1 **Single-cell and single-nucleus RNA-sequencing from paired normal-adenocarcinoma lung**
2 **samples provides both common and discordant biological insights**
3
4

5 Sébastien Renaut[1], Victoria Saavedra Armero[1], Dominique K. Boudreau[1], Nathalie Gaudreault[1], Patrice
6 Desmeules[1], Sébastien Thériault[1], Patrick Mathieu[1], Philippe Joubert[1], Yohan Bossé[1,2]
7
8 1) Institut universitaire de cardiologie et de pneumologie de Québec – Université Laval, Quebec City,
9 Canada
10
11 2) Department of Molecular Medicine, Université Laval, Quebec City, Canada
12
13
14 **Corresponding author**
15 Yohan Bossé, Ph.D.
16 Scientific Director of the Quebec Heart and Lung Institute
17 Professor, Laval University
18 Department of Molecular Medicine
19 Canada Research Chair in Genomics of Heart and Lung Diseases
20 Institut universitaire de cardiologie et de pneumologie de Québec – Université Laval
21 Pavillon Marguerite-d'Youville, Y2106
22 2725 chemin Sainte-Foy
23 Quebec City (Quebec)
24 Canada, G1V 4G5
25 Tel: 418-656-8711 ext. 3725
26 email: yohan.bosse@criucpq.ulaval.ca
27
28

## Abstract

Whether single-cell RNA-sequencing (scRNA-seq) captures the same biological information as single-nuclei RNA-sequencing (snRNA-seq) remains uncertain and likely to be context-dependent. Herein, a head-to-head comparison was performed in matched normal-adenocarcinoma human lung samples to assess biological insights derived from scRNA-seq versus snRNA-seq and better understand the cellular transition that occurs from normal to tumoral tissue. Here, the transcriptome of 160,621 cells/nuclei was obtained. In non-tumor lung, cell type proportions varied widely between scRNA-seq and snRNA-seq with a predominance of immune cells in the former (81.5%) and epithelial cells (69.9%) in the later. Similar results were observed in adenocarcinomas, in addition to an overall increase in cell type heterogeneity and a greater prevalence of copy number variants in cells of epithelial origin, which suggests malignant assignment. The cell type transition that occurs from normal lung tissue to adenocarcinoma was often discordant whether cells or nuclei were examined. In addition, we showed that the ligand-receptor interactome landscape of lung adenocarcinoma was largely different whether cells or nuclei were evaluated. Immune cell depletion in fresh specimens partly mitigated the difference in cell type composition observed between cells and nuclei. However, the extra manipulations affected cell viability and amplified the transcriptional signatures associated with stress responses. In conclusion, research applications focussing on mapping the immune landscape of lung adenocarcinoma benefit from scRNA-seq in fresh samples, whereas snRNA-seq of frozen samples provide a low-cost alternative to profile more epithelial and cancer cells, and yield cell type proportion that more closely match tissue content.

## Introduction

Single-cell transcriptomics (scRNA-seq) has the ability to inspect the cellular heterogeneity of tissue and cancer with unprecedented details, and as such provide important insights into the cellular origin and cell-specific molecular defects that play a role in disease pathogenesis[1–4]. However, given the pace at which the field is evolving, uncertainties remain with respect to the design and analysis of single-cell transcriptomic datasets in order to gain the most from priceless biological samples. Fresh biospecimens are generally prioritized for cell viability and greater yield of high-quality cells. For tissues, scRNA-seq requires disaggregating the tissue to release individual cells into a single-cell suspension. Differences in dissociation and sample preparation efficiency across cell types are known to affect RNA integrity and can skew cell type proportions. A well-known instance of dissociation bias is observed in human lung tissue, where dissociation of fresh tumor (biopsies or resected specimens) commonly results in a majority of immune cells being sequenced[5–7]. While the aforementioned cell-type dissociation bias can be partly alleviated by enriching the epithelial cell fraction using EPCAM-based cell sorting[6], single cell preparation protocols may also affect cell viability and introduce transcriptional signatures associated with dissociation and stress responses[6,8,9].

Analyzing nuclei (snRNA-seq) instead of cells has been proposed as an alternative for frozen samples and tissues that cannot be readily dissociated[10]. While cellular compositions recovered from scRNA-seq versus snRNA-seq can vary substantially[11], the transition from cell to nuclei sequencing may help to reduce the dissociation bias and transcriptional stress responses, facilitate the study of difficult-to-dissociate tissues and cell types, and allow the assessment of large cells that cannot pass through microfluidics systems. At the same time, reference databases and cell type-specific gene markers, which are readily used to annotate unknown cell populations, have been largely built from scRNA-seq datasets[4] and therefore may not be optimal for snRNA-seq. Cell types and gene expression differences between scRNA-seq and snRNA-seq have been observed in mouse kidneys[12,13] and

76  brain[14,15] as well as in human metastatic breast cancer and neuroblastoma[11]. However, head-to-head

77  comparisons between scRNA-seq and snRNA-seq are still scarce and to the best of our knowledge, this

78  direct comparison has never been evaluated in the context of patient-matched normal lung and tumor

79  tissues.

80      Lung cancer is highly prevalent and the number one cause of cancer mortality. It thus represents

81  a medically valuable case study to compare the biological signal recovered through cells and nuclei

82  sequencing. A variety of experimental designs and samples have been evaluated by scRNA-seq in

83  patients with lung cancer. This includes lung samples enriched (e.g. FACS-sorted) for immune cells[16,17],

84  lung tumor of mixed histological types[2,7], and non-small cell lung cancer (NSCLC) samples before and

85  after targeted therapy[18] or immunotherapy[19]. More specifically in lung adenocarcinomas (LUAD), the

86  most common histological subtype of lung cancer, which originates from epithelial cells that line the

87  inside of the lungs, resected specimens or biopsies from two to eleven[2,5–7,20] patients have been

88  evaluated, but with a very limited number of paired normal-adenocarcinoma lung samples. Compared

89  with normal lung samples, epithelial cells from lung adenocarcinomas were characterized by a

90  depletion of alveolar cells (AT1 and AT2)[2,6], lost cell identity and more cells annotated as mixed-

91  lineage[5,21], higher transcriptome complexity and cell heterogeneity[6,22], patient-specific cancer cell

92  clusters[18], transcriptional states associated with survival[20,21], and AT2 cells dedifferentiated into a stem-

93  like state[22]. The shift in immune cells from normal to LUAD samples observed in previous studies

94  were similarly informative. It unveiled an increase in B, plasma and T regulatory cells coupled with a

95  decline in natural killer cells as well as reduced signatures of cytotoxicity in T cells, antigen

96  presentation in macrophages, and inflammation in dendritic cells, which are all coherent features of an

97  immunosuppressive tumor microenvironment[6,16]. Finally, differentially enriched ligand-receptor

98  interactions promoting tumorigenesis were also observed between LUADs and normal tissues[6,20].

4

99    Herein, specimens derived from the same patients were tested using both scRNA-seq in fresh

100   tissues and snRNA-seq from flash frozen tissues using the 10x Genomics® workflows. The biology

101   captured by both methods was compared in the context of paired tumor-normal human lung samples

102   explanted from patients that underwent surgery for lung adenocarcinoma. This study design revealed

103   the cellular and molecular transition that occurs from normal lung to adenocarcinoma, and evaluated

104   the commonality and discordance in the stemming biological insights gained from cells versus nuclei.

105   In addition, we compared the same paired normal-adenocarcinoma human lung samples using an

106   immune cell depletion protocol that alleviates the cell-type dissociation bias, with the aim of recovering

107   a more representative biological signal.

108

# Results

**Single cell/Nucleus dataset preparation**

Four patients, two tissue type (Normal/Tumor) and three experimental methods (scRNA-seq, snRNA-seq & immune-depleted scRNA-seq, hereafter labelled as *Cell*, *Nucleus* and *Immune-depleted cell*) were processed for a total of twenty-four samples. 160,621 cells/nuclei passed quality control (53,286; 57,078 and 50,257 for *Cell*, *Nucleus* and *Immune-depleted cell* datasets respectively) with a mean of 6,692 cells per sample (6,661; 7,135 and *6,282* for *Cell*, *Nucleus* and *Immune-depleted cell* datasets respectively, **Fig. 2A**) and a mean of 2,214 genes per cell (1,868; 2,309 and 2,473 genes for *Cell*, *Nucleus* and *Immune-depleted cell* datasets respectively, **Fig. 2B**). The experimental design is presented in **Fig. 1A-B**, while the clinical and cellular characteristics are detailed in **Tables S1** & **S2**, respectively.

