

# Northumbria Research Link

Citation: Khaliq, Ateeq M, Kurt, Zeyneb, Grunvald, Miles W, Erdogan, Cihat, Turgut, Sultan Sevgi, Rand, Tim, Khare, Sonal, Borgia, Jeffrey A, Hayden, Dana M, Pappas, Sam G, Govekar, Henry R, Bhama, Anuradha R, Singh, Ajaypal, Jacobson, Richard A, Kam, Audery E, Zloza, Andrew, Reiser, Jochen, Catenacci, Daniel V, Turaga, Kiran, Radovich, Milan, Thyparambil, Sheeno, Levy, Mia A, Subramanian, Janakiraman, Kuzel, Timothy M, Sadanandam, Anguraj, Hussain, Arif, El-Rayes, Bassel, Salahudeen, Ameen and Masood, Ashiq (2021) Redefining tumor classification and clinical stratification through a colorectal cancer single-cell atlas. bioRxiv. p. 429256. (Submitted)

Published by: Cold Spring Harbor Laboratory

URL: <https://doi.org/10.1101/2021.02.02.429256>  
<<https://doi.org/10.1101/2021.02.02.429256>>

This version was downloaded from Northumbria Research Link:  
<http://nrl.northumbria.ac.uk/id/eprint/45403/>

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: <http://nrl.northumbria.ac.uk/policies.html>

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)



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

Ateeq M. Khaliq<sup>1,\*</sup>, Zeyneb Kurt<sup>2,\*</sup>, Miles W. Grunvald<sup>1,\*</sup>, Cihat Erdogan<sup>3</sup>, Sevgi S. Turgut<sup>3</sup>, Tim Rand<sup>4</sup>, Sonal Khare<sup>4</sup>, Jefferey A. Borgia<sup>1</sup>, Dana M. Hayden<sup>1</sup>, Sam G. Pappas<sup>1</sup>, Henry R. Govekar<sup>1</sup>, Anuradha R. Bhama<sup>1</sup>, Ajaypal Singh<sup>1</sup>, Richard A. Jacobson<sup>1</sup>, Audrey E. Kam<sup>1</sup>, Andrew Zloza<sup>1</sup>, Jochen Reiser<sup>1</sup>, Daniel V. Catenacci<sup>5</sup>, Kiran Turaga<sup>5</sup>, Milan Radovich<sup>6</sup>, Sheeno Thyparambil<sup>7</sup>, Mia A. Levy<sup>1</sup>, Janakiraman Subramanian<sup>8</sup>, Timothy M. Kuzel<sup>1</sup>, Anguraj Sadanandam<sup>9</sup>, Arif Hussain<sup>10</sup>, Bassel El-Rayes<sup>11</sup>, Ameen A. Salahudeen<sup>4</sup>✉, Ashiq Masood<sup>1</sup>✉

## Affiliations:

<sup>1</sup>Rush University Medical Center, Chicago, IL, USA.

<sup>2</sup>Northumbria University, Newcastle Upon Tyne, UK.

<sup>3</sup>Isparta University of Applied Sciences, Isparta, Turkey.

<sup>4</sup>Tempus Labs, Inc., Chicago, IL, USA.

<sup>5</sup>The University of Chicago, Chicago, IL, USA.

<sup>6</sup>Indiana University School of Medicine, Indianapolis, IN, USA.

<sup>7</sup>mProbe Inc. Rockville, Maryland, USA

<sup>8</sup>University of Missouri-Kansas City School of Medicine, Kansas City, MO, USA.

<sup>9</sup>Institute of Cancer Research, London, UK.

<sup>10</sup>University of Maryland Marlene and Stewart Greenebaum Comprehensive Cancer Center, Baltimore, MD USA.

<sup>11</sup>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](mailto:ashiq_masood@rush.edu); [ameen@tempus.com](mailto:ameen@tempus.com)

# ABSTRACT

Colorectal cancer (CRC), a disease of high incidence and mortality, exhibits a large degree of inter- and intra-tumoral heterogeneity. The cellular etiology of this heterogeneity is poorly understood. Here, we generated and analyzed a single-cell transcriptome atlas of 49,859 CRC cells from 16 patients, validated with an additional 31,383 cells from an independent CRC patient cohort. We describe subclonal transcriptomic heterogeneity of CRC tumor epithelial cells, as well as discrete stromal populations of cancer-associated fibroblasts (CAFs). Within CRC CAFs, we identify the transcriptional signature of specific subtypes (CAF-S1 and CAF-S4) that significantly stratifies overall survival in more than 1,500 CRC patients with bulk transcriptomic data. We also uncovered two CAF-S1 subpopulations, *ecm-myCAF* and *TGF $\beta$ -myCAF*, known to be associated with primary resistance to immunotherapies. We demonstrate that scRNA analysis of malignant, stromal, and immune cells exhibit a more complex picture than portrayed by bulk transcriptomic-based Consensus Molecular Subtypes (CMS) classification. By demonstrating an abundant degree of heterogeneity amongst these cell types, our work shows that CRC is best represented in a transcriptomic continuum crossing traditional classification systems boundaries. Overall, this CRC cell map provides a framework to re-evaluate CRC tumor biology with implications for clinical trial design and therapeutic development.



## Main

Colorectal cancer (CRC) is the third most commonly diagnosed cancer and a leading cause of cancer-related mortality worldwide<sup>1,2</sup>. Approximately 50% of patients experience disease relapse following curative-intent surgical resection and chemotherapy<sup>3,4</sup>. Despite the high incidence and mortality of advanced CRC, few effective therapies have been approved in the past several decades<sup>5</sup>. One barrier to the development of efficacious therapeutics is the biological heterogeneity of CRC and its variable clinical course<sup>6</sup>. 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<sup>7-12</sup>. 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)<sup>8-12</sup>.

Bulk gene expression analyses by several independent groups have identified distinct CRC subtypes<sup>13-15</sup>. Reflecting both the tumor and TME, an international consortium published the Consensus Molecular Subtypes (CMS), which proposed four distinct subtypes of CRC<sup>15</sup>. Unfortunately, the association between CMS and meaningful therapeutic response to specific agents have shown inconsistent results across studies and CMS lacks a concordance between primary and metastatic CRC tumors, limiting its utility thus far in clinical decision making<sup>16-23</sup>. As a result, an improved CMS classification or an alternative classification system is required 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 heterotropic cell states of tumor epithelia and stromal cells, including a cancer-associated fibroblast

(CAF) population. The CAF population's clinical and prognostic significance became apparent when CAF signatures were applied to large, independent CRC transcriptomic cohorts.

## RESULTS

We profiled sixteen primary colon tissue samples and eight adjacent non-malignant tissues (24 in total) using droplet-based, scRNA-seq. Altogether, we captured and retained 49,589 single cells after performing quality control for downstream analysis (**Fig. 1a, supplementary table 1**). All scRNA-seq data were merged and normalized to identify robust discrete clusters of epithelial cells (*EPCAM*<sup>+</sup>, *KRT8*<sup>+</sup>, and *KRT18*<sup>+</sup>), fibroblasts (*COL1A2*<sup>+</sup>), endothelial cells (*CD31*<sup>+</sup>), T cells (*CD3D*<sup>+</sup>), B cells (*CD79A*<sup>+</sup>), and myeloid cells (*LYZ*<sup>+</sup>) using canonical marker genes. Additionally, each cell type compartment was analyzed separately. Cluster v0.4.1 and manual review of differentially expressed genes in each subcluster were studied to choose the best cluster resolution without cluster destabilization (see methods)<sup>24</sup>. Cell population designation was chosen by specific gene expression, and SingleR was also utilized for unbiased cell type recognition (see methods)<sup>25–28</sup>.

In addition to cancer cells, we identified diverse TME cell phenotypes, including fibroblasts subsets (*CAF-S1* and *CAF-S4*), endothelial cells, CD4<sup>+</sup> subsets (*naïve/memory*, *Th17*, and *Tregs*), CD8<sup>+</sup> subsets (*naïve/memory*, *cytotoxic*, *tissue-resident memory*, and *Mucosa-Associated Invariant (MAIT) cells*), NK cells, innate lymphoid cell (ILC) types, B cell phenotypes (*naïve*, *memory*, *germinal center*, and *plasma cells*), and monocyte lineage phenotypes (*C1DC*<sup>+</sup> *dendritic cells*, *proinflammatory monocytes* [*IL1B*, *IL6*, *S100A8*, and *S100A9*]), and M2 polarized anti-inflammatory [*CD163*, *SEPP1*, *APOE*, and *MAF*]), tumor-associated macrophages (TAMs) (**Fig. 1b-d, Extended Data Figs. 1-3, and Extended Data Tables 1-4**)<sup>25–28</sup>.

For validation, we additionally profiled 31,383 high-quality, single cells from an independent cohort using stringent criteria to corroborate our findings (see methods)<sup>29</sup>. Thus, a total of 81,242 high-quality cells were profiled to produce a single-cell map of 39 colorectal cancer patients. The results of the primary CRC cohort (49,859 single-cells) are available at the Colon Cancer Atlas ([www.websiteinprogress.com](http://www.websiteinprogress.com)).

