

## **CITEseq analysis of non-small-cell lung cancer lesions reveals an axis of immune cell activation associated with tumor antigen load and *TP53* mutations**

*Andrew M. Leader*<sup>1,2,3</sup>, *John A. Grout*<sup>1,2,3</sup>, *Christie Chang*<sup>1,4</sup>, *Barbara Maier*<sup>1,2,3</sup>, *Alexandra Tabachnikova*<sup>1,2,3</sup>, *Laura Walker*<sup>1,4</sup>, *Alona Lansky*<sup>1,2,3</sup>, *Jessica LeBerichel*<sup>1,2,3</sup>, *Naussica Malissen*<sup>1,2,3</sup>, *Melanie Davila*<sup>1,4</sup>, *Jerome Martin*<sup>1,2,3,5</sup>, *Giuliana Magri*<sup>1,2,3,#</sup>, *Kevin Tuballes*<sup>1,4</sup>, *Zhen Zhao*<sup>6</sup>, *Francesca Petralia*<sup>7,8</sup>, *Robert Samstein*<sup>1,2,3,9</sup>, *Natalie Roy D'Amore*<sup>10</sup>, *Gavin Thurston*<sup>11</sup>, *Alice Kamphorst*<sup>1,2,3</sup>, *Andrea Wolf*<sup>d2</sup>, *Raja Flores*<sup>12</sup>, *Pei Wang*<sup>7,8</sup>, *Mary Beth Beasley*<sup>6</sup>, *Helene Salmon*<sup>1,2,3,&</sup>, *Adeeb H. Rahman*<sup>1,4,7</sup>, *Thomas U. Marron*<sup>1,2,13</sup>, *Ephraim Kenigsberg*<sup>1,7,8,§</sup>, *Miriam Merad*<sup>1,2,3,4,§</sup>

<sup>1</sup>The Precision Immunology Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>2</sup>The Tisch Cancer Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>3</sup>Department of Oncological Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>4</sup>Human Immune Monitoring Center, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>5</sup>Nantes Université, CHU Nantes, Inserm, Centre de Recherche en Transplantation et Immunologie, UMR 1064, F-44000 Nantes, France ; CHU Nantes, Nantes Université, Laboratoire d'Immunologie, F-44000 Nantes, France

<sup>6</sup>Department of Pathology, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>7</sup>Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>8</sup>Icahn Institute for Genomics and Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>9</sup>Department of Radiation Oncology, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>10</sup>Immuno-oncology Drug Discovery Unit, Millennium Pharmaceuticals, Inc. a wholly owned subsidiary of Takeda Pharmaceutical Company Limited.

<sup>11</sup>Department of Oncology & Angiogenesis, Regeneron Pharmaceuticals Inc., Tarrytown, NY, 10591, USA

<sup>12</sup>Department of Thoracic Surgery, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>13</sup>Division of Hematology/Oncology, Icahn School of Medicine at Mount Sinai, New York, NY, USA.

Correspondence should be sent to: [miriam.merad@mssm.edu](mailto:miriam.merad@mssm.edu), [ephraim.kenigsberg@mssm.edu](mailto:ephraim.kenigsberg@mssm.edu)

<sup>§</sup>Authors contributed equally to this work

<sup>#</sup>Present address: Program for Inflammatory and Cardiovascular Disorders, Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain

<sup>&</sup>Present address: INSERM U932, Institut Curie, 26 rue d'Ulm, 75005 Paris, France

1 **SUMMARY**

2 Immunotherapy is becoming a mainstay in the treatment of NSCLC. While tumor mutational  
3 burden (TMB) has been shown to correlate with response to immunotherapy, little is known about  
4 the relation of the baseline immune response with the tumor genotype. Here, we profiled 35 early  
5 stage NSCLC lesions using multiscale single cell sequencing. Unsupervised clustering identified  
6 in a subset of patients a key cellular module consisting of *PDCDI*+ *CXCL13*+ activated T cells,  
7 IgG+ plasma cells, and *SPPI*+ macrophages, referred to as the lung cancer activation module  
8 (LCAM<sup>hi</sup>). Transcriptional data from two NSCLC cohorts confirmed a subset of patients with  
9 LCAM<sup>hi</sup> enrichment, which was independent of overall immune cell content. The LCAM<sup>hi</sup> module  
10 strongly correlated with TMB, expression of cancer testis antigens, and with *TP53* mutations in  
11 smokers and non-smokers. These data establish LCAM as a key mode of immune cell activation  
12 associated with high tumor antigen load and driver mutations.

13

## 14 INTRODUCTION

15 Lung cancer is the most common cause of cancer-related death<sup>1</sup>, and the most common subgroup  
16 of lung cancer is non-small cell lung cancer (NSCLC)<sup>2</sup>. In recent years, immune checkpoint  
17 blockade (ICB) targeting the PD-1/PD-L1 axis has become first-line therapy for a majority of  
18 patients with metastatic and locally advanced disease<sup>3</sup>. Though ICB studies have achieved  
19 improved overall survival, fewer than half of patients achieve significant clinical benefit, though  
20 still may experience physical and financial toxicity. Biomarkers are lacking to determine optimal  
21 treatment regimens for patients, as our understanding of tumor-associated immune phenotypes and  
22 immune correlates of response to ICB remains incomplete.

23 While multiple studies have used single-cell assays to profile NSCLC tumor-infiltrating  
24 immune cells in comparison to patient-matched, non-involved lung (nLung)<sup>4,5</sup>, blood<sup>6</sup>, or both<sup>7,8</sup>,  
25 or characterized tumor-infiltrating lymphocytes (TIL)<sup>9-12</sup>, we continue to lack a comprehensive  
26 understanding of how immune cell phenotypes vary across patients. In particular, it remains  
27 unclear which immune cell populations and phenotypes are associated with robust, tumor-directed  
28 T cell responses and response to ICB, and how these features are connected to tumor-cell intrinsic  
29 characteristics such as tumor mutational burden (TMB)<sup>13,14</sup>. A deeper analysis is further required  
30 for uncovering the cell types and states associated with immunostimulatory versus  
31 immunoregulatory presentation of tumor-associated antigens, as well as parsing the tumor-related  
32 effects on tissue-resident and migratory innate cell types. Attempts to integrate these data across  
33 the innate and adaptive arms of the immune system are of crucial importance to optimizing rational  
34 design of immunotherapies. Furthermore, while response to ICB has been associated with specific  
35 patient groups, individual driver mutations, the degree of immune infiltrate, and TMB, the complex  
36 interplay between these factors remains poorly understood.

37 Here, we sought to define the molecular immune states induced in the tumor  
38 microenvironment by profiling an expanded patient cohort compared to previous related studies<sup>5,6</sup>  
39 via multiscale single-cell analyses. We integrated the results of single-cell RNA sequencing  
40 (scRNAseq) of immune cells with cellular indexing of transcriptomes and epitopes by sequencing  
41 (CITEseq)<sup>15</sup>, a method allowing for combined scRNAseq and multiplexed single-cell surface  
42 protein measurement. To further elucidate the TCR landscape across T cell phenotypes, we  
43 analyzed these results together with joint scRNAseq/TCRseq. We revealed a pattern of inter-tumor  
44 variability involving innate and adaptive immune responses which we validated in two bulk RNA  
45 datasets, allowing us to detect an association with TMB and tumor driver mutations.

46

## 47 **RESULTS**

48

49 **Integrative analyses unify phenotypic mappings across substrates and datasets.** To probe  
50 transcriptional states of immune cells in the lung cancer microenvironment, we set out to profile  
51 cells from a cohort of untreated, early-stage NSCLC patients undergoing resection with curative  
52 intent (Figure 1A). The cohort was diverse with respect to age, smoking status, sex, and  
53 histological subtype (Figure 1B). We generated three datasets integrating antibody profiling of  
54 surface marker proteins using CITEseq<sup>15</sup>, scRNAseq, and TCRseq with single cell resolution  
55 (Tables S1-S3). We performed CITEseq on matched tumor and non-involved lung (nLung) tissues  
56 from 7 patients, in addition to performing scRNAseq of matched tumor and nLung in 28 additional  
57 patients. Finally, to expand on our annotation of T cell clusters based on the distribution of clonally  
58 expanded populations, we performed paired single-cell TCRseq and scRNAseq on T cells isolated  
59 from 3 patients.

60           ScRNA profiles in tumor and nLung were clustered together using a batch-aware algorithm  
61 we recently developed in order to combine data across patients, while accounting for batch-specific  
62 background noise<sup>16</sup> (Figures 1C and S1A-C). To develop a general gene-expression model of  
63 clusters representing different cell types and states, we relied only on 19 nLung and 22 tumor  
64 samples processed with 10X Chromium 3' V2 chemistry. We then used this model to classify cells  
65 from additional samples processed with different protocols or from different datasets showing  
66 similar transcriptional profiles. (Figure S1B-D, methods). The RNA-based clustering identified  
67 49 immune clusters within 6 compartments including subsets of T cells, B cells, plasma cells, mast  
68 cells, plasmacytoid dendritic cells (pDC), and mononuclear phagocytes (MNP) consisting of  
69 macrophages (M $\Phi$ ), monocytes, putative monocyte-derived dendritic cells (MoDC), and  
70 conventional dendritic cells (cDC; Figure 1C and S1E, F). Overall, 377,549 single cells from 35  
71 tumors and 32 matched nLung samples from patients at Mount Sinai were classified into 6  
72 compartments and 30 annotated transcriptional states. CITEseq data further confirmed cell  
73 identities using well-established protein markers (Figure 1D). For example, annotation of pDC  
74 was based on expression of transcripts associated with this lineage (*LILRA4*, *IRF8*; Figure 1C) and  
75 high expression of known population-defining surface markers (CD123; Figure 1D). While cluster  
76 frequencies varied widely among patients, clusters mapped between 590 and 23812 cells, and all  
77 clusters included cells from multiple patients (Figure 1E and S1F).

78           We first compared the variability of samples from different regions within individual  
79 tumors to the variability between different patients' tumors with respect to immune cell type  
80 composition. To do this, we examined samples from a study that analyzed three regions per tumor  
81 in 8 patients<sup>5</sup>, mapping cells to the clusters produced with our expanded dataset. Clustering the  
82 samples by correlation of cell type frequencies among immune cells demonstrated that samples

83 almost always clustered by patient (Figure S1G), and similarly, the Euclidean distances between  
84 patient-matched samples of different tumor regions was strongly reduced compared to the  
85 distances between samples from different patients (Figure S1H). Therefore, while the total level  
86 of immune content may still vary regionally in and around tumors<sup>17</sup>, these analyses demonstrated  
87 that inter-tumor differences drive lung tumor immune variability in terms of the phenotypic  
88 makeup of the immune cells that are present.

89 To understand whether the immune changes between tumor and nLung were distinct across  
90 patients or, alternatively, globally similar, we estimated the immune diversity within tumor and  
91 nLung using the Euclidean distances between the log-transformed cluster-frequencies. This  
92 analysis indicated that nLung samples were significantly more homogeneous (Figure 1F; “nLung-  
93 nLung distances”) than tumor samples (“Tumor-Tumor distances”;  $t=8.3$ ;  $p<2.2e-16$ ). We further  
94 compared distances among nLung and among tumor to the distances between nLung and tumor.  
95 This analysis showed that the diversity between independent (unmatched) tumor and nLung  
96 samples was larger than the diversity within tumor samples ( $t=19.6$ ;  $p<2.2e-16$ ) and nLung  
97 samples ( $t=24.6$ ;  $p<2.2e-16$ ), suggesting that immune landscapes within the TME were  
98 significantly changed compared to non-involved tissues, and that most tumors harbored many  
99 conserved changes (Figure 1E-F).

100 We next sought to test if the differences between nLung and tumor could be observed in  
101 an independent cohort. The cell type frequencies of 8 matched tumor-nLung pairs described in  
102 ref.<sup>5</sup> indeed validated the distinct microenvironments we observed in our cohort (Figure 1G). This  
103 result demonstrated that the observed tumor signatures were robust and reproducible, encouraging  
104 us to further study the transcriptional states within it.

105

106 **The intratumoral dendritic cell compartment is characterized by expansion of monocyte-**  
107 **derived DC.** We next investigated the heterogeneity with the myeloid compartment, given that  
108 different myeloid cells have various important roles in generating or inhibiting tumor directed  
109 immune responses, including antigen presentation, T cell co-stimulation, and shaping the cytokine  
110 milieu within the TME<sup>13</sup>. We identified conventional DC1 (cDC1) expressing *IRF8*, *WDFY4*, and  
111 *CLEC9A* transcripts (Figure 2A) as well as CD141 and CD26 surface markers (Figure 2B), and  
112 cDC2 expressing high *CD1C* and *FCER1A* transcripts as well as CD1c and CD5 protein. We also  
113 detected a DC cluster expressing *FSCN1* and *CCR7* transcripts and elevated HLADR, CD86, PD-  
114 L1, and CD40 surface protein which we described as mature DC enriched in regulatory molecules  
115 (mregDC) in great detail elsewhere<sup>18</sup>; in this study we found mregDC were correlated with tumor-  
116 antigen uptake and thus help define antigen-charged DC<sup>18</sup>. This phenotype was also consistent  
117 with an activated DC phenotype detected in lung and liver tumors by others<sup>6,19</sup>. We furthermore  
118 identified clusters that expressed cDC2 markers such as CD1c and *CLEC10A*, but also expressed  
119 high levels of monocyte and MΦ genes including *SI00A8*, *SI00A9*, *CIQA*, and *CIQB*, lacked  
120 surface expression of the pre-DC surface marker CD5<sup>20,21</sup>, and exhibited increased expression of  
121 CD11b and CD14 (Figure 2A, B); we annotated such clusters as MoDC. Importantly, MoDC were  
122 distinct from MΦ based on higher levels of CD1c surface protein in addition to their upregulation  
123 of the DC2-like transcriptional signature (Figure 1C, clusters 52, 29, and 30). Overall, MoDC were  
124 the most prevalent DC subtype and were increased in tumors compared to nLung, whereas  
125 mregDC were the rarest (Figure 2C and S2A-C). As we had seen previously<sup>7</sup>, the fraction of cDC1  
126 were strongly reduced in tumors (Figure 2C).

127 Since the activation profile of mregDC is crucial for inducing tumor directed T cell  
128 responses<sup>18</sup>, we examined the mregDC distribution in tumors by multiplexed

129 immunohistochemical consecutive staining on a single slide (MICSSS)<sup>22</sup>. We stained for DC-  
130 LAMP and PD-L1, as the transcripts of these genes (*LAMP3* and *CD274*, respectively) were highly  
131 enriched in the mregDC cluster (Figure 2D). We found that mregDC expressing DC-LAMP and  
132 PD-L1 accumulated in tertiary lymphoid structures (TLS) in close proximity to T cells (Figure  
133 2E). CD3-negative areas of TLS, which are putatively analogous to lymph node B cell zones, were  
134 frequently populated by MYH11<sup>+</sup> follicular dendritic cells<sup>23</sup>, a stromal cell type commonly found  
135 in B cell zones (Figure 2F).

136 To better understand the relationship between MoDC and other MNP, we searched for  
137 genes that were mutually exclusive among CD14 monocytes, cDC2, and MΦ (Figure S2D, Table  
138 S4). Scoring MoDC using these gene lists in comparison with other MNP populations revealed  
139 that MoDC were distinct from MΦ and CD14<sup>+</sup> monocytes. Ordering cells within each of these  
140 compartments by the expression of these distinct monocyte- and cDC2 gene programs  
141 demonstrated anticorrelation of these gene sets among MoDC but not cDC2 (Figure 2G;  
142  $\rho = -0.33$ ,  $p < 2.2 \times 10^{-16}$ ;  $\rho = 0.016$ ,  $p = 0.29$ , respectively), demonstrating that MoDC inhabited a  
143 phenotypic spectrum between monocytes and cDC2-like cells. While some MΦ genes were  
144 expressed in MoDC higher than in cDC2 cells, MoDC were distinct from MΦ based on higher  
145 cDC2 gene expression and lower MΦ gene expression (score distributions are detailed in Figure  
146 S2E).

147 To further uncover transcriptional programs that were variable among MoDC and cDC  
148 without relying on specific cell classifications, we analyzed the covariance structure of variable  
149 genes among all DC. This approach resulted in distinct sets of co-expressed genes (gene  
150 “modules”) that varied together across cells, independent of cluster assignments (Figure S2F, G,  
151 Table S5). Gene module analysis across the DC compartment revealed upregulation in tumors of



152 multiple modules that were mainly restricted to MoDC and DC2 (Figure S2H, I). The gene  
153 modules most upregulated in tumors compared to nlung included genes associated with glycolysis  
154 (mod39) and cell cycle (mod38), which were mainly expressed in MoDC cluster 52 (Figure S2G-  
155 I). Frequent upregulation of many monocyte- or M $\Phi$ -like modules (7, 3, 4, 6, 5, 37, 10) was  
156 consistent with a higher frequency of MoDC compared to cDC in tumors.

157 We also identified a cDC2 module (mod34) which was enriched in tumor lesions compared  
158 to nLung (Figure S2G, H) and included *CD1A* and *CD207*. These genes mark the lesional cells of  
159 Langerhans cell histiocytosis (LCH), a myeloid inflammatory condition driven by enhanced ERK  
160 activation<sup>24</sup>; we therefore referred to this module as “LCH-like”. LCH cells produce many  
161 inflammatory cytokines that promote the accumulation of Tregs and activated T cells in LCH  
162 lesions<sup>25</sup>. Interestingly, *IL22RA2*, encoding the IL22 decoy receptor IL22-BP, was also included  
163 in this module (Table S5). IL22 modulates epithelial cell growth and plays a role in tissue  
164 protection through modulation of tissue inflammation and in promoting tumor growth through  
165 induction of tissue repair<sup>26</sup>. Expression of the IL22 receptor (*IL22RA1*), meanwhile, negatively  
166 correlated with survival in KRAS-mutated lung cancer lesions<sup>27</sup>. These genes were mainly induced  
167 in the *bona fide* cDC2 cluster, but were also upregulated in MoDC (Figure 2H). Probing DC  
168 expression in an independent scRNAseq dataset of NSCLC immune cells<sup>5</sup> confirmed upregulation  
169 of these genes in tumor associated DC transcriptomes (Figure S2J).

170

### 171 **Tumors are dominated by monocyte-derived M $\Phi$ that are distinct from alveolar M $\Phi$ .**

172 While previous studies have demonstrated phenotypic differences between M $\Phi$  populating nLung  
173 versus tumors<sup>5,7</sup>, they have been limited in their ability to parse specific M $\Phi$  subpopulations with  
174 potentially distinct ontogeny and function. Our data showed remarkable heterogeneity within the

175 M $\Phi$  compartment as demonstrated by the varying expression of classical marker genes among  
176 clusters (Figures 3A and S3A). This level of resolution allowed the identification of alveolar M $\Phi$   
177 (AM $\Phi$ ) clusters expressing *SERPINA1* and *PPARG* and a cluster expressing genes consistent with  
178 interstitial M $\Phi$  (IM $\Phi$ ), which thus far have only been defined to a limited extent in humans<sup>28,29</sup>. In  
179 contrast to AM $\Phi$  that self-renew locally independent of blood precursors<sup>30</sup>, IM $\Phi$  are thought to be  
180 maintained by circulating monocyte pools even in steady state, albeit at lower rates of turnover  
181 than in settings of overt inflammation. IM $\Phi$  lacked *PPARG* and expressed *MAF* family  
182 transcription factors, *MERTK*, *CSF1R*, *LYVE1*, and *CX3CR1*<sup>31</sup>. CD14<sup>+</sup> and CD16<sup>+</sup> monocytes  
183 were defined by the expression of *CD14* or *FCGR3A* respectively and the lack of M $\Phi$  markers  
184 *MRC1*, *VSIG4*, and *SIGLEC1*. Other M $\Phi$  clusters expressed genes such as *MAFB*, *CEBPD*,  
185 *FCGR2B*, and *CSF1R*, which are indicative of monocyte origin and shared by monocytes and  
186 IM $\Phi$ ; therefore, these clusters were annotated as MoM $\Phi$ . A remaining population of M $\Phi$   
187 expressed genes consistent with primary granule formation (*AZU1*, *ELANE*, *CTSG*) but distinct  
188 from bone marrow progenitors due to lack of *MPO*, and also lacked elevation of neutrophil marker  
189 genes<sup>6</sup> *CSF3R*, *LRG1*, *FFAR4*, and *VASP* compared to other myeloid cells (Figure S3A). This  
190 cluster was referred to as *AZU1*<sup>+</sup> M $\Phi$ .