From the 61 finest cell types annotations defined by Human Lung Cell Atlas (HLCA)[4], 35 were present in the current dataset at a frequency of >100 cells and we were able to annotate confidently 97.7% of cells at the coarsest level (*immune*, *epithelial*, *endothelial*, *stroma*, **Fig. 2C, Table S3**). This reference-based mapping and annotation approach is consistent with a marker-based approach for both the *Cell* and *Nucleus* datasets (**Fig. S1**). Nevertheless, cell type annotation scores were significantly and consistently lower (smaller fraction of annotated cells) in the *Nucleus* compared to the *Cell* dataset (two-way ANOVA, $p$-value < 2e-16), fine-level compared to high-level annotations ($p$-value < 2e-16) and Tumor compared to Normal tissue ($p$-value < 2e-16).

**Cell composition differs from Nucleus in Normal lung tissue.**

In **Fig. 3**, the UMAP visualisation showed that the *Cell* dataset from Normal lung tissue was largely dominated by immune cells, with 23,044 immune cells (81.5% of total, **Fig. 3A**). Conversely, the *Nucleus* dataset was dominated by epithelial cells, with 12,556 epithelial cells (69.9%, **Fig. 3B**). In

6

133　addition, the *Nucleus* dataset contained a larger fraction of unclassified cells compared to the *Cell*

134　dataset (7.3 % vs 0.1 %, Fisher Exact Test [FET], *p*-value < 2e-16).

135　　　　To further refine the immune community of cells, we sub-setted only the immune cells and

136　labelled the plots with a finer level (level 3) annotation (*Cell*, **Fig. 3C**; *Nucleus*, **Fig. 3D**). We observed

137　that the *Cell* dataset provided a better fine-grained classification as proportionally more cells could be

138　classified into specific cell types. To this effect, the *Nucleus* dataset contained a larger fraction of

139　unclassified cells (41.7 % vs 0.7 %, FET, *p*-value < 2e-16).

140　　　　We repeated this sub-setting approach for epithelial cells, given their primary role in the onset

141　of lung adenocarcinoma. We observed that *Cell* samples form distinct clusters mainly composed of

142　AT1, AT2 and multiciliated lineages (**Fig. 3E-F**). The *Nucleus* dataset, which had more than five times

143　more epithelial cells than the *Cell* dataset (12,556 versus 2,264), contained similar cell types and

144　mainly in similar proportions, except for a sizable fraction of unclassified cells that appeared largely

145　scattered in the UMAPs (10.9 % unclassified in *Nucleus* versus 1.29% in *Cell*, FET, *p*-value < 2e-16,

146　**Fig. 3E-F**).

147　　　　In **Fig. 4**, we present, for each cell type (level 3 annotation), the fraction of cells originating

148　from each patient (**Fig. 4A**), the number of cells (**Fig. 4B**) and the number of genes per cell (**Fig. 4C**).

149　In **Fig. 4D-F**, we present the same information for the *Nucleus* dataset and this visualization confirmed

150　that the *Nucleus* dataset has similar cellular composition, except for the over-representation of immune

151　cells in the *Cell* dataset. Both in *Cell* and *Nucleus* datasets, epithelial cell types were dominated by AT1

152　first and then AT2; endothelial cell types were dominated by capillary; and stromal cell types were

153　dominated by fibroblasts. With respect to the number of genes (transcripts) per cell (**Fig. 4 C, F**), we

154　observed many discordant patterns between *Nucleus* and *Cell* datasets, indicating that similar cell types

155　presented different overall transcriptional signatures based on the experimental method. For example,

156　in the *Cell* dataset, median numbers of genes per cell were low for monocytes (635), but high for T

7

157    cells (1,709), and the pattern was in the opposite direction for the *Nucleus* dataset (Monocytes = 2,729,

158    T cells = 1,055). For their part, alveolar cells AT1 and AT2 contained 50% more genes expressed in the

159    *Cell* dataset (AT1: 2,479 and AT2: 3,126) compared to the *Nucleus* (AT1: 1,639 and AT2: 2,004), and

160    fibroblast two times as much (2,101 vs 1,061).

161

162    **The cellular origin of tumoral cells**

163         In **Fig. 5A**, the UMAPs showed that *Cell* sequencing samples from lung tumor tissues were

164    largely dominated by immune cell types (20,410 immune cells vs 5,764 in *Nucleus* dataset), while in

165    **Fig. 5B**, the *Nucleus* dataset were dominated by epithelial cells (27,362 epithelial cells in *Nucleus* vs

166    1,220 in *Cell* dataset). For both Cell and Nucleus datasets, cells appeared more scattered (i.e., more

167    heterogeneous) in the tumor compared to normal lung (median *silhouette index* $_{(Normal)}$ = 0.69; median

168    *silhouette index* $_{(Tumor)}$ = 0.53; two-way ANOVA, *p*-value < 2e-16, **Fig. S2**). This shows a suboptimal

169    cell type assignment of heterogeneous tumor samples to the described lung cell types from the HLCA

170    reference.

171         In **Fig. 6A-C**, we present, for each level 3 annotation cell type, the fraction of cells from each

172    patient (**Fig. 6A**), the number of cells (**Fig. 6B**), the number of genes per cell (**Fig. 6C**) and in **Fig. 6E-**

173    **G**, we present the same information for the *Nucleus* dataset. First, we observed that within a coarse

174    level annotation, similar cell types and similar proportions are observed in *Cell* and *Nucleus* datasets.

175    For example, T cells largely dominated the immune cells, fibroblasts dominated the stroma cells and

176    endothelial cell types were relatively rare. With respect to epithelial cells, these were mainly composed

177    of unclassified and AT1 in both *Cell* and *Nucleus* datasets, and secretory epithelial cells appeared to be

178    mainly segregated to patient 3. However, rare cell types were much more common in the *Nucleus* than

179    the *Cell* datasets.

8

180    To distinguish malignant and non-malignant cells, we defined a genome-wide summary score

181    (CNV score) that relies on gene expression levels to identify gene deletion and duplication and serves

182    as a proxy to identify cancerous aneuploid cells[23]. This score was the highest for different epithelial cell

183    types depending whether we analysed the *Cell* dataset (rare, multiciliated lineage, AT1, unclassified,

184    **Fig. 6D)** or the Nucleus dataset (multiciliated lineage, secretory and unclassified, **Fig. 6H**). In addition,

185    we also noted that annotation scores were negatively correlated with CNV scores for *Cell* ($r^2 = 0.11$, *p*-

186    value $< 2e-16$) and *Nucleus* ($r^2 = 0.05$, *p*-value $< 2e-16$) datasets (**Fig. S3**).

187

**The cellular transition to lung adenocarcinoma**

189    Given the known epithelial origin of lung adenocarcinoma and the role of the immune system in

190    effectively controlling the growth of carcinoma cells, we analysed the transition in the proportions of

191    epithelial and immune cells from normal to adenocarcinoma tissue (**Fig. 7A-B**). Alveolar Type 1, AT2

192    and multiciliated cells decreased in relative abundance in adenocarcinomas, and this was consistent for

193    the *Cell* and *Nucleus* datasets. On the contrary, rare, secretory and unclassified epithelial cell types

194    increased in abundance in adenocarcinoma tissue in a consistent manner between *Cell* and *Nucleus*

195    datasets. For Immune cells, patterns were harder to interpret given the small number of immune cells in

196    the *Nucleus* dataset. Nevertheless, an augmentation of B and T cell lineages in adenocarcinoma was

197    found for both datasets, as well as a sharp drop in natural killer cells in the *Cell* dataset. For

198    macrophages and monocytes, a discordance in the transition from normal to tumor between scRNA and

199    snRNA was observed.

200

**The Ligand-receptor interactome differs between Cell and Nucleus**

202    In **Fig. 8A**, we visualised the incoming and outcoming interactions among 319 ligand-receptor

203    interactions (cell-cell contact) for the *Cell-Normal* dataset. The number of interactions between cell

204 types varies first according to the Cell vs. Nucleus method (two-way ANOVA, $F = 90.7$, $p$-value < 2e-

205 16) and then the Normal vs. Tumor tissue type ($F = 68.2$, $p$-value = 3.6e-16). In **Fig. 8B**, we show an

206 example of a typical pathway common in *Cell*, rare in *Nucleus* (Major Histocompatibility Complex-I)

207 and its interacting genes, which is more similar between *Normal* vs *Tumor* tissue of the same

208 experimental method (*Cell* vs *Nucleus*). An example pathway, rare in *Cell* but common in *Nucleus*

209 (Protein Tyrosine Phosphatase Receptor Type M) and its self interacting gene is presented in **Fig. 8C**.

