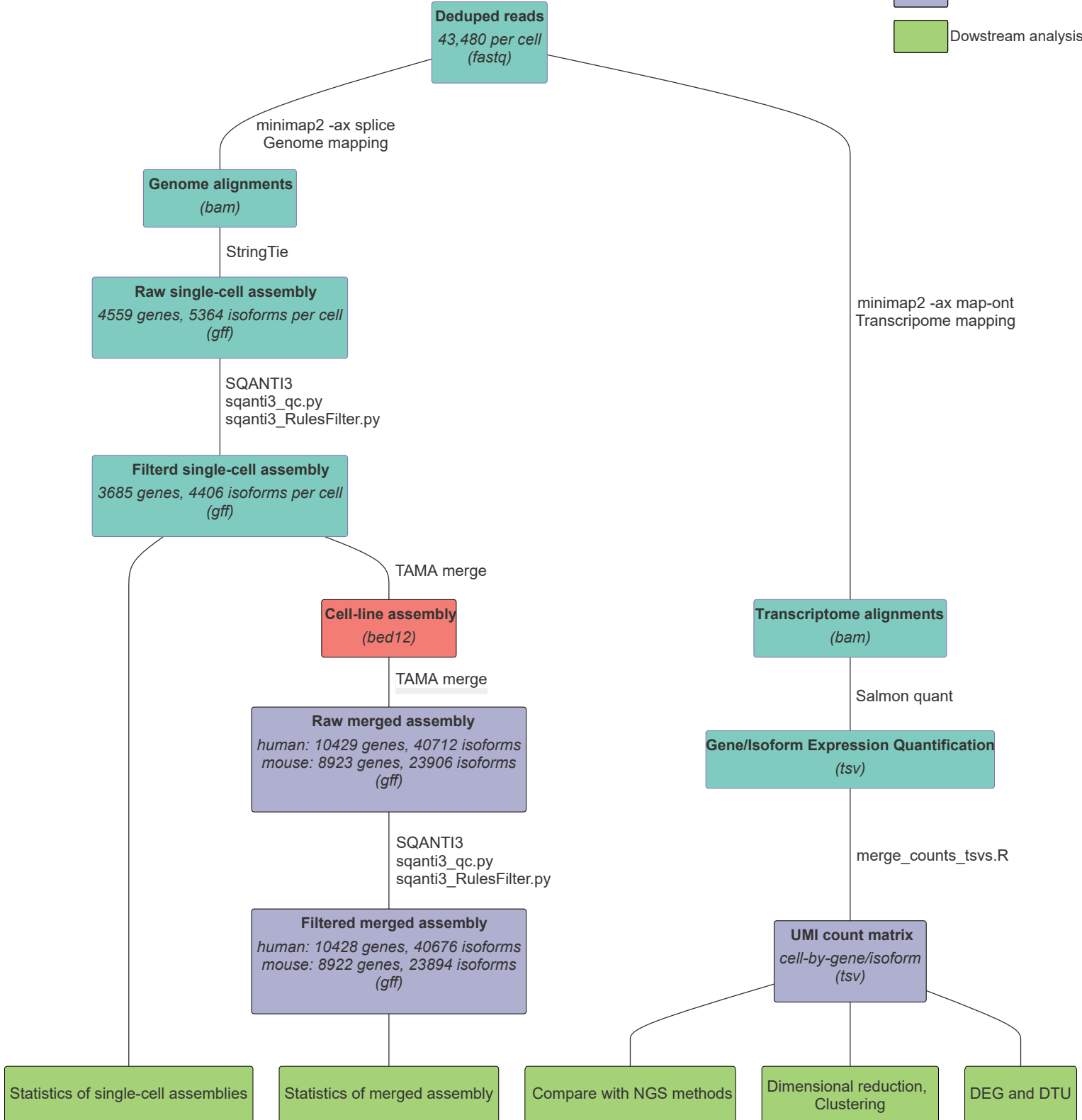
## a Reads QC and deduplication



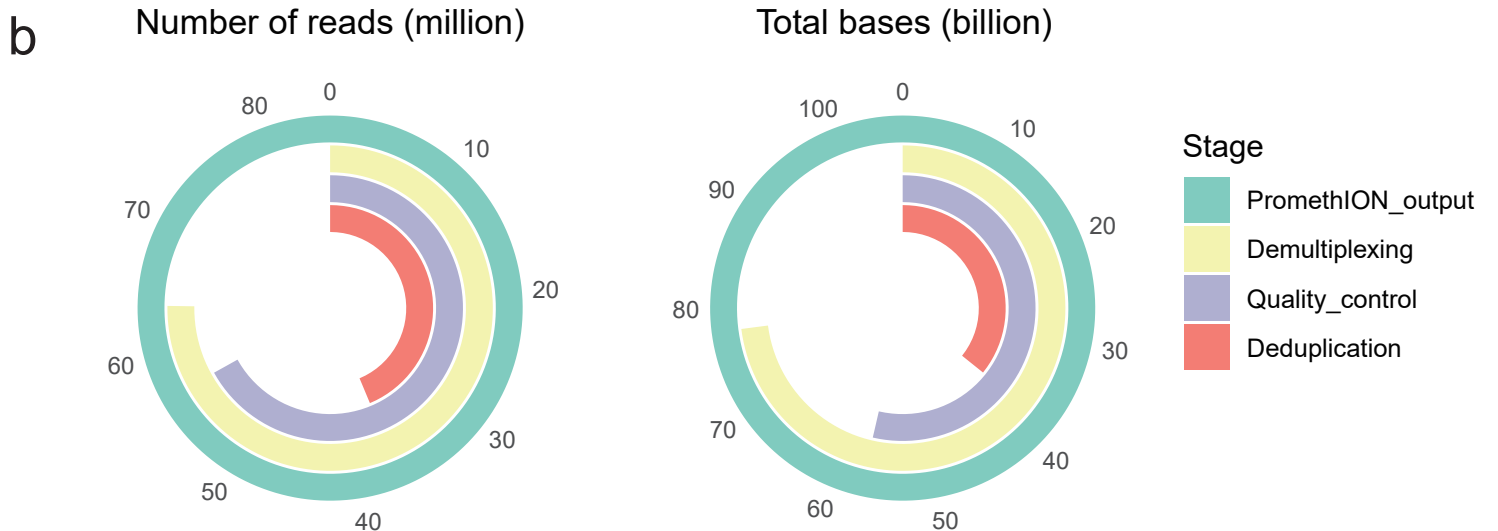
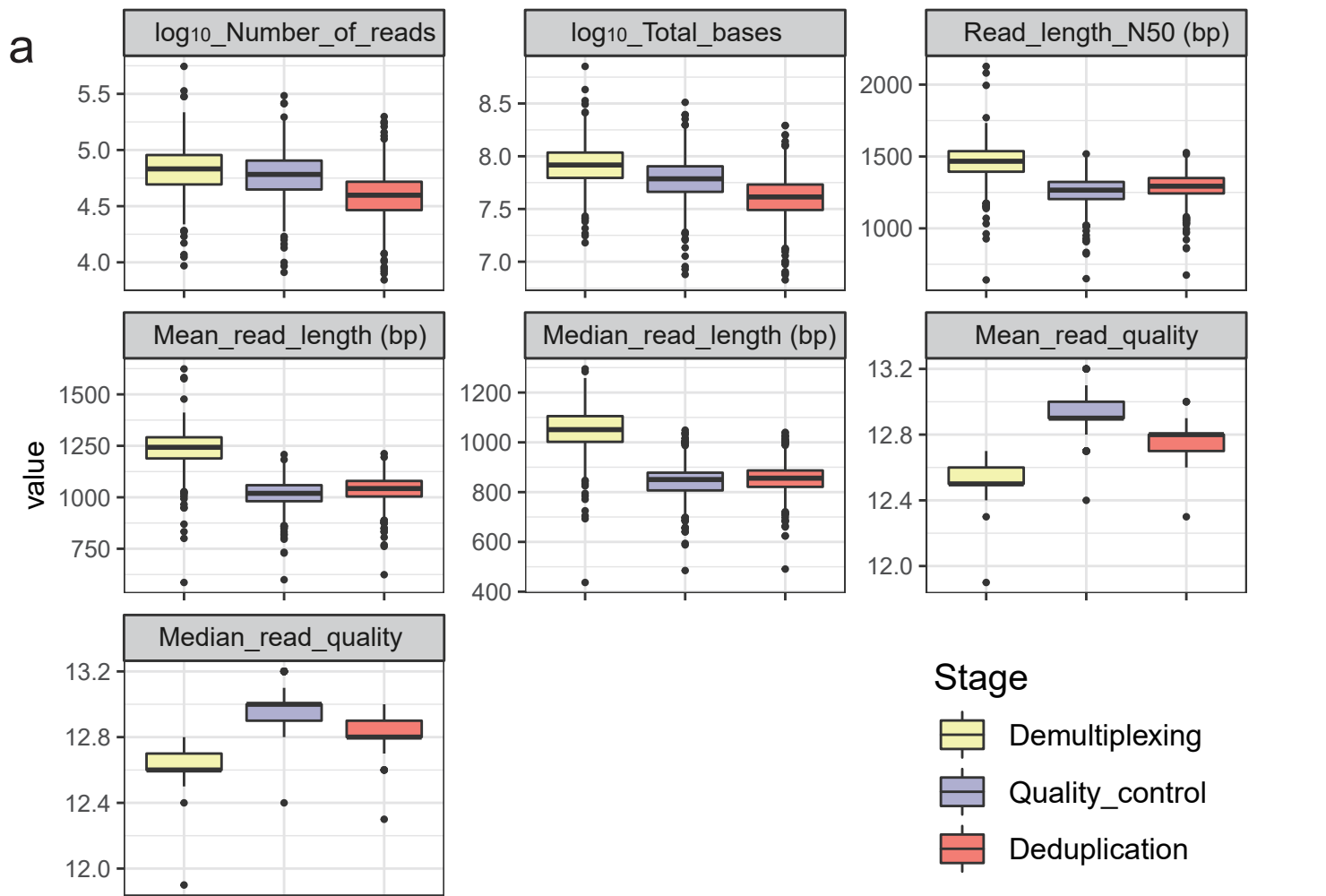
### Legend



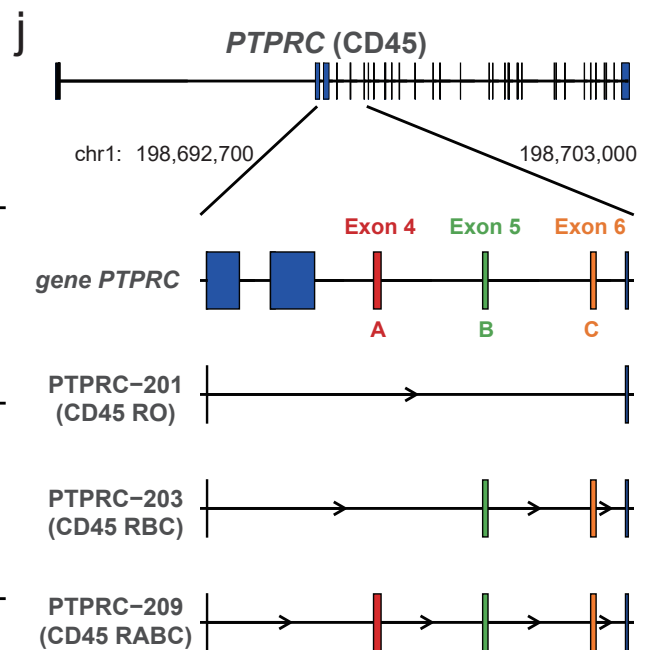
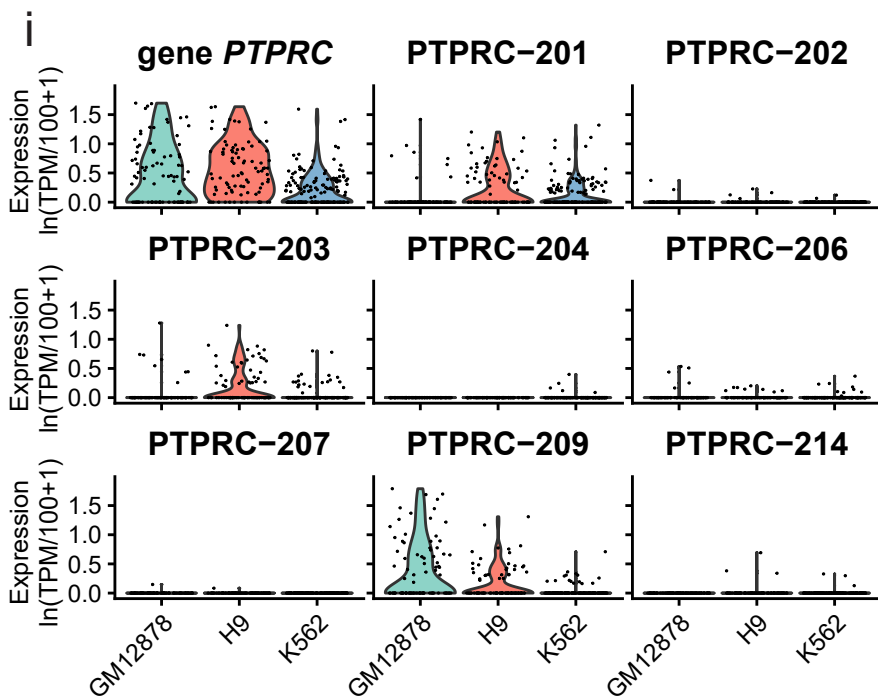
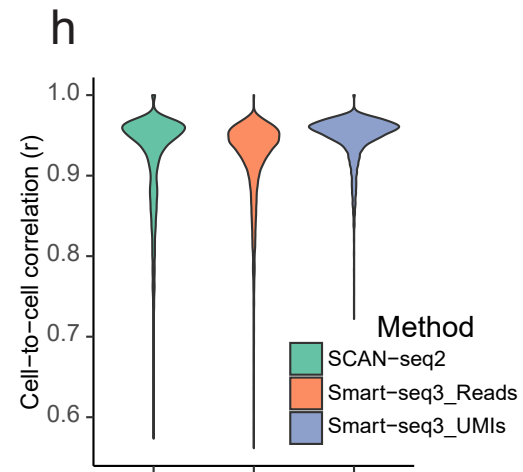
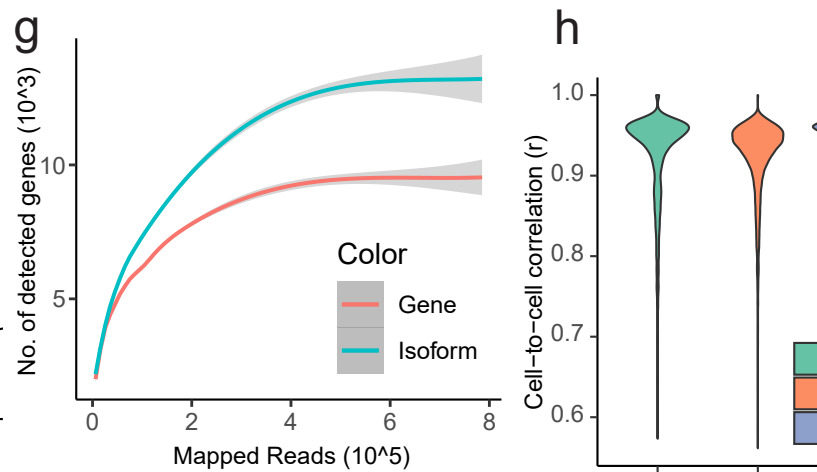
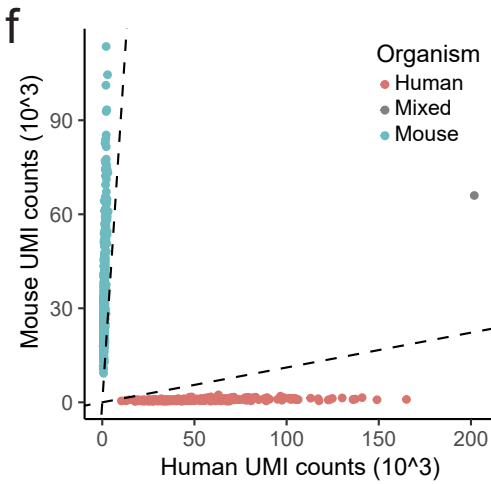
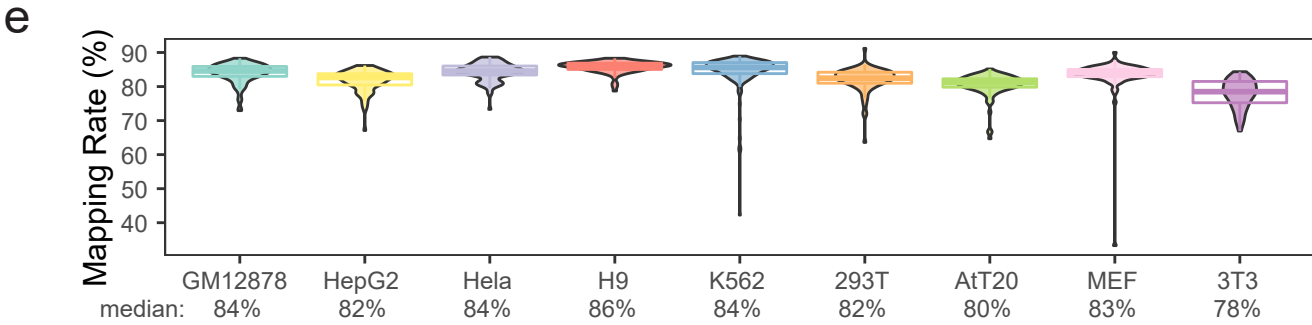
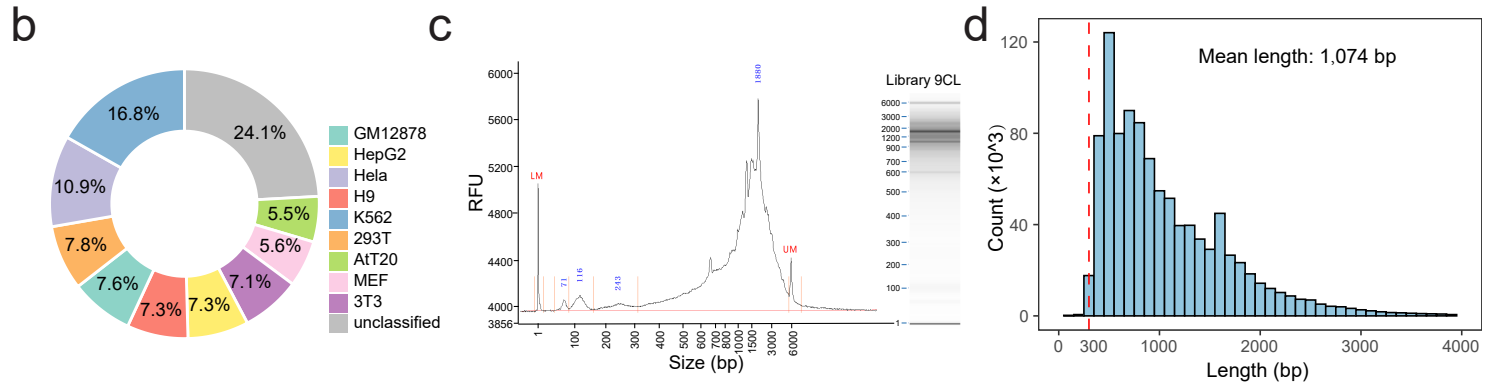
## b Transcriptome quantification and assembly