191 Using a CITEseq panel of established immune surface markers, we validated the  
192 transcription-based cluster annotations and associated new surface markers with the M $\Phi$   
193 subpopulations. For example, we found that CD10, not previously appreciated as a M $\Phi$  marker,  
194 could distinguish AM $\Phi$  from other lung myeloid populations (Figures 3B and S3B). This staining  
195 was consistent with RNA expression patterns (Figure S3A) and was verified by  
196 immunohistochemical staining (IHC) of airspace-residing AM $\Phi$  in nLung (Figure S3C). MoM $\Phi$   
197 expressed higher levels of CD11c and CD14 than other M $\Phi$  populations, whereas IM $\Phi$  were

198 notably CD14<sup>+</sup>/HLADR<sup>+</sup>/CD11c<sup>int</sup>/CD86<sup>-</sup>/CD10<sup>-</sup> (Figure 3B). Thus, CITEseq protein staining  
199 confirmed the main MΦ subpopulations identified by the transcriptional classification and defined  
200 potential sorting strategies (Figure S3B).

201 Gene module analysis across all monocyte and MΦ clusters (Figure S3D-G) revealed three  
202 broad signatures, consistent with genes that were highly expressed in both AMΦ and MoMΦ  
203 (module group I), MoMΦ or IMΦ (module group II), and monocytes (module group III; Figure  
204 S3D). Individual modules could be identifiably associated with cell type annotations as well as  
205 cell states reflecting, for example, interferon response (modules 32 and 19), heat shock genes  
206 (module 49) cell cycle (module 42), HLA class-II expression (module 28), and glycolysis (module  
207 47; Figure S3E). Examining the expression patterns of specific modules across MoMΦ clusters  
208 led us to divide them into MoMΦ subtypes I-IV: MoMΦ-II clusters expressed the highest levels  
209 of the tumor-enriched module 48, which was driven mainly by *SPPI* and also included IL-1  
210 receptor antagonist *ILIRN*, and module 47 consisting of genes indicating a glycolytically active  
211 state (*GAPDH*, *ENO1*, *LDHA*, *ALDOA*, *TPH1*), and lower levels than other MΦ of *CIQ* and HLA-  
212 class-II transcripts (Figures 3C and S3E, G). MoMΦ-III clusters were enriched in module 24  
213 (including *TREM2* and *LILRB4*) and module 25 (including *APOE* and *GPNMB*). MoMΦ-I and  
214 MoMΦ-IV were less distinctive than the other MoMac subtypes, but each comprised their own  
215 unique gene expression patterns. For example, MoMac-IV expressed the highest levels of module  
216 27 which included *CTSS*, *CFD*, and *ALDH1A1* and also expressed some genes otherwise confined  
217 to AMΦ (module 36), whereas MoMΦ-I was enriched in module 20 (including chemokine ligands  
218 *CCL13* and *CCL2*) while MoMac IV was not. Together, these analyses identify multiple tumor  
219 MoMΦ phenotypes with distinct metabolic and immunomodulatory gene programs that are  
220 enriched in the tumor milieu and likely contribute to defining the tumor microenvironment.

221 Gene set scores based on mutually exclusive, differentially expressed genes among CD14+  
222 monocytes, AM $\Phi$ , and MoM $\Phi$  (Figure S3H and Table S4) showed that AM $\Phi$  and MoM $\Phi$  were  
223 each distinct from CD14+ monocytes (Figure 3D) but that MoM $\Phi$  expressed a gradient of the  
224 CD14+ monocyte score (Figure 3E). Analysis of the gene expression patterns of hundreds of genes  
225 within the scores supported the general trends (Figure 3F). MoM $\Phi$  clusters were also distinct from  
226 the IM $\Phi$  cluster based on many transcripts and surface proteins (Figure 3A, B), although some  
227 MoM $\Phi$ , especially those that were the most distant from AM $\Phi$ , shared some IM $\Phi$  genes such as  
228 *CSF1R*, *FOLR2*, and *MERTK* (Figures 3A, F and S3A).

229 The predominant populations that increased in tumors were MoM $\Phi$ , while AM $\Phi$  were  
230 strongly depleted from tumors and IM $\Phi$  frequencies were unchanged (Figure 3G). Monocyte  
231 frequencies were also decreased, possibly reflecting their differentiation to MoM $\Phi$  or MoDC.  
232 Given that individual MoM $\Phi$ -subsets changed between nLung and tumor to different extents, we  
233 asked how these differences related to the underlying phenotypic heterogeneity within the MoM $\Phi$   
234 compartment, beyond signatures revealed by module analysis. Selecting for a set of highly  
235 expressed transcripts encoding secreted factors demonstrated strong differences between MNP  
236 subsets (Figure 3H). MoM $\Phi$ -II, the most tumor-enriched subset, expressed the highest levels of  
237 inflammatory cytokines *TNF* and *IL6*, transcripts encoding the pleiotropic factor *SPP1*, a broad  
238 collection of matrix metalloproteinases *MMP-7*, *-9*, and *-12*, as well as CCR2/5 ligands *CCL-2*, *-8*,  
239 and *-7*. By comparison, other MoM $\Phi$  populations expressed less distinct secretory profiles.  
240 Multiple MNP populations expressed the CXCR3-ligand chemotactic factors *CXCL-9*, *-10*, *-11*  
241 including MoM $\Phi$ , MoDC, and mregDC, while these ligands were distinctly absent from AM $\Phi$ ,  
242 IM $\Phi$ , AZU1+ M $\Phi$ , monocytes, cDC1, and cDC2. MregDC, meanwhile, expressed distinct  
243 cytokines and chemotactic factors associated with T cell engagement, including *IL12B*, *EBI3*,

244 *CCL17*, *CCL22*, and *CCL19*, which were expressed to a minimal or greatly reduced degree in  
245 monocytes, M $\Phi$ , or MoDC.

246

247 **TCRs limited to tumors mark T cells with distinct phenotypic features.** CITEseq  
248 characterization of T cells identified populations of CD8<sup>+</sup> cells that were characterized by an NK-  
249 like signature (T<sub>NK-like</sub>), high expression of *GZMK* (T<sub>GZMK</sub>), expression of genes related to tissue-  
250 residence such as *ITGAI* transcript and CD103 and CD69 protein (CD8<sup>+</sup> T<sub>rm</sub>), and a cluster  
251 consistent with activated T cells, expressing high levels of *IFNG*, *GZMB*, *LAG3*, *CXCL13*, and  
252 *HAVCR2* transcripts, as well as high PD-1, ICOS, and CD39 protein (T<sub>activated</sub>; Figure 4A, B). Other  
253 clusters, which mostly consisted of CD4<sup>+</sup> cells, could be separated into T<sub>reg</sub>, T<sub>rm</sub>, cells expressing  
254 a profile consistent with either central memory or naïve cells (T<sub>CM/Naïve-like-I</sub>; *TCF7*, *SELL*, *LEF1*,  
255 *MAL*, and surface expression of CD127), and a group of clusters expressing both intermediate  
256 levels of this signature as well as a tissue-residency signature (T<sub>CM/Naïve-like-II</sub>). Cells within the  
257 T<sub>CM/Naïve-like</sub> clusters did not otherwise segregate by signatures related to antigen experience, TCR  
258 engagement, activation, or exhaustion state.

259 While clustering cells using their transcriptional profiles did not result in complete  
260 separation of CD4<sup>+</sup> and CD8<sup>+</sup> cells, CITEseq allowed for the comparison of CD4<sup>+</sup> versus CD8<sup>+</sup>  
261 cells within otherwise transcriptionally similar groups. The T<sub>activated</sub> cluster could therefore be  
262 separated into CD4<sup>+</sup> and CD8<sup>+</sup> components (15.5% and 74.8%, respectively). Differential  
263 expression analysis between these subsets showed that, on average, CD4<sup>+</sup> cells in this cluster  
264 expressed increased levels of *CXCL13*, *CD40LG*, *BCL6*, and *IL21* (Figure S4A) consistent with a  
265 phenotype similar to T-follicular-helper. We next asked whether we could use profiles of CD8<sup>+</sup>  
266 and CD4<sup>+</sup> cells within this cluster to classify CD4 and CD8 cells from samples lacking CITEseq

267 surface staining, which was not available for the majority of our dataset. Learning a signature-  
268 based classifier from a training set consisting of the  $T_{\text{activated}}$  cells from 2 patients and testing this  
269 signature on the remaining patients with CITEseq staining demonstrated that transcriptional based  
270 classification guided by antibody signals was highly accurate (86% on test set; Figure S4B). This  
271 classification could further discriminate cells that uniquely expressed *CD8-A/B* transcripts or *CD4*  
272 transcripts across the remaining cells in the dataset (84% accuracy; Figure S4C). Applying this  
273 classification generally allowed for the separation of  $CD4 T_{\text{activated}}$  from  $CD8 T_{\text{activated}}$  across the  
274 dataset (Figure S4D). Similar to a recent report<sup>32</sup>, independent quantification of these cells  
275 separately and comparing their frequencies demonstrated a high correlation across tumors (Figure  
276 S4E;  $\rho=0.58$ ,  $p=2.7e-4$ ), so they were continued to be grouped for further analysis.

277 While  $T_{\text{activated}}$  and  $T_{\text{reg}}$  were the most increased T cell populations in tumors compared to  
278 nLung (Figure 4C), another cluster, characterized by high expression of cell-cycle genes *MKI67*  
279 and *STMN1*, and surface expression of HLA-DR and CD38, was also significantly increased in  
280 tumors ( $T_{\text{cycle}}$ ; Figure 4A-C). Other than expressing these hallmarks of proliferation, the  $T_{\text{cycle}}$   
281 cluster was diverse with respect to RNA and protein expression (Figure 4A, B). Analyzing the  
282 cells comprising  $T_{\text{cycle}}$  by gene scores constructed from genes differentially expressed among the  
283 other clusters demonstrated that  $T_{\text{cycle}}$  is a mixture of multiple T cell phenotypes that share the  
284 cycling state (Figure S4F and Table S4). While tumors expressed overall higher frequencies of  
285 cycling T cells (Figures 4C and S4G),  $T_{\text{activated}}$  and  $T_{\text{reg}}$  showed the highest frequencies of cycling  
286 cells compared to other phenotypes (Figure S4H).

287 To understand the clonal relationships among T cell phenotypes in tumor and nLung  
288 tissues, we performed paired scRNAseq and TCRseq using a nested PCR approach on paired  
289 tissues from 3 patients. Classification of the transcriptomes among the clusters and analysis of the

290 T cell repertoires among these phenotypes confirmed that cells mapping to the  $T_{\text{activated}}$  cluster were  
291 the most clonal population in tumors (Figure S4I). Furthermore, dividing clones into groups based  
292 on their expansion in nLung or tumor determined the phenotypes of shared clones compared to  
293 clones detected in either tissue specifically (Figures 4D, E,  
294 S4J, K). In nLung samples, the phenotypic distribution of T cells with TCRs either shared with  
295 tumor samples or only present only in nLung was similar (Figure 4Ei). In tumors, however, we  
296 observed differences in phenotypic distributions between cells with shared versus tissue-specific  
297 TCRs (Figure 4Eii). Specifically, no  $T_{\text{NK-like}}$  cells in tumors had TCRs that were uniquely expanded  
298 in tumors. Furthermore, the proportion of  $T_{\text{cycle}}$ ,  $T_{\text{reg}}$ , and  $T_{\text{activated}}$  among cells with TCRs uniquely  
299 expanded in tumors were all markedly increased compared to their proportions among cells with  
300 TCRs present in both tissues, and these relationships were not observed in nLung. By controlling  
301 for the distribution of cells with shared TCRs in the tumor, we found that the clonal enrichment in  
302 these populations was not simply due to enrichment of these phenotypes within tumor. Together,  
303 the finding that  $T_{\text{activated}}$  are enriched in clonally expanded and cycling T cells at the tumor suggests  
304 that their accumulation is at least in part due to local clonal expansion.

305

306 **B cells and plasma cells are increased in tumors, but the B:plasma cell ratio is conserved**  
307 **between tumor and nLung.**

308 B and plasma cells represented the most globally increased lineage among immune cells in tumors  
309 compared to nLung across multiple datasets (Figure 1C-E); B cells were increased as a proportion  
310 of immune cells by a median of 6.4-fold (IQR: 2.5-8.4), while plasma cells were similarly  
311 increased by a median of 4.1-fold (IQR: 2.2-9.4). B cells and plasma cells were strongly distinct  
312 both on the RNA and surface marker level (Figure 4F, G). Plasma cell clusters included rare  $\text{IgD}^+$

313 plasma cells, which were also the only CD38<sup>int</sup> population. B cell frequencies overwhelmed plasma  
314 cell frequencies with IgD<sup>+</sup> and IgM<sup>+</sup> plasma cells being the rarest, but lineage-normalized  
315 frequencies were not different between nLung and Tumor (Figures 1E and S4L). B cells and  
316 plasma cells were therefore found to increase in tumors without significant overall perturbation of  
317 the B:plasma cell ratio or plasma cell isotype ratios.

318  
319 **Ligand-receptor interactions identify potential drivers of an adaptive activation module.** In  
320 order to identify links between cellular phenotypes that may drive patient diversity, we performed  
321 correlation analyses across cell type frequencies in tumors normalized within lineage (Figure 5A).  
322 Among the most highly correlated cell types were T<sub>activated</sub>, IgG<sup>+</sup> plasma cells, and MoMΦ-II; we  
323 therefore called these cell types collectively the lung cancer activation module (LCAM). The cell  
324 types that were most anticorrelated to this module included B cells, T<sub>cm/naïve-II</sub>, AMΦ, resting cDC,  
325 and AZU1<sup>+</sup> MΦ. Sorting patients by these cell types revealed that patients could be broadly  
326 grouped into those with high or low frequencies of LCAM cell types (Figure 5B). We called these  
327 groups LCAM<sup>hi</sup> and LCAM<sup>lo</sup>, respectively. Including samples from external datasets in this  
328 stratification supported the overall pattern (Figure 5B and S5A). This stratification was not  
329 strongly associated with changes in lineage frequencies among total immune cells, and  
330 accordingly, samples from both LCAM<sup>hi</sup> and LCAM<sup>lo</sup> groups generally displayed lineage-  
331 population shifts in line with overall tumor versus nLung differences, such as decreased NK and  
332 increased B lineage frequencies (Figure 5C). Therefore, while LCAM cell types included some of  
333 the populations that were most enriched in tumor compared to nLung on average, the LCAM axis  
334 was importantly not a reflection of tumor sample purity.



335 To identify tumor-specific immune dysregulation that may contribute to shaping the  
336 LCAM<sup>hi</sup> vs. LCAM<sup>lo</sup> cellular organization among patients, we performed an unbiased analysis of  
337 ligand-receptor pairs between immune subsets, leveraging a dataset of secreted ligands and their  
338 experimentally validated receptors<sup>16,33</sup> (Figure S5B, C and Table S6), comparing differences in  
339 ligand-receptor (LR) intensity scores<sup>16</sup> between LCAM<sup>hi</sup> and LCAM<sup>lo</sup> groups, as well as between  
340 each group and their respective adjacent nLung tissues (Table S7). Overall, both LCAM<sup>hi</sup> and  
341 LCAM<sup>lo</sup> patients demonstrated correlated modes of LR activation in tumors compared to nLung  
342 (Figure 5D). In particular, tumors in both LCAM<sup>hi</sup> patients and LCAM<sup>lo</sup> patients exhibited strong  
343 intensity scores between T-cell derived CXCL13 and B cell CXCR5, which is likely contributing  
344 to the influx of B and plasma cells seen in tumors (Figure 5E). T cells in LCAM<sup>hi</sup> but not LCAM<sup>lo</sup>  
345 patients also produced other factors in tumors but not nLung capable of stimulating B cells through  
346 the IFNG-IFNGR1 axis and the BTLA-TNFRSF14 axis (Figure 5E). In addition, B cells from  
347 LCAM<sup>hi</sup> but not LCAM<sup>lo</sup> patients highly expressed TNFSF9 (41BBL), which ligates TNFRSF9  
348 (41BB), that we found highly expressed on T<sub>activated</sub> cells (Figure S5D), indicating B cells from  
349 LCAM<sup>hi</sup> patients participate in activation of T cells via TNFSF9-TNFRSF9 interaction.

350 In addition, we observed increased IFNG-IFNGR signaling between T cells and myeloid  
351 cells in LCAM<sup>hi</sup> patients (Figures 5F and S5E, F). Potentially in result, LCAM<sup>hi</sup> patients displayed  
352 higher activation of the CXCL9/10/11-CXCR3 axis between myeloid and T cells (Figure 5G-I).  
353 Whereas MΦ and MoDC demonstrated many conserved ligands upregulated in both LCAM<sup>hi</sup> and  
354 LCAM<sup>lo</sup> tumors such as IL10 and OSM (Figure 5G, H; note distribution of highlighted LR pairs  
355 along the diagonal), tumor cDC shared few ligands between the two groups, and rather upregulated  
356 CCL19 higher in LCAM<sup>hi</sup> patients compared to CCL17 in LCAM<sup>lo</sup> patients (Figure 5I); MoDC  
357 also demonstrated the latter pattern, selectively expressing CCL17 in LCAM<sup>lo</sup> patients (Figure

358 5H), suggesting that DC expression of CCL19 may be a unique feature of cDC necessary for  
359 activation of T cells as well as induction of a humoral immune response. Overall, differences in  
360 ligand-receptor intensity scores between LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patients supported such a patient  
361 stratification, and provide possible mechanistic insight into immune-cell crosstalk underlying the  
362 development of the LCAM axis, including IFN $\gamma$  signaling as a major driver.

363

364 **Projection of bulk-transcriptomic data onto scRNA-derived signatures reveals the presence**  
365 **of the LCAM<sup>hi</sup> module in two independent LUAD datasets.**

366 To identify tumor-related correlates of the LCAM module, we aimed to analyze a larger patient  
367 cohort in order to increase statistical power. Therefore, we implemented an unbiased method of  
368 scoring bulk transcriptomic signatures along the LCAM axis<sup>16,34</sup>. Specifically, we identified genes  
369 that were both differentially expressed between LCAM<sup>hi</sup> and LCAM<sup>lo</sup> tumor samples (Figure  
370 S6A), and also highly specific to the cell types enriched or depleted in the LCAM<sup>hi</sup> tumors (Figure  
371 S6B, see methods). Using a published tabulation on estimated immune content of 512 lung  
372 adenocarcinoma (LUAD) patients available from TCGA based on expression of immune genes of  
373 all lineages<sup>35</sup>, we saw as expected that scores generated with either gene set were highly correlated  
374 with estimates of overall immune content (Figure S6C, D), but an ensemble LCAM score  
375 computed by the difference of these scores (LCAM<sup>hi</sup> score – LCAM<sup>lo</sup> score) was not (Figure S6E).  
376 As predicted by our scRNAseq data, when controlling for the immune content, we generally  
377 observed negative correlations of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> gene scores among tumors except for  
378 samples with the 10% lowest immune content (Figure S6F, G), suggesting that the ensemble  
379 LCAM score might measure a mode of immune activation that is independent of the overall  
380 immune infiltration measured by immune content. We excluded the samples with the lowest 10%

381 immune content from further analysis because probing the immune signatures was likely less  
382 informative within these samples. Sorting the patients by the ensemble LCAM score revealed the  
383 presence of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patient groups within the cohort (Figure 6A).