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

We detected 8,965 tumor and benign epithelial cells (*EPCAM*<sup>+</sup>, *KRT8*<sup>+</sup>, and *KRT18*<sup>+</sup>) and, on reclustering, produced 17 epithelial clusters (designated C1 to C17) (**Fig. 2a and Supplementary Table 2**). Clusters were chiefly influenced by colonic epithelial markers, including those for stemness (*LGR5*, *ASCL2*, *OLFM4*, and *STMN1*), enterocytes (*FABP1* and *CA2*), goblet cells (*ZG16*, *MUC2*, *SPINK4*, and *TFF3*), and enteroendocrine cells (*PYY* and *CHGA*) (**Supplementary Fig. 2a and Supplementary Table 2**). Tumor cells exhibited a high degree of de-differentiated state of plasticity possibly accounting for lasting cancer growth (**Supplementary Fig. 2b**)<sup>30</sup>. Each distinct tumor-derived cluster was mostly patient-specific, reflecting a high degree of inter-patient tumoral cell heterogeneity. In contrast, epithelial populations derived from non-malignant tissue samples from multiple patients clustered together, a pattern observed in previous studies confirming both normal tissue homeostasis and limited sample batch effects (**Fig. 2a**)<sup>31,32</sup>.

We next aimed to identify gene expression programs shared across these clusters using hallmark pathway analysis<sup>33</sup>. A strong overlap was observed for multiple pathways such as activation of inflammatory, epithelial-mesenchymal transformation (EMT), immune response, and metabolic pathways (**Fig. 2b**). Interestingly, high microsatellite instability (MSI-H) and microsatellite stable (MSS) CRC tumors, considered clinically separate entities, demonstrated similar pathway

program activation within the tumor epithelial populations. Some clusters also showed activation of unique pathways such as activation of apical junctions and angiogenesis (C6), hypoxia and fatty acid metabolism (C11) and Notch signaling and DNA repair (C14), among others. However, MSI-H tumors differed from MSS tumors based on immune cell infiltration (**Extended Data Fig. 1**).

Since intratumoral heterogeneity is recognized as a key mechanism contributing to drug resistance, cancer progression, and recurrence, we next focused on dissecting potential transcriptomic states to identify heterogeneity within each patient's tumor<sup>34–37</sup>. We found that each tumor specimen contained 2–10 distinct tumor epithelial clusters (**Fig. 2c**). Gene set variation analysis (GSVA) was performed on cells from individual tumor samples and illustrated the subclonal transcriptomic heterogeneity within each specimen (**Supplementary Fig. 2c**)<sup>38</sup>. Clusters identified in individual pathway analysis demonstrated the up- or down-regulation of crucial metabolic and oncogenic pathways between samples, suggesting wide phenotype variations between cells from the same tumor<sup>39</sup>.

Given the evidence of intratumoral epithelial heterogeneity, we next performed trajectory inference using pseudotime analysis to identify potential alignments or lineage relationships (i.e., right versus left-sided CRC), CMS classification, or MSI status (**Fig. 2d**)<sup>40,41</sup>. This analysis also served as a control for inter-patient genomic heterogeneity and provided an orthogonal strategy to confirm the transcriptomic trends we identified. We detected five molecular states (S1n/t to S5n/t) with malignant and normal epithelial cells intermixed along a joint transcriptional trajectory. This observation is consistent with prior studies demonstrating that CRC cells recapitulate normal colon epithelia's multilineage differentiation process as each transcriptional state's pathway activation in both normal and tumor cells was related to normal

colon epithelial function of nutrient absorption or maintaining colon homeostasis (**Supplementary Fig. 3 and Supplementary Table 3**)<sup>42,43</sup>.

Additionally, tumor cells showed upregulation of embryogenesis (S2t), consistent with previous findings that tumor cells revert to their embryological states in cancer development (**Supplementary Table 4**)<sup>44</sup>. Interestingly, there were no significant associations with anatomic location, CMS classification, or MSI status within our dataset or an independent dataset of 31,383 single cells (**Supplementary Fig. 4**)<sup>29</sup>. Hence, in our analysis, CRC development mainly represents a hijacking of the normal epithelial differentiation program, coupled with the acquisition of embryogenic pathways (**Supplementary Fig. 3**)<sup>45</sup>.

**CRC-associated fibroblasts in the tumor microenvironment exhibit diverse phenotypes, and specific subtypes are associated with poor prognosis.**

We next focused on CRC TME subpopulations. High-quality 819 fibroblasts were re-clustered into eight clusters, and then phenotypically classified into two major subtypes to assess for further CAF heterogeneity. These phenotypic subtypes were found to be immunomodulatory CAF-S1 (*PDGFRA*<sup>+</sup> and *PDPN*<sup>+</sup>) and contractile CAF-S4 (*RGS5*<sup>+</sup> and *MCAM*<sup>+</sup>) (**Fig. 3 and Supplementary Table 4**)<sup>46</sup>. This fibroblast cluster dichotomy was also observed in the independent CRC patient scRNA-seq dataset of 31,383 cells (**Supplementary Fig. 5**)<sup>29</sup>.

The CAF-S1 and CAF-S4 subtypes showed striking resemblances to the mCAF (extracellular matrix) and vCAF (vascular) fibroblast subtypes, respectively, as previously described in a mouse breast cancer model<sup>47</sup>. Most clusters were found in multiple patients, albeit in varying proportions, signifying shared patterns in CAF transcriptomic programs between patients. Fibroblasts derived from MSI-H tumors were distributed similarly throughout these clusters.

CAF-S1 exhibited high chemokine expressions such as *CXCL1*, *CXCL2*, *CXCL12*, *CXCL14*, and immunomodulatory molecules including *TNFRSF12A* (**Supplementary Fig. 6**). Additionally, CAF-S1 expressed extracellular matrix genes including matrix-modifying enzymes (*LOXL1* and *LOX*)<sup>47</sup>. To determine this population's functional significance, we compared the CAF-S1 population transcriptomes to those described recently in breast cancer, lung cancer, and head and neck cancer<sup>48</sup>. We recovered five CAF subtypes, within the CAF-S1 population, including *ecm-myCAF* (extracellular; *GJB+*), *IL-iCAF* (growth factor, *TNF* and interleukin pathway; *SCARA5+*), *detox-iCAF* (detoxification and inflammation; *ADH1B+*), *wound-myCAF* (collagen fibrils and wound healing; *SEMA3C*), and *TGFβ-myCAF* (*TGF-β* signaling and matrisome; *CST1+*, *TGFβ1+*), which were previously divided into two major subtypes: iCAF and *myCAF* (**Fig. 3b**). Among these five subtypes, *ecm-myCAF* and *TGFβ-myCAF* are known to correlate with immunosuppressive environments and are enriched in tumors with high regulatory T lymphocytes (Tregs) and depleted CD8<sup>+</sup> lymphocytes. Additionally, these subtypes are associated with primary immunotherapy resistance in melanoma and lung cancer<sup>48</sup>.

The CAF-S4 population expressed pericyte markers (*RGS5+*, *CSPG4+*, and *PDGFRA+*), *CD248* (*endosialin*), and *EPAS1* (*HIF2-α*), that this particular CAF subtype is vessel-associated, with hypoxia potentially contributing to invasion and metastasis, as has been shown in another study<sup>47</sup>. CAF-S4 clustered into the immature phenotype (*RGS5+*, *PDGFRB+*, and *CD36+*) and the differentiated myogenic subtype (*TAGLN+* and *MYH11+*)(**Fig. 3c and Supplementary Table 4**)<sup>49</sup>.

Given the correlation between CMS4 and fibroblast infiltration, we next sought to test the existence of CAF-S1 and CAF-S4 signatures in bulk transcriptomic data and their association with clinical outcomes<sup>15</sup>. To this end, we interrogated and carried out a meta-analysis of eight colorectal cancer transcriptomic datasets comprising 1,584 samples. We detected a strong and

positive correlation between specific gene expressions characterizing each CAF subtype in CRC. We also confirmed the presence of CAF-S1 and CAF-S4 signatures in pancreatic adenocarcinoma (n=118) and non-small cell lung cancer (NSCLC, n=80) cohorts (**Fig. 4**) (see methods for datasets). The gene signatures were specific to each CAF-S1 and CAF-S4, thus confirming their existence in TME of CRC and other tumor types.

We found high CAF-S1 and CAF-S4 signatures associated with poor median overall-survival (HR>1, p<0.05), irrespective of CMS subtypes in three independent CRC datasets (**Supplementary Fig 7**). Additionally, CAF signatures stratified the CMS4 subtype into high- and low-risk overall survival in all datasets, thus identifying additional heterogeneity and providing prognostication in this aggressive patient subgroup (**Fig. 4b-d**). Here, using scRNA-seq, we show for the first time that high CAF infiltration in CRC is associated with poor prognosis across all molecular subtypes, and which further stratifies the CMS4 subgroup into high and low-risk clinical phenotypes in CRC cohorts.

# **Single-cell RNA sequencing reveals heterogeneity beyond Consensus Molecular Subtypes in colorectal cancers and offers therapeutic opportunities.**

The lack of association between tumor epithelia and CMS classification, as well as the survival differences between high- and low-risk CAF signatures across CRC molecular subtypes suggest CRCs are much more heterogeneous than the traditional classification systems have indicated (e.g. those systems defined by somatic alterations, epigenomic features, and bulk gene expression data)<sup>13,14,50,51</sup>.

