

High-throughput single-nucleus hybrid sequencing reveals genome-transcriptome correlations in cancer

Siran Li^{1*}, Joan Alexander¹, Jude Kendall¹, Peter Andrews¹, Elizabeth Rose¹, Hope Orjuela¹, Sarah Park¹, Craig Podszus¹, Liam Shanley¹, Rong Ma², Nissim Ranade¹, Michael Ronemus¹, Arvind Rishi³, David L. Donoho², Gary L. Goldberg⁴, Dan Levy¹, Michael Wigler^{1*}

¹ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

² Department of Statistics, Stanford University, Stanford, CA, USA

³ Department of Pathology and Laboratory Medicine, Zucker School of Medicine at Hofstra/Northwell, Hempstead, NY, USA.

⁴ Department of Obstetrics and Gynecology, Division of Gynecologic Oncology, Zucker School of Medicine at Hofstra/Northwell, Northwell Health, New Hyde Park, NY, USA

* Correspondence: siranli@cshl.edu; wigler@cshl.edu

Abstract

Single-cell genomic analyses can provide information on cellular mutation and tumor heterogeneity, whereas single-cell transcriptomic analyses can distinguish cell types and states. However, the disconnect between genomic and transcriptomic spaces limits our understanding of cancer development. To address this, we developed a novel high-throughput method that simultaneously captures both DNA and RNA from single nuclei and new algorithms for the quantitative clustering and filtering of single-cell data. We applied this hybrid protocol to 65,499 single nuclei extracted from frozen biopsies of five different endometrial cancer patients and separately clustered the genome and expression data. We also analyzed 34,651 and 21,432 nuclei using RNA-only and DNA-only protocols, respectively, from the same samples to verify the clustering. Multiple tumor genome and/or expression clusters were often present within an individual patient, and different tumor clones could project into distinct or shared expression states. Almost all possible genome-transcriptome correlations were observed in the cohort. Stromal clusters were largely shared between patients, but some patients possessed unique stromal components, or mutant stroma with a significant loss of the X chromosome. This study reveals the complex landscape involving genome and transcriptome interactions at single-cell level, and provides new insights into mutant stroma as a potential clinical biomarker.

Main

To enhance our understanding and treatment of cancer, it is important to understand its heterogeneity, its interaction with the host, and the role the host plays in assisting or inhibiting the invasive somatic clone. Single-cell analysis offers one possible route to improved understanding. Integrating genomic and transcriptomic information would enable us to better explore the stromal reaction to neoplasm and the diversity of transcriptional states within tumors, making it possible to understand cooperation and competition between cancer cells and stroma, to view the emergence of malignant from pre-malignant cells, and to discern the forces that drive particular expression states. To achieve this aim and demonstrate its potential, we developed and applied a high-throughput single-cell

42 analysis of DNA and RNA to five cases of uterine cancer. We show the method yields new information,
43 and provides insights into the nature of the cancer stroma.

44 Methods for high throughput single-cell DNA¹⁻³ or RNA⁴⁻¹¹ alone are well established. For
45 investigating both omics, there are methods¹²⁻¹⁶ for inferring copy number states from RNA-seq data
46 based on the assumed positive correlation between genome copies and gene expression. While there
47 are low-throughput methods¹⁷⁻²⁸ for capturing both nucleic acids from single cells, among the limited
48 high-throughput whole-genome and whole-transcriptome techniques²⁹, none have been applied to
49 tumor biopsies to systematically study the cancer heterogeneity and stromal mutations with the
50 throughput demonstrated in this study.

51 In this paper, we introduced an innovative high-throughput technology in which both DNA and
52 RNA sequences can be captured from the same cell nucleus. Post-sequencing, the DNA and RNA layers
53 are bioinformatically separated by their mapping properties and subsequently analyzed. We chose to
54 work with nuclei because nuclei from frozen biopsies were a more available and abundant clinical
55 sample source than cells. This multi-omics method is a progression from the single-omic BAG platform³⁰,
56 in which single cells were encapsulated into individual balls of acrylamide gel, with either DNA or RNA
57 captured by Acrydite primers that were copolymerized into the gel matrix. We applied our hybrid DNA-
58 RNA approach to frozen tissue biopsies of five patients with endometrial cancer. We found sufficient
59 transcriptional complexity in the nuclear RNA to cluster by cell type and state, and sufficient copy
60 number information in the DNA layer to readily distinguish stroma from tumor components by genomic
61 DNA analyses. These clustering patterns were confirmed by published RNA-only and DNA-only
62 protocols. In addition, we developed a novel multinomial algorithm, the “multinomial wheel,” to
63 quantitatively measure the deviation of each single cell from the main clusters. This allowed us to
64 effectively remove cell collisions, and gave us insights into the minority of cells showcasing tumor
65 genomes and stromal expression patterns and those with normal genomes but tumor expression
66 profiles.

67 Although tumor expression clusters are highly distinct between patients, within a given patient,
68 we often observed multiple tumor expression states. We found that one tumor DNA clone may project
69 into one or more of these distinct expression states, which may or may not be shared by another tumor
70 clone from the same patient. In five patients, we observed virtually every possible projection pattern
71 between genomic and transcriptomic states. Conversely, stromal expression clusters are largely shared
72 between patients, although in different proportions. Two patients exhibited unique stromal
73 components, and all five displayed instances of mutant stroma. In the patient with the worst clinical
74 outcome, almost half of her plasma cells lost one copy of the X chromosome. These observations
75 demonstrated the potential and immediate applicability of this multiomic method in the investigation of
76 cancer evolution, stromal mutations, and the intricate interplay between genome and transcriptome.

77

78 Results

79 The Hybrid Platform

80 We study DNA and RNA templates from the same individual nuclei using BAG technology, a
81 flexible platform open to design modifications³⁰. The BAG platform extends the capabilities of traditional
82 droplet methods^{3,4} by polymerizing droplets containing nuclei in the presence of Acrydite-modified
83 primers. The Acrydite-modified primers form a primer-template duplex with the nucleic acid contents of
84 the cell or nucleus. Once polymerized, the “balls of acrylamide gel” or BAGs are removed from oil and
85 processed in aqueous solution. For DNA-only or RNA-only protocols, the primers that capture nucleic
86 acids are extended using polymerase or reverse transcriptase, resulting in template copies from the
87 same individual nuclei being covalently bound together in the same bag. After processing, we randomly
88 distribute BAGs into 96 wells, affixing one of 96 unique well-barcode sequences to each template on

89 each BAG in each well. Pooling the BAGs together and randomly splitting them again is defined as a
90 “pool-and-split” process. We repeat this process three times to generate a unique signature of 96^3
91 about one million unique **BAG barcodes**. During the first two pool-and-split labeling steps, in addition to
92 the BAG barcodes, we also add eight random bases to each template. When combined with the four
93 bases from the genomic sequences, these twelve random bases form a **varietal tag** (or UMI) that
94 uniquely labels each template molecule.

95 In this paper, we used BAGs to capture DNA and RNA from the same individual nucleus or cell
96 under mild denaturing conditions. To capture both types of nucleic acids, we used a mixture of Acrydite-
97 modified random T/G and oligo d(T) primers (**Fig. 1a**). The primers are then extended using reverse
98 transcriptase and DNA polymerase under conditions favorable to both reactions (see **Methods** and
99 **Supplementary Protocol**). We apply the pool-and-split labeling process with varietal tags as described
100 above, and then combine all the BAGs into a single PCR reaction, followed by tagmentation to generate
101 a sequencing library.

102 The targets of our investigation are frozen biopsy samples from five uterine cancer patients. The
103 samples include tumor tissue, which we denote as **Tumor 1** through **Tumor 5**. For three patients (1, 2
104 and 4), we also have samples from normal adjacent tissue, which we denote as **Normal 1**, **Normal 2**, and
105 **Normal 4**. In addition to our new **hybrid** protocol for simultaneously analyzing both DNA and RNA, we
106 also generated sequencing libraries using traditional BAG sequencing of DNA alone (labeled as **DNA-**
107 **only**) and RNA alone (labeled as **RNA-only**) on the same samples. The main differences between the
108 three protocols are illustrated in **Supplementary Fig. 1a**. Additionally, we applied the 10x Genomic
109 Chromium v3 single-cell RNA sequencing method on Tumor 1 for comparison purposes.

110 By design, the paired-end reads in the sequencing library are asymmetric: one end (Read 2) of
111 the read-pair contains information about the template sequence, while the other end (Read 1) contains
112 the BAG barcode, the varietal tag, and some template sequence. We first confirm that the reads have
113 the correct structure with BAG barcodes and NLAIII cutting sites at the correct positions (**Fig. 1a**), and
114 then extract the BAG barcode, varietal tag, and sequence information. We then use HISAT2³¹, a sensitive
115 alignment program useful for mapping both DNA and RNA reads, to map the template sequences. On
116 average, 84% of reads have a high-quality map with a mapping quality score equal to 1 or 60 and are
117 only primary mappings. Since sequence amplification occurs before tagmentation, we often obtain
118 reads from different fragmented copies of the same initial template. These reads will have the same
119 BAG barcode, varietal tag, and similar mapping regions, but they will have different sequence
120 information on the opposite end of the read-pair. For this reason, we group reads with the same BAG
121 barcode and varietal tag into **templates**, and templates with the same BAG barcode are grouped into
122 **BAGs**.

123 We only consider BAGs with a sufficient number of templates, using the “elbow bend” in the
124 cumulative distribution to establish the cutoff. For nuclei from frozen tissue biopsies, on average, we
125 observe 20 reads per template, 8200 templates per BAG, and 3500 BAGs per experiment when using the
126 hybrid protocol. The sequence and its context are used to discern whether templates are derived from
127 DNA or RNA. For each template, we document all its mapping information, including overlap with exons,
128 introns, UTRs, and behavior at splice sites. We established four categories:

- 129 1. **Exonic**: If greater than 90% of template bases are within a single gene transcript **and**
130 either >50% of the bases mapped to exons of that gene or the template includes a known splice
131 junction of the target gene.
- 132 2. **Intergenic**: If all template bases are intergenic.
- 133 3. **Intronic**: If less than 10% of template bases map to exons and not intergenic.
- 134 4. **Uncategorized**: does not satisfy 1-3.

135 We apply these categorization rules to all seven datasets for all protocols. The full distribution of each
136 category per sample is shown in **Supplementary Fig. 1b**. From the DNA-only protocol, we find that 2.5%
137 of templates are Exonic, 52% are Intergenic and 45% are Intronic, aligning closely with the proportions

138 anticipated from random sampling of the genome. In contrast, the great majority (>80%) RNA-only data
139 templates are mapped to genes, and a sizable proportion are mapped to exons (about 18%) and only
140 17% are Intergenic. Templates from the hybrid protocol were intermediate to the two: 10% Exonic
141 templates, 30% Intergenic templates, and 55% Intronic templates. The hybrid protocol is composed of a
142 mixture of DNA-only and RNA-only distributions. To divide the templates into RNA or DNA components,
143 we apply a conservative rule: we restrict RNA data analysis to Exonic templates (**RNA layer**) and we
144 restrict DNA data analysis to Intergenic templates (DNA layer) (see Methods for detailed filtering
145 process). The distribution of total unique templates, unique Exonic templates, and unique genes for
146 each sample in the cancer cohort—as well as for a human-mouse mixture experiment—using the hybrid
147 protocol, is presented in **Supplementary Fig. 2**.
148

149 Splitting layers comparison

150 Because our informatics relies on a conservative segregation of templates into molecular layers,
151 we expect some degradation of signal from excluded templates. To measure the extent of this loss, we
152 apply this same layer-splitting method to the hybrid protocol, as well as to the DNA-only and RNA-only
153 datasets. Before conducting separate layer analyses, we found it necessary to introduce additional
154 measures to control bias in the hybrid data, as we observed interference of RNA with DNA profiling. For
155 DNA profiling and clustering, we excluded certain intergenic genomic hotspots that many nuclear RNA
156 sequences map to, which significantly increased the quality of copy number data from the hybrid
157 protocol (**Supplementary Fig. 3**). These hotspots were present in data from both tumor and normal
158 tissues. The copy number data quality from the hybrid protocol was substantially superior to that
159 derived from the RNA-only protocol, with a quantitative measure of the signal-to-noise ratio being
160 discussed in the following section.
161

162 DNA-layer results

163 Restricting the DNA-only sequence data to the DNA layer results in a 54% reduction in template
164 counts. First, to estimate the effect on the signal-to-noise ratio for single cells between different
165 protocols, we use **Tumor 1 (Fig. 2a)** which has a 100 MB deletion on chromosome 5 that spans 13 bins
166 for all tumor cells despite different tumor clones. For each single cell, we compare the average template
167 counts of a normal copy 2 region (13 bins from chromosome 10) to the average template counts in the
168 copy 1 region in chromosome 5. The ratio of these averages estimates the signal strength between copy
169 numbers 2 and 1. We compute the mean normalized standard deviations over the bins in these regions
170 to estimate the relative noise. **Supplementary Fig. 4** shows the mean versus standard deviation (SD) for
171 single cells from the DNA layer for DNA-only (green), RNA-only (blue), and the hybrid method (orange).
172 We also show the results of using all the templates from the DNA-only experiments (red). As expected,
173 we have the least noise using all the DNA-only data (SD of 0.148). The signal-noise-ratios (mean/SD) for
174 all-molecules DNA-only data, DNA layers of DNA-only data, hybrid data, and RNA-only data are 13.0, 9.7,
175 4.3, and 1.0, respectively.