384 To see whether a similar pattern was present in additional datasets, we probed the  
385 independent CPTAC: LUAD cohort, consisting of 110 treatment-naive LUAD patients undergoing  
386 surgical resection on whom bulk RNAseq and WES had been performed (Michael A. Gillette, et  
387 al. *Cell*, In press). Similar to the TCGA cohort, sorting the patients by the ensemble LCAM score  
388 revealed the presence of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patient groups (Figure S6H), further establishing  
389 the prevalence of this cellular module in a subset of LUAD patients.

390

### 391 **LCAM immune response correlates with tumor-genotype and expression of tumor-antigens** 392 **in LUAD lesions**

393 While the anti-tumor immune response can be modulated by many tumor-intrinsic and tumor-  
394 extrinsic factors, the tumor-infiltrating immune cells exist as part of a complex microenvironment  
395 that includes many other stromal populations<sup>13,36</sup>. To ask whether the ensemble LCAM score is  
396 associated with other non-tumor, non-immune stromal populations, we derived gene lists that were  
397 specific for individual stromal populations identified in a public dataset of 8 NSCLC patients<sup>5</sup>  
398 (Table S4), and used these genes to quantify enrichment of stromal populations in TCGA LUAD  
399 data. The ensemble LCAM score correlated with a cancer-associated fibroblast (CAF) enrichment  
400 score, anticorrelated with a normal fibroblast enrichment score, and strongly correlated with the  
401 difference of these scores (Figure 6B-D). Meanwhile, it exhibited weak or absent correlations with  
402 a tumor-associated blood endothelial cell (BEC) enrichment score, an nLung BEC enrichment  
403 score, and a lymphatic endothelial score (Figure S6I). These data suggest an intimate link between

404 development of the LCAM cellular module and a CAF-like fibroblast phenotype, which should be  
405 explored in further detail as CAF have been suggested to act as major regulators of TIL  
406 function<sup>36,37</sup>.

407 We hypothesized that variability in immune and stromal states captured by the LCAM and  
408 CAF signatures could be associated with different tumor properties. While the ensemble LCAM  
409 score demonstrated a small but significant increase in large tumors ( $t=2.60$ ,  $p=0.01$  between TNM  
410 T-stage=T1 and T-stage>T1), we observed variable LCAM presence among tumors of all stages  
411 (Figure S6J). Furthermore, while PD-L1 expression is the most commonly used biomarker guiding  
412 ICB treatment, we also observed a weak correlation between the ensemble LCAM score and total  
413 *CD274* expression ( $r=0.21$ ,  $p=2.7e-5$ ). TMB, meanwhile, has been demonstrated to be one of the  
414 most robust predictors of checkpoint response<sup>38</sup>, and is supported by the key mechanistic  
415 hypothesis that tumors with many mutations more easily activate and are targeted by the immune  
416 system via the generation of mutated peptides and damage-associated molecular patterns.  
417 Strikingly, the data showed that the ensemble LCAM score was strongly correlated with TMB  
418 both in TCGA (Figure 6E;  $r=0.47$   $p<2.2e-16$ ) and in CPTAC (Figure S6K) ( $r=0.53$   $p=2.3e-9$ ). By  
419 comparison, other scores measuring overall immune content (Immune ESTIMATE<sup>35</sup>) or specific  
420 aspects of immune state (T cell-inflamed gene expression profile (GEP) score<sup>39,40</sup>) had much  
421 weaker associations with TMB (Figure S6L, M). Importantly, correlation with TMB was observed  
422 broadly across LCAM<sup>hi</sup> genes expressed in multiple cell types, whereas conversely, anti-  
423 correlation with TMB was also observed broadly across LCAM<sup>lo</sup> genes (Figure S6N). The  
424 ensemble LCAM score correlated with TMB to similar extents among patients grouped within  
425 each TNM T-stage (Figure S6O).

426 In LUAD cases, TMB is strongly associated with smoking history. Consistent with this  
427 relationship, the ensemble LCAM score correlated with smoking pack-years ( $\rho=0.23$ ,  $p=4.4e-5$ ).  
428 Therefore, smoking history confounded the correlation we observed between TMB and the  
429 ensemble LCAM score, suggesting that the immune signature could be only indirectly related to  
430 mutations and specifically mutated neoantigens, but rather due to alternate modes of immune  
431 dysregulation related to smoking exposure. To test this hypothesis, we stratified tumors by the  
432 detection of the smoking-related mutational signature characterized by C>A de-aminations within  
433 specific trinucleotide contexts<sup>41,42</sup>. This approach removed uncertainty related to unreliability of  
434 patient-reported smoking statistics and missing clinical data. We observed that both tumors with  
435 and without detection of this signature exhibited significant correlations between TMB and the  
436 ensemble LCAM score ( $r=0.38$ ;  $p=9.2e-5$  in the undetected smoking signature group) despite  
437 having clearly distinct distributions of TMB (Figure 6E), suggesting that this relationship was  
438 independent of smoking-driven immunomodulation.

439 We then asked which additional features of the tumors may influence the ensemble LCAM  
440 score beyond the effect caused by differences in TMB. To perform this analysis, we regressed the  
441 ensemble LCAM score onto the LogTMB and correlated candidate variables with regression  
442 residuals, which quantify the difference between the observed and expected LCAM scores based  
443 on this relationship. For example, scores quantifying total predicted single-nucleotide-variant- or  
444 Insertion/deletion-induced neoantigens did not correlate with these differences (Figure S6P, Q),  
445 indicating that these neoantigen prediction scores did not provide more information regarding the  
446 LCAM immune modulation than TMB alone. However, consistent with the hypothesis that  
447 generation of tumor-associated antigens was the key mechanism connecting TMB to an LCAM  
448 response, we found that a score quantifying total tumor associated but not tumor-specific cancer-

449 testis antigens (CTA) was correlated with the regression residuals (Figure 6F;  $r=0.16$ ,  $p=3.4e-3$ ),  
450 suggesting that additional tumor-associated antigens beyond those directly caused by tumor  
451 mutations may also contribute to induction of the LCAM response.

452 Most adenocarcinoma patients have at least one of a small number of common driver  
453 mutations, including *KRAS*, *EGFR*, *STK11*, and *TP53*. Recently, it was shown that LUAD patients  
454 responsive to immune checkpoint blockade frequently have tumors harboring *TP53* mutations, and  
455 that *TP53* mutant status was associated with enrichment of CD8 T cells in the TME<sup>43,44</sup>. However,  
456 immune-related effects of individual mutations have generally not been considered independently  
457 given their correlation with TMB. Specifically, while *TP53*-mutant tumors had higher ensemble  
458 LCAM scores compared to *TP53*-WT/(*EGFR* or *KRAS* or *STK11*)-mut tumors in both TCGA and  
459 CPTAC datasets (Figures 6G and S6R), *TP53* was also most strongly associated with increased  
460 TMB (Figures 6H and S6S). In order to statistically test whether these mutations were associated  
461 with higher LCAM scores while controlling for TMB, we regressed the LCAM score onto the  
462 LogTMB and asked whether any individual mutations were correlated with the regression  
463 residuals. Interestingly, this analysis showed that *TP53*-mutant patients had higher LCAM scores  
464 than expected by a model assuming only correlation with TMB (Figure 6I;  $p=1.4e-3$ ). *KRAS*-  
465 mutant patients, meanwhile, had lower LCAM scores than expected by this model (Figure 6J;  
466  $p=1.6e-4$ ). There was no similar deviation seen in either *STK11*- or *EGFR*-mutant patients (Figure  
467 S6T). Overall, projection of bulk signatures onto axes defined by variation in our scRNAseq cohort  
468 suggested that expression of the LCAM cellular module is a marker of adaptive response against  
469 mutated and ectopically-expressed tumor-associated antigens that is independent from the overall  
470 level of immune infiltration.

471

472 **DISCUSSION**

473 The analysis of matched tumor and nLung tissues from 35 patients as described here provides the  
474 largest unbiased single-cell map of the immune response of early-stage lung cancer lesions to date.  
475 CITEseq analysis, combining phenotypic classifications based on surface protein expression with  
476 transcriptomic profile, serves here to help unite high-dimensional models of cellular classification  
477 and refine our understanding of the immune cellular landscape in disease lesions. By further  
478 integrating tumor and nLung samples from public datasets, we demonstrated the robustness of the  
479 reported signatures across platforms. Importantly, based on high levels of changes conserved  
480 across tumor lesions, these data support the notion that common immunotherapy treatment  
481 paradigms could be beneficial for large subsets of patients despite existing disease heterogeneity.

482         Among tumors, patients could, however, be stratified along a dominant LCAM axis that  
483 was independent of overall immune infiltration or changes in proportions of immune lineages.  
484 This axis was defined by a high level of IgG<sup>+</sup> plasma cells, activated T cells that were clonally  
485 enriched in the tumor and expressing a proliferation signature, and MoMac-II that expressed *SPPI*,  
486 a glycolysis signature, and a set of inflammatory secreted factors; this module of cell types  
487 anticorrelated with B cells, T cells with a Tcm/naïve-like phenotype, resting cDC, AMΦ, and MΦ  
488 expressing *AZUI*. We therefore propose that LCAM<sup>hi</sup> patients are undergoing a more vigorous  
489 antigen-specific antitumor adaptive immune response, whereas LCAM<sup>lo</sup> patients fail to mount an  
490 adaptive response to such a degree. Unbiased ligand-receptor analyses showed that, while both  
491 LCAM<sup>hi</sup> and LCAM<sup>lo</sup> tumors expressed similar patterns of ligand-receptor pairs among immune  
492 cells compared to nLung, LCAM<sup>hi</sup> status was specifically related to heightened *CXCL13*  
493 expression by T cells, IFNγ signaling from T cells to myeloid and B cells, and *CXCL-9,10,11*  
494 signaling from myeloid cells; cDC meanwhile expressed more *CCL19* in LCAM<sup>hi</sup> tumors

495 compared to more *CCL17* in LCAM<sup>lo</sup> tumors. These factors likely served to modify the immune  
496 response around a set of conserved changes compared to nLung observed in both LCAM<sup>hi</sup> and  
497 LCAM<sup>lo</sup> tumors.

498         When analyzed in the context of broader datasets with paired bulk transcriptomics and  
499 whole exome sequencing, an ensemble gene score learned from the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patients  
500 and associated cell types strongly correlated with measures indicative of high levels of tumor-  
501 associated antigens, namely TMB and a cancer testis antigen score. Interestingly, the LCAM score  
502 was not correlated with the overall immune infiltrate, and was independent of the T-cell inflamed  
503 gene expression profile score commonly used to reflect immune activation in tumor lesions<sup>39,40</sup>.  
504 While the ensemble LCAM score was correlated with smoking status and weakly correlated with  
505 stage, the relationship with TMB remained even after controlling for these possible confounders.  
506 Given that TMB has demonstrated predictive power with response to ICB response in NSCLC<sup>14,38</sup>,  
507 the relationship between TMB and the LCAM cellular module in treatment-naïve patients suggests  
508 that this effect may be mediated via a conditioning of the immune system that exists prior to  
509 treatment, and that measurement of this cellular module may provide a more direct indicator with  
510 respect to the immune system's propensity for ICB response. Specifically, the fact that many  
511 factors significantly influence the LCAM score, not just TMB, demonstrates how the immune  
512 system integrates multiple types of signals to establish its set point.

513         Importantly, while previously reported immune signatures have been proposed to reflect  
514 tumor cytolytic activity or T cell and IFN $\gamma$ -driven immune response in association with tumor  
515 antigens and immune evasion modes<sup>39,40,45</sup>, the LCAM axis presented here represents an integrated  
516 assessment of the immune cellular organization, based on all immune cell types as defined by  
517 scRNA across patients, likely arising as a direct response to tumor antigens.



518           An additional question of clinical interest relates to how different driver mutations affect  
519 the conditioning of the immune system and ICB response. The analysis presented here shows that  
520 the common LUAD driver mutations *EGFR* and *STK11* had little effect on the LCAM response  
521 beyond that explained by their association with TMB. While *STK11* mutation status has been  
522 shown to be the most prevalent genomic driver of primary resistance to ICB<sup>44,46</sup>, there were no  
523 patients with *STK11*-mutated tumors in our scRNAseq cohort, so this effect can therefore not be  
524 addressed here. Meanwhile, compared to what was expected based on each tumor's TMB alone,  
525 *TP53* mutation intensified the LCAM response and *KRAS* mutation blunted it. Interestingly, the  
526 latter result is consistent with a recent report demonstrating that pharmacological blockade of  
527 *KRAS*-G12C in preclinical studies resulted in a robust immune response and synergized with  
528 anti-PD1 treatment<sup>47</sup>. The mechanisms of these effects remain to be seen, and may relate to the  
529 expression of immunomodulatory factors by the tumor, or the re-shaping of the metabolic  
530 microenvironment, for example. To elucidate such pathways, close study of the tumor on a broader  
531 molecular scale, in conjunction with the immune cell composition and state, is necessary.

532           A further, surprising result from our bulk RNA analyses was that the LCAM axis was  
533 highly consistent with a change in fibroblast phenotype based on signatures derived from  
534 scRNAseq of NSCLC stromal clusters<sup>5</sup>. This association could suggest that the development of  
535 the tumor fibroblast phenotype is in response to overwhelming immune activation that may be  
536 instigated by an adaptive, antigen-specific response.

537           An important limitation of these findings relates to the site of initiation of the LCAM  
538 response; while the LCAM cellular module consists of cells undergoing an apparent active immune  
539 response, this study does not demonstrate the extent to which the module is instigated or  
540 perpetuated *in situ* at the tumor lesion. Specifically, despite evidence of clonally expanding T<sub>activated</sub>

541 cells, it remains unclear whether these lineages are primed *in situ* versus in tumor-draining lymph  
542 nodes (TdLN). While understanding the timescale of the tumor specific response will always be  
543 challenging due to variation in patient presentation timelines, it will nevertheless be important to  
544 correlate the cell types and states present in the TdLN in order to determine whether the LCAM  
545 response depends on lymph node priming, as well as to develop a deeper understanding of the  
546 spatial dynamics of the LCAM cell types.

547 Overall, the model presented here identifies an immune activation signature, derived from  
548 definitions of immune phenotypes defined by single-cell RNA and CITEseq, as an integrator of  
549 tumor-associated antigen load and driver mutation status that is not related to overall immune  
550 content. We believe that this axis, therefore, can serve as a more direct measure of antigen-specific,  
551 anti-tumor immune activation compared to previously suggested immune readouts.

552

553 **METHODS**

554

555 **Human subjects**

556 Samples of tumor and non-involved lung were obtained from surgical specimens of patients  
557 undergoing resection at Mount Sinai Hospital (New York, NY) after obtaining informed consent  
558 in accordance with a protocol reviewed and approved by the Institutional Review Board at the  
559 Icahn School of Medicine at Mount Sinai (IRB Human Subjects Electronic Research Applications  
560 10-00472 and 10-00135) and in collaboration with the Biorepository and Department of Pathology.

561

562 **Tissue processing**

563 Tissues were rinsed in PBS, minced and incubated for 40 minutes at 37°C in Collagenase IV  
564 0.25mg/ml, Collagenase D 200U/ml and DNase I 0.1mg/ml (all Sigma). Cell suspensions were  
565 then aspirated through a 18G needle ten times and strained through a 70-micron mesh prior to RBC  
566 lysis. Cell suspensions were enriched for CD45<sup>+</sup> cells by either bead positive selection (Miltenyi)  
567 per kit instructions or FACS sorting on a BD FACSAria flow sorter (as indicated in Table S1)  
568 prior to processing for scRNAseq or CITEseq.

569

570 **ScRNA- and TCR-seq**

571 For each sample, 10,000 cells were loaded onto a 10X Chromium single-cell encapsulation chip  
572 according to manufacturer instructions. Kit versions for each sample are indicated in Table S1.  
573 Libraries were prepared according to manufacturer instructions. QC of cDNA and final libraries  
574 was performed using CyberGreen qPCR library quantification assay. Sequencing was performed  
575 on Illumina sequencers to a depth of at least 80 million reads per library.

576 TCRseq was performed using the Chromium Single Cell 5' VDJ kit, following  
577 manufacturer's instructions. For patients 695 and 706, cells were subject to a CD2+ bead  
578 enrichment (Miltenyi) instead of CD45+ enrichment prior to encapsulation.

579

#### 580 **CITEseq**

581 For each sample, cell suspensions were split and barcoded using “hashing antibodies”<sup>48</sup> staining  
582 beta-2-microglobulin and CD298 and conjugated to “hash-tag” oligonucleotides (HTOs). Hashed  
583 samples were pooled and stained with CITEseq antibodies that had been purchased either from the  
584 Biologend TOTALseq catalog or conjugated using the Thunder-Link PLUS Oligo Conjugation kit  
585 (Expedeon). Sample hashing schemes and CITEseq panels are detailed in Tables S1 and S2,  
586 respectively. Stained cells were then encapsulated for single-cell reverse transcription using the  
587 10X Chromium platform and libraries were prepared as previously described<sup>15</sup> with minor  
588 modifications. Briefly, cDNA amplification was performed in the presence of 2pM of an antibody-  
589 oligo specific primer to increase yield of antibody derived tags (ADTs). The amplified cDNA was  
590 then separated by SPRI size selection into cDNA fractions containing mRNA derived cDNA  
591 (>300bp) and ADT-derived cDNAs (<180bp), which were further purified by additional rounds of  
592 SPRI selection. Independent sequencing libraries were generated from the mRNA and ADT cDNA  
593 fractions, which were quantified, pooled and sequenced together on an Illumina Nextseq to a depth  
594 of at least 80 million reads per gene expression library and 20 million reads per ADT library.

595

#### 596 **MICSSS**

597 FFPE tissues were stained using multiplexed immunohistochemical consecutive staining on a  
598 single slide as previously described<sup>22</sup>. Briefly, slides were baked at 37°C overnight, deparaffinized

599 in xylene, and rehydrated in decreasing concentrations of ethanol. Tissue sections were incubated  
600 in citrate buffer (pH6 or 9) for antigen retrieval at 95°C for 30 minutes, followed by incubation in  
601 3% hydrogen peroxide and in serum-free protein block solution (Dako, X0909) before adding  
602 primary antibody for 1 hour at room temperature. After signal amplification using secondary  
603 antibody conjugated to streptavidin-horseradish peroxidase and chromogenic revelation using 3-  
604 amino-9-ethylcarbazole (AEC), slides were counterstained with hematoxylin, mounted with a  
605 glycerol-based mounting medium and scanned for digital imaging (Pannoramic 250 Flash III  
606 whole-slide scanner, 3DHISTECH). Then the same slides were successively bleached and re-  
607 stained as previously described<sup>22</sup>. Primary antibodies were: anti-human CD10 (200103, R&D  
608 systems), DC-Lamp (1010E1.01, Novus biologicals), pan-cytokeratin (AE1/AE3, Dako), PDPN  
609 (D@-40, Ventana), CD163 (10D6, Novus Biologicals) and PD-L1 (E1L3N, Cell Signaling Tech).

610

### 611 **Analysis of Sequencing data**

612 Transcriptomic and TCR library reads were aligned to the GRCh38/84 reference genome and  
613 quantified using Cellranger (v3.1.0). CITEseq ADT and CITEseq HTO reads were queried for  
614 antibody- and cell-specific oligonucleotide sequence barcodes in the designated read positions,  
615 including antibody sequences within a Hamming distance of 1 from the reference, using the  
616 feature-indexing function of Cellranger. Resulting alignment statistics are reported in Table S3.  
617 TCR data was aligned using Cellranger *vdj* function with default parameters.