210 In this case, each network shows differences according to both the experimental method and tissue.

211

212 **The effect of immune depletion on Cell sequencing**

213 In order to remove the large fraction of immune cells, we performed immune depletion in

214 Normal and Tumor single-cell suspensions. We confirmed that the *Immune-depleted cell* dataset was

215 enriched in epithelial cells and depleted in immune cells (**Fig. 9A-B**). As such, both the Normal and

216 Tumor tissues resemble the *Nucleus* dataset in the fact that they harbor a majority of epithelial cells

217 (61.5% and 69.9% of total for the *Immune-depleted cell* and *Nucleus* dataset, respectively), yet they

218 differ given that immune depleted cells harbor proportionally more endothelial (17.8% vs 4%) and

219 stromal (18.4% vs 7.9%) cell types, but less immune cells (1.3% vs 13.0%). In addition, Normal tissues

220 were largely composed of epithelial AT1 and AT2, while Tumor tissues also harbored secretory, rare

221 and unclassified cell types, much like the *Nucleus* dataset (**Fig. 9C-D**). Finally, as we observed for the

222 non-depleted dataset, we saw an increase in the heterogeneity from Normal to Tumor datasets (median

223 Silhouette index for each level 3 cell type annotation: $s_{i\ (Normal)} = 0.56$, median $s_{i\ (Tumor)} = 0.2$, two-way

224 ANOVA, $p$-value < 2e-16, **Fig. S2**).

225 Finally, we downloaded a set of 512 heat shock and stress response genes that were previously

226 identified as affected by the scRNA-seq method[9]. Ninety four percent (482 genes) of the genes in this

227 core dataset were also present in our current dataset, with varying levels of expression. More

228    specifically, the percentage of cells expressing these genes was largely dependent on the method (**Fig.**

229    **9E**, two-way ANOVA, *p*-value < 2e-16). The *Immune-depleted cell* dataset showed the highest

230    expression of the stress response genes, whereas on average a cell from the *Immune-depleted cell*

231    dataset expressed 21% of the 482 genes, compared to 11.0% and 6.9% for the *Cell* and *Nucleus* dataset,

232    respectively. In addition, the proportions of cells expressing this core set of stress response genes was

233    slightly, but significantly (*p*-value = 9.7e-8) higher in Tumor than in Normal (12.4 % and 11.5 %,

234    respectively) tissue. In a similar manner, higher mitochondrial contamination is often considered a sign

235    of lower cell quality or viability[24] and we observed that the percentage of unique sequences (UMIs)

236    assigned to mitochondrial genes in the raw data prior to any filtering was significantly higher (two-way

237    ANOVA, *p*-value = 3.6e-5) in the *Immune-depleted cell* (mean = 15.2 %) and *Cell* (11.2 %) compared

238    to the *Nucleus* (2.6%) dataset, while the tissue type (*p*-value = 0.10) had no significant effect (**Fig. S4**).

239

240

241

242

243

## Discussion

244         In this study we generated a dataset of 160,621 cells/nuclei showing commonalities and

246    discordances in biological insights derived from single-cell and single-nucleus RNA-sequencing of

247    paired normal-adenocarcinoma human lung specimens. A distinct portrait of cellular composition was

248    observed per experimental methods that favors scRNA-seq of fresh samples to map the immune

249    landscape of lung adenocarcinoma. On the other hand, snRNA-seq of frozen samples surpassed the

250    relative merits of scRNA-seq to obtain a dataset with cell type proportion that match tissue content and

251    to provide a more cost-effective approach for research applications necessitating a higher number of

252    epithelial and cancer cells (see **Table S4** for a summary of the benefits of each method). In these paired

253    lung samples, we identified gene expression and cell type transitions from normal to tumoral tissue that

254    were not always concordant whether cells or nuclei were examined. The most striking difference was

255    the ligand-receptor interactions that varied more across methods (cells vs. nuclei) rather than tissue

256    types (normal vs. tumor). Immune cell depletion partly alleviated the difference in cell type

257    composition between cells and nuclei, but at the detriment of inducing a stress response. Finally, our

258    analysis revealed that the recently proposed five-level hierarchical cell type annotation system by the

259    Human Lung Cell Atlas[4] will require customization for assigning cell types from nuclei and tumor

260    samples.

261         Despite the fact that samples originated from the same patients' specimens, scRNA-seq and

262    snRNA-seq varied substantially in their recovered cellular compositions and transcriptional landscape,

263    thus highlighting the considerable impact of methodology on biological inference. While it has been

264    shown previously that cryopreservation of tissue sample (such as performed for snRNA-seq) results in

265    a major loss of epithelial cell types and an underrepresentation of T, B, and NK lymphocytes in the

266    single-nucleus libraries[11,13], it is not necessarily apparent which experimental method is more

267    biologically relevant. Slyper et al.[11] have suggested to analyse both fresh and frozen tissue, but this is

12

268 often unrealistic in practice. For their part, Denisenko et al.[13] indicate that the apparent discordance in

269 the recovered cellular composition between scRNA and snRNA might be due to either an under-

270 representation of immune cells in snRNA, or an under-representation of other cell types cells in scRNA

271 due to incomplete dissociation. Early pioneering work in lung histology would suggest the latter,

272 whereas cell staining and electron microscopy has revealed that the alveolar regions of normal human

273 lungs are comprised mainly of epithelial, endothelial and interstitial cells, while immune cells

274 (macrophages) comprised a small fraction (~5%) of all cells identified[25]. We thus conclude that in the

275 context of lung adenocarcinoma and patient-matched normal samples, snRNA-seq provides a dataset

276 comprising cell populations more closely matching tissue content.

277       In addition, we observed a decrease in cell viability in both depleted and non-depleted scRNA-

278 seq, likely due to the longer sample preparation times at room temperature. While this could be partly

279 alleviated by cold-activated proteases[9], it favors snRNA-seq as a experimental protocol to preserve

280 sample integrity. Although immune depletion works well for removing immune cells and therefore

281 might draw a more accurate representation of the lung cellular composition that is closer to snRNA-seq,

282 it requires extra laboratory manipulations and has the adverse effect of affecting both cell viability (**Fig.**

283 **S4**) and inducing a dissociation transcriptional stress response (**Fig. 9E**), as shown previously[12].

284       The reference-based annotation used here provides an attractive alternative to unsupervised

285 analysis[26]. We annotated the large majority of cells/nuclei in all tissue types, methods and patients (**Fig.**

286 **2, Fig. S5**) while showing that it performed as well as a marker-based approach, at least at the coarsest

287 annotation level (**Fig. S1**). Arguably, the confidence in this reference-based annotation approach

288 depends on several factors. Notably, the comprehensiveness of the reference, the quality and type of

289 query data and the level of cellular granularity required to answer the biological question of interest

290 will dictate the best approach to use. Nevertheless, an unsupervised-marker based approach also

291 depends on several factors such as the clustering algorithm, the gene markers used, and almost always,

13

292   the expertise and subjectivity of the person annotating the dataset[27,28]. Here, annotation and mapping

293   were done using the same analytical framework for all samples and therefore provided an objective

294   overview of the transcriptional cellular landscape. Fortunately, we were able to use a recently published

295   comprehensive atlas of the lung (HLCA)[4], although thorough cell atlases might not exist for all tissue

296   types, biological conditions and demographic states[29]. The lower annotation scores observed in nuclei

297   and tumor samples and consequently the greater number of unclassified cells, especially at the finer

298   annotation levels suggest that these cells or nuclei have a distinct signature from the current reference

299   cell types. A similar phenomenon was also observed in the HLCA for different disease states[4] and the

300   authors concluded themselves that the HLCA must be viewed as a live resource that will require

301   continuous updates in the future, including samples of diverse ethnic, clinical and experimental (e.g.

302   snRNA-seq) backgrounds.

303       During the transition from normal to tumoral tissue, we identified a drop in AT1, AT2 and NK

304   cells, concurrently with a rise in immune B and T cells, as previously identified[2,6,16]. In addition,

305   tumoral cells showed an increased transcriptomic heterogeneity and a greater prevalence of copy

306   number variants in epithelial cells. Similarly, it has been described that NSCLC exhibit important

307   interpatient histologic heterogeneity and inferred origin of tumor cells[30]. Here, we showed that

308   epithelial AT1, secretory and multiciliated lineages cell types had higher Copy Number Variants scores

309   than AT2, which suggests malignant assignment. Yet, the distinction between these epithelial cells is

310   not always straightforward, especially in a context of oncogenesis. Along those lines, we noted that

311   annotation scores were negatively correlated with CNV scores which implies that cells with high CNV

312   (likely carcinoma cells) loose their cellular identity and become harder to classify as distinct lung cell

313   types. During the construction of the HLCA, Sikkema *et al.*[4] also noted than a significant fraction of

314   cells from adenocarcinomas did not cluster into the specific fine level cell types. Similarly, Wang et

315   al.[22] argued that cancer cells originate from 'AT2-like' cells, but also nuanced this fact and stated that

14

316    these form a distinct cluster from regular AT2 cells and in fact, have a transcriptional profile closely

317    resembling other epithelial cells. Again, a more refined and thorough reference database will help to

318    solve these questions.