176 Second, we demonstrate that clustering copy number patterns using DNA-layer molecules or all
177 DNA molecules generates similar DNA clonal information. Using DNA-only data from Tumor 2 as an
178 example (**Supplementary Fig. 5**), we show that Seurat clustering generates the same number of clusters
179 (**Supplementary Fig. 5a,b**) and almost identical copy-number heatmaps (**Supplementary Fig. 5c,d**) using
180 either all molecules or only DNA-layer molecules as bin counts. Both clusterings in Supplementary Fig.
181 5a and 5b use the same nuclei. The heatmap in **Supplementary Fig. 5e** shows the number of cells in
182 each clone of the two clusterings, and we find that most of the cells belong to the same DNA clones in
183 both clusterings. To quantify this result, we examine all the nuclei pairs in both clusterings and check
184 whether each nuclei pair stays in the same or different DNA clones when clustered using all molecules or

185 only DNA-layer molecules. We present the result in **Supplementary Fig. 5f** and show that 97.6% of the
186 nuclei pairs are on the diagonal, indicating that the DNA layer preserves the DNA clonal information.
187 Third, we show the similarity between the DNA layer of the hybrid protocol and the DNA-only
188 protocol. We use tumor clones of Tumor2 as an example. Clusterings of both protocols generate the
189 same number of DNA clones and copy-number patterns (**Supplementary Fig. 6a,b**). To quantify this
190 result, we measure the proximity of every single tumor nucleus to the centroid of each tumor clone. We
191 compute the centroid of each tumor clone by averaging all the nuclei belonging to that clone as
192 determined by Seurat clustering. Given the centroid of each of the four tumor clones and 8 intermediate
193 linear combinations equally spaced between each pair of tumor clone centroids (for a total of 52
194 possible states), we calculated the most likely state for each nucleus based on multinomial distribution.
195 For each nucleus, the distance between this spot and its Seurat-assigned clone is called its “Distance
196 from home”. More details and validation of multinomial wheel analysis will be further discussed in the
197 last section. In **Supplementary Fig. 6c,d**, we plot the most likely state for each tumor nucleus from the
198 hybrid protocol and DNA-only protocol, respectively, with the colors marking its cluster identity. The
199 histograms on the right show the distribution of the distance from home for each nucleus. Compared to
200 the DNA-only protocol where 84.5% of the nuclei are within two units from the centroid of each clone,
201 this number drops to 77.2% for the hybrid protocol. Therefore, we estimate the resolution of the hybrid
202 protocol dropped by about 9% compared to the DNA-only protocol. Furthermore, we calculate the
203 distance of nuclei from the hybrid protocol to the tumor clone centroid determined by the DNA-only
204 protocol, and vice versa. As shown in **Supplementary Fig. 6e,f**, there is a 10.7% reduction of hybrid data
205 to DNA-only centroids compared to hybrid centroids, but there is no reduction of DNA-only data on
206 hybrid protocol centroids, showing the normalized averaged bin counts of each tumor clone for each
207 protocol are similar. Heatmaps of all five cases show clonal similarities between the hybrid protocol and
208 DNA-only protocol (**Supplementary Fig. 7.**)
209

210 **RNA-layer results**

211 Similar to the comparison between DNA layer and all DNA molecules, we use RNA-only data
212 from Tumor 2 to demonstrate that restricting the analysis to the RNA layer generates similar results to
213 using all of the RNA templates mapped within transcripts (**Supplementary Fig. 8**). UMAP clustering using
214 only the RNA layer templates (**Supplementary Fig. 8a**) or all RNA templates (**Supplementary Fig. 8b**)
215 generates the same number of clusters. Both clusterings were performed using the same group of
216 nuclei. The number of nuclei in each clustering is shown in **Supplementary Fig. 8c**, in which 94.3% of the
217 nuclei remain in the same expression clusters regardless of whether the RNA-layer or all-molecule
218 clustering was performed. We further quantify this result by examining all nuclei pairs to check whether
219 they reside in the same or different clusters under two clusterings **Supplementary Fig. 8d**. We found
220 that 95.2% of the nuclei pairs are on the diagonal, indicating that restricting the analysis to the RNA-
221 layer conserves much of the transcriptional clustering information.

222 Second, to demonstrate that the expression clustering results were similar between the hybrid
223 and RNA-only protocols, still using Tumor 2 as an example (**Supplementary Fig. 9**), we first clustered the
224 data from each protocol separately (**Supplementary Fig. 9c, 9g**). We then combined the data generated
225 from the two protocols and clustered the merged dataset together, adjusting for technical variations in
226 the protocols (see the **Methods** section). We show nuclei from each protocol in the merged clustering in
227 the panels of **Supplementary Fig. 9a, 9e**. We found that all clusters were populated by nuclei from both
228 protocols (**Supplementary Fig. 9d**). The merged-data clustering generated the same grouping of nuclei
229 as the clustering by either protocol alone. This was illustrated in heatmap matrices of identities
230 (**Supplementary Fig. 9b, 9f**). The columns of the heatmap show clusters from either the hybrid or the
231 RNA protocol alone, while the rows show clusters from the merged data. We found that most of the
232 nuclei fell on the diagonal of the heatmaps, meaning that nuclei of the same type were predominantly
233 grouped together whether clustered alone or in the merged data. For all five cases, the quality of

234 separations from the hybrid protocols (leftmost panel) and the RNA-only protocols (rightmost panel)
235 was very similar (**Supplementary Fig. 10**).

236 Finally, we compare the RNA layer of the hybrid protocol to the commonly used 10x Chromium
237 v3 single-cell RNA-seq method (**Supplementary Fig. 11**). We used Tumor 1 as an example since, in
238 addition to diverse stromal components, we observed two distinct tumor expression states (RNA1a and
239 RNA1b) that are independent of the two tumor DNA clones (discussed further in later sections and in
240 **Fig. 2a**). We find that the RNA clustering from 10x Chromium v3 also separates the tumor nuclei into
241 two expression states mainly based on collagen-related gene expressions, and with similar cell
242 proportions to the hybrid protocol (**Supplementary Fig. 11a,b**). To demonstrate the similarities between
243 these two protocols in distinguishing tumor cluster RNA1a from RNA1b, we study the correlation of the
244 fold change between RNA1a and RNA1b of all the genes for the two protocols (**Supplementary Fig. 11c**).
245 We restrict the analysis to genes that were detected in at least 10% of cells in either RNA1a or RNA1b
246 clusters for both protocols. Running an ordinary least squares (OLS) regression of y on x produces a
247 highly significant coefficient of 0.92 (**Supplementary Fig. 11c**). In addition, we plot the ratio of the
248 proportions of cells expressing these genes in RNA1a versus RNA1b for both protocols (**Supplementary**
249 **Fig. 11d**). Correlation tests under the null hypothesis $H_0: \rho=0$ and the alternative hypothesis $H_1: \rho>0$
250 confirm that the gene expressions and cell proportions in these two protocols have strong positive
251 correlations with both p -values less than $2.2e-16$. The top marker genes that either positively or
252 negatively distinguish RNA1a from RNA1b for both protocols, with p -values from the Wilcoxon Rank Sum
253 test, are listed in **Supplementary Fig. 11e**.

254

255

Doublets

256

257

258

259

260

261

262

263

264

265

266

We performed a human-mouse nuclei mixture experiment to study the levels of doublets and cross-contamination. Out of 1299 nuclei, after removing templates that mapped to both the mouse genome and human genome, we detected 619 nuclei with the majority (>85%) of the templates mapped to the human genome, and 662 nuclei with the majority (>85%) of the templates mapped to the mouse genome, and there were 18 doubles (1.39%) (**Fig. 1b**). We also show that the level of doublets is consistent between DNA layer and RNA layer, as shown in **Supplementary Fig. 12**, where only 5 out of 1299 nuclei (0.38%) nuclei showed disagreement of identities between two layers. The experiment also presented a low level of cross-contamination. As shown in **Fig. 1c**, the percentage of mouse templates in human nuclei had a median of 0.20%, and the human templates in mouse nuclei had a median percentage of 0.55%. This low level of cross-contamination also is preserved in both the DNA layer and RNA layer (**Supplementary Fig. 12**).

267

268

269

270

271

272

273

274

275

276

277

In addition, we performed two mixture experiments using nuclei from tumor biopsies. In each experiment, the frozen material from two patients was mixed prior to preparing single nuclei, sorting, and encapsulating them in BAGs. The source of each nucleus in the mixture experiments could be readily determined by its abundance of germline single nucleotide polymorphisms (SNPs), and collisions/doublets could be readily identified by having SNPs from both sources. One library was a mixture of nuclei from Patient 1 and Patient 5 (**Supplementary Fig. 13a**), and the other was from Patient 2 and Patient 5 (**Supplementary Fig. 13b**). If the SNPs in a BAG were contaminated by at least 15% of SNPs from the other patient, it was considered a collision. Based on this measure, 5.3% (out of 2526 nuclei) and 6.0% (out of 4244 nuclei) of BAGs were judged as collisions in these two experiments, respectively (**Supplementary Fig. 13c, 13d**).

278

Alluvial diagrams illustrating complex tumor projection patterns

279

280

281

Having validated the clustering patterns from the hybrid platform, we were able to confidently assess the projections of the genomic clusters onto the expression clusters. We first demonstrate this concept through a mixture experiment involving two cell lines: a normal male fibroblast, SKN1, and a

282 breast cancer cell line, SKBR3. The distributions of the numbers of total unique molecules, Exonic
283 molecules, and detected genes from this experiment are shown in **Fig. 1d-f**. We illustrate the clustering
284 results and heatmaps based on the copy number and gene expression in **Fig. 1g-k**. The alluvial diagram
285 (**Fig. 1i**) shows the projection of the genomic clones into the expression clusters. As expected, we
286 observed a good one-to-one correlation between the genome and transcriptome of each cell type.

287 After validating the alluvial diagram using cell lines, we used it to illustrate the genome-
288 transcriptome correlations in all five tumor samples (**Fig. 2**). In each panel of Fig. 2, we present the
289 clustering results for both the DNA layer and RNA layer, as well as the copy-number and gene-
290 expression heatmaps to illustrate the features that distinguish each DNA clone or RNA cluster in the five
291 cases. In general, the tumor genome clones projected distinctly from the normal genome clones,
292 although there were exceptions which we term “crossovers” and will discuss in more depth later.

293 Next, we focused mainly on the tumor genome projection patterns. To classify projections, we
294 used a set of letters and numbers to represent the tumor genome and tumor RNA clusters, respectively.
295 For example, {A:1,2; B:2} indicates tumor clone A projected into RNA clusters 1 and 2, whereas tumor
296 clone B from the same primary tumor tissue projected only into RNA cluster 2. Each of the five tumors
297 had a different projection pattern, and we observed almost all the possible patterns, which are defined
298 as follows: distinct tumor clones could each project into distinct expression clusters (e.g., {A:1; B:2} for
299 Tumor 5), into shared clusters (e.g., {A:1,2; B:1,2} for Tumor 1), or into a combination of distinct and
300 shared clusters (e.g., {A:1,2; B:1} for Tumor 4). Alternatively, multiple DNA clones could project into a
301 single RNA cluster (e.g., {A:1; B:1; C:2; D:2} for Tumor 2), or a single tumor clone could project into two
302 RNA clusters (e.g., {A:1,2} for Tumor 3). We discuss the special aspects of each case in the following
303 paragraphs.

