# Redefining colorectal cancer classification and clinical stratification through a single-cell atlas

Ateeq M. Khaliq[1,*], Zeyneb Kurt[2,*], Miles W. Grunvald[1,*], Cihat Erdogan[3], Sevgi S. Turgut[4], Tim Rand[5], Sonal Khare[5], Jeffrey A. Borgia[1], Dana M. Hayden[1], Sam G. Pappas[1], Henry R. Govekar[1], Anuradha R. Bhama[1], Ajaypal Singh[1], Richard A. Jacobson[1], Audrey E. Kam[1], Andrew Zloza[1], Jochen Reiser[1], Daniel V. Catenacci[6], Kiran Turaga[6], Milan Radovich[7], Vineet Gupta[1], Ram Al-Sabti[1], Sheeno Thyparambil[8], Mia A. Levy[1], Janakiraman Subramanian[9], Timothy M. Kuzel[1], Anguraj Sadanandam[10], Arif Hussain[12], Bassel El-Rayes[12], Ameen A. Salahudeen[5✉],  Ashiq Masood[1✉]


**Affiliations:**
[1]Rush University Medical Center, Chicago, IL, USA.

[2]Northumbria University, Newcastle Upon Tyne, UK.

[3]Isparta University of Applied Sciences, Isparta, Turkey.

[4]Yildiz Technical University, Istanbul, Turkey.

[5]Tempus Labs, Inc., Chicago, IL, USA.

[6]The University of Chicago, Chicago, IL, USA.

[7]Indiana University School of Medicine, Indianapolis, IN, USA.

[8]mProbe Inc. Rockville, Maryland, USA.

[9]University of Missouri-Kansas City School of Medicine, Kansas City, MO, USA.

[10]Institute of Cancer Research, London, UK.

[11]University of Maryland Marlene and Stewart Greenebaum Comprehensive Cancer Center, Baltimore, MD USA.

[12]Emory University Winship Cancer Institute, Atlanta, GA, USA.

*These authors contributed equally: Ateeq M. Khaliq, Zeyneb Kurt and Miles W. Grunvald.

✉e-mail: ashiq_masood@rush.edu; ameen@tempus.com

51 **ABSTRACT**

52 Colorectal cancer (CRC), a disease of high incidence and mortality, has had few treatment advances owing

53 to a large degree of inter- and intratumoral heterogeneity. Attempts to classify subtypes of colorectal cancer

54 to develop treatment strategies has been attempted by Consensus Molecular Subtypes (CMS) classification.

55 However, the cellular etiology of CMS classification is incompletely understood and controversial. Here,

56 we generated and analyzed a single-cell transcriptome atlas of 49,859 CRC cells from 16 patients, validated

57 with an additional 31,383 cells from an independent CRC patient cohort. We describe subclonal

58 transcriptomic heterogeneity of CRC tumor epithelial cells, as well as discrete stromal populations of

59 cancer-associated fibroblasts (CAFs). Within CRC CAFs, we identify the transcriptional signature of

60 specific subtypes (CAF-S1 and CAF-S4) in more than 1,500 CRC patients using bulk transcriptomic data

61 that significantly stratifies overall survival in multiple independent cohorts. We also uncovered two CAF-

62 S1 subpopulations, ecm-myCAF and TGFß-myCAF, known to be associated with primary resistance to

63 immunotherapies. We demonstrate that scRNA analysis of malignant, stromal, and immune cells exhibit a

64 more complex picture than portrayed by bulk transcriptomic-based Consensus Molecular Subtypes (CMS)

65 classification. By demonstrating an abundant degree of heterogeneity amongst these cell types, our work

66 shows that CRC is best represented in a transcriptomic continuum crossing traditional classification systems

67 boundaries. Overall, this CRC cell map provides a framework to re-evaluate CRC tumor biology with

68 implications for clinical trial design and therapeutic development.

69

70

71

72

73

74

75

76

## INTRODUCTION

Colorectal cancer (CRC) is the third most commonly diagnosed cancer and a leading cause of cancer-related mortality worldwide[1,2]. Approximately one-third of patients experience disease relapse following curative-intent surgical resection and chemotherapy[3]. Despite the high incidence and mortality of advanced CRC, few effective therapies have been approved in the past several decades[4]. One barrier to the development of efficacious therapeutics is the biological heterogeneity of CRC and its variable clinical course. While landmark studies from The Cancer Genome Atlas (TCGA) have defined the somatic mutational landscape within CRC, several studies have shown that stromal signatures, including fibroblasts and cytotoxic T cells, are likely the main drivers of clinical outcomes[5–9]. These findings suggest that the clinical phenotypes of CRC and by extension, its tumor biology is shaped by a complex niche of heterotypic cell interactions within the tumor microenvironment (TME).

Bulk gene expression analyses by several independent groups have identified distinct CRC subtypes[10–12]. Based on bulk transcriptomic signatures, an international consortium published the Consensus Molecular Subtypes (CMS) that classified CRC into CMS1 (MSI immune), CMS2 (canonical), CMS3 (metabolic), and CMS4 (mesenchymal) subtypes[12]. Unfortunately, associations between CMS and meaningful therapeutic responses to specific agents have been inconsistent across studies[13–16]. Further, there has been a lack of concordance between primary and metastatic CRC tumors within the CMS framework, limiting its overall utility in clinical decision making[16–18]. Thus, an improved CMS classification, or an alternative classification system, is needed to improve clinical utility.

To overcome the limitations of bulk-RNA sequence profiling, we utilized single-cell RNA sequencing (scRNA-seq) to more thoroughly evaluate the CRC subtypes at the molecular level, including within the context of the currently defined CMS classification. We dissected heterotopic cell states of tumor epithelia and stromal cells, including a cancer-associated fibroblast (CAF) population. The CAF population's clinical

102  and prognostic significance became apparent when CAF signatures were applied to large, independent CRC

103  transcriptomic cohorts.

104

105  **RESULTS**

106  We profiled sixteen primary CRC tumor tissue samples and eight adjacent, normal, colonic tissue samples

107  (24 in total) using droplet-based, scRNA-seq. Altogether, we captured and retained 49,589 high-quality

108  single cells after performing quality control for downstream analysis (**Fig. 1a, Supplementary Table 1**).

109  All scRNA-seq data were merged and normalized to identify robust discrete clusters of epithelial cells

110  (*EPCAM+, KRT8+,* and *KRT18+*), fibroblasts (*COL1A2+*), endothelial cells (*CD31+*), T cells (*CD3D+),*

111  B cells *(CD79A+),* and myeloid cells (*LYZ+*) using canonical marker genes **(Fig. 1b-c).** Additionally, each

112  cell type compartment was analyzed separately. Clustree (v0.4.1) and manual review of differentially

113  expressed genes in each subcluster were studied to choose the best cluster resolution without cluster

114  destabilization (see methods)[19]. Cell population designation was chosen by specific gene expression, and

115  SingleR was also utilized for unbiased cell type recognition (see methods)[20–23].

116

117  In addition to cancer cells, we identified diverse TME cell phenotypes, including fibroblasts subsets (*CAF-*

118  *S1* and *CAF-S4*), endothelial cells, CD4+ subsets (naïve/memory, Th17, and Tregs),  CD8+  subsets

119  (naïve/memory, cytotoxic, tissue-resident memory, and Mucosa-Associated Invariant (MAIT) cells),  NK

120  cells, innate lymphoid cell (ILC) types, B cell phenotypes (naïve, memory, germinal center, and plasma

121  cells), and monocyte lineage phenotypes (C1DC+ dendritic cells, proinflammatory monocytes [IL1B, IL6,

122  S100A8, and S100A9]), and M2 polarized anti-inflammatory [*CD163, SEPP1, APOE,* and *MAF]*), tumor-

123  associated macrophages (TAMs) (**Fig. 1d, Extended Data Figs. 1-3, Extended Data Tables 1-4**)[20–22,24].

124

125  We also profiled an independent CRC dataset for validation and retained 31,863 cells after quality control

126  and strict filtering of cells expressing hybrid markers[25]. The re-clustering of individual compartments

127  further refined our analysis, and cells that expressed hybrid or distinct lineage markers within a cluster were

128    removed from the downstream analysis (**Supplementary Fig. 4-5, 9-10**) (see methods). Thus, a total of

129    81,242 high-quality cells were profiled to produce a single-cell map of 39 colorectal cancer patients. The

130    results of the primary CRC cohort (49,859 single-cells) are available at the Colon Cancer Atlas

131    (*www.websiteinprogress.com*).

132

133    **Malignant colon cancer reveals tumor epithelial cell subclonal heterogeneity and stochastic behavior.**

134    We detected 8,965 tumor and benign epithelial cells (*EPCAM+, KRT8+,* and *KRT18+*) and, on re-

135    clustering, produced 17 epithelial clusters (designated C1 to C17) (**Fig. 2a-b**). Tumor cells were confirmed

136    to be of malignant origin by inferring chromosomal copy number alterations (**Supplementary Fig. 1**).

137    Clusters were chiefly influenced by colonic epithelial markers, including those for stemness (*LGR5, ASCL2,*

138    *OLFM4,* and *STMN1*), enterocytes (*FABP1* and *CA2*), goblet cells (*ZG16, MUC2, SPINK4,* and *TFF3*), and

139    enteroendocrine cells (*PYY* and *CHGA*) (**Supplementary Table 2**). Tumor cells exhibited a high degree of

140    de-differentiated state of plasticity, possibly accounting for lasting cancer growth (**Supplementary Fig.**

141    **2b**)[26]. Distinct tumor-derived clusters were predominantly patient-specific, reflecting a high degree of inter-

142    patient tumoral cell heterogeneity (**Fig. 1d**). In contrast, epithelial populations derived from normal colon

143    tissue samples across multiple patients clustered together, a pattern observed in previous studies confirming

144    both normal tissue homeostasis and limited sample batch effects (**Fig. 1d**)[27,28].

145

146    We next aimed to identify gene expression programs shared across these clusters using hallmark pathway

147    analysis[29]. A strong overlap was observed for multiple pathways such as activation of inflammatory,

148    epithelial-mesenchymal transformation (EMT), immune response, and metabolic pathways (**Fig. 2b**).

149    Interestingly, high microsatellite instability (MSI-H) and microsatellite stable (MSS) CRC tumors,

150    considered clinically separate entities, demonstrated similar pathway program activation within the tumor

151    epithelial populations **(Fig. 2b).** Some clusters also showed activation of unique pathways such as activation

152    of apical junctions and angiogenesis (C6), hypoxia and fatty acid metabolism (C11) and Notch signaling

153    and DNA repair (C14), among others **(Fig. 2b).** However, MSI-H tumors differed from MSS tumors based

154    on immune cell infiltration (**Extended Data Fig. 1**).

155

156    Since intratumoral heterogeneity is recognized as a key mechanism contributing to drug resistance, cancer

157    progression, and recurrence, we next focused on dissecting potential transcriptomic states to identify

158    heterogeneity within each  tumor[30–32]. We found that each tumor specimen contained 2-10 distinct tumor

159    epithelial clusters (**Fig. 1d**). Gene set variation analysis (GSVA) was performed on cells from individual

160    tumor samples and illustrated the sub-clonal transcriptomic heterogeneity within each specimen

161    (**Supplementary Fig. 2c**)[33]. Clusters identified in individual pathway analysis demonstrated the up- or

162    down-regulation of crucial metabolic and oncogenic pathways between samples, suggesting wide

163    phenotype variations between cells from the same tumor[34].

164

165    Given the evidence of intratumoral epithelial heterogeneity, we next performed trajectory inference using

166    pseudotime analysis to identify potential alignments or lineage relationships (i.e., right versus left-sided

167    CRC), CMS classification, or MSI status[35,36]. This analysis also served as a control for inter-patient

168    genomic heterogeneity and provided an orthogonal strategy to confirm the transcriptomic trends we

169    identified. We detected five molecular states (S1n/t to S5n/t) with malignant and normal epithelial cells

170    intermixed  and aligned along a common transcriptional trajectory **(Fig. 2c)**[37,38].  In both normal and

171    tumor cells, each transcriptional state pathway activation was shared and related to colon epithelial

172    function of nutrient absorption, maintaining the colon homeostasis, and the activation of cancer-related

173    pathways such as apoptosis and cell development[39]. However, these cell states showed differential gene

174    enrichment activation between normal and tumor cells (**Supplementary Fig. 3, Supplementary Table**

175    **3)**.

176

177    Additionally, tumor cells showed upregulation of embryogenesis (S2t), consistent with previous findings

178    that tumor cells revert to their embryological states in cancer development (**Supplementary Table 3**)[40].

179  Interestingly, there were no significant associations with anatomic location, CMS classification, or MSI

180  status within our dataset or an independent dataset of 31,383 single cells (**Fig. 2d**, **Supplementary Fig.**

181  **4**)[25]. Hence, in our analysis, CRC oncogenesis represent the hijacking of the normal epithelial

182  differentiation pathways coupled with the acquisition of  embryonic pathways[37,40].

183

184  **CRC-associated fibroblasts in the tumor microenvironment exhibit diverse phenotypes, and specific**

185  **subtypes are associated with poor prognosis.**

186  We next focused on CRC TME subpopulations. 819 high-quality fibroblasts were re-clustered into eight

187  clusters, and then phenotypically classified into two major subtypes to assess for further CAF heterogeneity.

188  These phenotypic subtypes were found to be immunomodulatory CAF-S1 (*PDGFRA+* and *PDPN+*) and

189  contractile CAF-S4 (*RGS5+* and *MCAM+*) (**Fig. 3a-b, Supplementary Table 4**)[41,42]. This fibroblast cluster

190  dichotomy was also observed in the independent CRC patient scRNA-seq dataset of 31,383 cells

191  (**Supplementary Fig. 5**)[25].

192

193  The CAF-S1 and CAF-S4 subtypes showed striking resemblances to the mCAF (extracellular matrix) and

194  vCAF (vascular) fibroblast subtypes, respectively, as previously described in a mouse breast cancer

195  model[43]. Most clusters were found in multiple patients, albeit in varying proportions, signifying shared

196  patterns in CAF transcriptomic programs between patients. Fibroblasts derived from MSI-H tumors were

197  distributed similarly throughout these clusters (**Fig. 1d**).

198

199  CAF-S1 cells exhibited high chemokine expressions such as *CXCL1, CXCL2, CXCL12, CXCL14,* and

200  immunomodulatory molecules including *TNFRSF12A* (**Supplementary Table S4**). Additionally, CAF-S1

201  cells expressed extracellular matrix genes including matrix-modifying enzymes (*LOXL1* and *LOX*)[43]. To

202  determine this population's functional significance, we compared the CAF-S1 population transcriptomes to

203  those described recently in breast cancer, lung cancer, and head and neck cancer[44]. We recovered five CAF-

204  S1 subtypes that included ecm-myCAF (extracellular; *GJB2*), IL-iCAF (growth factor, *TN*F and interleukin

205    pathway; *SCARA5*), detox-iCAF (detoxification and inflammation; *ADH1B*), wound-myCAF (collagen

206    fibrils and wound healing; *SEMA3C*), and TGFβ-myCAF (*TGF-β* signaling and matrisome; *CST1, TGFb1*),

207    which were previously divided into two major subtypes: iCAF and myCAF (**Fig. 3c**). Among these five

208    subtypes, ecm-myCAF and TGFβ-myCAF are known to correlate with immunosuppressive environments

209    and are enriched in tumors with high regulatory T lymphocytes (Tregs) and depleted CD8+ lymphocytes.

210    Additionally, these subtypes are associated with primary immunotherapy resistance in melanoma and lung

211    cancer[44].

212

213    The CAF-S4 population expressed pericyte markers (*RGS5+, CSPG4+,* and *PDGFRA+*), *CD248*

214    *(endosialin)*, and *EPAS1 (HIF2-α),* that this particular CAF subtype is vessel-associated, with hypoxia

215    potentially contributing to invasion and metastasis, as has been shown in another study **(Fig. 3a-b,**

216    **Supplementary Table 4)**[43]. CAF-S4 clustered into the immature phenotype (*RGS5+, PDGFRB+,* and

217    *CD36+)* and the differentiated myogenic subtype (*TAGLN+* and *MYH11+*) **(Supplementary Table 4)**[42].

218

219    Given the correlation between CMS4 and fibroblast infiltration, we next sought to test the existence of

220    CAF-S1 and CAF-S4 signatures in bulk transcriptomic data and their association with clinical outcomes[12].

221    To this end, we interrogated and carried out a meta-analysis of eight colorectal cancer transcriptomic

222    datasets comprising 1,584 samples and confirmed the presence of CAF-S1 and CAF-S4 gene signatures in

223    CRC and other cancer types (**Fig. 4a**). We detected a strong and positive correlation between specific genes

224    differentially expressed between each CAF subtype in CRC **(Fig. 4a-d)**[43]. We also confirmed the presence

225    of CAF-S1 and CAF-S4 signatures in pancreatic adenocarcinoma (n=118) and non-small cell lung cancer

226    (NSCLC, n=80) cohorts (**Fig. 4a**) (see methods for datasets). Gene signatures were specific to each CAF-

227    S1 and CAF-S4 in bulk transcriptomic datasets thus confirming their existence in TME of CRC and other

228    tumor types.

229

230    We found high CAF-S1 and CAF-S4 signatures associated with significantly poor median overall-survival,

231    irrespective of CMS in three independent CRC datasets (**Fig. 4b-d, Supplementary Fig. 6**). Additionally,

232    CAF signatures stratified the CMS4 subtype into high- and low-risk overall survival in all datasets, thus

233    identifying additional heterogeneity and providing prognostication in this aggressive patient subgroup (**Fig.**

234    **4b-d**).  Here, using scRNA-seq, we show for the first time that high CAF infiltration in CRC is associated

235    with poor prognosis across all molecular subtypes, and which further stratifies the CMS4 subgroup into

236    high and low-risk clinical phenotypes in CRC cohorts.