319        Ultimately, we hope to develop a comprehensive transcriptional resource for the identification

320    of cell-targeted biomarkers and therapeutic targets to treat and prevent LUAD and other ailing aspects

321    of the lung. Accordingly, this study may have clinical significance as immunotherapy is currently

322    revolutionizing the treatment of lung cancer. Response to immune checkpoint inhibitors relies on the

323    existing cell-cell interactions between tumor and T cells (e.g., commercial immunotherapy drugs

324    targeting the interaction between PD-1 in tumor cells and PD-L1 in T cells)[31] and identifying accurate

325    biomarkers of response to immunotherapy is a major challenge in the field of lung cancer[32].

326    Consequently, this seems like a clinical problem where single-cell genomics can provide a solution.

327    However, here we demonstrated that the ligand-receptor interactome landscape of lung

328    adenocarcinoma is largely different whether cells or nuclei are evaluated. This may lead to conflicting

329    prediction response to these novel immunotherapy agents. Accordingly, at least in the context of lung

330    cancer, the choice between scRNA-seq and snRNA-seq has important implications. Our results favor

331    scRNA-seq on fresh samples to provide a more comprehensive portray and granularity of the immune

332    cells diversity. On the other hand, it may not be representative of the true cellular community, and lead

333    to fewer difficult-to-dissociate tumor cells to assess relevant tumor-immune interactions. More studies

334    will be needed to assess the best methods as well as to overcome other barriers to move single-cell

335    genomics into the clinical setting[33].

336

15

# Materials and methods

**Patients and samples**

Lung samples were collected from four patients that underwent curative intent primary lung cancer surgery at the *Institut universitaire de cardiologie et de pneumologie de Québec – Université Laval* (IUCPQ-UL) in 2021-2023, henceforth referred to patient 1, 2, 3 and 4. The four patients were self-reported white French Canadian (European ancestry) with no prior chemotherapy and/or radiation therapy, and all patients were between the age of 59 and 69, former smokers with adenocarcinomas (See **Fig. 1** for overview of experimental design, and **Table S1** for detailed clinical characteristics of patients).

Following surgery, the explanted lobes were immediately transferred to the pathology department. For each patient, two ☐1 cm$^3$ fresh tumor samples and two ☐1 cm$^3$ non-tumor (normal) lung samples located distant from the tumor were harvested. The first set of tumor/non-tumor samples was transferred in dedicated tubes containing ice-cold RPMI (ThermoFisher, Cat. 11875093) for immediate cell dissociation and single-cells RNA sequencing (scRNA-seq) experiment. The second set of tumor/non-tumor samples was transferred in dedicated tubes, immediately snap-frozen in liquid nitrogen and stored at -80°C until the day of the single-nucleus RNA sequencing (snRNA-seq) experiment. A histologic slide of each specimen was stained (H&E) and reviewed by a pathologist. Staging was performed using the 8[th] edition of the TNM Classification of Malignant Tumours[34]. Lung tissue samples were obtained in accordance with the Institutional Review Board guidelines. All patients provided written informed consent, and the ethics committee of the IUCPQ-UL approved the study.

**Sample preparation for scRNA-seq**

Immediately after collection, the weight of each sample was recorded. Samples were transferred to 6-well cell culture plates, washed twice with 3 mL ice-cold PBS (Thermo Fisher, cat. 10010023) to

16

361 remove excess blood and transferred to a 5 mL glass beaker. Using a 1 mL syringe and 25G needle,

362 300 µL of Enzyme dissociation mix was injected in the tissue followed by mechanical mincing into

363 small fragments (<1mm³) using spring scissors for 2 minutes. Samples were then transferred to 50 mL

364 Falcon tubes containing 5,7 mL of Enzyme dissociation mix and pipette mixed 5 times using wide bore

365 1 mL tips. The enzymatic digestion was performed at 37°C, using a Vari-Mix™ test tube rocker at max

366 speed for 35 minutes. Samples were pipette mixed 20 times after 15 and 30 minutes using wide bore 1

367 mL tips. Enzyme dissociation mix contained: Pronase 1250 µg/mL (Sigma Aldrich, cat. 10165921001),

368 Elastase 18.4 µg/ml (Worthington Biochemical, cat. LS006363), DNase I 100 µg/mL (Sigma Aldrich,

369 cat. 11284932001), Dispase 100 µg/mL (Worthington Biochemical, cat. LS02100), Collagenase A

370 1500 µg/mL (Sigma Aldrich, cat.10103578001) and Collagenase IV 100 µg/mL (Worthington

371 Biochemical, cat. LS 004186) in HBSS (Thermo Fisher, cat. 14170112). Enzymatic digestion was

372 stopped by adding 1.5 mL of fetal bovine serum (FBS, ThermoFisher, cat. A3840301) followed by

373 pipette mix 5 times using wide bore 1 mL tips. Dissociated cells were filtered through a 70 µm strainer

374 and washed with 7.5 mL ice-cold PBS. Cells were then pelleted at 400g, 4°C for five minutes and

375 supernatant was removed. Three cycles of red blood cells removal were performed as follow: cell pellet

376 resuspended by manual agitation in 500 µL of ACK Lysis Buffer (ThermoFisher, cat. A1049201) and

377 incubated on ice one minute. One mL of ice-cold PBS was added and cells were centrifuged at 400g,

378 4°C for two minutes and the supernatant was removed. The final pellet was resuspended in 500 µL ice-

379 cold-PBS containing 0.04% Bovine Serum Albumin (BSA, Sigma Aldrich Cat. A7284) and 10% FBS.

380 Cell suspensions were successively passed through 100 µm, 70 µm and 40 µm strainer using quick spin

381 to reach 400g to filtrate each sample. Samples were transferred to 2.0 mL low binding tubes and kept at

382 4°C. Cell count and viability were performed using a 1:1 mix of cell suspension, Trypan blue

383 (ThermoFisher, cat. 15250061), haemocytometer and conventional light microscopy. Cells suspensions

384 meeting the following criteria were accepted for scRNA-seq library preparation: absence of aggregated

385    cells, a viability >80%, and a total cell count between 400 and 1200 cells/µL. $1 \times 10^5$ cells were

386    transferred to a low binding 2 mL tube and kept at 4°C (non-depleted fraction). The remaining cells

387    (from 2 to 5 $\times 10^6$ cells) were submitted to CD45 immune cell depletion protocol (single cells depleted

388    fraction) as described below. The characteristics of the lung specimen and the single cell suspension for

389    each sample are given in **Table S2**.

390

391    **CD45 immune cell depletion**

392    Cells (from 2 to 5 $\times 10^6$ cells) were centrifuged at 300g, 4°C, 10 minutes. The supernatant was

393    removed and the cell pellet was resuspended in 80 µL MACS buffer (0.5% BSA, 2 mM EDTA pH 8.0

394    in PBS) previously degassed for 1 hour at room temperature. Twenty µL of CD45 microbeads

395    (Miltenyi Cat. 130-045-801) were added and sample was incubated 15 minutes at 4°C followed by

396    addition of 1 mL MACS buffer and centrifugation 300g, 10 minutes at room temperature. Supernatant

397    was removed and pellet resuspended in 2-steps 100 µL + 400 µL MACS buffer. The total volume (500

398    µL) was applied to a LS Positive Selection Column (Miltenyi Cat. 130-042-401) previously rinsed with

399    3 mL MACS buffer and installed on a MidiMACS magnetic Separator with a collection tube. Column

400    was rinsed with 3 X 3 mL MACS buffer and all volumes (9.5 mL) were collected which contained the

401    CD45-negative fraction. CD45-negative cells were centrifuged 300g, 10 minutes at room temperature

402    followed by supernatant removal. Cells were washed twice with 1 mL PBS followed by centrifugation

403    at 300g, 10 minutes after each wash. Cells were finally resuspended in 100 µL BSA 0.04%, 10% FBS

404    in PBS and kept at 4°C. Cell count and viability were performed using a 1:1 mix of cell suspension,

405    Trypan blue, haemocytometer and conventional light microscopy. Cells suspensions meeting the

406    following criteria were accepted for scRNA-seq library preparation: absence of aggregated cells, a

407    viability >80%, and a total cell count between 400 and 1200 cells/µL.

408

**Sample preparation for snRNA-seq**

Nuclei suspension was prepared from ~30 mg snap frozen tissue using Chromium Nuclei
Isolation Kit as per manufacturer protocol (10x Genomics Cat. 1000494). Nuclei count and integrity
were performed using a 1:1 mix of nuclei suspension and methylene blue 0.25% (Ricca Chemical, Cat.
48504), haemocytometer and conventional light microscopy. Nuclei suspensions meeting the following
criteria were accepted for snRNA-seq library preparation: absence of aggregated nuclei, nuclei with
circular shape and intact membrane (without blebbing) >80%, and a total nucleus count between 400
and 1200 nuclei/µL. Nuclei suspension were kept at 4°C until proceeding with 10x Genomics snRNA-
Seq library preparation protocol.

