

1 **TITLE**

2 **Aberrant expression of collagen type X in solid tumor stroma is associated with EMT,**  
3 **immunosuppressive and pro-metastatic pathways, bone marrow stromal cell signatures,**  
4 **and poor survival prognosis**

5  
6 **AUTHORS**

7 Elliot H.H. Famili-Youth\*<sup>†1,2</sup>, Aryana Famili-Youth\*<sup>†1,2</sup>, Dongfang Yang<sup>2</sup>, Ayesha Siddique<sup>2</sup>,  
8 Elizabeth Y. Wu<sup>2</sup>, Wenguang Liu<sup>3#</sup>, Murray B. Resnick<sup>2</sup>, Qian Chen\*<sup>3</sup>, and Alexander S.  
9 Brodsky\*<sup>2,4</sup>

10

11 <sup>1</sup>Medical Scientist Training Program, Northwestern University Feinberg School of Medicine,  
12 Chicago, IL, USA

13 <sup>2</sup>Department of Pathology and Laboratory Medicine, Rhode Island Hospital and Lifespan  
14 Medical Center, Warren Alpert Medical School of Brown University, Providence, RI, USA

15 <sup>3</sup>Department of Orthopedics, Rhode Island Hospital, Warren Alpert Medical School of Brown  
16 University, Providence, RI, USA

17 <sup>4</sup>Center for Computational Molecular Biology, Brown University, Providence, RI, USA

18 #Current address: School of Food Science and Engineering, Shaanxi University of Science and  
19 Technology, Xi'an, Shaanxi 710021, China

20 \*Corresponding authors

21 †Equal contribution

22

23 **Research Highlights**

- 24 ● ColX highlights features of EMT in breast and pancreatic cancer  
25 ● ColX gene modules are immunosuppressive and pro-metastatic  
26 ● ColX-associated gene networks contribute to sex differences in pancreatic cancer

- 27 • ColX-positive fibroblasts define more aggressive tumors with poorer survival
- 28 • ColX is emerging as a biomarker for bone marrow-derived cells in cancer

29

## 30 **ABSTRACT**

31       **Background:** Collagen type X (ColX $\alpha$ 1, encoded by *COL10A1*) is expressed specifically  
32 in the cartilage-to-bone transition, in bone marrow cells, and in osteoarthritic (OA) cartilage. We  
33 have previously shown that ColX $\alpha$ 1 is expressed in breast tumor stroma, correlates with tumor-  
34 infiltrating lymphocytes, and predicts poor adjuvant therapy outcomes in ER<sup>+</sup>/HER2<sup>+</sup> breast  
35 cancer. However, the underlying molecular mechanisms for these effects are unknown. In this  
36 study, we performed bioinformatic analysis of *COL10A1*-associated gene modules in breast and  
37 pancreatic cancer as well as in cells from bone marrow and OA cartilage. These findings  
38 provide important insights into the mechanisms of transcriptional and extracellular matrix  
39 changes which impact the local stromal microenvironment and tumor progression.

40       **Methods:** Immunohistochemistry was performed to examine collagen type X expression  
41 in solid tumors. WGCNA was used to generate *COL10A1*-associated gene networks in breast  
42 and pancreatic tumor cohorts using RNA-Seq data from The Cancer Genome Atlas.  
43 Computational analysis was employed to assess the impact of these gene networks on  
44 development and progression of cancer and OA. Data processing and statistical analysis was  
45 performed using R and various publicly-available computational tools.

46       **Results:** Expression of *COL10A1* and its associated gene networks highlights  
47 inflammatory and immunosuppressive microenvironments, which identify aggressive breast and  
48 pancreatic tumors and contribute to metastatic potential in a sex-dependent manner. Both  
49 cancer types are enriched in stroma, and *COL10A1* implicates bone marrow-derived fibroblasts  
50 as drivers of the epithelial-to-mesenchymal transition (EMT) in these tumors. Heightened  
51 expression of *COL10A1* and its associated gene networks is correlated with poorer patient  
52 outcomes in both breast and pancreatic cancer. Common transcriptional changes and

53 chondrogenic activity are shared between cancer and OA cartilage, suggesting that similar  
54 microenvironmental alterations may underlie both diseases.

55 **Conclusions:** *COL10A1*-associated gene networks may hold substantial value as  
56 regulators and biomarkers of aggressive tumor phenotypes with implications for therapy  
57 development and clinical outcomes. Identification of tumors which exhibit high expression of  
58 *COL10A1* and its associated genes may reveal the presence of bone marrow-derived stromal  
59 microenvironments with heightened EMT capacity and metastatic potential. Our analysis may  
60 enable more effective risk assessment and more precise treatment of patients with breast and  
61 pancreatic cancer.