Among these the widely adopted CMS classification, which reflects both the malignant cell phenotypes and the TME, classified CRC into CMS1 (MSI immune), CMS2 (canonical), CMS3



(metabolic), and CMS4 (mesenchymal) subtypes based on bulk transcriptomic signatures<sup>15</sup>. To test our hypothesis, we estimated every cell type fraction using single-cell data from eight pooled datasets (>1,500 samples) with a machine-learning algorithm, CIBERSORTx<sup>52</sup>. When we compared epithelial, immune, and stromal cell populations among the CMS subtypes, we did not detect a distinct pattern of tumor, immune, or stromal cell signatures across the different CMS subtypes. Each CMS subtype was enriched in these cell types in varying proportions but without a clear distinction between the four subtypes, suggesting a lack of clear separation among the CMS subsets at the single-cell resolution (**Fig. 5a-b**). Upon analysis of four independent bulk RNA datasets, there was significant discordance in terms of cell phenotype enrichment with respect to each CMS subtype across the datasets except CMS4 which had predominant stromal enrichment (**Supplementary Figs. 8-11**)<sup>53</sup>. These discordant results could be due to intra-patient CMS heterogeneity, intratumoral variation in tumor purity, stromal and immune cell infiltration, and CMS's inability to address tumor/TME-to-tumor/TME variability, among others<sup>54</sup>. Thus, novel approaches that consider these factors are required to stratify patients for optimal biomarker and therapeutic development.

Based on the above findings, we postulate that CRC is more accurately represented in a transcriptomic continuum previously proposed by Ma et al.<sup>55</sup>. The authors analyzed bulk transcriptomic data using a novel computational framework in which denovo, unsupervised clustering methods (k-medoid, non-negative-matrix factorization, and consensus clustering) demonstrated the existence of CRC in a transcriptomic continuum<sup>56-58</sup>. They further carried out principal component analysis and robustly validated two principal components, PC Cluster Subtype Scores 1 and 2 (PCSS1 and PCSS2, respectively). We reasoned that using single-cell data could deepen our understanding of how each cell phenotype contributes to the CRC tumor microenvironment using continuous scores that inform CRC diversity beyond binning CRC into traditional classifications.



We evaluated every cell fraction (epithelial, stromal, and immune components) from our data with the Ma et al. algorithm on eight pooled bulk transcriptomic datasets, focusing on PCSS1 and PCSS2, since these were validated in the original study (**Fig. 5c-d**). On projecting single-cell expression profiles on quadrants corresponding to each of the four CMS, we noted a lack of separation of all cell phenotypes between CMS subtypes suggesting that CRC exists in a transcriptomic continuum<sup>55</sup>.

To test continuous scores reproducibility, we analyzed four bulk transcriptomic datasets separately; we found that transcriptional shifts were reproducible across datasets for all major cell types (**Supplementary Figs. 8-11**). The continuous scores showed no reliability in classifying CRC into immune–stromal rich (CMS1/CMS4) or immune-stromal desert (CMS2/CMS3) subtypes as proposed previously<sup>20</sup>. Thus, confirming continuous scores rather than discrete subtypes may improve classifying CRC tumors and may explain tumor-to-tumor variability, tumor/TME-to-tumor/TME and TME-to-TME variabilities<sup>55</sup>. Of note, CAF-S1 exhibited high PCSS1 and PCSS2 scores across independent datasets, correlating with the CMS4 subtype. Thus, our analysis identified CAF-S1 as a cell of origin for biological heterogeneity in CMS4 subtypes associated with poor prognosis (**Supplementary table 5**).

## DISCUSSION

In the present study, we evaluated the CMS classification of CRC that have been developed by bulk RNA-seq through single-cell resolution transcriptomic analysis. We find that stromal cells engender a more significant contribution to biological heterogeneity. Although previous studies employing bulk transcriptomics have demonstrated that the degree of stromal infiltration is associated with prognosis and a small scRNAseq study utilizing 26 fibroblasts demonstrated poor survival among CAF-enriched CRC tumors especially the CMS4 subtype. Here, using much larger

sample set of 1,182 high-quality fibroblasts, we identify the CAF-S1 subtype to be the cell of origin associated with poor prognosis across all CRC CMS subtypes and not just CMS4<sup>59</sup>. Further, we developed a novel signature of CAF infiltration and demonstrate that CMS4 can be stratified into risk groups associated with good or poor median overall survival. These findings are significant since the CMS4 subtype, is primarily stromal-driven and is enriched in more than 40% of metastatic CRC samples from patients with worse outcomes<sup>21</sup>. We also identified CAF-S1 subtypes associated with certain biological functions in other cancers, including the *ecm-myCAF* and *TGFβ-myCAF* subtypes (responsible for immunotherapy resistance in NSCLC and melanoma), and CAF-S4s known to play a role in inducing cancer cell invasion<sup>48,49</sup>. Thus, targeting CAFs to remodel the tumor microenvironment may lead to improved and much-needed therapeutic development for metastatic CRC patients<sup>21,22</sup>.

Targeting of CAFs in solid tumors is being explored in multiple clinical trials with variable results<sup>60–62</sup>. Such studies likely failed to address CAF heterogeneity and their complex interactions with the other cells of TME. Our study suggests that CRC may be intricately entwined with the stroma, and therefore may be amenable to stromal targeted combinatoric approaches, including monoclonal antibodies that abrogate CAF-S1 function. In future studies, the treatment of CRC patients should involve stroma targeted therapies and take the above aspects into consideration<sup>60,63</sup>. The scRNA-seq or bulk-RNA-seq signatures corresponding to CAF-S1 and CAF-S4 may serve as suitable biomarkers for tumors that are reliant on this axis. Immunotherapy responses in MSS CRC, which comprise almost 95% of metastatic CRC, are lacking; CAF subpopulations within the TME may be suppressing immune responses in these tumors<sup>64</sup>. Based on our analysis, we speculate that targeting CAF-derived chemokines and cytokines via biospecific antibodies, vaccines, or even cell-based therapies, may enhance current checkpoint blockade strategies<sup>60</sup>. Functional validation and clinical studies will be required to confirm the clinical utility of targeting these CAF populations in CRC.

More importantly, our study's single-cell resolution enables us to investigate whether tumor cell transcriptomes, and by extension, biological phenotypes, are the primary determinant of CMS classification. Based on our findings, it appears that bulk analysis may have been confounded by varying tumor microenvironment population enrichment, and that tumor cells within each patient do not segregate into static phenotypes but rather exhibit considerable plasticity. In contrast, the single-cell analysis uncovered the complex and mixed cellular-phenotypes among each cell specific subpopulation, which projected in a transcriptomic continuum across CMS subtypes. These findings were further supported by scRNA-seq CMS classification analysis that assigned each CRC sample to multiple CMS subtypes thereby suggesting CMS heterogeneity in each CRC tumor<sup>20,65</sup>. These findings may also explain why CMS-defined populations of tumors have not been readily observed in transcriptomic data from independent CRC cohorts<sup>21,53</sup>. Our data indicate that attempts to divide CRC phenotypes into the current discrete subtypes may undermine optimal patient stratification in the clinical trial setting. Intriguingly, by applying two independent algorithms, we demonstrate that CRC tumors and their ecosystems exist in a transcriptomic continuum and not only show tumor-to-tumor variability (as proposed by Ma et al.) but also demonstrate tumor/TME-to-tumor/TME transcriptional variability at the single-cell resolution<sup>55</sup>. The continuous scores are reproducible across transcriptomic datasets, thus allowing robust identification of patient subtypes. This may help to optimize CRC treatment in future studies.

In conclusion, our study lends strong support to the tumor biology models proposed by Ma et al. (and other groups) and represents a conceptual shift in our understanding of CRC pathogenesis, clinical management, and therapeutic development. Future studies will need to consider tumor-TME to tumor-TME heterogeneity which will be critical for optimizing biomarkers and treatment strategies for CRC.

350

## 351 **ACKNOWLEDGEMENTS**

352 This study was supported by the startup fund provided to A.M. by the Rush University Medical  
353 Center; the OCM grant to A.M. by Rush University Cancer Center. Part of A.H.'s time was  
354 supported by a Merit Review Award (I01 BX000545) from the Medical Research Service,  
355 Department of Veterans Affairs. We would also like to thank Dr. Levi Waldron for sharing code  
356 from his publication entitled, "Continuity of transcriptomes among colorectal cancer subtypes  
357 based on meta-analysis." Above all we want to thank our patients who participated in this study  
358 and their families.

359

## 360 **AUTHOR CONTRIBUTIONS**

361 A.M. devised, supervised the study, and wrote the manuscript. A.M.K. performed data analyses,  
362 wrote the manuscript, and created figures. Z.K. performed data analyses and wrote the  
363 manuscript. M.W.G. aided in analysis, wrote the manuscript, and generated figures. A.S.  
364 supervised study and wrote manuscript. C.E. and S.S.T. helped with bulk transcriptomic analysis.  
365 D.M.H, H.R.G., A.R.B helped with sample collection. All other authors contributed substantially  
366 to data interpretation, and manuscript editing. All authors read and approved this manuscript.

367

## 368 **CODE AVAILABILITY**

369 The code generated and utilized in the completion of this publication will be available in a Github  
370 repository specific to this project.

371

## 372 **DATA AVAILABILITY**

373 Sequencing data deposition is currently in progress. Ten bulk transcriptomic datasets were  
374 accessed from the Gene Expression Omnibus (GEO) database  
375 (<https://www.ncbi.nlm.nih.gov/geo/>).

376

377 **COMPETING INTERESTS**

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

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

380 A.S. holds patent PCT/IB2013/060416, '*Colorectal cancer classification with differential prognosis*

381 *and personalized therapeutic responses*' and patent number 2011213.2 '*Prognostic and*

382 *Treatment Response Predictive Method.*'

## METHODS

**Patient and tissue sample collection.** Patients with resectable untreated CRC who underwent curative colon resection at Rush University Medical Center (Chicago, IL, USA) were included in this Institutional Review Board (IRB)-approved study. CRC specimens from 16 patients including nine Caucasian, six African American and one Asian patient with corresponding 8 adjacent normal tissue samples were processed immediately after collection at Rush University Medical Center Biorepository and sent for scRNA-seq. Thus, our scRNA-seq atlas represent diverse patient population. The study was conducted in accordance with ethical standards and all patients provided written informed consent.