**10x Genomics sn/scRNA-seq library preparation**

For each sample, approximatively 15,000 nuclei or cells were loaded into each channel of a
Chromium Next Gel Beads-in-emulsion (GEM) Chip G (10x Genomics Cat. 1000127) as per
manufacturer instruction for GEM generation and barcoding. Given the cell capture efficiency of
around 65%, 10,000 cells per library were therefore expected. The Chip was run on the Chromium
Controller, GEMs were aspirated and transferred to a strip tube for cDNA synthesis, cDNA
amplification and library construction using Chromium Next GEM single-cell 3' Library Kit v3.1 (10x
Genomics Cat. 1000128) and Single Index Kit T Set A (10x Genomics Cat. 2000240) as per
manufacturer instruction. The library average fragment size and quantification was performed using
Agilent Bioanalyzer High Sensitivity DNA kit (Agilent Cat. 5067-4626) and a final concentration
determination was performed using NEBNext® Library Quant Kit for Illumina (New England Biolabs
Cat. E7630) prior to library sequencing.

**Next generation sequencing**

Libraries were individually diluted to 10 nM, pooled and sequenced on an Illumina NextSeq2000 system following manufacturer's recommendations. Sequencing was realized on a P3 (100 cycles) cartridge, aiming for 200 to 500 million reads per library (sample). Run parameters for paired-end sequencing were as follows: read 1, 28 nucleotides; read 2, 91 nucleotides; index 1, 8 nucleotides; and index 2, 0 nucleotide.

**Single cell/nucleus data preparation**

Demultiplexing, alignment and transcript counting was performed using the *Cellranger* software (v7.1.0, 10x Genomics) on our local server (Lenovo ThinkSystem SR650, 40 cores and 384GB RAM). The BCL files from the Illumina sequencing run were first demultiplexed into FASTQ files using the *cellranger mkfastq* command. Read alignment and UMI counting were then executed with the *cellranger count* command (see alignment and cell statistics in **Table S5**). We used GRCh38 as the reference transcriptome available on Gencode, release 43 (GRCh38.p13).

**Data quality control**

The most up-to-date bioinformatics procedure defined by the R (v4.3.0)[35] library *Seurat* (v4.3.0)[24] was used to create an object for each sample and calculate values for *nCount* (number of Unique Molecular Identifiers [UMI] per cell), *nFeatures* (number of genes expressed per cell) and *percent.mt* (fraction of UMIs aligning to mitochondrial genes) parameters. Using the R library *scuttle* (v1.10.1)[36], we determined outlier values for *nCount*, *nFeatures* and *percent.mt* based on the median absolute deviation and sub-setted each sample accordingly. Note that for the *percent.mt* parameter, if necessary, we further capped this outlier value at twenty-five percent per sample.

455    For each sample, we then performed normalization and variance stabilization using the function

456    *SCTransform*, which also has the benefit to regress out the *percent.mt* effect from the underlying count

457    data. Then, using the R library *DoubletFinder* (v2.0.3)[37], we identified and removed doublets

458    (assuming a five percent doublet rate), which occur when multiple cells are captured into a single oil

459    droplet during the GEM generation.

460

461    **Reference-based cell type annotation and mapping**

462    On each of these curated samples, cellular annotation was performed using the R library

463    *Azimuth* (v0.4.6)[26] and the most recent version of the Human Lung Cancer Atlas (HLCA v2)[4]. Note

464    that in the subsequent methodology, *cell* annotation refers to the annotation of a uniquely barcoded

465    GEM sample stemming from either a scRNA-seq or a snRNA-seq dataset.

466    The HLCA is a comprehensive and curated reference dataset constructed using a diverse set of

467    107 healthy lung samples (584,444 cells) and which allows to identify the transcriptional signature of

468    61 hierarchical cell types, from the coarsest possible annotations (level 1: *Immune*, *Epithelial*,

469    *Endothelial* and *Stroma*), recursively broken down into finer levels (levels 2-5). In addition, this

470    reference-based mapping approach allows to robustly and sensitively compare samples of broad

471    cellular compositions, while also identifying specific and rare cell populations[24,26,38]

472    Specifically, for each sample (query), the algorithmic approach first identifies anchors between

473    the reference and query (that is, pairs of cells from each dataset that are contained within each other's

474    neighborhoods) and uses these anchors to integrate the query dataset onto the reference. Then, the

475    embeddings of the query data onto the reference Principal Components (50 PCs) are calculated and

476    visualised directly onto the reference two-dimensional Uniform Manifold Approximation and

477    Projection (UMAP). Finally, annotation scores [0:1], which reflect the confidence in the annotation,

478    were used to label cell types, whereas cells with annotation scores < 0.5 were labelled as *unclassified*.

479

**Copy number variations analysis**

For each scRNA-seq and snRNA-seq tumor sample, we performed an analysis of Copy-Number Variants (CNVs) in order to identify malignant aneuploid cells based on the premise that gene CNVs can be identified using the difference between the mean log expression level of non-cancerous reference cells (here immune cells) and the log gene expression level of a cell of interest. This was performed using the R library *infercnv* (v1.17.0)[23] and a general index (CNV score) for each cell was then defined as the mean sum of square of scaled [-1;+1] standardized log fold-change values.

487

**Biological dataset comparisons**

We integrated twenty-four samples into six different datasets (*Cell-Normal*, *Nucleus-Normal*, *Cell-Tumor*, *Nucleus-Tumor*, *Immune-depleted cell-Normal*, *Immune-depleted cell-Tumor*), in order to quantify biological similarities and differences among datasets (see **Fig. 1F-G** for summary of comparisons and accompanying figures). Given that the same reference dimensionality reduction (PCA) and visualisation space (UMAP) was used for each sample, we could simply merge expression data, metadata and projections into objects that accounts for technical variation among sample in order to quantify patterns. For each individual cell, we also calculated a Silhouette index[39] to evaluate the goodness of fit of the clustering, whereas the index is calculated from the UMAP embeddings and the clusters correspond to specific cell type (level 3) annotations. We then tested the effect of the experimental method and tissue type on the Silhouette index using a two-way Analysis of Variance (ANOVA).

500

**Ligand-receptor analysis**

502    In order to infer and visualise the intercellular communication among cell populations, we used

503    the R library *cellchat* (v 1.6.1)[40]. We quantified the cell-cell interaction pathways in normal and tumor

504    tissue (*cell* and *nucleus* dataset) to describe the cellular transition during oncogenesis and quantify how

505    the experimental method and tissue type affected the results. We limited this analysis to level 3

506    annotation and excluded infrequent cell types (<500 cells in total) and cells that were unclassified at the

507    level 3 annotation. We quantified the number of interactions from and to each cell type and tested the

508    effect of the experimental method and tissue type using a two-way ANOVA.

509

510    **Stress-related genes**

511    To quantify the effect of our *Cell*, *Nucleus* and *Immune Depleted Cell* experimental methods on

512    the overall stress responses of the cell populations, we analysed the expression pattern of a core set of

513    512 heat shock and stress response genes that were previously identified to be affected by the scRNA-

514    seq sample preparation method[9]. We quantified the proportions of cells that expressed these genes for

515    each sample and tested the effect of the experimental method, tissue type and patient using a two-way

516    ANOVA.

517

# Supplementary Information

**Authors' contributions**

PD, ST, PM, PJ and YB conceived the study. PD and PJ oversaw the sample pathology. SR and YB

wrote the manuscript. VA, DB, NG conducted the single-cell experiments and sequencing. SR

analyzed the data. All authors read and approved the final manuscript.

**Data availability statement**

The datasets generated by *Cellranger* will be available as open-access downloadable files on Zenodo

upon acceptance (zenodo.org/records/10144050). All analytical codes used to produce the results of

this study will be made available at https://github.com/Yohan-Bosse-Lab/scRNA

539

# **References**

541    1.  Puram, S. V. *et al.* Single-Cell Transcriptomic Analysis of Primary and Metastatic Tumor

542        Ecosystems in Head and Neck Cancer. *Cell* **171**, 1611-1624.e24 (2017).

543    2.  Lambrechts, D. *et al.* Phenotype molding of stromal cells in the lung tumor microenvironment. *Nat.*

544        *Med.* **24**, 1277–1289 (2018).

545    3.  Wu, S. Z. *et al.* A single-cell and spatially resolved atlas of human breast cancers. *Nat. Genet.* **53**,

546        1334–1347 (2021).

547    4.  Sikkema, L. *et al.* An integrated cell atlas of the lung in health and disease. *Nat. Med.* **29**, 1563–

548        1577 (2023).

549    5.  Laughney, A. M. *et al.* Regenerative lineages and immune-mediated pruning in lung cancer

550        metastasis. *Nat. Med.* **26**, 259–269 (2020).

551    6.  Sinjab, A. *et al.* Resolving the Spatial and Cellular Architecture of Lung Adenocarcinoma by

552        Multiregion Single-Cell Sequencing. *Cancer Discov.* **11**, 2506–2523 (2021).

553    7.  Zilionis, R. *et al.* Single-Cell Transcriptomics of Human and Mouse Lung Cancers Reveals

554        Conserved Myeloid Populations across Individuals and Species. *Immunity* **50**, 1317-1334.e10

555        (2019).