62

### 63 **KEYWORDS**

64 collagen type X, tumor microenvironment, breast cancer, pancreatic cancer, osteoarthritis

65

### 66 **BACKGROUND**

67 The tumor microenvironment profoundly influences cancer progression and aggression  
68 through direct and indirect interactions between neoplastic cells and surrounding structures  
69 such as the extracellular matrix (ECM), whose remodeling plays a critical role in tumor growth,  
70 invasion, and metastasis<sup>1</sup>. The ECM encompasses a broad array of glycoproteins, collagens,  
71 and proteoglycans, as well as affiliated regulators and secreted factors; together, these exhibit a  
72 multitude of structural and signaling functions across diverse biological contexts and their  
73 alteration contributes to the development of pathological states ranging from fibrotic diseases to  
74 cancer<sup>2,3</sup>. Collagens are the most abundant protein component of the ECM and the composition  
75 of both major and minor collagens has been shown to vary substantially across different cancer  
76 types; additionally, numerous collagens have been identified as biomarkers associated with  
77 molecular alterations and overall survival in cancers of diverse primary tissues<sup>4</sup>. The  
78 composition and distribution of collagens within the local ECM is largely driven by fibroblasts,

79 which influence the processes of inflammation and angiogenesis through regulation of the  
80 ECM<sup>5</sup>, although other cells may also play a role in the production and degradation of collagens  
81 in disease states. Fibroblastic activity appears to be an important driver of disease across  
82 diverse tumor types, but the full composition and function of the ECM in cancer remains  
83 uncertain.

84       Collagen type X (*COL10A1*, ColX) is a non-fibrillar collagen synthesized specifically by  
85 hypertrophic chondrocytes to regulate matrix mineralization, stiffness, and metabolism<sup>6</sup>. ColX  
86 promotes the cartilage-to-bone transition in skeletal development, is highly expressed by bone  
87 marrow stromal cells (BMSCs), and becomes progressively elevated in articular cartilage during  
88 the development of osteoarthritis (OA)<sup>6-9</sup>. OA pathogenesis has been associated with  
89 senescence of mesenchymal stromal cells, chondrocyte death, calcification and degradation of  
90 the extracellular matrix, and angiogenic invasion<sup>10-12</sup>. Previously we found that ColX is  
91 expressed in breast tumor tissue, the first time that ColX was shown to be highly expressed in  
92 non-skeletal tissues<sup>13-15</sup>. Prior studies by our group have shown that ColX is not only expressed  
93 in many types of breast tumors, but is also associated with overall survival outcomes for  
94 ER<sup>+</sup>/HER2<sup>+</sup> breast tumors in particular<sup>13,15</sup>. However, the mechanism by which ColX is involved  
95 in tumor progression and treatment outcomes remains unknown. Pancreatic ductal  
96 adenocarcinoma (PDAC) remains a highly lethal malignancy due to its aggressive nature and  
97 the paucity of effective treatment options; such tumors are notable for their high fractions of  
98 desmoplastic stroma which contributes significantly to drug and immune resistance<sup>16</sup>. Complex  
99 interactions between PDAC cancer cells and surrounding stromal features such as activated  
100 fibroblasts and collagens play a major role in aggressive, treatment-refractory disease<sup>16,17</sup>.  
101 Thus, one approach to improve our understanding of stromal impacts in cancer is to identify key  
102 ECM features which drive the development, survival, and progression of such tumors.

103       Core environmental factors which influence tumor outcomes include stromal  
104 composition, blood vessel density, and infiltrating immune cells<sup>18</sup>. The complicated interplay

105 between resident and foreign host cells, the extracellular matrix, and molecular signals all  
106 contribute to primary tumor treatment responses. Recent studies have suggested that the  
107 stromal fractions of breast and pancreatic tumors feature significant proportions of cancer-  
108 associated fibroblasts (CAFs), which exhibit substantial heterogeneity, originate from both the  
109 local biome and differentiated bone marrow-derived mesenchymal stromal cells, and contribute  
110 significantly to patient prognosis and response to therapy<sup>19,20</sup>. Given ColX's important role in  
111 cartilage development and the bone marrow niche<sup>21</sup>, along with its dysregulated expression  
112 across both OA cartilage and solid tumors, we sought to characterize its pathophysiologic role in  
113 cancer through bioinformatic analysis of *COL10A1*-expressing cancer and non-cancer cells. We  
114 hypothesized that similar stromal microenvironments across certain cancers, bone marrow, and  
115 OA cartilage are defined by ColX and its associated gene networks, which may contribute to  
116 molecular mechanisms underlying tumor progression. In this study, we defined gene co-  
117 expression modules to characterize pathways and microenvironmental components related to  
118 and correlated with ColX expression in breast and pancreatic tumors from The Cancer Genome  
119 Atlas (TCGA). By analyzing ColX expression in BMSCs and OA cartilage cells, we found  
120 notable common biological pathways in both cancer and BMSCs, thereby linking these two  
121 types of stromal cells. Characterization of ColX expression networks and their pathological  
122 mechanisms will improve understanding of aggressive disease states and offer opportunities for  
123 devising future therapies.

124

## 125 **METHODS**

### 126 **Immunohistochemistry and ColX $\alpha$ 1 expression scoring**

127 Two PAAD samples were tested compared to breast tumor observations. One stage 3  
128 and one stage 4 sample were evaluated for ColX $\alpha$ 1 protein expression as follows. Four-micron  
129 sections were cut from formalin-fixed paraffin-embedded tissue blocks, heated at 60°C for 30  
130 minutes, deparaffinized, rehydrated, and subjected to antigen retrieval by heating the slides in

131 epitope retrieval buffer in a water bath at 95°C for 45 minutes. The slides were then incubated  
132 with either mouse monoclonal antibodies or rabbit polyclonal antibodies for 30 minutes at room  
133 temperature in a DAKO Autostainer. Anti-ColX $\alpha$ 1 (1:100, eBioscience/Affymetrix, Clone X53)  
134 was used for immunohistochemistry (IHC). Immunoreactivity was detected using the DAKO  
135 EnVision method according to the manufacturer's recommended protocol.