304

305 Special aspects of each case

306 Tumor 1 was a uterine carcinosarcoma with greater than 90% of cells in the tumor tissue having
307 copy number variations (**Fig. 2a**). The copy number heatmap showed two tumor DNA clones, with the
308 primary difference being that clone T1B had a lower copy number in a region on chromosome 13. The
309 alluvial plot showed that cells from each tumor DNA clone projected about equally into both tumor RNA
310 clusters. Cluster RNA1a had a high expression level of fibroblast-specific genes such as fibroblast growth
311 factor receptor genes (*FGFR3*) and collagen genes (*COL9A2*) (**Supplementary Fig. 14a**), consistent with
312 the pathological classification of this tumor as having a sarcomatous component. These fibroblast genes
313 had lower expression in cluster RNA1b. On the other hand, cluster RNA1b had higher expression of
314 *RSPO4*, a key regulator of the Wnt/ β -catenin signaling pathway, and *DUSP6*, a negative regulator of the
315 ERK signaling pathway (**Supplementary Fig. 14a**). When we projected the two tumor RNA clusters back
316 into DNA UMAP space, we found that the nuclei from both RNA clusters were randomly distributed in
317 DNA space and unrelated to DNA cluster patterns (**Supplementary Fig. 15a-c**). The two tumor RNA
318 clusters were also not separated by cell cycle, template, or gene counts (**Supplementary Fig. 15d**). The
319 normal nuclei from this patient's primary tumor made up a minor fraction of the total.

320 Tumor 2 was diagnosed as a uterine serous carcinoma, from which we observed four DNA tumor
321 clones and two tumor RNA clusters (**Fig. 2b**). The two RNA clusters to which these tumor clones
322 projected had many distinguishing gene sets (**Fig. 2b**). Tumor DNA clones T2A and T2B each projected to
323 RNA2a, while T2C and T2D each projected to RNA2b. The major feature shared by T2A and T2B, but not
324 T2C and T2D, was the loss of an entire X chromosome. In fact, one gene that significantly distinguished
325 RNA2a and RNA2b was *XIST*, which had low expression in RNA2a. Compared to tumor 1, the normal
326 genomes of tumor 2 projected to many more distinct RNA clusters. In this example, we observed two
327 normal DNA clusters, one with a single copy of the X chromosome named as DNA clone “Nx”. From the
328 alluvial plot in **Fig. 2b**, we found that most of the nuclei in this Nx DNA clone projected to plasma cells as
329 well as T cells. As can be seen in the zoomed image from **Supplementary Fig. 14b**, the projection of “Nx”

330 into the "Plasma cells" cluster also had a low expression of *XIST*. We quantify this observation by
331 comparing the *XIST* RNA-layer counts (Exonic templates) between nuclei with two copies and one copy
332 of the X chromosome. Compared to the clone with one copy of the X chromosome, where 95% of the
333 nuclei have 0 *XIST* RNA counts, there is a significantly higher RNA count (p-value < 1.7e-22, t-test) for the
334 other group of plasma cells where the median *XIST* RNA count is 3 (**Supplementary Fig. 14b**). As *XIST*
335 RNA is required for X-chromosome inactivation and is only expressed from the inactive X-
336 chromosome³², this result also verifies that nuclei from the "Nx" DNA clone lost their inactivated X
337 chromosomes instead of not being captured by the hybrid protocol.

338 Tumor 3 was an endometrial adenocarcinoma, and was the only example in which we did not
339 observe any discernible tumor subclones. Despite this, the tumor projected into two distinct RNA
340 clusters (**Fig. 2c**), which differed most notably in the elevated expression of the estrogen receptor gene
341 *ESR1* in RNA3a (**Supplementary Fig. 14c**). As assessed by immunostaining, about 50% of the tumor cells
342 from the primary tumor expressed the estrogen receptor (**Supplementary Fig. 14c**), which aligned with
343 what we saw from the RNA expression.

344 Tumor 4 was another uterine carcinosarcoma case. The copy number heatmap showed that the
345 nuclei with copy-number variations were clustered into two DNA clones with significant differences in
346 chromosomes 1, 8, and the X chromosome (**Fig. 2d**). However, for RNA clustering, there was only one
347 major RNA tumor cluster (RNA4a) containing nuclei from both DNA tumor clones. In addition, there
348 were two small RNA clusters, RNA4b and RNA4c, both close to the main cluster RNA4a in RNA space.
349 RNA4c had projections from both tumor DNA clones and had a high expression level of G2-phase marker
350 genes, *MKI67* and *CENPF*. Nuclei in cluster RNA4b were mostly from DNA clone T4A (**Fig. 2d**). Compared
351 to the main tumor RNA cluster RNA4a, RNA4b alone had a high expression level of many actin (*ACTG1*,
352 *ACTB*) and tubulin genes (*TUBA1B*, *TUBA1A*, *TUBB*). RNA4b also highly expressed *EEF2*, an essential
353 factor for protein synthesis, and *GAPDH*, a key enzyme in glycolysis, as shown in **Supplementary Fig.**
354 **14d**.

355 Tumor 5 was a uterine leiomyosarcoma. The biopsy sample from this patient had two sectors,
356 with one being more hemorrhagic than the other. Clustering from each sector was done separately, and
357 overall, the projections were similar. In **Fig. 2e**, we present only the results from Tumor 5-2 (with more
358 nuclei). Tumor 5 had two DNA tumor clones with significant differences in chromosomes 1, 2, 6, 7, 8, 12,
359 and the X chromosome. In RNA clustering, there were two large and one small tumor RNA clusters. The
360 two distinct tumor DNA clones projected mainly to distinct tumor RNA clusters (**Fig. 2e**). Both of the two
361 tumor RNA clusters, RNA5a and RNA5b, had high expression of fibroblast markers, which was
362 concordant with the immunohistochemical analyses showing that the tumor cells were positive for h-
363 caldesmon³³ (**Supplementary Fig. 14e**). RNA5a had higher expression of *TNNT3*, *PLXDC1*, and *MTMR11*,
364 whereas RNA5b had higher expression of *ADAM12*, (involved in skeletal muscle regeneration), *ZFHX4*
365 (related to muscle differentiation), *FN1* (which encodes fibronectin) and collagen genes such as *COL1A1*
366 and *COL6A2*. A small proportion of nuclei from both DNA clones went to the tumor RNA cluster RNA5c,
367 which was cell-cycle related. RNA5c specifically had high expression of typical G2-phase markers *MKI67*,
368 *TOP2A*, *CENPF*, and *CENPE* (**Supplementary Fig. 14e**).

369

370 Common stromal and distinct tumor clusters

371 To obtain a clearer picture of the data, we clustered single-nuclei hybrid data from all patients,
372 separately for RNA and DNA, using the "FindClusters" and "RunUMAP" functions of Seurat (**Fig. 3a** and
373 **Fig. 3b**, respectively). We examined a total of 35,369 nuclei from the sources indicated (**Fig. 3c**). For RNA
374 clustering, we downsampled all nuclei to 400 "Exonic" templates. In addition, and only for the RNA
375 clustering, we created a "garbage" cluster by including 3,500 nuclei from the DNA-only platform treated
376 as hybrid libraries. As expected, the 3,500 nuclei from DNA-only libraries clustered together and were
377 well separated from all other clusters (**Supplementary Fig. 16**). This "garbage" cluster also included 853
378 low-quality nuclei from the hybrid protocol (**Supplementary Fig. 16**). We removed these nuclei from

379 further analysis. The aggregated "Exonic" template counts for all the nuclei in each expression cluster
380 are shown in **Supplementary Table 1**.

381 Tumor-genome clusters were quite distinct from the normal-genome clusters (**Fig. 3b**), and
382 distinct clones within a patient mapped nearby to each other or merged into a single cluster at this
383 resolution of clustering. The intra-tumor clones of patient 2 (T2A, T2B, T2C, and T2D) and patient 5 (T5A
384 and T5B) were still preserved under this resolution. Cluster "N" was distinct from "Nx," the latter being
385 otherwise normal cells with only one copy of the X chromosome.

386 The projections of nuclei from six sample sources into the combined RNA space (**Fig. 3a**) are
387 shown in **Fig. 3d-i**, where nuclei from a given sample are highlighted either in blue or red, depending on
388 whether they are classified by DNA as tumor or normal genomes. In each panel, the nuclei from other
389 samples are colored in light grey. The projections of the tumor genomes are very distinct between
390 patients, well-separated from each other and the projections of the normal genomes. At this resolution,
391 the intra-tumor RNA expression sub-clusters generally merged together.

392 By contrast, the normal-genome cell projections from a given patient were quite distinct, and
393 different patients have overlapping normal-genome projections. We labeled these common elements by
394 their distinctive patterns of expression³⁴ and list the marker genes in **Supplementary Table 2**. The blood
395 components can be further distinguished into finer subtypes, as shown in the zoomed figure of **Fig. 3a**.
396 The counts for these projections of DNA profiles into RNA profiles are shown in **Fig. 4a**. In addition, the
397 hybrid protocol had good consistency between experimental replicates (**Supplementary Fig. 17**).

398 Exceptions to overlapping normal projections were seen for patient 1, where there was a cluster
399 mainly consisting of normal-genome cells, well-separated from the main stromal clusters. The epithelial-
400 like cluster "EP-T1" from the tumor tissue of patient 1 and the "EP-N1" cluster from an adjacent normal
401 site were distinct from each other or the main epithelial cluster "EP". We believe this distinct "EP-T1"
402 cluster was not due to batch effects, as it had very distinct and plausible gene expression patterns
403 (**Supplementary Fig. 18**), and other stromal-cell projections from this sample mostly overlapped well
404 with other samples.

405

406 Loss of the X chromosome in blood elements

407 The loss of a single X chromosome in some cancer cells (such as in patients 2 and 4)³⁵, as well as
408 in a small proportion of certain stromal components in cancer patients or the elderly³⁶, did not surprise
409 us. However, we unexpectedly saw that nearly half of the relatively abundant plasma cells showed
410 losses of the X chromosome in the patient (patient 2) with the worst clinical outcomes (**Supplementary**
411 **Table 3**). The loss of the X chromosome in somatic lineages was observed in all five endometrial cancer
412 cases. The summary of our data on the projection from normal genomes with and without two copies of
413 the X chromosome is shown in **Table 1**. Another patient (patient 4) with poor clinical outcomes had
414 about 15% loss of the X chromosome in the T cell components, further indicating that this might be a
415 potential biomarker for negative outcomes.

416

417 Crossovers and the multinomial wheel

418 We summarize the projection data for all nuclei from the hybrid protocol (**Fig. 4a**). As
419 represented there and even more visually in the alluvial plots (**Fig. 2**), we saw what we termed
420 "crossovers": nuclei that clustered as tumor or normal genomes, but then clustered in opposition as
421 normal or tumor expression patterns. 357 tumor genomes projected to normal expression clusters, and
422 401 normal genomes projected to tumor expression clusters. If correct, these crossovers could have
423 profound biological significance (**Discussion**), but it is well-known that in single-cell methods, two cells
424 or nuclei may 'collide' and create a merged profile. These collisions could occur at the droplet
425 generation stage, during informatic processing (e.g., barcode collisions), or even biologically by cell

426 fusion in the host. We developed a filter to eliminate such collisions by making a quantitative measure
427 of the deviation from a cluster.

428 To this end, we developed the "multinomial wheel." The idea behind this is to create a
429 "multinomial state" from each Seurat cluster. There are as many multinomial states as Seurat clusters,
430 one set of multinomial states for genomes and another set for expression. Each multinomial state is
431 represented as a K -length vector summing to one, where $K = 300$ bins for DNA data and $K = 29,637$
432 genes for RNA data. The value of a multinomial vector is the normalized centroid for the cluster. We
433 then computed the deviation of every nucleus from its Seurat-assigned multinomial state as follows.
434 Between every two multinomial states, we created nine new states that were equally spaced linear
435 combinations of the two multinomial vectors. Thus, for N "home" states, we created $9 \times N \times (N-1)/2$ new
436 multinomial states, making up the wheel (see **Fig. 4a** for the genome states, and **Supplementary Fig. 19**
437 for expression states). Every home state has $L = 9 \times (N-1)$ other states linked directly to it, each from one
438 to ten 'units' away. For each nucleus, we asked which of the $L+1$ multinomial states (including its home
439 state) would be the most likely to generate its observed template counts. The distance of that nucleus
440 from its Seurat-assigned multinomial state was the number of units to its closest multinomial state. This
441 was defined as the "distance from home" shown in the histograms in **Fig. 4b** and **Supplementary Fig. 19**.