237

238    **Single-cell RNA sequencing reveals heterogeneity beyond Consensus Molecular Subtypes in**

239    **colorectal cancers and offers therapeutic opportunities.**

240    The lack of association between tumor epithelia and CMS classification, as well as the survival differences

241    between high- and low-risk CAF signatures across CRC molecular subtypes suggest CRCs  are less well

242    defined than the traditional classification systems have indicated (e.g. those systems defined by somatic

243    alterations, epigenomic features, and bulk gene expression data)[10–12,45,46].

244

245    To test our hypothesis, we estimated every cell type fraction using single-cell data in a discovery cohort

246    (GSE39852[11], n=585) with a machine-learning algorithm, CIBERSORTx[47]. Cell type fractions were then

247    validated in two independent cohorts of 177 (GSE17536[48]) and 290 (GSE14333[49]) tumors[9]. When we

248    compared epithelial, immune, and stromal cell populations among the CMS subtypes, we did not detect a

249    distinct pattern of tumor, immune, or stromal cell abundance across the four subtypes (**Fig. 5a-b,**

250    **Supplementary Figs. 7a-b, 8a-b**). Although, we noted some trends in cell pattern enrichment amongst a

251    few cell types consistent with CMS classification (CAFs were enriched in CMS4; dendritic cells, monocytes

252    and TAMs were enriched in CMS1/CMS4; epithelial cells were enriched in CMS3) overall cellular

253    phenotypes were present in varying proportions without a clear distinction between the four subtypes. The

254    discovery and  validation cohorts showed significant discordance in terms of  cell phenotype enrichment

255    with respect to each CMS subtype (**Fig. 5a-b, Supplementary Figs. 7a-b, 8a-b**)[9]. These discordant results

256    could potentially be due to  intra-tumoral variations in tumor purity,  location of tumor biopsy, stromal and

257    immune cell infiltration, and/or CMS's inability to address tumor-TME to tumor-TME and TME to TME

258    variabilities, among other factors[9,50].

259

260    Based on the above findings, we hypothesized that CRC tumors may be more accurately represented as a

261    continuum as has been proposed by Ma et al.[51]. The authors analyzed bulk transcriptomic data using a novel

262    computational framework in which *denovo*, unsupervised clustering methods (k-medoid, non-negative-

263    matrix factorization, and consensus clustering) best classified CRC tumors  in a transcriptomic continuum[52–

264    54]. They further carried out principal component analysis and robustly validated two principal components,

265    PC Cluster Subtype Scores 1 and 2 (PCSS1 and PCSS2, respectively). Using this framework, we reasoned

266    that single-cell data could elucidate the biological underpinnings of a CRC continuum model, and resolve

267    stromal confounding seen using bulk transcriptomes[50,55].

268

269    We evaluated every cell fraction (epithelial, stromal, and immune components) using the Ma et al.

270    algorithm on our discovery and validation cohorts, focusing on the validated PCSS1 and PCSS2. Upon

271    projecting bulk transcriptomes onto the four CMS quadrants, we analyzed single cell-specific gene

272    signatures after deconvolution. Intriguingly, we noted that each cell type not only projected in the

273    expected quadrant but some of these same cell types also projected on other CMS quadrants **(Fig. 6a-b,**

274    **supplementary 9a-b, 10a-b, Supplementary table 5-6)**. For example, epithelial cells, in addition to

275    projecting on CMS2/CMS3 quadrants, also projected on other CMS quadrants and exhibited continuum.

276    These data reflect significant intra-tumoral heterogeneity among all the cellular components that makeup

277    the tumor ecosystem. Thus, it appears that CRC exists in a transcriptomic continuum not only with

278    respect to the tumor cells themselves but also all the other cell types that make up the TME. These aspects

279    would not have been apparent based on bulk transcriptomics analysis alone.

280

281    We found that transcriptional shifts were reproducible across discovery and validation datasets for major

282    cell types (**Fig. 6a-b, supplementary 9a-b, 10a-b**). Our analysis  show no reliability in classifying CRC

283    into immune–stromal rich (CMS1/CMS4) or immune-stromal desert (CMS2/CMS3) subtypes as proposed

284    previously[16]. Thus, confirming continuous scores rather than discrete subtypes may improve classifying

285    CRC tumors and may explain tumor-to-tumor variability, tumor/TME-to-tumor/TME and TME-to-TME

286    variabilities within CMS subgroups as has been previously described [51]. While CAF-S1 and CAF-S4 were

287    present across all CMS groups these CAFs exhibited high PCSS1 and PCSS2 scores correlating with

288    increased enrichment in the CMS4 group. Thus, our analysis identified CAF-S1 and CAF-S4 as the cells

289    of origin for biological heterogeneity in CMS4 subtypes associated with poor prognosis (**Supplementary**

290    **table 5-6**).

291

292    **DISCUSSION**

293    In the present study, we evaluated the CMS classification system of CRC that has been established using

294    bulk RNA-seq through the lens of single-cell transcriptomics. We identified significant intratumoral

295    transcriptomic heterogeneity of cell states within tumor epithelial cells. These findings are consistent with

296    a prior DNA-based multi-omics approach by Sottoriva et al. which proposed a 'big bang model' of

297    CRC[56]. This model suggests that at initiation, CRC is composed of a mixture of subclones, underscoring

298    the significant heterogeneity inherent in CRC biology and potentially reflecting non-linear pathways

299    during CRC evolution and progression. Furthermore, both our findings (using scRNA-seq and

300    pseudotime trajectory analysis) and Sottoriva et al. agree with clinical observations that tumors with

301    diverse clonal populations respond unpredictably to monolithic treatment strategies "targeted" to average

302    expression profiles based on bulk sequencing. Thus, targeting sub-clonal transcriptomic patterns is likely

303    to be more effective when designing personalized CRC therapies [56] [57]. Interestingly, a recent study in

304    gastroesophageal adenocarcinoma (GEA) that employed therapies to target dominant sub-clonal

305    populations in metastatic GEA patients demonstrated improved survival as compared to standard

306    monolithic treatment strategies[58].

307

308    We find that stromal cells are important contributors to the observed biological heterogeneity of CRC, in

309    agreement with emerging data that suggest stromal-derived signatures play key role in modulating CRC

310    biology[7–9,37,59,60]. These recent studies employing bulk transcriptomics demonstrated that the degree of

311    stromal infiltration is associated with prognosis, while a small scRNAseq study utilizing 26 fibroblasts

312    showed poor survival particularly among the CMS4 subtype CAF-enriched CRC tumors[61]. By using a much

313    larger sample set of 1,182 high-quality fibroblasts, our current work identifies the CAF-S1 and CAF-S4

314    subtypes to be the cells of origin associated with poor prognosis across all CRC patients and not just those

315    with CMS4 classified tumors[61]. CAF signatures are able to stratify patients by median overall survival (**Fig.**

316    **4b-d**). These findings are significant since the CMS4 subtype, is primarily stromal-driven and is enriched

317    in more than 40% of metastatic CRC samples from patients with worse outcomes[17,18]. Thus, targeting CAFs

318    to remodel the tumor microenvironment may lead to improved and much-needed therapeutic development

319    for metastatic CRC patients[17,18].

320

321    Targeting of CAFs in solid tumors is being explored in multiple clinical trials with variable results[62,63]. Such

322    studies likely failed to address CAF heterogeneity and their complex interactions with the other cells of

323    TME. Our study suggests that CRC may be intricately entwined with the stroma, and therefore may be

324    amenable to stromal targeted combinatoric approaches, including monoclonal antibodies that abrogate

325    CAF-S1 function. In future studies, the treatment of CRC patients should involve stroma targeted therapies

326    and take the above aspects into consideration[63,64]. The scRNA-seq or bulk-RNA-seq signatures

327    corresponding to CAF-S1 and CAF-S4 may serve as suitable biomarkers for tumors that are reliant on this

328    axis.

329

330    Immunotherapy responses in MSS CRC, which comprise almost 95% of metastatic CRC, are lacking[65];

331    CAF subpopulations within the TME may be suppressing immune responses in these tumors. Based on our

332    analysis we speculate that targeting ecm-myCAF and TGFβ-myCAF subtypes (responsible for

333   immunotherapy resistance in NSCLC and melanoma) via bispecific antibodies, vaccines, or even cell-based

334   therapies, may enhance current checkpoint blockade strategies[44,63,64]. Functional validation and clinical

335   studies will be required to confirm the clinical utility of targeting these CAF populations in CRC.

336

337   More importantly, our study's single-cell resolution enables us to investigate whether tumor cell

338   transcriptomes, and by extension, biological phenotypes, are the primary determinant of CMS

339   classification. Based on our findings, it appears that bulk analysis may have been confounded by varying

340   degrees of tumor microenvironment population enrichment, and that tumor cells within each patient do not

341   segregate into static phenotypes but rather exhibit considerable plasticity. Our single-cell analysis

342   uncovered complex and mixed cellular phenotypes among each cell-specific subpopulation, which

343   projected in a transcriptomic continuum across the CMS groups .These findings were further supported by

344   our scRNA-seq CMS classification analysis that assigned each CRC sample to multiple CMS groups,

345   thereby suggesting CMS heterogeneity within each CRC tumor (**Fig. 1d, Supplementary Table 1**)[16,66].

346   These findings may also explain why CMS-defined populations of tumors have not been readily observed

347   in transcriptomic data from independent CRC cohorts[9,17].  Our results suggest that attempts to divide CRC

348   phenotypes into the current discrete subtypes may undermine optimal patient stratification in the clinical

349   trial setting[9]. The only prospective study to date that utilized the CMS classification (specifically the CMS4

350   subtype) for patient selection based on dual PD-L1/TGF-ß expression signatures was halted due to futility,

351   suggesting CMS classification may adversely impact clinical trial design[67].

352

353   Previous studies relying on bulk transcriptomics had concerns that stromal content may conceal subtle

354   critical gene signals originating from other key cellular phenotypes within the CRC spectrum, and thus

355   affect CRC classification[8,9,50]. Our extensive single-cell analysis allowed for detailed evaluation of all

356   cellular subtypes of CRC simultaneously, thereby uncovering contributions of the different components

357   making up the complex cellular milieu within the CMS classification schema. The present analysis resolved

358    the issue of stromal confounding in CRC and showed that all cellular components contribute to a

359    transcriptomic continuum encompassing all the subtypes that together define the CMS system.

360

361    In conclusion, our analysis not only shows tumor-to-tumor variability (as proposed by Ma et al. and other

362    groups), but also demonstrates tumor-TME to tumor-TME as well as TME to TME transcriptional

363    variability at the single-cell resolution level[9,51]. These data contribute to the conceptual advances in

364    understanding CRC pathogenesis, clinical management, and therapeutic development. We suggest

365    revisiting CMS-like classifications that define CRCs into distinct static subtypes while, in reality CRC

366    tumors exist in a continuum. We caution in using CMS to stratify patients for drug development and

367    clinical trial design. Future studies should concentrate on developing biomarkers (such as CAF's) and

368    therapeutic agents that can stratify CRC patients beyond traditional classification system boundaries.

369    Ultimately, newer single-cell multiomic technologies will allow us to detect somatic mutations,

370    transcriptomes, proteomes, epigenomes and metabolomes in real-time at the single-cell resolution to

371    better guide more individualized and improved therapies.

372

382

383    **AUTHOR CONTRIBUTIONS**

384    A.M. devised, supervised the study, conducted data analysis, and wrote the manuscript. A.M.K. performed

385    data analyses, wrote the manuscript, and created figures. Z.K. performed data analyses and wrote the

386    manuscript. M.W.G. aided in analysis, wrote the manuscript, and generated figures. A.S. supervised study

387    and wrote manuscript. C.E. and S.S.T. helped with bulk transcriptomic analysis. D.M.H, H.R.G., A.R.B

388    helped with sample collection. All other authors contributed substantially to data interpretation, and

389    manuscript editing. All authors read and approved this manuscript.

390

391    **CODE AVAILABILITY**

392    The code generated and utilized in the completion of this publication will be available in a Github

393    repository specific to this project.

394

395    **DATA AVAILABILITY**

396    Sequencing data deposition is currently in progress. Ten bulk transcriptomic datasets were accessed from

397    the Gene Expression Omnibus (GEO) database (https://www.ncbi.nlm.nih.gov/geo/).

398

399    **COMPETING INTERESTS**

400    A.M. and J.A.B. received research funding from Tempus lab.

401    A.S. receives research funding from Bristol-Myers Squibb; Merck KGaA, Pierre Fabre. Further, A.S. holds

402    patent PCT/IB2013/060416, '*Colorectal cancer classification with differential prognosis and personalized*

403    *therapeutic responses*' and patent number 2011213.2 '*Prognostic and Treatment Response Predictive*

404    *Method.*'

## REFERENCES

1. Cancer of the Colon and Rectum - Cancer Stat Facts. SEER. Accessed December 29, 2020. https://seer.cancer.gov/statfacts/html/colorect.html

2. Arnold M, Sierra MS, Laversanne M, Soerjomataram I, Jemal A, Bray F. Global patterns and trends in colorectal cancer incidence and mortality. *Gut*. 2017;66(4):683-691. doi:10.1136/gutjnl-2015-310912

3. Kunst N, Alarid-Escudero F, Aas E, Coupé VMH, Schrag D, Kuntz KM. Estimating population-based recurrence rates of colorectal cancer over time in the United States. *Cancer Epidemiol Biomarkers Prev*. 2020;29(12):2710-2718. doi:10.1158/1055-9965.EPI-20-0490

4. Xie Y-H, Chen Y-X, Fang J-Y. Comprehensive review of targeted therapy for colorectal cancer. *Signal Transduction and Targeted Therapy*. 2020;5(1):1-30. doi:10.1038/s41392-020-0116-z

5. Pagès F, Mlecnik B, Marliot F, et al. International validation of the consensus Immunoscore for the classification of colon cancer: a prognostic and accuracy study. *The Lancet*. 2018;391(10135):2128-2139. doi:10.1016/S0140-6736(18)30789-X

6. The Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon and rectal cancer. *Nature*. 2012;487(7407):330-337. doi:10.1038/nature11252

7. Calon A, Lonardo E, Berenguer-Llergo A, et al. Stromal gene expression defines poor-prognosis subtypes in colorectal cancer. *Nature Genetics*. 2015;47(4):320-329. doi:10.1038/ng.3225

8. Isella C, Terrasi A, Bellomo SE, et al. Stromal contribution to the colorectal cancer transcriptome. *Nat Genet*. 2015;47(4):312-319. doi:10.1038/ng.3224

9. Dunne PD, McArt DG, Bradley CA, et al. Challenging the cancer molecular stratification dogma: intratumoral heterogeneity undermines consensus molecular subtypes and potential piagnostic value in colorectal cancer. *Clin Cancer Res*. 2016;22(16):4095-4104. doi:10.1158/1078-0432.CCR-16-0032

10. Sadanandam A, Lyssiotis CA, Homicsko K, et al. A colorectal cancer classification system that associates cellular phenotype and responses to therapy. *Nature Medicine*. 2013;19(5):619-625. doi:10.1038/nm.3175

11. Marisa L, Reyniès A de, Duval A, et al. Gene Expression Classification of Colon Cancer into Molecular Subtypes: Characterization, Validation, and Prognostic Value. *PLOS Medicine*. 2013;10(5):e1001453. doi:10.1371/journal.pmed.1001453

12. Guinney J, Dienstmann R, Wang X, et al. The consensus molecular subtypes of colorectal cancer. *Nature Medicine*. 2015;21(11):1350-1356. doi:10.1038/nm.3967

13. Lenz H-J, Ou F-S, Venook AP, et al. Impact of Consensus Molecular Subtype on Survival in Patients With Metastatic Colorectal Cancer: Results From CALGB/SWOG 80405 (Alliance). *J Clin Oncol*. 2019;37(22):1876-1885. doi:10.1200/JCO.18.02258

14. Stintzing S, Wirapati P, Lenz H-J, et al. Consensus molecular subgroups (CMS) of colorectal cancer (CRC) and first-line efficacy of FOLFIRI plus cetuximab or bevacizumab in the FIRE3 (AIO KRK-0306) trial. *Ann Oncol*. 2019;30(11):1796-1803. doi:10.1093/annonc/mdz387

15. Mooi JK, Wirapati P, Asher R, et al. The prognostic impact of consensus molecular subtypes (CMS) and its predictive effects for bevacizumab benefit in metastatic colorectal cancer: molecular analysis of the AGITG MAX clinical trial. *Ann Oncol*. 2018;29(11):2240-2246. doi:10.1093/annonc/mdy410

16. Sveen A, Cremolini C, Dienstmann R. Predictive modeling in colorectal cancer: time to move beyond consensus molecular subtypes. *Annals of Oncology*. 2019;30(11):1682-1685. doi:10.1093/annonc/mdz412

17. Fontana E, Eason K, Cervantes A, Salazar R, Sadanandam A. Context matters—consensus molecular subtypes of colorectal cancer as biomarkers for clinical trials. *Annals of Oncology*. 2019;30(4):520-527. doi:10.1093/annonc/mdz052

18. Khambata-Ford S, Garrett CR, Meropol NJ, et al. Expression of epiregulin and amphiregulin and K-ras mutation status predict disease control in metastatic colorectal cancer patients treated with cetuximab. *J Clin Oncol*. 2007;25(22):3230-3237. doi:10.1200/JCO.2006.10.5437

19. Zappia L, Oshlack A. Clustering trees: a visualization for evaluating clusterings at multiple resolutions. *Gigascience*. 2018;7(7). doi:10.1093/gigascience/giy083

20. Zhang L, Li Z, Skrzypczynska KM, et al. Single-cell analyses inform mechanisms of myeloid-targeted therapies in colon cancer. *Cell*. 2020;181(2):442-459.e29. doi:10.1016/j.cell.2020.03.048

21. Zhang L, Yu X, Zheng L, et al. Lineage tracking reveals dynamic relationships of T cells in colorectal cancer. *Nature*. 2018;564(7735):268-272. doi:10.1038/s41586-018-0694-x

