

1

2 **Single-cell analysis of human primary prostate cancer reveals the heterogeneity**
3 **of tumor-associated epithelial cell states**

4

5 Hanbing Song¹, Hannah N.W. Weinstein^{1#}, Paul Allegakoen^{1#}, Marc H. Wadsworth II^{2#},
6 Jamie Xie¹, Heiko Yang³, Felix Y. Feng⁴, Peter R. Carroll⁵, Bruce Wang⁶, Matthew R.
7 Cooperberg⁷, Alex K. Shalek², Franklin W. Huang^{1*}

8

9 1. Division of Hematology/Oncology, Department of Medicine; Helen Diller Family
10 Comprehensive Cancer Center; Bakar Computational Health Sciences Institute; Institute
11 for Human Genetics; University of California, San Francisco, San Francisco, CA 94143,
12 USA

13 2. The Ragon Institute of Massachusetts General Hospital, Massachusetts Institute of
14 Technology and Harvard University, 400 Technology Square, Cambridge, MA 02139,
15 USA; Institute for Medical Engineering and Science (IMES), Department of Chemistry,
16 Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA
17 02139, USA; Koch Institute for Integrative Cancer Research, Massachusetts Institute of
18 Technology, 77 Massachusetts Avenue, Cambridge, MA 02139, USA; Broad Institute of
19 Massachusetts Institute of Technology and Harvard, 415 Main St., Cambridge, MA 02142,
20 USA.

21 3. Department of Urology; Helen Diller Family Comprehensive Cancer Center; University
22 of California, San Francisco, 550 16th Street, 6th Floor, San Francisco, CA 94143, USA

23 4. Departments of Radiation Oncology, Urology, and Medicine, Helen Diller Family
24 Comprehensive Cancer Center; University of California, San Francisco, Helen Diller
25 Family Cancer Research Building, 1450 Third Street, Room 383, San Francisco, CA
26 94143, USA

27 5. Department of Urology; Helen Diller Family Comprehensive Cancer Center; University
28 of California, San Francisco, 550 16th Street, 6th Floor, San Francisco, CA 94143, USA

29 6. Division of Gastroenterology, Department of Medicine, University of California, San
30 Francisco, CA 94143, USA

31 7. Department of Urology; Epidemiology & Biostatistics; Helen Diller Family
32 Comprehensive Cancer Center; University of California, San Francisco, 550 16th Street,
33 6th Floor, San Francisco, CA 94143, USA

34 #These authors contributed equally.

35 * Correspondence: Franklin.Huang@ucsf.edu

36

37

38 **Abstract**

39 Prostate cancer is the second most common malignancy in men worldwide and
40 consists of a mixture of tumor and non-tumor cell types. To characterize the prostate
41 cancer tumor microenvironment, we performed single-cell RNA-sequencing on prostate
42 biopsies, prostatectomy specimens, and patient-derived organoids from localized
43 prostate cancer patients. We identify a population of tumor-associated club cells that may
44 act as progenitor cells and uncover heterogeneous cellular states in prostate epithelial
45 cells marked by high androgen signaling states that are enriched in prostate
46 cancer. *ERG*⁻ tumor cells, compared to *ERG*⁺ cells, demonstrate shared heterogeneity
47 with surrounding luminal epithelial cells and appear to give rise to common tumor
48 microenvironment responses. Finally, we show that prostate epithelial organoids
49 recapitulate tumor-associated epithelial cell states and are enriched with distinct cell types
50 and states from their parent tissues. Our results provide diagnostically relevant insights
51 and advance our understanding of the cellular states associated with prostate
52 carcinogenesis.

53

54 **Introduction**

55 The prostate consists of multiple cell types including epithelial, stromal, and
56 immune cells, each of which has a specialized gene expression profile. The
57 development of cancer from prostate tissue involves complex interactions of tumor cells
58 with surrounding epithelial and stromal cells and can occur multifocally, suggesting that
59 prostate epithelial cells may undergo cellular state transitions towards carcinogenesis¹⁻
60 ⁶. Previous studies on prostate cancer (PCa) molecular changes have focused on
61 unsorted bulk tissue samples, leaving a gap in our understanding of the adjacent
62 epithelial cell states.

63 The classification of prostate epithelial cells has been expanded over the past
64 few years from three types (basal epithelial cells, luminal epithelial cells, and
65 neuroendocrine)^{7,8} to include hillock cells and club cells⁹. The roles of these additional
66 cell types in the prostate are largely unknown. Most PCa are marked by the expansion
67 of malignant cells with luminal epithelial features and the absence of basal epithelial
68 cells. However, to date, the role of additional cell populations beyond the luminal and
69 basal types is not well known.

70 Another underexplored area is the tumor microenvironment changes that occur
71 based on dominant genomic drivers in PCa. PCa tumor cells are driven by a number of
72 oncogenic alterations that include highly prevalent gene fusion events involving ETS
73 family transcription factors, such as *TMPRSS2-ERG* and *ETV1/4/5*^{1,10-12}. Tumor cells
74 without ETS fusion events and non-malignant luminal cells, however, have not been
75 thoroughly characterized on a single-cell level, and uncertainty remains whether ETS
76 fusion events could evoke differential stromal and immune cell responses.

77 To characterize tumor cells and the surrounding epithelial, stromal, and immune
78 cell microenvironment and identify cell states that are associated with tumorigenesis,
79 we performed single-cell RNA-sequencing (scRNA-seq) on PCa samples. Furthermore,
80 we derived *in vitro* organoids from PCa tumor tissues followed by scRNA-seq to chart
81 molecular and cellular changes in prostate cancers from localized PCa patients. We
82 aimed to understand at single-cell resolution the tumor microenvironment and cellular
83 states associated with prostate carcinogenesis.

84

85 **Results**

86 To probe the diversity of cell types and transcriptional states of cells in localized
87 prostate cancer specimens, we isolated single cells from biopsies and surgical resection
88 specimens from men with localized prostate cancer for scRNA-seq (**Supplemental**
89 **Table 1**) using an improved Seq-well single-cell platform¹³. Altogether, 21,743 cells
90 were analyzed and a total of nine different major cell types were identified, marked by
91 specific gene expression profiles (**Methods, Figure 1a,b**).

92 Cell type identification was determined by examining differentially expressed
93 genes (DEGs) as well as signature scores from normal prostate and immune cell
94 population gene sets^{9,14}. Cells in the merged dataset were annotated as epithelial,
95 stromal (endothelial, fibroblast, and smooth muscle), and immune (T-cells, myeloid
96 cells, plasma cells, mast cells, and B-cells) cells based on established marker genes.
97 Epithelial cells (N = 13,322) were identified based upon the expression of luminal
98 epithelial (LE) markers *KLK3*, *ACPP*, and *MSMB*, consistent with LE cells found as the
99 dominant epithelial cell type in PCa samples. Immune cells were identified based on the

100 high-level expression of *PTPRC* in five clusters, of which one cluster was marked by
101 high-level expression of *IL7R*, *CD8A*, and *CD69*, indicating a mixture of both CD8 and
102 CD4 T-cells; a second cluster was characterized by the myeloid cell markers *APOE*,
103 *LYZ*, and *IL1B*^{15–18}. The third *PTPRC*+ cluster represented plasma cells marked by high
104 level expression of *MZB1* and *IGJ*. The other two remaining *PTPRC*+ clusters were
105 annotated as mast cells expressing *CPA3*, *KIT*, and *TPSAB1*, and a population of B-
106 cells expressing *MS4A1*, *CD22*, and *CD79A*. Stromal cells in our dataset consisted of
107 endothelial cells characterized by *CLDN5* and *SELE* expression, fibroblasts expressing
108 *C1S*, *DCN*, and *C7*, and smooth muscle cells expressing *ACTA2*, *MYH11*, and *RGS5*
109 **(Figure 1c)**.

110 As our samples consisted of prostate biopsies (N = 3 patients) and radical
111 prostatectomy (RP) specimens (N = 8 patients), half of which had matched benign-
112 appearing tissue **(Supplemental Table 1)**, we tested whether each sampling strategy
113 captured a similar distribution of different cell types across samples. All major cell types
114 were captured in each sample with epithelial cells comprising the largest population
115 **(Figure 1d)**. No significant difference was found among the three sample types ($p >$
116 0.05 , Mann Whitney U-test) **(Figure 1e)**. We also compared the cell type composition
117 among paired tumor (N = 4), paired normal (N = 4), and RP unpaired tumor tissues (N =
118 4) **(Supplemental Table 1)** and found no significant differences ($p > 0.05$, Mann-
119 Whitney U test). The main cell types identified were validated by SingleR annotation¹⁹
120 **(Supplemental Table 1)**. Furthermore, within each biopsied patient, we tested whether
121 biopsies from the two anatomical regions identified similar cell types and found that all

122 cell types were recovered in each biopsy sample with some sampling differences by
123 anatomical regions (**Supplemental Table 1**).

124

125 **Epithelial cell clusters reveal tumor cells and non-tumor surrounding epithelial** 126 **cell heterogeneity**

127 To identify the transcriptional cell states of epithelial cells associated with
128 prostate cancer, we performed a graph-based clustering analysis and identified 20
129 clusters (**Figure 2a**). We then conducted single-sample gene set enrichment
130 analyses^{20,21} (ssGSEA) using signature gene sets developed previously from single-cell
131 profiling of normal prostates (**Supplemental Table 2**) to determine the major cell
132 subtypes⁹. Clusters with *KRT5*, *KRT15*, *KRT17*, and *TP63* expression (**Figure 2b**) and
133 significantly upregulated basal epithelial (BE) signature scores were identified as BE
134 cells. Given that tumor cells predominantly express LE cell markers such as *KLK2*,
135 *KLK3*, *ACPP*, and *NKX3-1*, clusters with high LE signature scores could be either tumor
136 cells or non-malignant LE cells (**Figure 2b**). BE and LE signature feature plots also
137 revealed a cluster of cells (cluster 5) that we termed other epithelial (OE) cells (**Figure**
138 **2a,c**), with lower BE and LE signatures scores (**Supplemental Figure 1a**), and were
139 characterized by *PIGR*, *MMP7*, and *CP* expression. In previous studies, *PIGR* has
140 shown a role in promoting cell transformation and proliferation²²; *MMP7* may promote
141 prostate carcinogenesis through induction of epithelial-to-mesenchymal transition²³, and
142 serum *CP* levels have been used as a marker in PCa²⁴ (**Figure 2b**).

143 Approximately 50% of PCa tumors from European ancestry patients harbor
144 *TMPRSS2-ERG* fusion events and less frequently harbor other ETS fusion events

145 (*ETV1*, *ETV4*, *ETV5*)²⁵. To identify tumor cells, we tested cells for expression of *ERG*,
146 *ETV1*, *ETV4*, and *ETV5*. *ERG* expression was found upregulated in four clusters
147 (**Figure 2b, Supplemental Figure 1a**); therefore, we annotated these four clusters as
148 *ERG*+ tumor cells that comprised cells from six patients. Other than tumor cell clusters,
149 only endothelial cells showed high-level *ERG* expression. The identity of *ERG*+ tumor
150 cells was further supported by the upregulation of the SETLUR PROSTATE CANCER
151 *TMPRSS2 ERG FUSION UP* gene set signature score in these cells²⁶. Furthermore,
152 STAR-Fusion²⁷ identified potential fusion transcripts of *TMPRSS2-ERG* fusion in two
153 *ERG*+ patients. No clusters with *ETV1*, *ETV4*, or *ETV5* expression were detected
154 (**Supplemental Figure 1b**).

155 To identify tumor cells without ETS fusion events, we tested the LIU PROSTATE
156 CANCER UP and other PCa tumor marker gene set signature scores and identified
157 seven clusters in total with upregulated signature scores of at least one prostate cancer
158 gene set (**Supplemental Figure 1a,b**). Single sample gene set enrichment analysis
159 (ssGSEA) on these 11 clusters also showed at least one prostate cancer gene set that
160 scored in the top 1% of all C2 CGP gene set collection (N = 3,297) (**Supplemental**
161 **Table 3**). Therefore, we classified four clusters with *ERG* expression as *ERG*+ tumor
162 cell clusters and the other seven as *ERG*- tumor cells (**Figure 2c**). All *ERG*- tumor cell
163 clusters expressed tumor marker *SPON2*²⁸ (**Figure 2b**).

164 To validate our tumor cell assignments, we estimated copy number variants
165 (CNV) via InferCNV²⁹, using non-malignant LE cells as a reference. From the CNV
166 estimation visualization (**Supplemental Figure 1c**), we identified significantly different
167 CNV patterns in both *ERG*+ and *ERG*- tumor cells. Non-uniform CNV profiles were

168 detected within *ERG*+ and *ERG*- tumor cell populations, suggesting heterogeneity in
169 both tumor cell populations.

170 While we did not observe a separate neuroendocrine cell cluster, we tested for
171 prostate neuroendocrine (NE) cells^{9,30} using an established NE cell signature gene set⁹
172 and computed the NE signature scores for each epithelial cell. Taking the cells ranking
173 in the top 0.5% NE signature score, we detected 66 putative NE cells within the BE cell
174 population, characterized by *CHGB*, *KRT4*, and *LY6D* expression⁹ (**Supplemental**
175 **Figure 1c,d**).

176 To examine if our annotation method could accurately identify each cell type, we
177 computed the top 10 biomarkers for each cell type (**Figure 2d**). BE cells showed high
178 expression of established basal epithelial cell markers *KRT5*, *KRT15*, and *KRT17*. The
179 top biomarkers in the OE clusters were *PSCA*, *PIGR*, *MMP7*, *SCGB1A1*, and *LTF*, of
180 which *PSCA* is considered to be a prostate progenitor cell marker enriched in PCa^{31–33}
181 and *SCGB1A1* a marker for lung club cells³⁴. *ERG*+ and *ERG*- tumor cells and non-
182 malignant LE cells all showed high expression of luminal markers *KLK3*, *KLK2*, and
183 *ACPP*³⁵. *ERG*+ tumor cells were characterized by expression of *ERG* and tumor
184 markers *PCA3*, *AMAC*, and *TRPM8*^{35–37}; *ERG*- tumor cells were marked by the
185 expression of tumor markers *PCA3* and *TRPM8*^{35–37} (**Figure 2d**).

186 Since most PCa are androgen-responsive with tumor cell proliferation dependent
187 on the activity of the androgen receptor (*AR*)^{36–39}, we tested for androgen
188 responsiveness among the epithelial cell populations and identified LE cells and tumor
189 cells as the most androgen responsive due to significantly higher *AR* signature scores
190 compared to the other epithelial cell types (**Supplemental Figure 1a**). To identify

191 putative prostate cancer stem cells that may contribute to prostate cancer development,
192 we used an adult stem cell signature gene set³⁸ and found that 56.4% of the BE cell
193 population was enriched for the stem cell signature (**Supplemental Figure 1b**).

194 A previous single-cell study of normal human prostate reported two populations
195 of other epithelial cells: hillock cells characterized by *KRT13*, *SERPINB1*, *CLDN4*, and
196 *APOBEC2* expression and club cells characterized by the expression of *SCGB3A1*,
197 *PIGR*, *MMP7*, *CP*, and *LCN2*⁹. While we did not detect a separate hillock cell population
198 within our prostate cancer epithelial cells (**Supplemental Figure 1e**), we did detect a
199 distinct population representing 6.5% of all epithelial cells (872 of 13,322) characterized
200 by expression of *PIGR*, *MMP7*, *CP*, and *LTF* (**Figure 2d**) (FDR $q < 10e-20$). We
201 hypothesized that this epithelial cluster represented club cells that had previously been
202 described in lung³⁴ and normal prostate specimens⁹. To test this hypothesis, we applied
203 a normal prostate club cell signature gene set⁹ and projected the signature onto our
204 epithelial UMAP. We found that cells with high club cell signature scores largely
205 overlapped with this OE cluster (cluster 5) (**Figure 2e**). Furthermore, this cluster was
206 enriched for a lung club cell signature compared to other clusters ($p < 0.001$, Wilcoxon
207 rank sum test) (**Figure 2f**). Based on these results, we annotated this cluster as club
208 cells. We then conducted an ssGSEA analysis on all epithelial cells using the BE, LE,
209 and club cell signatures generated from the DEG profiles (**Supplemental Table 2**). All
210 three cell type signature scores were strongly correlated to the corresponding cell types,
211 supporting our annotation (**Supplemental Figure 1f**).

212

213 **Club and BE cells harbor PCa-enriched LE-like cell states that are upregulated in**
214 **AR signaling**

215 A recent study identified a luminal progenitor cell type in mouse and human
216 prostates characterized by high expression of LE markers *KRT8*, *KRT18*, and other
217 markers including *PSCA*, *KRT4*, *TACSTD2*, and *PIGR*³⁹. In both normal and PCa
218 epithelial cells datasets in our study, we could not identify a single cell type
219 distinguished by the co-expression of *KRT8*, *KRT18*, and *TACSTD2*; however, *PSCA*
220 and *PIGR* were expressed at higher levels in club cells compared to other epithelial cell
221 types (**Supplemental Figure 2a**), indicating that the luminal progenitor cells previously
222 identified are most similar to the club cells in our analysis.

223 Club cells in PCa have not been previously characterized. Since we exclusively
224 captured club cells but not hillock cells in our PCa samples, we hypothesized that club
225 cells play a role in carcinogenesis. To test this hypothesis, we integrated our prostate
226 cancer club cells (Club PCa) with normal club cells from a previous study from healthy
227 controls⁹ (Club Normal) and detected six cell states with distinct transcriptomic profiles
228 (**Figure 3a**) by selecting an optimal resolution to yield stable clusters (**Supplemental**
229 **Figure 2b**). Overall, compared to club cells from normal samples, PCa club cells
230 exhibited downregulation of genes including lipocalin 2 (*LCN2*) and growth-inhibitory
231 cytokine *SCGB3A*^{140,41} and upregulation of *LTF*, *AR*, and *AR* downstream members
232 including *KLK3*, *KLK2*, *ACPP*, and *NKX3-1* (**Figure 3b**), which we hypothesized could
233 be driven by the enrichment of one or more specific club cell states in the PCa samples.

234 For each of the six subclusters, a group of distinctive DEGs was identified
235 (**Figure 3c**) and each subcluster was detected in both Club PCa and Club Normal

236 **(Supplemental Figure 2c)**, of which, Club PCa was significantly enriched in cluster 0
237 by more than three-fold compared to Club Normal ($p < 0.001$, Fisher's exact test (FET))
238 **(Figure 3d)**. This cluster was distinguished by a higher level of expression of *LTF*,
239 luminal markers, and downstream *AR* pathway molecules *KLK2*, *KLK3*, *ACPP*,
240 *PLA2G2A*, and *NKX3-1* **(Figure 3e)**, suggesting a luminal-like and androgen-responsive
241 state³⁹.

242 To test the functional role of this cell state, we performed GSEA analysis using
243 C2 canonical pathways (N = 2,232) **(Supplemental Table 4)** and Hallmark (N = 50)
244 gene set collections on cluster 0 vs other cell states. Among the top significantly
245 upregulated gene sets in cluster 0 was the Hallmark Androgen Response pathway
246 (FDR $q < 10e-5$) **(Figure 3f)**. These results were consistent with the upregulation of
247 downstream *AR* pathway molecules in cluster 0.

248 Next, we tested whether this PCa-enriched cell state represented a luminal-like
249 cell state. We observed higher LE signature scores in cluster 0 compared to other cell
250 states ($p < 0.001$, Wilcoxon rank sum test) **(Figure 3g)**. Specifically, we compared the
251 expression levels of all LE markers among cluster 0, other club cells, and the LE
252 population within the PCa samples, and found that Club cell cluster 0 exhibited higher
253 expression of *KLK2*, *KLK3*, *ACPP*, *NKX3-1*, *KLK4*, *PLA2G2A*, *SPDEF* and *GOLM1* than
254 other club cells **(Figure 3h)**. While *AR* itself was not upregulated in cluster 0
255 **(Supplemental Figure 2d)**, *KLK4*, a regulator of androgen response signaling, was
256 upregulated in this cell cluster⁴².

257 Overall, the population of PCa club cells, compared to normal prostate clubs,
258 was characterized by higher androgen signaling and an enrichment of an *LTF^{high}* and
259 *NKX3-1^{high}* luminal-like cell state (**Figure 3i**).