# Supplementary Fig. S2

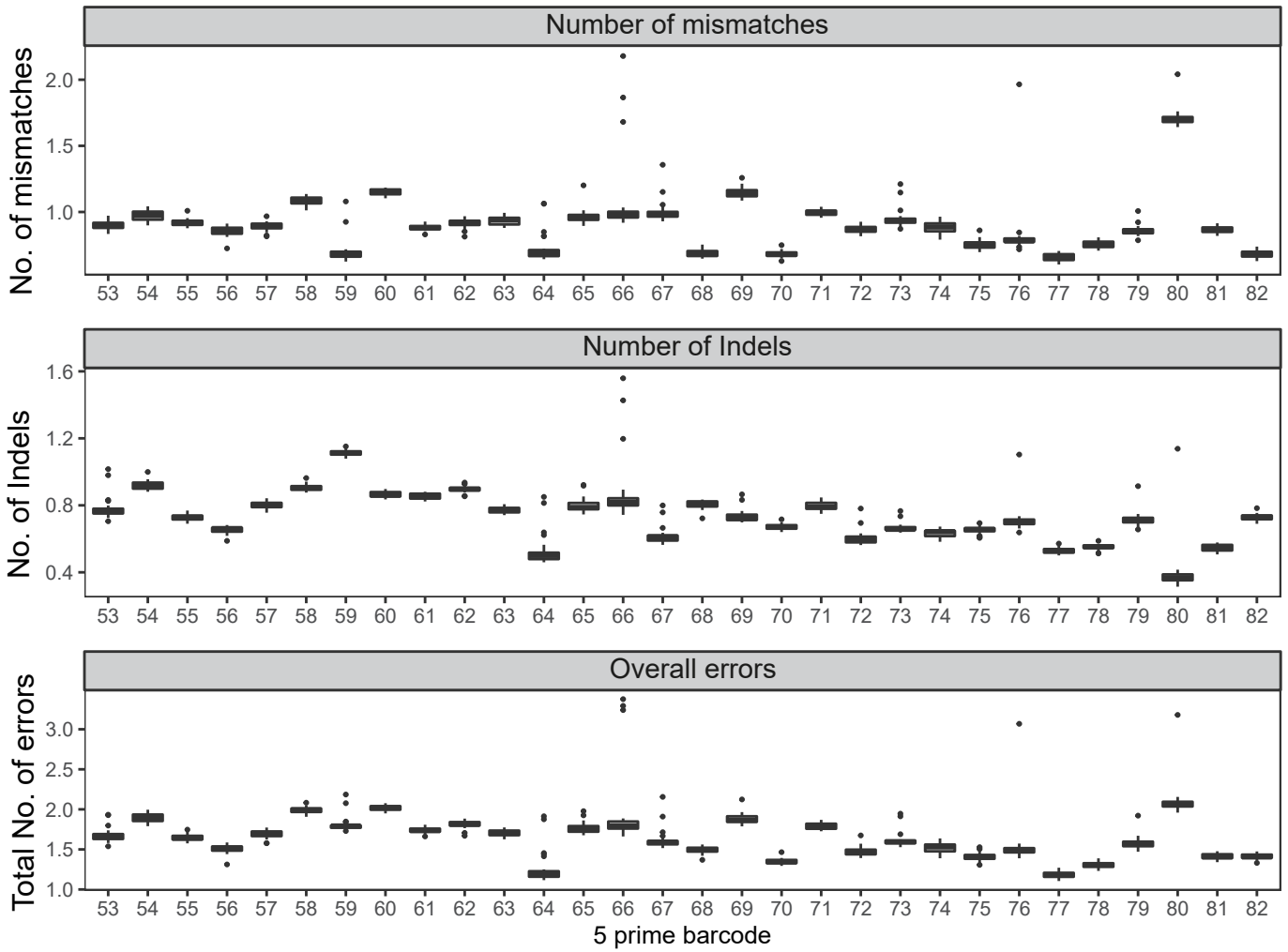


# Supplementary Fig. S3

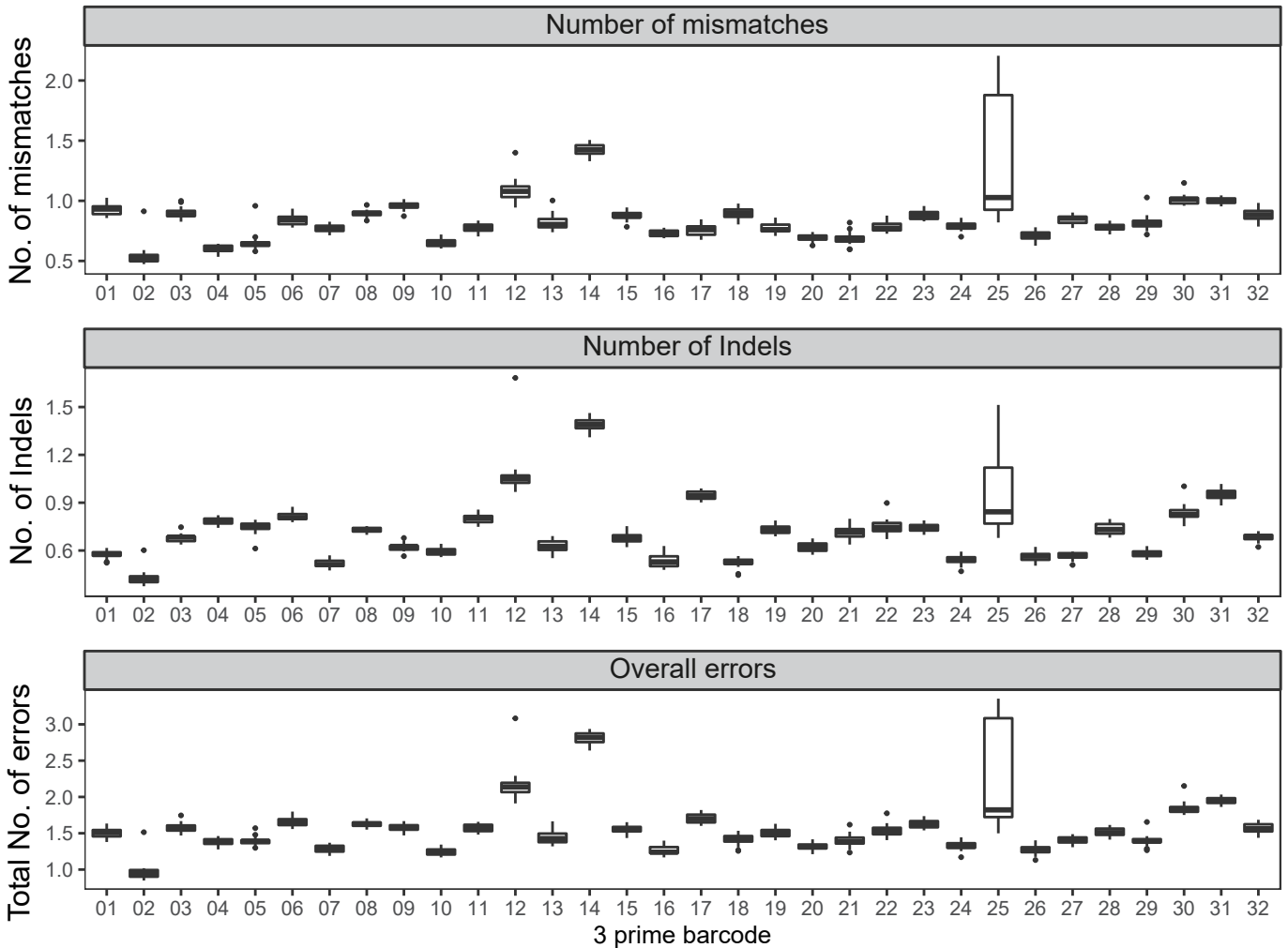


# Supplementary Fig. S4

**a**



**b**

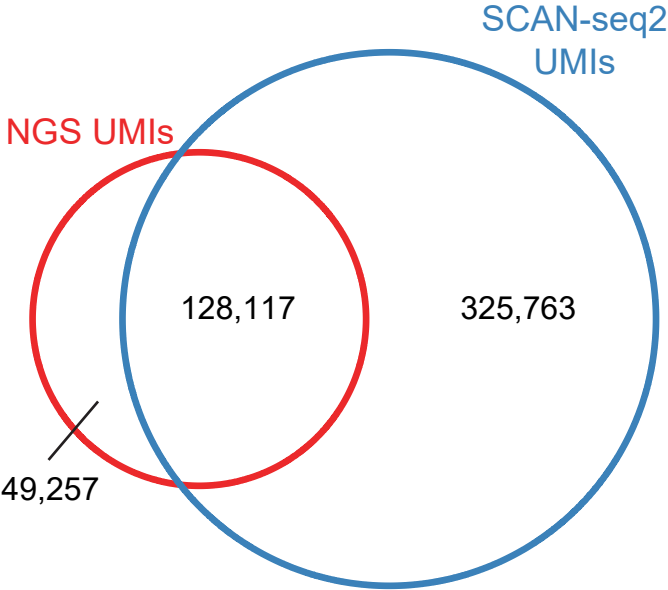




Supplementary Fig. S5

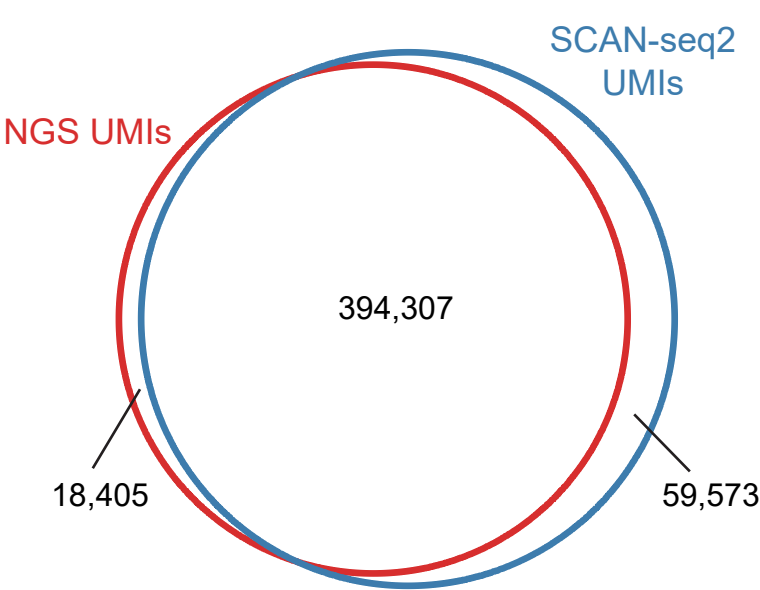