136

### 137 **Data analysis and visualization**

138 All data processing and analysis was performed in R (version 4.0.2)<sup>22</sup> unless otherwise  
139 stated. Visualizations were generated in R using the ggplot2 (version 3.3.6)<sup>23</sup>, gplots (version  
140 3.1.3)<sup>24</sup>, and eulerr (version 7.7.0)<sup>25,26</sup> packages. See **Figure 1C** for overview of tumor sample  
141 datasets and computational tools employed in this study.

142

### 143 **Data acquisition and pre-processing**

#### 144 *TCGA gene expression data*

145 Log-normalized expression for *COL10A1* across all TCGA cancer types was  
146 downloaded from the Broad Institute's Genome Data Analysis Center FireBrowse portal. Batch-  
147 corrected, normalized RNA-Seq-derived RSEM values for breast invasive carcinoma (BRCA,  $n$   
148 = 1,095 samples) and pancreatic adenocarcinoma (PAAD,  $n$  = 178 samples) TCGA cohorts  
149 were downloaded from the NIH National Cancer Institute Genomic Data Commons<sup>27,28</sup>. Genes  
150 with RSEM values  $< 1$  in  $\geq 50\%$  of samples and mean RSEM values  $< 50$  overall were defined  
151 as "low-expression" and excluded from downstream analysis.

#### 152 *Microarray gene expression data*

153 To assess whether RNA-Seq-derived ColX-related modules would be robust to different  
154 methodologies of gene expression quantification, microarray datasets with similar sample sizes  
155 were selected for comparison to results from the BRCA and PAAD cohorts from TCGA. For  
156 breast cancer, raw fluorescence intensity data from a previously-described collection of 12

157 studies on primary early-stage breast cancer in females was downloaded from the NCBI Gene  
158 Expression Omnibus<sup>29</sup> (GEO) database ( $n = 1,763$  samples in total), all of which were  
159 expression-profiled on the GPL570 (Affymetrix Human Genome U133 Plus 2.0 Array) platform  
160 (see **Table S1A** for comparison to TCGA data and **Table S1B** for GSE accessions and number  
161 of samples analyzed from each study)<sup>30–42</sup>. For pancreatic cancer, array-based gene expression  
162 data for the Australian pancreatic cancer cohort (PACA-AU,  $n = 269$  samples) was downloaded  
163 from the International Cancer Genome Consortium (ICGC) Data Portal<sup>43</sup>. Pre-processing of  
164 microarray data was carried out using the following packages in R: oligo (version 1.52.1)<sup>44</sup>,  
165 hgu133plus2.db (version 3.2.3)<sup>45</sup>, AnnotationDbi (version 1.50.3)<sup>46</sup>, tidyverse (version 1.3.0)<sup>47</sup>,  
166 WGCNA (version 1.69)<sup>48,49</sup>, and sva (version 3.36.0)<sup>50</sup>. Probe intensity values across each  
167 cancer were log-transformed and normalized using the Robust Multichip Average (RMA)  
168 quantile method. Probes were mapped to gene IDs based on the GPL570 annotation database,  
169 and unmapped or multi-mapping probes were removed. Expression values for multiple probes  
170 mapping to the same gene were consolidated using the collapseRows function. Gene  
171 expression values for the breast cancer samples were then combined across GEO studies and  
172 batch-corrected using the ComBat function. Finally, “low-intensity” expression thresholds were  
173 established for each cancer dataset, and all genes with expression values below these  
174 thresholds in  $> 80\%$  of samples were defined as “low-expression” and excluded from  
175 downstream analysis.

#### 176 *OA gene expression data*

177 Raw RNA-Seq-derived gene counts for 4 knee joint-derived OA cell types (normal  
178 cartilage stromal cells/NCSCs, OA mesenchymal stromal cells/OA-MSCs, OA  
179 chondrocytes/OACs, bone marrow stromal cells/BMSCs;  $n = 3$  each) were sourced from GEO  
180 accession GSE176199<sup>10</sup>. Genes exhibiting fewer than 5 counts in  $\geq 90\%$  of samples were  
181 defined as “low-expression” and excluded from downstream analysis. DESeq2 (version  
182 1.34.0)<sup>51</sup> was used to normalize raw gene counts and perform differential expression analysis.

183 For each OA cell type, cell type-specific genes were defined based on the definition of “tissue-  
184 enriched” employed by the Human Protein Atlas (at least four-fold higher mRNA level in a given  
185 tissue compared to any other tissues)<sup>52</sup>; i.e., all genes with  $\log_2$ -fold change  $\geq +2$  and adjusted  
186 p-value  $< 0.05$  relative to each other cell type.

187

## 188 **Characterization of ColX consensus modules**

### 189 *Generation of ColX modules*

190 Correlated gene network modules were generated for each TCGA dataset based on  
191 normalized, filtered RSEM values using WGCNA (version 1.69)<sup>48,49</sup>. Briefly, for each dataset,  
192 the soft thresholding power  $\beta$  was selected to ensure approximately scale-free topology of the  
193 gene co-expression network, and signed modules were generated using the `blockwiseModules`  
194 function with parameters `deepSplit = 2` and `minModuleSize = 30`. ColX modules were defined  
195 for each dataset separately, comprising all genes which co-modularized with *COL10A1*.