**Droplet based scRNA-seq - 10x library preparation and sequencing.** Single-cell RNA sequencing (scRNA-seq) was performed using 10X Genomics Single Cell 5' Platform. Tumors and non-malignant samples were enzymatically dissociated (*Miltenyi*), filtered through a 70-micron cell strainer, pelleted after centrifugation at 300 xg and resuspended in DAPI-FACS buffer (PBS, 0.04% BSA). Samples were sorted and viable singlets were gated on the basis of scatter properties and DAPI exclusion. Approximately 3000 cells were pelleted and resuspended in PBS, and cells underwent single cell droplet-based capture on 10X Chromium instruments according to the 10X Genomics Single Cell 5' Platform protocol. Transcriptome libraries post-fragmentation, end-repair, and A-tailing double-sided size selection, and subsequent adaptor ligation also followed the manufacturer's protocol. Illumina *NextSeq 550* was used for library sequencing and data were mapped and counted using Cellranger-v3.1.0 (*GRCh38/hg38*).

**scRNA-seq data quality control, gene-expression quantification, dimensionality reduction, and identification of cell clusters.** *Cell Ranger* was utilized to process the raw gene expression matrices per samples and all samples from multiple patients were combined in R package (v3.6.3

2020-02-29] -- "*Holding the Windsock*"). Seurat package (v3.2.2) was used in this integrative multimodal analysis<sup>66</sup>. Genes detected in fewer than three cells and cells expressing less than 200 detected genes were filtered out and excluded from analysis. In addition, cells expressing > 25% mitochondria were removed. Cell cycle scoring was performed, (for the S phase and the G2M phase) and the predicted cell cycle phases were calculated. Doublet detection and any higher-order multiplets that were not dissociated during sample preparation were removed via the *DoubletFinder* (v2.0.2) package using default settings<sup>67</sup>. Following quality control one non-malignant colon sample (B-cac13) was discarded due to poor data quality. Finally, 49,859 cells remained and were utilized for downstream analysis.

We adopted the general protocol described in Stuart et al. (2019) to group single cells into different cell subsets<sup>66</sup>. We employed the following steps: clustering the cells within each compartment (including the selection of variable genes for each dataset based on a variance stabilizing transformation [VST]), canonical correlation analysis (CCA) to remove batch effects among the samples, reduction of dimensionality, and projection of cells onto graphs<sup>68,69</sup>. Principal component analysis (PCA) was carried out on the scaled data of highly variable genes<sup>70</sup>. The first 30 principal components (PCs) were used to cluster the cells and to perform a subtype analysis by nonlinear dimensionality reduction (t-SNE) and to construct Uniform Manifold Approximation and Projection (UMAP) for cell embeddings<sup>71,72</sup>. We identified cell clusters under the optimal resolution by a shared nearest neighbor (SNN) modularity optimization-based clustering method. We implemented the *FindClusters* function of the Seurat package, which first calculated *k-nearest neighbors* and constructed the SNN graph. We implemented the original *Louvain algorithm* (algorithm = 1) for modularity optimization. Additionally, we utilized Clustree (v0.4.3) and manual review for identifying the best clustering resolution<sup>24</sup>.

**Major cell type detection and data visualization.** To identify all major cell types, we evaluated differentially expressed markers in each identity cell group by comparing them to other clusters

using the Seurat *FindAllMarkers* function. We used positively expressed genes with an average expression of  $\geq 2$ -fold higher in that subcluster than the average expression in the rest of the other subclusters. We used known marker genes, which have the highest fold expression in that cluster with respect to the other clusters. We also utilized SingleR ((v0.99.10, R Package), which leverage large transcriptomic datasets of well-annotated cell types and manual annotation for cell-type identification<sup>31,73–75</sup>. Depending on the presence of known marker genes the clusters were grouped as: epithelial cells (*EPCAM*, *KRT8*, and *KRT18*), fibroblasts (*COL1A1*, *DCN*, *COL1A2*, and *C1R*), endothelial cells (CD31+), myeloid cells (*LYZ*, *MARCO*, *CD68*, and *FCGR3A*), CD4 T cells (*CD4*), CD8 T cells (*CD8A* and *CD8B*), and B cells (*MZB1*),<sup>31,47,73,76–80</sup>. The cells were eventually assembled into DGE matrices within each compartment, containing all six cell types.

**Major-cell type subclustering and data visualization.** Each major cell type, including epithelial cells, endothelial cells, T cells, B cells, myeloid cells, and fibroblasts was reclustered and reanalyzed to study each compartment at a higher resolution to detect granular cellular heterogeneity in CRC. Clustree (v0.4.3) and manual review were utilized for optimal cluster detection. For cell annotation of each cell type, we utilized published literature gene expression signatures and manual review of differential genes among clusters. Additionally, we again utilized SingleR (v0.99.10, R Package) for unbiased cell annotation. Interestingly, reclustering of major compartments individually also detected clusters expressing hybrid markers as well as cell clusters expressing markers from distinct lineages (such as T cell clusters expressing B cells); these were removed and excluded for further analysis. We utilized UMAP for visualization purposes. For validation, we analyzed 65,362 cells from 23 patients and applied the same quality control metrics as outlined above, retaining 31,383 high-quality single cells for further analysis<sup>29</sup>. These high-quality cells were analyzed utilizing the same pipelines and parameters as that for our primary cohort (**Supplementary Figs. 4-5 and 12-13**).



The InferCNV (v1.2.1) package was used with default parameters to identify the evidence for somatic large-scale chromosomal copy number alteration in epithelial cells (*EPCAM*<sup>+</sup>, *KRT8*<sup>+</sup>, *KRT18*<sup>+</sup>)<sup>81</sup>. Non-malignant epithelial cells were used as the control group.

**Trajectory analysis.** We used Monocle v.2 (v2.14.0), a reverse graph embedding method to reconstruct single-cell trajectories in tumor and non-malignant epithelium<sup>82</sup>. In brief, we used UMI count matrices and the *negbinomial.size()* parameter to create a *CellDataSet* object in the default setting. We grouped projected cells on UMAP in default settings for visualization of monocle results. We defined the cumulative duration of the trajectory to show the average amount of transcriptional transition that a cell undergoes as it passes from the starting state to the end state. The cells were also ordered in pseudotime to explain the transition of cells from one state to another.

**Pathway- Gene set variation analysis (GSVA).** Pathway analysis was performed on the 50 hallmark gene sets downloaded from *Molecular Signatures Database* (v7.2). We used GSVA (v1.34.0), a non-parametric, unsupervised method to estimate the gene set variations and evaluation of pathway enrichment, and pathway scores were calculated for each cell using standard settings<sup>33,38</sup>.

**DNA and bulk RNA library construction.** DNA and bulk RNA sequencing was performed as previously described<sup>83</sup>. One hundred nanograms of DNA from each tumor was mechanically sheared to an average size of 200 bp. Using the *KAPA Hyper Prep Pack*, DNA libraries were packed, hybridized into the *xT probe* package, and amplified with the *KAPA HiFi HotStart ReadyMix*. For uniformity, each sample needed to have 95% of all targeted base pairs sequenced to a minimum depth of 300x. One hundred nanograms of RNA per tumor sample was heat fragmented to a mean size of 200 base pairs in the presence of magnesium. Using random

primers, the RNA was used for first-strand cDNA synthesis, followed by second-strand synthesis and A-tailing, adapter ligation, bead-based cleanup, and amplification of the library. After library planning, the *IDT xGEN Exome Test Panel* was hybridized with samples. Streptavidin-coated beads and target recovery were carried out, accompanied by amplification using the *KAPA HiFi* library amplification package. The RNA libraries were sequenced on an *Illumina HiSeq 4000* using patterned flow cell technology to achieve at least 50 million reads.

**Detection of somatic variation on DNA sequencing data.** The tumor and normal FASTQ files were paired. For quality management measurement, FASTQ files were evaluated using FASTQC and matched with Novoalign (Novocraft, Inc.)<sup>83,84</sup>. SAM files were generated and converted to BAM files. The BAM files were sorted, and duplicates were marked. Single nucleotide variations (SNVs) were called after alignment and sorting. For discovery of copy number alterations, the de-duplicated BAM files and the VCF generated from the variant calling pipeline were processed to compute read depth and variance of heterozygous germline SNVs between the tumor sample and normal sample. Binary circular segmentation was introduced and segments with strongly differential log<sub>2</sub> ratios between the tumor and its comparator were chosen. From a combination of differential coverage in segmented regions and estimation of stromal admixture provided by analysis of heterozygous germline SNVs, an estimated integer copy number was determined

**Microsatellite instability status.** Probes for 43 microsatellite regions were developed using *Tempus xT* assay<sup>83</sup>. Tumors were categorized into three groups by the MSI classification algorithm as described by Tempus: microsatellite instability-high (MSI-H), microsatellite stable (MSS) or microsatellite equivocal (MSE). MSI screening for paired tumor-normal patients used reads mapped to the microsatellite loci with at least 5 bps flanking the microsatellite. The sample was graded as MSI-H if there was a >70% chance of MSI-H classification. If the likelihood of MSI-H status was 30-70%, the test findings were too ambiguous to interpret and those samples were

listed as MSE. If there was a <30% chance of MSI-H status, the sample was called MSS. Additionally, IHC results were used to classify tumors into MSS or MSI molecular subtypes. Both of these modalities were concordant and produced the same results.