618

### 619 **CITEseq processing and normalization**

620 Doublets were removed based on co-staining of distinct sample-barcoding (“hashing”) antibodies  
621 ( $[\textit{maximum HTO counts}]/[\textit{2}^{nd} \textit{ most HTO counts}] < 5$ ) and cell barcodes with few HTO counts

622 (*maximum HTO counts* < 10) were also excluded. Cells were then assigned to samples based on  
623 their maximum staining HTO. HTO to sample associations are detailed in Table S1.

624 To normalize ADT counts across experimental batches given different CITEseq staining  
625 panels and sequencing runs, we performed a quantile-normalization on the ADT count values for  
626 each surface marker for the immune cells in each 10X encapsulation batch. To do this, the  
627 geometric average of the quantile function was computed across batches

$$628 \quad \overline{F_m^{-1}(p)} = \left( \prod_{b=1}^N [F_{m,b}^{-1}(p) + d] \right)^{\frac{1}{N}}$$

629 where  $F_{m,b}^{-1}(p)$  is the quantile function, or inverse cumulative distribution function, for counts of  
630 CITEseq marker  $m$  on immune cells in each of  $N$  10X encapsulation batches  $b$  and regularization  
631 factor  $d=1$  ADT count, evaluated at quantile  $p$  in interval  $[0,1]$ . This geometric average quantile  
632 function provided a reference function for a common mapping of cells based on their single-  
633 channel, batch-specific staining quantile  $p$  to a normalized staining intensity. Of note, this  
634 normalization method preserved the relationships between channels while constraining the  
635 observed differences in staining across experiments within individual channels.

636

### 637 **Unsupervised batch-aware clustering analysis**

638 Immune cells from tumor and nLung samples were filtered for cell barcodes recording > 500 UMI,  
639 with < 25% mitochondrial gene expression, and with less than defined thresholds of expression  
640 for genes associated with red blood cells and with epithelial cells (Table S4). Cells were clustered  
641 using an unsupervised batch-aware clustering method we have recently described<sup>16</sup> with minor  
642 adjustments. This EM-like algorithm, which was also based on earlier studies<sup>49,50</sup>, iteratively  
643 updates both cluster assignments and sample-wise noise estimates until it converges, using a

644 multinomial mixture model capturing the transcriptional profiles of the different cell-states and  
645 sample specific fractions of background noise. We clustered 19 nLung and 22 tumor samples  
646 jointly and 46 additional tumor and nLung samples were mapped onto the final model as described  
647 below.

648 The model definitions and estimation of model parameters were as described in <sup>(16)</sup>.  
649 Specifically, the probability of observing gene  $i$  in cell  $j$  is defined as:

$$650 \quad p_{ji} = \frac{1}{Z} \left[ K_{reg} + (1 - \eta_{bj}) \cdot \alpha_{i, map^j} + \eta_{bj} \cdot \beta_{i, bj} \right]$$

651 Where  $map^j$  and  $b^j$  are assignments of cells  $j$  to cell-type and batch respectively;  $\eta_{bj}$  is the fraction  
652 of UMIs contributed by background noise;  $\alpha_{i, map^j}$  is the probability that a molecule drawn from  
653 celltype  $map^j$  is of gene  $i$  (assuming no background noise)  $\beta_{i, bj}$  is the probability that a noise UMI  
654 drawn from batch  $b^j$  will be of gene  $i$ , and  $K_{reg}$  is a small regularization constant.

655 We also used here the pseudo expectation-maximization (EM) algorithm<sup>16</sup> to infer the  
656 model parameters with minor modifications: (1) training set size was 2000 instead of 1000 cells  
657 and (2) the best clustering initiation was selected from 1000 instead of 10000 kmeans+ runs. For  
658 this clustering we included barcodes with more than 800 UMIs and used  $K_{reg\_ds} = 0.2$  ;  $(P_1, P_2) =$   
659  $(0^{th}, 30^{th})$  percentiles;  $K_{reg} = 5 \cdot 10^{-6}$  ;  $k=60$ . Genes with high variability between patients were  
660 not used in the clustering. Those genes consisted of mitochondrial, stress, metallothionein genes,  
661 immunoglobulin variable chain genes, HLA class I and II genes and 3 specific genes with  
662 variable/noisy expression: *MALATI*, *JCHAIN* and *XIST* (Table S4). Ribosomal genes were  
663 excluded only from the k-means clustering (Step 2.D as described in <sup>(16)</sup>). Samples used to  
664 generate this model included only those that were enriched for CD45+ immune cells using bead  
665 enrichment and were processed with the 10X Chromium V2 workflow.

666

## 667 **Integration of additional single-cell data**

668 The resulting clustering model was used to analyze additional data that was both generated in-  
669 house or downloaded from public datasets. Single cells were mapped to clusters defined by the  
670 previously generated model  $\alpha$ . Similarly to the clustering iterations, this process associates single-  
671 cells of a sample with multinomial probability vectors defined by the model and estimates the  
672 noise fractions of the sample  $\eta_b$  by optimizing the likelihood function (<sup>16</sup>):

$$673 \quad f(\eta_b) = \sum_j \sum_i U_{ij} \log(p_{ji})$$

674 For  $p_{ji}$  as defined above, while  $\alpha_{i, \text{map}j}$  are updated using maximum likelihood.

675 Integrating inDrop data from (<sup>6</sup>) and 10X Chromium 5' data required addressing the  
676 systematic differences<sup>51</sup> in gene capture present between these technologies and 10X Chromium  
677 3' data that was used to develop the clustering model. Analysis of the differences in gene  
678 expression between the technologies suggested that a multiplicative correction factor  $C_i$  per each  
679 gene  $i$  could adjust for the capture efficiency differences. The following process was used to  
680 estimate the correction parameters:

- 681 1. Map cells to the original cluster models, as above, assuming absent noise in order to prevent  
682 the estimated noise term from being driven by error due to batch differences instead of true  
683 noise.
- 684 2. Re-calculate models using the average expression of the mapped cells for each cluster to  
685 form “data-based models”  $\alpha^D$ .
- 686 3. Calculate a weight matrix  $W$ , that weights individual genes for each cluster.  $W$  is calculated  
687 by

$$688 \quad W_{i,j} = \max(\alpha_{i,j}, \alpha^D_{i,j}) + w_{reg}$$



689 for original cluster model matrix  $\alpha$ , data-based cluster model  $\alpha^D$ , gene  $i$ , cluster  $j$ , and  
690 regularization constant  $w_{reg} = 10^{-10}$ . Since highly detected genes tend to dominate the  
691 mapping results, it is important to account for genes that are highly detected in either the  
692 original (10X Chromium V2) platform or the new platform

693 4. Construct a vector of gene-specific conversion factors that can operate between platforms:

$$694 \quad C_i = \sum_j W_{i,j} ([c + \alpha^D_{i,j}] / [c + \alpha_{i,j}])$$

695 for regularization factor  $c = 10^{-6}$ .

696 5. Generate transformed cluster models  $\alpha'_{i,j}$  by multiplying the original models by the  
697 conversion vector and dividing by a normalization factor  $Z$ :

$$698 \quad \alpha'_{i,j} = \frac{1}{Z} * C_i \alpha_{i,j}$$

699 6. Map cells to transformed models without fixing the noise.

700

701 Analysis of the gene expression profiles of the mapped cells in each cluster demonstrated  
702 correspondence between the model and the mapped samples across the different technologies.  
703 (Fig S1B).

704

#### 705 **Analysis of public datasets**

706 Fastqs of scRNAseq data of tumors and nLung from 8 NSCLC patients<sup>5</sup> acquired using 10X  
707 Chromium protocols was downloaded from ArrayExpress accessions E-MTAB-6149 and E-  
708 MTAB-6653. Sequencing reads were re-aligned using Cellranger as described above. Single-cells  
709 were mapped to clusters as described above. Tumor samples included 3 separate samples from the  
710 core, middle, and edge of each tumor. Regional tumor samples were considered separately for the

711 intra- versus inter-patient variability analyses (Figure S1G, H). For remaining analyses, cell counts  
712 of projected tumor samples were pooled by patient.

713 scRNAseq data of tumors from 7 NSCLC patients<sup>6</sup> acquired using inDrop was downloaded  
714 from GEO accession GSE127465. Since neutrophils were not detected in 10X Chromium data,  
715 cells that were annotated as neutrophils in the  
716 *GSE127465\_human\_cell\_metadata\_54773x25.tsv.gz* file were excluded from analysis. Cells were  
717 classified by projection as described above, using the modified procedure for inDrop data.

718 TCGA LUAD RNAseq data was downloaded using the *GDCquery* and *GDCdownload*  
719 functions from the *TCGAbiolinks* R package. *GDCquery* options included *project="TCGA-*  
720 *LUAD"*, *data.category="Transcriptome Profiling"*, *data.type="Gene Expression*  
721 *Quantification"*, *workflow.type="HTSeq – FPKM"*, *experimental.strategy="RNA-Seq"*, and  
722 *legacy=F*. Whole exome sequencing data was downloaded using the *GDCquery\_Maf* function  
723 with arguments *tumor="TCGA-LUAD"* and *pipelines="mutect2"*. Clinical data was downloaded  
724 using the *GDCquery\_clinic* function with arguments *project="TCGA-LUAD"* and  
725 *type="clinical"*.

726 Processed CPTAC lung adenocarcinoma data was downloaded from the CPTAC Data  
727 Portal <https://cptac-data-portal.georgetown.edu/cptacPublic/>.

728

## 729 **Determination of sample-sample distances**

730 Sample-sample distances were computed as the Euclidean distance between vectors consisting of  
731 the Log<sub>10</sub>-transformed cell type frequencies, where frequencies were computed as a fraction of  
732 total immune cells. A regularization factor of 10<sup>-3</sup> was applied prior to applying the log-transform.

733

### 734 **Determination of myeloid cell type-specific gene scores**

735 Lists of mutually-exclusive genes were used to compare monocytes, cDC2, and M $\Phi$  in Figure 2,  
736 and monocytes, AM $\Phi$ , and MoM $\Phi$  in Figure 3. For these analyses, genes were identified as  
737 “mutually exclusive” if the average expression was at least 2x greater in a given population than  
738 in the other comparison populations. To account for the large diversity of MoM $\Phi$  clusters, the  
739 maximum average expression of each MoM $\Phi$  subtype was used instead of the overall average  
740 expression. Resulting gene lists are presented in Table S4. Cells were scored according to the  
741 resulting gene lists as the Log-transformed fraction of UMI belonging to the gene list. Histograms  
742 were generated with the R function *density* using default parameters.

743

### 744 **Modules analyses**

745 Gene-gene correlation modules were generated using a similar method to that previously  
746 described. Briefly, cells were downsampled to 2000 UMI prior to selecting a set of variable genes,  
747 similar to the selection of genes in preparation for seeding the clustering<sup>16</sup>. The gene-gene  
748 correlation matrix for this gene set was then computed for each sample over the cell population of  
749 interest. Correlation matrices were averaged following a Fisher Z-transformation. Applying the  
750 inverse transformation then resulted in the best-estimate correlation coefficients of gene-gene  
751 interactions across the dataset. Genes were clustered into modules using complete linkage  
752 hierarchical clustering over correlation distance. Histograms of module expression scores were  
753 generated with the R function *density* using default parameters.

754

### 755 **Classification of CD4+ versus CD8+ T<sub>activated</sub> cells**

756 CITEseq staining on a subset of patients was used to build a gene-set-based classifier that could  
757 use mRNA UMI data to discriminate CD4<sup>+</sup> versus CD8<sup>+</sup> cells within the T<sub>activated</sub> cluster. To  
758 identify these gene sets, cells from 2 patients used as a training set were gated based on a Log<sub>2</sub>FC  
759 of raw ADT counts of raw CD4/CD8 > 1 and compared by differential expression. Genes were  
760 filtered by expression > 10<sup>-4</sup> and a Log<sub>2</sub>FC > 1, and nonspecific or noise-related genes such as  
761 those associated with cell-cycle, long-non-coding RNAs, heat shock proteins, immunoglobulin  
762 genes, ribosomal proteins, *XIST*, and histone transcripts. Resulting gene lists are reported in Table  
763 S4. Cells were scored based on the fraction of RNAs belonging to the resulting gene lists, and a  
764 discrimination threshold for the ratio of the CD4 vs. CD8 gene lists was determined based on the  
765 overall accuracy in discriminating between CITEseq-defined CD4<sup>+</sup> vs. CD8<sup>+</sup> cells in the training  
766 set. This gene score discriminator was validated using cells from a test set comprised of cells from  
767 4 additional patients analyzed by CITEseq (584, 593, 596, 630), and on cells with unique detection  
768 of either *CD4* or at least one of (*CD8A*, *CD8B*).

769

#### 770 **Analysis of cycling T cell cluster**

771 To analyze the phenotypic makeup of the cluster of T cells expressing cell-cycle genes, we  
772 generated gene sets based on the other T cell phenotypes described here to score each cell within  
773 the cluster. To do this, we pooled the cells of each other T cell phenotype to compute its average  
774 expression. We then identified a gene list for each phenotype defined by expression > 1e-5 and  
775 Log<sub>2</sub>FC > 0.25 compared to the maximum of the other phenotypes. From this list, we excluded  
776 variable TCR genes, and other genes associated with noise or cell stress. The gene lists for the  
777 T<sub>naive/CM-like</sub> cell types were grouped, since these phenotypes were very similar. Resulting gene lists  
778 are reported in Table S4.

779 For each cell in the cycling cluster, we then computed the fraction of UMI belonging to  
780 each gene signature after removing UMIs belonging to the list of genes associated with the cycling  
781 cluster that was calculated as above. We performed spherical k-means clustering using the function  
782 *skmeans()* in the *skmeans* R package on these signature fractions in order to group cells within the  
783 cycling cluster according to phenotypic subtype by spherical k-means cluster.

784

### 785 **Single-cell TCRseq analysis**

786 Single T cells were grouped by clonotype according to their precise combination of  $\alpha$  and  $\beta$  chains  
787 present (uniquely defined by CDR3 sequence and V, D, and J gene usage), with the following  
788 acceptations in order to filter for high quality singlets:

789 1. Cells with contigs encoding  $> 3$  productive  $\alpha$  and  $\beta$  chains were excluded as multiplets.

790 2. Cells with contigs encoding  $> 3$  productive  $\alpha$  and  $\beta$  chains that completely overlapped  
791 with observed cells within the multiplets were also excluded as multiplets.

792 3. Remaining cells with 3 unique  $\alpha$  and  $\beta$  chains that could be uniquely associated with  
793 similar cells displaying 2 unique  $\alpha$  and  $\beta$  chains were assumed to be clonally related, whereas cells  
794 that could be similarly associated with multiple distinct sets of cells expressing 2 unique  $\alpha$  and  $\beta$   
795 chains were excluded as doublets.

796 4. Cells in which a single TCR chain was observed were assumed to be clonally related to  
797 any cells with 2 unique  $\alpha$  and  $\beta$  chains to which they uniquely associated.

798 5. Remaining cells in which a single TCR chain was observed were excluded if they  
799 matched ambiguously to multiple cells with 2- or 3-chains.

800 Clonality scores were computed for each T cell type in each patient as *1-Peilon's evenness*  
801 over the set of unique TCRs as previously described<sup>52</sup>.

802

### 803 **Ligand-receptor analysis**

804 Ligand-receptor intensity scores for a set of secreted ligands (ref<sup>(33)</sup>) and Table S6) were calculated  
805 as previously reported<sup>16</sup>. Briefly, for each ligand-receptor interaction, for each source cell type and  
806 each receiver cell type, the intensity score was equal to the product of ligand generation from the  
807 source cell type relative to the total RNA with the expression of the receptor on the receiver cell  
808 type. Scores were independently calculated for LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patient sets in nLung and  
809 Tumor tissues. To determine these patient sets, patients were sorted by the geometric mean of  
810 lineage-normalized cellular frequencies of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> cell types, and the top half of  
811 patients were defined as LCAM<sup>hi</sup> with the bottom half defined as LCAM<sup>lo</sup>. Only patients analyzed  
812 using 10X Chromium V2 with immune cells purified with magnetic beads were used for this  
813 analysis. The patients included in these groups were: LCAM<sup>hi</sup>: (408, 403, 522, 371, 570, 714, 584,  
814 377, 406, 564, 630, 578, 514); LCAM<sup>lo</sup>: (571, 596, 393, 593, 626, 378, 370, 410, 572, 558, 581,  
815 596, 729).

816

### 817 **Identification of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> bulk-RNA gene signatures**

818 To define genes that could probe the presence of LCAM<sup>hi</sup> or LCAM<sup>lo</sup> cell types in bulk RNA  
819 data, we adopted a similar strategy to that used previously for the projection of bulk data onto  
820 signatures defined by cellular axes as measured with scRNA<sup>16,34</sup>. Cells were evenly sampled from  
821 LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patients (1409 cells per patient), and sampled cells were then pooled within  
822 the groups. Differentially expressed genes ( $FDR < 10^{-3}$  and  $\text{Log}_2\text{FC} > 1$ ) were retained. Genes that  
823 were expressed in the filtered epithelial cells  $> 2x$  higher than in immune cells on average were  
824 removed. Among the remaining differentially expressed genes, those that were expressed on

825 average within any LCAM<sup>hi</sup> or LCAM<sup>lo</sup> subtype with Log2FC > 3 compared to the highest  
826 expressing subtype in the opposite group were retained. These gene lists were further abbreviated  
827 to include no more than 10 genes per cell type, in order to balance the number of genes coming  
828 from any individual cell type. In order to increase the differential expression effect sizes observed,  
829 only the most extreme 6 LCAM<sup>hi</sup> and LCAM<sup>lo</sup> patients processed with CD45+ bead enrichment  
830 and 10X Chromium V2 were included in the differential expression analysis. These patients were  
831 LCAM<sup>hi</sup>: (408, 403, 714, 522, 371, 570), and LCAM<sup>lo</sup>: (571, 596, 393, 593, 626, 378).

832

### 833 **Calculation of LCAM<sup>hi</sup>, LCAM<sup>lo</sup>, and ensemble LCAM scores in bulk-RNA datasets**

834 Bulk RNA expression datasets were log-transformed and z-scored. For each cell type associated  
835 with LCAM<sup>hi</sup> or LCAM<sup>lo</sup>, the resulting z-scores of the associated genes were averaged and z-  
836 scored. A summary average of these values was then computed across all the cell types associated  
837 with either LCAM<sup>hi</sup> or LCAM<sup>lo</sup> cell types.

838

### 839 **Published statistics for TCGA Lung adenocarcinoma patients**

840 Estimates of total immune content present in each TCGA sample (ESTIMATE score)<sup>35</sup> were  
841 download from <https://bioinformatics.mdanderson.org/public-software/estimate/>.

842 Scores associating mutational signatures<sup>41</sup> with individual TCGA samples were  
843 downloaded from the mSignatureDB<sup>42</sup> website <http://tardis.cgu.edu.tw/msignaturedb/>. For the present  
844 study, detection of Signature 4 was used to indicate presence of smoking-related mutations.

845 Counts of Indel Neoantigens, SNV Neoantigens, and CTA score in TCGA cases were  
846 accessed from Table S1 of ref. (<sup>53</sup>).

847

### 848 **Generation of stromal cell type scores**

849 Fibroblast and endothelial cell count matrices from tumor and nLung of 8 NSCLC patients<sup>5</sup> were  
850 downloaded from  
851 [https://gbiomed.kuleuven.be/english/research/50000622/laboratories/54213024/scRNAseq-](https://gbiomed.kuleuven.be/english/research/50000622/laboratories/54213024/scRNAseq-NSCLC)  
852 [NSCLC](https://gbiomed.kuleuven.be/english/research/50000622/laboratories/54213024/scRNAseq-NSCLC). Previously-applied<sup>5</sup> cluster annotations were assumed, where endothelial cluster 6 was  
853 defined as “lymphatics”, and endothelial and fibroblast clusters were defined based on enrichment  
854 in tumor or nLung: endothelial clusters 3 and 4 were pooled as “Tumor BEC”, endothelial clusters  
855 1 and 5 were pooled as “Normal BEC”, fibroblast cluster 1 was defined as “Normal fibroblast”,  
856 and fibroblast clusters 1, 2, 3, 4, 5, and 7 were pooled as “CAF”. For each of these cell types, gene  
857 scores were defined based on a minimum average expression of  $10^{-4}$  and a minimum fold-change  
858 threshold of 4 compared to any other stromal cell type. Cell type gene-scores were defined in  
859 TCGA lung adenocarcinoma using the average z-scored gene expression of each stromal gene list.