442 We display each nucleus as a point on the wheel at its closest genome state (**Fig. 4b**, second
443 panel from left). If the distance of the nucleus exceeds 5, the point is colored in red; otherwise, it is
444 colored blue. We show a histogram of the distances (**Fig. 4b**, leftmost panel) and saw that most nuclei
445 were within one-unit distance away from the Seurat-assigned genome multinomial state.

446 To test the utility of the multinomial wheel to detect collisions, we utilized the two mixture
447 experiments using nuclei from frozen tumor biopsies where the collisions were verified by patient-
448 specific SNPs. In **Fig. 4b** (second panel from right), we show only the nuclei determined as collisions from
449 the first mixture experiment between patients 1 and 5. The BAGs judged by SNPs to be collisions usually
450 reside in the middle of two multinomial states, not close to either one of them (**Fig. 4b**, second panel
451 from right). We showed a histogram of their distances to Seurat-assigned clusters with a peak distance
452 of 4 (**Fig. 4b**, rightmost panel). Most existing doublets detection methods³⁷⁻⁴³ for single-cell RNA
453 sequencing start from the individual count vectors of single cells, and then make artificial doublets by
454 adding/averaging random droplet pairs and use these to train the model. However, unlike the existing
455 methods, the multinomial wheel acknowledges the major clusters determined by Seurat UMAP
456 clustering and then measures the deviation of every cell from the centroids of these major clusters.
457 Different from "DoubletDecon"⁴⁴ method which is also based on clustering information but decides if a
458 droplet resembles artificial droplets based on a deconvolution algorithm⁴⁵, we assume the multinomial
459 of a cluster could be viewed as a k -sided die, with each throw of the die landing on a face with a fixed
460 probability, with probabilities summing to 1. In our context, each face was either a gene (if an
461 expression multinomial) or a genomic bin (if a genomic multinomial), with its probability determined by
462 its relative frequency in the cluster of cells. Each cell in the cluster could be considered as the outcome
463 of N rolls of its multinomial, where N was the count of templates that were observed. This method
464 works for both DNA and RNA space (**Fig. 4b** and **Supplementary Fig. 19**), not only serving as a doublet
465 detector but also providing quantitative measurements of the cells in between different states.

466 These experiments justified using the multinomial wheel as a filter and removing the genome
467 'violators' from all BAGs. The violators were nuclei with a distance of ≥ 2 from their Seurat-assigned
468 home multinomial states. We plotted all violators from the combined DNA analysis in the respective
469 DNA and RNA multinomial wheels (**Supplementary Fig. 19b**). Importantly, the violators of the genome
470 wheel from the mixing experiment were also violators of the expression wheel (**Supplementary Fig. 19c-**
471 **d**). This filtration was stringent, as we removed about 20% of the BAGs, in excess of our expectation of
472 5% collisions. We then re-tabulated the projections of nuclei that passed the filter (**Fig. 4c**). In the table
473 of **Fig. 4c**, we highlighted the remaining crossovers in red. Filtration reduced crossovers from 357 to 24
474 for tumor genomes with normal expression, and from 401 to 24 for normal genomes with tumor

475 expression. Thus, removing 20% of the BAGs by filtration eliminated greater than 90% of all crossovers.
476 The alluvial plots for each case before and after violator removal are shown in **Supplementary Fig. 20**.

477 One case of note in **Fig. 4c** was that two nuclei in the adjacent normal tissue of patient 2
478 (Normal2) showed tumor genomes, and their respective copy-number profiles are exhibited in
479 **Supplementary Fig. 21b**. Using multinomial wheel analyses, we found both nuclei were within 1 unit
480 distance from the same tumor clone T2B, and in RNA space, both were within 1 unit distance from
481 tumor RNA cluster RNA2a. These two nuclei were not crossovers, but they represented tumor subclone
482 infiltration into the adjacent normal tissue for patient 2.

483 Patients 2 and 5 showed the largest number of residual crossovers, and they also had the largest
484 numbers prior to filtering. We checked all residual crossovers to make sure the Seurat "FindClusters"
485 assignments agreed with UMAP spatial assignments, and that the individual copy number profiles had
486 good quality. We plotted the crossovers from Tumor 5 with normal genomes and tumor expression on
487 the DNA and RNA multinomial wheel in **Supplementary Fig. 21c**, and the crossovers with tumor
488 genomes and normal expressions from Tumor 2 in **Supplementary Fig. 21d**, which gave a clearer
489 picture. We discuss the biological implication of these remaining crossovers in the next section.
490

491 Discussion

492 We sought to develop a high-throughput method for the assessment of both RNA and DNA from
493 individual cells of a population, and to begin to explore its utility in the description of the cellular
494 composition of primary cancer sites. Our experimental design incorporated four elements. First, we
495 chose BAG single-cell technology because of its flexibility and excellent performance for either RNA-only
496 or DNA-only protocols³⁰. Second, we collected both RNA and DNA from individual units at the same time
497 in a hybrid protocol because this was simpler than trying to capture the two nucleic acid types
498 sequentially. Third, we chose nuclei over cells because nuclear RNA was sufficient for the classification
499 of cell types⁴⁶⁻⁴⁸, and isolating intact cells from frozen or fixed biopsy samples is problematic. Finally, we
500 chose to examine nuclei from one target organ, in this case the uterus, so that we could better assess
501 the commonality of the stroma and the diversity of tumor expression⁴⁹. Although this hybrid method, in
502 its current format, has lower genome coverage per cell than the single-nucleus DNA-only or RNA-only
503 protocol³⁰, it still shows many advantages over plate-based, low-throughput methods that assess both
504 DNA and RNA together (**Supplementary Table 4**). For example, it enables the detection of mutant
505 stroma existing only in a small proportion of most cell types, which is challenging for low-throughput
506 methods. This new method not only enables the analysis of many more cells, but also overcomes some
507 of the challenges imposed by existing techniques, such as high labor intensity, limitation of sample
508 types, and a preference for entire cells over nuclei alone^{19,20,50}.

509 At low resolution, we found that the expression clusters of the cancer cells themselves were
510 quite distinct, well separated from each patient and from the normal clusters. At higher resolution, each
511 cancer had more than one expression cluster. The relationship between these tumor expression clusters
512 and the DNA subclonal populations of the cancers was not consistent from patient to patient. In some
513 cases, cells of distinct tumor subpopulations projected to distinct expression clusters; in some cases, the
514 cells of distinct subpopulations projected to the same expression cluster; and in some cases, cells from
515 the same subpopulation "split": they projected to distinct expression clusters. We believe that these
516 split expression patterns are consistent with epigenetic drift rather than being caused by genetic
517 variations. We take note of two special cases: in patient 3, the same cancer population projected to
518 estrogen receptor positive and receptor negative expression clusters; and in patient 1 with uterine
519 carcinosarcoma, each of two cancer subclones projected to high and low collagen expression clusters. In
520 both patients 1 and 3, the distinct cancer expression types were physically interspersed as determined
521 by histopathology. This was consistent with the idea of epigenetic variation, rather than genetic
522 variation, because we would expect the latter to show physical segregation.

523 Using this multiomic technology, we observed that a significant number of stroma cells lost one
524 copy of chromosome X. Especially in patient 2, almost half of the plasma cells showed loss of one copy of
525 chrX, suggesting extensive clonality in this lineage. These observations raise additional questions: does
526 somatic clonality indicate failure of the immune checkpoint mechanisms? Do these cells hinder or help
527 the tumor penetrate the host? Do the somatic elements travel with the cancer when it metastasizes?

528 Although aneuploid and diploid lineages generally projected to distinct expression clusters, we
529 initially observed many exceptions that we termed "crossovers." While it would not be surprising to see
530 early tumor lineages without copy number changes begin to express the tumor pattern, it would be
531 surprising to see the entire program expressed so early, before the selection within the host for the
532 predominant tumor clone. Also, if tumor cells can take on the expression pattern of normal cells, they
533 could possibly escape host surveillance or chemotherapy. Such crossovers could therefore be of
534 immense interest, provided they are not artifacts. We therefore refined our methods to minimize
535 possible artifacts. The most likely artifacts are from 'collisions,' BAGs that either report multiple nuclei or
536 with coincidental identifiers. Preparing intentionally mixed samples that were distinguishable by SNVs
537 enabled us to determine that collisions occurred in about 5% of the BAG data. After removing possible
538 collisions, the few remaining crossovers merit further future study with larger data sets. Some of these
539 remaining crossovers might be from an earlier cancer lineage or mutated stroma⁵¹⁻⁵⁴.

540 To better understand "crossovers," we needed a tool for quantifying the similarity of a cell to
541 others in its cluster. The widely used clustering program in Seurat was effective at finding clusters,
542 offering a manifest of marker genes that distinguished the clusters, and providing clear graphical
543 displays. However, the current clustering method failed to detect cells intermediate between clusters
544 and was highly dependent on parameter settings. We, therefore, experimented with a simpler
545 mathematical paradigm for clustering, the multinomial distribution and developed the "multinomial
546 wheel" method to filter the collisions. This multinomial wheel algorithm has extensive utility beyond this
547 specific case or method. It can assist in any clustering analyses to provide quantitative measurement of
548 how each single cell fits into each cluster, which would help identify outliers, collisions, and cells in
549 transition from one state to another in either genomic or transcriptomic space.

550 In summary, we have developed a high-throughput multiomic method that connects genotypes
551 and expression profiles at single-cell resolution. When tumors have copy number changes, it is now
552 possible to distinguish stromal expression patterns from tumor expression patterns. In exploring five
553 uterine tumors, we uncovered all possible patterns of connection between tumor subclones and
554 expression sub-clusters. We saw differences in the proportionate composition by stromal type, and
555 observed clear evidence of genomic variants in stromal subtypes. How these observations relate to
556 cancer biology in general, or to the classification of cancer subtypes and their relation to disease
557 outcome, await more extensive studies. We expect this pilot study opens a window to the complex
558 relationship between genome and transcriptome, and will lead to new insights into cancer biology, new
559 methods for monitoring cancer progression and evaluating clinical prospects, and possibly new
560 treatments.

561

562 **Methods**

563 **Pulverization of frozen biopsy samples in liquid Nitrogen**

564 All patient tissue biopsy samples were pulverized in Liquid Nitrogen (LN2) with a sterile mortar
565 and pestle prior to analysis. Mortar and pestles were submerged in LN2 and cooled to LN2 temperature.
566 The cooled vessels were then partially filled with fresh LN2 and transferred to a basin containing a
567 shallow pool of LN2. The presence of LN2 in both the mortar and basin helped maintain a constant
568 temperature during the pulverization process and prevent sample heating due to friction. The tissue
569 samples were then transferred to the sterile mortar, submerged in LN2, and pulverized until they were
570 mostly a fine, homogeneous powder. Once pulverized, residual tissue material was scraped off the

571 pestle back into the mortar with a sterile, LN2-cooled disposable spatula. The mortar was then removed
572 from the basin to allow for the LN2 to evaporate out of the mortar. Subsequently, pulverized tissue was
573 immediately collected with a fresh, sterile, LN2-cooled disposable spatula into 2.0 mL DNA LoBind
574 Eppendorf tubes submerged in LN2. Pulverized samples were placed on dry ice with the caps open to
575 allow for temperature equilibration before closing the tubes, and then stored at -80°C until further use.
576 All samples were pulverized with separate sterile mortar and pestles to avoid cross-contamination
577 between tumor and normal adjacent biopsy tissues.
578

579 The sample cohort

580 We studied samples from five patients (patient 1 – patient 5). The samples were from biopsies
581 of their uterine cancers (Tumor 1 – Tumor 5), and in three patients also from adjacent normal
582 endometrial sites (Normal 1, Normal 2, and Normal 4). For most samples (Normal 1, Tumor 1, Tumor 2,
583 Tumor 3, Normal 4, Tumor 4, Tumor 5), we sequenced single nuclei of the same sample on each of three
584 platforms: DNA-only, RNA-only, and the hybrid protocol. We performed comparison analyses and
585 showed the validity of the hybrid protocol mainly using the above five trio data sets.
586

587 Hybrid BAG generation

588 We dissolved the pulverized tissue in ice-cold NST detergent buffer⁴⁷ and stained with DAPI. We
589 performed single-nuclei sorting using DAPI-H vs. DAPI-A single-nuclei gate on a FACSAria II SORP cell
590 sorter to remove debris and clumps. We confirmed (data not shown) that single-nuclei sorting based on
591 ploidy would not be able to distinguish cancer cells from normal cells because the hypodiploid peak of
592 cancer cells often overlaps with the diploid peak of normal cells⁴⁷. Single nuclei were loaded into the
593 microfluidic device described in detail in a previous publication³⁰. Nuclei were encapsulated into
594 droplets with an average diameter of 120 microns. For the capture of nucleic acids, we used 5' Acrydite
595 oligonucleotides. All the Acrydite-modified oligonucleotides became covalently co-polymerized into the
596 gel ball matrix. They also all contained, at their 5' end, a universal PCR primer (UP1) for subsequent
597 amplification. For RNA-only protocol, we used oligo-dT; for DNA-only protocol, we used random T/G
598 primers, and followed their respective published protocols³⁰. To capture both RNA and DNA together in
599 the new hybrid protocol, we used both Acrydite primer designs, but we altered the protocol in two
600 important ways.