**Bulk RNA-seq and microarray analysis.** We downloaded and pooled eight colorectal gene expression datasets (GSE13067<sup>85</sup>, GSE13294<sup>85</sup>, GSE14333<sup>86</sup>, GSE17536<sup>87</sup>, GSE20916<sup>88</sup>, GSE33113<sup>89</sup>, GSE35896<sup>90</sup>, and GSE39582<sup>14</sup>), a pancreatic cancer dataset (GSE62165<sup>91</sup>) and a non-small cell lung cancer dataset (GSE33532<sup>92</sup>) to validate our findings from the single cell compartments by deconvoluting the bulk gene expression profiles into pseudo single-cell resolutions. We used Affy (v1.64.0) for the data analysis and for exploration of Affymetrix oligonucleotide array probe level data<sup>93</sup>. Batch correction was carried out using the removeBatchEffect (v3.42.2) function of the LIMMA program and CMScaller for the CMS classification (see below)<sup>94</sup>. Three datasets (GSE17536<sup>87</sup>, GSE33113<sup>89</sup>, and GSE39582<sup>14</sup>) were utilized for clinical outcome analysis<sup>94,95</sup>.

**Correlation patterns in bulk gene expressions for CAF compartments.** To identify the top correlated CAF-marker genes within the combined eight CRC datasets, four bulk CRC gene expression sets individually, pancreas cancer and lung cancer datasets. We first transformed the bulk gene expression sets with log<sub>2</sub> transformation. Next, marker genes with an average log<sub>2</sub> FC ≥ 0.5 and p < 0.05 obtained from the single cell data of CAF-S1 and CAF-S4 compartments were separately intersected with the bulk gene expression sets. Genes with an average Spearman correlation score greater than 0.8 were kept as the CAF signatures within the bulk gene expression. Heatmaps illustrating the correlation patterns within and between the CAF compartments were prepared with the heatmap.2 function from ggplot package (v3.1.1) utilizing the Pearson correlation coefficient. Heatmaps illustrating the correlation patterns within and

between the CAF compartments were prepared using the ggplot package (v3.1.1) utilizing the Pearson correlation coefficient<sup>96</sup>.

The Cox proportional hazard regression model was used to examine the significance of 20 cell types from scRNA-seq in bulk expression data. Each cell type's marker genes with an average logFC>1 and adjusted P<0.05 were intersected with the bulk expression datasets separately. We only kept the marker genes with a high correlation with each other in bulk, which provides an average correlation score of > 0.8. The average bulk expression of each cell type's remaining marker genes was calculated and used in the hazard regression model as the representative of this cell type. For analysis of relationships with patient outcome, univariate models were calculated using Cox proportional hazard regression (coxph function from survival R package)<sup>97</sup>.

#### ***Deconvoluting public bulk gene expression profiles into pseudo single-cell expressions.***

We used CIBERSORTx v1 to estimate composition of various cell populations in pooled eight microarray datasets<sup>52</sup>. Signature gene matrices were created using the expression profiles of 49,859 cells as the reference single cell profile. We ran the 'hires' module with default parameters except for the 'rmbatchBmode,' and the bulk-mode batch correction argument was set to true. After the deconvolution process, we normalized the gene expressions according to the cell fractions in each sample and calculated each gene's Z-transformed expression values. The average normalized expression of each cell type across all samples was plotted with the heatmap.3 R function of the GMD package (v0.3.3)<sup>98</sup>. A signature matrix highlighting marker genes of the different cell types was prepared with a heatmap.2 R function of ggplot (v3.1.1). We also applied the same parameters to deconvolute GSE14333<sup>86</sup>, GSE17536<sup>87</sup>, GSE33113<sup>89</sup>, and GSE39582<sup>14</sup> datasets individually.

**Consensus molecular subtyping of colorectal cancer (CMS Classification).** We used R package CMScaller(v0.9.2), a nearest template prediction (NTP) algorithm, for the classification of gene expression datasets<sup>95</sup>. We set the permutation number to 1000 to predict the CMS classes of the samples in the GEO datasets with a p-value < 0.05. We ran CMScaller with default parameters.

**Continuous subtype discovery using scRNA-seq analysis.** Bulk mRNA expression profiles of the combined and batch adjusted eight GEO datasets (GSE13067<sup>85</sup>, GSE13294<sup>85</sup>, GSE14333<sup>86</sup>, GSE17536<sup>87</sup>, GSE20916<sup>88</sup>, GSE33113<sup>89</sup>, GSE35896<sup>90</sup>, and GSE39582<sup>14</sup>), composed of 1584 samples in total, were deconvoluted into the pseudo single-cell expression profiles via CIBERSORTx utilizing the expression data consisting of 20 different cell types from our scRNA-seq dataset<sup>52</sup>. We transformed the deconvoluted expression matrix with log2 transformation. The principal components cluster subtype scores (PCSSs) of the CMS subtypes among the 1584 samples, were determined separately for each cell type using an algorithm published by Ma et al<sup>55</sup>. To obtain the PCSSs, the average loading vectors were used. The results obtained for 20 cell types were projected on the first two PCSSs (PCSS1 and PCSS2) as they were validated by Ma et al. in their analysis using 18 datasets. We also analyzed four datasets (GSE14333<sup>86</sup>, GSE17536<sup>87</sup>, GSE33113<sup>89</sup>, and GSE39582<sup>14</sup>) to independently confirm reproducibility of continuous scores.

**Statistics and reproducibility.** All statistical analyses and graphs were created in R (v3.6.3) and using a Python-based computational analysis tool. Schematic representations were made using the Inkscape (<https://inkscape.org/>) software. Dim plots, bar plots and box plots were generated using the dittoSeq (v1.1.7) package with default parameters<sup>99</sup>. Violin plots were generated using the patchwork (v1.1.0) package and ggplot2 (v3.3.2) package in R with default parameters. Heatmaps were generated using Morpheus.R with default parameters<sup>100,101</sup>. ANOVA and pair-

wised t-tests for the CMS classes across the deconvoluted expression profiles were performed in R using the ggpubr R (v0.4.0) package<sup>102</sup>. The Box and Whisker plots were generated using the boxplot function of the R base package at default parameters. The mean of the log<sub>2</sub> transformed deconvoluted expression value of the samples in each CMS group was demonstrated with a horizontal straight line within each box. The length of a boxplot corresponds to the interquartile range (IQR), which is defined as the range between the first and third quartiles (Q1 and Q3), whereas the whiskers are the upper and lower extreme values of the data (either data's extremum values, or the Q3+1.5\*IQR and Q1-1.5\*IQR values, whichever was less extreme).

**Survival analysis.** Survival curves were obtained according to the Kaplan-Meier method survfit (v3.2-7), and differences between survival distributions were assessed by Log-rank test. The patients were divided into two groups (high/poor and low/good risk) according to their median expression values (survminer (v0.4.8)). The surv\_cutpoint function uses the maximally selected rank statistics and implements standard methods for the approximation of the null distribution of maximally selected rank statistics (maxstat (v0.7-25)).

The proportional hazard assumption was tested to examine the fit of the model for survival of the samples in four GEO datasets (GSE14333<sup>86</sup>, GSE17536<sup>87</sup>, GSE33113<sup>89</sup>, and GSE39582<sup>14</sup>) with respect to the deconvoluted bulk mRNA expressions. For analysis of the relationships with patient outcome, multivariate models were calculated using the Cox proportional hazard regression (coxph survival R package)<sup>97</sup>.



# REFERENCES

1. Cancer of the Colon and Rectum - Cancer Stat Facts. *SEER* <https://seer.cancer.gov/statfacts/html/colorect.html>.
2. Arnold, M. *et al.* Global patterns and trends in colorectal cancer incidence and mortality. *Gut* **66**, 683–691 (2017).
3. Osterman, E. *et al.* Recurrence risk after radical colorectal cancer surgery—less than before, but how high is it? *Cancers* **12**, 3308 (2020).
4. Molinari, C. *et al.* Heterogeneity in colorectal cancer: a challenge for personalized medicine? *Int J Mol Sci* **19**, (2018).
5. Xie, Y.-H., Chen, Y.-X. & Fang, J.-Y. Comprehensive review of targeted therapy for colorectal cancer. *Signal Transduction and Targeted Therapy* **5**, 1–30 (2020).
6. Budinska, E. *et al.* Gene expression patterns unveil a new level of molecular heterogeneity in colorectal cancer. *J Pathol* **231**, 63–76 (2013).
7. Deschoolmeester, V. *et al.* Tumor infiltrating lymphocytes: an intriguing player in the survival of colorectal cancer patients. *BMC Immunol* **11**, 19 (2010).
8. Lee, W.-S., Park, S., Lee, W. Y., Yun, S. H. & Chun, H.-K. Clinical impact of tumor-infiltrating lymphocytes for survival in stage II colon cancer. *Cancer* **116**, 5188–5199 (2010).
9. Perez, E. A. *et al.* Association of stromal tumor-infiltrating lymphocytes with recurrence-free Survival in the N9831 adjuvant trial in patients with early-stage HER2-positive breast cancer. *JAMA Oncology* **2**, 56–64 (2016).
10. Nearchou, I. P. *et al.* Spatial immune profiling of the colorectal tumor microenvironment predicts good outcome in stage II patients. *npj Digital Medicine* **3**, 1–10 (2020).
11. Pagès, F. *et al.* International validation of the consensus Immunoscore for the classification of colon cancer: a prognostic and accuracy study. *The Lancet* **391**, 2128–2139 (2018).
12. The Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon and rectal cancer. *Nature* **487**, 330–337 (2012).
13. Sadanandam, A. *et al.* A colorectal cancer classification system that associates cellular phenotype and responses to therapy. *Nature Medicine* **19**, 619–625 (2013).
14. Marisa, L. *et al.* Gene Expression Classification of Colon Cancer into Molecular Subtypes: Characterization, Validation, and Prognostic Value. *PLOS Medicine* **10**, e1001453 (2013).
15. Guinney, J. *et al.* The consensus molecular subtypes of colorectal cancer. *Nature Medicine* **21**, 1350–1356 (2015).
16. Berg, I. van den *et al.* Improving clinical management of colon cancer through CONNECTION, a nation-wide colon cancer registry and stratification effort (CONNECTION II trial): rationale and protocol of a single arm intervention study. *BMC Cancer* **20**, 1–8 (2020).
17. Lenz, H.-J. *et al.* Impact of Consensus Molecular Subtype on Survival in Patients With Metastatic Colorectal Cancer: Results From CALGB/SWOG 80405 (Alliance). *J Clin Oncol* **37**, 1876–1885 (2019).
18. Stintzing, S. *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* **30**, 1796–1803 (2019).
19. Mooi, J. K. *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* **29**, 2240–2246 (2018).
20. Svein, A., Cremolini, C. & Dienstmann, R. Predictive modeling in colorectal cancer: time to move beyond consensus molecular subtypes. *Annals of Oncology* **30**, 1682–1685 (2019).
21. 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* **30**, 520–527 (2019).