556    8.  van den Brink, S. C. *et al.* Single-cell sequencing reveals dissociation-induced gene expression in

557        tissue subpopulations. *Nat. Methods* **14**, 935–936 (2017).

558    9.  O'Flanagan, C. H. *et al.* Dissociation of solid tumor tissues with cold active protease for single-cell

559        RNA-seq minimizes conserved collagenase-associated stress responses. *Genome Biol.* **20**, 210

560        (2019).

561    10. Krishnaswami, S. R. *et al.* Using single nuclei for RNA-seq to capture the transcriptome of

562        postmortem neurons. *Nat. Protoc.* **11**, 499–524 (2016).

563   11. Slyper, M. *et al.* A single-cell and single-nucleus RNA-Seq toolbox for fresh and frozen human

564        tumors. *Nat. Med.* **26**, 792–802 (2020).

565   12. Wu, H., Kirita, Y., Donnelly, E. L. & Humphreys, B. D. Advantages of Single-Nucleus over

566        Single-Cell RNA Sequencing of Adult Kidney: Rare Cell Types and Novel Cell States Revealed in

567        Fibrosis. *J. Am. Soc. Nephrol. JASN* **30**, 23–32 (2019).

568   13. Denisenko, E. *et al.* Systematic assessment of tissue dissociation and storage biases in single-cell

569        and single-nucleus RNA-seq workflows. *Genome Biol.* **21**, 130 (2020).

570   14. Bakken, T. E. *et al.* Single-nucleus and single-cell transcriptomes compared in matched cortical

571        cell types. *PloS One* **13**, e0209648 (2018).

572   15. Lake, B. B. *et al.* A comparative strategy for single-nucleus and single-cell transcriptomes confirms

573        accuracy in predicted cell-type expression from nuclear RNA. *Sci. Rep.* **7**, 6031 (2017).

574   16. Leader, A. M. *et al.* Single-cell analysis of human non-small cell lung cancer lesions refines tumor

575        classification and patient stratification. *Cancer Cell* **39**, 1594-1609.e12 (2021).

576   17. Guo, X. *et al.* Global characterization of T cells in non-small-cell lung cancer by single-cell

577        sequencing. *Nat. Med.* **24**, 978–985 (2018).

578   18. Maynard, A. *et al.* Therapy-Induced Evolution of Human Lung Cancer Revealed by Single-Cell

579        RNA Sequencing. *Cell* **182**, 1232-1251.e22 (2020).

580   19. Liu, B. *et al.* Temporal single-cell tracing reveals clonal revival and expansion of precursor

581        exhausted T cells during anti-PD-1 therapy in lung cancer. *Nat. Cancer* **3**, 108–121 (2022).

582   20. Kim, N. *et al.* Single-cell RNA sequencing demonstrates the molecular and cellular reprogramming

583        of metastatic lung adenocarcinoma. *Nat. Commun.* **11**, 2285 (2020).

584   21. Marjanovic, N. D. *et al.* Emergence of a High-Plasticity Cell State during Lung Cancer Evolution.

585        *Cancer Cell* **38**, 229-246.e13 (2020).

586    22. Wang, Z. *et al.* Deciphering cell lineage specification of human lung adenocarcinoma with single-

587        cell RNA sequencing. *Nat. Commun.* **12**, 6500 (2021).

588    23. Tickle, T., Tirosh, I., Georgescu, C., Brown, M. & Haas, B. *inferCNV of the Trinity CTAT Project.*

589        (Klarman Cell Observatory, Broad Institute of MIT and Harvard, 2019).

590    24. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29 (2021).

591    25. Crapo, J. D., Barry, B. E., Gehr, P., Bachofen, M. & Weibel, E. R. Cell Number and Cell

592        Characteristics of the Normal Human Lung.

593    26. Butler, A., Darby, C., Hao, Y., Hoffman, P. & Satija, R. *Azimuth: A Shiny App Demonstrating a*

594        *Query-Reference Mapping Algorithm for Single-Cell Data.* (2022).

595    27. Xie, B., Jiang, Q., Mora, A. & Li, X. Automatic cell type identification methods for single-cell

596        RNA sequencing. *Comput. Struct. Biotechnol. J.* **19**, 5874–5887 (2021).

597    28. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a tutorial.

598        *Mol. Syst. Biol.* **15**, e8746 (2019).

599    29. Snyder, M. P. *et al.* The human body at cellular resolution: the NIH Human Biomolecular Atlas

600        Program. *Nature* **574**, 187–192 (2019).

601    30. Chen, Z., Fillmore, C. M., Hammerman, P. S., Kim, C. F. & Wong, K.-K. Non-small-cell lung

602        cancers: a heterogeneous set of diseases. *Nat. Rev. Cancer* **14**, 535–546 (2014).

603    31. Garon, E. B. *et al.* Pembrolizumab for the Treatment of Non–Small-Cell Lung Cancer. *N. Engl. J.*

604        *Med.* **372**, 2018–2028 (2015).

605    32. Mino-Kenudson, M. *et al.* Predictive Biomarkers for Immunotherapy in Lung Cancer: Perspective

606        From the International Association for the Study of Lung Cancer Pathology Committee. *J. Thorac.*

607        *Oncol. Off. Publ. Int. Assoc. Study Lung Cancer* **17**, 1335–1354 (2022).

608    33. Lim, J. *et al.* Transitioning single-cell genomics into the clinic. *Nat. Rev. Genet.* **24**, 573–584

609        (2023).

610 34. Brierley, J. D., Gospodarowicz, M. K. & Wittekind, C. *TNM Classification of Malignant Tumours*.

611 (John Wiley & Sons, 2017).

612 35. R Core Team. *R: A Language and Environment for Statistical Computing*. (R Foundation for

613 Statistical Computing, 2023).

614 36. McCarthy, D. J., Campbell, K. R., Lun, A. T. L. & Wills, Q. F. Scater: pre-processing, quality

615 control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* **33**,

616 1179–1186 (2017).

617 37. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-Cell

618 RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst.* **8**, 329-337.e4 (2019).

619 38. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).

620 39. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis.

621 *J. Comput. Appl. Math.* **20**, 53–65 (1987).

622 40. Jin, S. *CellChat: Inference and analysis of cell-cell communication from single-cell and spatial*

623 *transcriptomics data*. (2023).

624

625

626

627



628
629
630 **Figure 1 | Overview of the experimental design.** For each patient (**A**), a tumor specimen and a
631 normal (non-malignant) lung specimen harvested from a site distant from the tumor were resected (**B**).
632 The research specimens were immediately divided into smaller fragments. For both normal and tumor
633 lung specimens, a fragment was frozen in liquid nitrogen and stored at -80$^{\circ}$C until further processing
634 for snRNA-seq. For fresh specimens, the fragments proceeded directly to dissociation into single-cell
635 suspensions. A subsample of the dissociation mix underwent immune cell depletion (**C**). The final set
636 of samples (**D**) were then loaded in wells of the microfluidic chip (**E**) in order to generate the
637 transcriptome of approximately 10,000 cells or nuclei per sample (**F**). Dataset comparisons performed
638 with accompanying figures (**G**).
639
640

641



**Figure 2** | **Overview of the 160,621 cells/nuclei that passed quality control obtained from lung tumors and distal normal lung samples. A.** Number of cells retained after quality control for each patient, each experimental method (*Cell*, *Nucleus*, *Immune-depleted cell*) and tissue type (*Normal*, *Tumor*). **B.** Mean number of genes per cell, per patient, method and tissue type. **C.** The fraction of annotated cells for each of the five-level HLCA hierarchical cell annotation reference framework, per method and tissue type.