22. Corridoni D, Antanaviciute A, Gupta T, et al. Single-cell atlas of colonic CD8 + T cells in ulcerative colitis. *Nature Medicine*. 2020;26(9):1480-1490. doi:10.1038/s41591-020-1003-4

23. Aran D, Looney AP, Liu L, et al. Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nature Immunology*. 2019;20(2):163-172. doi:10.1038/s41590-018-0276-y

24. Luoma AM, Suo S, Williams HL, et al. Molecular Pathways of Colon Inflammation Induced by Cancer Immunotherapy. *Cell*. 2020;182(3):655-671.e22. doi:10.1016/j.cell.2020.06.001

25. Lee H-O, Hong Y, Etlioglu HE, et al. Lineage-dependent gene expression programs influence the immune landscape of colorectal cancer. *Nature Genetics*. 2020;52(6):594-603. doi:10.1038/s41588-020-0636-z

26. Mills JC, Sansom OJ. Reserve stem cells: Differentiated cells reprogram to fuel repair, metaplasia, and neoplasia in the adult gastrointestinal tract. *Sci Signal*. 2015;8(385):re8. doi:10.1126/scisignal.aaa7540

27. Lambrechts D, Wauters E, Boeckx B, et al. Phenotype molding of stromal cells in the lung tumor microenvironment. *Nat Med*. 2018;24(8):1277-1289. doi:10.1038/s41591-018-0096-5

28. Izar B, Tirosh I, Stover EH, et al. A single-cell landscape of high-grade serous ovarian cancer. *Nature Medicine*. 2020;26(8):1271-1279. doi:10.1038/s41591-020-0926-0

29. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst*. 2015;1(6):417-425. doi:10.1016/j.cels.2015.12.004

30. Malikic S, Jahn K, Kuipers J, Sahinalp SC, Beerenwinkel N. Integrative inference of subclonal tumour evolution from single-cell and bulk sequencing data. *Nat Commun*. 2019;10(1):2750. doi:10.1038/s41467-019-10737-5

31. Lim SB, Yeo T, Lee WD, et al. Addressing cellular heterogeneity in tumor and circulation for refined prognostication. *Proc Natl Acad Sci U S A*. 2019;116(36):17957-17962. doi:10.1073/pnas.1907904116

32. Wang R, Dang M, Harada K, et al. Single-cell dissection of intratumoral heterogeneity and lineage diversity in metastatic gastric adenocarcinoma. *Nat Med*. 2021;27(1):141-151. doi:10.1038/s41591-020-1125-8

33. Hänzelmann S, Castelo R, Guinney J. GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics*. 2013;14:7. doi:10.1186/1471-2105-14-7

34. Sathe A, Grimes SM, Lau BT, et al. Single-cell genomic characterization reveals the cellular reprogramming of the gastric tumor microenvironment. *Clin Cancer Res*. 2020;26(11):2640-2653. doi:10.1158/1078-0432.CCR-19-3231

35. Trapnell C, Cacchiarelli D, Grimsby J, et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol*. 2014;32(4):381-386. doi:10.1038/nbt.2859

36. Qiu X, Mao Q, Tang Y, et al. Reversed graph embedding resolves complex single-cell trajectories. *Nat Methods*. 2017;14(10):979-982. doi:10.1038/nmeth.4402

37. Dalerba P, Kalisky T, Sahoo D, et al. Single-cell dissection of transcriptional heterogeneity in human colon tumors. *Nat Biotechnol*. 2011;29(12):1120-1127. doi:10.1038/nbt.2038

38. Vermeulen L, Todaro M, Mello F de S, et al. Single-cell cloning of colon cancer stem cells reveals a multi-lineage differentiation capacity. *PNAS*. 2008;105(36):13427-13432. doi:10.1073/pnas.0805706105

39. Kim N, Kim HK, Lee K, et al. Single-cell RNA sequencing demonstrates the molecular and cellular reprogramming of metastatic lung adenocarcinoma. *Nat Commun*. 2020;11(1):2285. doi:10.1038/s41467-020-16164-1

40. Monk M, Holding C. Human embryonic genes re-expressed in cancer cells. *Oncogene*. 2001;20(56):8085-8091. doi:10.1038/sj.onc.1205088

41. Costa A, Kieffer Y, Scholer-Dahirel A, et al. Fibroblast heterogeneity and immunosuppressive environment in human breast cancer. *Cancer Cell*. 2018;33(3):463-479.e10. doi:10.1016/j.ccell.2018.01.011

42. Wu SZ, Roden DL, Wang C, et al. Stromal cell diversity associated with immune evasion in human triple-negative breast cancer. *The EMBO Journal*. 2020;39(19). doi:10.15252/embj.2019104063

43. Bartoschek M, Oskolkov N, Bocci M, et al. Spatially and functionally distinct subclasses of breast cancer-associated fibroblasts revealed by single cell RNA sequencing. *Nat Commun*. 2018;9(1):5150. doi:10.1038/s41467-018-07582-3

44. Kieffer Y, Hocine HR, Gentric G, et al. Single-cell analysis reveals fibroblast clusters linked to immunotherapy resistance in cancer. *Cancer Discov*. 2020;10(9):1330-1351. doi:10.1158/2159-8290.CD-19-1384

45. Sadanandam A, Wang X, de Sousa E Melo F, et al. Reconciliation of classification systems defining molecular subtypes of colorectal cancer. *Cell Cycle*. 2014;13(3):353-357. doi:10.4161/cc.27769

46. De Sousa E Melo F, Wang X, Jansen M, et al. Poor-prognosis colon cancer is defined by a molecularly distinct subtype and develops from serrated precursor lesions. *Nat Med*. 2013;19(5):614-618. doi:10.1038/nm.3174

47. Newman AM, Steen CB, Liu CL, et al. Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nature Biotechnology*. 2019;37(7):773-782. doi:10.1038/s41587-019-0114-2

48. Jorissen RN, Gibbs P, Christie M, et al. Metastasis-associated gene expression changes predict poor outcomes in patients with Dukes stage B and C colorectal cancer. *Clin Cancer Res*. 2009;15(24):7642-7651. doi:10.1158/1078-0432.CCR-09-1431

49. Smith JJ, Deane NG, Wu F, et al. Experimentally derived metastasis gene expression profile predicts recurrence and death in patients with colon cancer. *Gastroenterology*. 2010;138(3):958-968. doi:10.1053/j.gastro.2009.11.005

50. Dunne PD, Alderdice M, O'Reilly PG, et al. Cancer-cell intrinsic gene expression signatures overcome intratumoural heterogeneity bias in colorectal cancer patient classification. *Nature Communications*. 2017;8(1):15657. doi:10.1038/ncomms15657

51. Ma S, Ogino S, Parsana P, et al. Continuity of transcriptomes among colorectal cancer subtypes based on meta-analysis. *Genome Biology*. 2018;19(1):142. doi:10.1186/s13059-018-1511-4

52. Clustering rules: a comparison of partitioning and hierarchical clustering algorithms. Accessed January 22, 2021. https://link.springer.com/article/10.1007/s10852-005-9022-1

53. Wilkerson MD, Hayes DN. ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics*. 2010;26(12):1572-1573. doi:10.1093/bioinformatics/btq170

54. Gaujoux R, Seoighe C. A flexible R package for nonnegative matrix factorization. *BMC Bioinformatics*. 2010;11(1):367. doi:10.1186/1471-2105-11-367

55. Isella C, Brundu F, Bellomo SE, et al. Selective analysis of cancer-cell intrinsic transcriptional traits defines novel clinically relevant subtypes of colorectal cancer. *Nature Communications*. 2017;8(1):15107. doi:10.1038/ncomms15107

554  56.  Sottoriva A, Kang H, Ma Z, et al. A Big Bang model of human colorectal tumor growth. *Nature*
555        *Genetics*. 2015;47(3):209-216. doi:10.1038/ng.3214
556  57.  Loeb LA, Kohrn BF, Loubet-Senear KJ, et al. Extensive subclonal mutational diversity in human
557        colorectal cancer and its significance. *PNAS*. 2019;116(52):26863-26872.
558        doi:10.1073/pnas.1910301116
559  58.  Catenacci DVT, Moya S, Lomnicki S, et al. Personalized Antibodies for Gastroesophageal
560        Adenocarcinoma (PANGEA): A Phase II Study Evaluating an Individualized Treatment Strategy
561        for Metastatic Disease. *Cancer Discov*. 2021;11(2):308-325. doi:10.1158/2159-8290.CD-20-1408
562  59.  Tauriello DVF, Palomo-Ponce S, Stork D, et al. TGFβ drives immune evasion in genetically
563        reconstituted colon cancer metastasis. *Nature*. 2018;554(7693):538-543. doi:10.1038/nature25492
564  60.  Calon A, Espinet E, Palomo-Ponce S, et al. Dependency of colorectal cancer on a TGF-β-driven
565        program in stromal cells for metastasis initiation. *Cancer Cell*. 2012;22(5):571-584.
566        doi:10.1016/j.ccr.2012.08.013
567  61.  Li H, Courtois ET, Sengupta D, et al. Reference component analysis of single-cell transcriptomes
568        elucidates cellular heterogeneity in human colorectal tumors. *Nat Genet*. 2017;49(5):708-718.
569        doi:10.1038/ng.3818
570  62.  Gascard P, Tlsty TD. Carcinoma-associated fibroblasts: orchestrating the composition of
571        malignancy. *Genes Dev*. 2016;30(9):1002-1019. doi:10.1101/gad.279737.116
572  63.  Sahai E, Astsaturov I, Cukierman E, et al. A framework for advancing our understanding of
573        cancer-associated fibroblasts. *Nat Rev Cancer*. 2020;20(3):174-186. doi:10.1038/s41568-019-
574        0238-1
575  64.  Liu T, Han C, Wang S, et al. Cancer-associated fibroblasts: an emerging target of anti-cancer
576        immunotherapy. *J Hematol Oncol*. 2019;12(1):86. doi:10.1186/s13045-019-0770-1
577  65.  André T, Shiu K-K, Kim TW, et al. Pembrolizumab in Microsatellite-Instability-High Advanced
578        Colorectal Cancer. *N Engl J Med*. 2020;383(23):2207-2218. doi:10.1056/NEJMoa2017699
579  66.  Laurent-Puig P, Marisa L, Ayadi M, et al. Colon cancer molecular subtype intratumoral
580        heterogeneity and its prognostic impact: An extensive molecular analysis of the PETACC-8.
581        *Annals of Oncology*. 2018;29:viii18. doi:10.1093/annonc/mdy269.058
582  67.  Mehrvarz Sarshekeh A, Lam M, Zorrilla IR, et al. Consensus molecular subtype (CMS) as a novel
583        integral biomarker in colorectal cancer: A phase II trial of bintrafusp alfa in CMS4 metastatic
584        CRC. *JCO*. 2020;38(15_suppl):4084-4084. doi:10.1200/JCO.2020.38.15_suppl.4084
585  68.  Stuart T, Butler A, Hoffman P, et al. Comprehensive integration of single-cell data. *Cell*.
586        2019;177(7):1888-1902.e21. doi:10.1016/j.cell.2019.05.031
587  69.  McGinnis CS, Murrow LM, Gartner ZJ. DoubletFinder: Doublet detection in single-cell RNA
588        sequencing data using artificial nearest neighbors. *Cell Syst*. 2019;8(4):329-337.e4.
589        doi:10.1016/j.cels.2019.03.003
590  70.  Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*.
591        2010;11(10):R106. doi:10.1186/gb-2010-11-10-r106
592  71.  Thompson B. Canonical correlation analysis. In: *Encyclopedia of Statistics in Behavioral Science*.
593        American Cancer Society; 2005. doi:10.1002/0470013192.bsa068
594  72.  Jolliffe I. Principal component analysis. In: Lovric M, ed. *International Encyclopedia of Statistical*
595        *Science*. Springer; 2011:1094-1096. doi:10.1007/978-3-642-04898-2_455
596  73.  Maaten L van der. Accelerating t-SNE using tree-based algorithms. *Journal of Machine Learning*
597        *Research*. 2014;15(93):3221-3245.
598  74.  McInnes L, Healy J, Saul N, Großberger L. UMAP: uniform manifold approximation and
599        projection. *Journal of Open Source Software*. 2018;3(29):861. doi:10.21105/joss.00861
600  75.  Zilionis R, Engblom C, Pfirschke C, et al. Single-cell transcriptomics of human and mouse lung
601        cancers reveals conserved myeloid populations across individuals and species. *Immunity*.
602        2019;50(5):1317-1334.e10. doi:10.1016/j.immuni.2019.03.009
603  76.  Helmink BA, Reddy SM, Gao J, et al. B cells and tertiary lymphoid structures promote
604        immunotherapy response. *Nature*. 2020;577(7791):549-555. doi:10.1038/s41586-019-1922-8

605 77.  Guo X, Zhang Y, Zheng L, et al. Global characterization of T cells in non-small-cell lung cancer
606     by single-cell sequencing. *Nat Med*. 2018;24(7):978-985. doi:10.1038/s41591-018-0045-3

607 78.  Szabo PA, Levitin HM, Miron M, et al. Single-cell transcriptomics of human T cells reveals tissue
608     and activation signatures in health and disease. *Nat Commun*. 2019;10(1):4706.
609     doi:10.1038/s41467-019-12464-3

610 79.  Nirschl CJ, Suárez-Fariñas M, Izar B, et al. IFNγ-dependent tissue-Immune homeostasis Is co-
611     opted in the tumor microenvironment. *Cell*. 2017;170(1):127-141.e15.
612     doi:10.1016/j.cell.2017.06.016

613 80.  Shi Z, Zhang Q, Yan H, et al. More than one antibody of individual B cells revealed by single-cell
614     immune profiling. *Cell Discov*. 2019;5:64. doi:10.1038/s41421-019-0137-3

615 81.  Ramesh A, Schubert RD, Greenfield AL, et al. A pathogenic and clonally expanded B cell
616     transcriptome in active multiple sclerosis. *Proc Natl Acad Sci U S A*. 2020;117(37):22932-22943.
617     doi:10.1073/pnas.2008523117

618 82.  Puram SV, Tirosh I, Parikh AS, et al. Single-cell transcriptomic analysis of primary and metastatic
619     tumor ecosystems in head and neck cancer. *Cell*. 2017;171(7):1611-1624.e24.
620     doi:10.1016/j.cell.2017.10.044

621 83.  Qiu X, Hill A, Packer J, Lin D, Ma Y-A, Trapnell C. Single-cell mRNA quantification and
622     differential analysis with Census. *Nat Methods*. 2017;14(3):309-315. doi:10.1038/nmeth.4150

623 84.  Beaubier N, Tell R, Lau D, et al. Clinical validation of the tempus xT next-generation targeted
624     oncology sequencing assay. *Oncotarget*. 2019;10(24):2384-2396. doi:10.18632/oncotarget.26797

625 85.  Andrews S. *FastQC: A Quality Control Tool for High Throughput Sequence Data.* Babraham
626     Institute; 2012. http://www.bioinformatics.babraham.ac.uk/projects/fastqc

627 86.  Jorissen RN, Lipton L, Gibbs P, et al. DNA copy-number alterations underlie gene expression
628     differences between microsatellite stable and unstable colorectal cancers. *Clin Cancer Res*.
629     2008;14(24):8061-8069. doi:10.1158/1078-0432.CCR-08-1431

630 87.  Skrzypczak M, Goryca K, Rubel T, et al. Modeling oncogenic signaling in colon tumors by
631     multidirectional analyses of microarray data directed for maximization of analytical reliability.
632     *PLoS One*. 2010;5(10). doi:10.1371/journal.pone.0013091

633 88.  Kemper K, Versloot M, Cameron K, et al. Mutations in the Ras-Raf Axis underlie the prognostic
634     value of CD133 in colorectal cancer. *Clin Cancer Res*. 2012;18(11):3132-3141. doi:10.1158/1078-
635     0432.CCR-11-3066

636 89.  Schlicker A, Beran G, Chresta CM, et al. Subtypes of primary colorectal tumors correlate with
637     response to targeted treatment in colorectal cell lines. *BMC Med Genomics*. 2012;5:66.
638     doi:10.1186/1755-8794-5-66

639 90.  Janky R, Binda MM, Allemeersch J, et al. Prognostic relevance of molecular subtypes and master
640     regulators in pancreatic ductal adenocarcinoma. *BMC Cancer*. 2016;16:632. doi:10.1186/s12885-
641     016-2540-6

642 91.  Meister M, Belousov A, EC X, et al. Intra-tumor heterogeneity of gene expression profiles in early
643     stage non-small cell lung cancer. *Journal of Bioinformatics Research Studies*. 2014;1.

644 92.  Gautier L, Cope L, Bolstad BM, Irizarry RA. affy—analysis of Affymetrix GeneChip data at the
645     probe level. *Bioinformatics*. 2004;20(3):307-315. doi:10.1093/bioinformatics/btg405

646 93.  Ritchie ME, Phipson B, Wu D, et al. limma powers differential expression analyses for RNA-
647     sequencing and microarray studies. *Nucleic Acids Res*. 2015;43(7):e47. doi:10.1093/nar/gkv007

648 94.  Eide PW, Bruun J, Lothe RA, Sveen A. CMScaller: an R package for consensus molecular
649     subtyping of colorectal cancer pre-clinical models. *Sci Rep*. 2017;7(1):16618. doi:10.1038/s41598-
650     017-16747-x

651 95.  Wickham H. *Ggplot2*. 2nd ed. Springer

652 96.  Therneau T. *A Package for Survival Analysis in R*.; 2020.