601 The first critical change was an incubation step at 85°C for 5 minutes instead of 95°C for 12
602 minutes for DNA denaturation in the DNA-only protocol. Otherwise, we observed significant destruction
603 of the RNA.

604 The second critical change took place after the BAGs were formed. The RNA and genomic DNA
605 trapped in the BAGs were used as templates to make covalently bound copies, and in the new hybrid
606 protocol, both reverse transcriptase and DNA polymerase were used. Template-switch-oligos were also
607 introduced in the hybrid protocol so that the cDNA products which were covalently linked to the BAG
608 matrix ended with a double-stranded region. This double-stranded DNA region included an NLA-III
609 cleavage site. Subsequently, DNA polymerase (Klenow) was added to extend the captured genomic DNA
610 from primers, forming a copy that was also covalently linked to the BAGs. Some, perhaps most, of the
611 cDNA-mRNA sequence was further partially converted to double-stranded cDNA. BAGs were pooled and
612 the covalently captured DNA and cDNA were cleaved with NLAIII leaving a sticky end used for
613 subsequent extensions.

614 BAG barcodes and varietal tags were added to the 3' ends of the covalently captured nucleic
615 acids in split-and-pool reactions. The BAG barcodes were present on both the genomic-DNA and RNA
616 copies. The varietal tags were used for counting. The first BAG barcode and varietal tag were added by
617 ligation extension (described in detail in the supplementary experimental protocol), leaving a common
618 3' sequence identical across all the molecules and BAGs. The second BAG barcode and varietal tag were

619 added by hybridization extension of the common 3' sequence, along with a second common sequence
620 adapter for the third split-and-pool step. The third barcode was added by a split PCR, using the first
621 universal PCR primer (UP1) and the second common sequence adapter as part of the PCR primer
622 sequences.

623 These amplified products were pooled and converted by tagmentation into paired-end Illumina
624 sequencing libraries. One end of the reads contained BAG barcode and varietal tag, as well as genomic
625 or transcriptomic sequence information. The other end from random tagmentation was mostly genomic
626 or transcriptomic sequence information.

628 Initial data processing

629 Sequencing libraries were sequenced in paired-end 150 bp format using an Illumina NovaSeq
630 6000. Briefly, each processing step is described in more detail in the immediately following sections. We
631 first checked the structure of each read pair in the fastq files. For the good read pairs with the correct
632 structure as shown in **Fig. 1a**, we extracted the BAG barcode, varietal tag, and genomic sequences from
633 both reads. We then mapped the genomic (including transcriptomic) sequence to the reference genome
634 with gene transcript information. Finally, we combined the mapping information from all reads
635 belonging to each varietal tag for each BAG barcode. In the end, we obtained a template data table with
636 each row containing the information of an original template/molecule. In the following section, we
637 explain each processing step from the fastq file to the template table in detail.

639 *Step 1 – Check sequence structure*

640 First, we filtered out reads from the fastq files where either Read 1 or Read 2 were less than 100
641 bases. Second, we examined if the sequences from the expected BAG barcode positions exactly matched
642 one of the 96×96×96 barcodes, and if the "CATG" cutting site was in the expected location, allowing for
643 one base mismatch. We removed read pairs that did not satisfy these requirements. Third, from Read 1
644 which started with barcodes and varietal tags, we trimmed away the first 80 bases containing the BAG
645 barcode, varietal tag, and adapter sequences, and also checked if the reverse complementary sequence
646 of the universal primer ("CCAAACACACCCAA") or oligo-dT ("AAAAAAAAAAAAAAAA") was present. If
647 present, it meant we had reached the end of the template, so these primer-related sequences were
648 trimmed off for downstream mapping. Similarly, for Read 2, the tagmentation end, we checked and
649 removed the adapter sequence ("GAGCGGACTCTGCG") from the first split-and-pool if it existed. After
650 trimming, we required both Read 1 and Read 2 to be at least 30 bases long. All the bases from Read 1
651 and Read 2 after trimming were then used for paired-end mapping (Step 3).

653 *Step 2 – Extract BAG barcode and varietal tags*

654 If a read pair passed Step 1, we extracted the BAG barcode and varietal tag information from the
655 first 80 bases of Read 1, and this information was appended to the read ID. The 17 base BAG barcodes
656 came from three cycles of the split-and-pool procedure, of which five bases came from the 1st-split, six
657 bases came from the 2nd-split, and six bases came from the 3rd-split. There were 96 different barcodes
658 for each split, so there were altogether 96×96×96 (≈ 1 million) varieties. The 12 base varietal tag came
659 from both the split-and-pool primers and the genomic sequence. Out of these twelve bases, four bases
660 came from the 1st-split, four bases came from the 2nd-split, and four bases came from the genomic
661 sequence that was two bases away from the "CATG" cutting site. These twelve bases provide 4¹² (≈ 16
662 million) varieties for each BAG.

664 *Step 3 – Map to the human genome*

665 After steps 1 and 2 above read pairs were mapped to the UCSC hg19 human genome using
666 HISAT2 version 2.1.0³¹. The reference genome we used included the primary chromosomes and

667 unlocalized and unplaced contigs. Alternate haplotypes were not included in the genome index. HISAT2
668 can take a file with known splice sites to use for alignment. This file was generated using a gtf formatted
669 file extracted from the NCBI refSeq gene annotation table from the UCSC genome browser and the
670 HISAT2 program, `hisat2_extract_splice_sites.py`. The bam files were then sorted and indexes using
671 samtools. In subsequent data analysis steps we designate by mapped reads the reads that HISAT2 marks
672 as being part of a proper pair and a primary mapping having a read mapping quality score greater than
673 zero.

674
675 *Step 4 – Combine read information with original template information*

676 We grouped the mapped reads based on their BAG barcodes. For the reads with the same BAG
677 barcode, we sorted the varietal tags by the number of reads associated with each tag in descending
678 order. We performed a "rollup" algorithm on the sorted varietal tags, and discarded varietal tags within
679 a Hamming distance of one from a more abundant varietal tag having at least ten times more reads. We
680 assumed the eliminated varietal tags originated from the tags with more abundant reads but contained
681 sequencing or PCR errors.

682 Using the varietal tags from the above "rollup" step, we aggregated the mapped segments for all
683 the reads with the same varietal tag. We checked the total coverage of each varietal tag against all
684 exons and transcript boundaries from the NCBI refSeq gene annotation file downloaded from the UCSC
685 genome browser, and wrote out one line per varietal tag with all the useful information into a "template
686 table". Each line of the template table contains the following information: BAG barcode, varietal tag,
687 chromosome, start mapping position, end mapping position, start and end mapping position for each
688 fragment if there was more than one continuous fragments, total bases covered by this template,
689 number of reads, number of genes, gene list, bases overlapping with the transcript of the best-matched
690 genes, bases overlapping with exons, number of splice junctions, number of unspliced sites, bases
691 overlapping with the coding regions, 5'UTR, and 3'UTR of the gene. The downstream data analyses were
692 mainly based on the information from this table. The best-mapped gene was deemed to be the gene
693 from the annotated transcript file having the highest overlap to the transcript. If more than one
694 transcript had the same overlap then best was determined by overlap to exons, then overlap to coding
695 sequence, then the number of splice junctions, then the fewest unspliced sites. If more than one gene
696 tied for all these criteria, then all genes are listed in the template table.

697
698 **Template processing**

699 *Sequence classification*

700 Starting from the template data table described above, each initial molecule was classified as
701 one of the four categories: "Exonic", "Intronic", "Intergenic", and "Uncategorized". This process was
702 applied uniformly regardless of the protocol types (RNA-only, DNA-only, or hybrid). We classified a
703 template as an "Exonic" template if over 90% of its bases were mapped within one gene. Furthermore,
704 we refined the "Exonic" classification only if 50% or more covered bases from this template were exonic,
705 or if 20% or more covered bases were exonic and at least one splicing event was observed (RNA layer). If
706 all the bases from a template were mapped to intergenic regions, we classified it as "Intergenic". If a
707 template was not classified as "Intergenic", but less than 10% of its covered bases were exonic and no
708 splicing events were observed, this template was classified as "Intronic". Only a small proportion of
709 templates failed to be classified into the above three categories, and these templates were classified as
710 "Uncategorized".

711 For expression clustering, we only used "Exonic" templates assigned to a single gene regardless
712 of protocols. For copy number clustering, we tested four versions of template choices on all the libraries,
713 and presented the comparative results on the two normal tissue samples (Normal 1 and Normal 4) in
714 **Supplementary Fig. 3**, which we will discuss in the next section.

715

716

Copy number plot varieties

717

We demonstrated four progressively improved versions of copy number estimation, named "all_molecules", "no_exon", "no_gene", and "no_gene.avoid50closeTN". The method "all_molecules" simply used all molecules for each retained nucleus for copy number as the name would imply. The method "no_exon" used molecules both classified as "Intronic" and "Intergenic" in the previous paragraph. The method "no_gene" only used "Intergenic" templates with no bases covering a transcript.

722

The method "no_gene.avoid50closeTN" (DNA layer) only retained the "Intergenic" molecules from the "no_gene" method that were at least 50 bases distant from RNA hotspots. We defined an RNA hotspot as the genomic region between two "Intergenic" templates that were within 50 bases of each other in RNA-only libraries from all tumor and normal samples in the cohort. RNA hotspots were expected to be some combination of actual unannotated transcripts and regions of DNA that were prone to being copied by reverse transcriptase. As these hotspot sequences distorted copy number profiles in normal and tumor biopsy specimens, we eliminated certain intergenic regions when determining copy number profiles for the hybrid protocol.

730

731

Empirical bin boundary generation for copy number

732

Separately for each of the four copy number molecule selection methods above, we used the genomic positions of all molecules from normal DNA samples to determine empirical bin boundaries for 300 bins with approximately equal molecule counts per bin. Excluding any molecules mapping to chromosome Y, we assigned to each chromosome 1-22 and X a number of bins in proportion to its fraction of total molecule counts. Within each chromosome, bin boundaries were assigned greedily from the start of the chromosome so that all but the final bin contained at least the same number of molecules that was equal to the total counts of molecules (or referred to as templates) divided by the number of bins for that chromosome. The observed count of molecules per bin was recorded as a normalization factor for later use during per-sample copy number estimation. This normalization factor could vary by up to 30% between chromosomes because a small number of bins (300) can only be imperfectly allocated by chromosomal molecule counts.

743

744

Copy number estimation

745

For each copy number variant separately, each selected molecule incremented a bin based on the established bin boundaries for that method. Each bin count was then divided by the per-bin normalization factor, and the result was multiplied by 2 divided by the median value over all bins. Assuming a mostly diploid sample, this process resulted in a copy number profile for the sample that was centered at a value of 2. Circular binary segmentation (DNACopy version 1.50.1)⁵⁵ in R was then performed on the copy number profile using parameters $\alpha=0.02$, $nperm=1000$, $undo.SD=0.5$ and $min.width=2$. For each profile, we also computed a quantity we call 'terrain' which was the sum of the absolute value of adjacent bin copy number differences. To produce copy number input for the "CreateSeuratObject" function of Seurat, the per-bin normalization factor was applied to each raw bin count for each cell, and a second per-cell normalization factor was then used so that each cell's total normalized count was set equal to its total unnormalized count.

756

757

RNA clustering

758

The RNA clustering was performed using Seurat package (version 3.1.5), and using the standard Seurat clustering pipeline⁵⁶. The gene names were also appended with the chromosome information to distinguish any ambiguous locations. We removed the ribosomal protein genes for clustering. For comparing expression clustering between the hybrid protocol and RNA-alone protocol, we normalized the gene-template matrix by cell, and removed the PCA components that most significantly

759

760

761

762

763 distinguished protocol differences. We normally used at least 15 PCA components for clustering. This
764 approach gave us similar clustering results as the "IntegrateData" function in Seurat v4. For the
765 combined RNA clustering of all the hybrid data, we downsampled the gene matrix to 400 Exonic
766 templates per nucleus, and included nuclei with more than 300 Exonic templates for clustering. In the
767 clustering process, we only used genes that showed up in at least 30 nuclei, and nuclei with at least 150
768 genes; we used the top 5,000 variable gene features for PCA analysis and used the first 50 PCA
769 components for subsequent UMAP and FindCluster functions.