260 The finding of a luminal-like club cell state led us to investigate if a similar cell
261 state existed in the BE cell population of prostate cancer samples. Therefore, we
262 integrated BE cells in the PCa samples (BE PCa) with BE cells from normal samples
263 (BE Normal) and identified nine cell states (**Figure 4a, Supplemental Figure 3a**) with
264 distinctive DEGs (**Supplemental Figure 3b**). While all nine cell states were represented
265 in both BE PCa and BE Normal cells (**Figure 4b**), BE PCa was found to be significantly
266 enriched in cluster 6 (31.8% vs 0.2%, PCa vs Normal) while BE Normal was enriched in
267 cluster 4 (0.8% vs 15.9%, PCa vs Normal) (FDR $q < 10e-5$, FET; **Figure 4b**,
268 **Supplemental Figure 3c**). This BE cluster 6 was marked by higher expression of
269 downstream *AR* pathway members *KLK3*, *KLK2*, and *ACPP* (**Supplemental Figure**
270 **3b**). Compared to other BE cells in the PCa samples, BE cluster 6 also showed
271 significant upregulation of *AR* ($p < 0.01$, Wilcoxon rank sum test, **Supplemental Figure**
272 **3d**). Among the top significantly upregulated gene sets were the Hallmark Androgen
273 Response pathway within the Hallmark gene set collection (**Supplemental Table 4**), as
274 well as androgen response pathways, estrogen pathways, the insulin signaling pathway,
275 and Kegg pathways in cancer within the C2 CP gene set collection (FDR $q < 0.1$,
276 Wilcoxon rank sum test) (**Figure 4d**)^{42–44}. As *AR* pathway members were among the
277 top biomarkers for cluster 6 (**Figure 4e**), we hypothesized that BE cluster 6 may
278 represent an intermediate BE/LE cell state, even though it did not cluster separately
279 from other BE cells on the epithelial cell UMAP (**Figure 4f**). Therefore, we compared the

280 expression levels of LE markers in BE cluster 6 and found that BE cluster 6 was
281 upregulated in multiple LE markers compared to other BE cell states (**Figure 4g**),
282 though at lower levels compared to the PCa LE cell population. Moreover, we found that
283 BE cluster 6 was significantly upregulated in the Hallmark Androgen Response
284 signature ($p < 0.001$, Wilcoxon rank sum test) and LE signature score (**Figure 4h**),
285 supporting that this cell state may be a luminal-like state associated with prostate
286 cancer.

287 Similarly, we identified eight cell states within the integrated LE dataset
288 (**Supplemental Figure 3e**). Unlike BE and club cells, we observed a clear separation
289 between LE PCa and LE Normal (**Supplemental Figure 3e**). LE PCa was significantly
290 enriched in four cell states and LE Normal significantly enriched in two ($p < 0.001$, FET)
291 (**Supplemental Figure 3f**). Cluster 5 was marked by co-expression of club cell markers
292 such as *PIGR*, *MMP7* and *CP*, suggesting an intermediate population of LE and club
293 cells. Cluster 1 was characterized by the overexpression of the *AR*-regulated gene
294 *TMEFF2* and insulin-like growth factor *IGFBP5* compared to other cell states, and
295 cluster 2 was upregulated in *AR* expression (**Supplemental Figure 3g**).

296 We then tested if the PCa-enriched cell states in BE and club cells (Club cell
297 cluster 0 and BE cluster 6) could be distinguished from other cell states in the
298 differentiation trajectory. Given that BE cells showed upregulated stem cell signature
299 scores (**Supplemental Figure 1a**), we used BE cells as the starting point and plotted
300 the diffusion pseudotime trajectory on the partition-based graph abstraction (PAGA)
301 initialized embedding with a list of cell type specific markers as well as proliferation
302 markers *MKI67* and *TOP2A* (**Supplemental Figure 4a,b**). We observed that *KRT5*+ BE

303 cells gave rise to all other epithelial cells and tumor cells, with tumor cells and LE cells
304 (*KLK3+*) appearing later than club cells (*PIGR+*, *LTF+* and *PSCA+*) in the pseudotime
305 trajectory (**Supplemental Figure 4c**), consistent with a previous analysis⁹. We ran
306 Monocle3⁴⁵ to compute the pseudotime trajectory for PCa club cells (**Supplemental**
307 **Figure 4d,e**). Club cells with higher LE signature scores were more differentiated in
308 pseudotime (**Supplemental Figure 4e**). This finding was further supported by
309 increasing expression levels of LE markers *ACPP* and *KLK3* along the trajectory
310 compared to club cell markers (**Supplemental Figure 4f,g**), suggesting that LE-like
311 club cells in PCa samples could be transitioning to LE cells or tumor cells.

312

313 **Integrated epithelial cell analysis reveals upregulated AR signaling in PCa** 314 **samples**

315 As PCa samples in this study included four paired tumor and normal samples, we
316 tested if PCa-enriched cell states in BE, LE, and club cells were enriched in the
317 surrounding epithelial cells of the PCa biopsies and in radical prostatectomy tissue
318 samples containing tumor cells. We compared the percentage composition of each BE
319 and club cell state within all BE and club cells in all five sample types respectively
320 (Normal, biopsy, RP paired tumor, RP paired normal, and RP unpaired tumor). The
321 PCa-enriched cell states of Club cell cluster 0 and BE cluster 6 were similarly
322 represented in the four paired tumor and normal samples ($p = 0.43$, Mann-Whitney U
323 test).

324 To identify the overall epithelial cell transcriptional programs in PCa samples, we
325 integrated all PCa epithelial cells (Epithelial PCa) with prostate epithelial cells from

326 normal healthy controls (Epithelial Normal)⁹ (**Figure 5a**). We identified differentially
327 expressed genes between tumor and normal samples across all three major types of
328 epithelial cells (LE, BE, and club cells). We found ATF transcription factors *FOS* and
329 *JUN*, members of the EGFR pathway that mediate gene regulation in response to
330 cytokines and growth factors⁴⁶, and prostate acid phosphatase (*PSAP*)⁴⁷ as commonly
331 upregulated across these cell types (**Figure 5b**). However, the DEGs could not be
332 recapitulated when comparing between paired tumor and normal samples
333 (**Supplemental Table 5**), suggesting that compared to normal prostate samples,
334 epithelial cells in the paired normal tissues were more similar to those from paired tumor
335 tissues taken from different anatomical regions within the same radical prostatectomy
336 specimen. Since the two PCa-enriched cell states in BE and club cells showed
337 upregulated *AR* signaling compared to other BE or club cells respectively, we tested *AR*
338 expression in the integrated dataset and found that in PCa epithelial cells, 21.4% of BE
339 cells (458 of 2,145 cells), 28.6% of club cells (249 of 872), 52.7% of LE cells (2,974 of
340 5,647 cells) and 43.2% of tumor cells (1,993 of 4,658 cells) were *AR*+, in which
341 significantly higher percentages of PCa BE, LE, and club cells were *AR*+ compared to
342 the same cell types from normal samples ($p < 0.001$, FET) (**Figure 5c**). We also
343 computed the Hallmark Androgen Response pathway signature scores for all cells and
344 found that the three major epithelial cell types in PCa samples were all upregulated in
345 *AR* signaling compared to normal samples ($p < 0.001$, Wilcoxon rank sum test) (**Figure**
346 **5c**).

347 To validate the two PCa-enriched epithelial cell states we identified in BE and
348 club cells and test their correlation with upregulated *AR* signaling, we ran ssGSEA on all

349 BE and club cells on the Hallmark Androgen Response pathway. The *AR* signature
350 score of BE cells was only significantly positively correlated to BE cluster 6 (information
351 coefficient (IC) = 0.499, FDR $q < 1e-5$), and the *AR* signature score in club cells was
352 significantly positively correlated to Club cell cluster 0 (IC = 0.385, FDR $q < 1e-5$)
353 (**Figure 5d**). Furthermore, to test if this correlation between a PCa-enriched cell state
354 and *AR* signaling could be replicated in other PCa datasets, we projected all BE and
355 club cell states across the TCGA²⁵ (N = 499) and SU2C⁴⁸ (N = 266) castration resistant
356 prostate cancer (CRPC) bulk RNA-seq datasets (**methods**). In both bulk RNA-seq
357 datasets, *AR* signature scores were positively correlated with BE cluster 6 (IC = 0.756,
358 FDR $q < 1e-5$) and Club cell cluster 0 (IC = 0.233, FDR $q < 1e-5$) (**Figure 5e**),
359 supporting our identification of cell states within BE cells and club cells that were more
360 androgen responsive and associated with prostate cancer.

361

362 **Transcriptomic profiles of *ERG*⁺ tumor cells are patient-specific while *ERG*⁻** 363 **tumor cells overlap with surrounding LE cells**

364 While *ERG*⁺ tumor cells clustered separately from non-malignant LE cells, *ERG*⁻
365 tumor cells resided more closely to non-malignant LE cells (**Figure 2c**). To investigate
366 this further, we first analyzed the sub-structure of *ERG*⁺ and *ERG*⁻ tumor cells
367 separately to identify distinct underlying cell states (**Figure 6a,b**). *ERG*⁺ tumor cells
368 clustered in a patient-specific manner, whereas no such pattern was seen for *ERG*⁻
369 tumor cells as most *ERG*⁻ tumor cell states were comprised of more than one patient
370 (**Figure 6c**).

371 One possibility for the different distribution patterns between *ERG*⁺ and *ERG*⁻
372 tumor cells is that *ERG*⁺ tumor cells for each patient represented a distinctive cell state
373 driven by a dominant oncogenic alteration, though no such distinction was seen in *ERG*⁻
374 tumor cells, suggesting more overlapping cell states between *ERG*⁻ tumor cells and
375 adjacent non-malignant LE cells. To test this hypothesis, we integrated *ERG*⁺ tumor
376 cells and *ERG*⁻ tumor cells separately with LE cells and performed sub-clustering
377 analyses. Overall, we found 1,244 genes significantly varied between *ERG*⁺ tumor cells
378 and LE cells (FDR $q < 0.01$, Wilcoxon rank sum test), while only 314 genes were
379 significantly varied between *ERG*⁻ tumor cells and LE cells (FDR $q < 0.01$, Wilcoxon
380 rank sum test). Fourteen and seventeen cell states were recovered in the *ERG*⁺ and
381 *ERG*⁻ integrated datasets, respectively (**Supplemental Figure 5a-b**). We observed a
382 clear separation between *ERG*⁺ tumor cells and non-malignant LE cells while *ERG*⁻
383 tumor cells were not clearly distinguishable from non-malignant LE cells in the analysis
384 (**Figure 6d**). From the cell state composition comparison, we observed three cell states
385 with more than 400 cells each that were almost exclusively detected in the *ERG*⁺ tumor
386 cells, with each cell state largely attributed to one specific patient (**Supplemental**
387 **Figure 5a**). In contrast, no such patient specificity was observed for *ERG*⁻ tumor cells
388 (**Figure 6e**) (**Supplemental Figure 5b**). In our dataset, *ERG*⁺ tumor cells were
389 predominantly found in tumor samples while *ERG*⁻ tumor cells were found in paired
390 tumor and normal samples (**Supplemental Figure 5c**). Using the DEGs between *ERG*⁺
391 and *ERG*⁻ tumor cells (**Supplemental Figure 5d**), we generated signature gene sets for
392 both types of tumor cells and tested if the signatures of *ERG*⁺ and *ERG*⁻ tumor cells
393 generated from this dataset were correlated with *TMPRSS2-ERG* fusion status in

394 TCGA²⁵ and SU2C⁴⁸ castration resistant prostate cancer (CRPC) bulk RNA-seq
395 datasets. *TMPRSS2-ERG* fusion status was significantly positively correlated with an
396 *ERG*+ tumor cell signature score in both datasets (TCGA: information coefficient (IC) =
397 0.673, FDR $q < 1e-5$; SU2C: IC = 0.407, FDR $q < 1e-5$) and the absence of *TMPRSS2-*
398 *ERG* fusion was significantly correlated with *ERG*- tumor signature scores (TCGA: IC =
399 -0.554, FDR $q < 1e-5$; SU2C: IC = -0.211, FDR $q < 1e-5$) (**Figure 6f**). These results
400 supported the tumor cell signatures and our use of *ERG* expression as a classification
401 in annotating tumor cells.

402 Furthermore, we compared the numbers of *ERG*+ tumor cell and *ERG*- tumor
403 cells in each patient. Tumor cells in five patients were over 90% *ERG*- and over 90%
404 *ERG*+ in two patients (**Supplemental Figure 5e**) Tumor cells in four patients harbored
405 both types of tumor cells. Using non-tumor epithelial cells as reference, we found
406 significantly different CNV profiles from the reference for each patient, further validating
407 our tumor cell identification (**Supplemental Figure 5f**). For our downstream analyses,
408 we classified patients based on ERG status by annotating the five patients with almost
409 exclusive *ERG*- tumor cells as *ERG*- patients and the other six patients (exclusive
410 *ERG*+ tumor cells and mixtures) as *ERG*+ patients.

411

412 **T-cell and stromal cell analysis reveals common signaling in *ERG*- patients**

413 The transcriptional differences between *ERG*+ and *ERG*- tumor cells suggested
414 that they might give rise to differential responses in the tumor microenvironment. To
415 identify tumor-related immune cells and whether specific immune cell types were
416 differentially enriched in either *ERG*+ or *ERG*- samples, we analyzed the T-cell

417 population and identified CD4 and CD8 T-cells, regulatory T-cells (Treg), and NK cells
418 based on differentially expressed genes (**Figure 7a**). We then stratified the T-cell
419 populations based on *ERG* status and found two CD4 T-cell clusters that were
420 differentially enriched. Between the two CD4+ T-cells we identified, CD4 T-cell cluster 1
421 was enriched in *ERG*+ patients with a 2.73 fold difference (20.5% vs 7.5%) (**Figure 7b**)
422 and was characterized by a higher level expression of immune response regulators
423 including AP-1 transcription factors⁴⁹ *FOSB* (log₂FC = 1.79, FDR q = 5e-30), *FOS*
424 (log₂FC = 1.78, FDR q = 6.2e-26) and *JUN* (log₂FC = 1.55, FDR q = 5.5e-22). CD4 T-
425 cell cluster 2 was enriched in *ERG*- patients with a 5.6 fold change (9.5% vs 1.7%)
426 (**Figure 7b**) ($p < 0.001$, Fisher's exact test) and was marked by higher expression of
427 *DUSP4* (log₂FC = 1.30, FDR q = 1.4e-20) and *CXCR6* (log₂ fold change (log₂FC) =
428 1.31, FDR q = 1.5e-22), which was previously shown to be expressed in the type-1
429 polarized T-cell subset and to contribute to tumor progression⁵⁰. We noted that the
430 DEGs between the two T-cell clusters were consistent with the DEGs identified between
431 *ERG*+ and *ERG*- tumor cells, with *FOSB*, *FOS*, and *JUN* overexpressed in *ERG*+ tumor
432 cells while *CXCR6* and *DUSP4* were overexpressed in *ERG*- tumor cells
433 (**Supplemental Figure 5d**). No other T-cell populations (CD8 T-cells, Treg, and NK
434 cells) showed a significant difference in cell type abundance between *ERG*+ and *ERG*-
435 patients.

436 Similarly, we stratified the stromal population based on the *ERG* status of
437 patients and identified three distinct clusters consistent with endothelial cells, smooth
438 muscle cells, and fibroblasts (**Figure 7c**). Of these three stromal cell types, fibroblasts
439 showed an enrichment in *ERG*+ patients ($p < 0.001$, FET).

440 To test if the differences between *ERG*- and *ERG*+ tumor cells could potentially
441 drive distinct and common stromal and immune responses, we ran independent GSEA
442 analyses between *ERG*- and *ERG*+ tumor cells, CD4 T-cells and stromal cells and
443 computed the intersection of significantly upregulated gene sets in *ERG*- patients (FDR
444 $q < 0.1$). Fourteen upregulated gene sets were identified that were commonly
445 upregulated in *ERG*- tumor cells, CD4 T-cells and stromal cells ($p < 10e-20$, multi-set
446 intersection exact test⁵¹) (**Figure 7f**). However, we did not detect any common pathway
447 changes in the other epithelial populations (**Figure 6g**). The fourteen common
448 upregulated gene sets in *ERG*- patients included Reactome PD-1 and Reactome
449 interferon gamma signaling (**Figure 7g**), which have both been reported to be
450 upregulated in advanced prostate cancers^{52,53}. Within these two gene sets, we found
451 that *ERG*- patient-enriched CD4 T-cells, tumor cells, and stromal cells showed
452 significantly higher expression of a family of HLA genes compared to *ERG*+ cell
453 populations ($p < 0.05$, Wilcoxon rank sum test) (**Figure 7h**). Within the T-cells, while
454 there was no difference in the cell composition of CD8 T-cells based on *ERG* status, the
455 *ERG*- CD8 T-cell population was also found to be upregulated in the Reactome PD-1
456 and Reactome interferon gamma signaling signatures (FDR $q < 0.1$, **Supplemental**
457 **Table 4**). To test if *ERG*- tumor cell-associated CD4 and CD8 T-cells could represent a
458 distinct immune cell niche, we tested a series of exhausted, cytotoxic markers⁵⁴ as well
459 as genes in the PD-1 and Reactome interferon gamma signaling pathway
460 (**Supplemental Table 6**). We found that *ERG*- CD4 T-cells were significantly
461 upregulated in exhausted T-cell markers *PDCD1* ($\log_2FC = 0.52$, $p < 0.01$, Wilcoxon
462 rank sum test) and *CTLA4* ($\log_2FC = 1.79$, $p < 0.001$, Wilcoxon rank sum test) and

463 cytotoxic markers *GZMA* (log2FC = 1.54, $p < 0.001$, Wilcoxon rank sum test) and *GZMB*
464 (log2FC = 1.09, $p < 0.05$, Wilcoxon rank sum test) compared to *ERG+* CD4 T-cells.
465 Additionally, *ERG-* CD8 T-cells were upregulated in exhausted T-cell markers *HAVCR2*
466 (log2FC = 0.68, $p < 0.05$, Wilcoxon rank sum test) and *LAG3* (log2FC = 0.86, $p < 0.001$,
467 Wilcoxon rank sum test) compared to *ERG+* CD8 T-cells (**Supplemental Figure 6a,b**).
468 These results suggested that CD4 and CD8 T-cells associated with *ERG-* tumor cells
469 represented a more exhausted and cytotoxic phenotype. Then, using CD4 phenotype
470 markers from a previous analysis⁵⁵, we tested the frequency of expression for these
471 markers in both *ERG+* and *ERG-* CD4 T-cells and found a significantly higher proportion
472 of *CCR7+* central memory CD4 T-cells, *CD69+* activated CD4 T-cells, *GZMB+* cytotoxic
473 CD4 T-cells, and *TOX+* exhausted CD4 T-cells⁵⁵ associated with *ERG-* patients
474 (**Supplemental Figure 6c**).

475 After T-cells, myeloid cells comprised the second largest immune cell population.
476 Annotation of the myeloid cell population with SingleR¹⁹ yielded four cell types:
477 neutrophils, eosinophils, macrophages, and monocytes (**Supplemental Table 7**;
478 **Supplemental Figure 7a-b**). Within the myeloid cell population, we did not detect any
479 significant composition differences in monocytes or macrophages between RP paired
480 tumor and paired normal samples or between *ERG+* and *ERG-* patients ($p > 0.05$, FET)
481 (**Supplemental Figure 7c**).

482 To investigate the subtypes of monocytes and macrophages that are associated
483 with tumor-related responses, we identified monocytes and macrophages with high
484 expression of cell cycle markers *MKI67* and *TOP2A*, indicating a cluster of proliferating
485 myeloid cells (**Supplemental Figure 7d**) that we termed *MKI67+* myeloid cells.

486 Monocytes were further classified by the expression of *CD14* (**Supplemental Figure**
487 **7d**). Within the macrophage population, we used previously established signatures^{56–59}
488 of dichotomous phenotypes to classify macrophages into M0, M1, and M2 types, of
489 which M1 macrophages have been described as pro-inflammatory and M2
490 macrophages as anti-inflammatory and associated with tumor progression⁶⁰. We
491 computed the signature scores of M1 and M2 macrophages and annotated the two
492 subtypes accordingly, based on signature scores as well as M1 specific markers, such
493 as *IL1A*, *CXCL3*, and *PTGS2*, and M2 specific markers, such as *ARG1*, *CCL22*, and
494 *FLT1*. Neither M1 nor M2 macrophages were clustered separately from normal M0
495 macrophages, consistent with a previous analysis of macrophage subtypes⁵⁸
496 (**Supplemental Figure 7d,e**).

497 A recent study on macrophages categorized macrophages into resident tissue
498 macrophages enriched in normal tissues (RTM) and tumor associated macrophages
499 (TAM) enriched in tumor tissues, which did not fit the M1/M2 phenotypes^{61,62}. We did
500 not detect RTMs within the PCa samples (**Supplemental Figure 7f**). In contrast, TAMs
501 were described as either *C1QC+* or *SPP1+*. These TAMs were reported to derive from
502 *FCN1+* monocyte-like macrophages, which was consistent with the detection of *FCN1*
503 in a cluster of PCa myeloid cells where we saw a mixture of monocytes and
504 macrophages (**Supplemental Figure 7f**). In total, 713 TAMs were identified but no
505 significant difference in composition was detected between paired tumor and normal
506 samples (77.9% vs 69.0%, $p = 0.58$, FET) (**Supplemental Figure 7g**).