653 97.  Zhao X, Valen E, Parker BJ, Sandelin A. Systematic clustering of transcription start site
654     landscapes. *PLOS ONE*. 2011;6(8):e23409. doi:10.1371/journal.pone.0023409

655  98.    Bunis DG, Andrews J, Fragiadakis GK, Burt TD, Sirota M. dittoSeq: universal user-friendly
656        single-cell and bulk RNA sequencing visualization toolkit. *Bioinformatics*. Published online
657        December 12, 2020. doi:10.1093/bioinformatics/btaa1011
658  99.    *Morpheus: Interactive Heat Maps Using "morpheus.Js" and "Htmlwidgets."*; 2021.
659        https://software.broadinstitute.org/morpheus
660  100.   Pedersen TL. *Patchwork: The Composer of Plots. R*.; 2020. https://CRAN.R-
661        project.org/package=patchwork
662  101.   Kassambara A. *Ggpubr: "ggplot2" Based Publication Ready Plots*. Based Publication; 2020.
663        https://CRAN.R-project.org/package=ggpubr

664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699

700    **FIGURE LEGENDS**

701

702    **Figure 1. Identification and clustering of single cells**. **a,** Workflow of sample collection, sorting, and
703    sequencing (methods contain full description for each step). **b,** UMAP characterization of the 49,859 cells
704    profiled. Coloring demonstrates clusters, tumor vs. normal sample origin (condition), and individual
705    sample origin. **c,** Identification of various cell types based on expression of specified marker genes. **d,**
706    Characterization of the proportion of cell types identified in tumor vs. normal colon tissue, sidedness
707    (right vs. left), microsatellite instability (MSI) status, single-cell Consensus Molecular Subtypes (scCMS)
708    classification, Consensus Molecular Subtypes (CMS) of bulk RNA-seq data, and origin of sample. The
709    transcription counts of tumor and normal tissue cell types are demonstrated at the bottom with boxplot
710    representation. The graph represents total clusters and cell types identified after re-clustering of each cell
711    compartment depicting global heterogeneous landscape of colorectal cancers.

712

713    **Figure 2. Reclustering and characterization of the epithelial compartment. a,** UMAP of tumor and
714    non-malignant epithelial reclustering demonstrating 17 distinct clusters. **b,** Bar plot representation of cell
715    proportions by sample, tissue type, MSI status, tumor location, scCMS score, and bulk CMS score. **c,**
716    Heatmap of Hallmark pathway analysis within the epithelial cell compartment. **d,** Trajectory analysis of
717    cells colored by tumor location, scCMS, MSI status, and bulk CMS status.

718

719    **Figure 3. Fibroblast clusters in colon and colorectal tumors. a,** UMAP of 819 fibroblasts colored by
720    distinct clusters, CAF status, tissue status and origin of sample.  UMAP of fibroblasts colored by specific
721    CAF-S1 subtypes. **b,** Heatmap showing the variable expression of fibroblast specific marker genes across
722    CAF-S1, CAF-S4, and normal fibroblasts. **c,** UMAP color-coded for marker genes for five CAF-S1
723    subtypes as indicated.

724

725    **Figure 4. Correlation of CAF-S1 and CAF-S4 gene profiles across human bulk transcriptomic data.**
726    **a,** Pearson's correlation of genes from CAF-S1 and CAF-S4 profiles in colorectal cancer (n= 1584; CAF-
727    S1 and CAF-S4 r > 0.8), pancreatic cancer (n= 118; CAF-S1 r = 0.70; CAF-S4 r= 0.60), non-small cell
728    lung cancer (n = 80; CAF-S1 r  = 0.69; CAF-S4 r = 0.67). **b-d,** Pearson correlation plots, Kaplan-Meyer
729    survival curves, and bar plots of CMS status assessing CAF expression in individual CRC datasets. Plots
730    b-c are generated from single GEO datasets; GSE17536 (n = 177), GSE39582 (n= 585) and GSE33113
731    (n= 96), respectively. Note: High CAF-S1 and CAF-S4 gene signatures are associated with poor survival
732    across all CMS . r = coefficient correlation. Hazard ratio > 1 in figures b-d.

733

734    **Figure 5. Average cell type abundance from CRC dataset GSE39582 and sorted by bulk CMS**
735    **status. a,** Boxplots show the distribution of cell types within tumors with varying CMS status. The
736    whiskers depict the 1.5 x IQR. The p-values for pairwise t-tests comparisons (with Benjamani-Hochberg
737    correction) and ANOVA tests of cell abundance across CMS are shown in the figure. **b,** Deconvolution
738    heatmap of different cell types by average expression using CIBERSORTx demonstrating cell type
739    distribution within each CMS category. ILCs = Innate lymphoid cells, GC Cells = Germinal Center B
740    Cells, NK Cells = Natural Killer Cells, TAMs = Tumor Associated Macrophages.

741

742

743    **Figure 6. Continuous scores for CRC dataset GSE39582. a,** Principal component analysis plot
744    showing PCSS1 and PCSS2 continuous scores reported by CMS classification across 19 cell types show
745    minimal separation in the top 2 principal components. **b,** All cell types projected on four quadrants
746    representing CMS1-4 using PCSS1 and PCSS2 scores. Note that the cell types largely form a continuum
747    along CMS status and are not clustered in discrete quadrants separate from one another. Cells and
748    markers are colored by bulk CMS status accordingly to the tumor sample of origin.  ILCs = Innate
749    lymphoid cells, GC Cells = Germinal Center B Cells, NK Cells = Natural Killer Cells, TAMs = Tumor
750    Associated Macrophages.

751   **METHODS**
752
753   ***Patient and tissue sample collection.*** Patients with resectable untreated CRC who underwent
754   curative colon resection at Rush University Medical Center (Chicago, IL, USA) were included in
755   this Institutional Review Board (IRB)-approved study. CRC specimens from 16 patients including
756   nine Caucasian, six African American and one Asian patient with corresponding 8 adjacent normal
757   tissue samples were processed immediately after collection at Rush University Medical Center
758   Biorepository and sent for scRNA-seq.  Thus, our scRNA-seq atlas represent diverse patient
759   population. The study was conducted in accordance with ethical standards and all patients provided
760   written informed consent.
761
762   ***Droplet based scRNA-seq - 10× library preparation and sequencing.*** Single-cell RNA sequencing
763   (scRNA-seq) was performed using 10X Genomics Single Cell 5' Platform. Tumors and normal
764   colon samples were enzymatically dissociated (*Miltenyi*), filtered through a 70-micron cell strainer,
765   pelleted after centrifugation at 300 x*g* and resuspended in DAPI-FACS buffer (PBS, 0.04% BSA).
766   Samples were sorted and viable singlets were gated on the basis of scatter properties and DAPI
767   exclusion. Approximately 3000 cells were pelleted and resuspended in PBS, and cells underwent
768   single cell droplet-based capture on 10X Chromium instruments according to the 10X Genomics
769   Single Cell 5' Platform protocol. Transcriptome libraries post-fragmentation, end-repair, and A-
770   tailing double-sided size selection, and subsequent adaptor ligation also followed the
771   manufacturer's protocol. Illumina *NextSeq 550* was used for library sequencing and data were
772   mapped and counted using Cellranger-v3.1.0 (*GRCh38/hg38*).
773
774   ***scRNA-seq data quality control, gene-expression quantification, dimensionality reduction, and***
775   ***identification of cell clusters.*** *Cell Ranger* was utilized to process the raw gene expression matrices
776   per samples and all samples from multiple patients were combined in R package (v3.6.3 2020-02-
777   29] -- "*Holding the Windsock*"). Seurat package (v3.2.2) was used in this integrative multimodal
778   analysis[68]. Genes detected in fewer than three cells and cells expressing less than 200 detected
779   genes were filtered out and excluded from analysis. In addition, cells expressing > 25%
780   mitochondria were removed. Cell cycle scoring was performed, (for the S phase and the G2M
781   phase) and the predicted cell cycle phases were calculated. Doublet detection and any higher-order
782   multiplets that were not dissociated during sample preparation were removed  via the
783   *DoubletFinder* (v2.0.2) package using default settings[69]. Following quality control one normal
784   colon sample (B-cac13) was discarded due to poor data quality. Finally, 49,859 cells remained and
785   were utilized for downstream analysis.
786
787   We adopted the general protocol described in Stuart et al. (2019) to group single cells into different
788   cell subsets[68]. We employed the following steps: clustering the cells within each compartment
789   (including the selection of variable genes for each dataset based on a variance stabilizing
790   transformation [VST]), canonical correlation analysis (CCA) to remove batch effects among the
791   samples, reduction of dimensionality, and projection of cells onto graphs [70,71]. Principal component
792   analysis (PCA) was carried out on the scaled data of highly variable genes[72]. The first 30 principal
793   components (PCs) were used to cluster the cells and to perform a subtype analysis by nonlinear
794   dimensionality reduction (t-SNE) and to construct Uniform Manifold Approximation and
795   Projection (UMAP) for cell embeddings[73,74].  We identified cell clusters under the optimal
796   resolution by a shared nearest neighbor (SNN) modularity optimization-based clustering method.
797   We implemented the *FindClusters* function of the Seurat package, which first calculated *k-nearest*
798   *neighbors* and constructed the SNN graph. We implemented the original *Louvain algorithm*
799   (algorithm = 1) for modularity optimization. Additionally, we utilized Clustree (v0.4.3) and manual
800   review for identifying the best clustering resolution[19].
801

802   ***Major cell type detection and data visualization.*** To identify all major cell types, we evaluated
803   differentially expressed markers in each identity cell group by comparing them to other clusters
804   using the Seurat *FindAllMarkers* function. We used positively expressed genes with an average
805   expression of >/= 2-fold higher in that subcluster than the average expression in the rest of the other
806   subclusters. We used known marker genes, which have the highest fold expression in that cluster
807   with respect to the other clusters. We also utilized SingleR ((v0.99.10, R Package), which leverage
808   large transcriptomic datasets of well-annotated cell types and manual annotation  for cell-type
809   identification[27,75–77]. Depending on the presence of known marker genes the clusters were grouped
810   as: epithelial cells (*EPCAM, KRT8,* and *KRT18*), fibroblasts (*COL1A1, DCN, COL1A2,* and *C1R*),
811   endothelial cells (CD31+), myeloid cells (*LYZ, MARCO, CD68,* and *FCGR3A*), CD4 T cells (*CD4*),
812   CD8 T cells (*CD8A* and *CD8B*),  and  B cells (*MZB1*), [27,39,43,75,78–81]. The cells were eventually
813   assembled into DGE matrices within each compartment, containing all six cell types.
814
815   ***Major-cell type subclustering and data visualization.*** Each major cell type, including epithelial
816   cells, endothelial cells, T cells, B cells, myeloid cells, and fibroblasts was reclustered and
817   reanalyzed to study each compartment at a higher resolution to detect granular cellular
818   heterogeneity in CRC. Clustree (v0.4.3) and manual review were utilized for optimal cluster
819   detection. For cell annotation of each cell type, we utilized published literature gene expression
820   signatures and manual review of differential genes among clusters. Additionally, we again utilized
821   SingleR (v0.99.10, R Package) for unbiased cell annotation. Interestingly, reclustering of major
822   compartments individually also detected clusters expressing hybrid markers as well as cell clusters
823   expressing markers from distinct lineages (such as T cell clusters expressing B cells); these were
824   removed and excluded for further analysis. We utilized UMAP for visualization purposes. For
825   validation, we analyzed 65,362 cells from 23 patients and applied the similar quality control metrics
826   as outlined above. In addition, we also applied *vars.to.regress* function to remove low quality cells
827   in an unbiased manner. We retained   31,383 high-quality single cells for further analysis[25].  These
828   high-quality cells were analyzed utilizing the same pipelines and parameters as that for our primary
829   cohort (**Supplementary Figs. 4-5 and 11-12**).
830
831   The InferCNV (v1.2.1) package was used with default paramets  to identify the evidence for
832   somatic large-scale chromosomal copy number alteration in epithelial cells (*EPCAM+, KRT8+,*
833   *KRT18+*)[82].  Normal epithelial cells were used as the control group.
834
835   ***Trajectory analysis.*** We used Monocle v.2 (v2.14.0), a reverse graph embedding method to
836   reconstruct single-cell trajectories in tumor and normal epithelium[83]. In brief, we used UMI count
837   matrices and the *negbinomial.size*() parameter to create a *CellDataSet* object in the default setting.
838   We grouped projected cells on UMAP in default settings for visualization of monocle results. We
839   defined the cumulative duration of the trajectory to show the average amount of transcriptional
840   transition that a cell undergoes as it passes from the starting state to the end state. The cells were
841   also ordered in pseudotime to explain the transition of cells from one state to another.
842
843   ***Pathway- Gene set variation analysis (GSVA).*** Pathway analysis was performed on the 50
844   hallmark gene sets downloaded from *Molecular Signatures Database (v7.2).* We used GSVA
845   (v1.34.0),  a non-parametric, unsupervised method to estimate the gene set variations and
846   evaluation of pathway enrichment, and pathway scores were calculated for each cell using standard
847   settings [29,33].
848
849   ***DNA and bulk RNA library construction.*** DNA and bulk RNA sequencing was performed as
850   previously described[84]. One hundred nanograms of DNA from each tumor was mechanically
851   sheared to an average size of 200 bp. Using the *KAPA Hyper Prep Pack*, DNA libraries were
852   packed, hybridized into the *xT probe* package, and amplified with the *KAPA HiFi HotStart*

853 *ReadyMix*. For uniformity, each sample needed to have 95% of all targeted base pairs sequenced
854 to a minimum depth of 300x. One hundred nanograms of RNA per tumor sample was heat
855 fragmented to a mean size of 200 base pairs in the presence of magnesium. Using random primers,
856 the RNA was used for first-strand cDNA synthesis, followed by second-strand synthesis and A-
857 tailing, adapter ligation, bead-based cleanup, and amplification of the library. After library
858 planning, the *IDT xGEN Exome Test Panel* was hybridized with samples. Streptavidin-coated beads
859 and target recovery were carried out, accompanied by amplification using the *KAPA HiFi* library
860 amplification package. The RNA libraries were sequenced on an *Illumina HiSeq 4000* using
861 patterned flow cell technology to achieve at least 50 million reads.

863 ***Detection of somatic variation on DNA sequencing data.*** The tumor and normal FASTQ files
864 were paired. For quality management measurement, FASTQ files were evaluated using FASTQC
865 and matched with Novoalign (Novocraft, Inc.)[84,85]. SAM files were generated and converted to
866 BAM files. The BAM files were sorted, and duplicates were marked. Single nucleotide variations
867 (SNVs) were called after alignment and sorting. For discovery of copy number alterations, the de-
868 duplicated BAM files and the VCF generated from the variant calling pipeline were processed to
869 compute read depth and variance of heterozygous germline SNVs between the tumor sample and
870 normal sample. Binary circular segmentation was introduced and segments with strongly
871 differential $\log_2$ ratios between the tumor and its comparator were chosen. From a combination of
872 differential coverage in segmented regions and estimation of stromal admixture provided by
873 analysis of heterozygous germline SNVs, an estimated integer copy number was determined

875 ***Microsatellite instability status.*** Probes for 43 microsatellite regions were developed using *Tempus*
876 *xT* assay[84]. Tumors were categorized into three groups by the MSI classification algorithm as
877 described by Tempus: microsatellite instability-high (MSI-H), microsatellite stable (MSS) or
878 microsatellite equivocal (MSE). MSI screening for paired tumor-normal patients used reads
879 mapped to the microsatellite loci with at least 5 bps flanking the microsatellite. The sample was
880 graded as MSI-H if there was a >70% chance of MSI-H classification. If the likelihood of MSI-H
881 status was 30-70%, the test findings were too ambiguous to interpret and those samples were listed
882 as MSE. If there was a <30% chance of MSI-H status, the sample was called MSS. Additionally,
883 IHC results were used to classify tumors into MSS or MSI molecular subtypes. Both of these
884 modalities were concordant and produced the same results.

886 ***Bulk transcriptomics analysis.*** We downloaded and pooled eight colorectal gene expression
887 datasets (GSE13067[86], GSE13294[86], GSE14333[48], GSE17536[49], GSE20916[87], GSE33113[88],
888 GSE35896[89], and GSE39582[11]), a pancreatic cancer dataset (GSE62165[90]) and a non-small cell
889 lung cancer dataset (GSE33532[91]) to validate our findings from the single cell compartments by
890 deconvoluting the bulk gene expression profiles into pseudo single-cell resolutions. We used Affy
891 (v1.64.0) for the data analysis and for exploration of Affymetrix oligonucleotide array probe level
892 data[92]. Batch correction was carried out using the removeBatchEffect (v3.42.2) function of the
893 LIMMA program and CMScaller for the CMS classification (see below)[93]. Three datasets
894 (GSE17536[49], GSE33113[88], and GSE39582[11]) were utilized for clinical outcome analysis[93,94].

896 ***Correlation patterns in bulk gene expressions for CAF compartments.*** To identify the top
897 correlated CAF-marker genes within the combined eight CRC datasets, three bulk CRC gene
898 expression sets individually, pancreas cancer and lung cancer datasets. We first transformed the
899 bulk gene expression sets with $\log_2$ transformation. Next, marker genes with an average $\log_2$ FC>/=
900 0.5 and p<0.05 obtained from the single cell data of CAF-S1 and CAF-S4 compartments were
901 separately intersected with the bulk gene expression sets. Genes with an average Spearman
902 correlation score greater than 0.8 were kept as the CAF signatures within the bulk gene expression.
903 Heatmaps illustrating the correlation patterns within and between the CAF compartments were

904 prepared with the heatmap.2 function from ggplot package (v3.1.1) utilizing the Pearson correlation
905 coefficient. Heatmaps illustrating the correlation patterns within and between the CAF
906 compartments were prepared using the ggplot package (v3.1.1) utilizing the Pearson correlation
907 coefficient[95].