a

No error tolerated

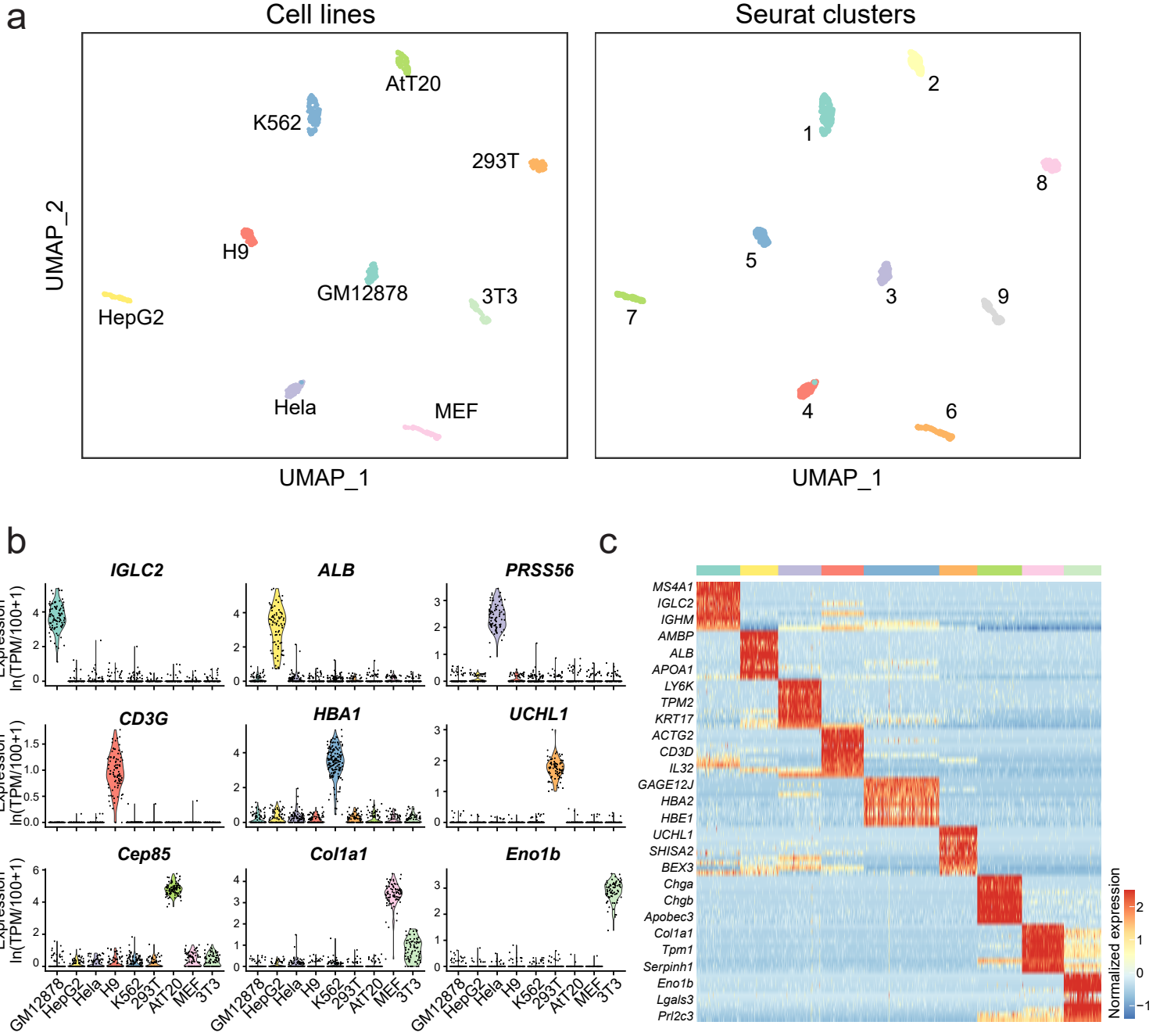


b

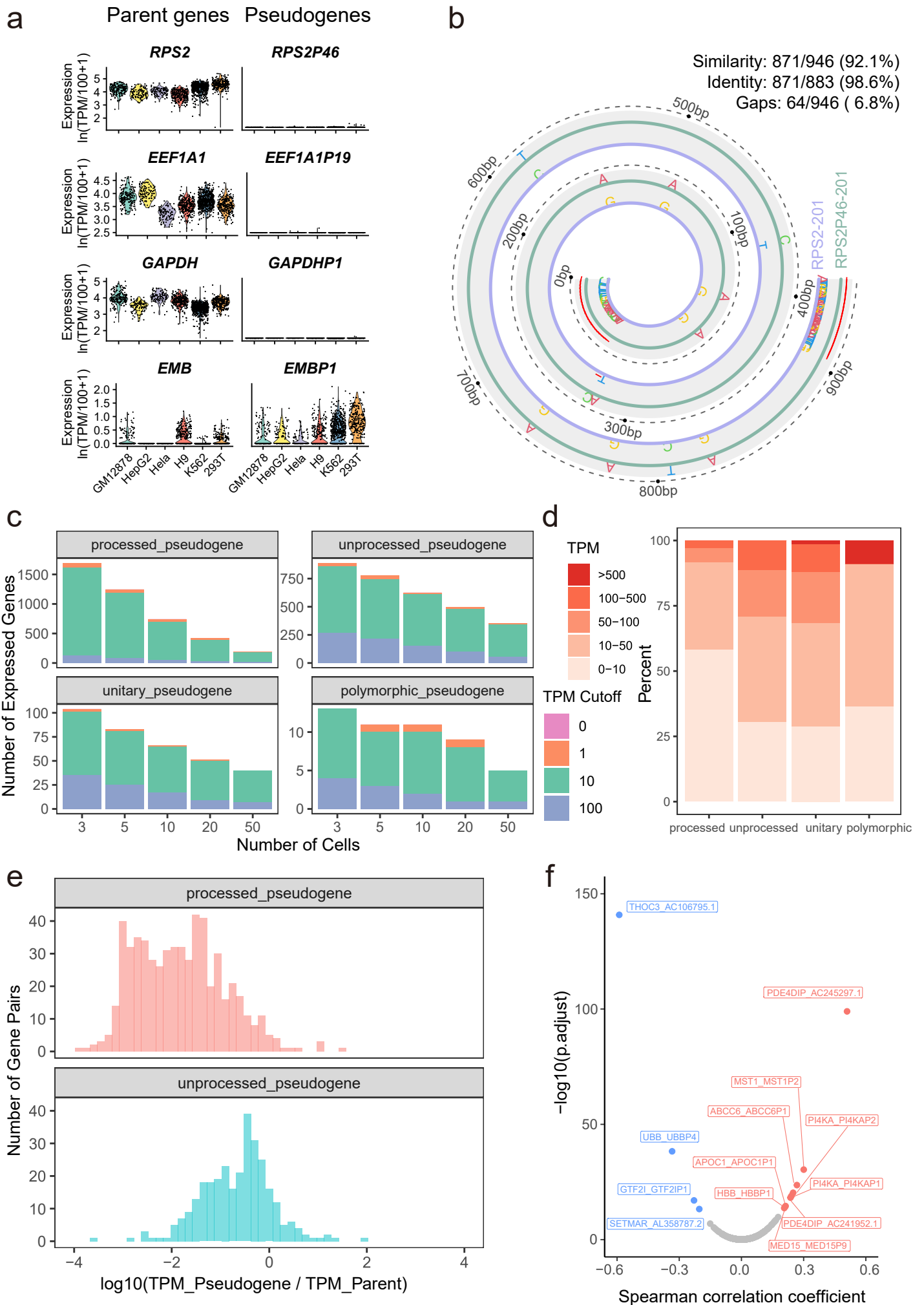
Tolerate edit distance of 1



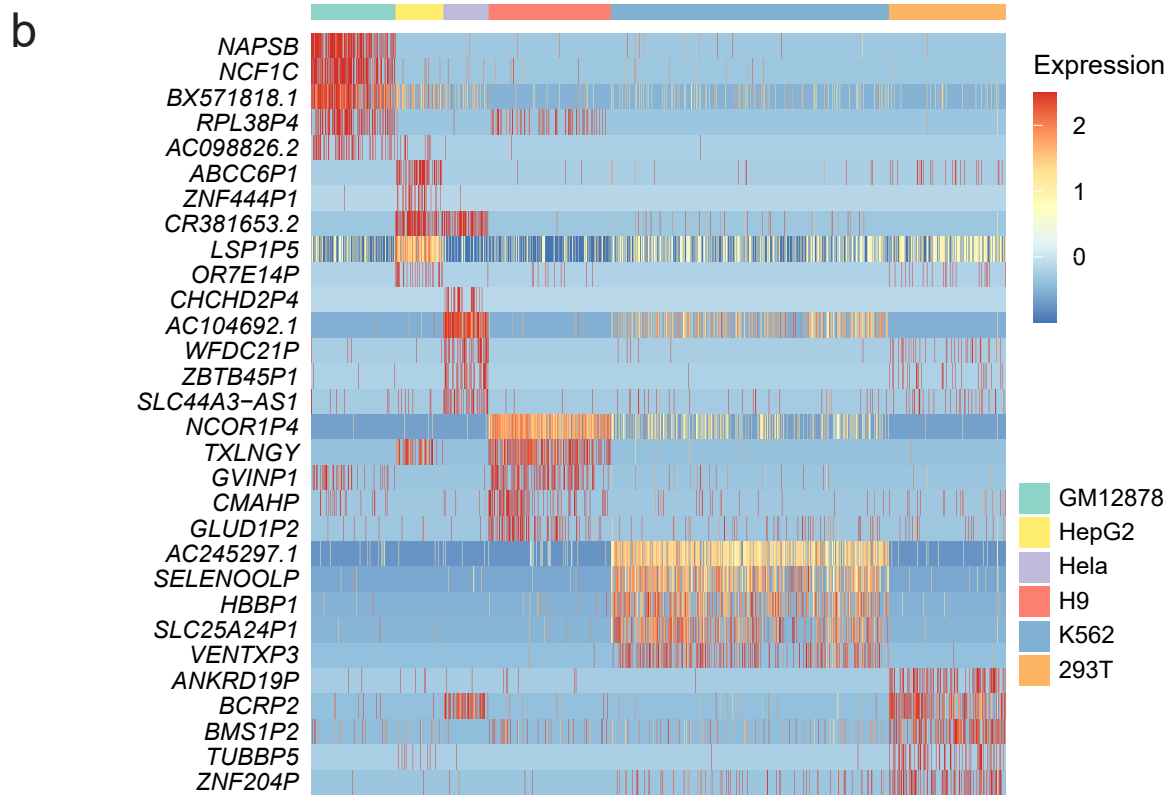
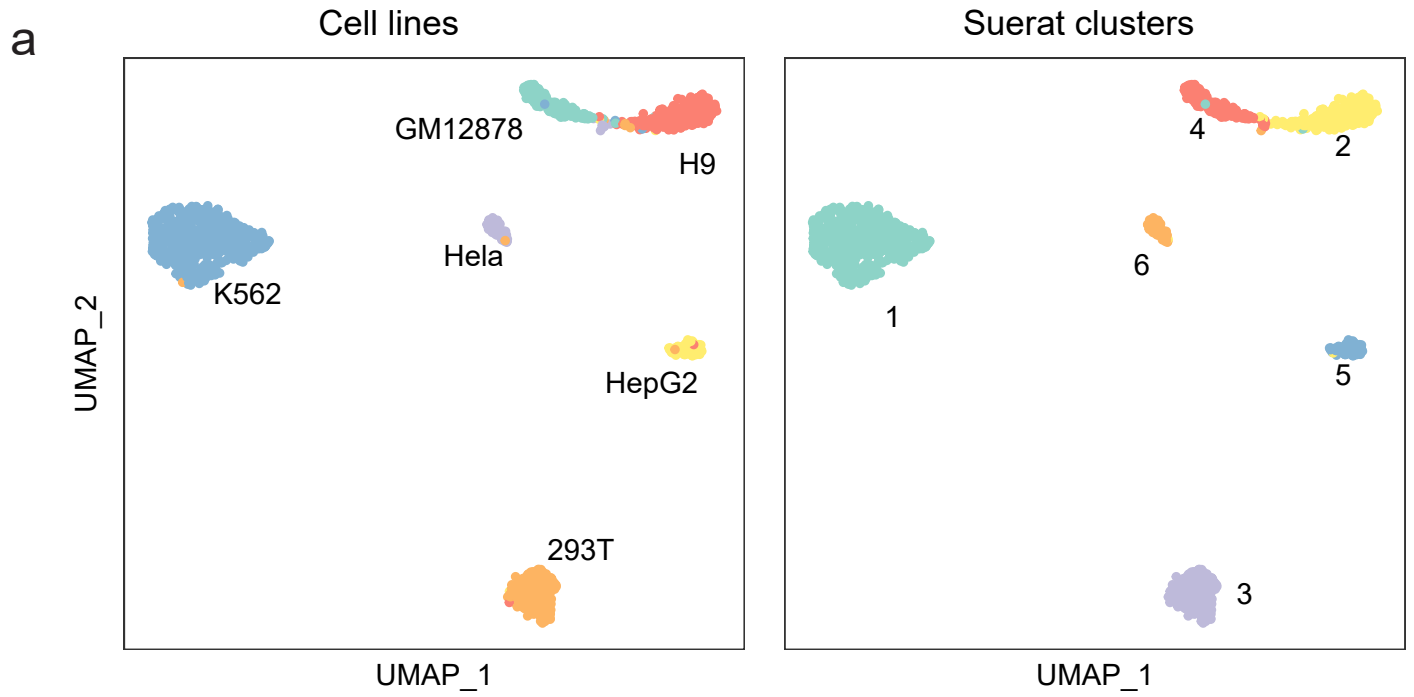
# Supplementary Fig. S6