### 196 *Enrichment analysis*

197 Enrichment analysis of ColX modules was performed using the ConsensusPathDB web  
198 tool (release 35)<sup>53</sup>. For each cancer, the *over-representation analysis* tool was run to identify all  
199 Reactome pathways and gene ontology (GO) terms enriched in the respective ColX module,  
200 using as background all genes which were retained after pre-processing the dataset from which  
201 the module was generated. Significantly-enriched pathways/terms were identified by FDR-  
202 corrected p-value  $< 0.05$ .

### 203 *Gene set overlap analysis*

204 Matrisome, hallmark pathway, and Gene Transcription Regulation Database (GTRD)  
205 transcription factor target (TFT) gene sets were downloaded from MSigDB<sup>2,54,55</sup>. Overlaps with  
206 breast and pancreatic cancer ColX modules were calculated using Fisher’s exact test. Adjusted  
207 p-values (p.adj) were computed using the Benjamini-Hochberg (BH) method<sup>56</sup>. Significantly-  
208 enriched gene sets were identified by  $p.\text{adj} < 0.05$  ( $p.\text{adj} < 0.10$  for candidate discovery of



209 enriched TFTs) and odds ratio > 1. Protein-protein interactions of specific transcription factors  
210 (TFs) of interest were queried using the STRING database<sup>57</sup>.

### 211 *Module preservation analysis*

212 Preservation of TCGA BRCA and PAAD RNA-Seq-derived ColX modules in  
213 corresponding cancer microarray datasets of comparable sample size was assessed using the  
214 modulePreservation function following standard WGCNA methodology<sup>58</sup>. The  $Z_{summary}$  statistic,  
215 defined as the mean of summarized density preservation statistics and connectivity preservation  
216 statistics, was computed for each TCGA-generated module in order to assess the relative  
217 preservation of the ColX module. A  $Z_{summary}$  value > 10 was considered to indicate significant  
218 module preservation.

219

## 220 **Differential pathway analysis**

### 221 *Gene set analysis*

222 Hallmark pathway and GTRD TFT gene sets were downloaded from MSigDB as  
223 described above. Immunome gene sets were obtained from The Cancer Immunome Atlas<sup>59</sup>.  
224 Cancer-associated fibroblast gene sets were extracted from previously-published datasets and  
225 defined as all genes exhibiting  $\geq 2$ -fold increased expression (with  $p_{adj} < 0.05$ ) in each  
226 fibroblast phenotype of interest relative to all others<sup>19</sup>.

### 227 *Definition of the G.A.M.E. metric*

228 To compare ColX module expression across all tumor samples within each dataset, a  
229 ranking metric was defined based on the sample-wise percentage of ColX module genes  
230 expressed above their respective median values across all samples (percentage of **Genes**  
231 **Above-Median Expression**, “%G.A.M.E.”). This effectively transformed the unimodal ColX  
232 module eigengene (ME) distribution into a bimodal G.A.M.E. distribution with a fixed range  
233 between 0 and 1, facilitating clustering of samples into discrete groups (**Figure S3**). The  
234 getJenksBreaks function from the BAMMtools package (version 2.1.10)<sup>60</sup> was used to divide

235 samples into “low”, “medium”, and “high” G.A.M.E. groups, which were used to proxy ColX  
236 module expression for comparisons between subgroups.

### 237 *QuSAGE analysis*

238 Differential activity of gene sets between high and low G.A.M.E. tumor samples were  
239 assessed using qusage (version 2.22.0)<sup>61–63</sup>. Log-fold change was used to quantify association  
240 with differential ColX module expression (i.e., enrichment). Activation and inhibition of individual  
241 gene sets/pathways were defined as positive and negative enrichment relative to the  
242 background gene set, respectively. Significance associated with ColX module expression was  
243 determined by BH-adjusted p-value < 0.05 and absolute magnitude of enrichment ≥ 25% of the  
244 magnitude of activity of the ColX module (both measured relative to the background gene set)  
245 within that dataset.

246

### 247 **Survival analysis**

#### 248 *Cox proportional hazards model analysis*

249 Survival analysis was carried out in R using the survival (version 3.3-1)<sup>64,65</sup> and  
250 survminer (version 0.4.9)<sup>66</sup> packages. Multivariate Cox proportional hazards models were  
251 constructed to assess the statistical dependence of overall survival (OS) and disease-free  
252 interval (DFI) on patient age, gender, tumor stage (binarized as stage 1-2 vs. stage 3-4), and  
253 ColX tumor signal (assessed as either normalized gene expression, ColX module eigengene  
254 (ME), or %G.A.M.E.) for breast and pancreatic cancer cohorts as well as gender-specific  
255 pancreatic cancer subcohorts. The proportional hazards assumption was validated for each  
256 signal variable by a test of the scaled Schoenfeld residuals. Significant  $\beta$  coefficients were  
257 determined by a p-value < 0.05, and the per-unit contribution of each significant variable to the  
258 overall hazard risk was computed as  $e^{\beta}$ .