507 Another group of tumor-associated myeloid cells termed myeloid-derived
508 suppressor cells (MDSC) has been characterized with roles in inflammation,

509 establishing host immune homeostasis, and driving castration resistance in prostate
510 cancer^{63–66}. These MDSCs can inhibit anti-tumor reactivity of T-cells and NK-cells and
511 the enrichment of MDSCs was correlated with tumor progression and worse clinical
512 outcomes⁶⁷. Two types of MDSCs have been described: monocytic MDSC (M-MDSC)
513 characterized by high expression of *CD11* and *CD14* and low expression of *HLA* and
514 *CD15* and granulocytic or polymorphonuclear MDSC (PMN-MDSC) characterized by
515 high expression of *CD11* and *CD15* and low expression of *CD14*. To test for the
516 presence of these MDSCs in our PCa samples, we used the co-expression of these
517 markers and identified 137 M-MDSCs within the 790 *CD14*⁺ monocytes and 11 PMN-
518 MDSCs within 974 *CD14*⁻ monocytes (**Supplemental Figure 7g**). M-MDSCs were
519 enriched in the paired tumor samples compared to paired normal (19.9% vs 3.6% of
520 total monocytes, $p = 0.0035$, FET).

521

522 **Prostate cancer organoids recapitulate epithelial cell types with uniquely** 523 **expanded cell states in BE and club cells**

524 To develop models to examine the cellular state heterogeneity revealed by
525 single-cell analysis and to determine if we could reconstitute and propagate prostate
526 cancer-associated club cells, we used established methods^{68,69} to generate localized
527 prostate cancer organoids from single cells from six patients who underwent radical
528 prostatectomies (four patients included in the tissue sample dataset) and characterized
529 them using scRNA-seq within three passages (**Figure 8a**). PCA-based clustering of
530 organoid samples yielded 23 clusters from a total of 15,073 cells. We identified a total of
531 six epithelial cell types with distinctive DEGs, based on the cell type signatures we

532 generated from the PCa tissue samples and the established signatures from normal
533 samples (**Supplemental Table 2**) (**Figure 8a**). The epithelial cell types included BE
534 cells characterized by high expression of *DST*, *KRT15*, *KRT5*, *KRT17*, and *TP63*, club
535 cells characterized by *PIGR*, *MMP7*, *CP*, and *CEACAM6*, hillock cells, consistent with
536 those in normal prostates showing high level expression of *KRT13*, *CLCA4*, and
537 *SERPINB3*, a mesenchymal stem cell (MSC) population expressing known MSC
538 markers⁷⁰⁻⁷² *LAMC2*, *VIM*, *MMP1*, and *KLK7* and a population with high level
539 expression of cell cycle markers *MKI67* and *TOP2A* termed *MKI67+* epithelial cells
540 (**Supplemental Figure 8a**). Notably, within these early passage organoids we identified
541 a tumor cell population expressing a high level of LE cell markers (*KLK3*, *KLK2*, and
542 *ACPP*) and tumor markers (*PCA3*, *TRPM8*, and *ERG*) (**Supplemental Figure 8a**). Cell
543 type annotation was supported by ssGSEA, which showed that the MSC population was
544 upregulated in the MSC signature gene set developed from a previous analysis⁷¹ and
545 that the *MKI67+* cluster was upregulated in a KEGG cell cycle signature indicating
546 proliferating cells (**Supplemental Figure 8b**). The identification of tumor cells was
547 further validated by InferCNV²¹ estimation (**Supplemental Figure 8c**). To validate our
548 recovery of the cell type diversity in the organoids, we performed immunofluorescence
549 staining for *KRT8+* luminal and *KRT5+* basal cells (**Figure 8b**). We validated club cell
550 proliferation *in vitro* by staining for *SCGB1A1*, an established club cell marker in the
551 lung and prostate⁹, and lactoferrin (*LTF*), which was upregulated in the PCa club cells
552 identified by scRNA-seq (**Figure 8b**).

553 To test the fidelity of the organoids as models for tumor tissues, we integrated
554 the cells in the early-passage (P0-P3) organoid samples (N = 10,990) with the epithelial

555 cells from the four RP specimens from which the organoids were derived (N = 8,719)
556 (**Figure 8c**). Compared to PCa tissue samples, LE cell markers or signature scores
557 could not identify a distinctive LE cell cluster in the organoid samples (**Supplemental**
558 **Figure 8b**), consistent with a previous study that LE cells were rarely captured in *in vitro*
559 organoid cultures analyzed by scRNA-seq⁷³. For the four patient-derived organoids,
560 only a small number of tumor cells were captured compared to the parent tissues
561 (tissue samples vs organoids, 34.11% vs 0.11%). However, hillock cells, MSCs and a
562 population of *MKI67+* epithelial cells were exclusive to the organoid samples and were
563 not observed in PCa tissue samples (**Figure 8d**).

564 As BE and club cells were the two primary overlapping cell types between tissue
565 and organoid samples (representing 11.9% and 29.0% of all cells, respectively, in the
566 organoid samples), we took the subset of BE cells and club cells in tissue and organoid
567 samples from the integrated dataset and computed the DEGs. BE markers *KRT5*, *DST*,
568 and *KRT15* were expressed in BE populations from tissue and organoids and club cell
569 genes *MMP7*, *LCN2*, and *CP* were expressed in both club cell populations (**Figure 8e**),
570 suggesting similarities between tissue and organoid BE and club cells.

571 We then investigated BE and club cell populations by integrating organoids with
572 tissue samples, respectively, to identify cell state differences in the organoid samples.
573 We identified nine clusters in the integrated BE cell dataset with distinctive groups of
574 DEGs (**Supplemental Figure 8d**). Compared to BE cells in PCa tissue samples,
575 significantly higher percentages of BE cells in organoids expressed *KRT6A* (organoid vs
576 tissue, 77.4% vs 0.56%, $p < 0.001$, FET), *KRT14* (organoid vs tissue, 71.2% vs 18.6%,
577 $p < 0.001$, FET), and *KRT23* (organoid vs tissue, 78.8% vs 20.2%, $p < 0.001$, FET)

578 **(Supplemental Figure 8e)**, suggesting that BE cells in organoid samples may be more
579 representative of a progenitor cell state.

580 Similarly, when analyzing the organoid club cells with club cells from PCa
581 tissues, we identified a total of eight clusters with distinctive DEGs **(Supplemental**
582 **Figure 8f)** and observed an expansion of cell states in the organoid samples **(Figure**
583 **8f)**. Among the eight clusters, five were predominantly comprised of organoid club cells,
584 while club cells from prostate tissue were only found in clusters 3, 4 and 7. By
585 comparing the expression levels of the top differentially expressed genes for these three
586 clusters split by tissue and organoid club cells, we found that in cluster 3, hillock cell
587 marker *KRT13* was expressed in tissue and organoid club cells, suggesting an
588 intermediate hillock-club cell state. In cluster 4, PCa club cell marker *PIGR* was
589 detected in 47% of organoid club cells (16 of 34) and 71% of tissue club cells (325 of
590 653). *LTF* was expressed in 15% of organoid club cells (5 of 34) compared to 50%
591 tissue club cells (326 of 653), suggesting that *LTF* may be a PCa tissue-specific club
592 cell marker. In contrast, the top DEGs for cluster 7 included LE markers such as *ACPP*,
593 *NKX3-1*, *KLK2* and *KLK3*, consistent with the profile of the previously-identified PCa-
594 enriched club cell state **(Figure 8g)**. In cluster 7, we observed approximately 20% of
595 organoid club cells expressing at least one LE cell marker. This cluster scored higher for
596 the PCa-enriched club cell state compared to all other clusters of organoid club cells,
597 suggesting that PCa-enriched club cell states were recapitulated in organoid samples.
598 Overall, we found that organoid samples harbored cell states found in tumor tissues and
599 an enrichment of progenitor-like cell states and intermediate cell states. The plasticity of
600 these organoid-enriched cell states within BE and club cells suggests that *in vitro*

601 organoid models may provide useful models to study cell state differences and identify
602 lineage relationships to tumorigenesis.

603

604 **Discussion**

605 Studies of localized prostate cancer have been extensively performed with bulk
606 RNA-seq and WES/WGS approaches that have provided key insights into the molecular
607 features of prostate cancer^{9,12,63–66}. Here, we performed single-cell analyses of localized
608 PCa biopsies and radical prostatectomy specimens to characterize the heterogeneity of
609 tumor cells and subpopulations of epithelial cells, stromal cells, and tumor
610 microenvironments.

611 Of note, we identified a distinctive epithelial cell population of club cells that has
612 not been previously observed in human prostate cancer samples. While club cells have
613 been noted in normal prostates^{9,77,78}, a population of club cells associated with prostate
614 cancer suggests they may play a previously unappreciated role in carcinogenesis.
615 Recent studies have identified a progenitor-like *CD38*^{low}/*PIGR*^{high}/*PSCA*^{high} luminal
616 epithelial cell sub-population with regenerative potential^{39,78,79}. Based on the similarity of
617 highly expressed genes including *PIGR*, *MMP7*, *CP*, and *LTF*, we believe those cells
618 are consistent with their identity as club cells. In our analysis, prostate cancer club cells
619 are characterized by the markedly lower expression of *SCGB3A1* and *LCN2* compared
620 to club cells from normal healthy controls⁹. Based on our gene signature analyses, our
621 results suggest that PCa club cells are more androgen responsive overall and harbor a
622 highly androgen-responsive cell state that may be a potential progenitor cancer cell or

623 function to support the overall androgen responsive cellular milieu of prostate
624 cancer^{80,81}.

625 *SCGB3A1*, a marker for club cells, was one of the top downregulated genes in
626 prostate cancer club cells compared to club cells from normal healthy control prostates.
627 *SCGB3A1* may play a tumor suppressor role in a number of cancers including breast,
628 prostate, and lung as its expression has been noted to be markedly lower in cancer
629 tissues compared to normal tissues⁸². We speculate that prostate club cells in the
630 normal epithelia may play a tumor suppressor role through secretion of *SCGB3A1*
631 which is then downregulated in concert with prostate cancer progression, as marked by
632 our finding of *SCGB3A1*^{low} club cells in prostate cancer tissues that can be propagated
633 in organoids. We did not find a distinct population of hillock cells in prostate cancer
634 tissues so it is possible that hillock cells may be depleted in prostate cancer
635 progression.

636 Consistent with other cancer single cell studies in which tumor cells cluster
637 separately, *ERG*⁺ tumor cells clustered separately by patient from non-malignant
638 epithelial clusters^{14,54,83–85}. However, our analysis of *ERG*⁻ tumor cells unexpectedly
639 found that *ERG*⁻ tumors did not cluster by patient and we observed a shared
640 heterogeneity for *ERG*⁻ tumor cells with non-malignant luminal cells.

641 Treating prostate cancer with immune checkpoint inhibition has had limited
642 efficacy to date and these therapies have largely focused on advanced castration-
643 resistant tumors^{43,44,86–90}. Our single-cell analysis reveals new insights into the tumor
644 immune microenvironment of localized prostate cancer based on *ERG* status. We
645 hypothesized that *ERG*⁻ tumor cells might evoke similar tumor microenvironment

646 responses and found common transcriptional pathways that were upregulated in the
647 tumor, stroma, and CD4 T cell populations of *ERG*- patients, including the PD-1 and
648 interferon gamma signaling pathway, suggesting that *ERG*- tumor cells may give rise to
649 a distinct immune cell niche and tumor microenvironment.

650 We note a potential limitation of our analysis in the identification of *ERG*- tumor
651 cells as we also found evidence for *ERG*- tumor cells in paired grossly normal
652 specimens. This could be attributed to tumor cells also being present in the seemingly
653 normal tissues from radical prostatectomy specimens^{14,85,91,92}. Analysis of somatic
654 mutations or structural variants on a single-cell level will contribute to the identification
655 of *ERG*- tumor cells and inform our understanding of tumor heterogeneity.

656 Furthermore, we showed that *in vitro* organoid cultures grown from tumor
657 specimens can recapitulate cell states found in tumor tissues. We identified a number of
658 new cell types that emerged in the organoid samples including hillock cells, MSC and
659 *MKI67*+ epithelial cells. The mechanisms by which hillock cells can propagate in
660 organoid cultures but not be found in the localized tumor tissue specimens are still to be
661 delineated. An expansion of cell states in BE and club cells in the organoids suggests a
662 broader view for their capacity for cell state transitions. Our results suggest that prostate
663 cancer epithelial organoids harbor many major cell types from tissue and provide a
664 useful model to investigate cell state plasticity in the context of selective pressures and
665 genetic perturbations. However, in contrast to previous studies on organoids generated
666 from prostate samples, we did not observe a distinctive *NKX3-1+/KLK3+/AR+* luminal
667 cell population^{68,93,94}. This might be due to a limitation of detection using single cell
668 sequencing technology or that we could not robustly grow differentiated luminal cells⁷³.

669 Comparing epithelial cells from PCa samples with those from normal healthy
670 controls revealed distinct high androgen-signaling cell states that were enriched in PCa
671 samples. We found that epithelial cells from PCa tissues were generally upregulated in
672 *AR* signaling. Given our identification of shared luminal-like, highly androgen-responsive
673 cell states across basal and club cell populations, we posit that these cell types may be
674 primed for tumor cell transformation and may also promote prostate tumorigenesis.
675 Further studies with lineage tracing and dissection of single cell somatic alterations
676 within these specific cell states will be informative for further characterization of their
677 potential tumorigenic roles. The identification of a tumor-associated club cell population
678 raises the possibility that these cells contribute to the interactions between tumor cells
679 and their surrounding epithelial microenvironment. Furthermore, our analyses identify
680 cell type specific signature gene sets within prostate cancer samples that should
681 contribute to a more precise and thorough classification of cells during prostate
682 carcinogenesis. In summary, we provide a single-cell transcriptomic blueprint of
683 localized prostate cancer that identifies and highlights the multicellular milieu and
684 cellular states associated with prostate tumorigenesis. Our results provide new insights
685 into the epithelial microenvironment and the cellular state changes associated with
686 prostate cancer toward improved PCa diagnosis.

687

688 **Methods**

689

690 **Experimental Details**

691

692 **Samples selection**

693 We obtained a total of six prostate biopsies from three different patients (two
694 biopsies for patient 1-3, obtained at the same time point), four radical prostatectomies
695 with tumor-only samples from four patients (patients 4-7) and four radical
696 prostatectomies with matched normal samples from four patients (patients 8-11,
697 matched normal samples were taken from adjacent seemingly normal regions).
698 Clinical/pathological data available for the samples is in **Supplemental Table 1**.

699

700 **Study Approval**

701 The UCSF Institutional Review Board (IRB) committee approved the collection of
702 the patient data included in this study.

703

704 **Tissue Dissociation**

705 Tissue samples were minced with surgical scissors and washed with RP-10
706 (RPMI + 10% FBS). Each sample was centrifuged at 1200 rpm for five minutes,
707 resuspended in 10 mL digestive media (HBSS + 1% HEPES) with Liberase TM (Roche,
708 Cat: 5401119001) or 1000 U/mL collagenase type IV (Worthington, Cat: LS004188),
709 and rotated for 30 minutes at 37 °C. Samples were triturated by pipetting ten times after
710 every ten minutes during the incubation or by pipetting 15 times at the end of the
711 incubation. Each sample was filtered through a 70 µm filter (Falcon, Cat: 352350),
712 washed with RP-10, centrifuged at 1200 rpm for five minutes, washed again with RP-10,
713 and resuspended in RP-10. A hemocytometer was used to count the cells.

714

715 **Single-cell RNA sequencing**

716 Sequencing was largely based on the Seq-Well S³ protocol^{13,95}. One to four
717 arrays were used per sample. Each array was loaded as previously described with
718 approximately 110,000 barcoded mRNA capture beads (ChemGenes, Cat: MACOSKO-
719 2011-10(V+)) and with 10,000-20,000 cells. Arrays were sealed with functionalized
720 polycarbonate membranes (Sterlitech, Cat: PCT00162X22100) and were incubated at
721 37°C for 40 minutes.

722 After sealing, each array was incubated in lysis buffer (5 M Guanidine
723 Thiocyanate, 1 mM EDTA, 0.5% Sarkosyl, 1% BME). After detachment and removal of
724 the top slides, arrays were rotated at 50 rpm for 20 minutes. Each array was washed
725 with hybridization buffer (2 M NaCl, 4% PEG8000) and then rocked in hybridization
726 buffer for 40 minutes. Beads from different arrays were collected separately. Each array
727 was washed ten times with wash buffer (2 M NaCl, 3 mM MgCl₂, 20 mM Tris-HCl pH
728 8.0, 4% PEG8000) and scraped ten times with a glass slide to collect beads into a
729 conical tube.

730 For each array, beads were washed with Maxima RT buffer (ThermoFisher, Cat:
731 EP0753) and resuspended in reverse transcription mastermix with Maxima RT buffer,
732 PEG8000, Template Switch Oligo, dNTPs (NEB, Cat: N0447L), RNase inhibitor (Life
733 Technologies, Cat: AM2696), and Maxima H Minus Reverse Transcriptase
734 (ThermoFisher, Cat: EP0753) in water. Samples were rotated end-to-end, first at room
735 temperature for 15 minutes and then at 52°C overnight. Beads were washed once with
736 TE-SDS and twice with TE-TW. They were treated with exonuclease I (NEB), rotating
737 for 50 minutes at 37°C. Beads were washed once with TE-SDS and twice with TE-TW,

738 and once with 10 mM Tris-HCl pH 8.0. They were resuspended in 0.1 M NaOH and
739 rotated for five minutes at room temperature. They were subsequently washed with TE-
740 TW and TE. They were taken through second strand synthesis with Maxima RT buffer,
741 PEG8000, dNTPs, dN-SMRT oligo, and Klenow Exo- (NEB, Cat: M0212L) in water.
742 After rotating at 37°C for one hour, beads were washed twice with TE-TW, once with
743 TE, and once with water.

744 KAPA HiFi Hotstart Readymix PCR Kit (Kapa Biosystems, Cat: KK2602) and
745 SMART PCR Primer were used in whole transcriptome amplification (WTA). For each
746 array, beads were distributed among 24 PCR reactions. Following WTA, three pools of
747 eight reactions were made and were then purified using SPRI beads (Beckman
748 Coulter), first at 0.6x and then at a 0.8x volumetric ratio.

749 For each sample, one pool was run on an HSD5000 tape (Agilent, Cat: 5067-
750 5592). The concentration of DNA for each of the three pools was measured via the
751 Qubit dsDNA HS Assay kit (ThermoFisher, Cat: Q33230). Libraries were prepared for
752 each pool, using 800-1000 pg of DNA and the Nextera XT DNA Library Preparation Kit.
753 They were dual-indexed with N700 and N500 oligonucleotides.

754 Library products were purified using SPRI beads, first at 0.6x and then at a 1x
755 volumetric ratio. Libraries were then run on an HSD1000 tape (Agilent, Cat: 50675584)
756 to determine the concentration between 100-1000 bp. For each library, 3 nM dilutions
757 were prepared. These dilutions were pooled for sequencing on a NovaSeq S4 flow cell.

758 The sequenced data were preprocessed and aligned using the
759 dropseq_workflow on Terra (app.terra.bio). A digital gene expression matrix was

760 generated for each sample, parsed and analyzed following a customized pipeline.

761 Additional details are provided below.

762

763 **Organoid culture**

764 Isolated single cells not used for single-cell sequencing were additionally frozen
765 in FBS + 10% DMSO, flash frozen on dry ice, or plated in Matrigel to grow as 3D
766 prostate organoid cultures. Organoid cultures were established by plating 20,000 cells
767 in 25uL Matrigel (Corning, Cat: 356231) in 48-well flat-bottom plates (Corning, Cat: EK-
768 47102). Prostate-specific serum-free culture media contained 500 ng/mL human
769 recombinant R-spondin (R&D Systems, Cat: 10820-904), 10uM SB202190 (Sigma, Cat:
770 S7076), 1uM Prostaglandin E3 (Tocris, CA: 229610), 1nM FGF10 (Peprotech, Cat:
771 100-26), 5 ng/mL FGF2 (Peprotech, CA: 100-18B), 10 ng/mL 5alpha-
772 Dihydrotestosterone (Sigma, Cat: D-073-1ML), 100 ng/mL human Noggin (Peprotech,
773 Cat: 102-10C), 500nM A83-01 (Fischer, Cat: 29-391-0), 5 ng/mL human EGF
774 (Peprotech, Cat: AF-100-15), 1.25mM N-acetyl-cysteine (Sigma, Cat: A9165), 10mM
775 Nicotinamide (Sigma, Cat: N3376), 1X B-27 (Gibco, Cat: 17504044), 1X P/S (Gibco,
776 Cat: 15140122), 10mM HEPES (Gibco, CA: 15630080), and 2mM GlutaMAX (Gibco,
777 Cat: 35050061)⁶⁹. Additionally, 10uM Y-27 (Biogems, Cat: 1293823) was included
778 during the first 2 weeks of growth and after passaging to promote growth⁶⁹. Generally,
779 organoid growth was apparent within two to three days and robust after two weeks.
780 250uL media was refreshed every two to four days using media stored at 4°C for a
781 maximum of ten days. Organoid growth was monitored using an EVOS-FL microscope.