22. Khambata-Ford, S. *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* **25**, 3230–3237 (2007).
23. Aderka, D., Stintzing, S. & Heinemann, V. Explaining the unexplainable: discrepancies in results from the CALGB/SWOG 80405 and FIRE-3 studies. *The Lancet Oncology* **20**, e274–e283 (2019).
24. Zappia, L. & Oshlack, A. Clustering trees: a visualization for evaluating clusterings at multiple resolutions. *Gigascience* **7**, (2018).
25. Zhang, L. *et al.* Single-cell analyses inform mechanisms of myeloid-targeted therapies in colon cancer. *Cell* **181**, 442–459.e29 (2020).
26. Zhang, L. *et al.* Lineage tracking reveals dynamic relationships of T cells in colorectal cancer. *Nature* **564**, 268–272 (2018).
27. Corridoni, D. *et al.* Single-cell atlas of colonic CD8 + T cells in ulcerative colitis. *Nature Medicine* **26**, 1480–1490 (2020).
28. Oh, D. Y. *et al.* Intratumoral CD4+ T Cells Mediate Anti-tumor Cytotoxicity in Human Bladder Cancer. *Cell* **181**, 1612–1625.e13 (2020).
29. Lee, H.-O. *et al.* Lineage-dependent gene expression programs influence the immune landscape of colorectal cancer. *Nature Genetics* **52**, 594–603 (2020).
30. Mills, J. C. & Sansom, O. J. Reserve stem cells: Differentiated cells reprogram to fuel repair, metaplasia, and neoplasia in the adult gastrointestinal tract. *Sci Signal* **8**, re8 (2015).
31. Lambrechts, D. *et al.* Phenotype molding of stromal cells in the lung tumor microenvironment. *Nat Med* **24**, 1277–1289 (2018).
32. Izar, B. *et al.* A single-cell landscape of high-grade serous ovarian cancer. *Nature Medicine* **26**, 1271–1279 (2020).
33. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* **1**, 417–425 (2015).
34. Jamal-Hanjani, M. *et al.* Tracking the evolution of non-small-cell lung cancer. *N Engl J Med* **376**, 2109–2121 (2017).
35. Malikic, S., Jahn, K., Kuipers, J., Sahinalp, S. C. & Beerenwinkel, N. Integrative inference of subclonal tumour evolution from single-cell and bulk sequencing data. *Nat Commun* **10**, 2750 (2019).
36. Lim, S. B. *et al.* Addressing cellular heterogeneity in tumor and circulation for refined prognostication. *Proc Natl Acad Sci U S A* **116**, 17957–17962 (2019).
37. Wang, R. *et al.* Single-cell dissection of intratumoral heterogeneity and lineage diversity in metastatic gastric adenocarcinoma. *Nat Med* **27**, 141–151 (2021).
38. Hänzelmann, S., Castelo, R. & Guinney, J. GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* **14**, 7 (2013).
39. Sathe, A. *et al.* Single-cell genomic characterization reveals the cellular reprogramming of the gastric tumor microenvironment. *Clin Cancer Res* **26**, 2640–2653 (2020).
40. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381–386 (2014).
41. Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat Methods* **14**, 979–982 (2017).
42. Dalerba, P. *et al.* Single-cell dissection of transcriptional heterogeneity in human colon tumors. *Nat Biotechnol* **29**, 1120–1127 (2011).
43. Vermeulen, L. *et al.* Single-cell cloning of colon cancer stem cells reveals a multi-lineage differentiation capacity. *PNAS* **105**, 13427–13432 (2008).
44. Monk, M. & Holding, C. Human embryonic genes re-expressed in cancer cells. *Oncogene* **20**, 8085–8091 (2001).
45. Fearon, E. R. & Vogelstein, B. A genetic model for colorectal tumorigenesis. *Cell* **61**, 759–767 (1990).