259 *Kaplan-Meier survival curve analysis*

260 Differential survival outcomes were assessed using Kaplan-Meier analysis. For each  
261 comparison, samples were divided into two groups: a “low” %G.A.M.E. group (defined as  
262 above), and a “high” %G.A.M.E. group (all other samples).

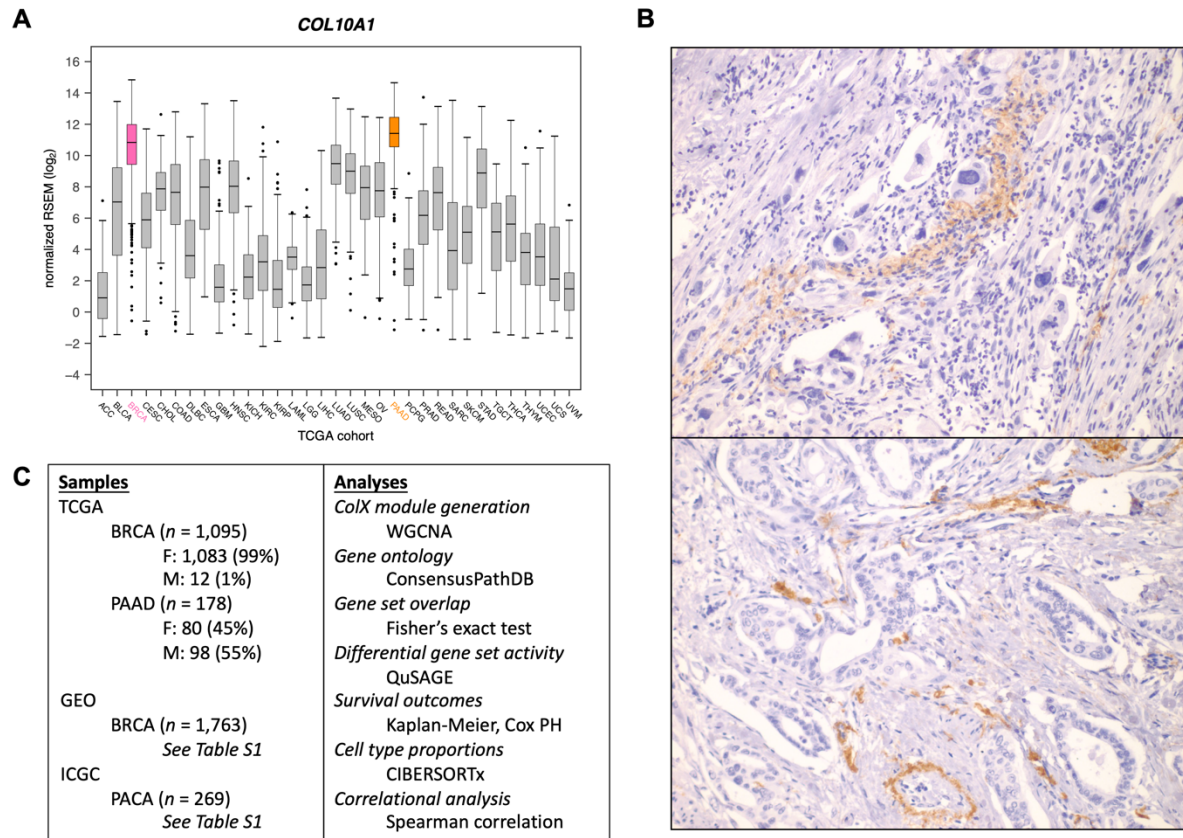
263

#### 264 **Tumoral osteoarthritic cell type proportion analysis**

265 CIBERSORTx was utilized to quantify OA cell type proportions for all TCGA tumor  
266 samples<sup>67</sup>. The *Create Signature Matrix* tool was used to generate a dimensionally-reduced  
267 expression signature (1,978 genes) to differentiate the 4 OA cell types profiled (NCSC, OA-  
268 MSC, OAC, BMSC). The *Impute Cell Fractions* tool was subsequently employed to infer the  
269 approximate proportions of each OA cell type present in each tumor sample, quantified on an  
270 absolute scale. Absolute proportions were scaled by sample-wise total OA cell type proportions  
271 to obtain sample-wise relative proportions.

272

273 **RESULTS**



274

275 **Figure 1: COL10A1 is highly expressed in breast and pancreatic tumors. (A)** Expression of  
 276 *COL10A1* across all TCGA sample cohorts. **(B)** Representative 400x ColX $\alpha$ 1 immunohistochemistry  
 277 staining in pancreatic tumors. **(C)** Outline of study samples and analyses. See Methods and Results  
 278 sections for details.