782 To passage prostate organoid cultures every 7-14 days, culture media was
783 replaced with 300 uL TrypLE (1X, Gibco, Cat: 12604013). Individual domes were
784 collected into 15mL Falcon tubes, disrupted by pipetting with wide-orifice tips and
785 incubated at 37°C for 30 minutes. Following incubation, the dissociation media was
786 neutralized using 10mL wash media: adDMEM/F12 containing 5% FBS, P/S, 10mM
787 HEPES (1M, Gibco, Cat: 15630080) and 2mM GlutaMAX (100x, Gibco, Cat:
788 35050061)⁶⁹. Cells were spun down at 500 G for five minutes and resuspended in 2mL
789 wash media. Finally, the media was aspirated, cells were resuspended in Matrigel, and
790 25 uL/dome were plated per well.

791 Organoids were accessed using single-cell sequencing at an early passage (P0-
792 4). To isolate single cells from Matrigel, organoids were collected in 500uL Trypsin
793 (0.25%, Gibco, Cat: 25-200-056) and incubated at 37°C for 30-45 minutes until few
794 clumps were visible. Throughout incubation, cells were triturated every five minutes.
795 Single cells were resuspended in 9mL DMEM + 5% FBS + 0.05mM EDTA and passed
796 through a 40 µM filter, followed by an additional wash of the filter with 1mL DMEM + 5%
797 FBS + 0.05mM EDTA. Cells were spun down at 300 G for 5 minutes, resuspended in
798 10mL of the same media, spun down again and finally, resuspended in 1-2mL media.
799 Cells were counted using a hemocytometer and loaded on to arrays for single-cell
800 sequencing as described for patient tissues.

801

802 **Immunofluorescence**

803 Organoids were passaged into 8-well Nunc Lab-Tek II Chamber Slides (Thermo
804 Scientific, Cat: 154453) and allowed to grow in prostate-specific media. Following seven

805 days, the media was removed, domes were washed twice with 300uL PBS and fixed in
806 4% paraformaldehyde (Electron Microscopy Sciences, Cat: 15710-S) at room
807 temperature for 20 minutes. Individual domes were washed 3x with IF Buffer (0.02%
808 Triton + 0.05% Tween + PBS) and blocked for one hour at room temperature with 0.5%
809 Triton X100 + 1% DMSO + 1% BSA + 5% Donkey Serum + PBS. Following the block,
810 domes were washed once with IF Buffer and incubated overnight with monoclonal
811 mouse anti-Lactoferrin (Abcam, Cat: ab101110, 1ug/mL), monoclonal rat anti-
812 Uteroglobulin/SCGB1A1 (R&D Systems, Cat: MAB4218-SP, 1ug/mL), polyclonal guinea
813 pig anti-Cytokeratin 8 + 18 (Fitzgerald, Cat: 20R-CP004, 1:100), and polyclonal chicken
814 anti-Keratin 5 (Biolegend, Cat: 905901, 1:100). Subsequently, domes were washed 3x
815 with IF Buffer and counterstained with Alexa Fluor 488-AffiniPure Donkey Anti-Chicken
816 IgY (IgG) (H+L) (Jackson ImmunoResearch, Cat: 703-545-155, 1:500), Donkey anti-
817 Mouse IgG (H+L) Cross-Adsorbed Secondary Antibody, DyLight 550 (Thermo Fisher
818 Scientific, Cat: SA5-10167, 1:500), Donkey anti-Rat IgG (H+L) Cross-Adsorbed
819 Secondary Antibody, DyLight 680 (Thermo Fisher Scientific, Cat: SA5-10030, 1:500),
820 and Alexa Fluor 790 AffiniPure Donkey Anti-Guinea Pig IgG (H+L) (Jackson
821 ImmunoResearch, Cat: 706-655-148, 1:500) containing DAPI (Sigma, Cat: D9542-5MG,
822 1:1000). Finally, wells were washed 3x with IF Buffer for five minutes and sealed with
823 Prolong Gold antifade mountant (Fischer Sci, Cat: P36930). Z-stack images were
824 captured on a Leica DCF9000 GT using Leica Application System X software.

825

826 **Quantification and Statistical Analysis**

827

828 **Sequencing and Alignment**

829 Sequencing results were returned as paired FASTQ reads and processed with
830 FastQC⁹⁶ for general quality checks in order to further improve our experimental
831 protocol. Then, the paired FASTQ files were aligned against the reference genome
832 using a STAR aligner⁹⁷ in the dropseq workflow
833 (https://cumulus.readthedocs.io/en/latest/drop_seq.html). The aligning pipeline output
834 included aligned and corrected bam files, two digital gene expression (DGE) matrix text
835 files (a raw read count matrix and a UMI-collapsed read count matrix where multiple
836 reads that matched the same UMI would be collapsed into one single UMI count) and
837 text-file reports of basic sample qualities such as the number of beads used in the
838 sequencing run, total number of reads, alignment logs. For each sample, the average
839 number of reads was 4,875,9687, and the mean read depth per barcode was 48,586.
840 The median and average number of genes per barcode were 767 and 1079. The
841 median and average number of UMI were 1,335 and 2,447. The mean percentage of
842 mitochondrial content per cell was 13.65%.

843

844 **Single-cell clustering analysis**

845 Cells in the samples were clustered and analyzed using customized codes based
846 on the Seurat V3.0 package on R²⁰. Cells with less than 300 genes, 500 transcripts, or a
847 mitochondrial level of 20% or greater, were filtered out as the first QC process. Then, by
848 examining the distribution histogram of the number of genes per cell in each sample, we
849 set the upper threshold for the number of genes per cell in each individual sample in
850 order to filter potential doublets. A total of 22,037 cells were acquired using these

851 thresholds. Since merging with and without integration of the samples showed no major
852 difference in the clustering of each cell type, in the subsequent analysis of these
853 samples we used the merged dataset without integration.

854 Doublets were removed by two steps: first we used DoubletFinder⁹⁸ and a
855 theoretical doublet rate of 5% to locate doublets in our dataset. 294 cells marked by
856 DoubletFinder as true positive were removed from further analysis. 21,743 cells were
857 used in the following cell clustering analysis. Then, after clustering, we removed cells
858 expressing biomarkers from more than one major cell type (epithelial, stromal, and
859 immune) as they were more likely to be doublets. In this step, we removed 276 cells
860 from our dataset and the follow-up analysis, leaving 21,467 cells in total.

861 UMI-collapsed read counts matrices for each cell were loaded in Seurat for
862 analysis²⁰. We followed the standard workflow by using the “LogNormalize” method that
863 normalized the gene expression for each cell by the total expression, multiplying by a
864 scale factor 10,000 and log-transforming the results. For downstream analysis to
865 identify different cell types, we then calculated and returned the top 2,000 most variably
866 expressed genes among the cells before applying a linear scaling by shifting the
867 expression of each gene in the dataset so that the mean expression across cells was 0
868 and the variance was 1. This way, the gene expression level could be comparable
869 among different cells and genes. PCA was run using the previously determined most
870 variably expressed genes for linear dimensional reduction and the first 100 principal
871 components (PCs) were stored which accounted for 25.42% of the total variance. To
872 determine how many PCs to use for the clustering, a JackStraw resampling method was
873 implemented by permutation on a subset of data (1% by default) and rerunning PCA for

874 a total of 100 replications to select the statistically significant principle component to
875 include for the K-nearest neighbors clustering. For graph-based clustering, the first 100
876 PC and a resolution of 3 were selected yielding a total of 46 cell clusters. We eliminated
877 the clustering side effect due to overclustering by constructing a cluster tree of the
878 average expression profile in each cluster and merging clusters together based on their
879 positions in the cluster tree. As a result, we ensured that each cluster would have at
880 least ten unique differentially expressed genes (DEGs). Differentially expressed genes
881 in each cluster were identified using the FindAllMarker function within Seurat package
882 and a corresponding p-value was given by the Wilcoxon's test followed by a Bonferroni
883 correction. Top differentially expressed gene markers were illustrated in a stacked violin
884 plot using a customized auxiliary function. Dot plots were generated as an alternative
885 way of visualization using the top ten differentially expressed genes in each cluster. Top
886 tier cell type clustering was also validated by the automated singleR annotation
887 **(Supplemental Table 1)**

888

889 **Cell type annotation by signature scores**

890 In order to annotate each cell type from the previous clustering, we took the
891 established studies and the signatures for each cell type (**Supplemental Table 2**).
892 Treating the signature score of each cell type as a pseudogene, we evaluated the
893 signature score for each cell in our dataset using the AddModuleScore function²⁰. Each
894 cluster in our dataset was assigned with an annotation of its cell type by top signature
895 scores within the cluster.

896

897 **Epithelial sub-clustering analysis and tumor cell inference**

898 All epithelial cells were clustered using the analytical workflow described above,
899 yielding 20 clusters. To compare the transcriptomic profiles between PCa samples and
900 normal prostates, a previous study on normal prostate single-cell RNA-seq was
901 downloaded and imported. Mean basal, luminal, hillock, and club signature scores were
902 calculated for each cluster, based on the top differentially expressed genes from a
903 previous scRNA-seq study on the normal prostate. A One-way ANOVA test was then
904 conducted to determine if the signature score of each cluster was significantly different
905 from the rest. We annotated the clusters with significantly upregulated basal epithelial
906 cell (BE) signature scores as BE. Cells in clusters with high luminal epithelial (LE)
907 signature scores could be either non-malignant luminal epithelial cells or tumor cells.
908 The clusters with low signature scores of both BE and LE were annotated as other
909 epithelial cells (OE). To efficiently identify tumor cells, we took the digital gene
910 expression matrix and conducted a single set gene set enrichment analysis on
911 GenePattern (<https://gsea-msigdb.github.io/ssGSEA-gpmodule/v10/index.html>) testing
912 against the C2 gene set collection curated on MSigDB ([https://www.gsea-](https://www.gsea-msigdb.org/gsea/msigdb/index.jsp)
913 [msigdb.org/gsea/msigdb/index.jsp](https://www.gsea-msigdb.org/gsea/msigdb/index.jsp)). Under the notion that tumor cells should have
914 higher expression of one or more tumor markers overlapping existing prostate cancer
915 gene sets, we projected the signatures of these prostate cancer gene sets on to our
916 epithelial clusters and annotated tumor cell clusters as the clusters with significantly
917 higher ($p < 0.05$ in one-way ANOVA test) signature scores of at least one prostate
918 cancer gene sets.

919 Approximately ~50% of prostate cancer cells from men of European ancestry
920 harbor *TMPRSS2-ERG* fusion events, indicating high gene expression of *ERG*^{99,100}.
921 Therefore, we hypothesized a high signature score of SETLUR PROSTATE CANCER
922 *TMPRSS2 ERG FUSION UP* gene set²⁶ would be a strong indicator of *ERG*+ tumor
923 cells. All the other tumor cell clusters were then annotated as *ERG*- tumor cell clusters
924 as they showed little to no *ERG* gene expression. All of the epithelial clusters with high
925 luminal signature scores and high expression of luminal markers such as *KLK3*, *KLK2*,
926 *ACPP*, *KRT8*, and *KRT18* were annotated as non-malignant luminal epithelial cells
927 (non-malignant LE). Compared to non-malignant cells, tumor cells harbor more single-
928 nucleotide variants and copy number variants, leading to distinctive patterns. To
929 validate our tumor cell annotation, we ran InferCNV on *ERG*+ and *ERG*- tumor clusters
930 with non-malignant LEs as reference²⁹ for an estimation of copy number alterations. We
931 classified tumor cells based on *ERG* gene expression. Then we defined patients
932 harboring *ERG*+ tumor cells as *ERG*+ patients and the other patients as *ERG*- patients.
933 This way, we were able to classify all the other cells based on the *ERG* status
934 (epithelial, stromal, and immune cells) as either *ERG*+ or *ERG*-.

935 To determine if common functional changes were present in more than one cell
936 type, we conducted gene set enrichment analysis (GSEA) for each cell type first and
937 imported the significantly changed gene sets to take the intersections. Statistical
938 significance of multi-set intersection was evaluated and visualized using the
939 SuperExacTest package⁵¹.

940

941 **Cell state analysis**

942 Gene expression profile differences in epithelial cells between PCa sample and
943 normal prostate samples were identified by integrating our PCa dataset with an
944 established dataset on normal prostates⁹. We utilized the integration method based on
945 commonly-expressed anchor genes by following the Seurat integration vignette²⁰ in
946 order to remove batch effects of samples sequenced with different technologies and
947 possible artifacts so that the cells were comparable.

948 In order to better characterize the transcriptomic profile and transition of cell
949 states among identified epithelial cells, both the tumor and paired normal samples were
950 integrated together and separately with the epithelial cells from a normal prostate
951 scRNA-seq dataset⁹ for *KRT5+* and *KRT15+* basal epithelial (BE), *KLK3+* and *ACPP+*
952 luminal epithelial (LE), and *PIGR+* and *MMP7+* club cell population together and
953 separately. An optimal resolution value was tested using the Clustree¹⁰¹ package.
954 Heatmaps of DEGs were generated to validate the cell state differentiation.
955 Compositions for each cell state were computed and compared between PCa samples
956 and normal samples using Fisher's exact test.

957 To assess the functional roles of the PCa-enriched cell states identified within the
958 integrated dataset, we ran GSEA analysis between the PCa-enriched cell state and all
959 the other cell states as a whole. The top 20 downregulated and upregulated gene sets
960 were visualized in terms of gene counts and ratio for each gene set. Using the DEGs
961 from each cell state, we generated signature gene sets for all the cell states in BE, LE,
962 and club cells. To validate the functional implications for the PCa-enriched cell states,
963 we conducted ssGSEA on PCa BE and club cells to compute the signature scores of
964 the upregulated gene sets using the ssGSEA module on GenePattern (

965 msigdb.github.io/ssGSEA-gpmodule/v10/index.html). Then, we computed the
966 information coefficient (IC) and corresponding p-values followed by FDR correction to
967 evaluate the correlation between these gene sets and cell states.

968

969 **Pseudotime analysis**

970 To evaluate the epithelial cell states with respect to their order in the
971 differentiation trajectory, we conducted pseudotime analysis on all epithelial and tumor
972 cells identified in the PCa samples. We first calculated a PAGA (partition-based graph
973 abstraction) graph using SCANPY's `sc.tl.paga()` function¹⁰² and then used
974 `sc.tl.draw_graph()` to generate the PAGA initialized single-cell embedding of the cell
975 types (**Supplemental Figure 4a**). The diffusion pseudotime for each cell was calculated
976 using SCANPY's `sc.tl.diffmap()` and `sc.tl.dpt()` with the root cell chosen from the stem
977 cell upregulated BE cluster and then was plotted on the PAGA initialized embedding.
978 (**Supplemental Figure 4b**). We then visualized the gene marker changes along the
979 pseudotime by cell type using `sc.pl.paga_path()` (Supplemental Figure 4c).

980 Furthermore, to test whether or not the luminal-like cell state within the club cell
981 population was more differentiated compared to other club cells, we utilized Monocle3⁴⁵
982 on club cells. Monocle3 object was generated using the count matrix for all club cells
983 and the pseudotime trajectory was computed following the standard Monocle3 workflow.
984 The starting point of the trajectory was identified using the cell with the highest adult
985 stem cell signature score (**Supplemental Figure 4d**) and the luminal-like club cells
986 were highlighted using the luminal epithelial cell signature (**Supplemental Figure 4d**).

987 Expression levels along the pseudotime trajectory for club cell markers *LTF* and *PIGR*,
988 and luminal markers *ACPP* and *KLK3* were then plotted.

989

990 **scRNA-seq Fusion detection**

991 Fusion transcripts were detected using STAR-Fusion²⁷ ([https://github.com/STAR-](https://github.com/STAR-Fusion/STAR-Fusion/wiki)
992 [Fusion/STAR-Fusion/wiki](https://github.com/STAR-Fusion/STAR-Fusion/wiki)) version 1.6.0. STAR-Fusion was run from a Docker container
993 using the following options: *--FusionInspectorvalidate*, *--examine_coding_effect*, and *--*
994 *denovo_reconstruct*. Due to the low coverage of scRNA-seq samples both filtered
995 fusion detection results and preliminary results were combined and processed, and we
996 only filtered for potential *TMPRSS2-ERG* fusion events.

997

998 **Signature analyses of bulk RNA-sequencing datasets**

999 Two publicly available bulk RNA-sequencing PCa datasets were used to test the
1000 correlation between the PCa-enriched cell state signatures and AR signaling, including
1001 Prostate Adenocarcinoma (TCGA²⁵, Firehose Legacy) dataset (N = 499, available at
1002 https://www.cbioportal.org/study/summary?id=prad_tcga) and SU2C/PCF Dream Team
1003 (SU2C⁴⁸, PNAS 2019) dataset (N = 266, available at
1004 https://www.cbioportal.org/study/summary?id=prad_su2c_2019). For each dataset,
1005 mRNA expression was downloaded and normalized. Signature scores of *AR* signaling
1006 (Hallmark androgen response pathway), BE, LE, and club cell states as well as *ERG+*
1007 and *ERG-* tumor cell signature scores were computed for each sample via ssGSEA.
1008 Samples in each dataset were rank ordered by the *AR* signature scores and heatmaps
1009 were generated using the customized scripts. To test the correlation between *AR*

1010 signature scores and each cell state signature score, we computed the information
1011 coefficient and corresponding p-values followed by FDR correction to evaluate the
1012 correlation. For tumor cell signatures, we computed the correlations between the *ERG*
1013 fusion status from each dataset and the signature scores of *ERG+* and *ERG-* tumor cell
1014 gene sets we had previously generated. We rank ordered the bulk RNA-seq samples
1015 according to whether or not the *TMPRSS2-ERG* fusion was detected and plotted the
1016 *ERG+* and *ERG-* tumor cell signature score heatmaps. Information coefficient (IC), p-
1017 values, and FDR q-values were computed.

1018

1019 **Immune cell analysis**

1020 T-cell and myeloid cell populations were sub-clustered separately following a
1021 similar pipeline as described above. For T-cells, 23 PCs and a resolution of 1.5 were
1022 selected for the clustering. For myeloid cells, 27 PCs and a resolution of 1.5 were
1023 selected. Cell clusters were annotated by a dot plot showing the top ten most expressed
1024 genes in each cluster.

1025 Monocytes, macrophages, neutrophils, and eosinophils were identified and
1026 annotated based on the automated SingleR analysis¹⁹. M1/M2 macrophage
1027 phenotypes, tumor associated macrophages, and two types of myeloid-derived
1028 suppressor cells were identified using documented markers from previous studies.

1029

1030 **Materials Availability**

1031 This study did not generate new unique reagents.

1032

1033 **Data and Code Availability**

1034 Processed single-cell RNA sequencing data that support this study will be
1035 deposited in the NCBI GEO database and available upon request to the corresponding
1036 author. All software algorithms used for analysis are available for download from public
1037 repositories. All code used to generate figures in the manuscript are available upon
1038 request.

1039

1040 **Acknowledgments**

1041 This work was in part supported by: Searle Scholars Program (A.K.S.), Beckman
1042 Young Investigator Program (A.K.S.), Sloan Fellowship in Chemistry (A.K.S.), Pew-
1043 Stewart Scholars Program for Cancer Research (A.K.S.), and the Prostate Cancer
1044 Foundation (F.W.H.).

1045 We thank Travis Hughes and Matthew Hellmann for helpful discussions.

1046

1047 **Author Contributions**

1048 Conceptualization, H.S and F.W.H.; Methodology, H.S, H.N.W, P.A, M.H.W,
1049 B.W., and F.W.H.; Investigation, H.S. and J.X., Writing – Original Draft, H.S. and
1050 F.W.H.; Writing – Review & Editing, H.S., H.N.W., P.A., J.X., M.H.W, F.Y.F. M.R.C., and
1051 A.K.S., F.W.H.,; Resources, M.R.C., P.C., B.W., H.Y., A.K.S. and F.W.H.; Supervision,
1052 F.W.H.