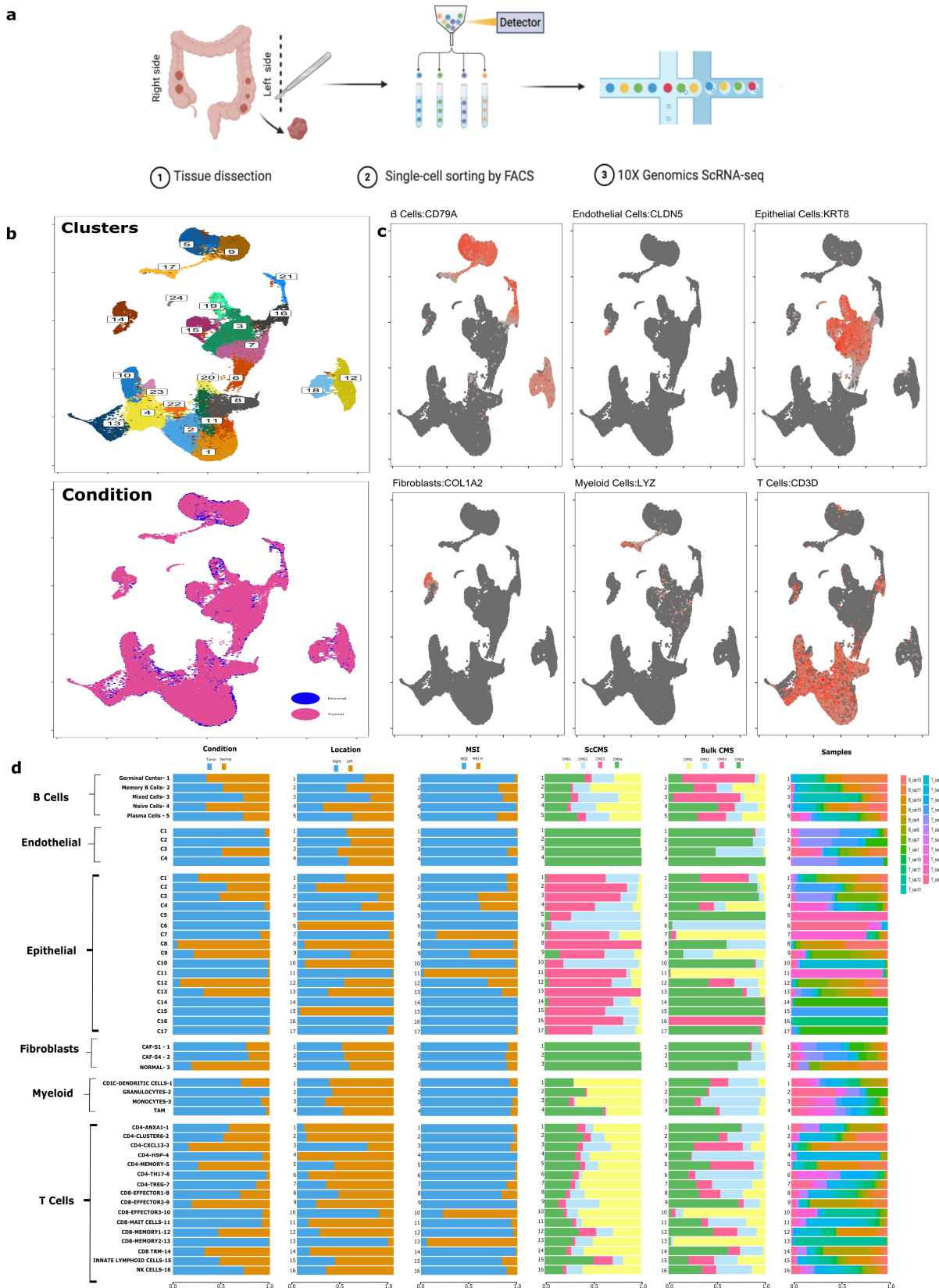


46. Costa, A. *et al.* Fibroblast heterogeneity and immunosuppressive environment in human breast cancer. *Cancer Cell* **33**, 463–479.e10 (2018).
47. Bartoschek, M. *et al.* Spatially and functionally distinct subclasses of breast cancer-associated fibroblasts revealed by single cell RNA sequencing. *Nat Commun* **9**, 5150 (2018).
48. Kieffer, Y. *et al.* Single-cell analysis reveals fibroblast clusters linked to immunotherapy resistance in cancer. *Cancer Discov* **10**, 1330–1351 (2020).
49. Wu, S. Z. *et al.* Stromal cell diversity associated with immune evasion in human triple-negative breast cancer. *The EMBO Journal* **39**, (2020).
50. Sadanandam, A. *et al.* Reconciliation of classification systems defining molecular subtypes of colorectal cancer. *Cell Cycle* **13**, 353–357 (2014).
51. De Sousa E Melo, F. *et al.* Poor-prognosis colon cancer is defined by a molecularly distinct subtype and develops from serrated precursor lesions. *Nat Med* **19**, 614–618 (2013).
52. Newman, A. M. *et al.* Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nature Biotechnology* **37**, 773–782 (2019).
53. Dunne, P. D. *et al.* Challenging the cancer molecular stratification dogma: intratumoral heterogeneity undermines consensus molecular subtypes and potential piagnostic value in colorectal cancer. *Clin Cancer Res* **22**, 4095–4104 (2016).
54. Roepman, P. *et al.* Colorectal cancer intrinsic subtypes predict chemotherapy benefit, deficient mismatch repair and epithelial-to-mesenchymal transition. *International Journal of Cancer* **134**, 552–562 (2014).
55. Ma, S. *et al.* Continuity of transcriptomes among colorectal cancer subtypes based on meta-analysis. *Genome Biology* **19**, 142 (2018).
56. Clustering rules: a comparison of partitioning and hierarchical clustering algorithms. <https://link.springer.com/article/10.1007/s10852-005-9022-1>.
57. Wilkerson, M. D. & Hayes, D. N. ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics* **26**, 1572–1573 (2010).
58. Gaujoux, R. & Seoighe, C. A flexible R package for nonnegative matrix factorization. *BMC Bioinformatics* **11**, 367 (2010).
59. Li, H. *et al.* Reference component analysis of single-cell transcriptomes elucidates cellular heterogeneity in human colorectal tumors. *Nat Genet* **49**, 708–718 (2017).
60. Liu, T. *et al.* Cancer-associated fibroblasts: an emerging target of anti-cancer immunotherapy. *J Hematol Oncol* **12**, 86 (2019).
61. Barnett, R. M. & Vilar, E. Targeted therapy for cancer-associated fibroblasts: are we there yet? *JNCI: Journal of the National Cancer Institute* **110**, 11–13 (2018).
62. Gascard, P. & Tlsty, T. D. Carcinoma-associated fibroblasts: orchestrating the composition of malignancy. *Genes Dev.* **30**, 1002–1019 (2016).
63. Sahai, E. *et al.* A framework for advancing our understanding of cancer-associated fibroblasts. *Nat Rev Cancer* **20**, 174–186 (2020).
64. Roth, A. D. *et al.* Prognostic role of KRAS and BRAF in stage II and III resected colon cancer: results of the translational study on the PETACC-3, EORTC 40993, SAKK 60-00 trial. *J Clin Oncol* **28**, 466–474 (2010).
65. Laurent-Puig, P. *et al.* Colon cancer molecular subtype intratumoral heterogeneity and its prognostic impact: An extensive molecular analysis of the PETACC-8. *Annals of Oncology* **29**, viii18 (2018).
66. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902.e21 (2019).
67. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet detection in single-cell RNA sequencing data using artificial nearest neighbors. *Cell Syst* **8**, 329–337.e4 (2019).
68. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol* **11**, R106 (2010).

69. Thompson, B. Canonical correlation analysis. in *Encyclopedia of Statistics in Behavioral Science* (American Cancer Society, 2005). doi:10.1002/0470013192.bsa068.
70. Jolliffe, I. Principal component analysis. in *International Encyclopedia of Statistical Science* (ed. Lovric, M.) 1094–1096 (Springer, 2011). doi:10.1007/978-3-642-04898-2\_455.
71. Maaten, L. van der. Accelerating t-SNE using tree-based algorithms. *Journal of Machine Learning Research* **15**, 3221–3245 (2014).
72. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: uniform manifold approximation and projection. *Journal of Open Source Software* **3**, 861 (2018).
73. Zilionis, R. *et al.* Single-cell transcriptomics of human and mouse lung cancers reveals conserved myeloid populations across individuals and species. *Immunity* **50**, 1317–1334.e10 (2019).
74. Helmink, B. A. *et al.* B cells and tertiary lymphoid structures promote immunotherapy response. *Nature* **577**, 549–555 (2020).
75. Guo, X. *et al.* Global characterization of T cells in non-small-cell lung cancer by single-cell sequencing. *Nat Med* **24**, 978–985 (2018).
76. Szabo, P. A. *et al.* Single-cell transcriptomics of human T cells reveals tissue and activation signatures in health and disease. *Nat Commun* **10**, 4706 (2019).
77. Nirschl, C. J. *et al.* IFN $\gamma$ -dependent tissue-immune homeostasis is co-opted in the tumor microenvironment. *Cell* **170**, 127–141.e15 (2017).
78. Kim, N. *et al.* Single-cell RNA sequencing demonstrates the molecular and cellular reprogramming of metastatic lung adenocarcinoma. *Nat Commun* **11**, 2285 (2020).
79. Shi, Z. *et al.* More than one antibody of individual B cells revealed by single-cell immune profiling. *Cell Discov* **5**, 64 (2019).
80. Ramesh, A. *et al.* A pathogenic and clonally expanded B cell transcriptome in active multiple sclerosis. *Proc Natl Acad Sci U S A* **117**, 22932–22943 (2020).
81. Puram, S. V. *et al.* Single-cell transcriptomic analysis of primary and metastatic tumor ecosystems in head and neck cancer. *Cell* **171**, 1611–1624.e24 (2017).
82. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with Census. *Nat Methods* **14**, 309–315 (2017).
83. Beaubier, N. *et al.* Clinical validation of the tempus xT next-generation targeted oncology sequencing assay. *Oncotarget* **10**, 2384–2396 (2019).
84. Andrews, S. *FastQC: a quality control tool for high throughput sequence data.* (Babraham Institute, 2012).
85. Jorissen, R. N. *et al.* DNA copy-number alterations underlie gene expression differences between microsatellite stable and unstable colorectal cancers. *Clin Cancer Res* **14**, 8061–8069 (2008).
86. Jorissen, R. N. *et al.* Metastasis-associated gene expression changes predict poor outcomes in patients with Dukes stage B and C colorectal cancer. *Clin Cancer Res* **15**, 7642–7651 (2009).
87. Smith, J. J. *et al.* Experimentally derived metastasis gene expression profile predicts recurrence and death in patients with colon cancer. *Gastroenterology* **138**, 958–968 (2010).
88. Skrzypczak, M. *et al.* Modeling oncogenic signaling in colon tumors by multidirectional analyses of microarray data directed for maximization of analytical reliability. *PLoS One* **5**, (2010).
89. Kemper, K. *et al.* Mutations in the Ras-Raf Axis underlie the prognostic value of CD133 in colorectal cancer. *Clin Cancer Res* **18**, 3132–3141 (2012).
90. Schlicker, A. *et al.* Subtypes of primary colorectal tumors correlate with response to targeted treatment in colorectal cell lines. *BMC Med Genomics* **5**, 66 (2012).
91. Janky, R. *et al.* Prognostic relevance of molecular subtypes and master regulators in pancreatic ductal adenocarcinoma. *BMC Cancer* **16**, 632 (2016).

92. Meister, M. *et al.* Intra-tumor heterogeneity of gene expression profiles in early stage non-small cell lung cancer. *Journal of Bioinformatics Research Studies* **1**, (2014).
93. Gautier, L., Cope, L., Bolstad, B. M. & Irizarry, R. A. affy—analysis of Affymetrix GeneChip data at the probe level. *Bioinformatics* **20**, 307–315 (2004).
94. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* **43**, e47 (2015).
95. Eide, P. W., Bruun, J., Lothe, R. A. & Sveen, A. CMScaller: an R package for consensus molecular subtyping of colorectal cancer pre-clinical models. *Sci Rep* **7**, 16618 (2017).
96. Wickham, H. *ggplot2*. (Springer).
97. Therneau, T. *A package for survival analysis in R*. (2020).
98. Zhao, X., Valen, E., Parker, B. J. & Sandelin, A. Systematic clustering of transcription start site landscapes. *PLOS ONE* **6**, e23409 (2011).
99. Bunis, D. G., Andrews, J., Fragiadakis, G. K., Burt, T. D. & Sirota, M. dittoSeq: universal user-friendly single-cell and bulk RNA sequencing visualization toolkit. *Bioinformatics* (2020) doi:10.1093/bioinformatics/btaa1011.
100. *morpheus: Interactive heat maps using 'morpheus.js' and 'htmlwidgets'*. (2021).
101. Pedersen, T. L. *Patchwork: The Composer of Plots*. R. (2020).
102. Kassambara, A. *ggpubr: 'ggplot2' Based Publication Ready Plots*. (Based Publication, 2020).