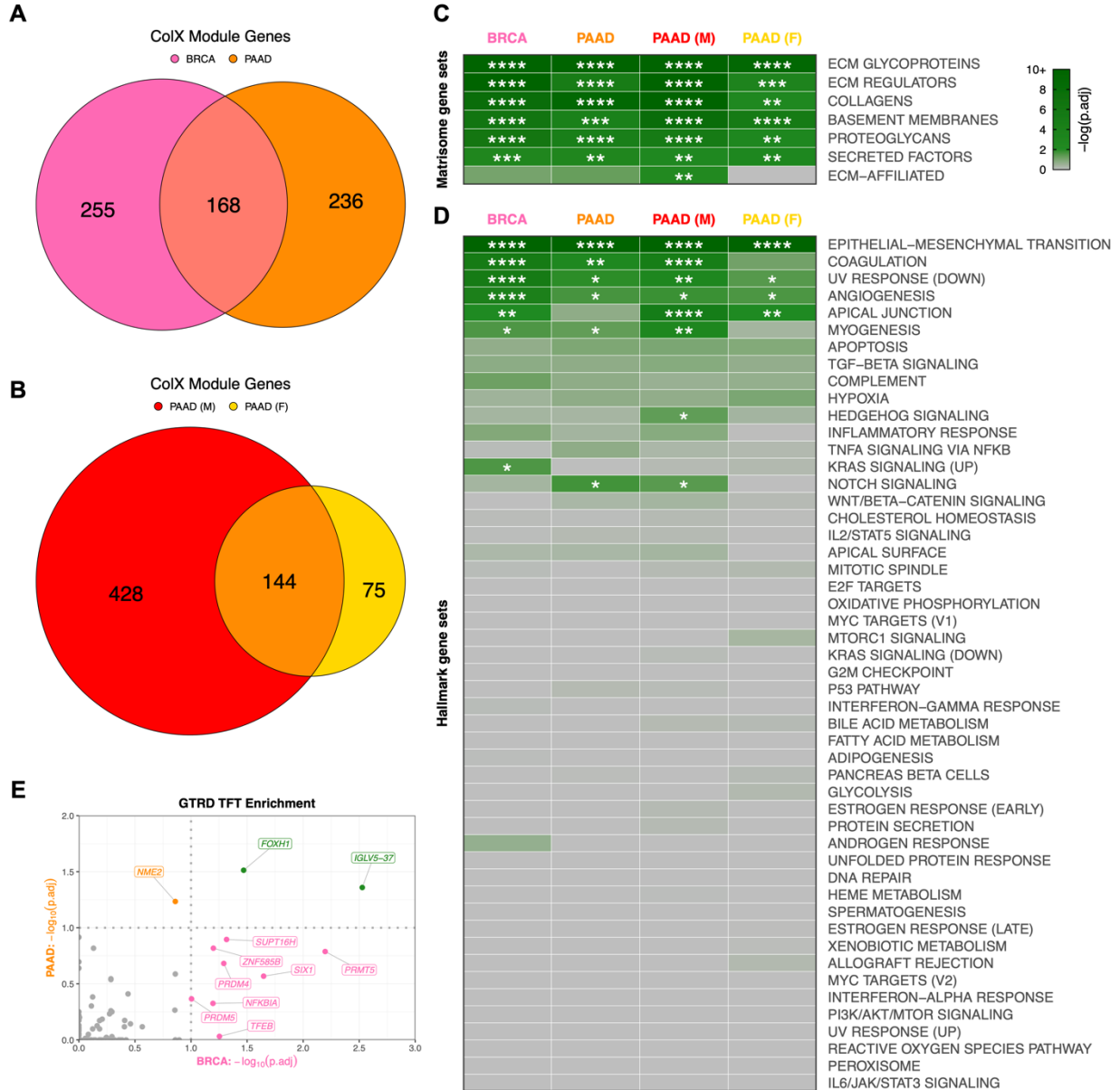
279 **ColX-associated gene modules are conserved across different cancers**

280 *COL10A1* mRNA expression was found to be the highest in breast invasive carcinoma  
 281 (BRCA) and pancreatic adenocarcinoma (PAAD) tumors in TCGA (Figure 1A). We evaluated  
 282 expression of ColX at the protein level by IHC to determine relative levels of expression. Strong  
 283 expression was observed in BRCA, as we have previously reported<sup>13–15</sup>. Similar to the patterns  
 284 observed in breast tumors, ColX $\alpha$ 1 IHC revealed a range of mild to strong expression of ColX  
 285 protein only in the stromal regions of pancreatic tumors (Figure 1B). Consistent with much lower  
 286 expression of *COL10A1* mRNA, no protein expression was observed in colon adenocarcinoma

287 or stomach adenocarcinoma by IHC. In accordance with these observations, we focused on  
288 BRCA and PAAD tumors to evaluate the impact of *COL10A1* on cancer pathophysiology.

289 To define the role of ColX in breast and pancreatic tumors, we analyzed multiple  
290 datasets by a variety of statistical methods (Figure 1C). To identify possible roles and gene  
291 networks associated with *COL10A1* in breast and pancreatic tumors, we used weighted gene  
292 co-expression network analysis (WGCNA) to generate cancer-specific ColX gene modules in  
293 the TCGA datasets for each cancer type. We defined “ColX-associated” genes as all those  
294 which were co-modularized with ColX in each dataset, signifying genes whose expression was  
295 broadly associated with that of ColX across the patient cohorts. This process yielded ColX  
296 modules of 423 and 404 genes across the breast and pancreatic cancer datasets, respectively,  
297 with 168 (approximately 40%) of these genes co-modularizing with ColX in both cancers (Figure  
298 2A, Table S2A; see Table S2B for full list of WGCNA modules for each dataset).