1053

1054 **Competing Interests**

1055 A.K.S. reports compensation for consulting and/or SAB membership from Merck,
1056 Honeycomb Biotechnologies, Cellarity, Repertoire Immune Medicines, Orche Bio, and
1057 Dahlia Biosciences.

1058 F.Y.F. reports compensation for consulting and/or SAB membership from
1059 Astellas, Bayer, Blue Earth Diagnostics, Celgene, Genentech, Janssen Oncology,
1060 Myovant, Roivant, Sanofi, PFS Genomics, and SerImmune.

1061

1062 **Reference**

- 1063 1. Tiwari, R., Manzar, N. & Ateeq, B. Dynamics of Cellular Plasticity in Prostate Cancer
1064 Progression. *Front. Mol. Biosci.* **7**, (2020).
- 1065 2. Blau, H. M. *et al.* Plasticity of the differentiated state. *Science* **230**, 758–766 (1985).
- 1066 3. Yuan, S., Norgard, R. J. & Stanger, B. Z. Cellular Plasticity in Cancer. *Cancer*
1067 *Discov.* **9**, 837–851 (2019).
- 1068 4. Varga, J. & Greten, F. R. Cell plasticity in epithelial homeostasis and tumorigenesis.
1069 *Nat. Cell Biol.* **19**, 1133–1141 (2017).
- 1070 5. Boutros, P. C. *et al.* Spatial genomic heterogeneity within localized, multifocal
1071 prostate cancer. *Nat. Genet.* **47**, 736–745 (2015).
- 1072 6. Cooper, C. S. *et al.* Analysis of the genetic phylogeny of multifocal prostate cancer
1073 identifies multiple independent clonal expansions in neoplastic and morphologically
1074 normal prostate tissue. *Nat. Genet.* **47**, 367–372 (2015).
- 1075 7. DeMarzo, A. M., Nelson, W. G., Isaacs, W. B. & Epstein, J. I. Pathological and
1076 molecular aspects of prostate cancer. *Lancet Lond. Engl.* **361**, 955–964 (2003).
- 1077 8. Shen, M. M. & Abate-Shen, C. Molecular genetics of prostate cancer: new prospects
1078 for old challenges. *Genes Dev.* **24**, 1967–2000 (2010).
- 1079 9. Henry, G. H. *et al.* A Cellular Anatomy of the Normal Adult Human Prostate and
1080 Prostatic Urethra. *Cell Rep.* **25**, 3530-3542.e5 (2018).
- 1081 10. Tomlins, S. A. *et al.* Recurrent fusion of TMPRSS2 and ETS transcription factor
1082 genes in prostate cancer. *Science* **310**, 644–648 (2005).
- 1083 11. Bhatia, V. & Ateeq, B. Molecular Underpinnings Governing Genetic Complexity of
1084 ETS-Fusion-Negative Prostate Cancer. *Trends Mol. Med.* **25**, 1024–1038 (2019).

- 1085 12. Abeshouse, A. *et al.* The Molecular Taxonomy of Primary Prostate Cancer. *Cell* **163**,
1086 1011–1025 (2015).
- 1087 13. Hughes, T. K. *et al.* Second-Strand Synthesis-Based Massively Parallel scRNA-Seq
1088 Reveals Cellular States and Molecular Features of Human Inflammatory Skin
1089 Pathologies. *Immunity* **53**, 878-894.e7 (2020).
- 1090 14. Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by
1091 single-cell RNA-seq. *Science* **352**, 189–196 (2016).
- 1092 15. Zilionis, R. *et al.* Single-Cell Transcriptomics of Human and Mouse Lung Cancers
1093 Reveals Conserved Myeloid Populations across Individuals and Species. *Immunity*
1094 **50**, 1317-1334.e10 (2019).
- 1095 16. Kapellos, T. S. *et al.* Human Monocyte Subsets and Phenotypes in Major Chronic
1096 Inflammatory Diseases. *Front. Immunol.* **10**, (2019).
- 1097 17. Tang-Huau, T.-L. *et al.* Human in vivo-generated monocyte-derived dendritic cells
1098 and macrophages cross-present antigens through a vacuolar pathway. *Nat.*
1099 *Commun.* **9**, 2570 (2018).
- 1100 18. Hadadi, E. *et al.* Differential IL-1 β secretion by monocyte subsets is regulated by
1101 Hsp27 through modulating mRNA stability. *Sci. Rep.* **6**, (2016).
- 1102 19. Aran, D. *et al.* Reference-based analysis of lung single-cell sequencing reveals a
1103 transitional profibrotic macrophage. *Nat. Immunol.* **20**, 163–172 (2019).
- 1104 20. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-
1105 1902.e21 (2019).

- 1106 21. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
1107 transcriptomic data across different conditions, technologies, and species. *Nat.*
1108 *Biotechnol.* **36**, 411–420 (2018).
- 1109 22. Yue, X. *et al.* Polymeric immunoglobulin receptor promotes tumor growth in
1110 hepatocellular carcinoma. *Hepatol. Baltim. Md* **65**, 1948–1962 (2017).
- 1111 23. Zhang, Q. *et al.* Interleukin-17 promotes prostate cancer via MMP7-induced
1112 epithelial-to-mesenchymal transition. *Oncogene* **36**, 687–699 (2017).
- 1113 24. Fotiou, K. *et al.* Serum ceruloplasmin as a marker in prostate cancer. *Minerva Urol.*
1114 *E Nefrol. Ital. J. Urol. Nephrol.* **59**, 407–411 (2007).
- 1115 25. The molecular taxonomy of primary prostate cancer. *Cell* **163**, 1011–1025 (2015).
- 1116 26. Setlur, S. R. *et al.* Estrogen-dependent signaling in a molecularly distinct subclass of
1117 aggressive prostate cancer. *J. Natl. Cancer Inst.* **100**, 815–825 (2008).
- 1118 27. Haas, B. J. *et al.* Accuracy assessment of fusion transcript detection via read-
1119 mapping and de novo fusion transcript assembly-based methods. *Genome Biol.* **20**,
1120 213 (2019).
- 1121 28. Qian, X. *et al.* Spondin-2 (SPON2), a More Prostate-Cancer-Specific Diagnostic
1122 Biomarker. *PLoS ONE* **7**, (2012).
- 1123 29. Kenny, P. A. InferCNV, a python web app for copy number inference from discrete
1124 gene-level amplification signals noted in clinical tumor profiling reports.
1125 *F1000Research* **8**, 807 (2019).
- 1126 30. Ellis, L. & Loda, M. Advanced neuroendocrine prostate tumors regress to stemness.
1127 *Proc. Natl. Acad. Sci.* **112**, 14406–14407 (2015).

- 1128 31. Goto, K. *et al.* Proximal prostatic stem cells are programmed to regenerate a
1129 proximal-distal ductal axis. *Stem Cells Dayt. Ohio* **24**, 1859–1868 (2006).
- 1130 32. Tsujimura, A. *et al.* Proximal location of mouse prostate epithelial stem cells: a
1131 model of prostatic homeostasis. *J. Cell Biol.* **157**, 1257–1265 (2002).
- 1132 33. Reiter, R. E. *et al.* Prostate stem cell antigen: a cell surface marker overexpressed in
1133 prostate cancer. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 1735–1740 (1998).
- 1134 34. Montoro, D. T. *et al.* A revised airway epithelial hierarchy includes CFTR-expressing
1135 ionocytes. *Nature* **560**, 319–324 (2018).
- 1136 35. Thorek, D. L., Evans, M. J., Carlsson, S. V., Ulmert, D. & Lilja, H. Prostate Specific
1137 Kallikrein-related Peptidases and Their Relation to Prostate Cancer Biology and
1138 Detection; Established Relevance and Emerging Roles. *Thromb. Haemost.* **110**,
1139 484–492 (2013).
- 1140 36. Hessels, D. & Schalken, J. A. Urinary biomarkers for prostate cancer: a review.
1141 *Asian J. Androl.* **15**, 333–339 (2013).
- 1142 37. Zhang, L. & Barritt, G. J. TRPM8 in prostate cancer cells: a potential diagnostic and
1143 prognostic marker with a secretory function? *Endocr. Relat. Cancer* **13**, 27–38
1144 (2006).
- 1145 38. Smith, B. A. *et al.* A Human Adult Stem Cell Signature Marks Aggressive Variants
1146 across Epithelial Cancers. *Cell Rep.* **24**, 3353-3366.e5 (2018).
- 1147 39. Guo, W. *et al.* Single-cell transcriptomics identifies a distinct luminal progenitor cell
1148 type in distal prostate invagination tips. *Nat. Genet.* **52**, 908–918 (2020).
- 1149 40. Gurioli, G. *et al.* Methylation pattern analysis in prostate cancer tissue: identification
1150 of biomarkers using an MS-MLPA approach. *J. Transl. Med.* **14**, 249 (2016).

- 1151 41. Tung, M.-C. *et al.* Knockdown of lipocalin-2 suppresses the growth and invasion of
1152 prostate cancer cells. *The Prostate* **73**, 1281–1290 (2013).
- 1153 42. Nelson, P. S. *et al.* The program of androgen-responsive genes in neoplastic
1154 prostate epithelium. *Proc. Natl. Acad. Sci.* **99**, 11890–11895 (2002).
- 1155 43. Feng, Q. & He, B. Androgen Receptor Signaling in the Development of Castration-
1156 Resistant Prostate Cancer. *Front. Oncol.* **9**, 858 (2019).
- 1157 44. Fujita, K. & Nonomura, N. Role of Androgen Receptor in Prostate Cancer: A Review.
1158 *World J. Mens Health* **37**, 288–295 (2019).
- 1159 45. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis.
1160 *Nature* **566**, 496–502 (2019).
- 1161 46. Guérin, O., Fischel, J. L., Ferrero, J.-M., Bozec, A. & Milano, G. EGFR Targeting in
1162 Hormone-Refractory Prostate Cancer: Current Appraisal and Prospects for
1163 Treatment. *Pharmaceuticals* **3**, 2238–2247 (2010).
- 1164 47. Varma, M., Berney, D. M., Jasani, B. & Rhodes, A. Technical variations in prostatic
1165 immunohistochemistry: need for standardisation and stringent quality assurance in
1166 PSA and PSAP immunostaining. *J. Clin. Pathol.* **57**, 687–690 (2004).
- 1167 48. Abida, W. *et al.* Genomic correlates of clinical outcome in advanced prostate cancer.
1168 *Proc. Natl. Acad. Sci.* **116**, 11428–11436 (2019).
- 1169 49. Atsaves, V., Leventaki, V., Rassidakis, G. Z. & Claret, F. X. AP-1 Transcription
1170 Factors as Regulators of Immune Responses in Cancer. *Cancers* **11**, (2019).
- 1171 50. Darash-Yahana, M. *et al.* The Chemokine CXCL16 and Its Receptor, CXCR6, as
1172 Markers and Promoters of Inflammation-Associated Cancers. *PLoS ONE* **4**, (2009).

- 1173 51. Wang, M., Zhao, Y. & Zhang, B. Efficient Test and Visualization of Multi-Set
1174 Intersections. *Sci. Rep.* **5**, 16923 (2015).
- 1175 52. Fay, A. P. & Antonarakis, E. S. Blocking the PD-1/PD-L1 axis in advanced prostate
1176 cancer: are we moving in the right direction? *Ann. Transl. Med.* **7**, (2019).
- 1177 53. Castro, F., Cardoso, A. P., Gonçalves, R. M., Serre, K. & Oliveira, M. J. Interferon-
1178 Gamma at the Crossroads of Tumor Immune Surveillance or Evasion. *Front.*
1179 *Immunol.* **9**, (2018).
- 1180 54. Chung, W. *et al.* Single-cell RNA-seq enables comprehensive tumour and immune
1181 cell profiling in primary breast cancer. *Nat. Commun.* **8**, 15081 (2017).
- 1182 55. Oh, D. Y. *et al.* Intratumoral CD4+ T Cells Mediate Anti-tumor Cytotoxicity in Human
1183 Bladder Cancer. *Cell* **181**, 1612-1625.e13 (2020).
- 1184 56. Orecchioni, M., Ghosheh, Y., Pramod, A. B. & Ley, K. Macrophage Polarization:
1185 Different Gene Signatures in M1(LPS+) vs. Classically and M2(LPS-) vs.
1186 Alternatively Activated Macrophages. *Front. Immunol.* **10**, (2019).
- 1187 57. Jablonski, K. A. *et al.* Novel Markers to Delineate Murine M1 and M2 Macrophages.
1188 *PLOS ONE* **10**, e0145342 (2015).
- 1189 58. Siefert, J. C. *et al.* *Human Prostate Cancer-Associated Macrophage Subtypes with*
1190 *Prognostic Potential Revealed by Single-cell Transcriptomics.*
1191 <http://biorxiv.org/lookup/doi/10.1101/2020.06.19.160770> (2020)
1192 doi:10.1101/2020.06.19.160770.
- 1193 59. Azizi, E. *et al.* Single-Cell Map of Diverse Immune Phenotypes in the Breast Tumor
1194 Microenvironment. *Cell* **174**, 1293-1308.e36 (2018).

- 1195 60. Martinez, F. O. & Gordon, S. The M1 and M2 paradigm of macrophage activation:
1196 time for reassessment. *F1000prime Rep.* **6**, 13 (2014).
- 1197 61. Zhang, L. *et al.* Single-Cell Analyses Inform Mechanisms of Myeloid-Targeted
1198 Therapies in Colon Cancer. *Cell* **181**, 442-459.e29 (2020).
- 1199 62. Chakarov, S. *et al.* Two distinct interstitial macrophage populations coexist across
1200 tissues in specific subtissular niches. *Science* **363**, (2019).
- 1201 63. Ouzounova, M. *et al.* Monocytic and granulocytic myeloid derived suppressor cells
1202 differentially regulate spatiotemporal tumour plasticity during metastatic cascade.
1203 *Nat. Commun.* **8**, 1–13 (2017).
- 1204 64. Youn, J.-I. & Gabrilovich, D. I. The biology of myeloid-derived suppressor cells: The
1205 blessing and the curse of morphological and functional heterogeneity. *Eur. J.*
1206 *Immunol.* **40**, 2969–2975 (2010).
- 1207 65. Gabrilovich, D. I. Myeloid-derived suppressor cells. *Cancer Immunol. Res.* **5**, 3–8
1208 (2017).
- 1209 66. Calcinotto, A. *et al.* IL-23 secreted by myeloid cells drives castration-resistant
1210 prostate cancer. *Nature* **559**, 363–369 (2018).
- 1211 67. Fleming, V. *et al.* Targeting Myeloid-Derived Suppressor Cells to Bypass Tumor-
1212 Induced Immunosuppression. *Front. Immunol.* **9**, (2018).
- 1213 68. Karthaus, W. R. *et al.* Identification of Multipotent Luminal Progenitor Cells in Human
1214 Prostate Organoid Cultures. *Cell* **159**, 163–175 (2014).
- 1215 69. Drost, J. *et al.* Organoid culture systems for prostate epithelial and cancer tissue.
1216 *Nat. Protoc.* **11**, 347–358 (2016).

- 1217 70. Mani, S. A. *et al.* The epithelial-mesenchymal transition generates cells with
1218 properties of stem cells. *Cell* **133**, 704–715 (2008).
- 1219 71. Medeiros Tavares Marques, J. C. *et al.* Identification of new genes associated to
1220 senescent and tumorigenic phenotypes in mesenchymal stem cells. *Sci. Rep.* **7**,
1221 (2017).
- 1222 72. Hashimoto, S. *et al.* Comprehensive single-cell transcriptome analysis reveals
1223 heterogeneity in endometrioid adenocarcinoma tissues. *Sci. Rep.* **7**, (2017).
- 1224 73. McCray, T., Moline, D., Baumann, B., Vander Griend, D. J. & Nonn, L. Single-cell
1225 RNA-Seq analysis identifies a putative epithelial stem cell population in human
1226 primary prostate cells in monolayer and organoid culture conditions. *Am. J. Clin.*
1227 *Exp. Urol.* **7**, 123–138 (2019).
- 1228 74. Fraser, M. *et al.* Genomic hallmarks of localized, non-indolent prostate cancer.
1229 *Nature* **541**, 359–364 (2017).
- 1230 75. Barbieri, C. E. *et al.* The mutational landscape of prostate cancer. *Eur. Urol.* **64**,
1231 567–576 (2013).
- 1232 76. Baca, S. C. *et al.* Punctuated evolution of prostate cancer genomes. *Cell* **153**, 666–
1233 677 (2013).
- 1234 77. Manyak, M. J., Kikukawa, T. & Mukherjee, A. B. Expression of a uteroglobin-like
1235 protein in human prostate. *J. Urol.* **140**, 176–182 (1988).
- 1236 78. Liu, X. *et al.* Low CD38 Identifies Progenitor-like Inflammation-Associated Luminal
1237 Cells that Can Initiate Human Prostate Cancer and Predict Poor Outcome. *Cell Rep.*
1238 **17**, 2596–2606 (2016).

- 1239 79. Karthaus, W. R. *et al.* Regenerative potential of prostate luminal cells revealed by
1240 single-cell analysis. *Science* **368**, 497–505 (2020).
- 1241 80. Goldstein, A. S., Huang, J., Guo, C., Garraway, I. P. & Witte, O. N. Identification of a
1242 cell-of-origin for human prostate cancer. *Science* **329**, 568–571 (2010).
- 1243 81. Wang, X. *et al.* A luminal epithelial stem cell that is a cell of origin for prostate
1244 cancer. *Nature* **461**, 495–500 (2009).
- 1245 82. Krop, I. *et al.* Frequent HIN-1 promoter methylation and lack of expression in
1246 multiple human tumor types. *Mol. Cancer Res. MCR* **2**, 489–494 (2004).
- 1247 83. Ji, A. L. *et al.* Multimodal Analysis of Composition and Spatial Architecture in Human
1248 Squamous Cell Carcinoma. *Cell* **0**, (2020).
- 1249 84. Marjanovic, N. D. *et al.* Emergence of a High-Plasticity Cell State during Lung
1250 Cancer Evolution. *Cancer Cell* **38**, 229-246.e13 (2020).
- 1251 85. Puram, S. V. *et al.* Single-cell transcriptomic analysis of primary and metastatic
1252 tumor ecosystems in head and neck cancer. *Cell* **171**, 1611-1624.e24 (2017).
- 1253 86. Denmeade, S. R. & Isaacs, J. T. A history of prostate cancer treatment. *Nat. Rev.*
1254 *Cancer* **2**, 389–396 (2002).
- 1255 87. Teo, M. Y., Rathkopf, D. E. & Kantoff, P. Treatment of Advanced Prostate Cancer.
1256 *Annu. Rev. Med.* **70**, 479–499 (2019).
- 1257 88. Karantanos, T., Corn, P. G. & Thompson, T. C. Prostate cancer progression after
1258 androgen deprivation therapy: mechanisms of castrate-resistance and novel
1259 therapeutic approaches. *Oncogene* **32**, 5501–5511 (2013).
- 1260 89. Tiwari, R. *et al.* Androgen deprivation upregulates SPINK1 expression and
1261 potentiates cellular plasticity in prostate cancer. *Nat. Commun.* **11**, 384 (2020).

- 1262 90. Zhang, Y. *et al.* Androgen deprivation promotes neuroendocrine differentiation and
1263 angiogenesis through CREB-EZH2-TSP1 pathway in prostate cancers. *Nat.*
1264 *Commun.* **9**, 4080 (2018).
- 1265 91. Maynard, A. *et al.* Therapy-Induced Evolution of Human Lung Cancer Revealed by
1266 Single-Cell RNA Sequencing. *Cell* **182**, 1232-1251.e22 (2020).
- 1267 92. Team, C. Is the playing field level in prostate cancer? *Wellcome Sanger Institute*
1268 *Blog* [https://sangerinstitute.blog/2015/04/01/is-the-playing-field-level-in-prostate-](https://sangerinstitute.blog/2015/04/01/is-the-playing-field-level-in-prostate-cancer/)
1269 [cancer/](https://sangerinstitute.blog/2015/04/01/is-the-playing-field-level-in-prostate-cancer/) (2015).
- 1270 93. Pietrzak, K. *et al.* TIP5 primes prostate luminal cells for the oncogenic
1271 transformation mediated by PTEN-loss. *Proc. Natl. Acad. Sci.* **117**, 3637–3647
1272 (2020).
- 1273 94. Chua, C. W. *et al.* Single luminal epithelial progenitors can generate prostate
1274 organoids in culture. *Nat. Cell Biol.* **16**, 951–4 (2014).
- 1275 95. Waldman, B. S. *et al.* Identification of a Master Regulator of Differentiation in
1276 *Toxoplasma*. *Cell* **180**, 359-372.e16 (2020).
- 1277 96. Trivedi, U. H. *et al.* Quality control of next-generation sequencing data without a
1278 reference. *Front. Genet.* **5**, (2014).
- 1279 97. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinforma. Oxf. Engl.*
1280 **29**, 15–21 (2013).
- 1281 98. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in
1282 Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst.* **8**,
1283 329-337.e4 (2019).