## 854 **FIGURES**

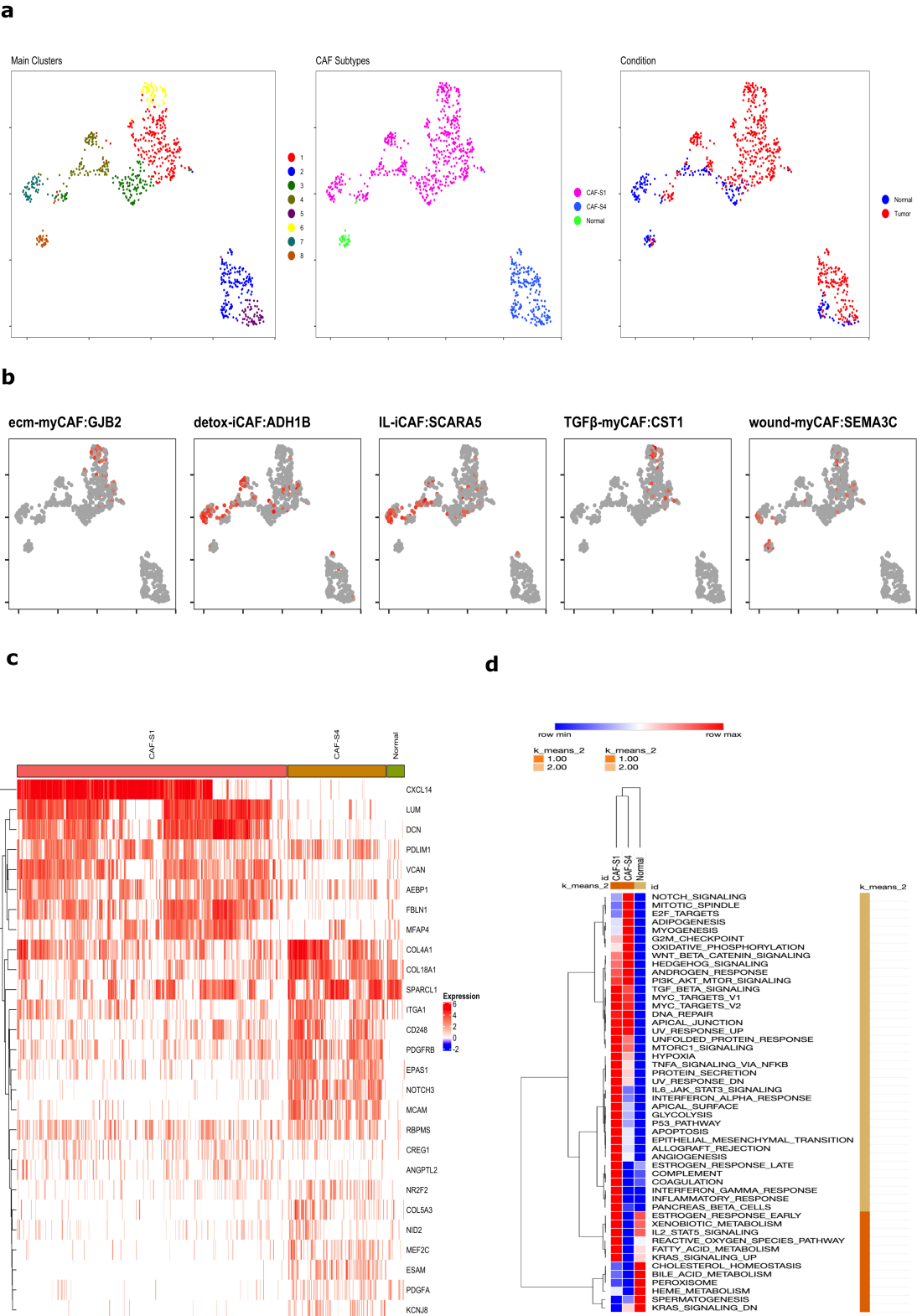


**Figure 1. Identification and clustering of single cells.** **a**, Workflow of sample collection, sorting, and sequencing (methods contain full description for each step). **b**, UMAP characterization of the 49,859 cells profiled. Coloring demonstrates clusters, tumor vs. non-malignant sample origin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



869 **Figure 2. Reclustering and characterization of the epithelial compartment.** **a**, UMAP of  
870 tumor and non-malignant epithelial reclustering demonstrating 17 distinct clusters. **b**, Heatmap  
871 of Hallmark pathway analysis within the epithelial cell compartment. **c**, Bar chart representation  
872 of cell proportions by sample, tissue type, MSI status, colonic location of sample, scCMS score,  
873 and bulk CMS score. **d**, Trajectory analysis of cells colored by colonic location, scCMS, MSI  
874 status, and bulk CMS status.

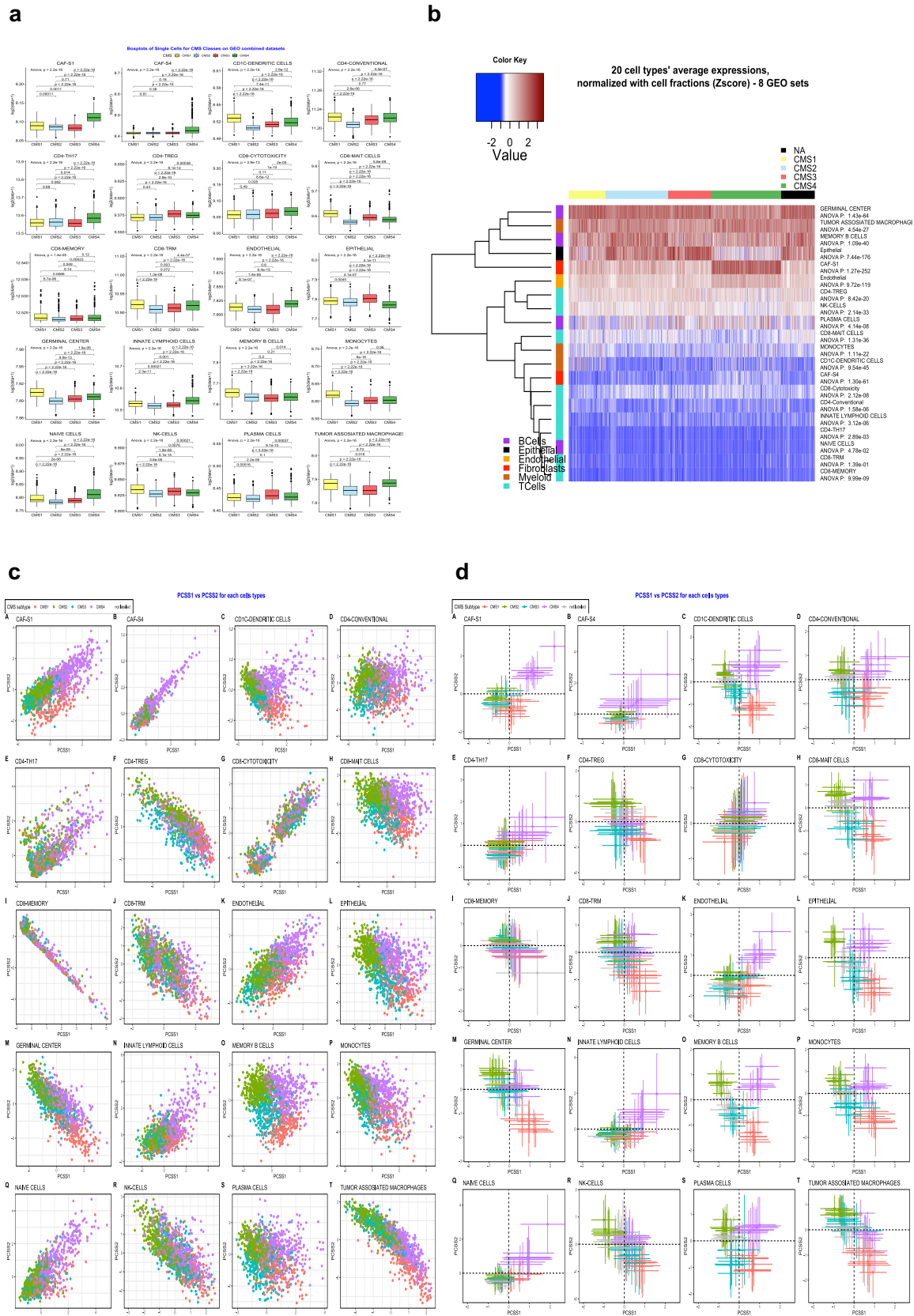




876 **Figure 3. Fibroblast clusters in colon and colorectal tumors. a**, UMAP of 819 fibroblasts  
 877 colored by distinct clusters, CAF status, tissue status and origin of sample. UMAP fibroblasts  
 878 colored by specific CAF-S1 subtypes. **b**, UMAP color-coded for marker genes for five CAF-S1  
 879 subtypes as indicated. **c**, Heatmap showing the variable expression of fibroblast specific marker  
 880 genes across CAF-S1, CAF-S4, and normal fibroblasts. **d**, Heatmap of Hallmark pathway  
 881 analysis of CAF-S1, CAF-S4, and normal cluster.



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



**Figure 5. Average cell type abundance from eight pooled CRC datasets and sorted by bulk CMS status.** **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 one-way ANOVA are shown in the figure. **b**, Deconvolution heatmap of different cell types by average expression using CIBERSORTx demonstrating cell type distribution (based on individual datasets) within each CMS category. **c**, All 20 cell types show no to little separation reported by CMS. **d**, All cell types projected on four quadrants representing CMS1-4 using PCSS1 and PCSS2 scores. Markers are colored by the bulk CMS status. Note that the cell types largely form a continuum along CMS status and are not clustered in discrete quadrants separate from one another. Cells are colored by bulk CMS status accordingly to origin of sample.