# Supplementary Fig. S7

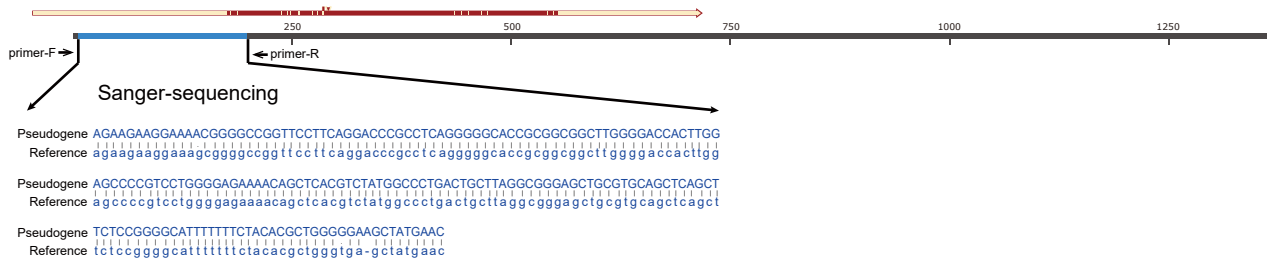


# Supplementary Fig. S8

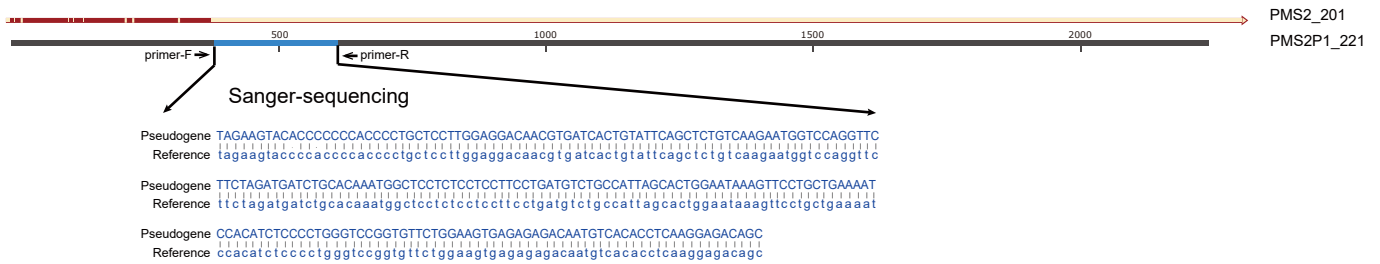


# Supplementary Fig. S9

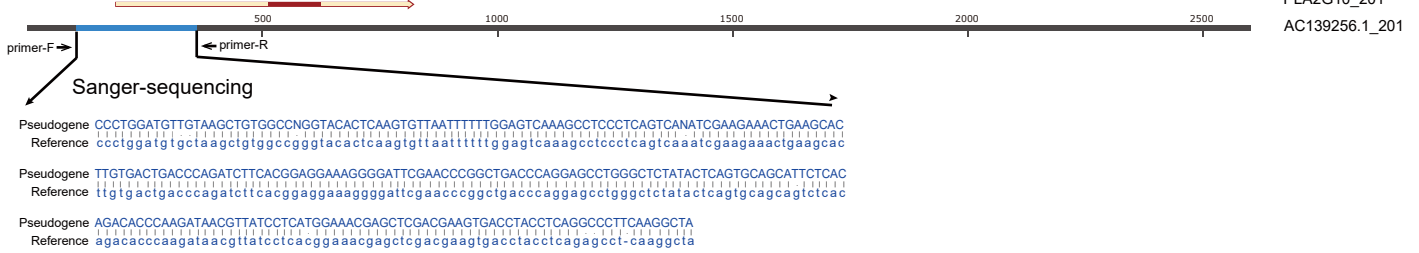
**a**



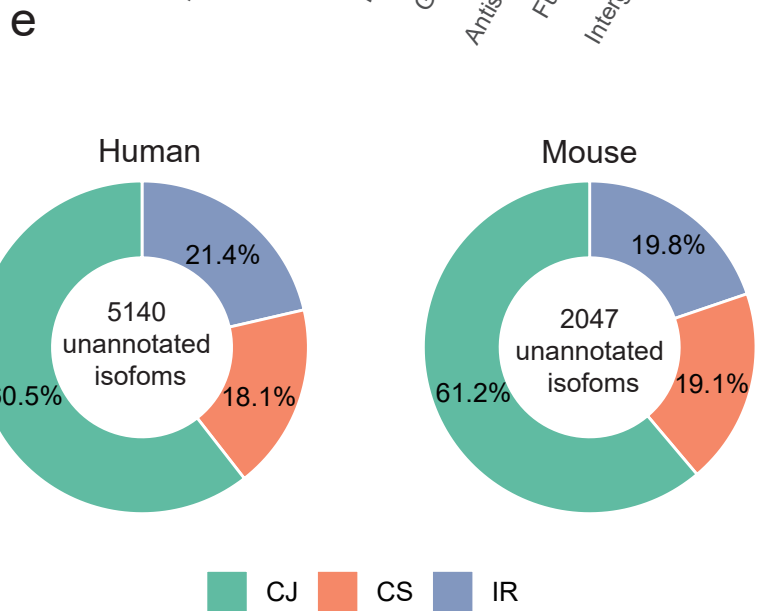
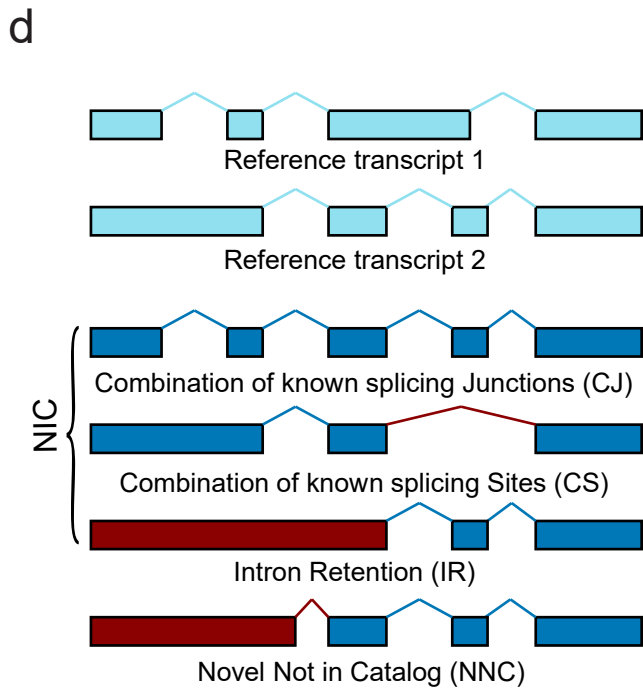
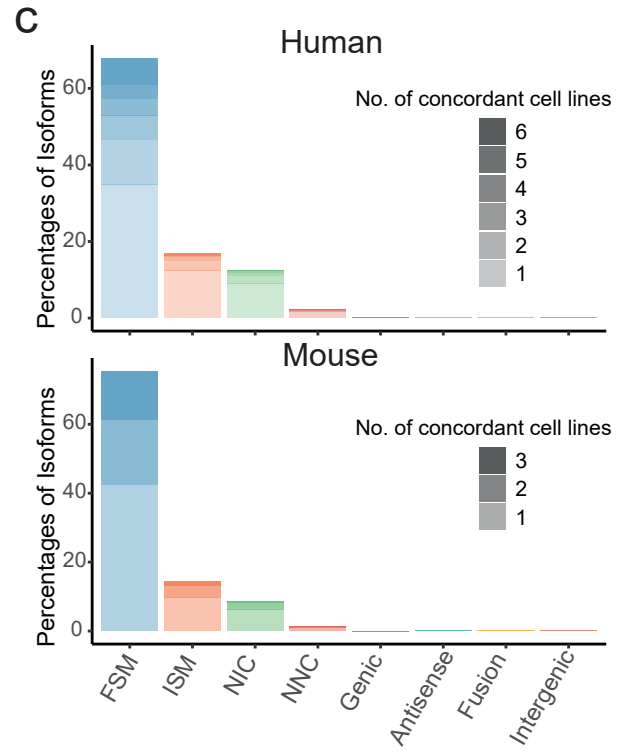
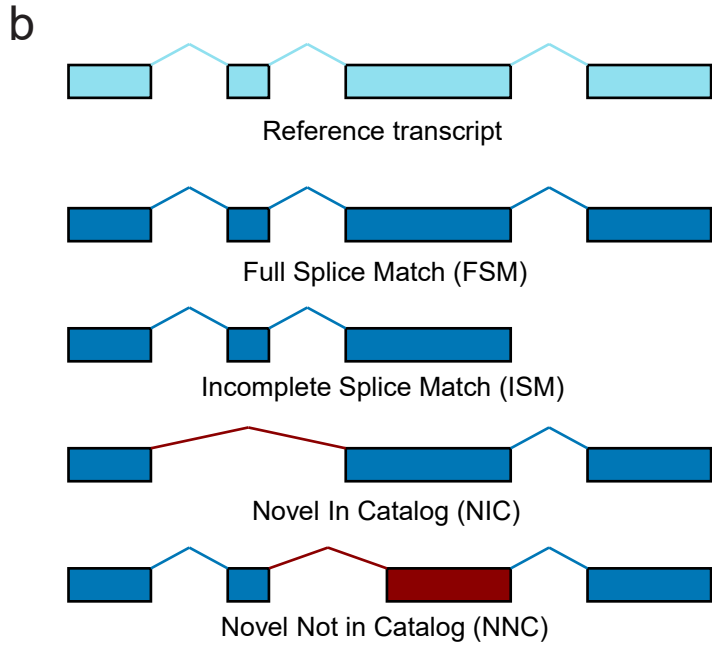
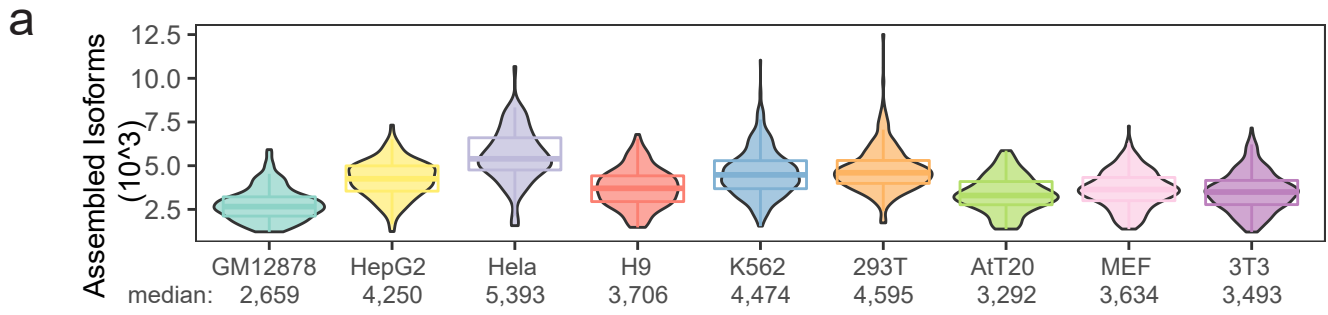
**b**