- 1284 99. Koga, Y. *et al.* Genomic Profiling of Prostate Cancers from Men with African and
1285 European Ancestry. *Clin. Cancer Res.* (2020) doi:10.1158/1078-0432.CCR-19-4112.
- 1286 100. Huang, F. W. *et al.* Exome sequencing of African-American prostate cancer reveals
1287 loss-of-function ERF mutations. *Cancer Discov.* **7**, 973–983 (2017).
- 1288 101. Zappia, L. & Oshlack, A. Clustering trees: a visualization for evaluating clusterings
1289 at multiple resolutions. *GigaScience* **7**, (2018).
- 1290 102. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene
1291 expression data analysis. *Genome Biol.* **19**, 15 (2018).

1292

1293

1294

1295

1296 **Figures and Legends**

1297

1298 **Figure 1. PCa sample single-cell RNA-sequencing overview and identification of**

1299 **major cell types in localized prostate cancer. a.** Single-cell RNA-sequencing

1300 workflow on PCa biopsies, radical prostatectomy (RP) specimens, and *in vitro* organoid

1301 cultures grown from RP tumor specimens using the Seq-Well platform. **b.** Overview of

1302 major cell types identified within the combined dataset consisting of 21,743 cells from all

1303 biopsies (N = 6) and RP specimens (N = 12). Cell types are labeled in colors from

1304 corresponding clusters in the UMAP. **c.** Heatmap for the top 10 differentially expressed

1305 genes in each cell type. **d.** Cell type composition stacked bar chart by sample. Cell

1306 counts for each sample are normalized to 100%. Sample type is annotated (top) and

1307 patients are labeled below the x-axis. **e.** Cell composition comparison for each cell type

1308 among three sample types: biopsy patients (N = 3), RP tumor specimens (N = 8), and

1309 RP paired normal tissues (N = 4). Mean and confidence interval for each cell type are

1310 indicated in the grouped bar chart.

1311

1312 **Figure 2. Identification of tumor cells and major epithelial cell types including**

1313 **club cells. a.** UMAP projection of all 20 clusters identified in the epithelial cells. Clusters

1314 are labeled in the UMAP. **b.** Violin plots of representative marker genes across the

1315 clusters. **c.** UMAP of epithelial cells annotated by cell types. **d.** Heatmap for the top 10

1316 differentially expressed genes in each cell type. **e.** Club cell signature scores of each

1317 epithelial cell projected on the UMAP and signature score violin plots across all clusters.

1318 **f.** Box plots of club cell signature scores from normal club cells and lung club cells
1319 across epithelial cell types (***: $p < 0.001$, Wilcoxon rank sum test).

1320

1321 **Figure 3. Identification of PCa-enriched club cell states with upregulated**
1322 **androgen response signature. a.** UMAP of integrated club cells from PCa samples

1323 (Club PCa) and club cells from normal samples (Club Normal), color coded by cell
1324 states with differential gene expression profiles (left) and sample type (right). **b.** Violin

1325 plots of representative marker genes between the two types of club cells. **c.** Heatmap
1326 for the top 10 differentially expressed genes in each cell state. **d.** Grouped bar chart

1327 comparison of 6 cell state compositions between Club PCa and Club Normal.

1328 Significance levels are labeled (***: FDR $q < 0.001$, Wilcoxon rank sum test). **e.** Volcano

1329 plots of the overexpressed genes in Club cell cluster 0 and other cell states within the
1330 PCa samples. **f.** Top 20 upregulated signaling pathways between Club cell cluster 0 and

1331 the other club cells on Hallmark gene set collection ($N = 50$) within the PCa samples.

1332 Gene counts for the corresponding gene set indicated by marker radius. Statistical

1333 significance levels (FDR) are shown by color gradient. **g.** Comparison of LE signature
1334 scores between Club cluster 0 and other club cells (***: $p < 0.001$, Wilcoxon rank sum

1335 test) within the PCa samples. **h.** Violin plot comparison between Club cluster 0, other
1336 club cells and LE for multiple LE markers within the PCa samples. **i.** Schematic marker

1337 of gene expression changes between Club Normal and Club PCa. Gene downregulation
1338 and upregulation in Club PCa compared to Club Normal represented by red and green

1339 arrows. Proportion of Club cell cluster 0 within all club cells represented by area in blue
1340 and characterized by its LE-like state and high-level expression of *LTF* and *NKX3-1*.

1341
1342 **Figure 4. Integration of BE and LE cells identifies tumor-associated cell states**
1343 **enriched in the PCa samples. a.** UMAP of integrated BE cells labeled by cell states
1344 (left) or samples type (BE PCa and BE Normal) (right). **b.** Cell composition comparison
1345 between BE PCa and BE Normal. **c.** PCa and normal enriched cell states 4 and 6
1346 highlighted in the integrated BE UMAP. **d.** Top 20 upregulated signaling pathways
1347 between cluster 6 and the other BE on C2 canonical gene set (C2CP) collection (N =
1348 2,332). Gene counts for the corresponding gene set are indicated by marker radius.
1349 Statistical significance levels (FDR) are shown by color gradient. Pathways associated
1350 with PCa tumor progression and invasiveness are highlighted in red. **e.** Volcano plots of
1351 the overexpressed genes in BE cluster 6 and other BE cell states within the PCa
1352 samples. **f.** Distribution of BE cluster 6, other BE and LE on the overall epithelial cell
1353 UMAP. **g.** Violin plot comparison between BE cluster 6, other BE and LE for multiple LE
1354 markers within the PCa samples. **h.** Comparison of Hallmark AR pathway signature and
1355 LE signature scores within the PCa samples (***: $p < 0.001$, Wilcoxon rank sum test).

1356
1357 **Figure 5. Integration of PCa and normal epithelial cells reveals common AR**
1358 **signaling upregulation driven by PCa-enriched BE and club cell states. a.** UMAP
1359 of integrated epithelial cells annotated by cell types and sample type (PCa and Normal),
1360 then separated by the origin (either previous normal epithelial cells or epithelial cells in
1361 the PCa samples). **b.** Heatmaps of top 20 differentially expressed genes between PCa
1362 samples and normal prostates for adjacent cell types (left: BE PCa, BE Normal. Middle:
1363 Club Normal, Club PCa. Right: LE PCa, LE Normal). Commonly upregulated genes in
1364 the PCa samples are labeled in red, and commonly upregulated genes in the normal

1365 samples are labeled in green. **c.** Top, *AR* expression percentages in all epithelial cell
1366 types within the integrated dataset. Significance levels are labeled in each comparison
1367 (***: $p < 0.001$, FDR). Bottom, Comparison of Hallmark *AR* pathway signature scores of
1368 each epithelial cell type. Significance levels are labeled for each common cell type (*: p
1369 < 0.05 , ***: $p < 0.001$, Wilcoxon rank sum test). **d.** The association of *AR* signature with
1370 BE and club cell state. Each cell is labeled (grey: 0, not in the cell state; black: 1, in the
1371 cell state). Information coefficient, accompanied p -values and FDR q values are labeled
1372 next to each cell state. **e.** The association of *AR* signature with BE and club cell state
1373 signature scores in the TCGA datasets ($N = 491$). Information coefficient, accompanied
1374 p -values and FDR q values are labeled next to each cell state.

1375

1376 **Figure 6. Comparison of *ERG+* and *ERG-* tumor cells reveals patient-specific cell**
1377 **states and intra-patient heterogeneity. a.** UMAP of *ERG+* tumor cells labeled by
1378 clusters with differential gene expression profiles (top). Heatmap of the top 10
1379 differentially expressed genes for each cluster (bottom). **b.** UMAP of *ERG-* tumor cells
1380 labeled by clusters with differential gene expression profiles (top). Heatmap of the top
1381 10 differentially expressed genes for each cluster (bottom). **c.** Patient composition in
1382 each cluster for *ERG+* tumor cells (top) and *ERG-* tumor cells (bottom). Cell counts in
1383 each cluster are normalized to 100%. **d.** UMAP of *ERG+* and *ERG-* tumor cells when
1384 integrated with non-malignant LE cells respectively. **e.** UMAP of *ERG+* and *ERG-* tumor
1385 cells when integrated with non-malignant LE cells labeled by patients. **f.** The association
1386 of *TMPRSS2-ERG* fusion status in the TCGA ($N = 290$) and SU2C ($N = 266$) datasets
1387 with *ERG+* and *ERG-* tumor cell signature (red: *TMPRSS2-ERG* fusion detected; blue:

1388 *TMPRSS2-ERG* fusion not detected). Information coefficient, accompanied p-values
1389 and FDR q values are labeled. **g.** Visualization of the intersection amongst significant
1390 GSEA results for BE, LE and club cells. The color intensity of the bars represents the p-
1391 value significance of the intersections.

1392

1393 **Figure 7. CD4 T subsets associated with *ERG* status and common upregulation of**

1394 **PD-1 and interferon gamma signaling in the *ERG*- tumor microenvironment. a.**

1395 UMAP of T-cells labeled by different cell types (left) and *ERG*+ or *ERG*- patients (right).

1396 **b.** Cell composition comparison between *ERG*+ and *ERG*- patients for all T-cell cell

1397 types. Significance levels are labeled in differentially enriched clusters. **c.** UMAP of

1398 stromal cells labeled by different cell types (left) and *ERG*+ or *ERG*- patients (right). **d.**

1399 Cell composition comparison between *ERG*+ and *ERG*- patients for all stromal cell

1400 types. Significance levels are labeled in differentially enriched clusters. **e.** Visualization

1401 of the intersections amongst significantly upregulated (top) and downregulated (bottom)

1402 gene sets within C2 CP gene set collection for tumor cells, two clusters of differentially

1403 enriched CD4 T-cell clusters and stromal cells. Significant GSEA results are

1404 represented by circle below bar chart with individual blocks showing “presence” (green)

1405 or “absence” (grey) of the gene sets in each intersection. P-value significance of the

1406 intersections are represented by color intensity of the bars. **f.** GSEA results for the

1407 *ERG*- patient-enriched CD4 T-cell cluster compared to the *ERG*+ patient-enriched

1408 cluster on the common upregulated gene sets (N = 14). Gene counts for the

1409 corresponding gene set are indicated by marker radius. Statistical significance levels

1410 (FDR) are shown by color gradient. Reactome PD-1 and Interferon gamma signaling

1411 pathways are highlighted in red. **g.** Gene expression heatmaps of genes in the
1412 Reactome PD-1 and Interferon gamma signaling pathways for tumor cells, CD4 T-cells
1413 and stromal cells in both *ERG+* and *ERG-* patients.

1414

1415 **Figure 8. *In vitro* organoid samples recapitulate PCa-enriched BE and club cell**

1416 **states. a.** UMAP of cells from organoid samples labeled by different cell types.

1417 Organoid culture snapshots are depicted in the upper right panel. **b.**

1418 Immunofluorescence staining for LE marker (*KRT8*), BE marker (*KRT5*) and club cell

1419 markers (*SCGB1A1*, *LTF*) of the organoid samples. **c.** UMAP of integrated dataset of

1420 cells from the organoid samples and epithelial cells from matching parent tissue

1421 samples, labeled by cell types. **d.** UMAP of integrated dataset, labeled by sample types

1422 (tissue or organoid samples). **e.** Heatmaps for the top 20 differentially expressed genes

1423 for BE and club cells between tumor tissues and organoid samples. **f.** UMAP of

1424 integrated club cell dataset of tumor tissue and organoid samples. Cell composition

1425 comparison is shown in the grouped bar charts. **g.** Dot plots of the top 10 differentially

1426 expressed genes in cluster 3, 4 and 7 in tissue and organoid club cells. Dot size

1427 represents proportions of gene expression in cells and expression levels are shown by

1428 color shading (low to high reflected as light to dark).