908

909 The Cox proportional hazard regression model was used to examine the significance of 20 cell
910 types from scRNA-seq in bulk expression data. Each cell type's marker genes with an average
911 logFC>1 and adjusted P<0.05 were intersected with the bulk expression datasets separately. We
912 only kept the marker genes with a high correlation with each other in bulk, which provides an
913 average correlation score of $> 0.8$. The average bulk expression of each cell type's remaining marker
914 genes was calculated and used in the hazard regression model as the representative of this cell type.
915 For analysis of relationships with patient outcome, univariate models were calculated using Cox
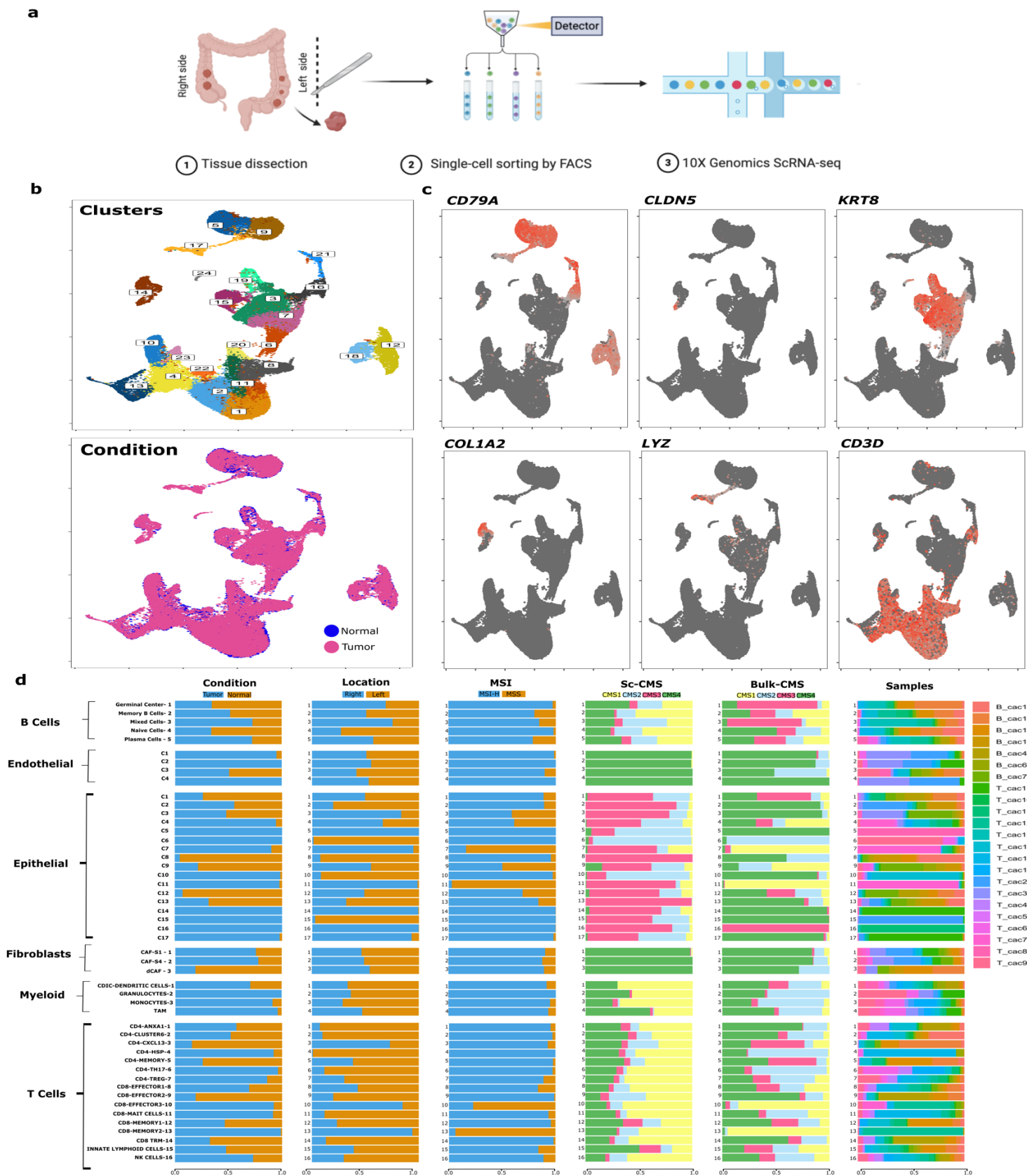916 proportional hazard regression (coxph function from survival R package)[96].

917

918 ***Deconvoluting public bulk gene expression profiles into pseudo single-cell expressions.*** We used
919 CIBERSORTx v1 to estimate composition of various cell populations in GSE39582[11],
920 GSE14333[48], and GSE17536[49] datasets[47]. Signature gene matrices were created using the
921 expression profiles of 49,859 cells as the reference single cell profile. We ran the 'hires' module
922 with default parameters except for the 'rmbatchBmode,' and the bulk-mode batch correction
923 argument was set to true. After the deconvolution process, we normalized the gene expressions
924 according to the cell fractions in each sample and calculated each gene's Z-transformed expression
925 values. The average normalized expression of each cell type across all samples was plotted with
926 the heatmap.3 R function of the GMD package (v0.3.3)[97]. A signature matrix highlighting marker
927 genes of the different cell types was prepared with a heatmap.2 R function of ggplot (v3.1.1).

928

929 ***Consensus molecular subtyping of colorectal cancer (CMS Classification).*** We used R package
930 CMScaller(v0.9.2), a nearest template prediction (NTP) algorithm, for the classification of gene
931 expression datasets[94]. We set the permutation number to 1000 to predict the CMS classes of the
932 samples in the GEO datasets with a p-value $< 0.05$. We ran CMScaller with default parameters.

933

934 ***Continuous subtype discovery using scRNA-seq analysis.*** Bulk mRNA expression profiles of the
935 GEO dataset (GSE39582[11]), composed of 585 samples in total, were deconvoluted into the pseudo
936 single-cell expression profiles via CIBERSORTx utilizing the expression data consisting of 20
937 different well known cell types from our scRNA-seq dataset[47]. We transformed the deconvoluted
938 expression matrix with log2 transformation. The principal components cluster subtype scores
939 (PCSSs) of the CMS among the samples, were determined separately for each cell type using an
940 algorithm published by Ma et al[51]. To obtain the PCSSs, the average loading vectors were used.
941 The results obtained for 20 cell types were projected on the first two PCSSs (PCSS1 and PCSS2)
942 as they were validated by Ma et al. in their analysis using 18 datasets. We also analyzed two datasets
943 (GSE14333[48], GSE17536[49]) to independently confirm reproducibility of continuous scores.

944

945 ***Statistics and reproducibility.*** All statistical analyses and graphs were created in R (v3.6.3) and
946 using a Python-based computational analysis tool. Schematic representations were made using the
947 Inkscape (https://inkscape.org/) software. Dim plots, bar plots and box plots were generated using
948 the dittoSeq (v1.1.7) package with default parameters[98]. Violin plots were generated using the
949 patchwork (v1.1.0) package and ggplot2 (v3.3.2) package in R with default parameters. Heatmaps
950 were generated using Morpheus.R with default parameters[99,100]. ANOVA and pair-wised t-tests for
951 the CMS classes across the deconvoluted expression profiles were performed in R using the ggpubr
952 R (v0.4.0) package[101]. The Box and Whisker plots were generated using the boxplot function of the
953 R base package at default parameters. The mean of the $\log_2$ transformed deconvoluted expression
954 value of the samples in each CMS group was demonstrated with a horizontal straight line within

955    each box. The length of a boxplot corresponds to the interquartile range (IQR), which is defined as
956    the range between the first and third quartiles (Q1 and Q3), whereas the whiskers are the upper and
957    lower extreme values of the data (either data's extremum values, or the Q3+1.5*IQR and Q1-
958    1.5*IQR values, whichever was less extreme). To test for differential score based on CIBERSORTx
959    cell abundances between different molecular subtypes pairwise t-test with Benjamini-Hochberg
960    correction was used.
961
962    ***Survival analysis.*** Survival curves were obtained according to the Kaplan-Meier method survfit
963    (v3.2-7), and differences between survival distributions were assessed by Log-rank test. The
964    patients were divided into two groups (high/poor and low/good risk) according to their median
965    expression values (survminer (v0.4.8)). The surv_cutpoint function uses the maximally selected
966    rank statistics and implements standard methods for the approximation of the null distribution of
967    maximally selected rank statistics (maxstat (v0.7-25)).
968
969    The proportional hazard assumption was tested to examine the fit of the model for survival of the
970    samples in three GEO datasets (GSE17536[49], GSE33113[88], and GSE39582[11]) with respect to the
971    deconvoluted bulk mRNA expressions. For analysis of the relationships with patient outcome,
972    multivariate models were calculated using the Cox proportional hazard regression (coxph survival
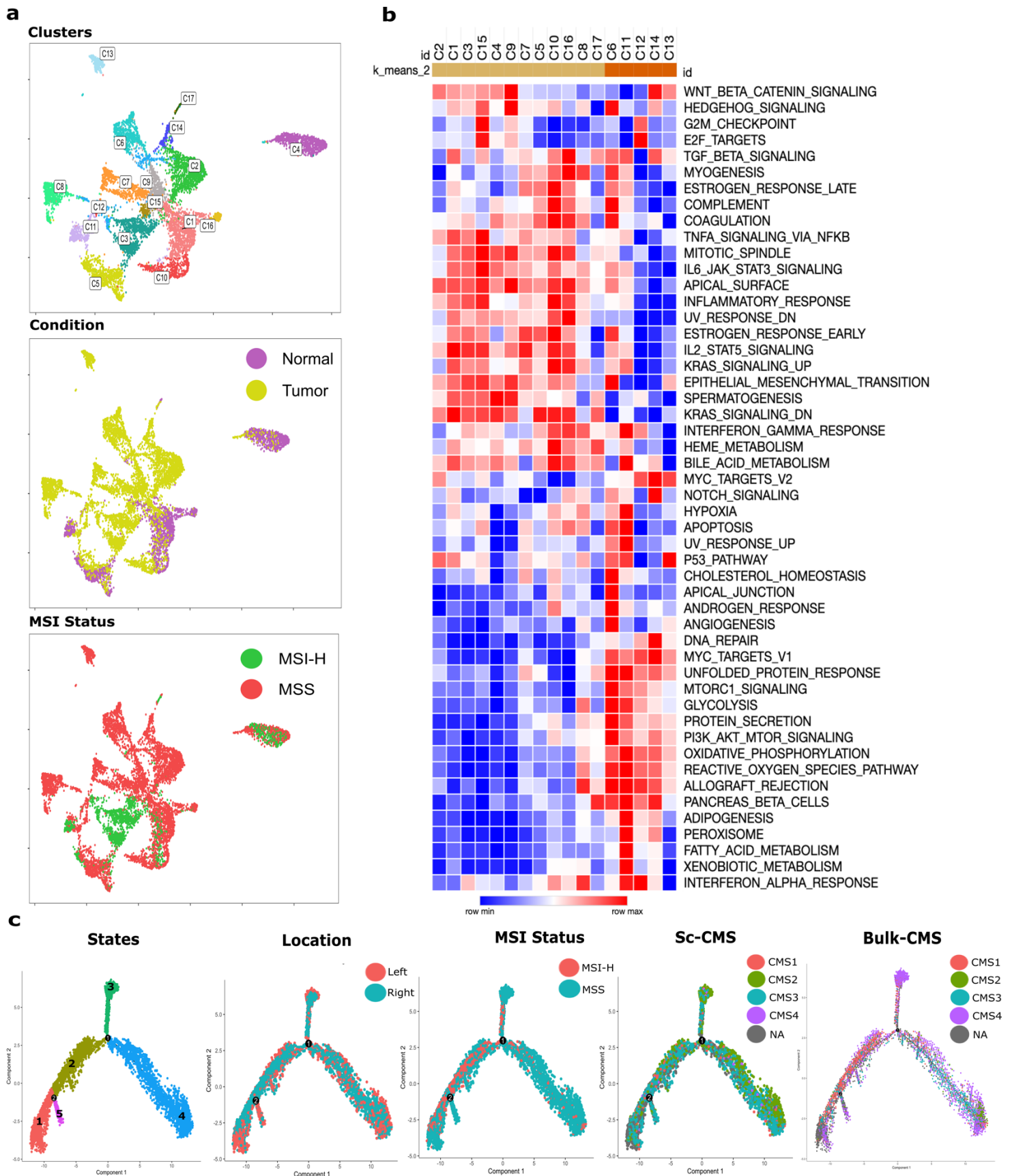973    R package)[96].
974