860

## 861 **DATA AVAILABILITY**

862 Human scRNAseq, TCRseq, and CITEseq data is available at GEO accession GSE154826.

863

## 864 **ACKNOWLEDGMENTS**

865 This work was supported by National Institutes of Health (NIH) grants 5T32CA078207 (to  
866 A.M.L.). We thank A. Magen, P. Hamon, M. Casanova-Acebes for critical comments on the  
867 manuscript; and the Mount Sinai flow cytometry core, Human Immune Monitoring Center and  
868 Mount Sinai Biorepository for support. Research reported in this paper was supported by the Office  
869 of Research Infrastructure of the National Institutes of Health under award numbers  
870 S10OD018522 and S10OD026880. The content is solely the responsibility of the authors and does  
871 not necessarily represent the official views of the National Institutes of Health. Data used in this



872 publication were generated by the TCGA Research Network, and the National Cancer Institute  
873 Clinical Proteomic Tumor Analysis Consortium (CPTAC). Research support was provided by  
874 Regeneron and Takeda. We recognize the patients and their families for their important  
875 contributions and sacrifices.

876

#### 877 **AUTHOR CONTRIBUTIONS**

878 MM and EK conceived the project. AML, AR, and MM designed the experiments. AML, EK, and  
879 MM wrote the manuscript. AML, EK, and MD performed computational analysis. TM, MB, AW,  
880 and RF facilitated access to human samples. AML, JG, CC, BM, AT, LW, JL, NM, GM, and KT  
881 performed experiments. NRD and GT funded part of the study. AL conducted patient consents and  
882 facilitated regulatory items. JG, JM, GM, ZZ, FP, RS, AK, PW, HS, and TM provided further  
883 intellectual input.

884

#### 885 **DECLARATION OF INTERESTS**

886 Research support for this work was provided by Regeneron and Takeda. The authors declare no  
887 other competing financial interests.

888

889 **FIGURE LEGENDS**

890

891 **Figure 1. scRNA- and CITE-seq establish the diversity of transcriptional states in the tumor**  
892 **microenvironment.**

893 **A**, Study overview. Resected specimens of tumor tissue and non-involved lung (nLung) were  
894 digested to single cell suspensions, enriched for CD45+ cells, and subjected to single cell assays  
895 including CITEseq and TCRseq.

896 **B**, Clinical data of patients undergoing resection indicating summary pathological stage, smoking  
897 history, histological diagnosis, and sex.

898 **C**, Expression of cell type marker genes across scRNAseq clusters of immune cells, grouped by  
899 lineage annotation (MNP: mononuclear phagocyte; pDC: plasmacytoid dendritic cell). Heatmap  
900 shows the number of unique molecular identifiers (UMI) per cell. Clusters are shown using an  
901 even number of randomly selected cells from 7 matched tumor and nLung sample pairs who were  
902 analyzed by CITEseq. Cells were downsampled to 2000 UMI/cell.

903 **D**, Expression of lineage-defining surface markers on single cells, as measured by CITEseq. Single  
904 cells correspond directly to cells shown in (C). CITEseq count values were first quantile  
905 normalized across patients, then row-normalized across cells in the heatmap.

906 **E**, Cells per cluster as a percent of total immune cells across 35 tumor and 32 matched nLung  
907 samples. Clusters correspond directly to those shown in (C) and (D).

908 **F**, Box-plots of Euclidean distances between pairs of samples among nLung only (nLung-nLung),  
909 tumor only (tumor-tumor), or between nLung samples and tumor samples (Tumor-nLung).

910 Distances between pairs of patient-matched samples were excluded from the Tumor-nLung

911 distribution to prevent confounding due to patient-specific effects. \*\*\*  $P < 0.001$ , Wilcoxon rank-  
912 sum test.

913 **G**, Log-ratios between cell type frequencies in tumor and nLung. Clusters were grouped by cell  
914 type annotation. Crosses represent error bars showing the mean  $\pm$  SEM of  $\text{Log}_2\text{FC}$  estimates of  
915 differences in cell type frequency between tumor and nLung using the cohort collected in the  
916 present study (Mount Sinai; x-axis) or the cohort in ref. <sup>5</sup> (y-axis).

917

918 **Figure 2. Intratumoral DC comprise expanded MoDC and express an LCH-like signature.**

919 **A**, Expression of key genes discriminating scRNAseq clusters of DC, grouped by cell type  
920 annotation (MoDC: monocyte-derived DC; cDC: classical DC). Heatmap shows the number of  
921 UMI per cell. Clusters are shown using an even number of randomly selected cells from each,  
922 drawing from patients who were analyzed by CITEseq with the DC panel shown in **(B)** (4 matched  
923 tumor-nLung tissue pairs). Cells were downsampled to 2000 UMI/cell.

924 **B**, Expression of DC surface markers on single cells, as measured by CITEseq. Single cells  
925 correspond directly to cells shown in **(A)**. CITEseq count values were first quantile normalized  
926 across patients, then row-normalized across cells in the heatmap.

927 **C**, Differences between tumor and nLung of DC frequencies normalized to total DC; \* $P < 0.05$ ,  
928 \*\* $P < 0.01$ , \*\*\* $P < 0.001$  (Wilcoxon signed-rank test with Bonferroni correction;  $N = 25$  matched  
929 tissue pairs with  $> 50$  DC observed in each tissue).

930 **D**, Barplots showing average expression of *LAMP3* (*DC-LAMP*) and *CD274* (*PD-L1*) across DC  
931 clusters.

932 **E**, MICSSS imaging showing spatial distribution of DC-LAMP<sup>+</sup>/PD-L1<sup>+</sup> DC in proximity to T  
933 cells in a TLS.

934 **F**, Expression of follicular dendritic cell marker MYH11 in TLS in an adjacent section to that  
935 shown in **(E)**.

936 **G**, Expression among CD14<sup>+</sup> monocytes and DC of monocyte, cDC2, and MΦ cell type specific  
937 gene signatures (See Figure S2D, E). Heatmaps show expression of 20 genes from each score  
938 among single-cells evenly sampled by cell type (left) and as corresponding summary scores. Cells  
939 were ordered by the ratio of monocyte:cDC2 summary scores and were downsampled to 2000  
940 UMI.

941 **H**, Boxplots showing average expression of LCH-like signature genes across DC populations in  
942 distinct nLung and tumor samples.

943

944 **Figure 3. Tumors exclude AMΦ and exhibit a diversity of MoMΦ populations.**

945 **A**, Average cluster expression of lineage-defining monocyte and MΦ clusters based on literature  
946 review, grouped by cell type annotation.

947 **B**, Expression of myeloid surface markers on single cells, as measured by CITEseq. Clusters are  
948 shown using an even number of randomly selected cells from each, from patients who were  
949 analyzed with the panel shown (4 matched tumor-nLung pairs). CITEseq count values were first  
950 quantile normalized across patients, then row-normalized across cells in the heatmap.

951 **C**, Histograms of gene module scores per cell type (see also Figure S3D-J).

952 **D-F**, Expression among CD14<sup>+</sup> monocytes, MoMΦ, and AMΦ of cell type specific gene scores.  
953 Gene scores were generated based on sets of mutually exclusive, differentially expressed genes  
954 among AMΦ, MoMΦ, and CD14<sup>+</sup> monocytes (see Figure S3I). Cells are plotted by AMΦ and  
955 MoMΦ score, and cell-annotations are indicated by colored dots or contour plots **(D)**. Cells are

956 plotted on similar axes and colored by CD14<sup>+</sup> monocyte score (**E**), or by expression of individual  
957 genes (**F**).

958 **G**, Differences between tumor and nLung of lineage-normalized monocyte and MΦ frequencies;  
959 \*P<0.05; \*\*P<0.01, \*\*\*P,0.001 (Wilcoxon signed-rank test with Bonferroni correction; N=32  
960 matched tissue pairs).

961 **H**, Average cell type expression of secreted factors across MNP cell types.

962

963 **Figure 4. CITEseq and TCR analysis of the adaptive immune compartment.**

964 **A**, Expression of key genes discriminating scRNAseq clusters of T cells, grouped by cell type  
965 annotation. Heatmap shows the number of UMI per cell. Clusters are shown using an even number  
966 of randomly selected cells from each, drawing from patients who were analyzed by CITEseq with  
967 the T cell panel shown in (**B**) (2 matched tumor-nLung tissue pairs). Cells were downsampled to  
968 2000 UMI/cell.

969 **B**, Expression of T cell surface markers on single cells, as measured by CITEseq. Single cells  
970 correspond directly to cells shown in (**A**). CITEseq count values were first quantile normalized  
971 across patients, then row-normalized across cells in the heatmap.

972 **C**, Differences between tumor and nLung of population frequencies normalized by total NK and  
973 T cells; \*P<0.05; \*\*P<0.01, \*\*\*P<0.001 (Wilcoxon signed-rank test with Bonferroni correction,  
974 N=32 matched tissue pairs).

975 **D, E**, Phenotypic distribution of T cells among tissue-stratified clonotypes. Frequencies of unique  
976 TCRs observed by scTCRseq in nLung (x-axis) or tumor in a representative patient (**D**). In (**E**),  
977 cells were first grouped by TCR tissue tropism categories as defined in (**D**); for 3 patients, the  
978 phenotypic makeup of the cells with unique TCRs, tissue-specific TCRs, or TCRs shared across

979 tissues is plotted for nLung (i) and tumor tissues (ii) is plotted as a percent of cells with similarly  
980 tissue-distributed TCRs. Each patient is indicated by shape.

981 **F**, Expression of key genes discriminating scRNAseq clusters of B and plasma cells, grouped by  
982 cell type annotation. Heatmap shows the number of UMI per cell. Clusters are shown using an  
983 even number of randomly selected cells from each, drawing from patients who were analyzed by  
984 CITEseq with the B cell panel shown in (G) (4 matched tumor-nLung tissue pairs). Cells were  
985 downsampled to 2000 UMI/cell.

986 **G**, Expression of B and plasma cell surface markers on single cells, as measured by CITEseq.  
987 Single cells correspond directly to cells shown in (F). CITEseq count values were first quantile  
988 normalized across patients, then row-normalized across cells in the heatmap.

989

990 **Figure 5. Cell-cell interactions drive an axis of adaptive activation.**

991 **A**, Spearman correlation of cell type frequencies after normalization within lineage. Analysis  
992 includes 23 tumors that were processed similarly using 10X Chromium V2 and CD45+ magnetic  
993 bead enrichment.

994 **B**, Lineage-normalized cell type frequencies of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> cell types among pooled  
995 nLung and Tumor samples from Mount Sinai and refs. (<sup>5,6</sup>) (50 tumor patients with 40 matched  
996 nLung samples). nLung samples are ordered to match the order of tumor samples based on  
997 frequencies of LCAM celltypes.

998 **C**, Immune lineage frequencies of nLung and Tumor samples; with columns corresponding to  
999 patient ordering in (B).

1000 **D-I**, Log2 Ratio of ligand-receptor (LR) intensity scores between tumor and nLung of LCAM<sup>hi</sup>  
1001 patients, (“LR ratio”; y-axis) and LCAM<sup>lo</sup> patients (x-axis). All interactions among T cells, B cells,

1002  $M\Phi$ , MoDC, cDC, and monocytes, colored by indication of significance (permutation test, **D**).

1003 Dashed diagonal line indicates unity.

1004 **E-I**, Showing same data as in (**D**), but highlighting in bold LR ratios for interactions between T

1005 cell ligands and B cell receptors (**E**), T cell ligands and cDC receptors (**F**),  $M\Phi$  ligands and T cell

1006 receptors (**G**), MoDC ligands and T cell receptors (**H**), and cDC ligands and T cell receptors (**I**).

1007 Labelled interactions are plotted in red.

1008

1009 **Figure 6. Tumor features related to the LCAM immune response.**

1010 **A**, Normalized expression of  $LCAM^{hi}$  and  $LCAM^{lo}$  bulk-RNA signature genes, determined as

1011 shown in Figure S6A, B and as described in the methods, in TCGA lung adenocarcinoma

1012 dataset. Cell type association with sets of genes for each signature is shown. Patients are sorted

1013 along y-axis by ensemble LCAM score.

1014 **B-D**, Scatter plots of the ensemble LCAM score (y-axis) with signature scores based on genes

1015 that are specific for CAFs (**B**), normal fibroblasts (**C**), or the difference between these scores (**D**)

1016 in TCGA lung adenocarcinoma data. Stromal signatures are based on the stromal data reported in

1017 ref. <sup>5</sup>.

1018 **E**, Scatter plot of LogTMB and ensemble LCAM score. Patients are divided into those with

1019 presence of a smoking-related mutational signature (black) and those without presence of the

1020 signature (red). Black and red lines indicate linear regression relationships computed over each

1021 group of patients independently ( $r=0.38$ ;  $p=9.2e-5$  in the undetected smoking signature group;

1022  $r=0.34$ ;  $p=1.1e-12$  in the detected signature group).

1023 **F**, Scatter plot of Cancer testes antigen expression score (CTA score), as computed in ref. <sup>53</sup>, and

1024 the residuals of the regression of the ensemble LCAM score on the LogTMB.

1025 **G and H**, Boxplots showing either the ensemble LCAM score (**G**), or TMB (**H**) among TCGA  
1026 lung adenocarcinoma patients, divided by combinations of driver mutations.

1027 **I and J**, Histograms of residuals of the regression of the ensemble LCAM score on the LogTMB,  
1028 with patients stratified by *TP53* (**I**) or *KRAS* (**J**) mutational status (Two-sided t-test).

1029

1030 **Figure S1. Integration of scRNA samples and datasets for common cell type analysis.**

1031 **A**, Comparison of per-sample estimated noise levels in the training set of cells used for clustering  
1032 and model formation (x-axis) compared to the per-sample estimated noise in a withheld test set of  
1033 cells that were mapped to the model clusters by probabilistic projection.

1034 **B and C**, Illustration of how incorporating a fit noise component improves the concordance  
1035 between predicted expression and of cells mapped to the  $T_{reg}$  cluster and observed expression. Y-  
1036 axis shows the predicted expression of  $T_{regs}$  in individual samples without accounting for noise (**B**)  
1037 or accounting for noise (**C**), against the observed average expression (X-axis). Genes were color-  
1038 coded by the ratio between the observed expression and the model without accounting for noise.  
1039 Estimation of the noise component is detailed in the methods.

1040 **D**, Per-sample estimated noise levels in 10X chromium V2 samples that were used for clustering,  
1041 10X chromium V2 samples that were analyzed by projection onto the clustering model and not  
1042 used in the clustering, 10X chromium 5' samples that were analyzed by projection, and samples  
1043 from external datasets<sup>5,6</sup> that were analyzed by projection.

1044 **E**, Boxplots showing the distribution of UMI per cell in each cluster.

1045 **F**, Barplots showing number of cells in Mount Sinai dataset mapping to each cluster.



1046 **G**, Heatmap showing row-normalized cell type frequencies in a public dataset<sup>5</sup> with samples  
1047 spanning 3 regions each in a cohort of 8 NSCLC patients. Samples are clustered by spearman  
1048 correlation distance. Sample names are colored by patient.

1049 **H**, Box plots of Euclidean distances based on log-transformed cluster frequencies between samples  
1050 of different patients or from the same patient, as in (**G**), from ref. <sup>5</sup>. \*\*\* P < 0.001, Wilcoxon rank-  
1051 sum test.

1052

1053 **Figure S2. Module analysis of DC.**

1054 **A**, Barplots showing total number of cells mapped to each individual DC cluster in the Mount  
1055 Sinai cohort.

1056 **B**, Boxplots showing number of cells mapped to each individual DC cluster per tumor sample in  
1057 the Mount Sinai cohort.

1058 **C**, Differences between tumor and nLung of DC frequencies normalized to total MNP; \*P<0.05,  
1059 \*\*P<0.01, \*\*\*P,0.001 (Wilcoxon signed-rank test with Bonferroni correction; N=26 matched  
1060 tissue pairs with >250 MNP observed in each tissue).

1061 **D**, Log<sub>2</sub>FC and expression level distributions of gene sets that are mutually exclusively expressed  
1062 in CD14+ monocytes, MΦ, and cDC2 (See Figure 2G).

1063 **E**. Histograms of cDC2 and MΦ scores, using gene lists generated as shown in (**D**).

1064 **F-I**, Gene module analysis of DC clusters. Correlation of gene module expression across all DC  
1065 (**F**), five example genes from each module, ranked by correlation to the other genes in the module  
1066 and colored by total expression in DC (**G**), boxplots showing Log<sub>2</sub>FC of module expression among  
1067 all DC between patient matched tumor and nLung samples (**H**), and normalized average cluster  
1068 expression of modules (**I**).

1069 **J**, Boxplots showing expression of LCH-like signature genes across DC populations in distinct  
1070 nLung and tumor samples from ref. <sup>5</sup>.

1071

1072 **Figure S3. Diversity of nlung and tumor-infiltrating MΦ populations.**

1073 **A**, Expression of key genes discriminating scRNAseq clusters of monocytes and MΦ, grouped by  
1074 cell type annotation. Heatmap shows the number of UMI per cell. Clusters are shown using an  
1075 even number of randomly selected cells from each, drawing from 35 tumor and 32 nLung samples.  
1076 Cells were downsampled to 2000 UMI/cell.

1077 **B**, Scatter plots showing normalized CITEseq CD10 and CD206 surface marker counts on AMΦ,  
1078 MoMΦ, and CD14+ monocytes in nLung of a representative patient.

1079 **C**, IHC of CD10 staining AMΦ in the airspaces of nLung tissue.

1080 **D-G**, Gene module analysis of monocyte and MΦ clusters. Correlation of gene module expression  
1081 across all monocytes and MΦ (**D**). Module groups illustrate groups of correlated modules which  
1082 are expressed most specifically on AMΦ, MoMΦ, and monocytes (see **G**). Five example genes  
1083 from each module, ranked by correlation to the other genes in the module and colored by total  
1084 expression in monocytes and MΦ (**E**), boxplots showing Log<sub>2</sub>FC of module expression among all  
1085 monocytes and MΦ between patient matched tumor and nLung samples (**F**), and normalized  
1086 average cluster expression of modules (**G**).

1087 **H**, Log<sub>2</sub>FC and expression level of gene sets that are mutually exclusively expressed in CD14+  
1088 monocytes, AMΦ, and MoMΦ (See Figure 3D-F).

1089

1090 **Figure S4. Phenotypic dissection of activated, cycling, and clonally expanded T cells.**

1091 **A**, Differential expression within the  $T_{\text{activated}}$  cluster of cells staining for CD4 versus CD8 by  
1092 CITEseq (y-axis) vs. average  $T_{\text{activated}}$  expression (x-axis).

1093 **B**, Classification of  $T_{\text{activated}}$  cells as CD4+ or CD8+ based on the ratio of CD4-related or CD8-  
1094 related gene signatures learned from cells of 2 patients (training set; open circles) and validated on  
1095 cells of 4 additional validation patients (test set; black dots). Red line indicates gene ratio threshold  
1096 learned from the training set. Only cells where the CITEseq CD4:CD8 count ratio is  $>2$  or  $<1/2$   
1097 are considered.

1098 **C**, Validation of CD4/CD8 classification scheme shown in **(B)** for cells without CITEseq staining.  
1099 Cells were considered to be CD4+ or CD8+ based on unique RNA detection of either *CD4* (blue  
1100 points) or at least one *CD8A* or *CD8B* transcript (green points). The discriminant line is equivalent  
1101 to the gene ratio threshold learned from CITEseq data, shown in **(B)**.