1429

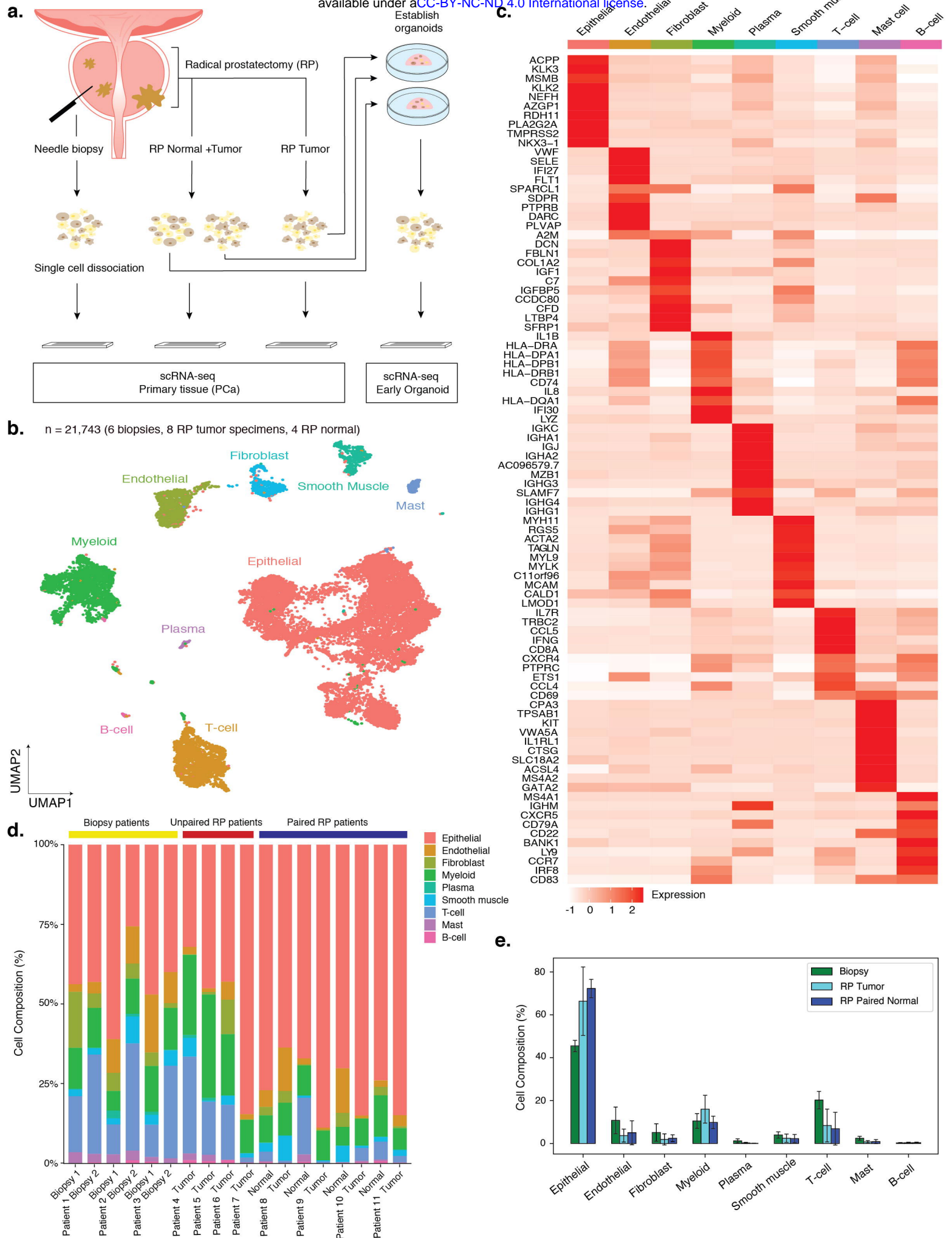


Figure 2

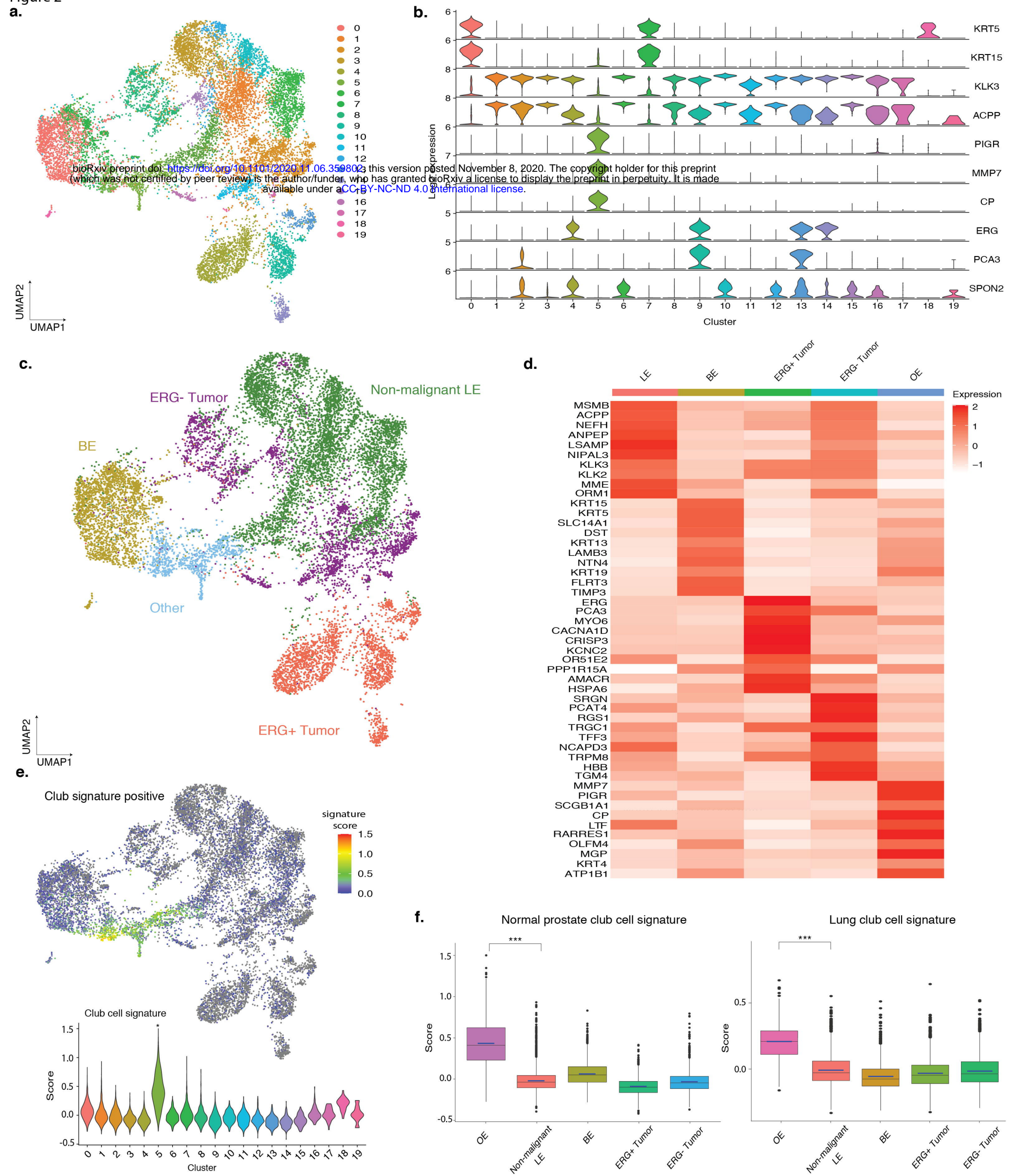


Figure 3

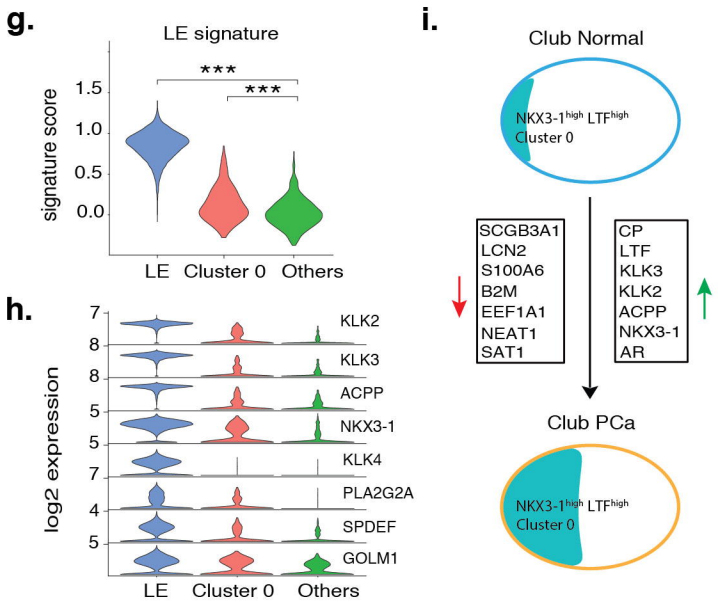
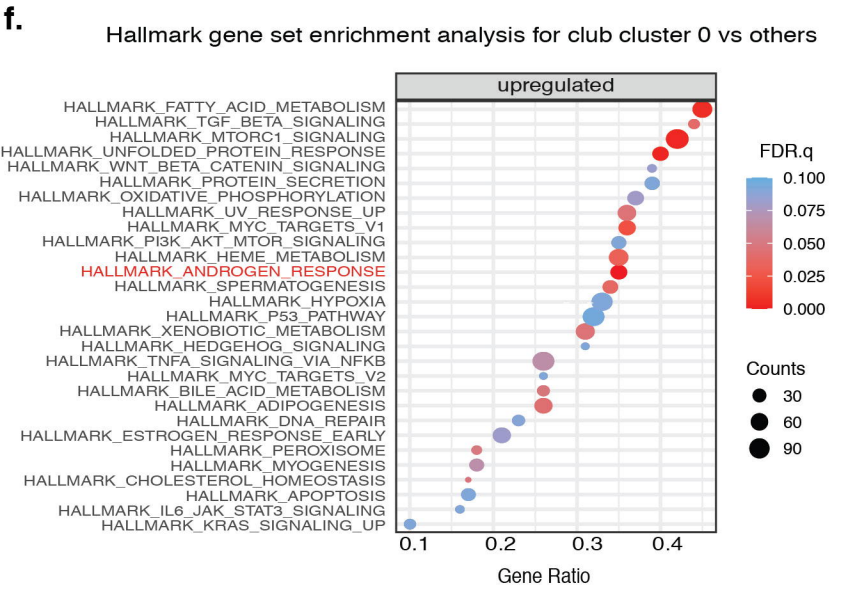
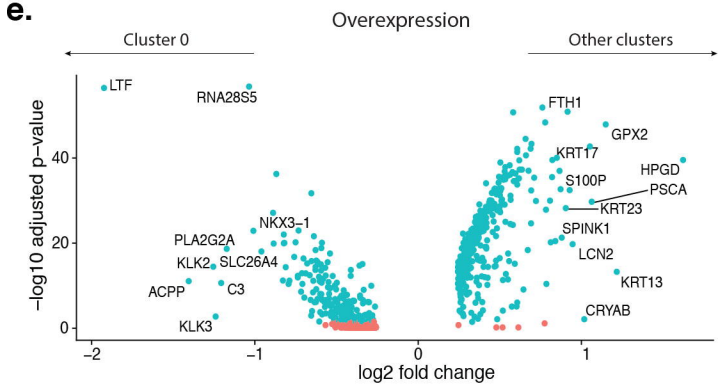
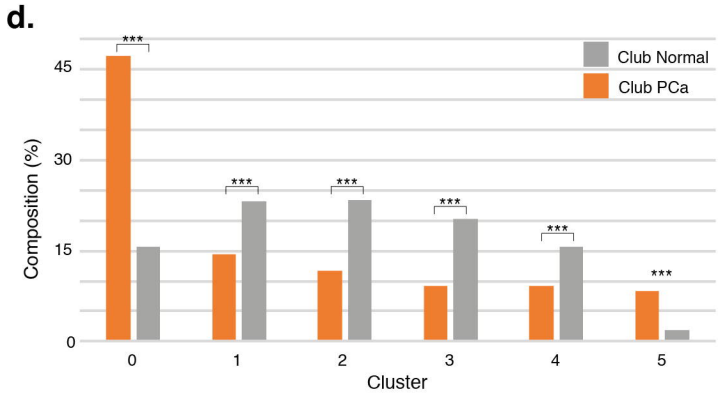
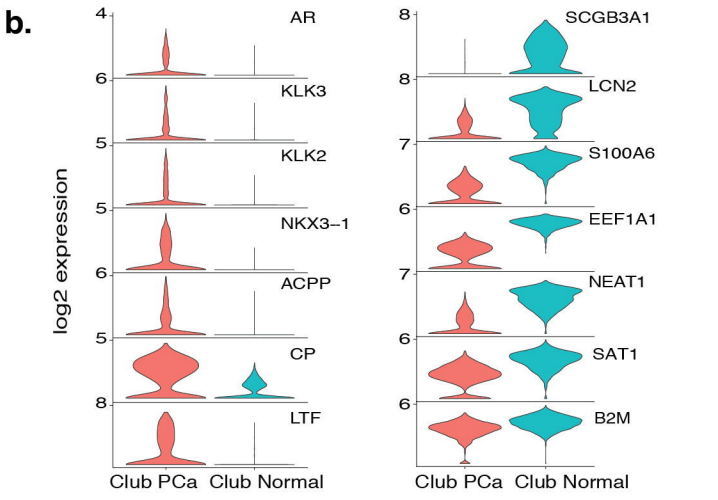
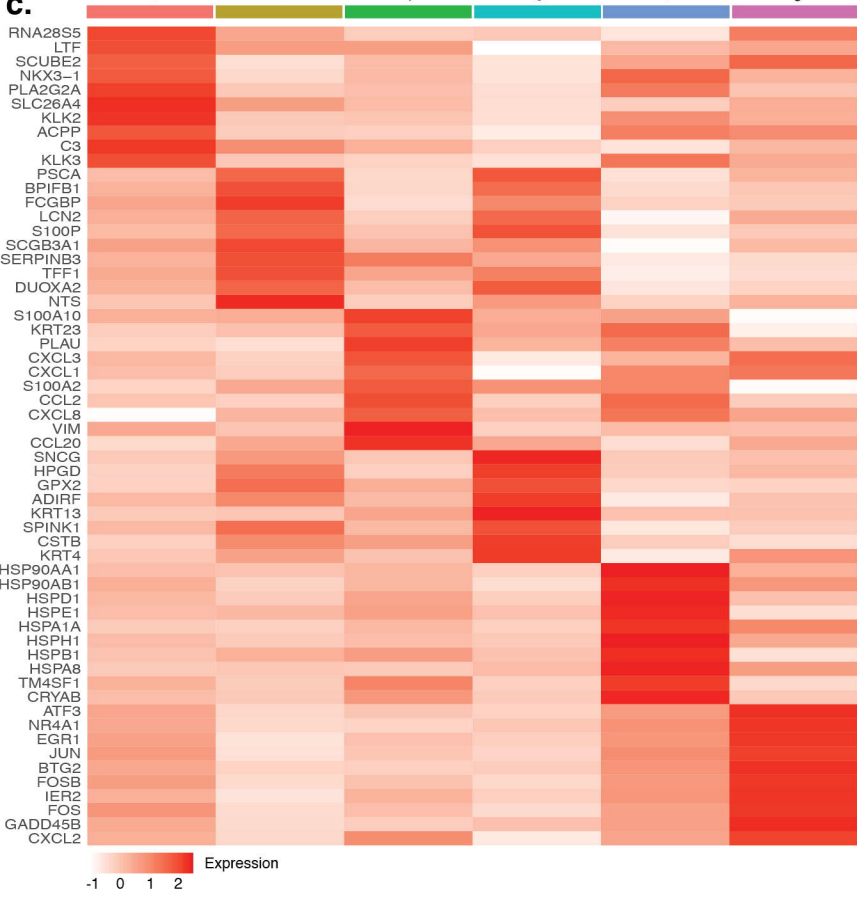
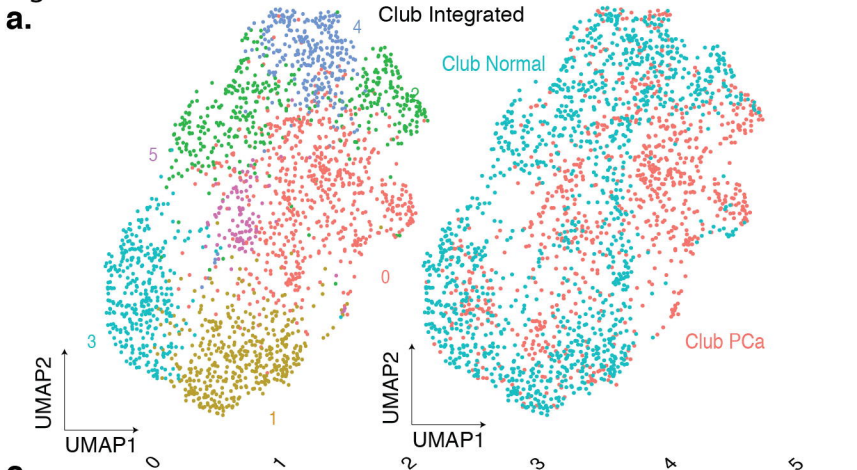
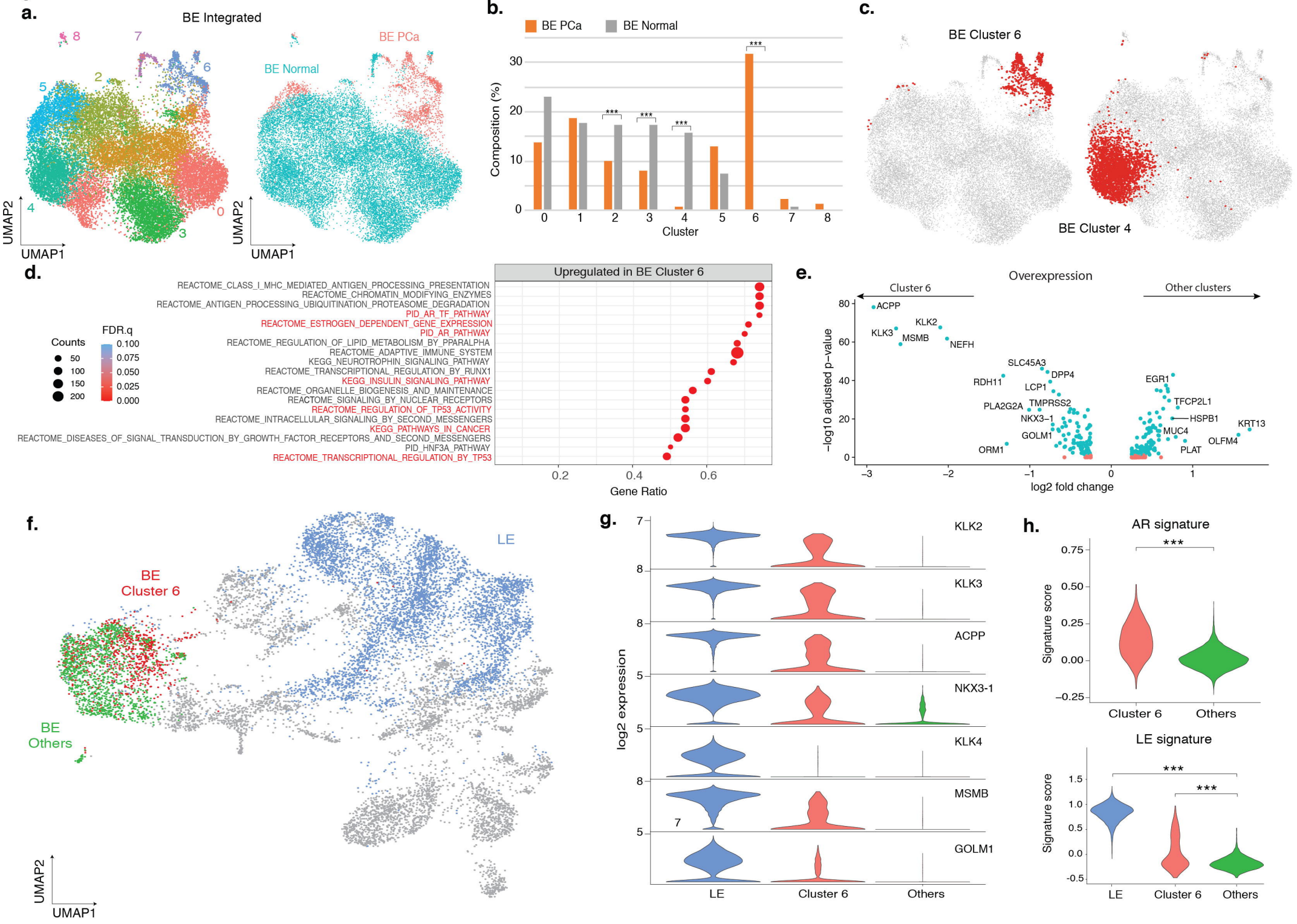
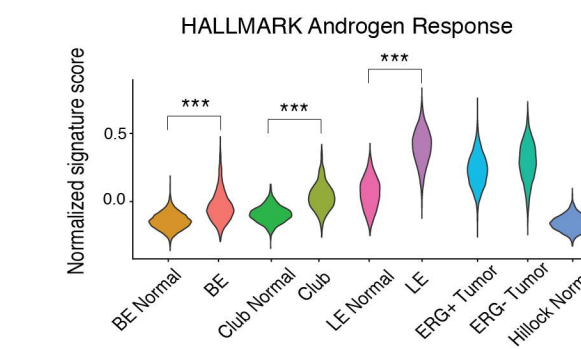
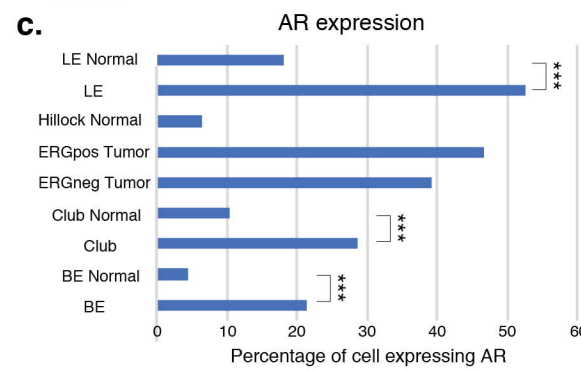
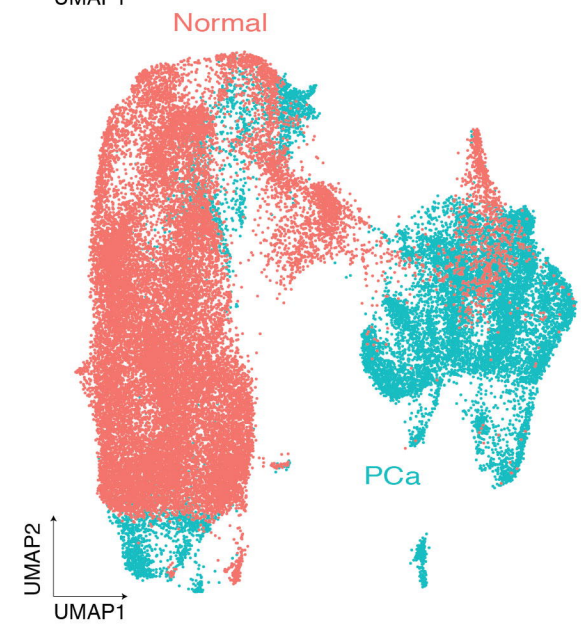
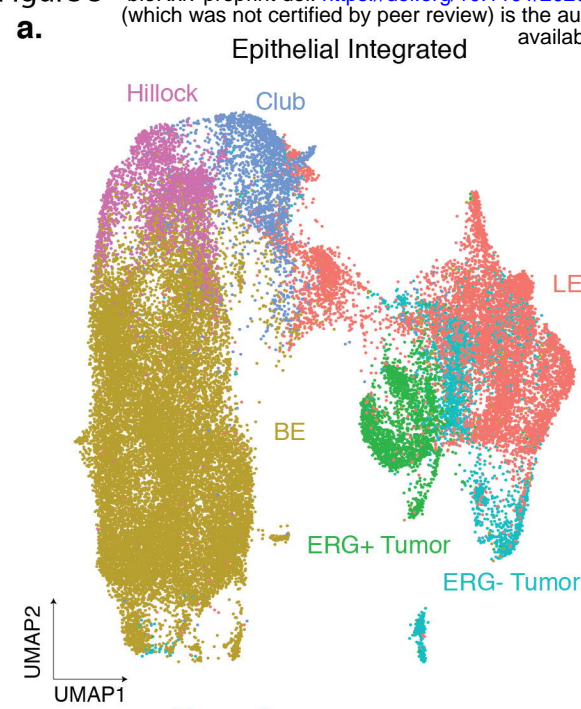