# Figure 1
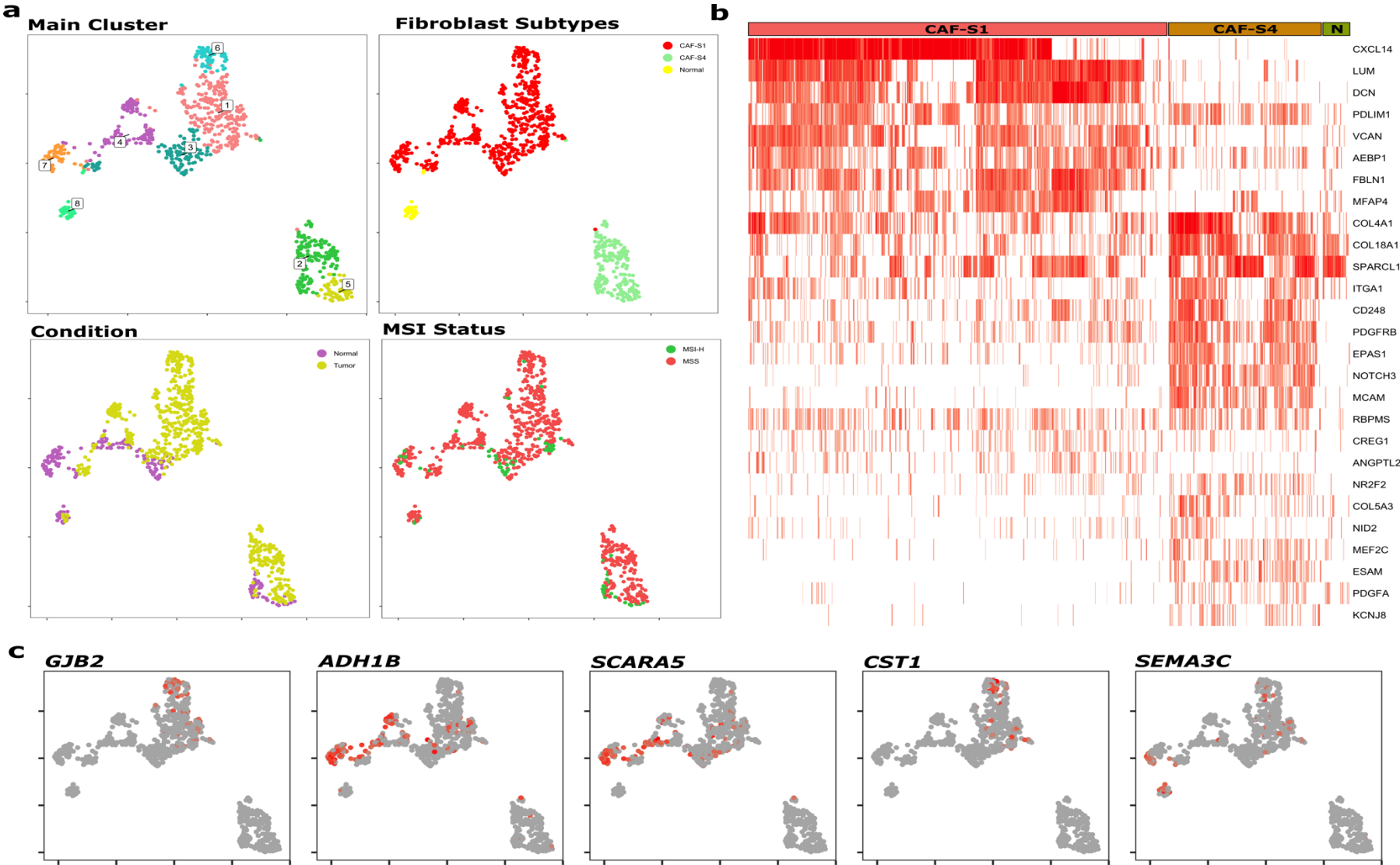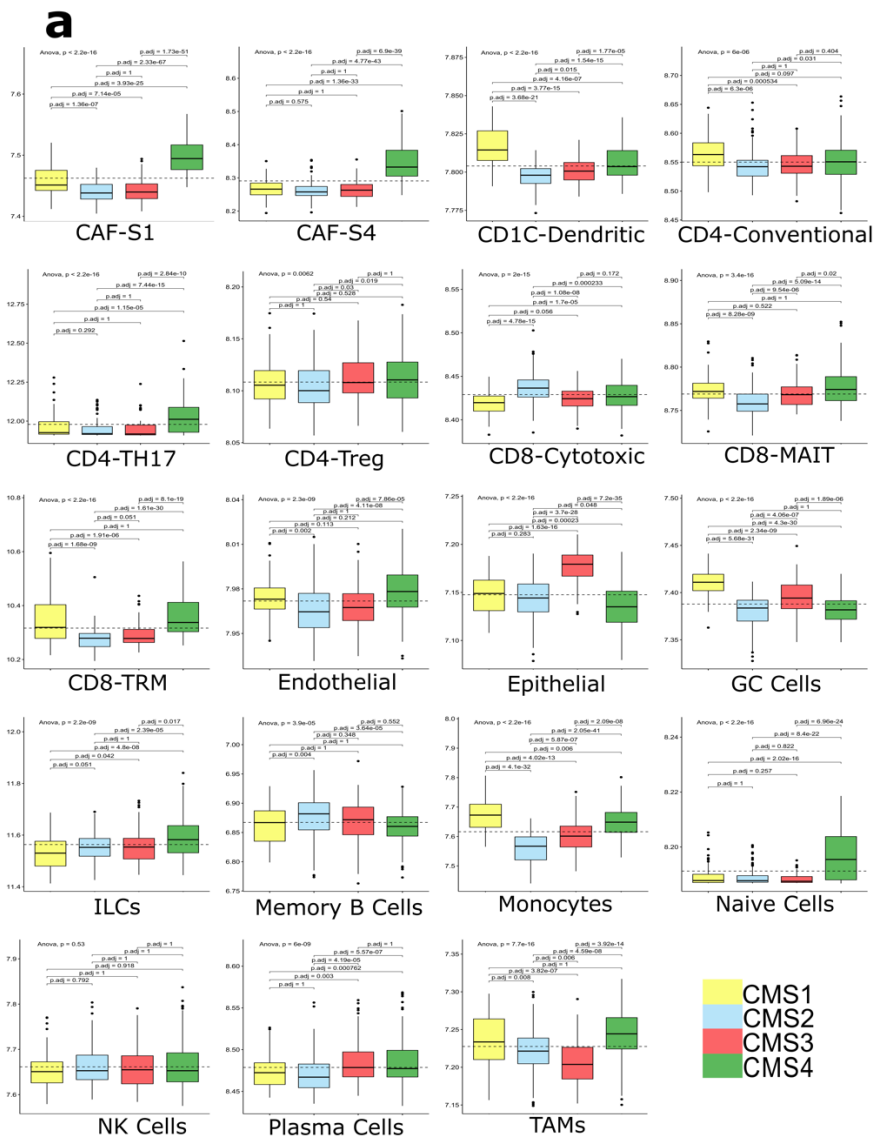
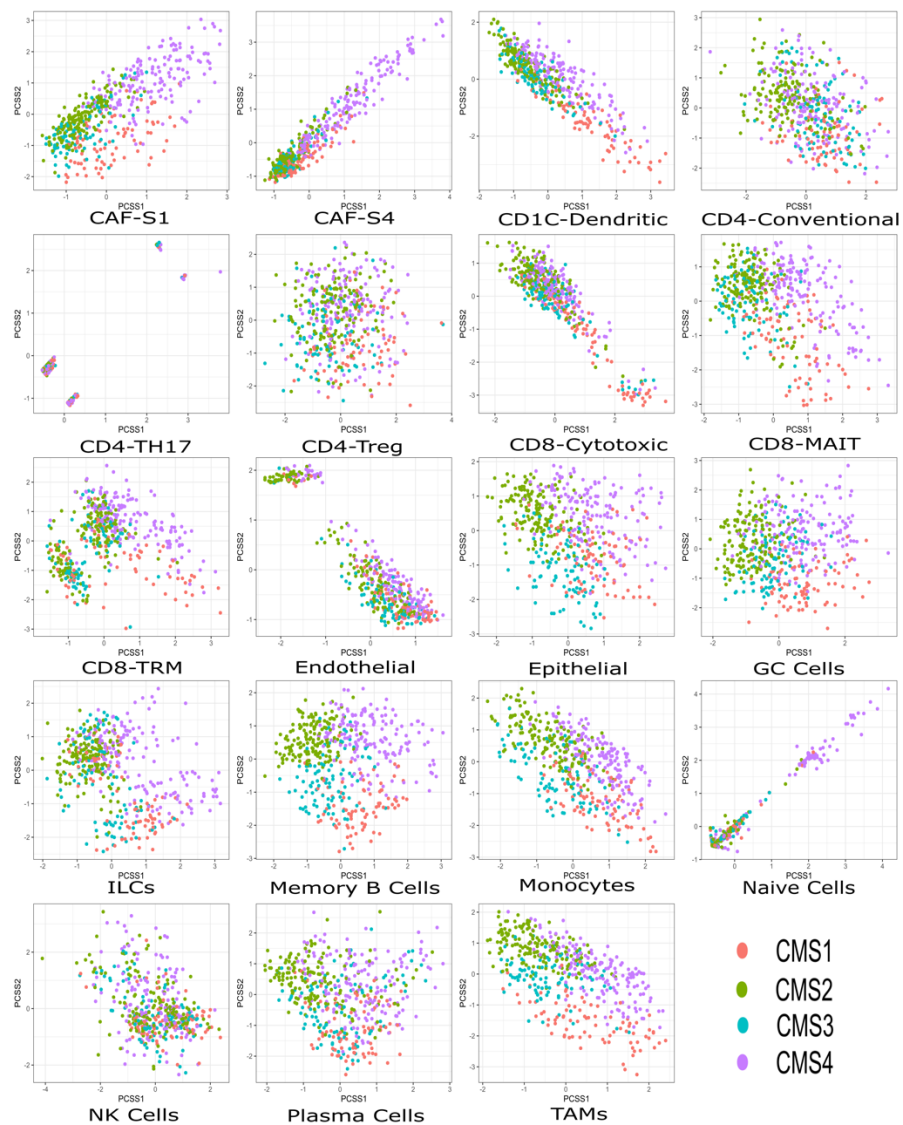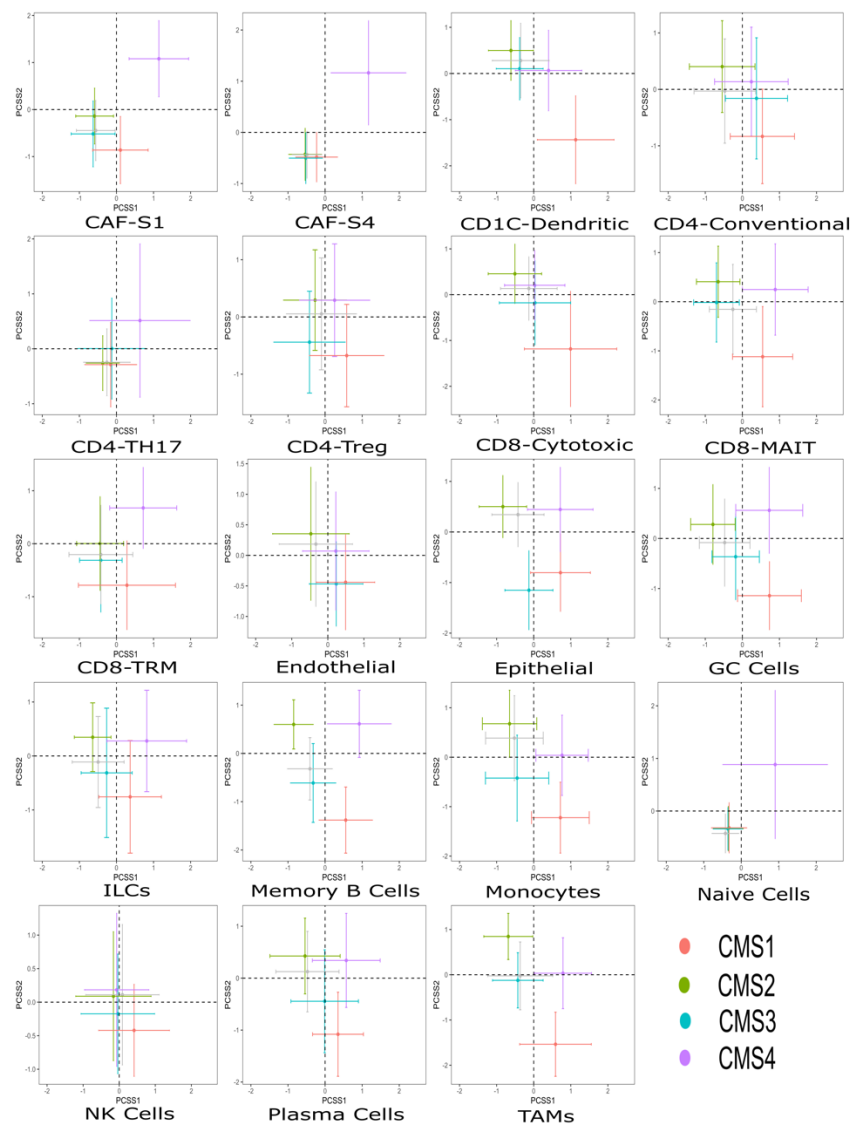# Figure 2

# Figure 3

# Figure 4
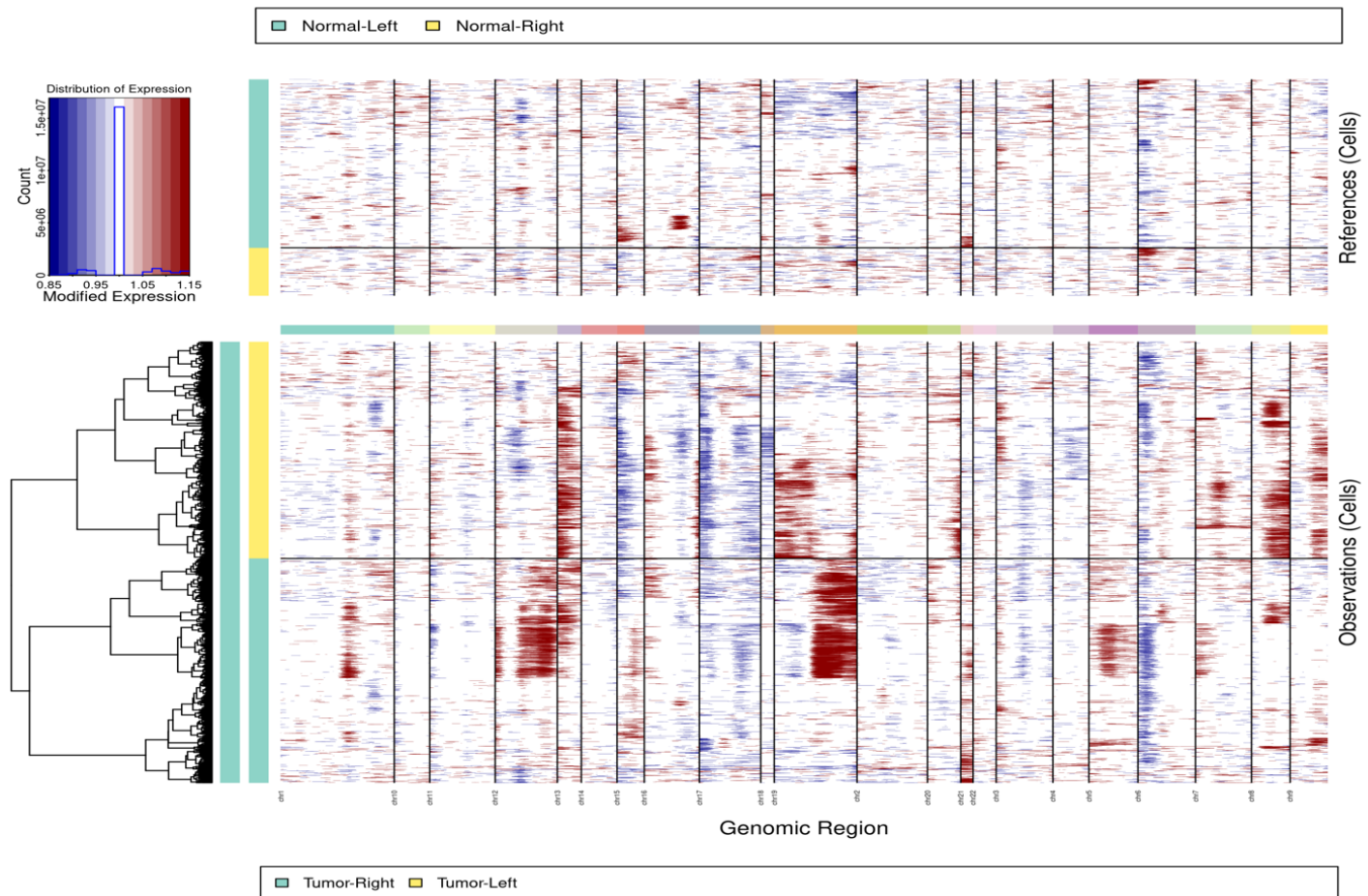
# Figure 5

# Figure 6

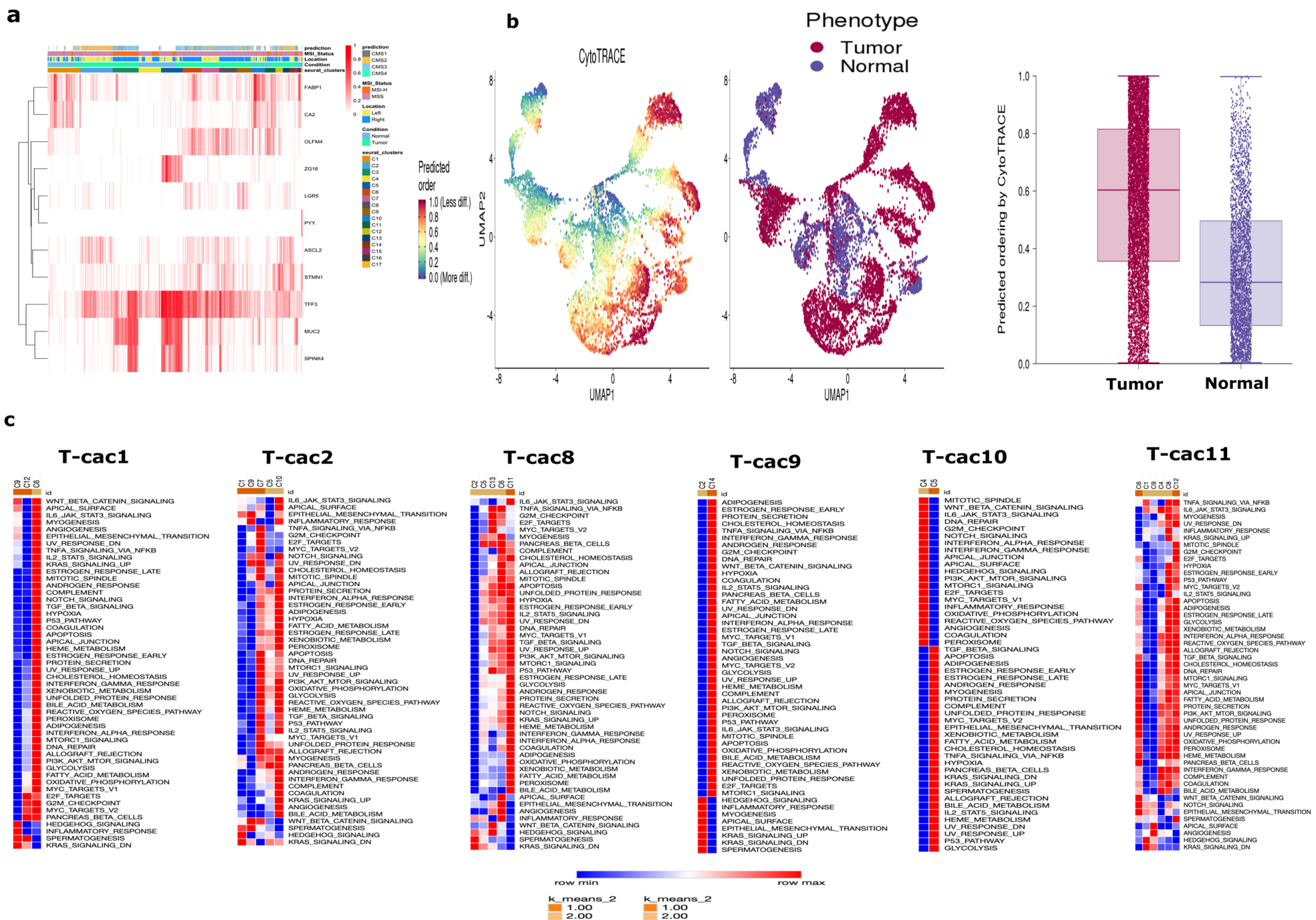**Supplementary Materials**

Khaliq et al.