770

771 Copy number clustering

772 Similar to RNA clustering, we used "RunUMAP" and "FindClusters" functions of Seurat to cluster
773 nuclei based on copy number. For each library, we had a bin-counts matrix, similar to the gene matrix
774 for RNA clustering. There were 300 rows in the matrix, representing 300 genomic bins. Each column
775 represented a nucleus. Each element of the 2D matrix represented the tag counts of the corresponding
776 bin in the corresponding nucleus. We first normalized the matrix by columns: for each nucleus, we
777 divided each bin count by the mean of 300 bins and then multiplied by 2. We not only used these 300
778 normalized single bin counts for clustering; additionally, we also included the median normalized bin
779 counts of every two and three adjacent bins, as long as these adjacent bins were within the same
780 chromosome. The reason for this step was that copy number segmentation usually requires similar
781 amplification or deletion patterns in at least two contiguous bins. By doing this, we appended another
782 277 rows from the two adjacent bins and 254 rows from the three adjacent bins onto the original 300-
783 row normalized bin-count matrix.

784 We performed clustering using the new matrix with 831 rows. We used a workflow similar to
785 that for RNA clustering, but we did not use "NormalizedData" function since the matrix had already
786 been normalized. For "FindVariableFeatures" function, we used the top 500 features by inputting
787 "selection.method = "vst", nfeatures = 500".

788

789 Copy number heatmap

790 The single-nucleus copy-number heatmap was plotted using Seurat "DoHeatmap" function. Each
791 row represented the median normalized counts of two adjacent bins, except for the first bin of each
792 chromosome, in which we used the normalized count of that single bin. The total of 300 rows were
793 sorted in genomic order, with chromosome Y eliminated.

794

795 Multinomial Wheel

796 To build a multinomial wheel in DNA space, we first computed a multinomial vector to represent
797 each Seurat cluster. Each multinomial vector had 300 elements, representing 300 genomic bins. Each
798 element was the total bin counts from all the nuclei in that cluster. We normalized each vector to sum
799 to one, serving as the multinomial probability vector representing that cluster. Next, we computed the
800 linear combination of multinomial probability vectors of every two Seurat clusters, and created 9 equally
801 spaced sampling states $C_{1,2,\dots,9} = pA + (1 - p)B$, for $p = (0.1, 0.2, \dots, 0.9)$, where A and B are the two
802 original states. We then assigned the nucleus to the state with the highest likelihood. In R language, we
803 used the "dmultinom" function to compute multinomial probabilities.

804 We applied a similar idea to create the RNA multinomial wheel. Different from the DNA
805 multinomial vector where each element was a genomic bin, in RNA space, each element represented
806 one of the 29,637 genes. We computed the sum of gene counts for each Seurat cluster $V_{1,2,\dots,n}$ (n is the
807 number of Seurat clusters, and V_i is a 29,637-element vector, $i = 1, 2, \dots, n$), but unlike DNA, there were
808 many elements still being zero which could not be used as a multinomial probability vector. We solved
809 the problem by adding a small value to each element that was proportional to the total expression level
810 of every gene, so that each vector V_i^* does not contain zero elements. For each gene element j, we did

811 the following transformation: $V_i^*[j] = V_i[j] + (0.05 \times (\sum_j V_i[j])) \times (\sum_i V_i[j]) \div (\sum_{i,j} V_i[j])$. We then
812 normalized each vector V_i^* to obtain the multinomial probability vector for cluster i.
813

814 **Supplementary Materials**

815 Supplementary Fig. 1 to Supplementary Fig. 21
816 Supplementary Table 1 to Supplementary Table 4
817 Supplementary Protocol
818

819 **Acknowledgments**

820 We thank P. Moody for cell sorting, Q. Gao for histology assistance, A. Runnels and E. Ghiban for
821 Illumina sequencing assistance, C. Danyko, A. Stepanyk, and D. Stauder for technical assistance, D.
822 Tuveson for helpful discussion, A. Kapedani, M. A. Green, and S. Chin from the Clinical Research Team in
823 the Department of Obstetrics and Gynecology of Long Island Jewish Medical Center for clinical
824 information assistance. Additionally, we extend our gratitude to S. S. Fox, M. Heywood, K. Quinn, B. M.
825 Weil, J. Jacob from Biobanking/Anatomic Pathology Team in the Northwell Health Biospecimen
826 Repository (NHBR) in Northwell Health Cancer Institute for sample transfer and pathological information
827 assistance. This work was supported by a grant from the Simons Foundation, Life Sciences Founders
828 Directed Giving-Research (award number 519054 to M.W.); The Breast Cancer Research Foundation
829 (BCRF, to M.W.); and by support from the Cold Spring Harbor Laboratory and Northwell Health
830 Affiliation (awarded to M.W.).
831

832 **Author contributions**

833 S.L. and M.W. conceived the idea and designed the study; S.L., J.K., P.A., and M.W. developed
834 bioinformatic analysis programs; S.L. and M.W. developed the experimental protocol; S.L., J.A., E.R, H.O.,
835 S.P., and L.S performed single-cell BAG-seq experiments; J.A. and C.P. assisted in sample transfer and
836 sample preparation; A.R. performed histology analyses; G.L.G provided samples and clinical guidance;
837 S.L., J.K., P.A., R.M., N.R., M.R., D.L.D., D.L., and M.W. performed informatics analyses and results
838 interpretation; S.L., D.L., and M.W. wrote the manuscript with input from all coauthors.
839

840 **Declaration of interests**

841 The authors declare that they have no competing interests.
842

843 **Data availability**

844 Illumina sequencing data for all the single-nucleus libraries are available at NCBI Sequencing
845 Read Archive (SRA) with accession code (PRJNA773107).
846

847 **Code availability**

848 Code is available through https://github.com/siranli01/DNA_RNA.
849

850 **References**

- 851
- 852 1. Lan, F., Demaree, B., Ahmed, N. & Abate, A.R. Single-cell genome sequencing at ultra-high-
853 throughput with microfluidic droplet barcoding. *Nature biotechnology* **35**, 640-646 (2017).
 - 854 2. Vitak, S.A. *et al.* Sequencing thousands of single-cell genomes with combinatorial indexing. *Nature*
855 *methods* **14**, 302-308 (2017).

- 856 3. Andor, N. *et al.* Joint single cell DNA-Seq and RNA-Seq of cancer reveals subclonal signatures of
857 genomic instability and gene expression. *Biorxiv*, 445932 (2018).
- 858 4. Macosko, E.Z. *et al.* Highly parallel genome-wide expression profiling of individual cells using
859 nanoliter droplets. *Cell* **161**, 1202-1214 (2015).
- 860 5. Klein, A.M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells.
861 *Cell* **161**, 1187-1201 (2015).
- 862 6. Gierahn, T.M. *et al.* Seq-Well: portable, low-cost RNA sequencing of single cells at high throughput.
863 *Nature methods* **14**, 395-398 (2017).
- 864 7. Cao, J. *et al.* Comprehensive single-cell transcriptional profiling of a multicellular organism.
865 *Science* **357**, 661-667 (2017).
- 866 8. Rosenberg, A.B. *et al.* Single-cell profiling of the developing mouse brain and spinal cord with split-
867 pool barcoding. *Science* **360**, 176-182 (2018).
- 868 9. Habib, N. *et al.* Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons.
869 *Science* **353**, 925-928 (2016).
- 870 10. Gao, R. *et al.* Nanogrid single-nucleus RNA sequencing reveals phenotypic diversity in breast
871 cancer. *Nature communications* **8**, 1-12 (2017).
- 872 11. Zheng, G.X. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nature*
873 *communications* **8**, 1-12 (2017).
- 874 12. Gao, R. *et al.* Delineating copy number and clonal substructure in human tumors from single-cell
875 transcriptomes. *Nature biotechnology* **39**, 599-608 (2021).
- 876 13. Fan, J. *et al.* Linking transcriptional and genetic tumor heterogeneity through allele analysis of
877 single-cell RNA-seq data. *Genome research* **28**, 1217-1227 (2018).
- 878 14. Elyanow, R., Zeira, R., Land, M. & Raphael, B.J. STARCH: Copy number and clone inference from
879 spatial transcriptomics data. *Physical Biology* **18**, 035001 (2021).
- 880 15. Patel, A.P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma.
881 *Science* **344**, 1396-1401 (2014).
- 882 16. Erickson, A. *et al.* Spatially resolved clonal copy number alterations in benign and malignant tissue.
883 *Nature* **608**, 360-367 (2022).
- 884 17. Dey, S.S., Kester, L., Spanjaard, B., Bienko, M. & Van Oudenaarden, A. Integrated genome and
885 transcriptome sequencing of the same cell. *Nature biotechnology* **33**, 285-289 (2015).
- 886 18. Macaulay, I.C. *et al.* G&T-seq: parallel sequencing of single-cell genomes and transcriptomes.
887 *Nature methods* **12**, 519-522 (2015).
- 888 19. Han, K.Y. *et al.* SIDR: simultaneous isolation and parallel sequencing of genomic DNA and total
889 RNA from single cells. *Genome research* **28**, 75-87 (2018).
- 890 20. Hou, Y. *et al.* Single-cell triple omics sequencing reveals genetic, epigenetic, and transcriptomic
891 heterogeneity in hepatocellular carcinomas. *Cell research* **26**, 304-319 (2016).
- 892 21. Bian, S. *et al.* Single-cell multiomics sequencing and analyses of human colorectal cancer. *Science*
893 **362**, 1060-1063 (2018).
- 894 22. Han, L. *et al.* Co-detection and sequencing of genes and transcripts from the same single cells
895 facilitated by a microfluidics platform. *Scientific reports* **4**, 1-9 (2014).
- 896 23. Van Strijp, D. *et al.* Complete sequence-based pathway analysis by differential on-chip DNA and
897 RNA extraction from a single cell. *Scientific reports* **7**, 1-9 (2017).
- 898 24. Kong, S.L. *et al.* Concurrent single-cell RNA and targeted DNA sequencing on an automated
899 platform for comeasurement of genomic and transcriptomic signatures. *Clinical chemistry* **65**,
900 272-281 (2019).
- 901 25. Cheow, L.F. *et al.* Single-cell multimodal profiling reveals cellular epigenetic heterogeneity. *Nature*
902 *methods* **13**, 833-836 (2016).
- 903 26. Rodriguez-Meira, A. *et al.* Unravelling intratumoral heterogeneity through high-sensitivity single-
904 cell mutational analysis and parallel RNA sequencing. *Molecular cell* **73**, 1292-1305. e8 (2019).