1102 **D**, Expression of key genes in CD4-related and CD8-related gene signatures for discriminating  
1103 CD4+ and CD8+ activated T cells. Cells are sorted by ratio of these gene signatures, and the line  
1104 is drawn to indicate the cells discriminated based on the threshold in panel **(B)**. Heatmap shows  
1105 the number of UMI per cell. Cells represent  $T_{\text{activated}}$  cells from 35 tumors, and were downsampled  
1106 to 2000 UMI/cell.

1107 **E**, Frequency of CD8+ or CD4+  $T_{\text{activated}}$  cells across 35 patients, as determined by gene signature  
1108 scores learned from CITEseq (as in **A-D**).

1109 **F-H**, Spherical k-means sub-clustering on cell type scores of cells within the cycling T cell cluster  
1110 18 based on gene scores generated from other T cell clusters. Heatmap of single-cell expression of  
1111 cell type scores, grouped by sub-cluster **(F)**, number of cells in each sub-cluster **(G)**; nLung shown  
1112 in blue, tumor in brown; lines dividing bars horizontally discriminate groups of cells from distinct

1113 patients), and the frequency of cycling T cells of each T cell phenotype (**H**; data points represent  
1114 samples with at least 50 cells of the given phenotype).

1115 **I**, TCR clonality score of phenotypic groups in nLung (blue) and tumor (brown). Dots represent  
1116 individual samples with at least 30 cells of indicated phenotype. N=3 patients with tumor-nLung  
1117 pairs.

1118 **J**, Number of cells within each TCR category, determined as in Figure 5D, in matched nLung and  
1119 tumor samples of 3 patients, each patient indicated by shape.

1120 **K**, Number of unique TCRs represented in each TCR category, determined as in Figure 5D.

1121 **L**, Differences between tumor and nLung of lineage-normalized B and plasma cell type  
1122 frequencies. All comparisons were not significant ( $P > 0.05$ , Wilcoxon signed-rank test, N=32  
1123 matched tissue pairs).

1124

1125 **Figure S5. Ligand-receptor interactions in LCAM<sup>hi</sup> and LCAM<sup>lo</sup> tumors.**

1126 **A**, Lineage-normalized cell type frequencies of all cell types among pooled nLung and Tumor  
1127 samples from Mount Sinai and refs. <sup>5,6</sup> (50 tumor patients with 40 matched nLung samples).

1128 **B and C**, Column-normalized expression of highly expressed secreted ligands (**B**) and associated  
1129 receptors (**C**) across all immune cell types, connected by lines linking ligands to receptors.  
1130 Connectors are colored by association with LCAM<sup>hi</sup> patients (purple), LCAM<sup>lo</sup> patients (green),  
1131 or all tumors (orange).

1132 **D-F**, Log<sub>2</sub> Ratio of ligand-receptor (LR) intensity scores between tumor and nLung of LCAM<sup>hi</sup>  
1133 patients, (“LR ratio”; y-axis) and LCAM<sup>lo</sup> patients (x-axis) as in Figure 5D, highlighting in bold  
1134 LR ratios for interactions between B cell ligands and T cell receptors (**D**), T cell ligands and MoDC  
1135 receptors (**E**), and T cell ligands and MΦ receptors (**F**). Labelled interactions are plotted in red.

1136

1137 **Figure S6. Projection of bulk RNA samples onto signatures defined by the LCAM scRNA**

1138 **axis.**

1139 **A and B**, Derivation of the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> gene signatures for scoring bulk RNA samples.

1140 Identification of differentially-expressed genes between averaged scRNAseq samples of LCAM<sup>hi</sup>

1141 and LCAM<sup>lo</sup> patients (**A**), and identification of differentially expressed genes that are specific to

1142 genes in the LCAM<sup>hi</sup> or LCAM<sup>lo</sup> cell types (**B**).

1143 **C-E**, Scatter plots of immune ESTIMATE score<sup>35</sup> with the LCAM<sup>hi</sup> signature score (**C**), the

1144 LCAM<sup>lo</sup> signature score (**D**), or the difference between the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> signature scores

1145 (i.e. the ensemble LCAM score; **E**).

1146 **F**, Spearman correlation of the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> signature scores among the deciles of

1147 immune content. Error bars represent the 95% confidence interval around the estimate of the

1148 spearman correlation.

1149 **G**, Scatter plots of the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> signature scores, showing the 1<sup>st</sup> (black), 3<sup>rd</sup> (green),

1150 and 10<sup>th</sup> (red) deciles of immune content. Labelled trend lines are shown for other deciles.

1151 **H**, Normalized expression of LCAM<sup>hi</sup> and LCAM<sup>lo</sup> bulk-RNA signature genes in the CPTAC lung

1152 adenocarcinoma dataset.

1153 **I**, Scatter plots of the ensemble LCAM score (y-axis) with signature scores based on genes that

1154 are specific for tumor blood endothelial cells (BEC; left), normal BEC (center), and lymphatic

1155 endothelial cells. Stromal signatures are based on the stromal data reported in ref. <sup>5</sup>.

1156 **J**, Boxplots showing ensemble LCAM scores among TCGA lung adenocarcinoma patients by

1157 TNM T-stage.

1158 **K**, Scatter plot of LogTMB and ensemble LCAM score in CTPAC lung adenocarcinoma patients,  
1159 with linear regression line.

1160 **L and M**, Scatter plots of the LogTMB and immune ESTIMATE score<sup>35</sup> (**L**) and the T-cell  
1161 inflamed gene expression profile (GEP<sup>39,40</sup>; **M**) in TCGA lung adenocarcinoma patients.

1162 **N**, Correlation between individual genes comprising the LCAM<sup>hi</sup> and LCAM<sup>lo</sup> bulk gene  
1163 signatures and LogTMB in TCGA lung adenocarcinoma patients.

1164 **O**, Scatter plots of LogTMB and the LCAM ensemble score for patients by T-stage.

1165 **P and Q**, Scatter plots of the number of indel-induced neoantigens (**P**) and SNV-induced  
1166 neoantigens (**Q**) as computed in ref. <sup>53</sup>, and the residuals of the regression of the ensemble LCAM  
1167 score on the LogTMB.

1168 **R and S**, Boxplots showing either the ensemble LCAM score (**G**), or TMB (**H**) among CPTAC  
1169 lung adenocarcinoma patients, divided by combinations of mutated driver mutations.

1170 **T and U**, Histograms of residuals of the regression of the ensemble LCAM score on the LogTMB,  
1171 with patients stratified by *STK11* (**T**) or *EGFR* (**U**) mutational status (Two-sided t-test).

1172

## 1173 **SUPPLEMENTAL TABLES**

1174 **Table S1.** Sample table, with information about patient, tissue, 10X loading, and QC metrics

1175 **Table S2.** CITEseq panels used

1176 **Table S3.** QC table of GEX, HTO, and ADT libraries

1177 **Table S4.** Gene lists used in paper

1178 **Table S5.** Gene modules

1179 **Table S6.** Ligand-receptor pairs used in the analysis

1180 **Table S7.** Ligand-receptor statistics

1181

## 1182 REFERENCES

- 1183 1. Siegel, R.L., Miller, K.D. & Jemal, A. Cancer statistics, 2019. *CA Cancer J Clin* **69**, 7-34 (2019).
- 1184 2. Travis, W.D., Brambilla, E., Burke, A.P., Marx, A. & Nicholson, A.G. Introduction to The 2015  
1185 World Health Organization Classification of Tumors of the Lung, Pleura, Thymus, and Heart. *J*  
1186 *Thorac Oncol* **10**, 1240-1242 (2015).
- 1187 3. Remon, J., *et al.* Immune Checkpoint Inhibitors in Thoracic Malignancies: Review of the Existing  
1188 Evidence by an IASLC Expert Panel and Recommendations. *Journal of Thoracic Oncology* (2020).
- 1189 4. Laughney, A.M., *et al.* Regenerative lineages and immune-mediated pruning in lung cancer  
1190 metastasis. *Nat Med* **26**, 259-269 (2020).
- 1191 5. Lambrechts, D., *et al.* Phenotype molding of stromal cells in the lung tumor microenvironment.  
1192 *Nat Med* **24**, 1277-1289 (2018).
- 1193 6. Zilionis, R., *et al.* Single-Cell Transcriptomics of Human and Mouse Lung Cancers Reveals  
1194 Conserved Myeloid Populations across Individuals and Species. *Immunity* **50**, 1317-1334 e1310  
1195 (2019).
- 1196 7. Lavin, Y., *et al.* Innate Immune Landscape in Early Lung Adenocarcinoma by Paired Single-Cell  
1197 Analyses. *Cell* **169**, 750-765.e717 (2017).
- 1198 8. Kargl, J., *et al.* Neutrophils dominate the immune cell composition in non-small cell lung cancer.  
1199 *Nature Communications* **8**, 14381 (2017).
- 1200 9. Guo, X., *et al.* Global characterization of T cells in non-small-cell lung cancer by single-cell  
1201 sequencing. *Nature Medicine* **24**, 978-985 (2018).
- 1202 10. Thommen, D.S., *et al.* A transcriptionally and functionally distinct PD-1+ CD8+ T cell pool with  
1203 predictive potential in non-small-cell lung cancer treated with PD-1 blockade. *Nature Medicine* **24**,  
1204 994-1004 (2018).
- 1205 11. Li, H., *et al.* Dysfunctional CD8 T Cells Form a Proliferative, Dynamically Regulated  
1206 Compartment within Human Melanoma. *Cell* **176**, 775-789 e718 (2019).
- 1207 12. Wu, T.D., *et al.* Peripheral T cell expansion predicts tumour infiltration and clinical response.  
1208 *Nature* **579**, 274-278 (2020).
- 1209 13. Binnewies, M., *et al.* Understanding the tumor immune microenvironment (TIME) for effective  
1210 therapy. *Nat Med* **24**, 541-550 (2018).
- 1211 14. Rizvi, N.A., *et al.* Mutational landscape determines sensitivity to PD-1 blockade in non-small cell  
1212 lung cancer. *Science* **348**, 124 (2015).
- 1213 15. Stoeckius, M., *et al.* Simultaneous epitope and transcriptome measurement in single cells. *Nature*  
1214 *methods* **14**, 865-868 (2017).

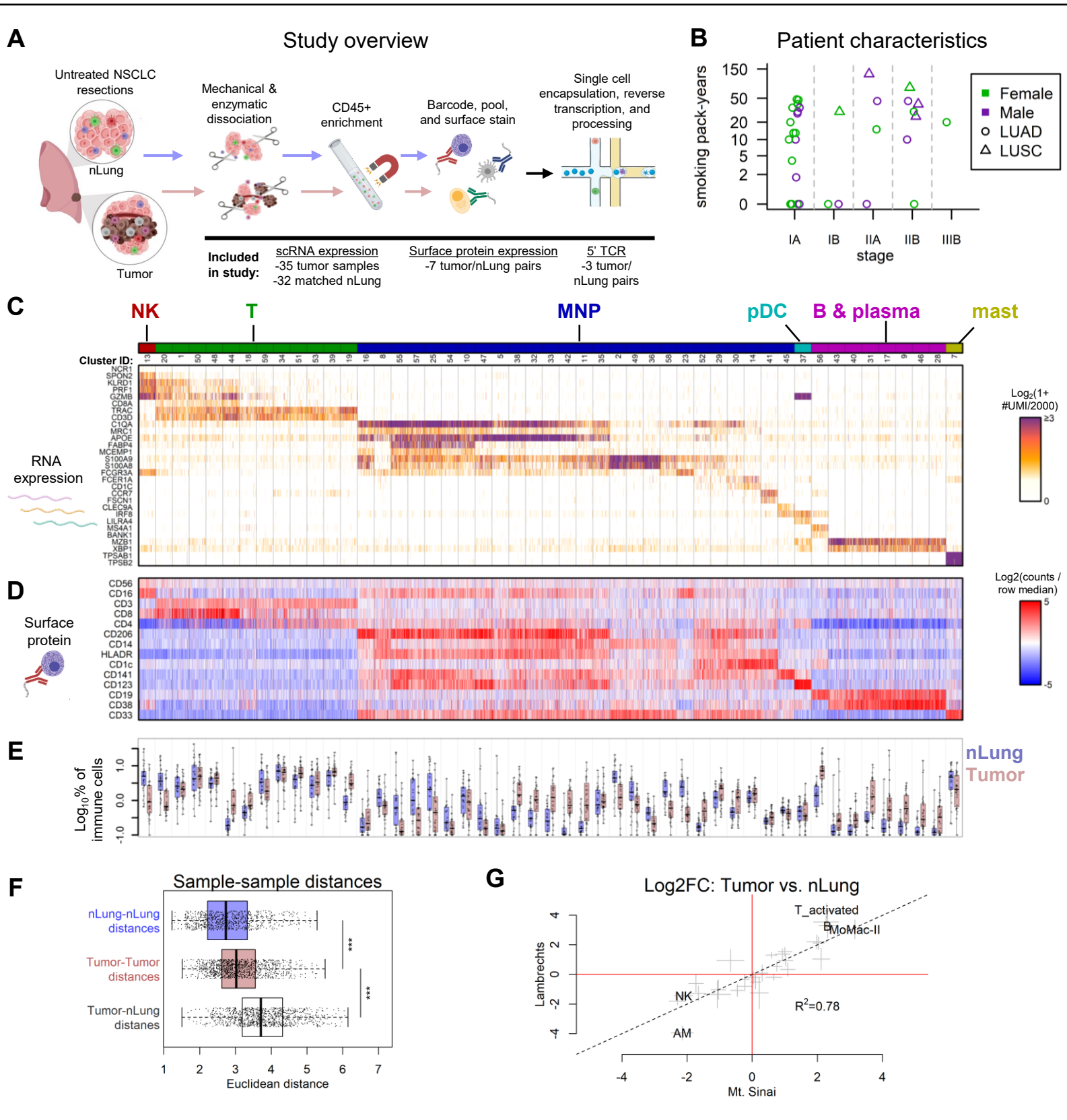
- 1215 16. Martin, J.C., *et al.* Single-Cell Analysis of Crohn's Disease Lesions Identifies a Pathogenic Cellular  
1216 Module Associated with Resistance to Anti-TNF Therapy. *Cell* **178**, 1493-1508.e1420 (2019).
- 1217 17. Rosenthal, R., *et al.* Neoantigen-directed immune escape in lung cancer evolution. *Nature* **567**,  
1218 479-485 (2019).
- 1219 18. Maier, B., *et al.* A conserved dendritic-cell regulatory program limits antitumour immunity. *Nature*  
1220 **580**, 257-262 (2020).
- 1221 19. Zhang, Q., *et al.* Landscape and Dynamics of Single Immune Cells in Hepatocellular Carcinoma.  
1222 *Cell* **179**, 829-845.e820 (2019).
- 1223 20. See, P., *et al.* Mapping the human DC lineage through the integration of high-dimensional  
1224 techniques. *Science* **356**(2017).
- 1225 21. Dutertre, C.A., *et al.* Single-Cell Analysis of Human Mononuclear Phagocytes Reveals Subset-  
1226 Defining Markers and Identifies Circulating Inflammatory Dendritic Cells. *Immunity* **51**, 573-589  
1227 e578 (2019).
- 1228 22. Remark, R., *et al.* In-depth tissue profiling using multiplexed immunohistochemical consecutive  
1229 staining on single slide. *Sci Immunol* **1**, aaf6925 (2016).
- 1230 23. Ioannidis, I. & Laurini, J.A. Use of Smooth Muscle Myosin Heavy Chain as an Effective Marker  
1231 of Follicular Dendritic Cells. *Appl Immunohistochem Mol Morphol* **27**, 48-53 (2019).
- 1232 24. Allen, C.E., Merad, M. & McClain, K.L. Langerhans-Cell Histiocytosis. *New England Journal of*  
1233 *Medicine* **379**, 856-868 (2018).
- 1234 25. Senechal, B., *et al.* Expansion of regulatory T cells in patients with Langerhans cell histiocytosis.  
1235 *PLoS Med* **4**, e253 (2007).
- 1236 26. Dudakov, J.A., Hanash, A.M. & van den Brink, M.R. Interleukin-22: immunobiology and  
1237 pathology. *Annu Rev Immunol* **33**, 747-785 (2015).
- 1238 27. Khosravi, N., *et al.* IL22 Promotes Kras Mutant Lung Cancer by Induction of a Pro-Tumor Immune  
1239 Response and Protection of Stemness Properties. *Cancer Immunology Research*,  
1240 canimm.0655.2017 (2018).
- 1241 28. Chakarov, S., *et al.* Two distinct interstitial macrophage populations coexist across tissues in  
1242 specific subtissular niches. *Science* **363**(2019).
- 1243 29. Leach, S.M., *et al.* Human and mouse transcriptome profiling identifies cross-species homology in  
1244 pulmonary and lymph node mononuclear phagocytes. *bioRxiv*, 2020.2004.2030.070839 (2020).
- 1245 30. Hashimoto, D., *et al.* Tissue-Resident Macrophages Self-Maintain Locally throughout Adult Life  
1246 with Minimal Contribution from Circulating Monocytes. *Immunity* **38**, 792-804 (2013).
- 1247 31. Gibbings, S.L., *et al.* Three Unique Interstitial Macrophages in the Murine Lung at Steady State.  
1248 *Am J Respir Cell Mol Biol* **57**, 66-76 (2017).
- 1249 32. Singh, D., *et al.* CD4<sup>+</sup> follicular helper-like T cells are key players in anti-tumor immunity. *bioRxiv*  
1250 (2020).