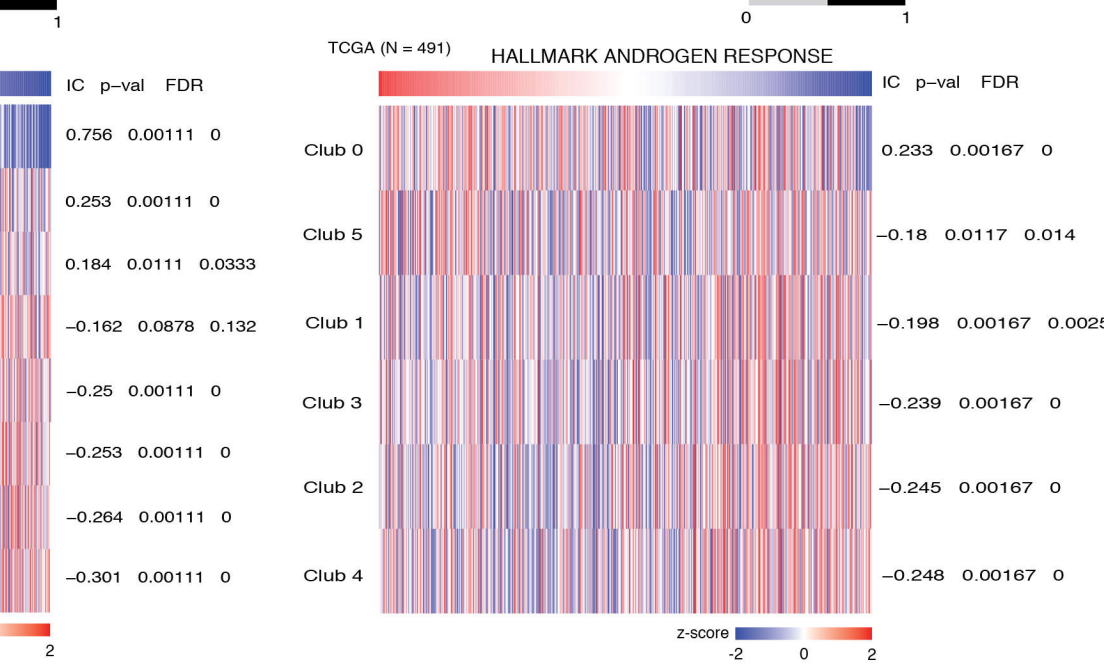
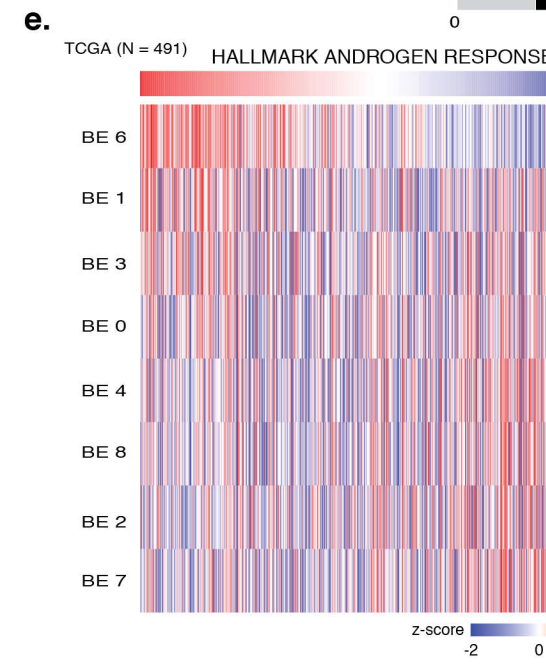
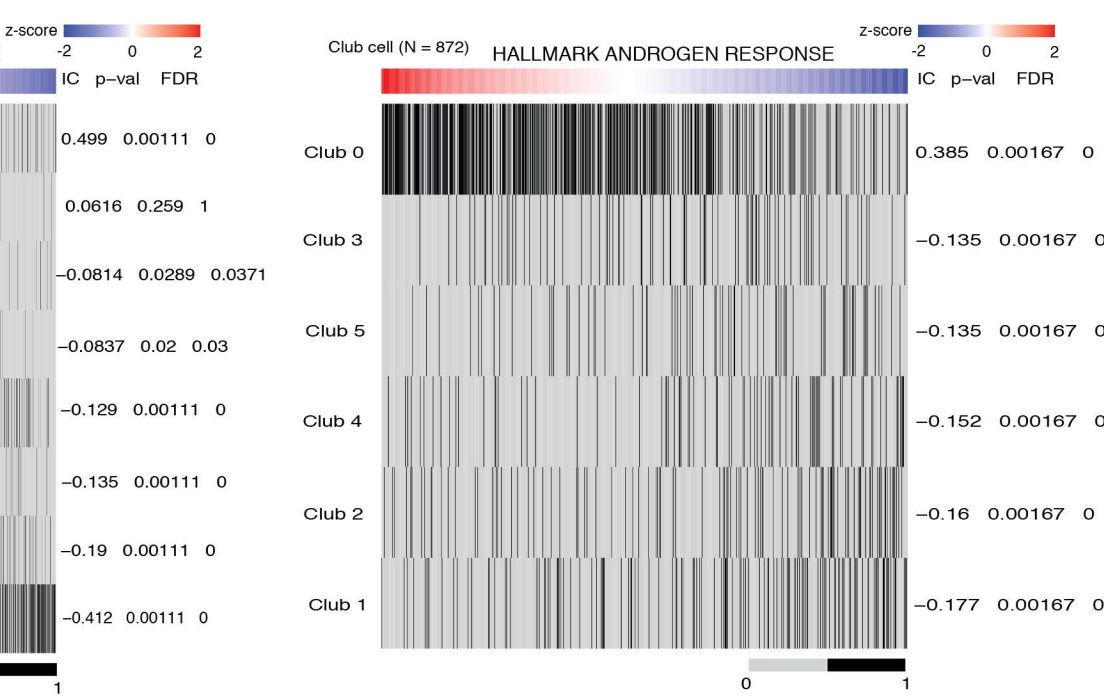
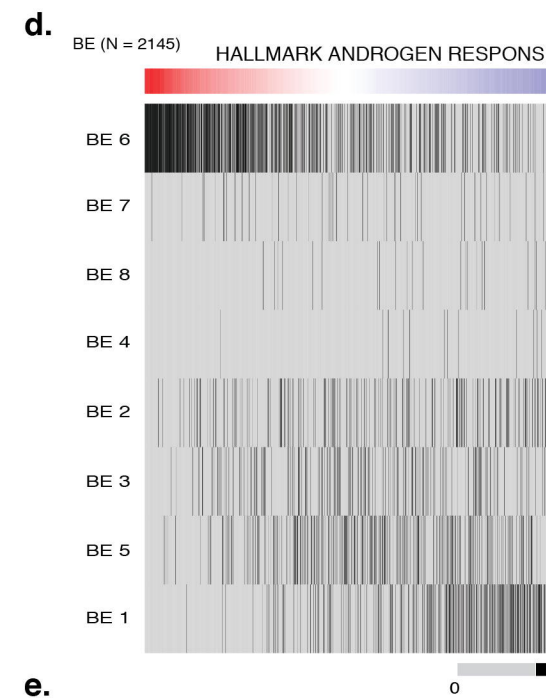
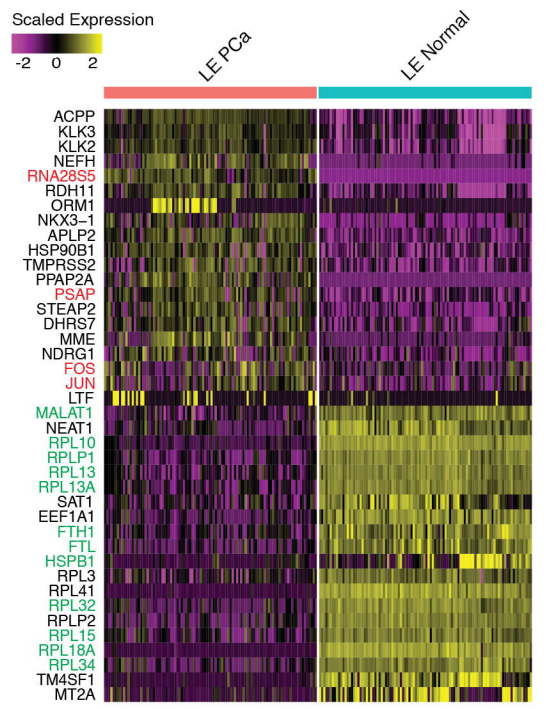
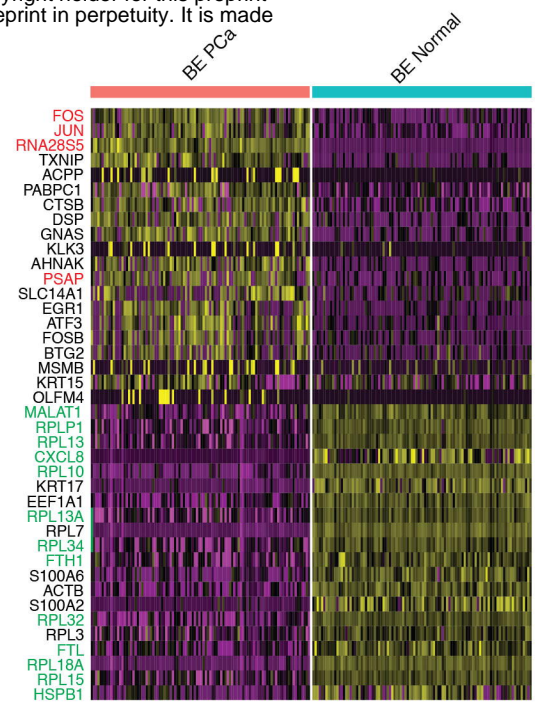
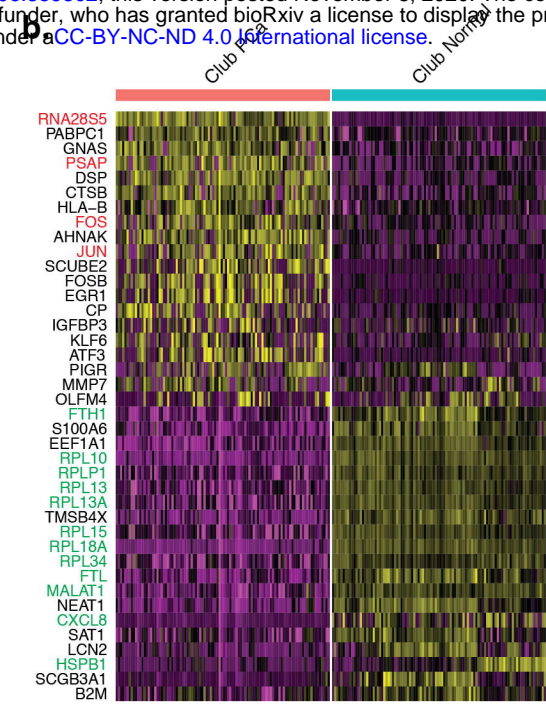
Figure 4

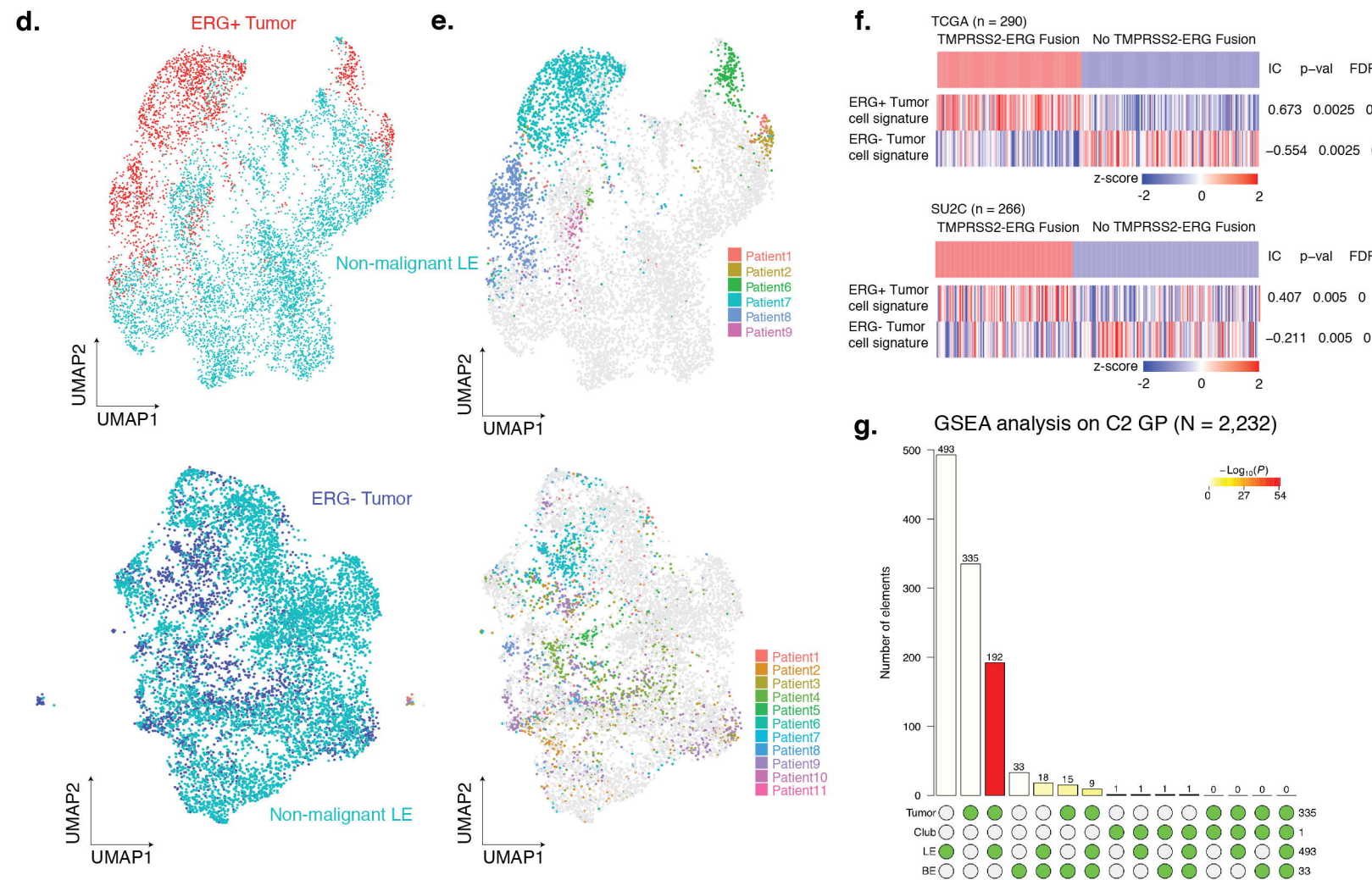
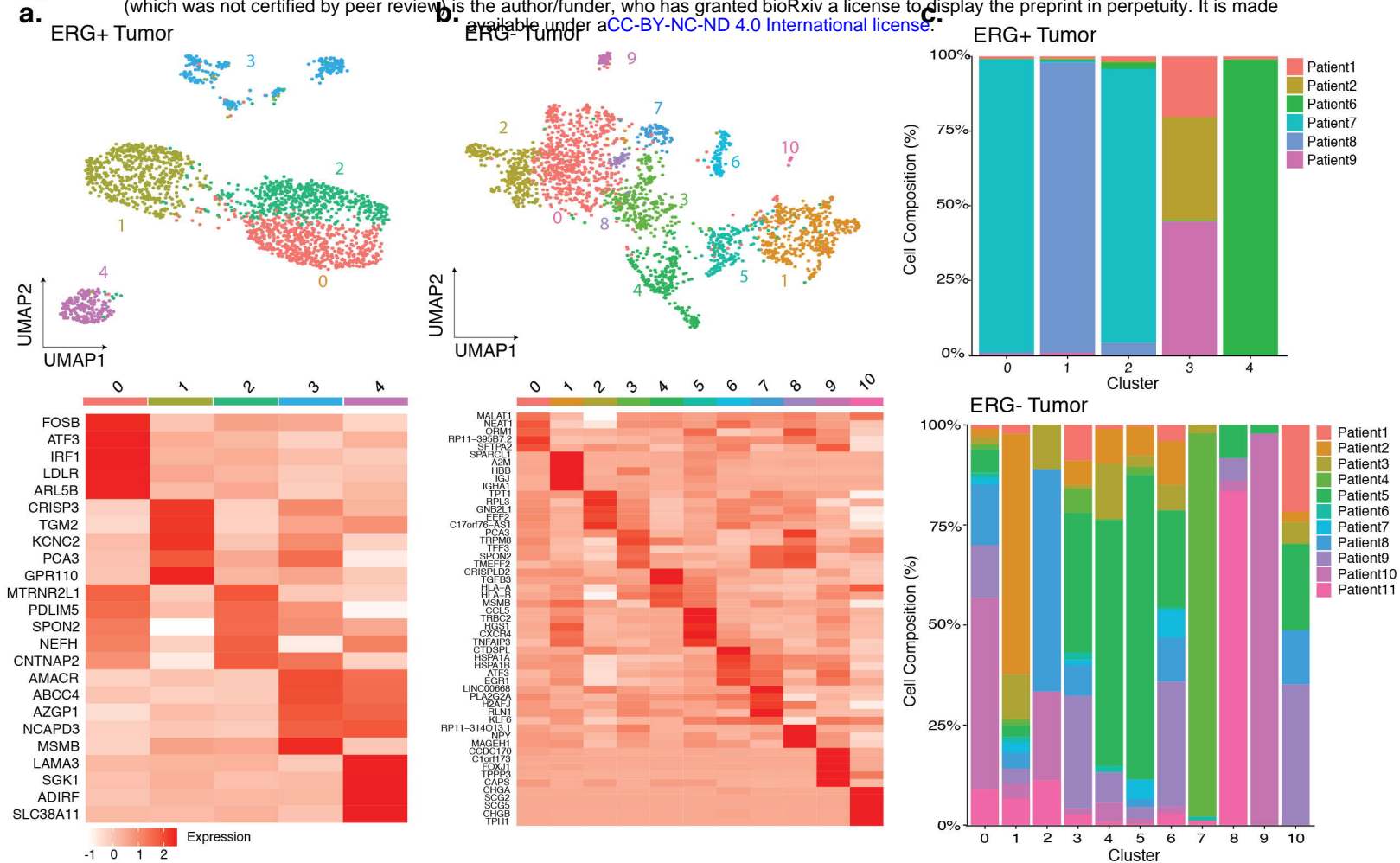


a.

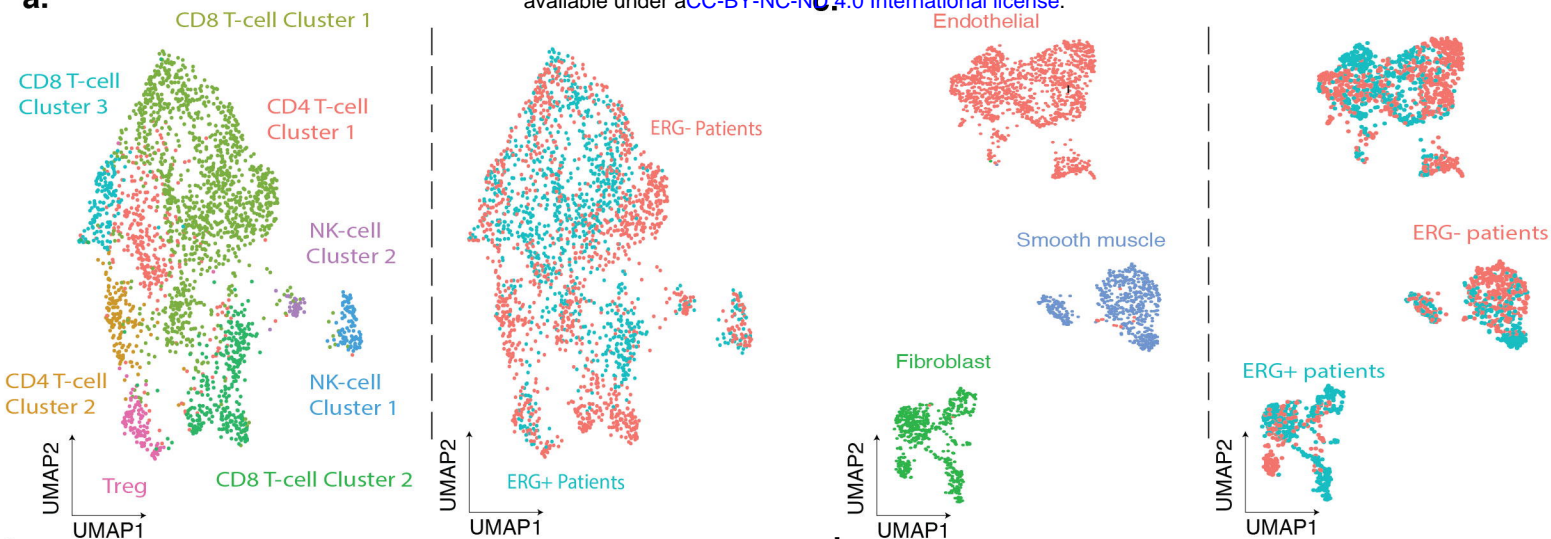


b.

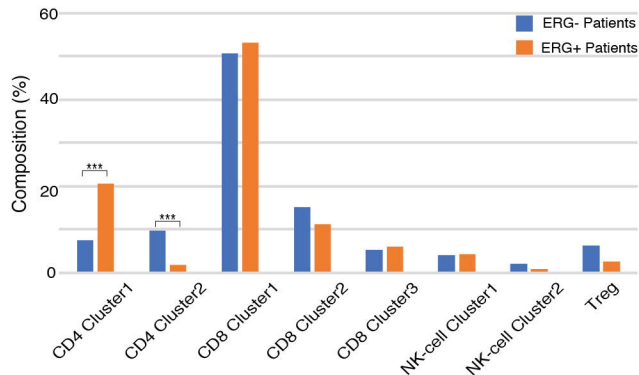




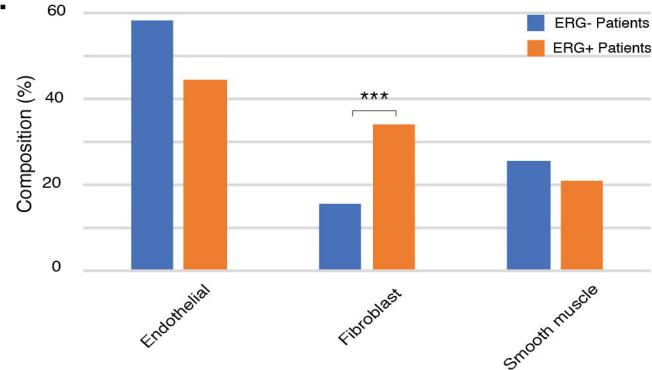
a.



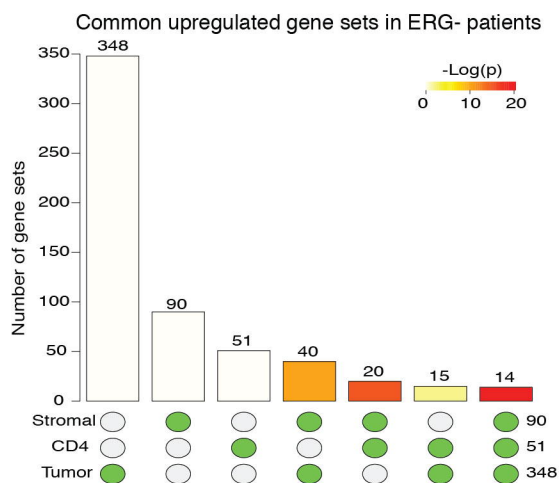
b.



d.

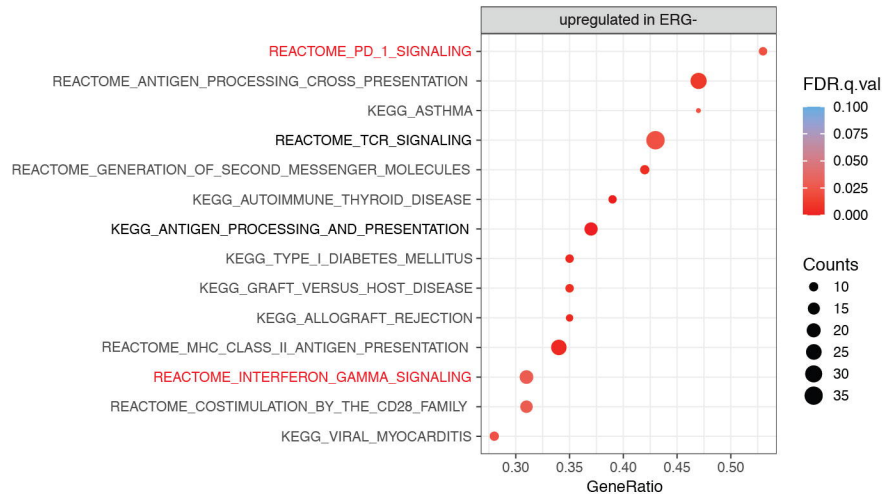


e.



f.

C2CP GSEA Results for ERG- enriched vs ERG+ enriched CD4 T-cell



g.