299 To verify the reproducibility of these ColX modules, we performed module preservation  
300 analysis on a comparably-sized microarray dataset for each cancer type: a collection of 1,763  
301 primary early-stage breast cancer samples sourced from GEO as well as 269 pancreatic cancer  
302 samples from the PACA-AU cohort. Both TCGA RNA-Seq-derived ColX modules were highly  
303 conserved in their corresponding microarray cancer datasets; the BRCA ColX module was the  
304 #2 most highly preserved among all 31 BRCA gene modules and the PAAD ColX module was  
305 the #1 most highly preserved among all 42 PAAD gene modules. This signifies the relevance of  
306 these cancer-specific gene sets across diverse patient cohorts (Figure S1).



307

308 **Figure 2: ColX modules are enriched for ECM genes, pro-metastatic pathways, and**  
309 **developmental/regulatory transcription factor targets. (A and B)** Overlap of genes within ColX  
310 WGCNA modules from TCGA **(A)** breast and pancreatic cancer and **(B)** gender-segregated pancreatic  
311 cancer datasets. See **Table S2A** for lists of ColX module genes, **Table S2B** for full list of WGCNA  
312 modules for each dataset, and **Table S3** for gene ontology and Reactome pathway enrichment analysis of  
313 cancer-specific and overlapping ColX-associated genes. **(C and D)** Enrichment of **(C)** human *in silico*  
314 matrisome gene sets<sup>2</sup> and **(D)** MSigDB hallmark pathway gene sets<sup>54</sup> within ColX modules. Gene sets  
315 within each block are ordered by mean significance rank (by Fisher's exact test) across all 4 modules.  
316 See **Figure S2A–D** and **Table S4** for hallmark pathway enrichment analysis of all WGCNA-inferred  
317 modules for each dataset. Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*,  
318 p.adj < 0.0001. **(E)** Enrichment of Gene Transcription Regulation Database (GTRD) transcription factor  
319 **(TF)** targets<sup>55</sup> within breast and pancreatic cancer ColX modules. Of note, the TFT gene set attributed to  
320 *IGLV5-37* in the GTRD database actually represents targets of the fusion oncoprotein *SS18-SSX*, as  
321 described in the text. Dotted lines correspond to p.adj = 0.10. Green labels/points indicate TFs whose  
322 targets are enriched in both breast and pancreatic cancer ColX modules. See **Table S5A** for reported TF  
323 functions and lists of overlapping TFTs. Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj <  
324 0.001; \*\*\*\*, p.adj < 0.0001.

## 325 **ColX-associated gene networks are enriched for ECM and developmental ontologies**

326 To probe the biological importance of these ColX modules in the context of breast and  
327 pancreatic cancer, we performed gene ontology (GO) and Reactome pathway enrichment  
328 analysis to identify functions of ColX modules in each dataset (**Tables S3, S4**). As anticipated,  
329 GO terms relating to extracellular matrix and collagen function were highly enriched across both  
330 the breast and pancreatic cancer ColX modules, as well as in their overlapping gene sets (**Table**  
331 **S3A–C**). Notably, GO categories including “cell migration” and “cell motility” were enriched in  
332 both cancer types, as were “ossification,” “cartilage development,” “skeletal system  
333 development,” and several terms related to Wnt signaling. Reactome pathways relating to  
334 collagen organization and ECM structure, function and degradation were similarly enriched,  
335 along with the role of *RUNX2* in regulating chondrocyte maturation (**Table S3F–H**).

336 In addition to these shared GO categories, breast and pancreatic cancer-specific ColX  
337 modules were individually enriched for several categories associated with the epithelial-to-  
338 mesenchymal transition (EMT). Enriched GO categories specific to the BRCA ColX module  
339 included “epithelial cell migration” and “mesenchymal cell proliferation,” and several Reactome  
340 pathways relating to signaling through *FGFR2* were also significantly enriched (**Table S3A,**

341 **S3F**). PAAD-specific GO hits included several developmental regulators such as Frizzled and  
342 Smoothed, which were corroborated by enrichment of Reactome pathways relating to Wnt,  
343 Hedgehog, and *RUNX2* signaling (**Table S3B, S3G**).

344

### 345 **ColX-associated gene networks are more strongly linked to OA and pro-metastatic** 346 **processes in male pancreatic cancer cohorts**

347 To assess the possible differential contribution of gender to the ColX-related pathology  
348 of pancreatic cancer, we also used WGCNA to define ColX modules in male and female patient  
349 subsets of the PAAD dataset, yielding 572 and 219 genes respectively, of which 144 were  
350 shared between male and female cohorts (**Figure 2B**). As male patients comprised only 1% of  
351 the BRCA dataset, we did not perform gender-specific analysis for breast cancer.

352 While both gender-specific pancreatic cancer modules were enriched for numerous GO  
353 terms and Reactome pathways which were also significant in their combined analysis (e.g.,  
354 various ECM terms, “ossification,” “cartilage development,” and several terms related to  
355 chondroitin sulfate metabolism), the male PAAD ColX module was additionally enriched for the  
356 GO terms “bone morphogenesis,” “Wnt signaling,” and “epithelial to mesenchymal transition,” as  
357 well as Reactome pathways relating to *RUNX2* signaling and platelet responses (**Table S3D,**  
358 **S3I**). In contrast, the female PAAD ColX module was not significantly enriched for these terms  
359 or pathways; the primary hits were largely similar to those of the full PAAD ColX module (**Table**  
360 **S3E, S3J**), with inference of a more functionally-restricted ColX-associated genetic network  
361 possibly due to the smaller sample size.