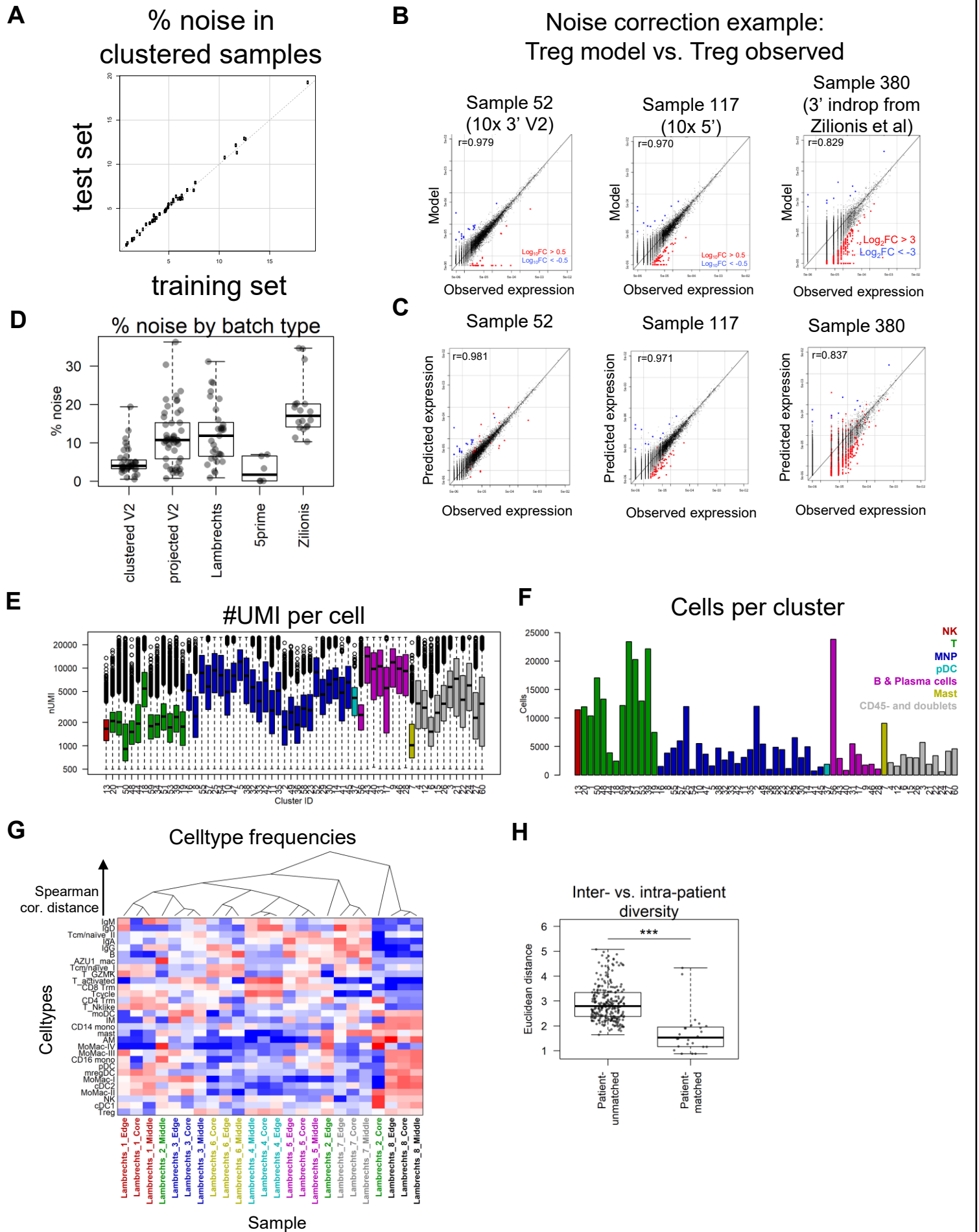
- 1251 33. Ramilowski, J.A., *et al.* A draft network of ligand–receptor-mediated multicellular signalling in  
1252 human. *Nature Communications* **6**, 7866 (2015).
- 1253 34. Neftel, C., *et al.* An Integrative Model of Cellular States, Plasticity, and Genetics for Glioblastoma.  
1254 *Cell* **178**, 835-849.e821 (2019).
- 1255 35. Yoshihara, K., *et al.* Inferring tumour purity and stromal and immune cell admixture from  
1256 expression data. *Nature Communications* **4**, 2612 (2013).
- 1257 36. Salmon, H., Remark, R., Gnjatic, S. & Merad, M. Host tissue determinants of tumour immunity.  
1258 *Nat Rev Cancer* **19**, 215-227 (2019).
- 1259 37. Salmon, H., *et al.* Matrix architecture defines the preferential localization and migration of T cells  
1260 into the stroma of human lung tumors. *J Clin Invest* **122**, 899-910 (2012).
- 1261 38. Samstein, R.M., *et al.* Tumor mutational load predicts survival after immunotherapy across  
1262 multiple cancer types. *Nat Genet* **51**, 202-206 (2019).
- 1263 39. Ayers, M., *et al.* IFN- $\gamma$ -related mRNA profile predicts clinical response to PD-1 blockade. *The*  
1264 *Journal of Clinical Investigation* **127**, 2930-2940 (2017).
- 1265 40. Cristescu, R., *et al.* Pan-tumor genomic biomarkers for PD-1 checkpoint blockade-based  
1266 immunotherapy. *Science* **362**, eaar3593 (2018).
- 1267 41. Alexandrov, L.B., *et al.* Mutational signatures associated with tobacco smoking in human cancer.  
1268 *Science* **354**, 618 (2016).
- 1269 42. Huang, P.-J., *et al.* mSignatureDB: a database for deciphering mutational signatures in human  
1270 cancers. *Nucleic Acids Research* **46**, D964-D970 (2017).
- 1271 43. Dong, Z.-Y., *et al.* Potential Predictive Value of TP53 and KRAS Mutation Status for Response to  
1272 PD-1 Blockade Immunotherapy in Lung Adenocarcinoma. *Clinical Cancer Research*,  
1273 clincanres.2554.2016 (2017).
- 1274 44. Biton, J., *et al.* TP53, STK11, and EGFR mutations predict tumor immune profile and the response  
1275 to anti-PD-1 in lung adenocarcinoma. *Clinical Cancer Research* **24**, 5710-5723 (2018).
- 1276 45. Rooney, Michael S., Shukla, Sachet A., Wu, Catherine J., Getz, G. & Hacohen, N. Molecular and  
1277 Genetic Properties of Tumors Associated with Local Immune Cytolytic Activity. *Cell* **160**, 48-61  
1278 (2015).
- 1279 46. Skoulidis, F., *et al.* STK11/LKB1 mutations and PD-1 inhibitor resistance in KRAS-mutant lung  
1280 adenocarcinoma. *Cancer discovery* **8**, 822-835 (2018).
- 1281 47. Canon, J., *et al.* The clinical KRAS(G12C) inhibitor AMG 510 drives anti-tumour immunity.  
1282 *Nature* **575**, 217-223 (2019).
- 1283 48. Stoeckius, M., *et al.* Cell Hashing with barcoded antibodies enables multiplexing and doublet  
1284 detection for single cell genomics. *Genome Biol* **19**, 224 (2018).
- 1285 49. Paul, F., *et al.* Transcriptional Heterogeneity and Lineage Commitment in Myeloid Progenitors.  
1286 *Cell* **163**, 1663-1677 (2015).

- 1287 50. Jaitin, D.A., *et al.* Massively parallel single-cell RNA-seq for marker-free decomposition of tissues  
1288 into cell types. *Science* **343**, 776-779 (2014).
- 1289 51. Ding, J., *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods.  
1290 *Nature Biotechnology* (2020).
- 1291 52. Kirsch, I., Vignali, M. & Robins, H. T-cell receptor profiling in cancer. *Mol Oncol* **9**, 2063-2070  
1292 (2015).
- 1293 53. Thorsson, V., *et al.* The Immune Landscape of Cancer. *Immunity* **48**, 812-830.e814 (2018).
- 1294