652



653
654
655 **Figure 3 | UMAP representations and cell types annotations (Normal tissue)** for *Cell* (**A**) and
656 *Nucleus* (**B**) datasets with general cell types (level 1) annotation. Finer-grained annotation (level 3) for
657 the subset of immune cells (**C**) or nuclei (**D**) and for the subset of epithelial cells (**E**) or nuclei (**F**). To
658 the right of each UMAP, stacked bar plots indicate the proportion of each cell type in the specific
659 dataset. Cell types present at < 1% are labelled as others.
660
661

662



663
664
665 **Figure 4 | Cell types characteristics (Normal tissue).** For each of the four coarse (level 1) cell types
666 annotation (*Immune*, *Epithelial*, *Endothelial*, *Stroma*) further refined into finer categories (level 3): the
667 fraction of cells (**A**: *Cell* dataset, **D**: *Nucleus*) and the number of cells (**B**: *Cell*, **E**: *Nucleus*) originating
668 from each patient. Box plots of the number of genes expressed per cell (**C**: *Cell*, **F**: *Nucleus*), with plot
669 center, box and whiskers corresponding to median, IQR and 1.5□×□IQR, respectively. Note that only
670 cell types with > 20 cells were retained for clarity in this visual representation.
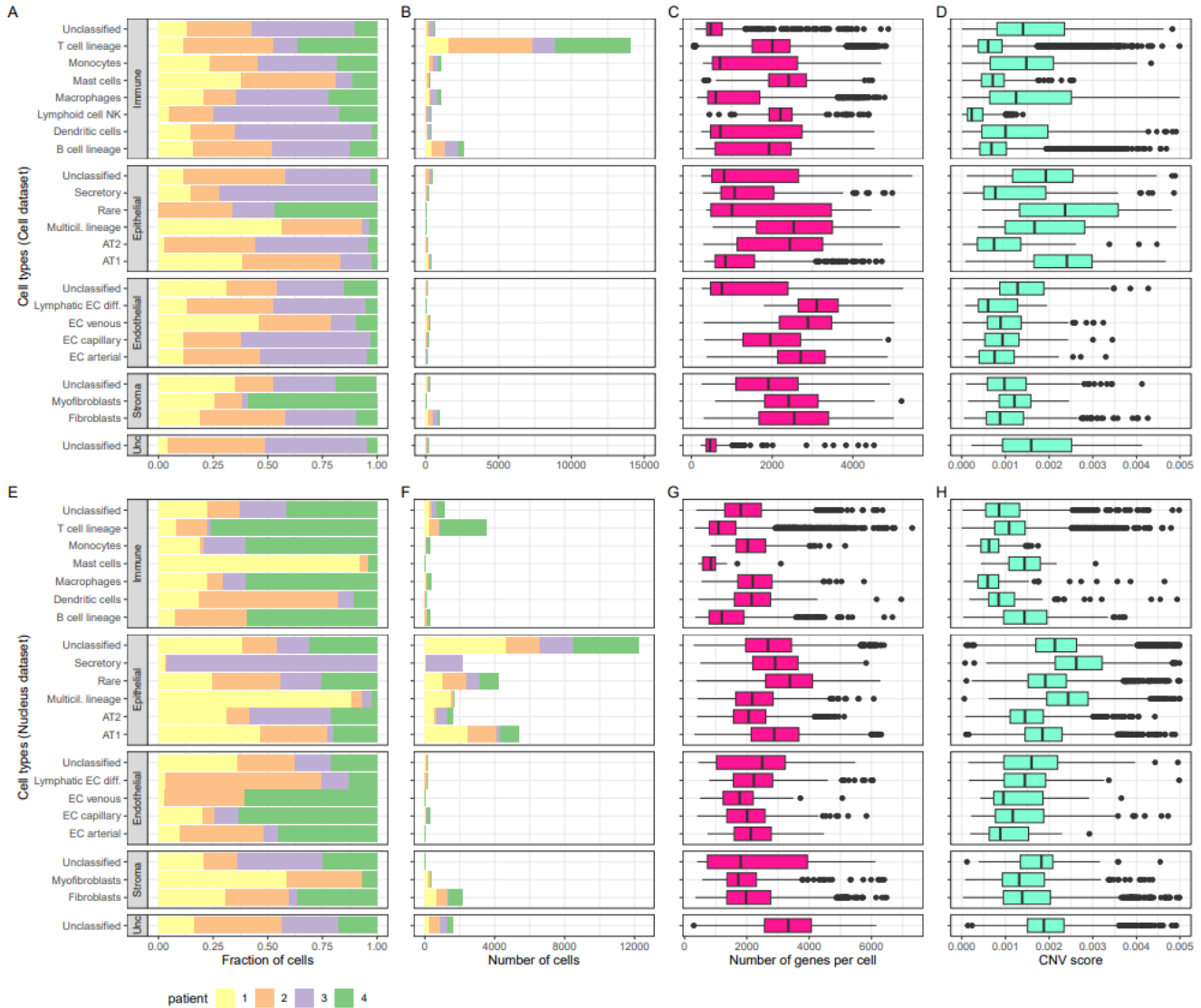
671
672

32

673



674
675

**Figure 5 | UMAP representations and cell types annotations (Tumor tissue)** for *Cell* (**A**) and *Nucleus* (**B**) datasets with general cell types (level 1) annotation. Tumor samples are overlaid on top of Normal samples (in gray). To the right of each UMAP, stacked bar plots indicate the proportion of each cell type in the specific dataset.
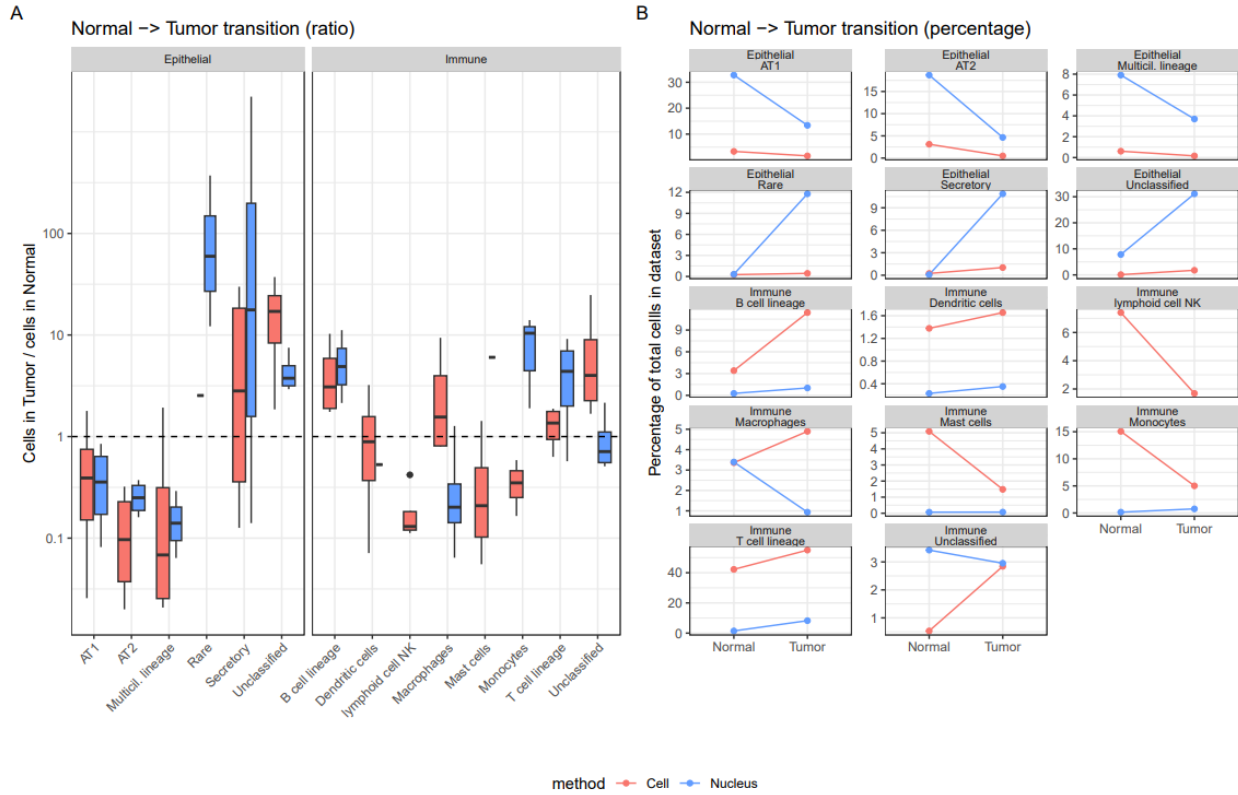
680
681