Title: Redefining tumor classification and clinical stratification through a colorectal cancer single-cell atlas
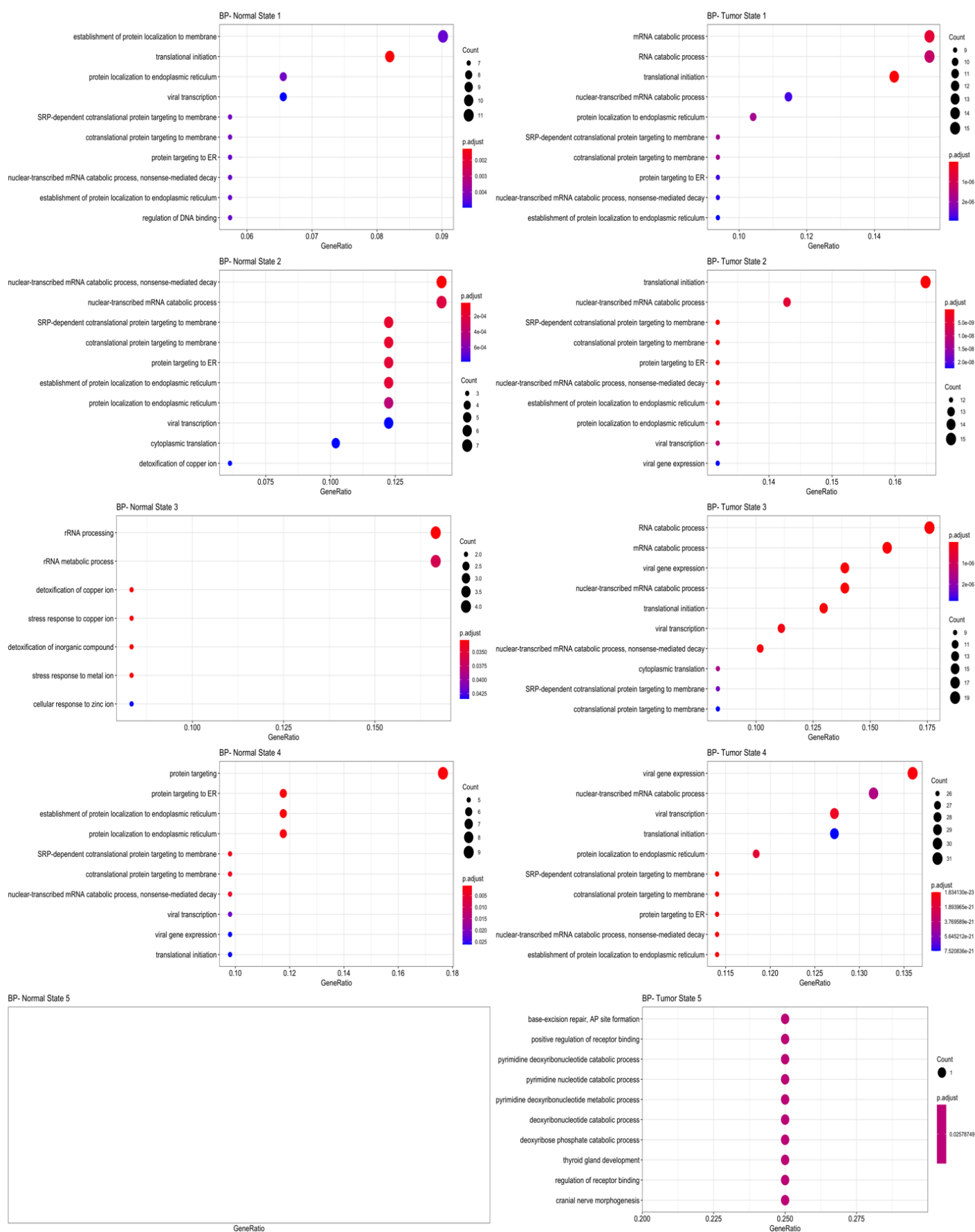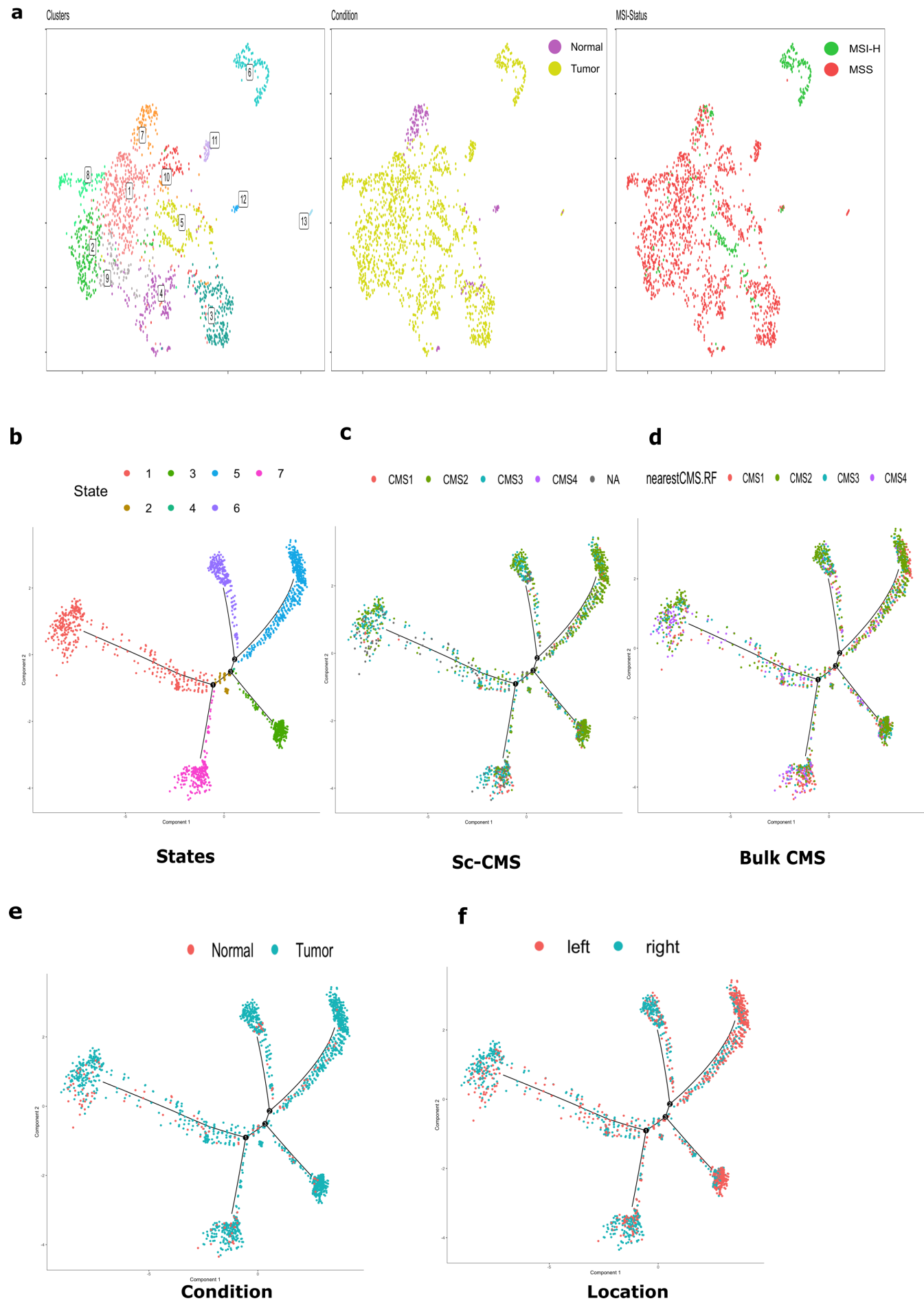
*Supplementary figures 1-12*

**Supplementary figure 1. Analysis of copy number variation amongst epithelial cells.** CNV analysis was conducted on the epithelial cell compartment. Increase CNV is seen in tumor derived epithelial cells (observation). Non-malignant derived epithelial cells were used as control (references).

**Supplementary figure 2**. **a,** Heatmap of marker gene expression for the epithelial and tumor cells. **b,** UMAP visualization of computational analysis of differentiation status using CytoTRACE (see methods). **c,** Heatmap representation of Hallmark pathway analysis of epithelial phenotypes within six different tumor samples demonstrating subclonal intratumoral heterogeneity.
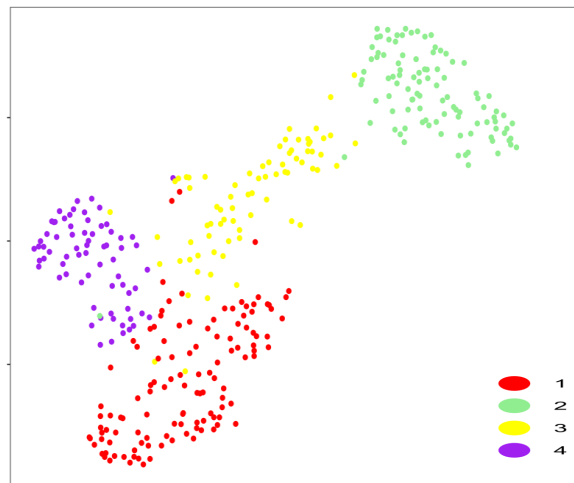
**Supplementary figure 3.** Pathway analysis (GO ontology) of gene-expressions specific to each cell-state in trajectory analysis stratified by malignant and non-malignant sample of origin.

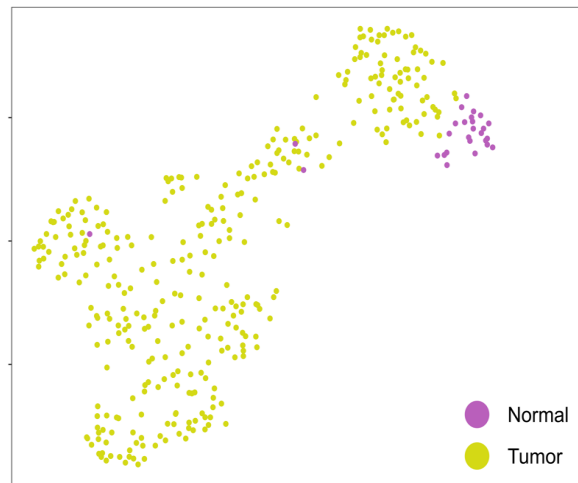**Supplementary figure 4 Epithelial clustering and differentiation trajectories from a validating cohort (Lee et al. data). a,** UMAP clustering of epithelial cells colored by cluster, sample of origin tumor vs. normal tissue status and MSI status. **b-f,** Differentiation trajectories of epithelial cells colored by differentiation state, CMS status, tumor status and colonic location.

**a**

Clusters

Conditions

MSI Status

- 1 (red)
- 2 (green)
- 3 (yellow)
- 4 (purple)

- Normal
- Tumor

- MSI-H
- MSS

**b**

Fibroblast: CAF-S1

Fibroblast: CAF-S4

Identity

Features: FBLN1, MFAP4, DCN, LUM, VCAN, AEBP1, CREG1, CXCL14

Features: SPARCL1, COL5A3, NID2, COL4A1, CD248, KCNJ8, RBPMS, MEF2C, PDLIM1, MCAM, NOTCH3, EPAS1, COL18A1, NR2F2, PDGFA, ANGPTL2, ESAM, ITGA1, PDGFRB

Scaled Average Expression
1.0
0.5
0.0
-0.5
-1.0

Percent Expressed
0
25
50
75

**Supplementary figure 5. Characterization of fibroblasts and their transcriptomic expression patterns (Lee et al. data)**. **a,** Fibroblasts colored by distinct groups, tumor vs normal sample, and sample specimen. **b,** Dot plot demonstrating variable expression patterns of subtypes of CAF-S1 and CAF-S4 confirming their relevance in colorectal cancer.

**Supplementary Figure 6 Association between relative cell abundance and patient survival from microarray-based datasets.** a. GSE17536 (n=177). b GSE39582 (n=585). c. GSE33113 (n=96). Note that CAF-S4 is not significant in GSE33113 (p=0.093).

**Supplementary Figure 7. Cell abundance prediction for each sample and projected on bulk CMS on GSE17536 (n= 177). a,** Boxplots show the distribution of cell types within tumors with varying CMS status. The whiskers depict the 1.5 x IQR. The p-values for pairwise t-tests comparisons (with Benjamani-Hochberg correction) and ANOVA tests of cell abundance across CMS are shown in the figure. **b,** Deconvolution heatmap of different cell types by average expression using CIBERSORTx demonstrating cell type distribution within each CMS category. ILCs = Innate lymphoid cells, GC Cells = Germinal Center B Cells, NK Cells = Natural Killer Cells, TAMs = Tumor Associated Macrophages.

**a**

CAF-S1, CAF-S4, CD1C-Dendritic, CD4-Conventional, CD4-TH17, CD4-Treg, CD8-Cytotox, CD8-Mait, CD8-TRM, Endothelial, Epithelial, GC Cells, ILCs, Memory B Cells, Monocytes, Naive Cells, NK Cells, Plasma Cells, TAMs

CMS1, CMS2, CMS3, CMS4

**b**

CMS1  CMS2  CMS3  CMS4  NA

Memory B Cells, CAF-S1, TAMs, CD4-Treg, GC Cells, Monocytes, NK Cells, Endothelial, Plasma Cells, Epithelial, CAF-S4, Naive Cells, CD1C-Dendritic, CD4-Conventional, CD8-Memory, CD8-Cytotox, CD8-Mait, CD8-TRM, ILCs, CD4-TH17

Z-Score
-4    +4

**Supplementary Figure 8. Cell abundance prediction for each sample and projected on bulk CMS on GSE14333 (n= 290). a,** Boxplots show the distribution of cell types within tumors with varying CMS status. The whiskers depict th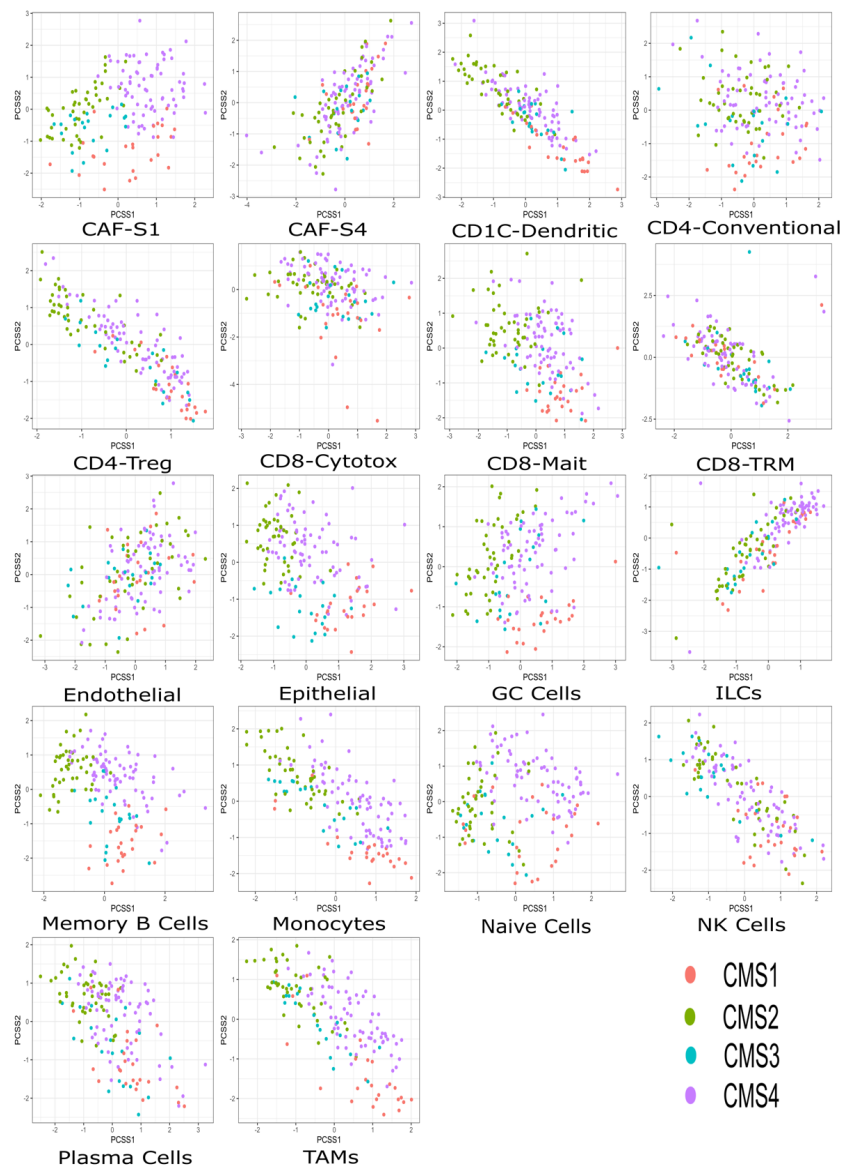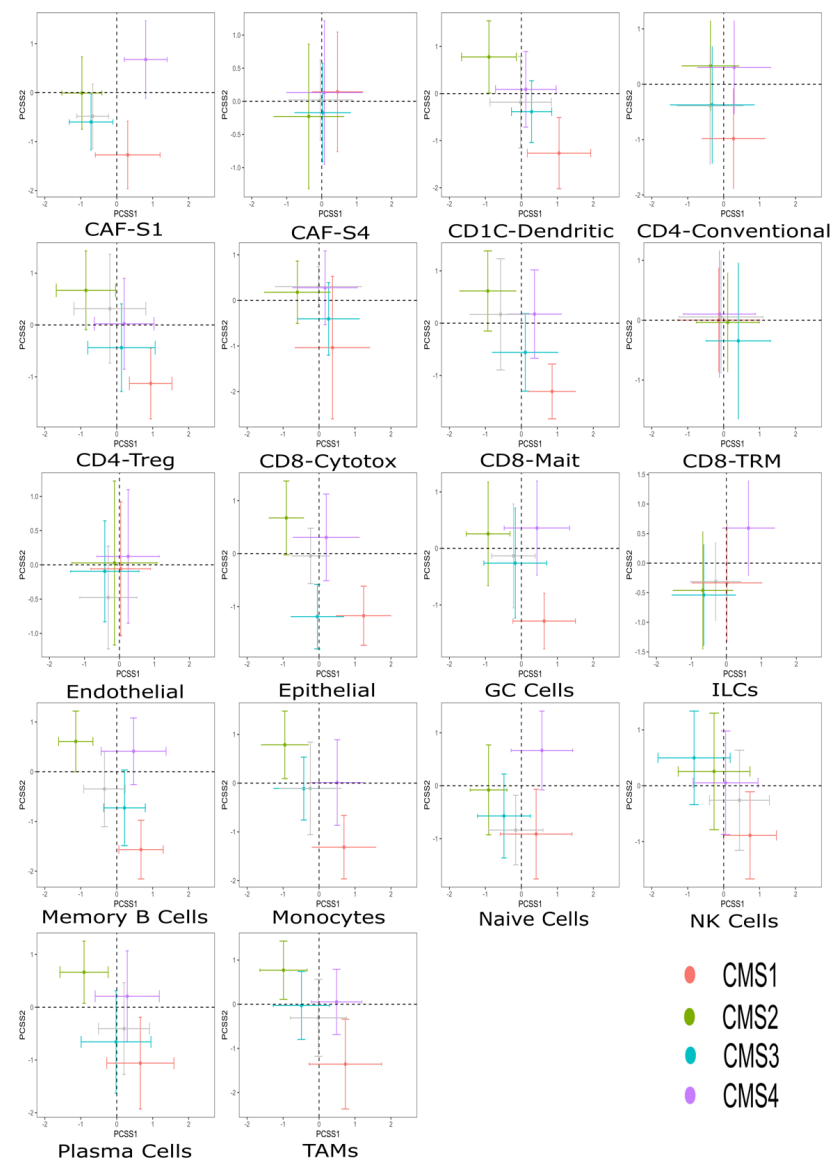e 1.5 x IQR. The p-values for pairwise t-tests comparisons (with Benjamani-Hochberg correction) and ANOVA tests of cell abundance across CMS are shown in the figure. **b,** Deconvolution heatmap of different cell types by average expression using CIBERSORTx demonstrating cell type distribution within each CMS category. ILCs = Innate lymphoid cells, GC Cells = Germinal Center B Cells, NK Cells = Natural Killer Cells, TAMs = Tumor Associated Macrophages.