**c**

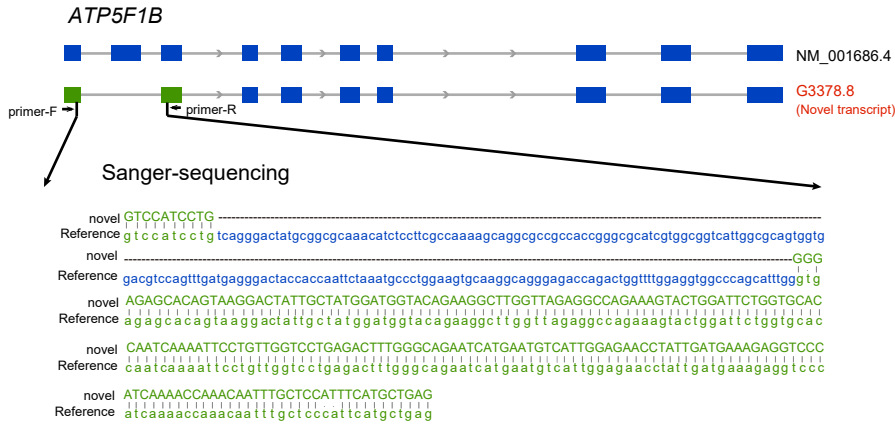


# Supplementary Fig. S10

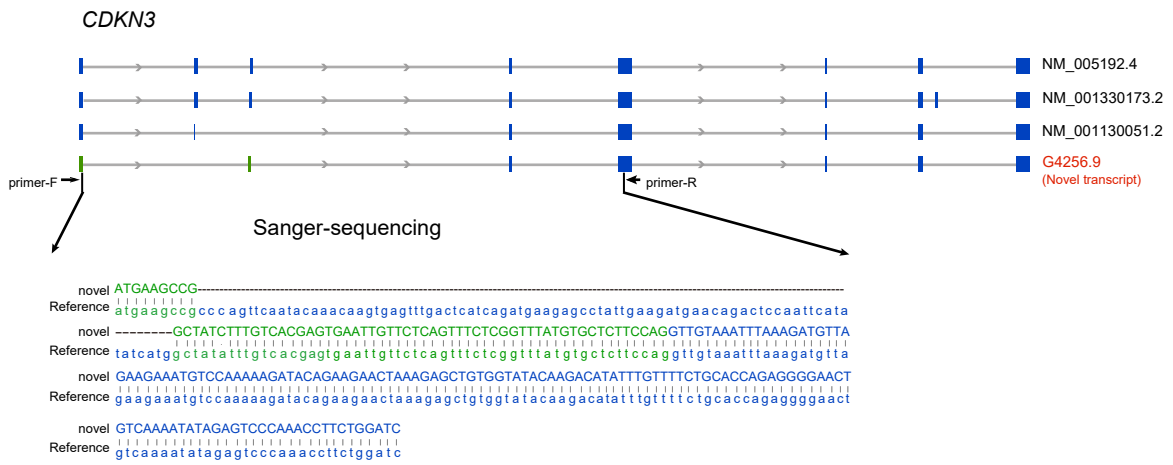


# Supplementary Fig. S11

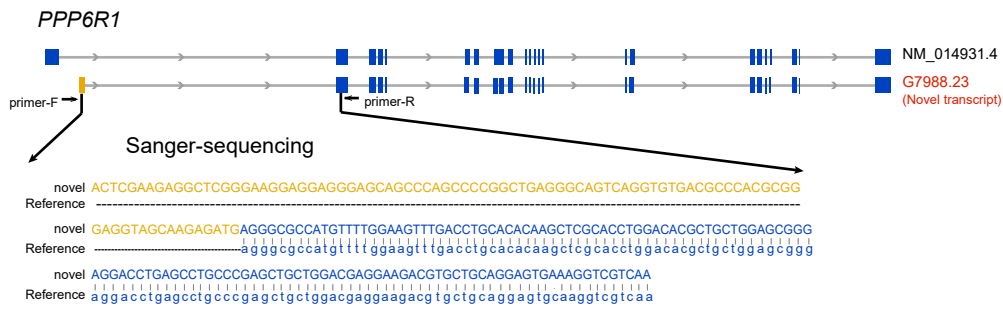
a



b



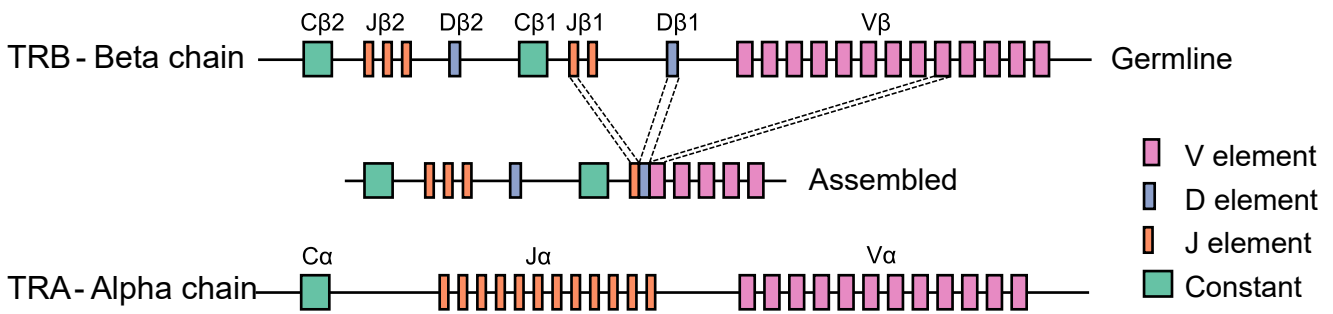
c



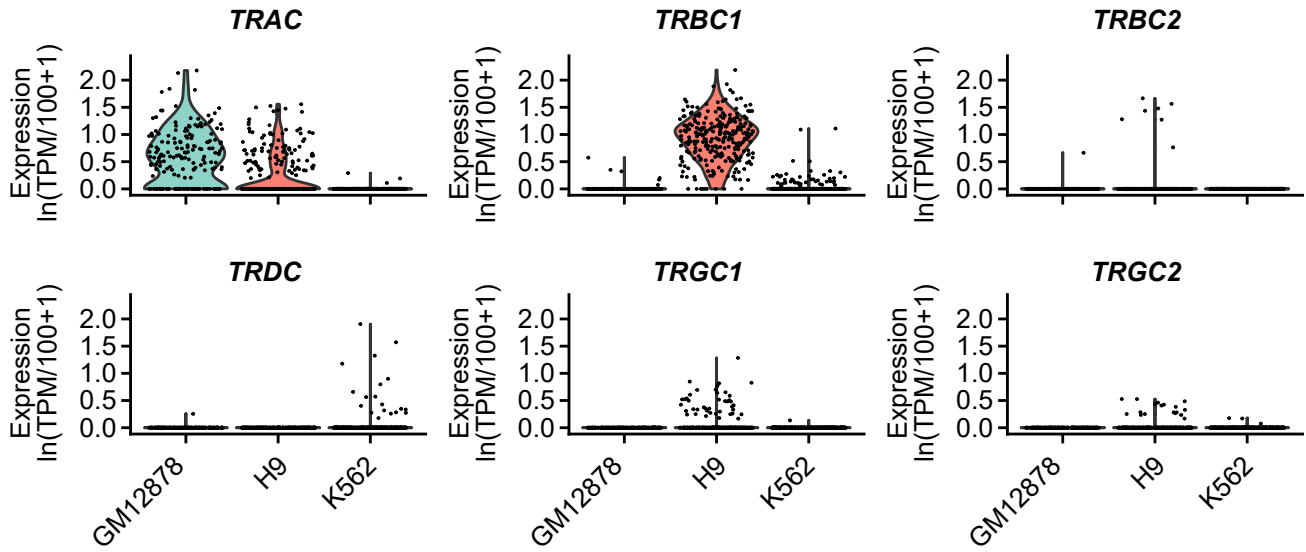
# Supplementary Fig. S12

a

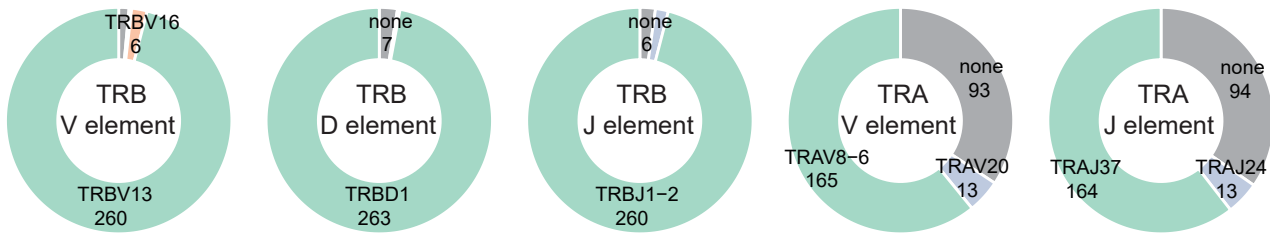
T cell receptor (TCR)



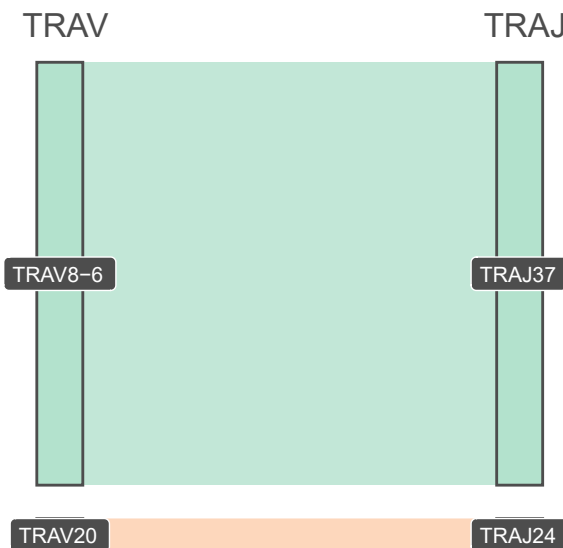
b



c



d

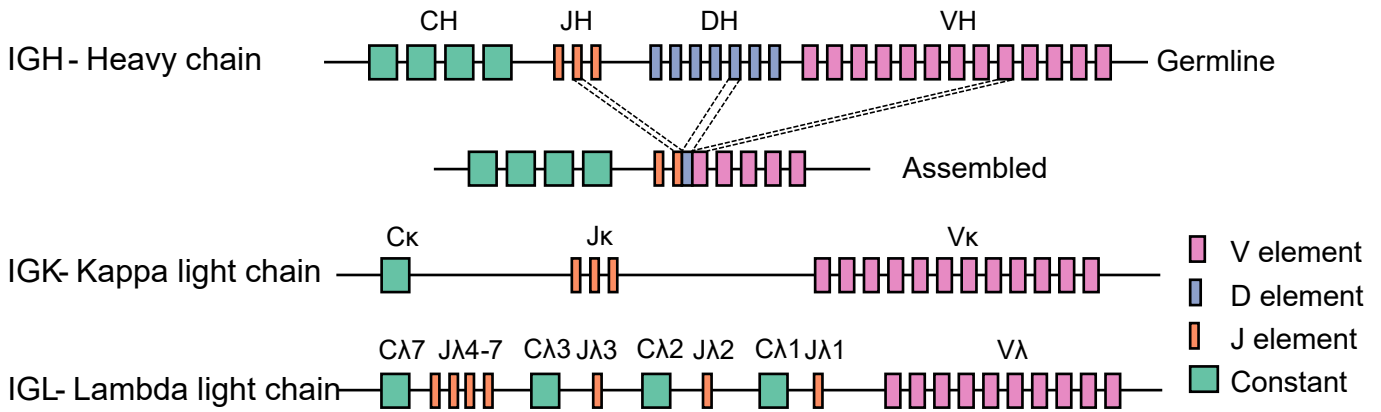




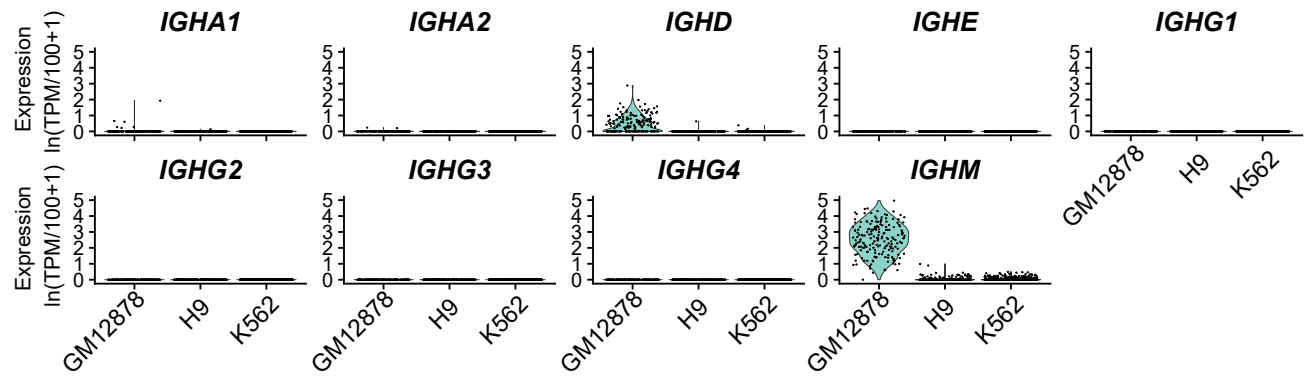
# Supplementary Fig. S13

**a**

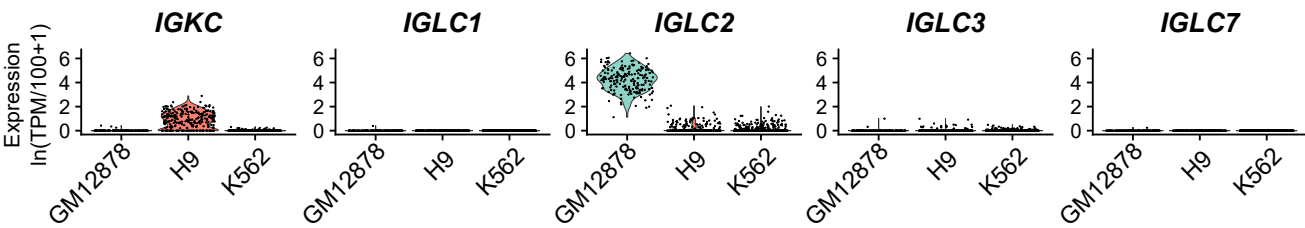
Immunoglobulin (B cell receptor, BCR)