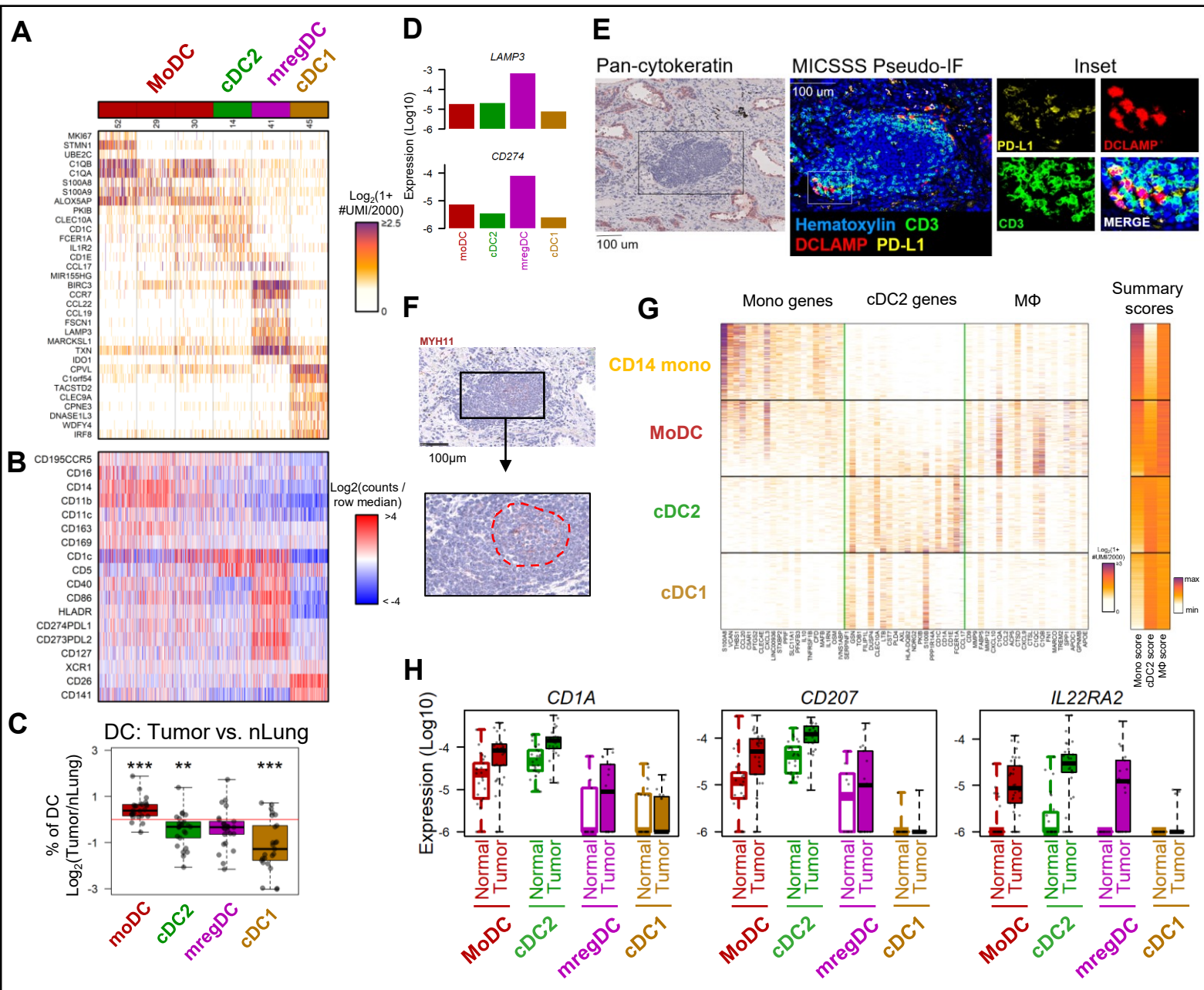
# Figure 1



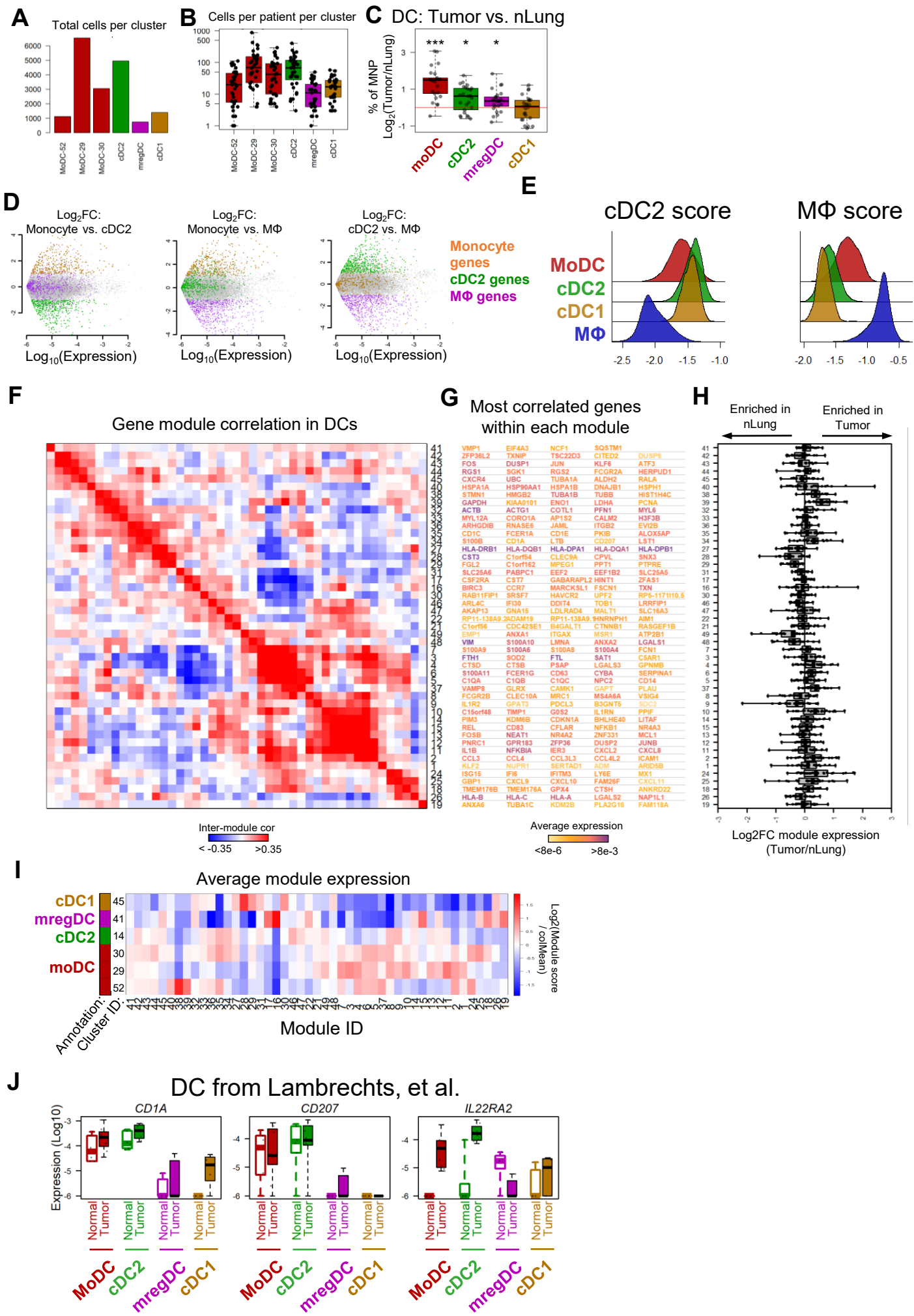
# Figure S1



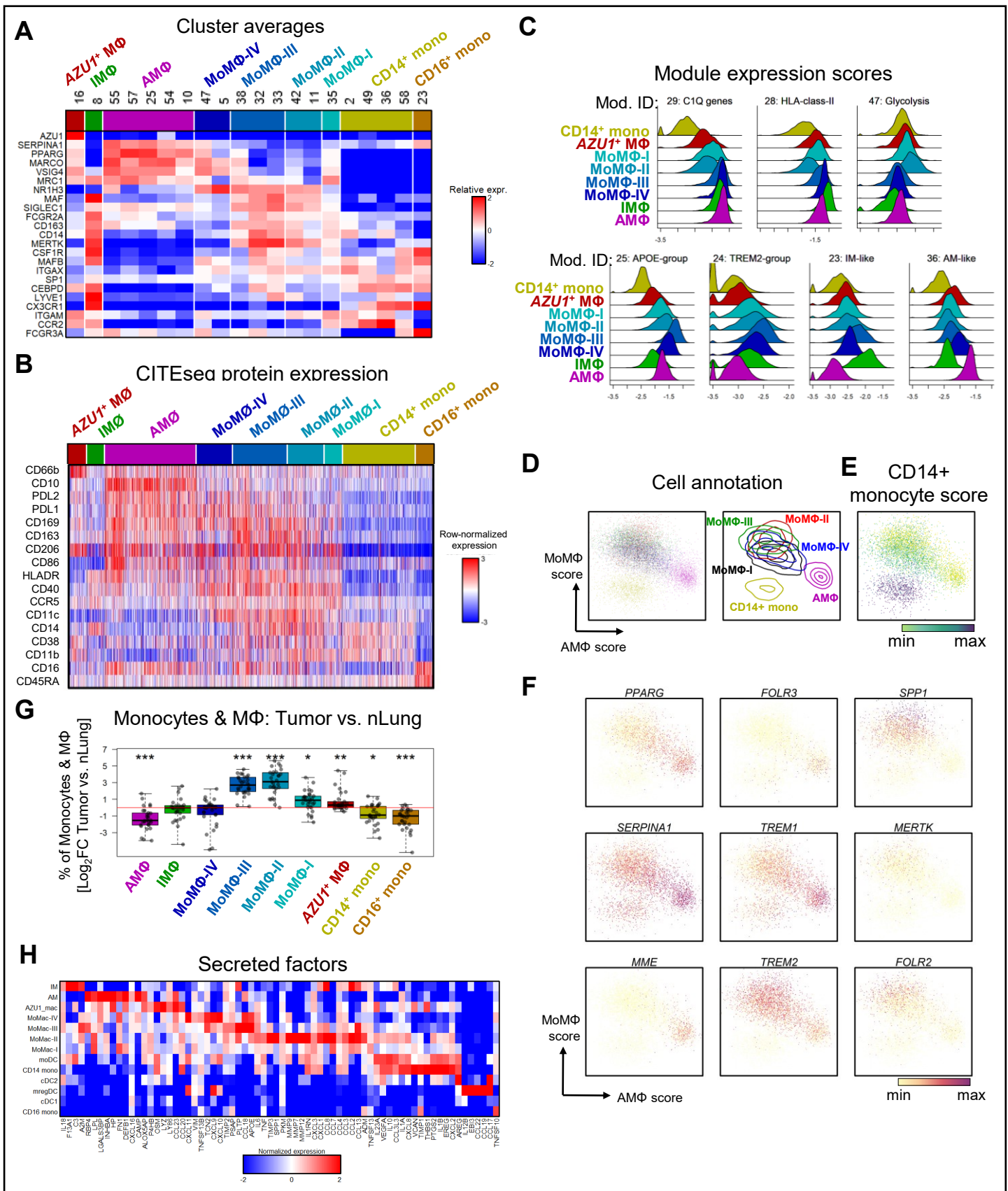
## Figure 2



# Figure S2

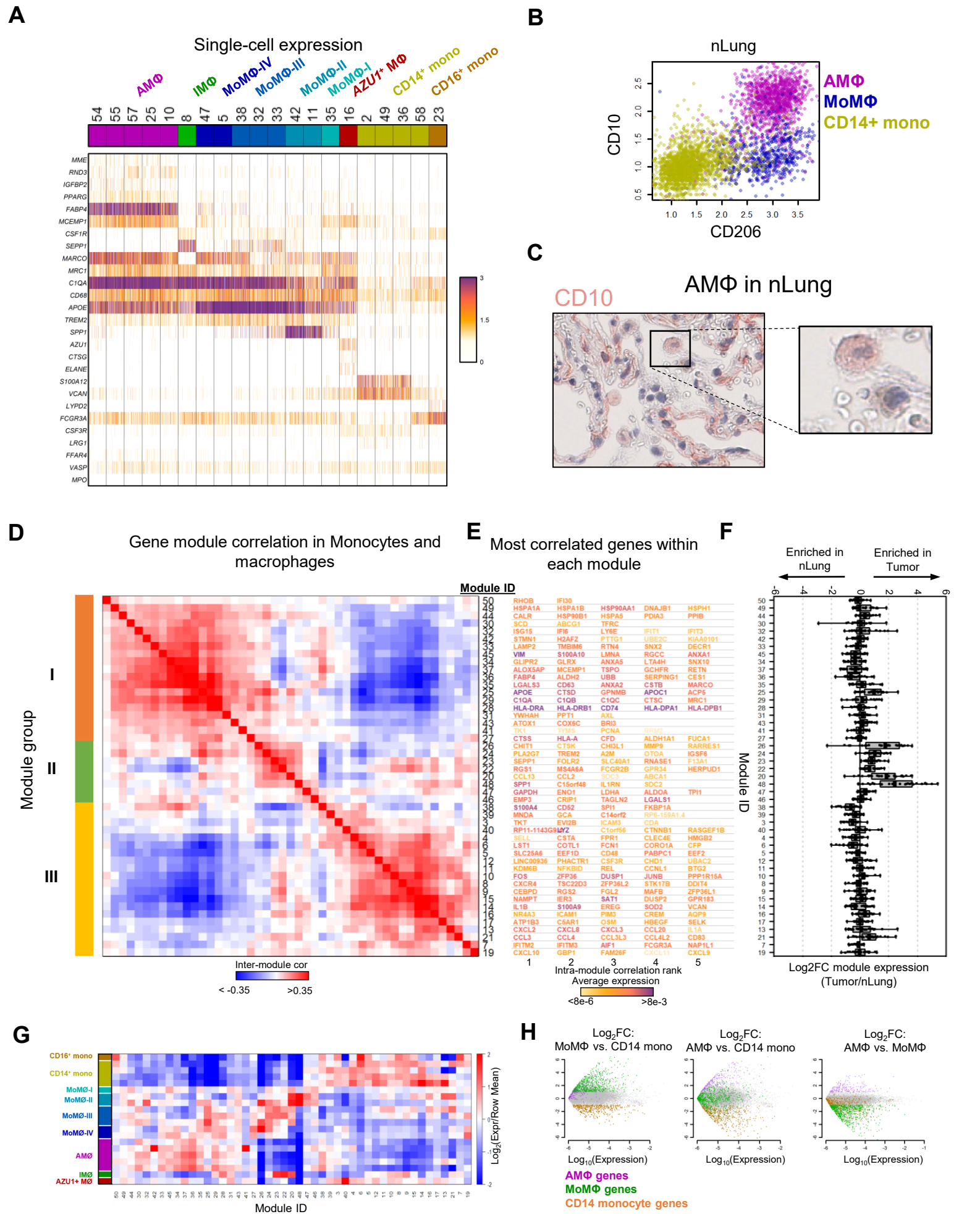


# Figure 3

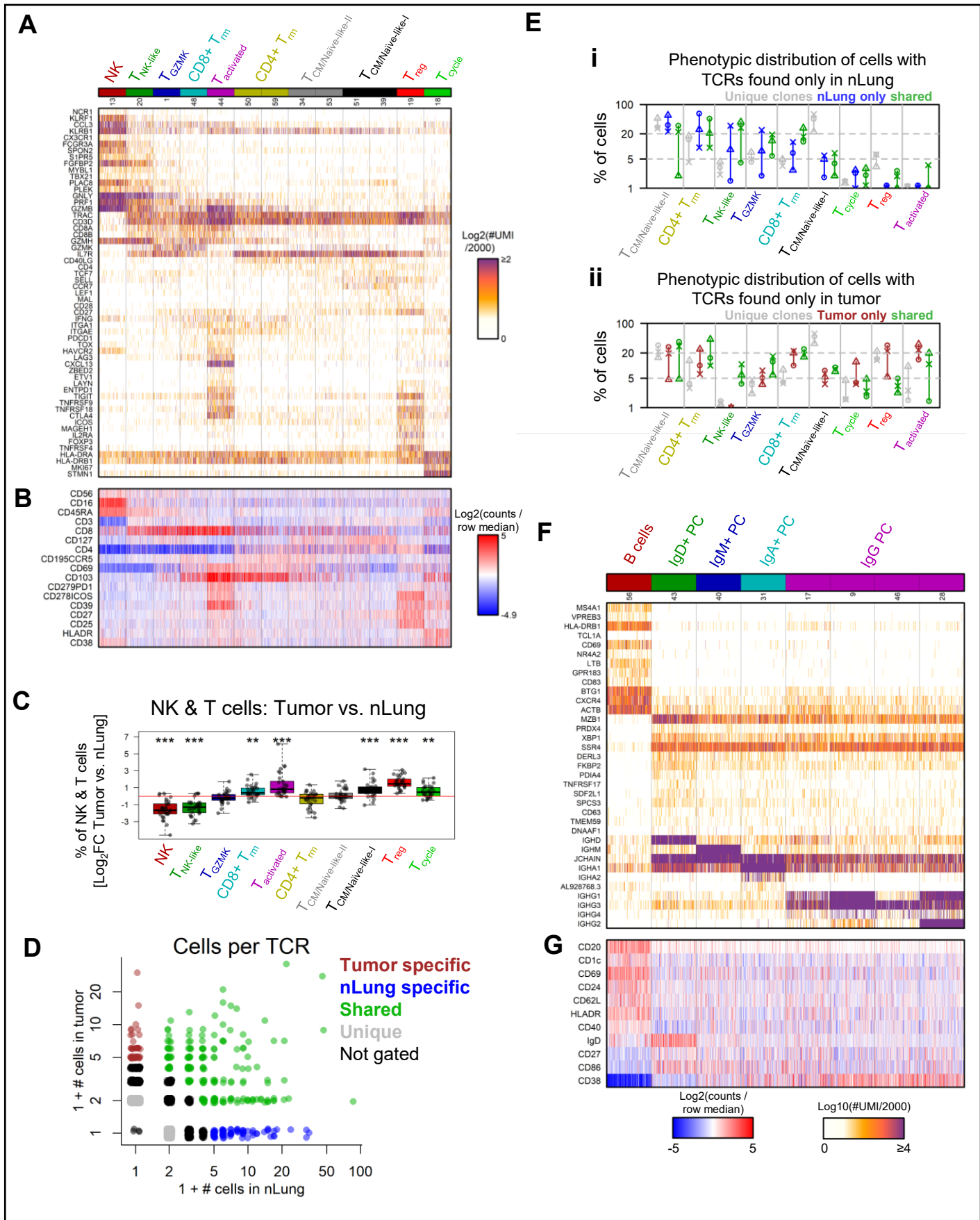


16  
8  
55  
57  
25  
54  
10  
47  
5  
38  
32  
42  
11  
35  
2  
49  
36  
58  
23

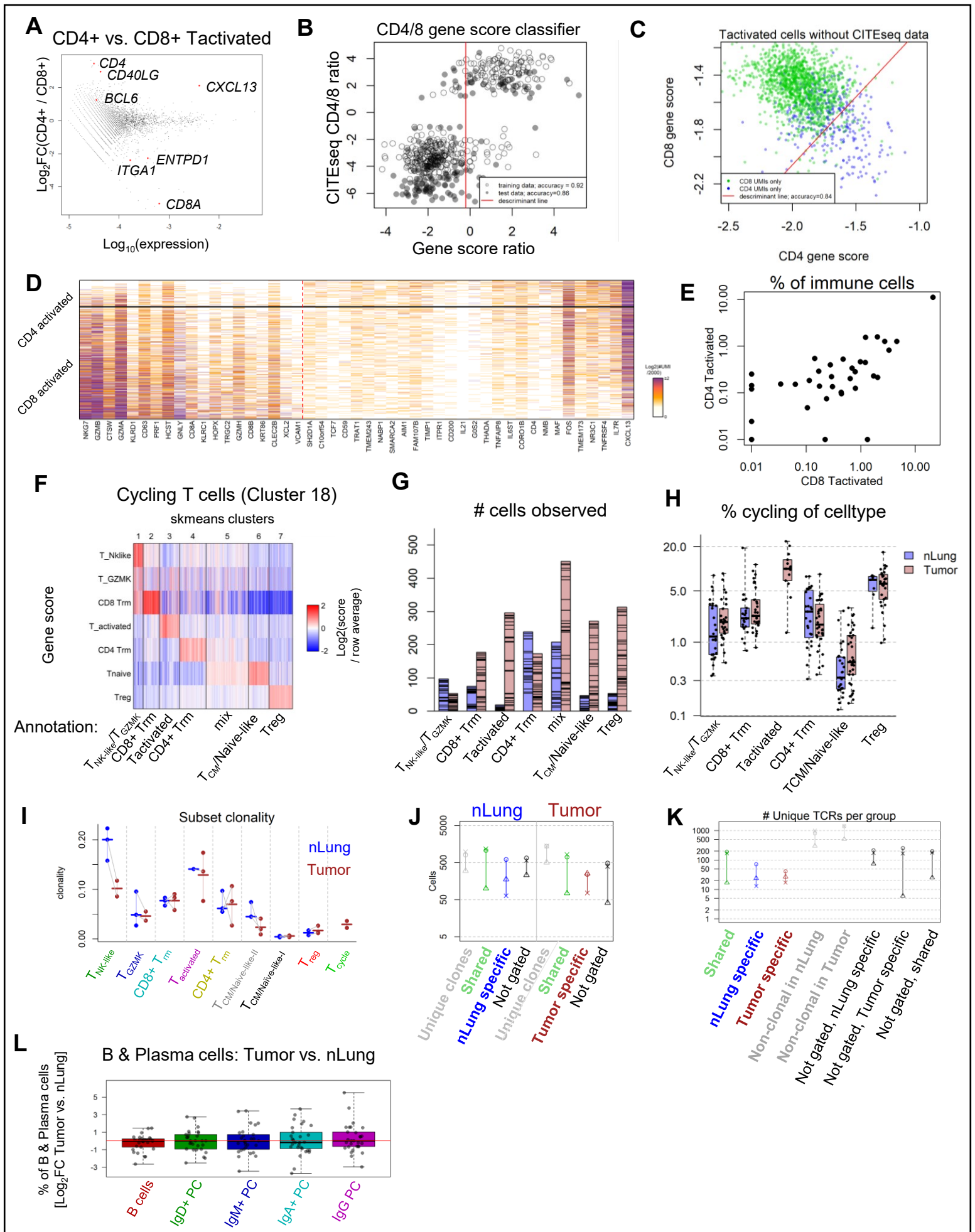
# Figure S3



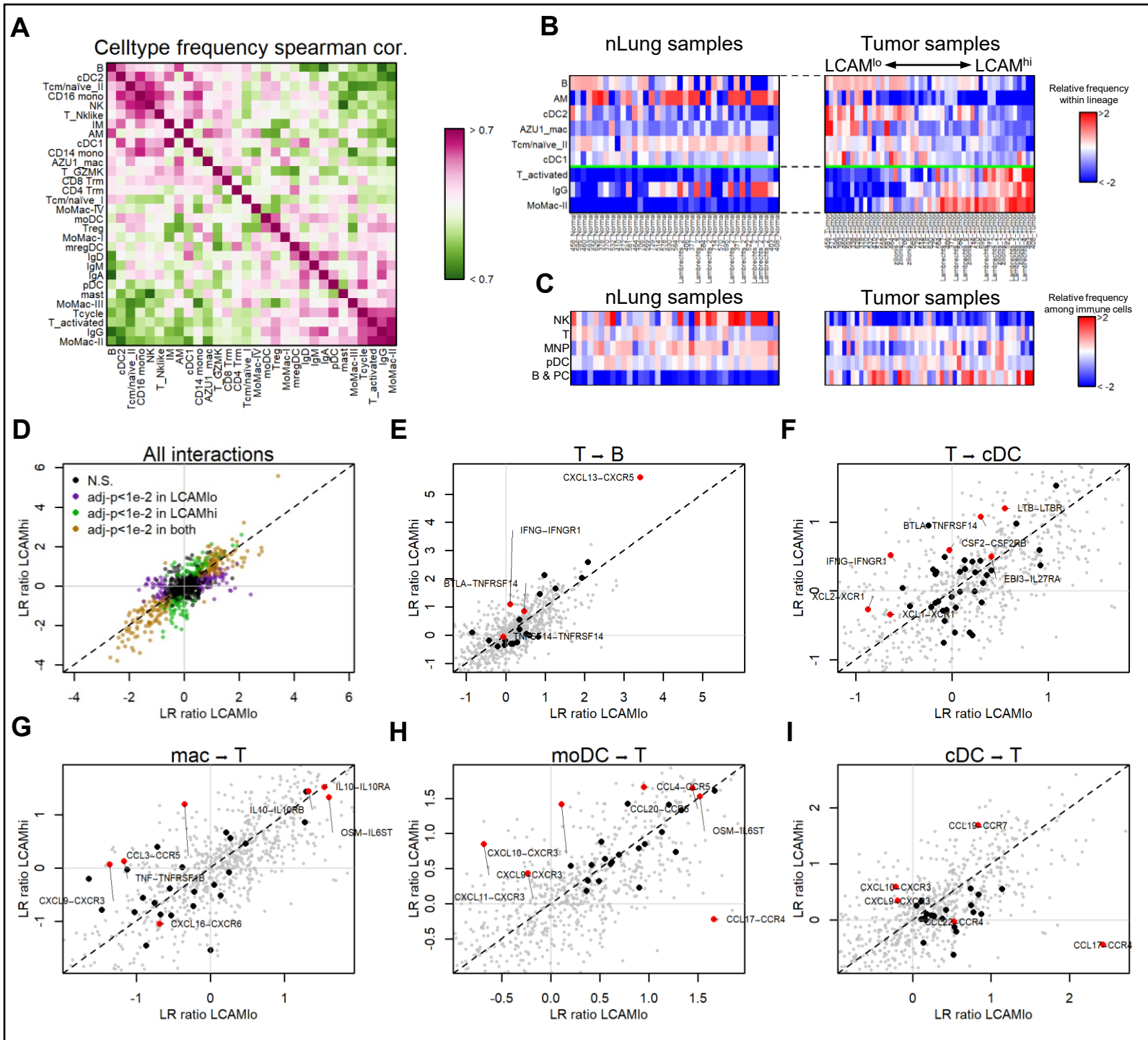




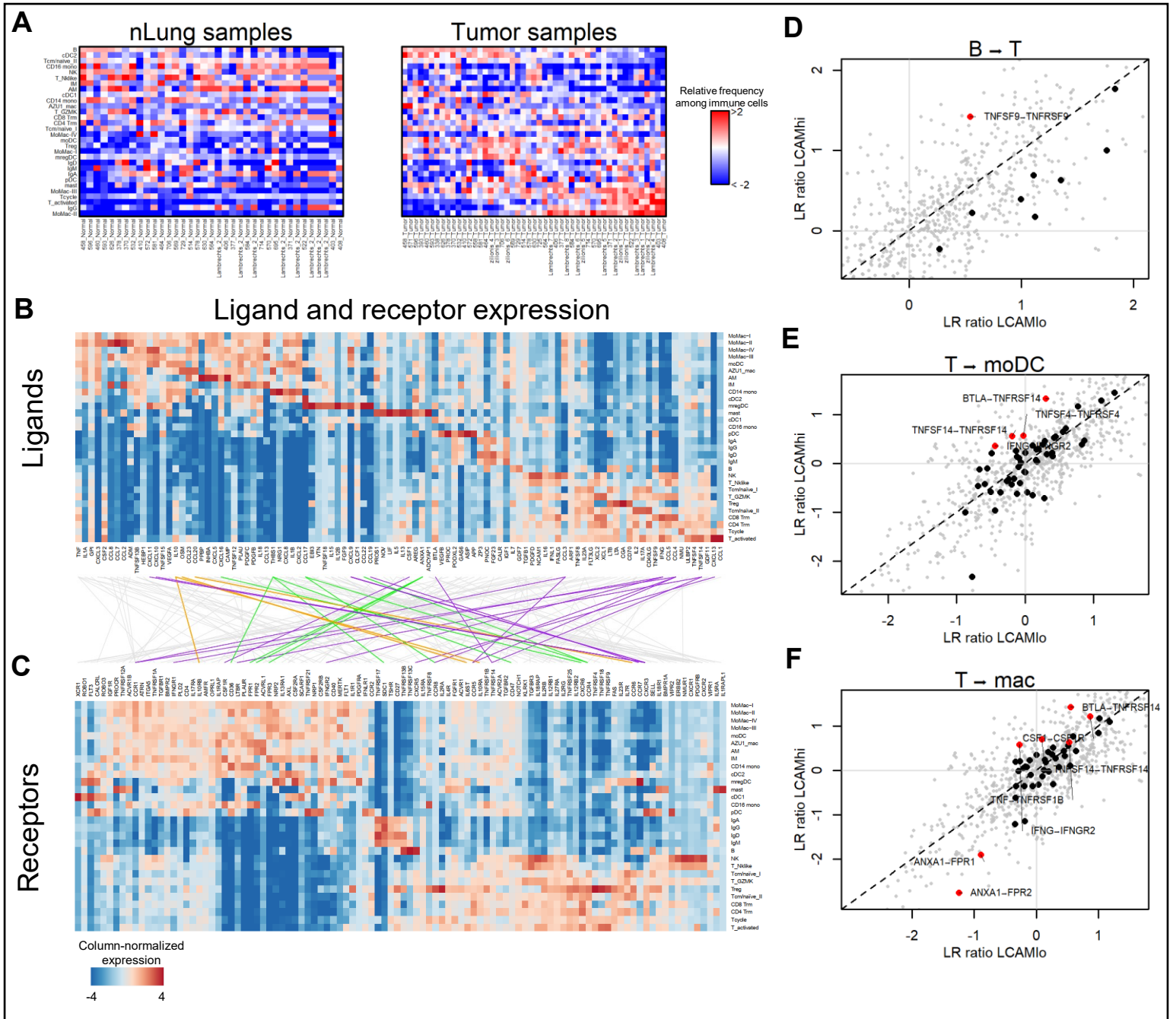
# Figure S4



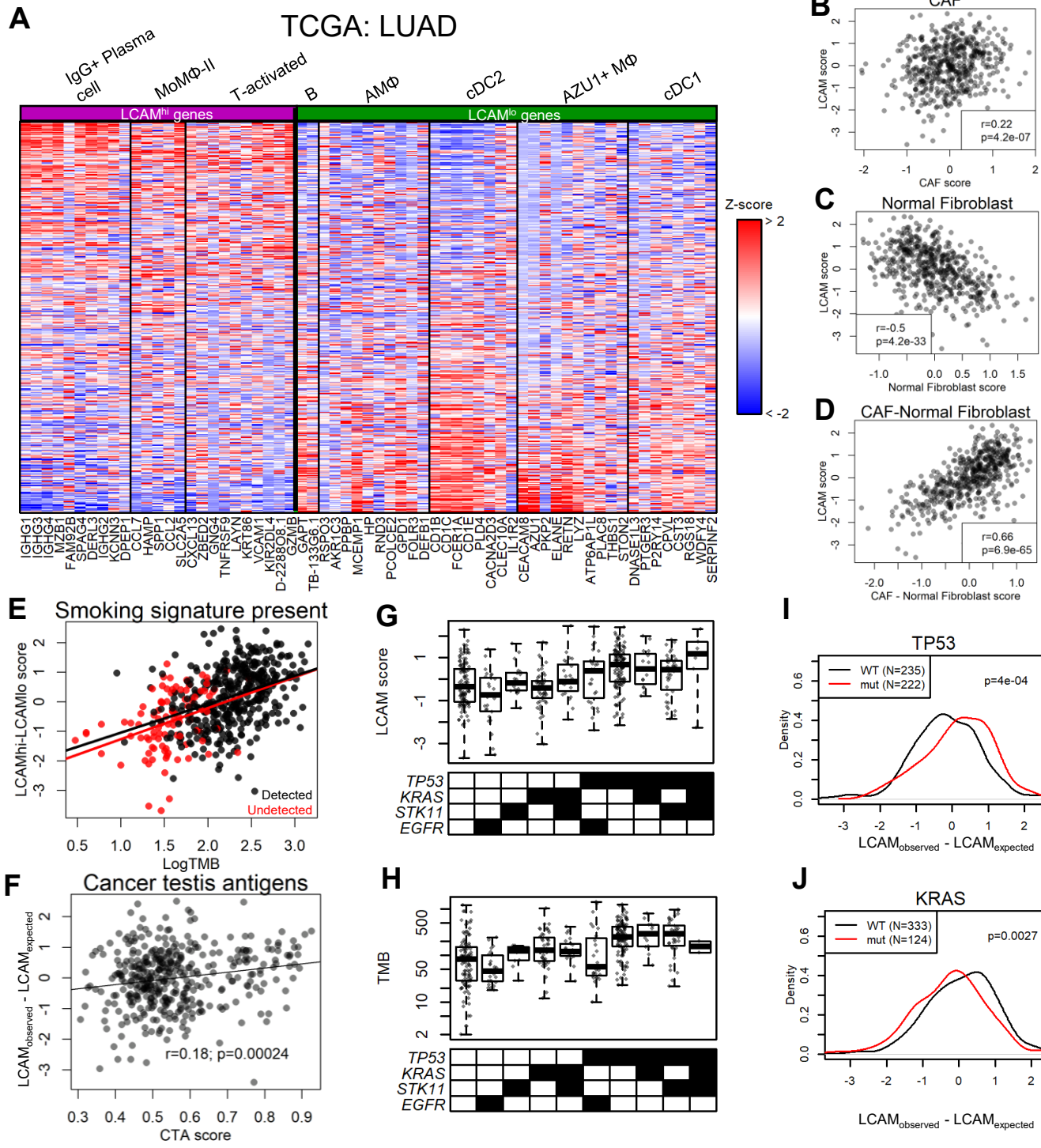
# Figure 5



# Figure S5



# Figure 6



# Figure S6