682



683
684
685 **Figure 6 | Cell types characteristics (tumor tissue).** For each of the four coarse (level 1) cell types
686 annotations (*Immune*, *Epithelial*, *Endothelial*, *Stroma*) and unclassified (*unc*), further refined into finer
687 categories (level 3 cell types): the fraction of cells (**A**: *cell* samples, **E**: *nuclei* samples) and the number
688 of cells (**B**: *cell*, **F**: *nucleus*) originating from each patient. Box plots of the number of genes expressed
689 (**C**: *cell*, **G**: *nucleus*) and the CNV score (**D**: *cell*, **H**: *nucleus*), with plot center, box and whiskers
690 corresponding to median, IQR and 1.5□×□IQR, respectively. Note that only cell types with > 20 cells
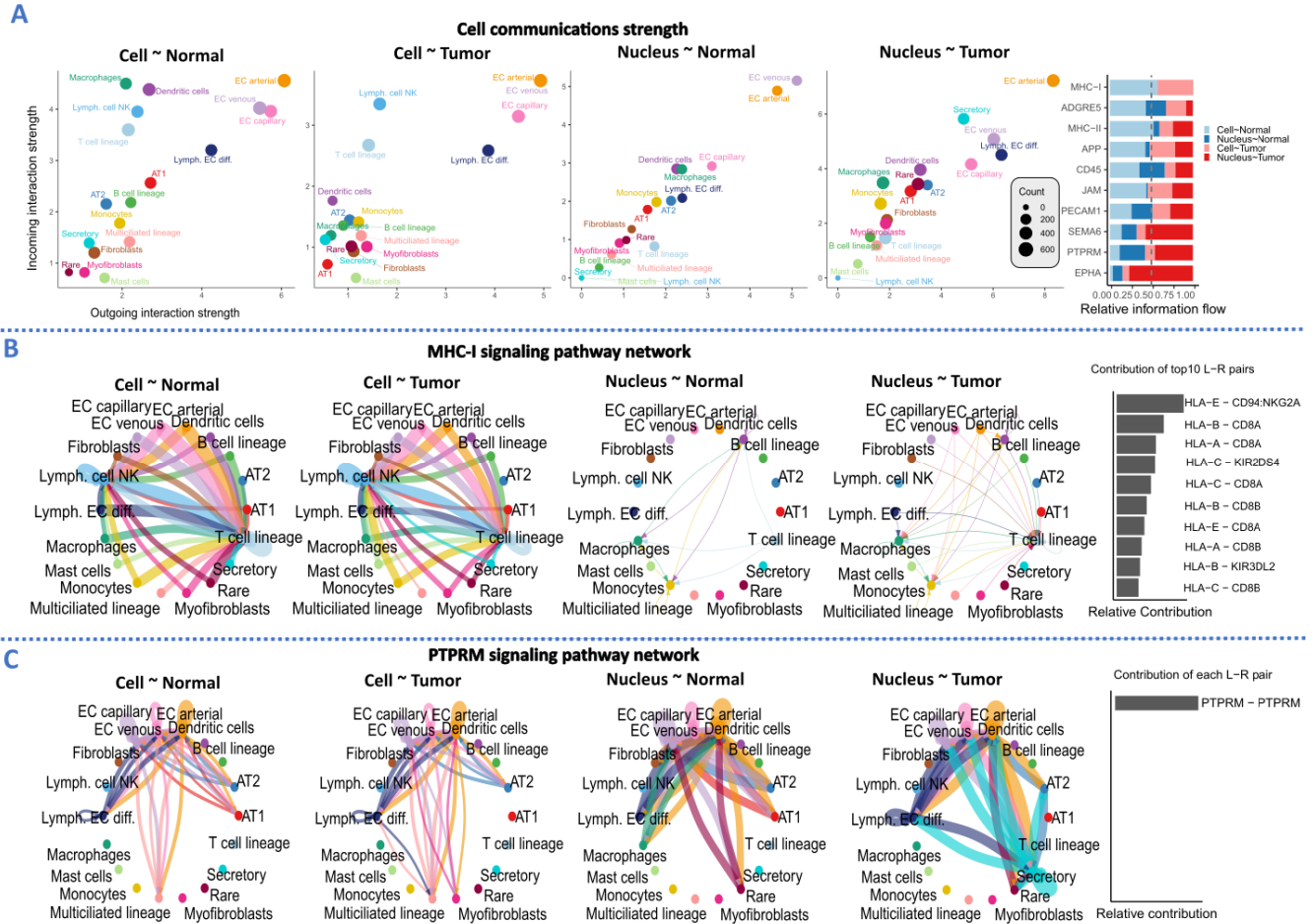691 were retained for clarity in this visual representation.
692
693

694



695
696
697 **Figure 7 | Normal - tumor transition**. **A.** For each specific (level 3) Epithelial or Immune cell type,
698 the fraction of cells they represent in the *Tumor* dataset divided by the fraction of cells they represent in
699 the *Normal* dataset (ratios above 1 represent an increase in the Tumor dataset), with plot center, box
700 and whiskers corresponding to median, IQR and $1.5 \times$ IQR, respectively **B.** The percentage of
701 specific (level 3) Epithelial or Immune cell types in *Tumor* and *Normal* dataset. Note that only cell
702 types with > 20 cells were retained for clarity in this visual representation.
703
704

705



706
707

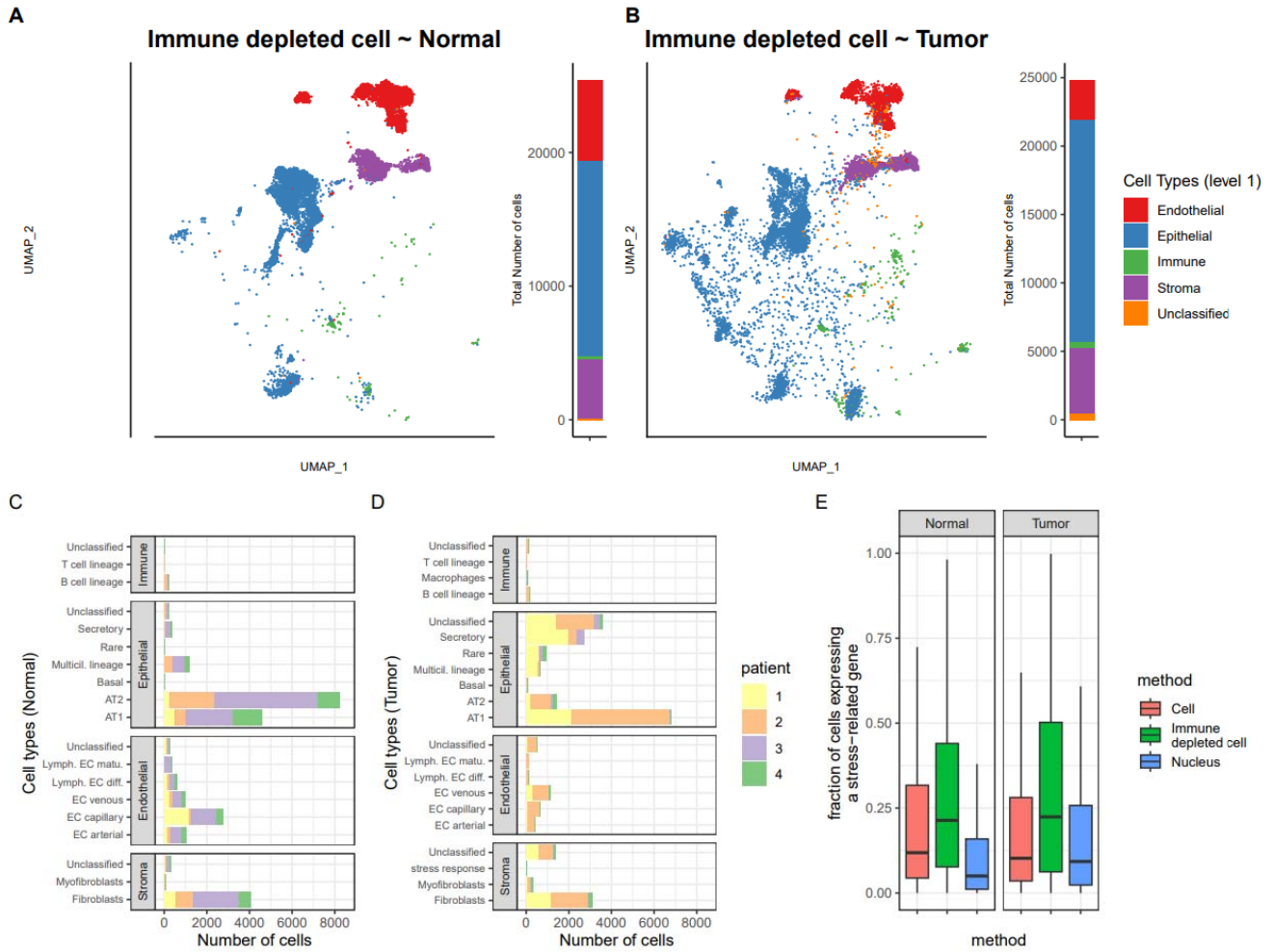**Figure 8** | **The ligand-receptor interactome**. **A.** Scatter plots of ingoing and outgoing interactions per
tissue type and method for common cell types (see methods) among all comparisons. To the right are
the top 10 interacting pathways**. B**: An example of pathway common in cell, rare in nucleus (MHC-I)
with the contribution of the top10 ligand-receptor interacting genes (bar plot to the right). **C:** An
example of pathway rare in cell, common in nucleus (PTPRM) with the ligand-receptor interacting
gene (bar plot to the right).

714

715



**Figure 9 | UMAP representations and cell types annotations (immune depleted cells)** for Normal (**A**) and Tumor (**B**) tissue samples with general cell types (level 1) annotation. To the right of each UMAP, stacked bar plots indicate the proportion of each cell type in the specific dataset. Number of cells in the Normal (**C**) and Tumor (**D**) tissue, per patient. **E**: The percentage of cells expressing a stress-related gene signature as a function of the experimental method and tissue type.