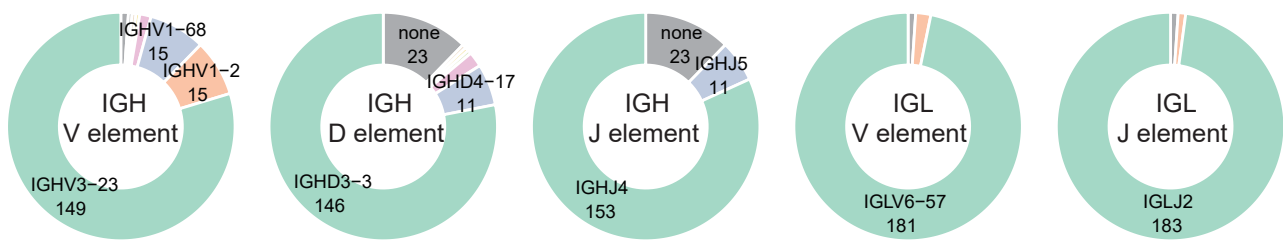
**b**



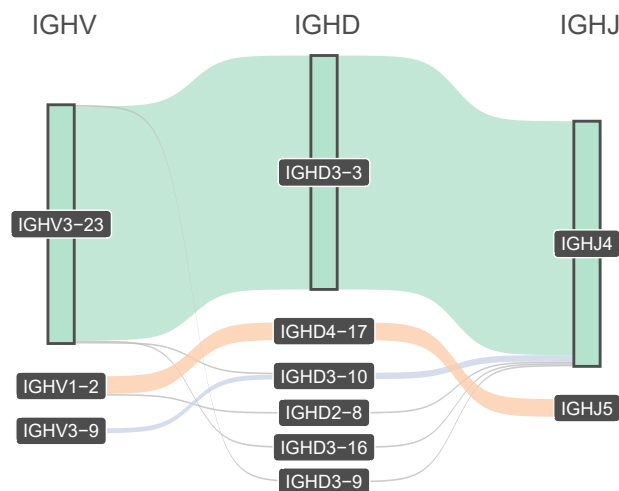
**c**



**d**



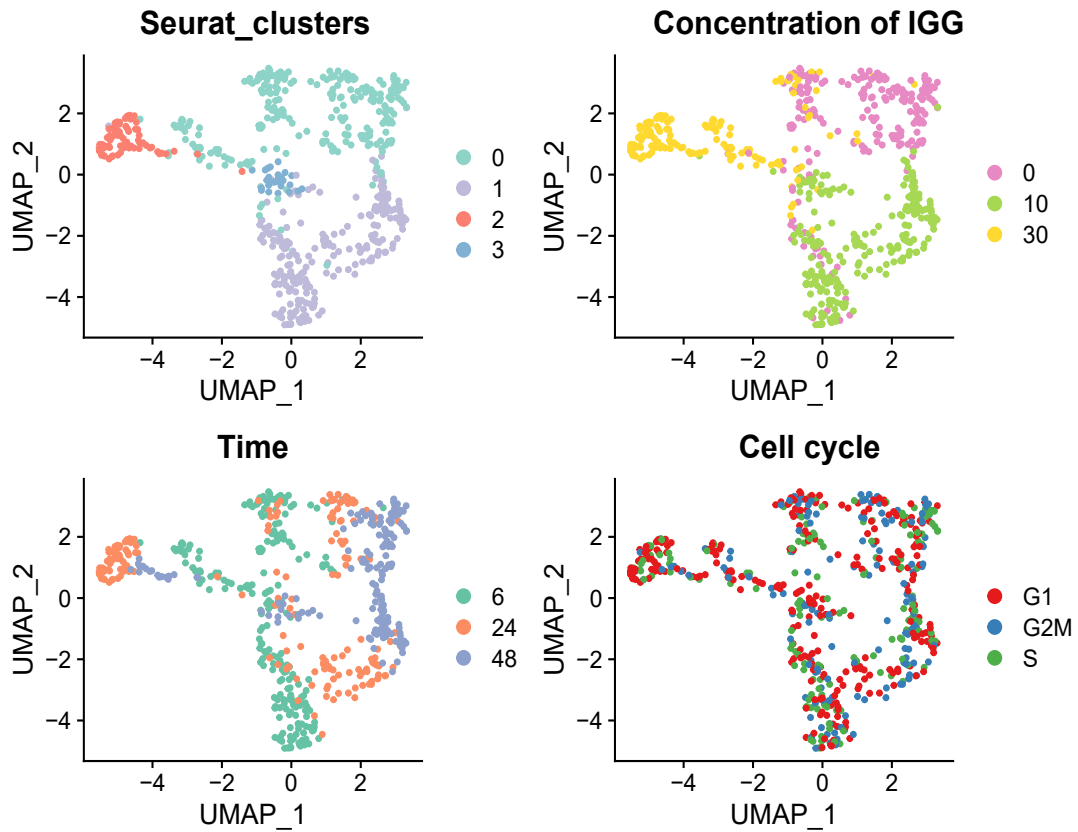
**e**



# Supplementary Fig. S14

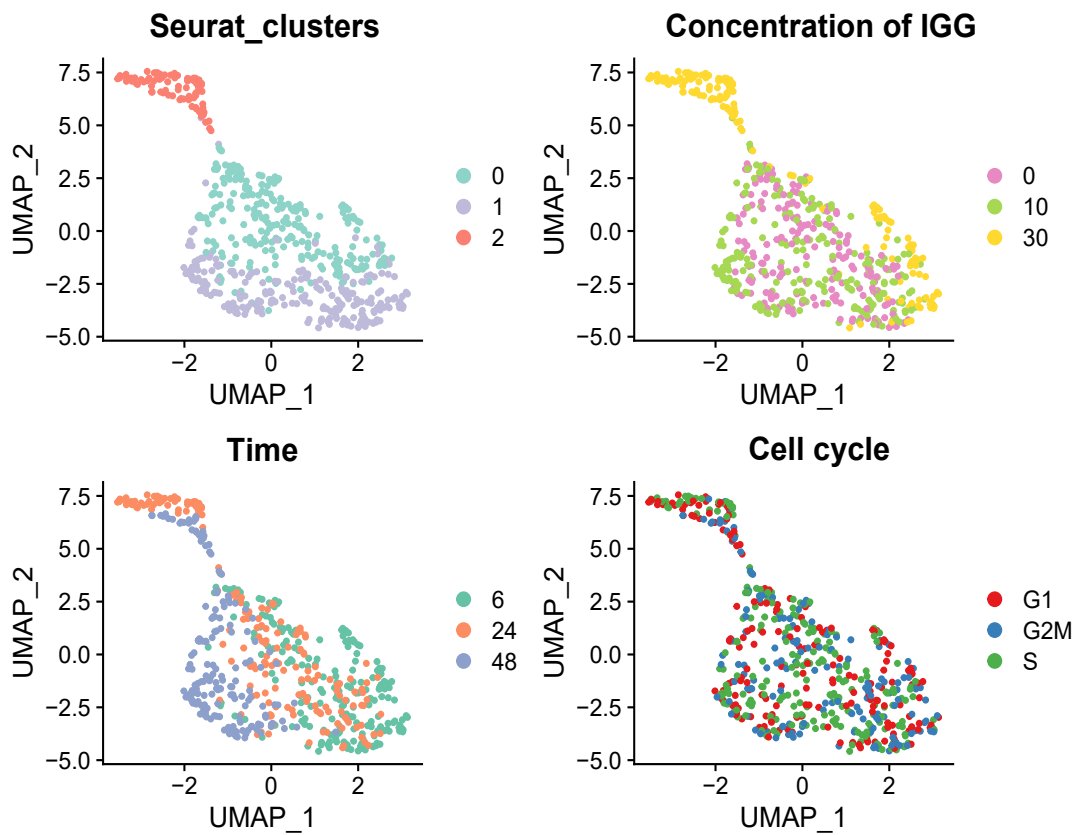
a

## SCAN-seq2

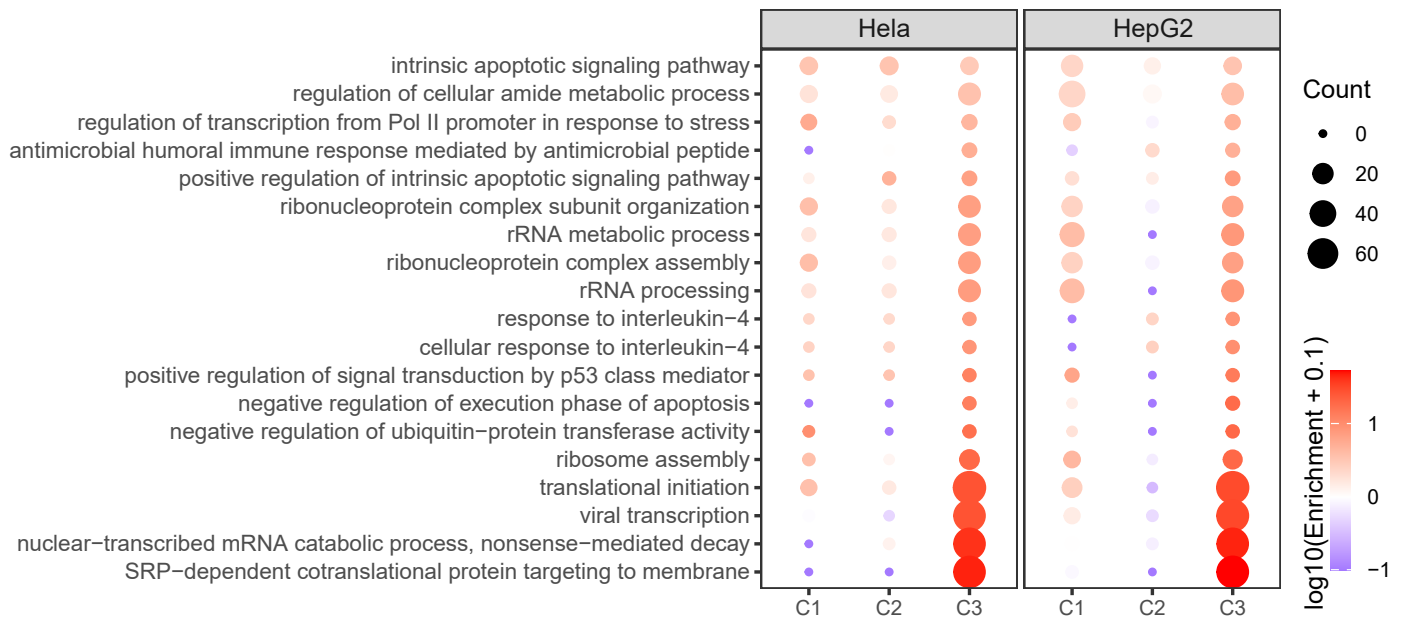


b

## NGS methods (STRT-seq)

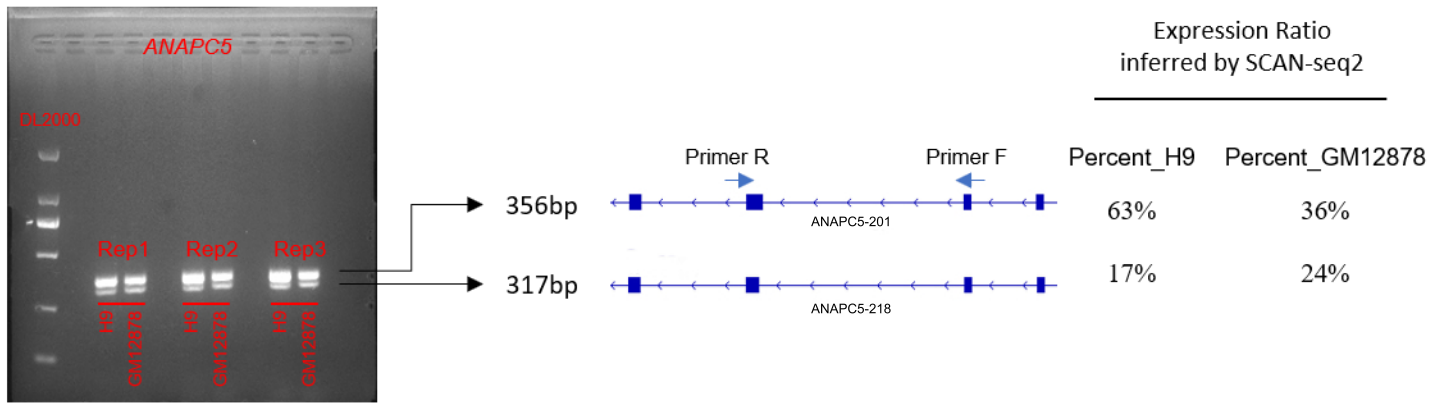


# Supplementary Fig. S15

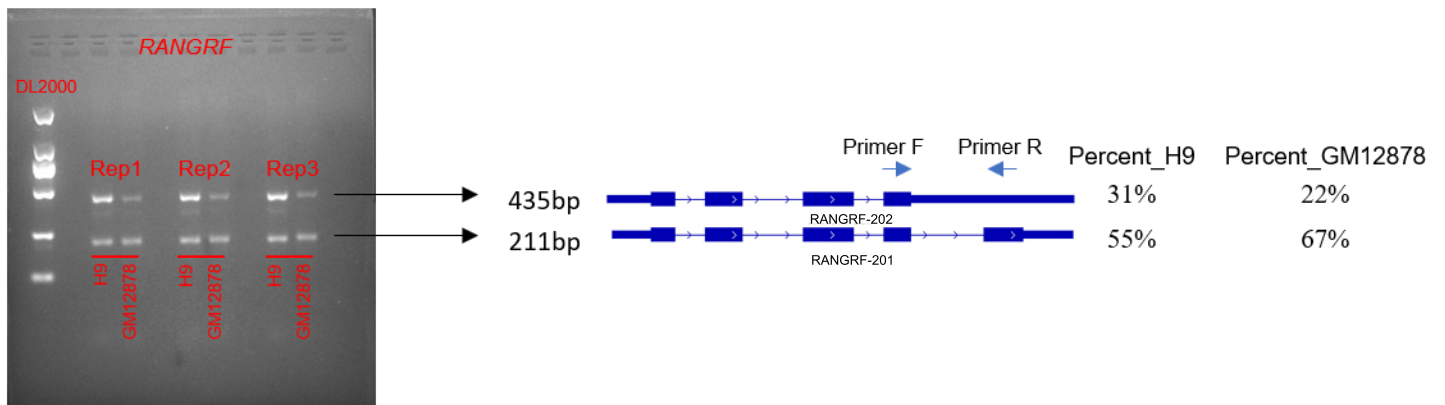


# Supplementary Fig. S16

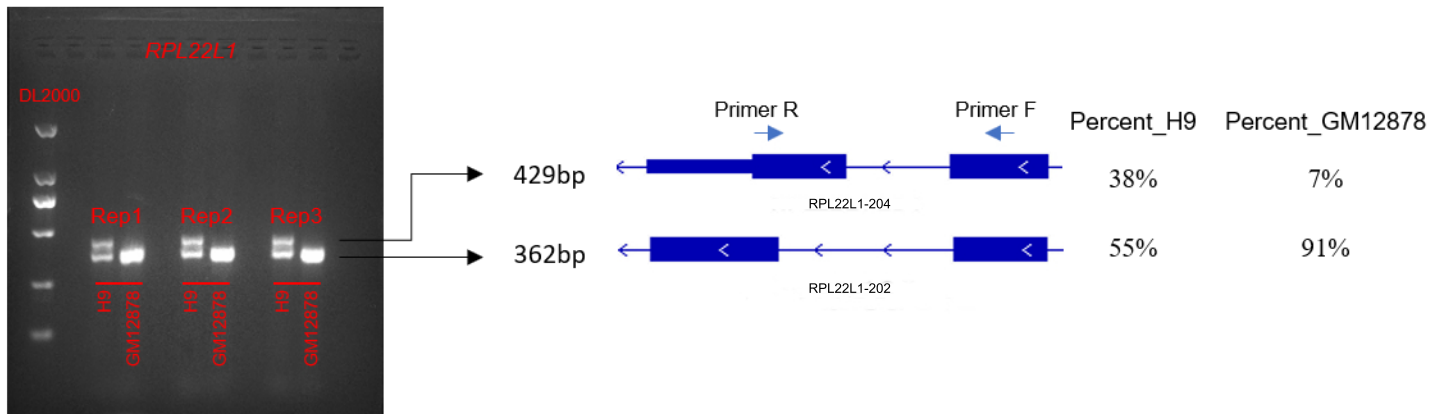
a



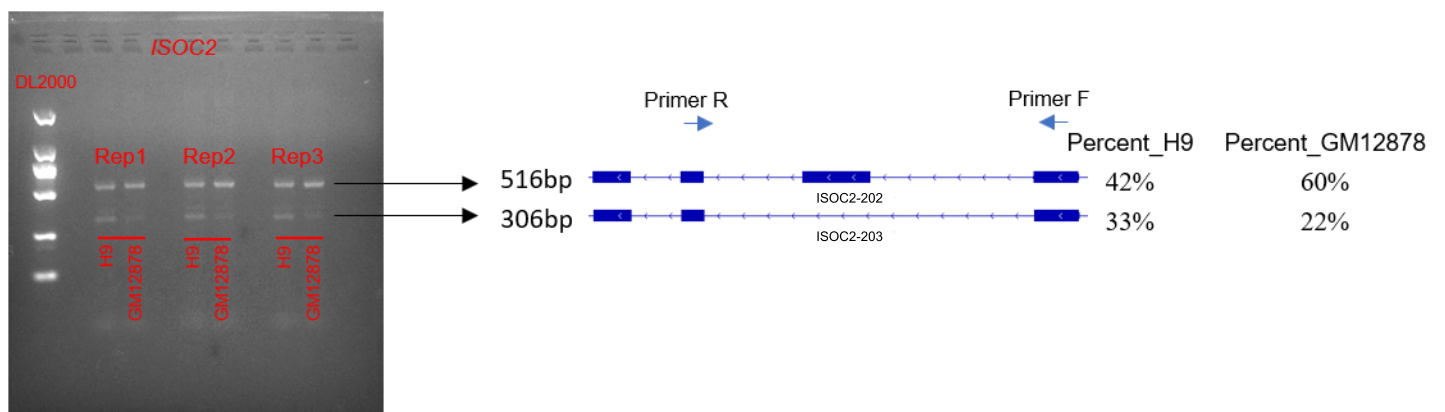
b



c



d



1 **High-throughput and high-sensitivity full-length single-cell RNA-seq**  
2 **analysis on third-generation sequencing platform**

3 Yuhan Liao<sup>1,5</sup>, Zhenyu Liu<sup>1,5</sup>, Yu Zhang<sup>1,5</sup>, Ping Lu<sup>1</sup>, Lu Wen<sup>1</sup> & Fuchou Tang<sup>1,2,3,4\*</sup>

4

5 **Materials and Methods**

6 **Experimental design**

7 We set five groups to evaluate stability and reliability of SCAN-seq2. Two groups of 960 cells  
8 were from 9 cell lines (K562, HepG2, HeLa, 293T, H9, GM12878 for human and MEF, 3T3, AtT20  
9 for mouse). One of them was sorting individual cells from each cell line into 96-well plates with  
10 known identity for each cell (Library 9CL). The other were first mixing these nine cell lines together  
11 and then sorting individual cells into 96-well plates (Library 9CL-mix, parallel group #1 to #3). One  
12 group of 96 cells (Library UMI-100) and one group of 192 cells (Library UMI-200) were both  
13 derived from human K562 cell line and mouse 3T3 cell line. The last group was from 4 cell lines  
14 (human K562 & 293T cells and mouse MEF & 3T3 cells), and we mixed every 16 human cells and  
15 16 mouse cells together after reverse transcription step and used the same 5' barcode primer to  
16 amplify them (Library 4CL).