- 905 27. Li, W., Calder, R.B., Mar, J.C. & Vijg, J. Single-cell transcriptogenomics reveals transcriptional
906 exclusion of ENU-mutated alleles. *Mutation Research/Fundamental and Molecular Mechanisms*
907 *of Mutagenesis* **772**, 55-62 (2015).
- 908 28. Yu, L. *et al.* scONE-seq: A single-cell multi-omics method enables simultaneous dissection of
909 phenotype and genotype heterogeneity from frozen tumors. *Science Advances* **9**, eabp8901
910 (2023).
- 911 29. Yin, Y. *et al.* High-throughput single-cell sequencing with linear amplification. *Molecular cell* **76**,
912 676-690. e10 (2019).
- 913 30. Li, S. *et al.* Copolymerization of single-cell nucleic acids into balls of acrylamide gel. *Genome*
914 *research* **30**, 49-61 (2020).
- 915 31. Kim, D., Paggi, J.M., Park, C., Bennett, C. & Salzberg, S.L. Graph-based genome alignment and
916 genotyping with HISAT2 and HISAT-genotype. *Nature biotechnology* **37**, 907-915 (2019).
- 917 32. Panning, B., Dausman, J. & Jaenisch, R. X chromosome inactivation is mediated by Xist RNA
918 stabilization. *Cell* **90**, 907-916 (1997).
- 919 33. Watanabe, K., Kusakabe, T., Hoshi, N., Saito, A. & Suzuki, T. h-Caldesmon in leiomyosarcoma and
920 tumors with smooth muscle cell-like differentiation: its specific expression in the smooth muscle
921 cell tumor. *Human pathology* **30**, 392-396 (1999).
- 922 34. Cao, J. *et al.* A human cell atlas of fetal gene expression. *Science* **370**, eaba7721 (2020).
- 923 35. Spatz, A., Borg, C. & Feunteun, J. X-chromosome genetics and human cancer. *Nature Reviews*
924 *Cancer* **4**, 617-629 (2004).
- 925 36. Zhou, Y. *et al.* Single-cell multiomics sequencing reveals prevalent genomic alterations in tumor
926 stromal cells of human colorectal cancer. *Cancer cell* **38**, 818-828. e5 (2020).
- 927 37. Xi, N.M. & Li, J.J. Benchmarking computational doublet-detection methods for single-cell RNA
928 sequencing data. *Cell systems* **12**, 176-194. e6 (2021).
- 929 38. Wolock, S.L., Lopez, R. & Klein, A.M. Scrublet: computational identification of cell doublets in
930 single-cell transcriptomic data. *Cell systems* **8**, 281-291. e9 (2019).
- 931 39. Lun, A.T., McCarthy, D.J. & Marioni, J.C. A step-by-step workflow for low-level analysis of single-
932 cell RNA-seq data with Bioconductor. *F1000Research* **5**(2016).
- 933 40. Bais, A.S. & Kostka, D. scds: computational annotation of doublets in single-cell RNA sequencing
934 data. *Bioinformatics* **36**, 1150-1158 (2020).
- 935 41. Gayoso, A., Shor, J., Carr, A.J., Sharma, R. & Pe'er, D. DoubletDetection (Version v2. 4). *Zenodo*,
936 *DOI* **10**(2018).
- 937 42. McGinnis, C.S., Murrow, L.M. & Gartner, Z.J. DoubletFinder: doublet detection in single-cell RNA
938 sequencing data using artificial nearest neighbors. *Cell systems* **8**, 329-337. e4 (2019).
- 939 43. Bernstein, N.J. *et al.* Solo: doublet identification in single-cell RNA-Seq via semi-supervised deep
940 learning. *Cell systems* **11**, 95-101. e5 (2020).
- 941 44. DePasquale, E.A. *et al.* DoubletDecon: deconvoluting doublets from single-cell RNA-sequencing
942 data. *Cell reports* **29**, 1718-1727. e8 (2019).
- 943 45. Gong, T. & Szustakowski, J.D. DeconRNASeq: a statistical framework for deconvolution of
944 heterogeneous tissue samples based on mRNA-Seq data. *Bioinformatics* **29**, 1083-1085 (2013).
- 945 46. Slyper, M. *et al.* A single-cell and single-nucleus RNA-Seq toolbox for fresh and frozen human
946 tumors. *Nature medicine* **26**, 792-802 (2020).
- 947 47. Navin, N. *et al.* Tumour evolution inferred by single-cell sequencing. *Nature* **472**, 90 (2011).
- 948 48. Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods.
949 *Nature biotechnology* **38**, 737-746 (2020).
- 950 49. Puram, S.V. *et al.* Single-cell transcriptomic analysis of primary and metastatic tumor ecosystems
951 in head and neck cancer. *Cell* **171**, 1611-1624. e24 (2017).

- 952 50. Zachariadis, V., Cheng, H., Andrews, N. & Enge, M. A highly scalable method for joint whole-
953 genome sequencing and gene-expression profiling of single cells. *Molecular Cell* **80**, 541-553. e5
954 (2020).
- 955 51. Navin, N. *et al.* Inferring tumor progression from genomic heterogeneity. *Genome research* **20**,
956 68-80 (2010).
- 957 52. Yokoyama, A. *et al.* Age-related remodelling of oesophageal epithelia by mutated cancer drivers.
958 *Nature* **565**, 312-317 (2019).
- 959 53. Yizhak, K. *et al.* RNA sequence analysis reveals macroscopic somatic clonal expansion across
960 normal tissues. *Science* **364**, eaaw0726 (2019).
- 961 54. Lee-Six, H. *et al.* The landscape of somatic mutation in normal colorectal epithelial cells. *Nature*
962 **574**, 532-537 (2019).
- 963 55. Olshen, A.B., Venkatraman, E., Lucito, R. & Wigler, M. Circular binary segmentation for the
964 analysis of array-based DNA copy number data. *Biostatistics* **5**, 557-572 (2004).
- 965 56. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888-1902. e21 (2019).
- 966

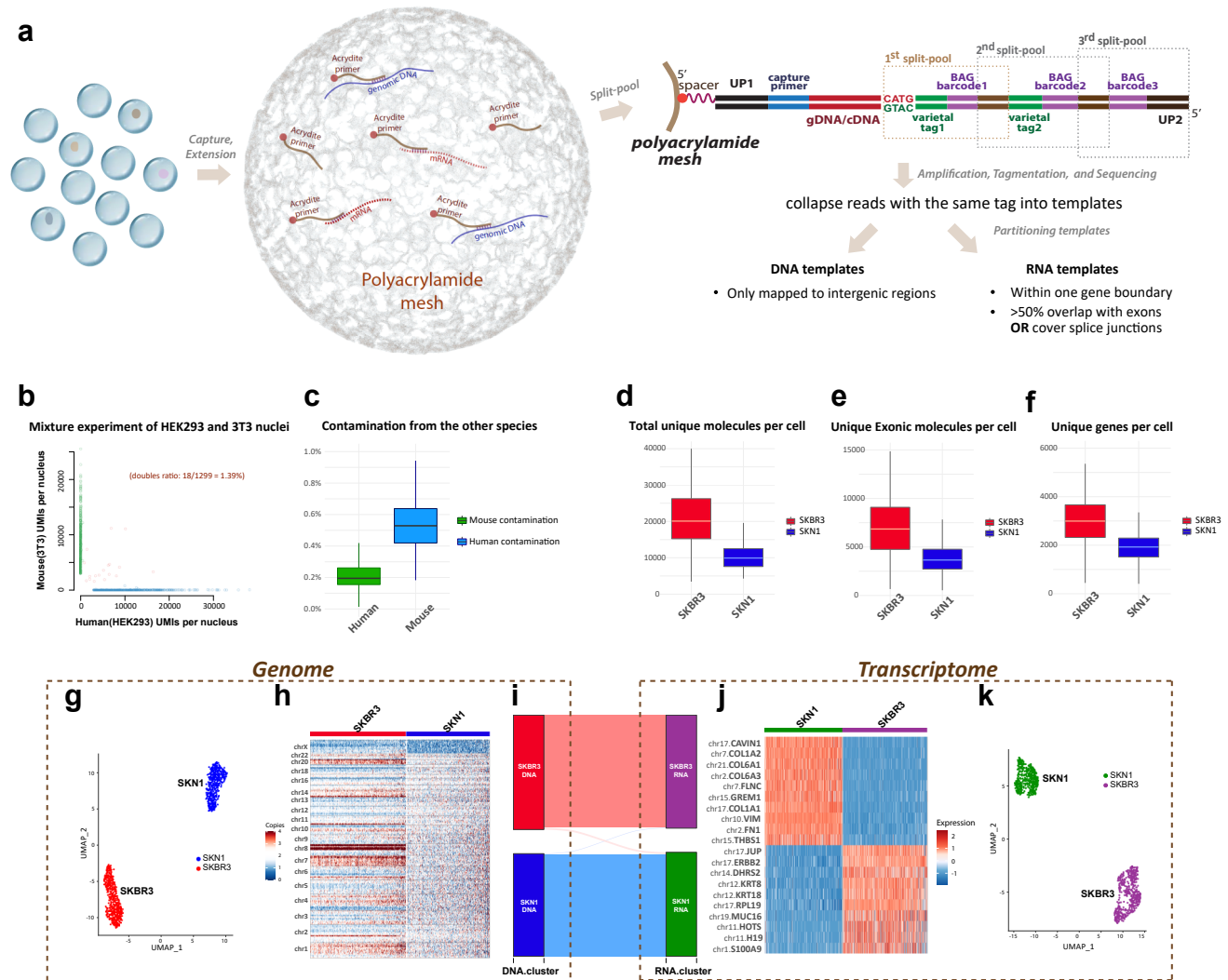


Fig. 1. Overview and basic performance of the single-nucleus/cell hybrid sequencing protocol to simultaneously capture and analyze the genome and transcriptome. **a**, Workflow showing the major steps of the single-nucleus/cell hybrid protocol. **b-c**, Doublets ratio (1.39% out of 1299 nuclei) (**b**) and cross-contamination level (0.2% mouse templates in human nuclei and 0.5% human templates in mouse nuclei) (**c**) from the other species of a HEK293-3T3 nuclei mixture experiment. **d-k**, Performance and genome-transcriptome correlation from a SKN1-SBKR3 mixture single-cell hybrid sequencing experiment; total unique molecules (**d**), exonic templates (**e**), and gene counts per cell (**f**); clustering based on DNA copy number (**g**) and copy number heatmap (**h**); clustering based on gene-count matrix (**k**) and heatmap of marker genes (**j**), and the correlation between two genomic clusters and two expression clusters (**i**).

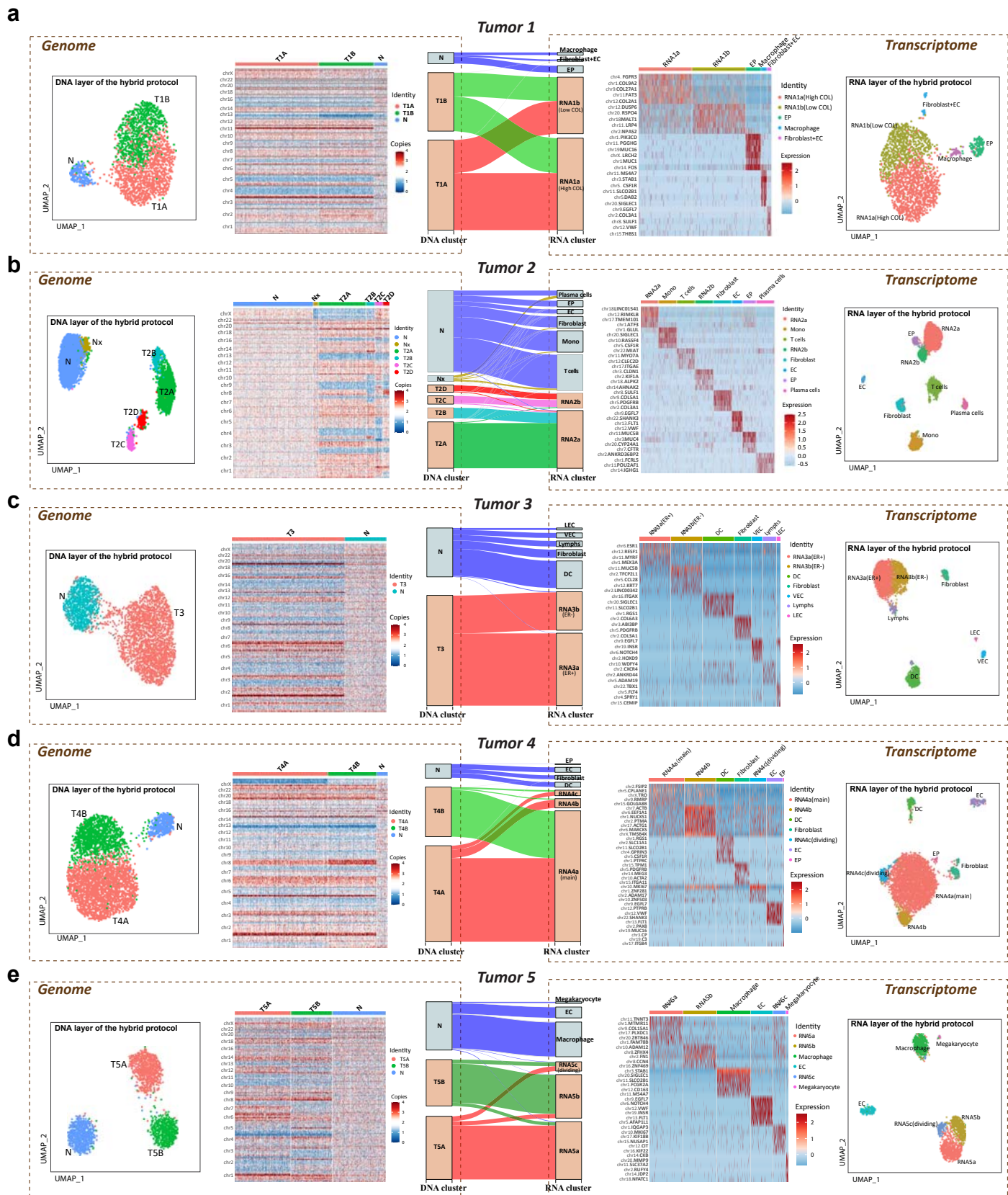


Fig. 2. Alluvial diagrams showing the genome-transcriptome correlations of the five tumor samples.

a-e, For each one of the five tumor samples Tumor1-5, we show the genomic clustering (leftmost panel) and copy number heatmap (second panel from the left), expression clustering (rightmost panel) and marker gene heatmap (second panel from the right), and the alluvial plots connecting the genome clones with RNA expression clusters (middle panel).