**a**

CAF-S1    CAF-S4    CD1C-Dendritic    CD4-Conventional

CD4-Treg    CD8-Cytotox    CD8-Mait    CD8-TRM

Endothelial    Epithelial    GC Cells    ILCs

Memory B Cells    Monocytes    Naive Cells    NK Cells

Plasma Cells    TAMs

CMS1
CMS2
CMS3
CMS4

**b**

CAF-S1    CAF-S4    CD1C-Dendritic    CD4-Conventional

CD4-Treg    CD8-Cytotox    CD8-Mait    CD8-TRM

Endothelial    Epithelial    GC Cells    ILCs

Memory B Cells    Monocytes    Naive Cells    NK Cells

Plasma Cells    TAMs

CMS1
CMS2
CMS3
CMS4

**Supplementary Figure 9. Continuous scores for CRC dataset GSE17536 (n= 177). a,** Principal component analysis plot showing PCSS1 and PCSS2 continuous scores reported by CMS classification across cell types show minimal separation in the top 2 principal components. **b,** All cell types projected on four quadrants representing CMS1-4 using PCSS1 and PCSS2 scores. Note that the cell types largely form a continuum along CMS status and are not clustered in discrete quadrants separate from one another. Cells and markers are colored by bulk CMS status accordingly to the tumor sample of origin. ILCs = Innate lymphoid cells, GC Cells = Germinal Center B Cells, NK Cells = Natural Killer Cells, TAMs = Tumor Associated Macrophages.

**a**

CAF-S1    CAF-S4    CD1C-Dendritic    CD4-Treg

CD8-Cytotox    CD8-Mait    Endothelial    Epithelial

GC Cells    Memory B Cells    Monocytes    Naive Cells

NK Cells    Plasma Cells    TAMs

CMS1
CMS2
CMS3
CMS4

**b**

CAF-S1    CAF-S4    CD1C-Dendritic    CD4-Treg

CD8-Cytotox    CD8-Mait    Endothelial    Epithelial

GC Cells    Memory B Cells    Monocytes    Naive Cells

NK Cells    Plasma Cells    TAMs

CMS1
CMS2
CMS3
CMS4

**Supplementary Figure 10.  Continuous scores for CRC dataset GSE14333 (n= 290). a,** Principal component analysis plot showing PCSS1 and PCSS2 continuous scores reported by CMS classification across cell types show minimal separation in the top 2 principal components. **b,** All cell types projected on four quadrants representing CMS1-4 using PCSS1 and PCSS2 scores. Note that the cell types largely form a continuum along CMS status and are not clustered in discrete quadrants separate from one another. Cells and markers are colored by bulk CMS status accordingly to the tumor sample of origin. ILCs = Innate lymphoid cells, GC Cells = Germinal Center B Cells, NK Cells = Natural Killer Cells, TAMs = Tumor Associated Macrophages.

**a**



**b**
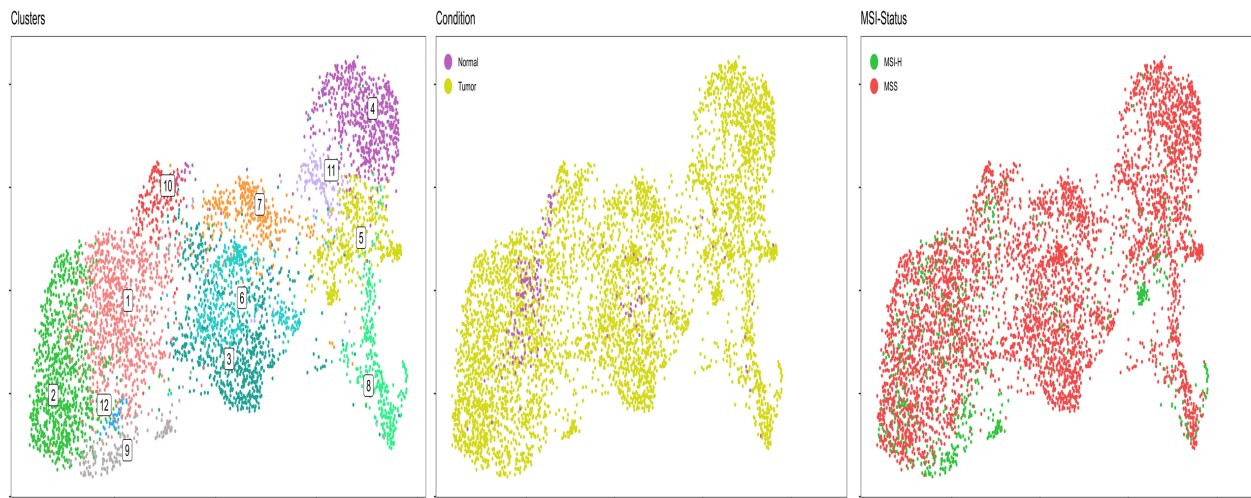


**Supplementary figure 9. T/NK cell data from Lee et al. a**, Reclustering of T/NK cells and coloring by clusters, tumor and normal status, MSI status and samples. **b**, Bar chart representation of cells colorized by our samples of origin, MSI status, tumor status, colonic location, CMS status.

**Supplementary figure 10. T cell expression analyses across T/NK cells . a,** SingleR heatmap cell type identification within each cluster. Note that genes associated with T cells show increased expressions, confirming the quality of the data. Note doublets were removed from the further analysis. **b,** T/NK cells gene specific expression to identify T cell heterogeneity using published literature (see methods).

**Extended Data Figure 1**. **T cell identification and characterization. a**, UMAP reclustering of T cells colored by cell phenotype, tissue malignancy status, and sample of origin. **b**, Violin plots showing the differential expression of T cell-specific marker genes between CD4 and CD8 phenotypes. **c**, Heatmap of the Hallmark pathway analysis for the T cell compartment. scCMS, single-cell consensus molecular subtyping; bulk CMS, consensus molecular subtyping on Bulk RNA-seq data.

**Description:** We reclustered and analyzed 22,525 cells from both tumors and adjacent, normal tissue samples and identified 11 CD4+ T cell and 12 CD8+ T cell clusters. We used known canonical markers

and published expression signatures to identify T cell states for further analysis (see methods). We identified conventional CD4+ T cells, CD4+ Tregs, CD8+ (naive/memory, cytotoxic, resident memory, and MAIT cells), NK cells, and innate lymphoid cells (ILC). Among the conventional CD4+ T cells, we identified the central memory/naive like state (CCR7+, SELL+, and TCF7+) enriched in non-malignant samples. In contrast, Th17 cells expressing IL-17, known as critical anti-tumor effectors, were enriched in tumor samples. CD4+ Tregs (FOXP3+) expressing immune checkpoint markers and costimulatory molecules were among the most abundant T cells in the colorectal TME compared to non-malignant tissue. Among the CD8+ T cell states, CD8+ cytotoxic cells were distributed across three clusters that we labeled CD8+ effector 1, CD8+ effector 2, and CD8+ effector 3. CD8+ effector clusters expressed cytotoxicity genes and chemokines as previously described in other tumor types. CD8 effector3 was predominantly enriched in MSI-H CRC patients and represented 77% of the total CD8+ effector 3 population among the 2 MSI-H CRC samples. This cluster expressed ITGAE, LAYN, CXCL13, and T cell exhaustion markers (LAG3, HAVCR2, and CD96), possibly explaining this CD8+ cell state's role in the response to immune checkpoint inhibitors in MSI-H colorectal tumors. Gene-set enrichment of CD8+ effectors further confirmed their distinct states. CD8+ effector 2 was a proliferative cluster with MYC activity, NOTCH activation, and EF2 targets. CD8+ effector(s), CD8+ MAIT cells, and NK cells were enriched in tumors, whereas Tissue-resident memory (Trm) cells were depleted in tumor tissue. Trm induction was recently seen to enhance cancer vaccine efficacy in other tumors, suggesting a possible therapeutic target in CRC.
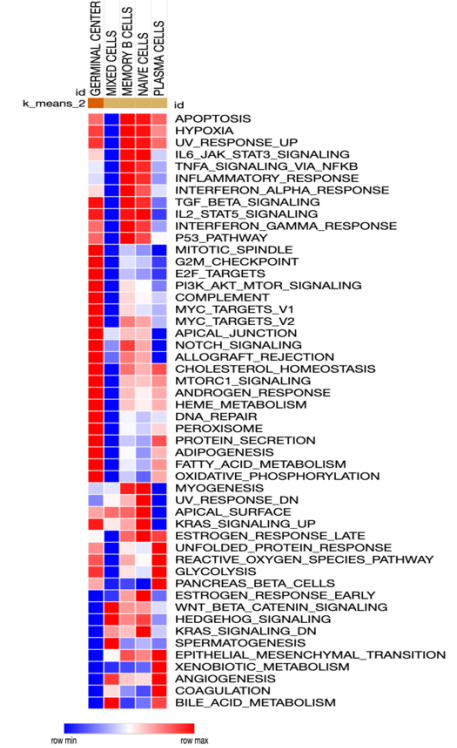
**a**

Clusters

Condition
- Normal
- Tumor

MSI-Status
- MSI-H
- MSS

**b**

Clusters
1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15

Endothelial_cells:blood_vessel
Endothelial_cells:lymphatic:KSHV
Endothelial_cells:lymphatic
Endothelial_cells:lymphatic:TNFa_48h
Endothelial_cells:HUVEC:Borrelia_burgdorferi
Endothelial_cells:HUVEC:IFNg
Endothelial_cells:HUVEC:VEGF
Endothelial_cells:HUVEC:FPV-infected
Endothelial_cells:HUVEC:IL-1b
Endothelial_cells:HUVEC:Serum_Amyloid_A
Endothelial_cells:HUVEC:PR8-infected
Endothelial_cells:HUVEC:H5N1-infected
Endothelial_cells:HUVEC
Endothelial_cells:HUVEC:B._anthracis_LT
Fibroblasts:foreskin
Tissue_stem_cells:BM_MSC:BMP2
Tissue_stem_cells:BM_MSC:TGFb3
Neurons:Schwann_cell
Chondrocytes:MSC-derived
iPS_cells:CRL2097_foreskin
Osteoblasts
Embryonic_stem_cells
Hepatocytes
Keratinocytes:IFNg
Keratinocytes:IL22
Keratinocytes:IL24
Keratinocytes:KGF
Keratinocytes
Keratinocytes:IL26
Keratinocytes:IL19
Keratinocytes:IL20
Epithelial_cells:bladder
Epithelial_cells:bronchial
NK_cell
NK_cell:CD56hiCD62L+
Pre-B_cell_CD34-
Pro-B_cell_CD34+
CMP
HSC_CD34+
T_cell:CD8+
GMP
T_cell:gamma-delta
B_cell:Memory
B_cell:immature
B_cell:Naive
B_cell:Plasma_cell
B_cell
B_cell:CXCR4+_centroblast
B_cell:CXCR4-_centrocyte
B_cell:Germinal_center

**c**

GERMINAL CENTER
MIXED CELLS
MEMORY B CELLS
NAIVE CELLS
PLASMA CELLS

k_means_2    id

APOPTOSIS
HYPOXIA
UV_RESPONSE_UP
IL6_JAK_STAT3_SIGNALING
TNFA_SIGNALING_VIA_NFKB
INFLAMMATORY_RESPONSE
INTERFERON_ALPHA_RESPONSE
TGF_BETA_SIGNALING
IL2_STAT5_SIGNALING
INTERFERON_GAMMA_RESPONSE
P53_PATHWAY
MITOTIC_SPINDLE
G2M_CHECKPOINT
E2F_TARGETS
PI3K_AKT_MTOR_SIGNALING
COMPLEMENT
MYC_TARGETS_V1
MYC_TARGETS_V2
APICAL_JUNCTION
NOTCH_SIGNALING
ALLOGRAFT_REJECTION
CHOLESTEROL_HOMEOSTASIS
MTORC1_SIGNALING
ANDROGEN_RESPONSE
HEME_METABOLISM
DNA_REPAIR
PEROXISOME
PROTEIN_SECRETION
ADIPOGENESIS
FATTY_ACID_METABOLISM
OXIDATIVE_PHOSPHORYLATION
MYOGENESIS
UV_RESPONSE_DN
APICAL_SURFACE
KRAS_SIGNALING_UP
ESTROGEN_RESPONSE_LATE
UNFOLDED_PROTEIN_RESPONSE
REACTIVE_OXYGEN_SPECIES_PATHWAY
GLYCOLYSIS
PANCREAS_BETA_CELLS
ESTROGEN_RESPONSE_EARLY
WNT_BETA_CATENIN_SIGNALING
HEDGEHOG_SIGNALING
KRAS_SIGNALING_DN
SPERMATOGENESIS
EPITHELIAL_MESENCHYMAL_TRANSITION
XENOBIOTIC_METABOLISM
ANGIOGENESIS
COAGULATION
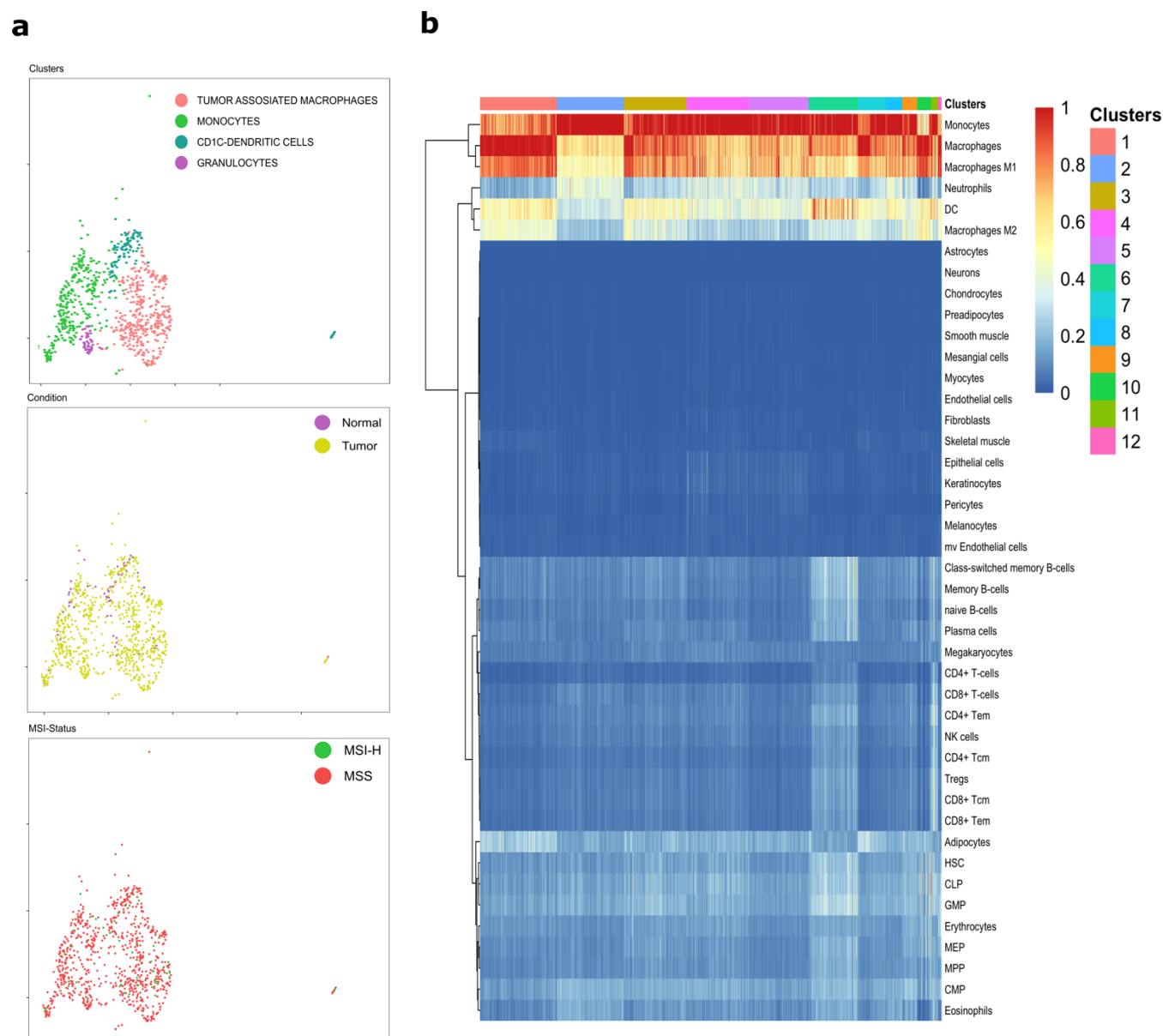BILE_ACID_METABOLISM

row min — row max

**Extended Data Figure 2. Reclustering of B cells with characterization of clusters and phenotypes**. **a**, UMAP depiction of B cell reclustering colored by cluster, malignancy status, and sample of origin. **b**, SingleR heatmap demonstration of B cell distribution within each cluster. **c**, B cell Hallmark pathway analysis by phenotype. scCMS, single-cell consensus molecular subtyping; bulk CMS, consensus molecular subtyping on Bulk RNA-seq data.

**Description:** To illustrate characteristics of B cells in CRC we reclustered 9,289 B cells that clearly identified naive cells, memory cells, plasma cells, and germinal center (GC) B cells. All B cell subtypes from the CRC TME and the non-malignant colonic tissue clustered together exhibiting transcriptional similarity among non-tumor and tumor-derived cells. Memory B cells and plasma cells were enriched in tumors, while naive and GC B cells were enriched in nonmalignant tissue.

**Extended Data Figure 3**. **Reclustering of the myeloid cell compartment. a**, UMAP depiction of myeloid cell reclustering colored by subtype, malignancy status, and sample of origin. **b**, SingleR heatmap demonstration of myeloid cell distribution within each cluster. scCMS, single-cell consensus molecular subtyping; bulk CMS, consensus molecular subtyping on Bulk RNA-seq data.

**Description**: We reclustered 819 myeloid cells and identified CD1C+ dendritic cells, tumor associated macrophages (TAM and MRC1+), monocytes (S100A8+), and granulocyte clusters. We recovered key cell types including M2 polarized macrophages, as seen in other tumor types. Monocytes revealed proinflammatory phenotypes (1L1B, S100A8, and S100A9), while TAM showed anti-inflammatory signatures (APOE, SEPP1, and CD163) consistent with the role of TAM in immune suppression and cancer progression (see methods).