17 **Cell culture**

18 K562 cells were maintained in RPMI 1640 medium, supplemented with 10% fetal bovine  
19 serum (FBS), 1% penicillin-streptomycin and 1% L-glutamine. GM12878 and H9 cells were both  
20 maintained in RPMI 1640 medium, supplemented with 10% FBS and 1% penicillin-streptomycin.  
21 HEK293T and MEF cells were both maintained in DMEM medium, supplemented with 10% FBS,  
22 1% penicillin-streptomycin and 1% L-glutamine. HeLa, HepG2, and 3T3 cells were maintained in  
23 DMEM medium, supplemented with 10% FBS, 1% penicillin-streptomycin. AtT20 cells were  
24 maintained in F12 medium, supplemented 15% horse serum (HS) and 2.5% FBS. All cell culture  
25 reagents were purchased from Gibco.

26 **Single cell isolation**

27 AtT20 cells were collected and washed with DPBS once, centrifuged at 300 rcf for 5 minutes,  
28 then resuspended with 1 mL medium. Same for K562 cells, H9 cells and GM12878 cells. For other  
29 5 cells lines, cells were washed with DPBS, then digested with 0.05% trypsin at 37°C for 1-2  
30 minutes. Cells were then centrifuged at 300 rcf for 5 minutes, resuspended with 1 mL medium.  
31 Followed by staining with 7-AAD, live single cells were sorted into individual wells of 96-well  
32 plates by FACS.

33 **SCAN-seq2 single cell amplification**

34 Cells from each line were sorted into 96-well plates containing lysis buffer. The lysis buffer  
35 comprised 2U RNase Inhibitor (Takara, Cat. 2313B), 0.475% Triton X-100 (Sigma-Aldrich, Cat.  
36 X100), 2.5 mM dNTP mixture (Thermo, Cat. R0193), 0.75 μM RT primer  
37 (TCAGACGTGTGCTCTTCCGATC-XXXXXXXXXXXXXXXXXXXXXXXXXXXX-N8-T25, with X  
38 representing the nucleotide of cell-specific barcode, N8 representing unique molecular identity),  
39 and 0.025% ERCC spike-in. Plates were thoroughly vortexed for 60 seconds and incubated at 72°C  
40 for 10 minutes, so that the linearized RNA molecules were released. Then they were immediately  
41 transferred on ice. Next, we added 2.85 μL RT mixture into each well, which consisted of 10U  
42 Maxima H-minus RT enzyme (Thermo, EP0752), 2.5U RNase Inhibitor, 40 mM DTT, 12.5 mM  
43 MgCl<sub>2</sub> (Sigma-Aldrich, Cat. 63020), 150 mM NaCl, 125 mM Tris-HCl pH8.3, 25% PEG 4000, 5  
44 mM GTP, and 10 μM TSO primer (5' biotin-AAGCAGTGGTATCAACGCAGAGTACATrGrG+G,  
45 with rG representing riboguanosines and +G representing the locked nucleic acid (LNA)-modified  
46 guanosine). The RT reaction was carried out at 42°C for 90 minutes, 11 cycles for 50°C for 2 minutes,

47 42°C for 2 minutes, and 85°C for 15 minutes to deactivate the enzyme. After that, plates were spun  
48 down. We pooled mRNA-cDNA hybrid strands of every 32 cells with different 3' barcodes together  
49 and purified with 0.8X Ampure XP beads (Beckman, Cat. A63882) once. PCR mixture that included  
50 2× KAPA HiFi Hot-Start Ready Mix, 266 nM 3' P2 primer  
51 (GTGACTGGAGTTCAGACGTGTGCTCTTCCGATC), 266 nM IS PCR oligo (ATGC-  
52 XXXXXXXXXXXXXXXXXXXXXXXXXXXX-AAGCAGTGGTATCAACGCAGAGT, with X  
53 representing the nucleotide of cell-specific barcode) was added into each tube. The amplification  
54 was performed by the following program: 4 cycles at 98°C for 20 seconds, 62°C for 30 seconds, and  
55 72°C for 5 minutes, followed by 15 cycles at 98°C for 20 seconds, 67°C for 15 seconds, and 72°C  
56 for 5 minutes, with a final cycle at 72°C for 5 minutes. Then, we pooled the cDNAs with different  
57 5' barcodes together and purified twice with 0.6X and 0.8X Ampure XP beads, respectively. And  
58 we quantified cDNA products with Equalbit® dsDNA HS Assay Kit (Vazyme, Cat. EQ111-01/02).  
59 Up to 1 µg cDNA products were used for further library construction.

## 60 **SCAN-seq2 library preparation and sequencing**

61 We constructed the library for Nanopore sequencing with Ligation Sequencing Kit 1D (ONT,  
62 Cat. SQK-LSK109). Briefly, the cDNA fragments were end-repaired and added dA-tailed using the  
63 Ultra II End Prep module (NEB, Cat. E7546). Then 1D adapters were tethered to above products by  
64 Quick Ligation Module (NEB, Cat. E6056). After that, each cDNA library was loaded into one R9.4  
65 chip and sequenced on PromethION 48.

## 66 **IGG treatment**

67 Isoginkgetin (IGG) (MCE, Cat. HY-N2117) was dissolved in DMSO and stored in -80°C. HeLa  
68 and HepG2 cells were pre-cultured to achieve sufficient amount. Then cells were equally divided  
69 into 96-well plates and each well contained 3,000 cells and 100 µL medium. We set three IGG  
70 concentration gradients (0 µM, 10 µM, and 30 µM) and added them to the culture medium  
71 correspondingly. Both cell lines had at least three parallel groups for each gradient. Then we  
72 digested cells at 6 h, 24 h, and 48 h after dosing treatment, respectively. Living cells were sorted by  
73 FACS after 7-AAD staining, and 2,304 cells were collected in total.

## 74 **SCAN-seq2 data processing**

75 Base calling was performed on the electric signals of nanopore sequencing to generate fastq  
76 files using Guppy (v4.0.1). Single-cell barcodes were extracted from 150 bp on both ends of the  
77 reads by nanoplexer (v0.1.2, <https://github.com/hanyue36/nanoplexer>). Reads with low quality ( $q < 7$ )  
78 and short length ( $< 100$  bp) were filtered out using NanoFilt (v2.7.1) [1]. Pychopper (v2.5.0,  
79 <https://github.com/nanoporetech/pychopper>) was utilized to trim adaptors, identify and orient full-  
80 length cDNA sequence. UMIs were extracted from 3' end of reads and added to header by UMI-  
81 tools (v1.0.1) [2] extract command. Poly-A sequences were trimmed using cutadapt (v3.2) [3] with  
82 parameters '-a "A{10}" -e 0.2'. Clean reads were mapped to a merged reference transcriptome of  
83 human (GRCh38.p13) and mouse (GRCm38.p6) cDNAs from Ensembl release 101 using minimap2  
84 (v 2.17-r941) [4]. PCR duplications were filtered based on UMI sequence and mapping position  
85 using dedup command of UMI-tools (v1.0.1) [2] with parameters "--method=directional --edit-  
86 distance-threshold=1 --per-gene --per-contig --buffer-whole-contig". Expression levels of each  
87 transcript and each gene were quantified using Salmon (v1.3.0) [5] in alignment-based mode.  
88 Transcript per million (TPM) of each gene and transcript was calculated as UMI count per million  
89 unique UMIs.

## 90 **Pseudogene analysis**

91 In the human and mouse genomes, a number of pseudogenes involving in complex gene  
92 regulatory networks and with potential as cancer biomarkers were identified. There are four major  
93 types of pseudogenes in human genome - the processed (~10,000), unprocessed (~3,500), unitary  
94 (~200) and polymorphic (~50) pseudogenes with different mechanisms of origin [6]. We identified  
95 pseudogenes based on Ensembl annotations. Genes with following biotypes were defined as  
96 pseudogenes in this research: processed pseudogene, unprocessed pseudogene, unitary pseudogene,  
97 polymorphic pseudogene. Other pseudogenes, including rRNA pseudogenes and pseudogenes of  
98 immunoglobulin and T cell receptor were not included. For processed pseudogene and unprocessed  
99 pseudogene, corresponding parent genes were also identified based on Ensembl gene annotation.

100 The expression level of pseudogenes was defined as the average TPM of top 100 cells with  
101 highest expression level of this gene. Correlation between pseudogene and parent gene was  
102 evaluated by Spearman's rank-order correlation. The absolute value of correlation coefficient larger  
103 than 0.2 was considered significant.

## 104 **Dimensionality reduction and clustering**

105 Further analysis based on expression matrix was performed with Seurat package (v4.0.3) [7].  
106 Cells with less than 2000 genes detected or more than 15% mitochondrial UMIs were discarded.  
107 Top 2,000 highly variable genes were selected with Seurat FindVariableFeatures function. Principle  
108 component analysis (PCA) was then performed on highly variable genes. Unsupervised clustering  
109 of cells was performed using original Louvain algorithm of Seurat FindClusters function. UMAP  
110 embedding was calculated to visualize cluster and cell type information.

## 111 **Differential gene expression (DGE) and differential transcript usage (DTU) analysis**

112 Differential gene expression analysis was performed using Wilcoxon rank sum test on  
113  $\log_2(\text{transcript per million})$  value. Genes with absolute  $\log_2$ -transformed fold change of  $>1$ , and an  
114 adjusted P value of  $P < 0.05$  were considered as differentially expressed.

115 Differential transcript usage analysis was performed with R package DTUrtle (v1.0.2) [8]. First,  
116 transcripts with less than 5 UMI support or detected in less than 25 cells were removed to reduce  
117 multiple testing. The Dirichlet-multinomial model was used to estimate the precision parameter.  
118 Next, a group-wise maximum likelihood estimation of transcript proportions was calculated. A  
119 likelihood ratio test for transcript proportions was used to identify DTU. Genes with false discovery  
120 rate (FDR)  $< 0.05$  were considered significant.

## 121 **Transcriptome assembly**

122 Unique reads after deduplication were mapped to reference genome of human (GRCh38.p13)  
123 or mouse (GRCm38.p6) from Ensembl release 101 using minimap2 (v 2.17-r941) [4]. Reads with  
124 MAPQ  $< 30$  were discarded. We performed transcriptome assembly for each single cell using  
125 StringTie (v2.1.7) [9] in long reads processing mode. Gene annotation gtf files from Ensembl were  
126 used to guide the assembly process. Single-cell assemblies were classified using the sqanti3\_qc.py  
127 script of SQANTI3 (v3.4.1) [10] with the parameters "--skipORF --report pdf" and then filtered  
128 using the sqanti3\_RulesFilter.py script of SQANTI3 (v3.4.1) with default parameters.

129 Single-cell transcriptome assemblies were integrated into merged assembly in a hierarchical  
130 manner using the merge subcommand of TAMA (v0.0) [11]. Briefly, single-cell assemblies were  
131 first merged into 9 cell-line assemblies with the parameter "--a 100 -m 20 -z 50". Only transcript



132 supported by more than 5 cells were retained in cell line assemblies. Next, 6 cell-line assemblies for  
133 human and 3 cell-line assemblies for mouse were merged respectively with the parameter “-a 10 -  
134 m 10 -z 10”. The merged assemblies were then compared with Ensembl gene annotation (release  
135 101) and filtered using SQANTI3 (v3.4.1) with identical parameters as used for single-cell  
136 assemblies.

### 137 **Consensus and polishing of TCR and immunoglobulin (BCR) sequences**

138 In each GM12878 or H9 cell, reads mapped to immunoglobulin or T cell receptor gene locus  
139 (IGH: chr14:105,586,437-106,879,844; IGL: chr22:22,026,076-22,922,913; TRA:  
140 chr14:21,621,904-22,552,132; TRB: chr7:142,299,011-142,813,287) were extracted. For each  
141 single cell, extracted reads were initially clustered using usearch (v.11.0.667) [12] -cluster\_fast -id  
142 0.75 -sizeout -centroids. The centroid read of the largest group was selected as representative  
143 sequence and used as the template for 4 rounds of polishing using all reads from the same cluster  
144 with minimap2 (v 2.17-r941) [4] -x map-ont followed by racon (v1.5.0) [13] -w 200 -m 8 -x -6 -g -  
145 8 -q 7. The racon-polished sequence was further corrected using all reads of the same cluster with  
146 Medaka (<https://github.com/nanoporetech/medaka>) consensus -m r941\_min\_high\_g360. The  
147 medaka-corrected sequences were regarded as the TCR/immunoglobulin transcripts and were utilized  
148 for the identification of V(D)J recombination. This pipeline was reported to produce sequence of  
149 ~99.995% accuracy at 25X coverage during amplicon sequencing [14].

### 150 **Identify V(D)J recombination for TCR and immunoglobulin (BCR)**

151 Genes encoding variable regions of B- and T- lymphocyte antigen receptors are assembled by  
152 recombination of variable (V), diversity (D), and joining (J) gene segments [15-16]. The V(D)J  
153 elements in the corrected TCR/immunoglobulin transcripts were identified with Igbblast (v 1.17.1) [17]  
154 with parameters “-organism human -show\_translation -outfmt 19” for immunoglobulin and an extra  
155 parameter “-ig\_seqtype TCR” for TCR. Annotation for human VDJ elements were downloaded  
156 from the international immunogenetics information system (IGMT). Hits with the smallest E value  
157 were retained for each cell.

158 Subclones of GM12878 and H9 cells were identified based on the V(D)J combination of each  
159 cell. Briefly, GM12878 cells with the same VDJ elements for immunoglobulin heavy chain and same  
160 VJ elements for light chain were considered as the same subclone. H9 cells with the same VDJ  
161 elements for  $\beta$  chain and same VJ elements for  $\alpha$  chain were considered as the same subclone.  
162 Subclones with more than 5 cells were considered solid and retained for further analysis.