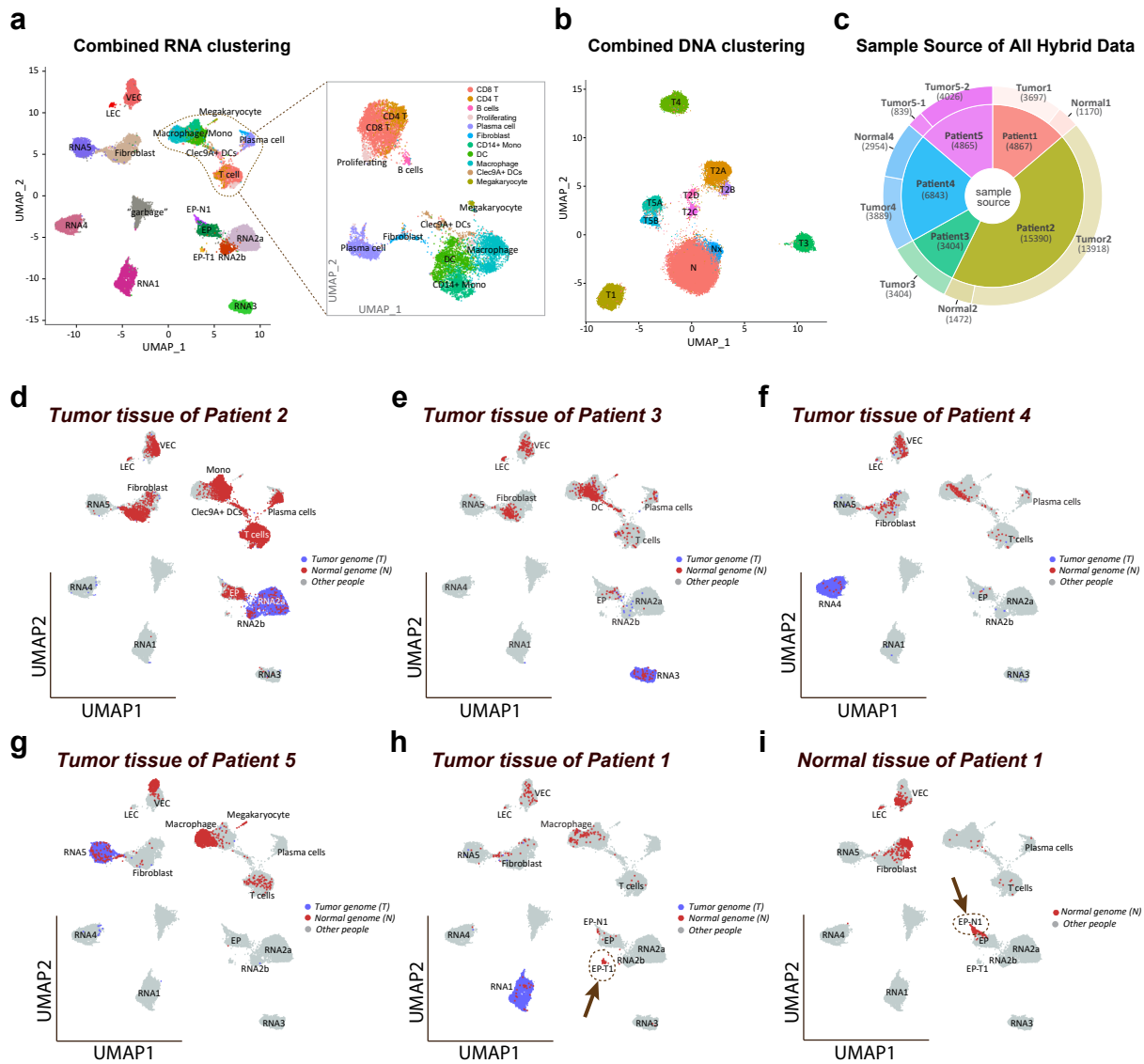


Fig. 3. RNA and DNA clustering using a combination of all nuclei from the hybrid protocol, distinguishing largely common stromal and distinct tumor clusters.

a, Combined RNA clustering using all nuclei from the hybrid protocols and 3,500 nuclei from seven DNA-only libraries. A zoomed plot in the same figure, circled by dash line, presents the subtypes of white blood cells when clustered separately. **b**, Combined DNA clustering after removing nuclei clustered to the “garbage” RNA state in (a). **c**, The sample sources of nuclei from the hybrid protocol in the combined analyses. **d-i**, The projections of tumor-genome (blue) and normal-genome (red) nuclei into the RNA UMAP space for six biopsy samples. The tumor-genome or normal-genome information is determined by the combined DNA analysis in (b); **d**, Tumor2; **e**, Tumor3; **f**, Tumor4; **g**, Tumor5; **h**, Tumor1; **i**, Normal1; in **h-i**, Unique stromal components are circled in dashed lines and indicated by arrows.

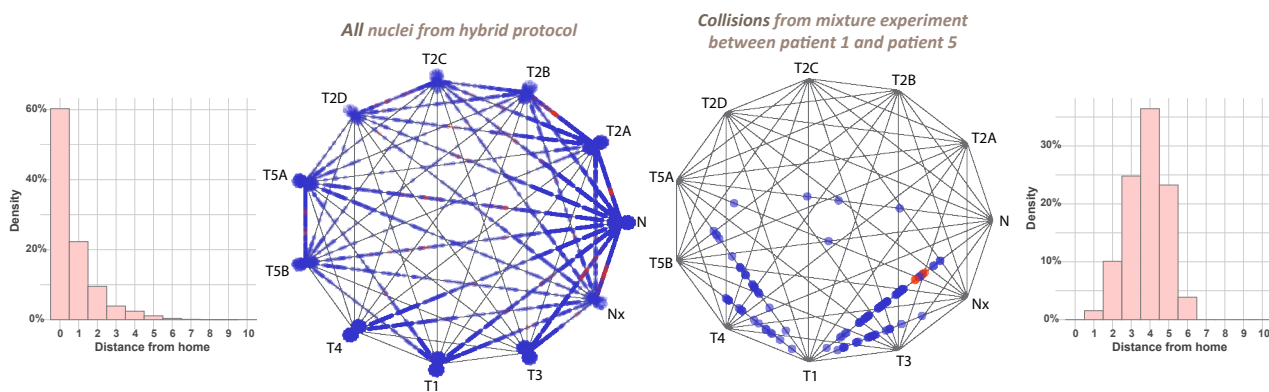
a

	RNA1	RNA2 (RNA2a+RNA2b)	RNA3	RNA4	RNA5	Macrophage/Mono	T cells	Clec9A+ DCs	Plasma cells	Fibroblast	EC (VEC+LEC)	EP	EP-N1	EP-T1	Megakaryocyte
Normal1_N	1			1		8	7	1	1	560	240	13	306		
Normal1_T															
Tumor1_N	25	3	3	2	8	75	6	1		23	46	8	2	118	
Tumor1_T	3358			1	3	1				2				3	
Normal2_N	1	2	2	1		32	38	4	2	625	173	572			
Normal2_T		3				2				9	1	2			
Tumor2_N	2	163	10	3	15	1929	3474	143	638	1072	575	606	1	2	2
Tumor2_T	1	5017	7	10	1	40	81	4	19	38	16	24	1		
Tumor3_N		3	41		4	483	42	46	19	151	130	16			7
Tumor3_T		8	2415			12	5	1	1			2			
Normal4_N		3		1	5	136	221	62	656	1140	484	189	5		
Normal4_T															
Tumor4_N		1		13	5	104	16	6	10	109	118	3			
Tumor4_T	1	3	4	3437	1	2			2	13	8	2			
Tumor5_N					83	1609	89	2		24	441	1			90
Tumor5_T		2		15	2393	27	1	1		17	19				1

Tumori_N means nuclei from the tumor tissue of patient i (Tumori) with normal genome (_N).

Normalj_T means nuclei from the adjacent normal tissue of patient j (Normalj) with tumor genome (_T).

b



c

	RNA1	RNA2 (RNA2a+RNA2b)	RNA3	RNA4	RNA5	Macrophage/Mono	T cells	Clec9A+ DCs	Plasma cells	Fibroblast	EC (VEC+LEC)	EP	EP-N1	EP-T1	Megakaryocyte
Normal1_N				1		7	8	1	1	489	226	10	270		
Normal1_T															
Tumor1_N						61	5	1		14	33	5	2	88	
Tumor1_T	3009									2					
Normal2_N		1				30	31	4	2	497	142	494			
Normal2_T		2													
Tumor2_N		9	2		1	1531	2745	115	486	825	456	495	1	2	
Tumor2_T		4399		2		2	6	1		2	3	4			
Tumor3_N			2		1	448	39	45	17	137	120	7			7
Tumor3_T			1882												
Normal4_N				1		120	207	57	612	1047	444	153	4		
Normal4_T															
Tumor4_N						84	10	4	5	83	104	2			
Tumor4_T	1		1	3053						1					
Tumor5_N					6	1333	76	2		16	328	1			63
Tumor5_T				3	2096					3					

Fig. 4. Multinomial wheel analyses quantifying the deviation of each cell to the major clusters and removing most of the cross-overs.

a, The projection of normal-genome (_N) or tumor-genome (_T) nuclei from each biopsy sample into RNA clusters based on clustering results in panels (a) and (b) of Fig. 3. **b**, We show the most-likely position for each nucleus on the DNA multinomial wheel based on the multinomial analysis (second panel from the left), and display in histogram the distance between the multinomial-wheel assignment and the Seurat assignment for all nuclei (first panel from the left). In addition, in the second panel from the right, we only show the nuclei with mixed identities (collisions) from the first mixture experiment of Tumor 1 and Tumor 5 on the same DNA multinomial wheel, and display in histogram the distance between the multinomial-wheel assignment and the Seurat assignment for all collisions (first panel from the right). **c**, The projection table in (a) after removing the suspected collisions which are between 2 and 10 units away from the major clusters .

The projection of normal-genome nuclei with 2 or 1 copy of chrX onto stroma cell types

	Macrophage /Mono	T cells	Cle9A+ DC	Plasma BC	Fibroblast	VEC	LEC	EP	EP-N1	EP-T1	Megakaryocyte
Normal1_N_2chrX	5	4	1	1	241	109	7	5	182		
Normal1_N_1chrX					17 (6.6%)	4	1				
Tumor1_N_2chrX	39	2	1		6	17	2	1	1	40	
Tumor1_N_1chrX										1	
Normal2_N_2chrX	23	21	2	2	385	103	10	392			
Normal2_N_1chrX	1	2			5	1		6			
Tumor2_N_2chrX	1119	2130	88	221	613	319	34	406	1		2
Tumor2_N_1chrX	23 (2.0%)	104 (4.7%)	4	149 (40.3%)	16 (2.5%)	4		2			
Tumor3_N_2chrX	212	29	31	10	76	35	17	3			2
Tumor3_N_1chrX	2		1	1	1						
Normal4_N_2chrX	61	112	36	286	508	222	16	90	2		
Normal4_N_1chrX		20 (15.2%)		15 (5.0%)	13 (2.5%)	3	1	1			
Tumor4_N_2chrX	41	4	3	3	36	50	2	2			
Tumor4_N_1chrX	1	1			3	1					
Tumor5_N_2chrX	897	56	2		12	209	2	1			40
Tumor5_N_1chrX	11 (1.2%)	2				1					

Normal_i_N_2chrX means the adjacent normal tissue from patient i (Normal_i) with normal genome (_N) with both copies of chrX (_2chrX).

Tumor_j_N_1chrX means the tumor tissue from patient j (Tumor_j) with otherwise normal genome (_N) but only one copy of chrX (_1chrX).

Nuclei with tumor DNA genome are not shown in this table.

Table 1. The projection of normal-genome cells with two or one copy of chrX into stromal expression clusters.