163

164 **Availability of data and materials**

165 All relevant data are available from the Gene Expression Omnibus (GEO) database (accession  
166 number: GSE203561).

167

168 **References**

- 169 [1] De Coster, W. et al. NanoPack: visualizing and processing long-read sequencing data.  
170 *Bioinformatics* **34**, 2666-2669 (2018).
- 171 [2] Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in Unique Molecular  
172 Identifiers to improve quantification accuracy. *Genome Research* **27**, 491-499 (2017).
- 173 [3] Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
174 *EMBnet.journal* **17**, 10-12 (2011)
- 175 [4] Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100  
176 (2018).
- 177 [5] Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-  
178 aware quantification of transcript expression. *Nature Methods* **14**, 417-419 (2017).
- 179 [6] Qian, S. H., Chen, L., Xiong, Y.-L. & Chen, Z.-X. Evolution and function of developmentally  
180 dynamic pseudogenes in mammals. *Genome Biology* **23** (2022).
- 181 [7] Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e3529  
182 (2021).
- 183 [8] Tekath, T., Dugas, M. & Boeva, V. Differential transcript usage analysis of bulk and single-cell  
184 RNA-seq data with DTUrtle. *Bioinformatics* **37**, 3781-3787 (2021).
- 185 [9] Kovaka, S. et al. Transcriptome assembly from long-read RNA-seq alignments with StringTie2.  
186 *Genome Biology* **20** (2019).
- 187 [10] Tardaguila, M. et al. SQANTI: extensive characterization of long-read transcript sequences for  
188 quality control in full-length transcriptome identification and quantification. *Genome Research* **28**,  
189 396-411 (2018).
- 190 [11] Kuo, R. I. et al. Illuminating the dark side of the human transcriptome with long read transcript  
191 sequencing. *BMC Genomics* **21** (2020).
- 192 [12] Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**,  
193 2460-2461 (2010).
- 194 [13] Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly  
195 from long uncorrected reads. *Genome Research* **27**, 737-746 (2017).
- 196 [14] Karst, S. M. et al. High-accuracy long-read amplicon sequences using unique molecular  
197 identifiers with Nanopore or PacBio sequencing. *Nature Methods* **18**, 165-169 (2021).
- 198 [15] Schatz, D. G. & Ji, Y. Recombination centres and the orchestration of V(D)J recombination.  
199 *Nature Reviews Immunology* **11**, 251-263 (2011).
- 200 [16] De Simone, M., Rossetti, G. & Pagani, M. Single Cell T Cell Receptor Sequencing: Techniques  
201 and Future Challenges. *Frontiers in Immunology* **9** (2018).
- 202 [17] Ye, J., Ma, N., Madden, T. L. & Ostell, J. M. IgBLAST: an immunoglobulin variable domain  
203 sequence analysis tool. *Nucleic Acids Research* **41**, W34-W40 (2013).
- 204

205 **Figure legend**

206 **Fig. 1 SCAN-seq2 technical performances and analysis of Isoginkgetin (IGG) responses in cell**  
207 **lines.** a. Schematic diagram of SCAN-seq2 library construction. N different single cells are labeled  
208 with 3' barcode during reverse transcription and pooled into the same tube for PCR amplification.  
209 M different tubes are pooled together and sequenced with Nanopore platform, allowing parallel  
210 sequencing of N×M cells. b. Number of detected genes (top) and isoforms (bottom) in 852 cells  
211 from library 9CL. Median values are labeled under each cell line. c. Pearson correlation of ERCC  
212 concentration and sequenced UMIs in each library. d. Correlation between gene expression  
213 quantification of SCAN-seq2 and Smart-seq3 in 293T cells. Single cells are aggregated into pseudo-  
214 bulk for comparison. e-f. UMAP embeddings of HeLa (e) and HepG2 (f) cells after IGG treatment  
215 at different concentration and time. Cells are colored by unsupervised clustering results (top left),  
216 cell cycle phase (top right), IGG concentration (bottom left) and time of treatment (bottom right).  
217 IGG-responsive clusters are highlighted in red circles. g. Venn diagram showing the overlap in  
218 upregulated differentially expressed genes (DEGs) and differential transcript usage (DTU) in IGG-  
219 responsive cluster of HeLa and HepG2 cells. h. Venn diagram of DEGs and genes with significant  
220 DTU in IGG-responsive cluster of HeLa cells. i. Fraction of each subcategory for NIC transcripts in  
221 different clusters. P values are calculated by two-tailed Wilcoxon rank-sum test. \*:  $p < 0.05$ , \*\*:  $p$   
222  $< 0.01$ , \*\*\*:  $p < 0.001$ . j. Examples of genes with significant differential transcript usage in IGG-  
223 responsive cluster. Exons with different usage are highlighted in red.

224

225 **Supplementary Fig. S1 Flowchart of SCAN-seq2 data processing.** a. Demultiplexing, trimming,  
226 filtering, and deduplication of raw Nanopore reads. b. Isoform expression quantification and  
227 reference-guided transcriptome assembly. Software utilized in each step is indicated next to the lines.  
228 The nodes are colored by the scale of the data, either single-cell level (green), cell line level (red),  
229 or all cells merged (purple). The number of reads, genes, and transcripts is also labeled in each node.  
230 For single cell level analysis, the number is calculated as the average value of all cells. The file  
231 format of each node is labeled in brackets.

232

233 **Supplementary Fig. S2 Statistics of SCAN-seq2 reads after each processing step.** a. Boxplot of  
234 reads statistics after demultiplexing, quality control, and deduplication in every single cell.  
235 Mean\_read\_quality (Q-score) b. Circle plot showing number of reads and total bases retained after  
236 each processing step. The total yield of a single PromethION is defined as a whole circle.

237

238 **Supplementary Fig. S3 Overview of SCAN-seq2.** a. Structure of SCAN-Seq2 library. Barcodes  
239 are introduced to both ends of cDNA for massively parallel analysis of thousands of single cells. b.

240 Donut chart showing the distribution of sequenced reads from library 9CL. 75.9% of reads have  
241 complete library structure and can be used for downstream analysis. c. Fragment analysis (FA) of  
242 library 9CL. d. Histogram indicating the length of cDNA sequence after adaptor trimming. e.  
243 Mapping rate of reads from each cell line. The median values are labeled below. f. Barnyard plot of  
244 710 cells from library 4CL. Dotted line indicates specificity of 90%. Red dots indicate human-  
245 specific barcodes. Blue dots indicate mouse-specific barcodes. One barcode associated with both  
246 human and mouse transcript (Gray). g. Saturation analysis of SCAN-seq2. Number of detected  
247 genes and isoforms reach plateau at 400,000 reads per cell. h. Cell-to-cell correlation of SCAN-seq2  
248 gene quantification in 293T cells, comparing with Smart-seq3. i. Expression level of *PTPRC* gene  
249 (top left) and its protein-coding isoforms in 3 immune cell lines. GM12878 cells mainly utilize  
250 PTPRC-209 while K562 cells mainly utilize PTPRC-201. j. Alternative splicing at exon 4 (A), 5  
251 (B), 6 (C) of *PTPRC* gene. PTPRC-201 uses none of these three exons, encoding CD45 RO. PTPRC-  
252 209 uses all three exons, thus encoding CD45 RABC.

253

254 **Supplementary Fig. S4 Evaluation of sequencing errors in 24-bp barcode sequences.** Errors of  
255 Nanopore sequencing, including mismatched (top), indel (middle), and overall errors (bottom) were  
256 counted in each demultiplexed barcode. a. 5 prime barcodes. b. 3 prime barcodes.

257

258 **Supplementary Fig. S5 Comparison of SCAN-seq2 UMIs with those in NGS methods.** UMIs  
259 of ERCC reads identified by SCAN-seq2 are compared to those identified by Illumina sequencing  
260 of the same library (but the original SCAN-seq2 sequencing library has been fragmented into about  
261 300bp fragments and the fragments containing 3' ends of the original library are captured by biotin-  
262 streptavidin affinity strategy and further ligated into NGS sequencing adaptor pairs). a. No sequence  
263 error is tolerated. b. Tolerated sequence error at an edit distance of 1.

264

265 **Supplementary Fig. S6 SCAN-seq2 analysis of 9 different cell lines.** a. UMAP embedding of  
266 library 9CL. Cells are colored by cell line (left) and unsupervised clusters (right). b. Violin plot  
267 showing selected markers of each cell line. c. Heatmap of cell-line specific markers. For each cell  
268 line, 5 genes with highest fold change against other cells are included.

269

270 **Supplementary Fig. S7 Systematic evaluation of pseudogene expression in human cell lines**  
271 **using SCAN-seq2.** a. Expression measurements in SCAN-seq2 of selected pseudogenes and  
272 corresponding parent genes. 4 gene pairs with identity 95% are selected. b. Spiral chart showing  
273 pairwise sequence alignment of RPS2-201 (purple) and RPS2P46-201 (green) transcript. Gaps are

274 highlighted in red. Mismatched bases and gap sequence are labeled by corresponding nucleobases.  
275 c. Number of expressed genes under different TPM cutoff and number of supported cells.  
276 Pseudogenes are grouped by gene type from Ensembl annotation. d. Fraction of expressed  
277 pseudogenes with different expression level. TPM of each gene is calculated as the average  
278 expression level of top 100 cells with highest expression level. e. Comparison of expression levels  
279 of pseudogenes and corresponding protein-coding parent genes. Pseudogenes generally possess  
280 significantly lower expression level comparing with parent gene. f. Correlation analysis of  
281 pseudogene expression and parent gene expression in 6 human cell lines. Pairs with absolute value  
282 of spearman correlation coefficient  $> 0.25$  are considered significant. Significant pairs are labeled  
283 with the symbol of parent genes and pseudogene. Positively correlated pairs are colored in red.  
284 Negatively correlated pairs are colored in blue.

285

286 **Supplementary Fig. S8 Unsupervised clustering on pseudogene expression distinguishes**  
287 **different cell lines.** a. UMAP embedding of 6 human cell line. Reductions are calculated on  
288 pseudogenes only. Cells are colored by cell line (left) and unsupervised clusters (right). b. Heatmap  
289 of cell-line specific pseudogene markers. For each cell line, 5 pseudogenes with highest fold change  
290 against other cells are included.

291

292 **Supplementary Fig. S9 Examples of pseudogenes identified in H9 and GM12878 cells.**  
293 Schematic displays 3 pseudogene-parent gene pairs (pseudogene *CRYBB2P1* and parent gene  
294 *CRYBB2*, pseudogene *PMS2P1* and parent gene *PMS2*, pseudogene *AC139256.1* and parent gene  
295 *PLA2G10*) (a-c). The red parts indicate regions of identical sequences between the pseudogene and  
296 corresponding parent gene. The blue parts are pseudogene-specific regions, which were amplified  
297 by RT-PCR followed by Sanger-sequencing. Details are showed in Supplementary Table. S8.

298

299 **Supplementary Fig. S10 Transcriptome assembly using SCAN-seq2.** a. Number of assembled  
300 isoforms in each single cell. The median values of each cell line are labeled below. b. Schematic  
301 diagram of isoform classification by comparing with gene annotations. Disagreements with  
302 annotations are highlighted in red. Full splice match (FSM) and incomplete splice match (ISM)  
303 indicate splicing events conforming to annotations. Novel in catalog (NIC) indicates combination  
304 of known splicing site and junctions. Novel not in catalog (NNC) indicates novel splicing sites. c.  
305 Fraction of each isoform classification in 6 human cell lines (top) and 3 mouse cell lines (bottom).  
306 Transparency of bars indicates number of concordant cell lines. Isoforms detected in all cell lines  
307 are colored with lowest transparency. d. Schematic diagram of unannotated transcript classification.

308 Disagreements with annotations are highlighted in red. CJ indicates unannotated combination of  
309 splicing junctions. CS utilizes known splicing sites to create unannotated junctions. IR indicates  
310 retained intron. NNC includes unannotated splicing sites. e. Donut chart demonstrating fraction of  
311 each NIC subcategories in 6 human cell lines (left) and 3 mouse cell lines (right).

312

313 **Supplementary Fig. S11 Examples of unannotated transcripts identified in H9 and GM12878**  
314 **cells.** a. Transcript G3378.8 matches partially to *ATP5F1B* gene sequences, but lacks the second  
315 exon. b. *CDKN3* gene has three annotated transcripts whereas transcript G4256.9 lacks the second  
316 exon compared with the known three transcripts. c. Transcript G7988.23 shares a portion of  
317 sequences with *PPP6R1* gene, but uses an alternative transcription start site (TSS). Details are  
318 showed in Supplementary Table. S8.

319

320 **Supplementary Fig. S12 T cell receptor (TCR) analysis of H9 cell line.** a. Schematic diagram of  
321 V(D)J recombination in TCR beta chain (TRB) and alpha chain (TRA). In the assembled genome,  
322 only a single VDJ element is retained for expression. b. Expression of all TCR constant region genes  
323 in 3 immune cell lines. H9 cells mainly utilize *TRBC1* for beta chain. c. Identification of VDJ  
324 elements in TRB transcript and VJ in TRA transcript in each single cell. Elements with more than 5  
325 cells are labeled. Cells with insufficient reads or no significant hits are labeled as none. d. Subclone  
326 analysis of H9 cell line. Two distinct clones with different TRAV and TRAJ elements are identified.

327

328 **Supplementary Fig. S13 Immunoglobulin (BCR) analysis of GM12878 cell line.** a. Schematic  
329 diagram of V(D)J recombination in Immunoglobulin heavy chain (IGH), kappa light chain (IGK) and  
330 lambda light chain (IGL). b. Expression of all immunoglobulin heavy chain constant region genes in  
331 3 immune cell lines. Transcripts of gene *IGHM* and *IGHD* are detected in GM12878 cells, indicating  
332 that most GM12878 cells express immunoglobulin M (IgM) whereas a small proportion of  
333 GM12878 cells express immunoglobulin D (IgD). c. Expression of all immunoglobulin light chain  
334 constant region genes in 3 immune cell lines. GM12878 cells only utilize *IGLC2* from lambda light  
335 chain locus for light chain construction. d. Identification of VDJ elements in IGH transcript and VJ  
336 in IGL transcript in each single cell. The D element of IGL is not detected, and VJ rearrangement is  
337 same in essentially all GM12878 cells. Elements with more than 5 cells are labeled. Cells with  
338 insufficient reads or no significant hits are labeled as none. E. Subclone analysis of GM12878 cell  
339 line. Three distinct clones with different VDJ elements are identified.

340

341 **Supplementary Fig. S14 Comparison of clustering on IGG-treated Hela cells by SCAN-seq2**

342 **and NGS methods.** UMAP embeddings of HeLa cells after IGG treatment at different  
343 concentrations and times. scRNA-seq was performed with both SCAN-seq2 (a) and NGS-based  
344 STRT method (b). Cells are colored by unsupervised clustering results (top left), IGG concentration  
345 (top right), time of treatment (bottom left) and cell cycle phase (bottom right).

346

347 **Supplementary Fig. S15 GO enrichment analysis in IGG-responsive cluster in 2 cell lines.** Dots  
348 are colored by fold of enrichment and sized by the number of relevant genes.

349

350 **Supplementary Fig. S16 Validation of DTU events between H9 and GM12878 cell lines.** DTU  
351 events in 4 genes (*ANAPC5*, *RANGRF*, *RPL22L1*, *ISOC2*) are shown (a-d). For each DTU event,  
352 two different cDNA products from transcripts of different lengths from a specific gene are amplified  
353 by one pair of primers. The expression ratios inferred by SCAN-seq2 are listed on the right.

354



355 **Acknowledgments**

356 This work was supported by the Beijing Advanced Innovation Center for Genomics at Peking  
357 University. We thank members in the Tang laboratory for discussions. Thank Dr. Ping Lu in  
358 construction of bioinformatics analysis process.

359

360 **Funding**

361 The work was supported by the National Key R&D Program of China (2021ZD0200102,  
362 2022YEF0203200 and 2018YFA0107601) and National Natural Science Foundation of China  
363 (32288102 and 31871457).

364 **Author contributions**

365 FT conceived the project. YL and YZ developed the protocol and contributed to cell culture, flow  
366 cytometry plus the cDNA sample preparation. ZL was in charge of the bioinformatic analysis. YL,  
367 YZ, ZL and FT wrote the manuscript with help from all authors.

368

369 **Ethics declarations**  
370 **Ethics approval and consent to participate**  
371 Not applicable.  
372 **Consent for publication**  
373 Not applicable.  
374 **Conflict of interest**  
375 The authors declare no competing interests.