362

### 363 **Breast and pancreatic cancer-specific ColX modules overlap with key matrisome and** 364 **oncogenic pathway gene sets**

365 To assess processes and gene functions captured within these four dataset-specific  
366 ColX modules, we performed statistical overlap analyses with multiple well-characterized gene



367 sets encompassing a broad range of physiological and clinical functions. Both breast and  
368 pancreatic ColX modules were found to overlap significantly with matrisome gene sets,  
369 including ECM glycoproteins, collagens, ECM regulators, basement membranes, proteoglycans,  
370 and secreted factors<sup>2</sup> (Figure 2C). The ColX modules include numerous matrisome genes  
371 previously implicated in tumor progression and metastasis, including ECM regulators *CTSB*,  
372 *LOXL2*, and *SERPINF1*, and ECM glycoprotein *SNED1*, which have been reported as markers  
373 of highly metastatic breast carcinomas that promote tumor invasiveness across a variety of  
374 models<sup>68</sup> (Table S2A). *CTSB* is also upregulated in pancreatic cancer and may indicate  
375 increased activity of *CSTB*, which enhances later metastatic extravasation in PDAC;  
376 additionally, several collagens that are highly expressed in PDAC relative to its precursor  
377 pancreatic intraepithelial neoplasia (*COL6A1*, *COL6A2*, and *COL11A1*) are present in the  
378 pancreatic ColX modules<sup>69,70</sup> (Table S2A). Thus, *COL10A1* is associated with common  
379 matrisome features that drive cancer progression, suggesting that ColX-associated genes may  
380 broadly play important roles in the development, maintenance, and pro-metastatic function of  
381 the tumor microenvironment<sup>1,71</sup>.

382 Multiple hallmark gene sets were significantly enriched in both breast and pancreatic  
383 ColX modules. The epithelial-to-mesenchymal transition (EMT) was the most significantly-  
384 overlapping hallmark gene set in both cancer types ( $p_{\text{adj}} = 5.2 \times 10^{-43}$  and  $3.2 \times 10^{-31}$  for breast  
385 and pancreatic cancer, respectively;  $p_{\text{adj}} = 8.0 \times 10^{-40}$  and  $3.0 \times 10^{-17}$  for male and female  
386 pancreatic cancer cohorts, respectively) (Figure 2D; Table S4). We observed substantial co-  
387 modularization of *COL10A1* with numerous collagen-binding integrin genes which play critical  
388 roles in invasion and blood vessel remodeling, including *ITGA1* (BRCA), *ITGA11* (BRCA, PAAD,  
389 and male PAAD), *ITGA5* (PAAD and male PAAD), *ITGAV* (BRCA), and *ITGB1* (BRCA and  
390 female PAAD) (Table S2A). The BRCA and male PAAD ColX modules additionally contain  
391 *DDR2*, a *COL10A1* receptor which has been shown to sustain the EMT phenotype to promote  
392 metastasis in breast cancer as well as induce EMT to accelerate the progression of pancreatic

393 cancer<sup>72,73</sup>, consistent with fibrosis and metastasis-associated EMT pathway genes being most  
394 strongly enriched in those two ColX modules. Notable overlap was also observed across both  
395 cancer types with numerous other hallmark gene sets related to cancer cell motility and aberrant  
396 cell repair, including angiogenesis, coagulation, apical junction, myogenesis, and decreased UV  
397 repair response. The lists of EMT hallmark genes present in both breast and pancreatic ColX  
398 modules includes not only multiple collagen genes but also various genes coding for ECM  
399 proteins involved in non-collagenous network formation and cartilage and bone development  
400 (*ADAM12*, *COMP*, *MATN3*, *FN1*) (Figure S2E). Several of these genes are similarly present in  
401 both male and female pancreatic ColX modules, indicating conservation of pro-metastatic  
402 function across gender (Figure S2F). The pairwise concordance in EMT hallmark genes overlap  
403 between these groups is statistically significant by Fisher's exact test ( $p = 6.2 \times 10^{-17}$  between  
404 BRCA and PAAD;  $p = 4.6 \times 10^{-8}$  between male and female PAAD), suggesting conservation of  
405 ColX module function across cancers and genders.

406 Several additional cancer type-specific hallmarks emerged as significant, including  
407 genes related to KRAS in breast cancer, and Notch signaling in pancreatic cancer. Interestingly,  
408 there were gender-specific differences within the pancreatic ColX modules, with the male ColX  
409 module significantly overlapping with genes involved in Hedgehog and Notch signaling, while  
410 neither of these gene sets overlapped significantly with the female ColX module (Figure 2D). In  
411 the female pancreatic cancer cohort, Hedgehog signaling was not significantly enriched in any  
412 module and Notch signaling was only enriched in module #12 (Table S4D), suggesting a  
413 stronger association between ColX expression and developmental pathways in male pancreatic  
414 tumors. Together, these observations provide evidence that ColX-associated gene signatures  
415 may confer more aggressive tumor features through activity of specific developmental  
416 pathways. Combined with the IHC observation that ColX $\alpha$ 1 is expressed in tumor stroma  
417 (Figure 1B), these findings suggest that *COL10A1* expression is associated with an oncogenic  
418 fibroblast environment.

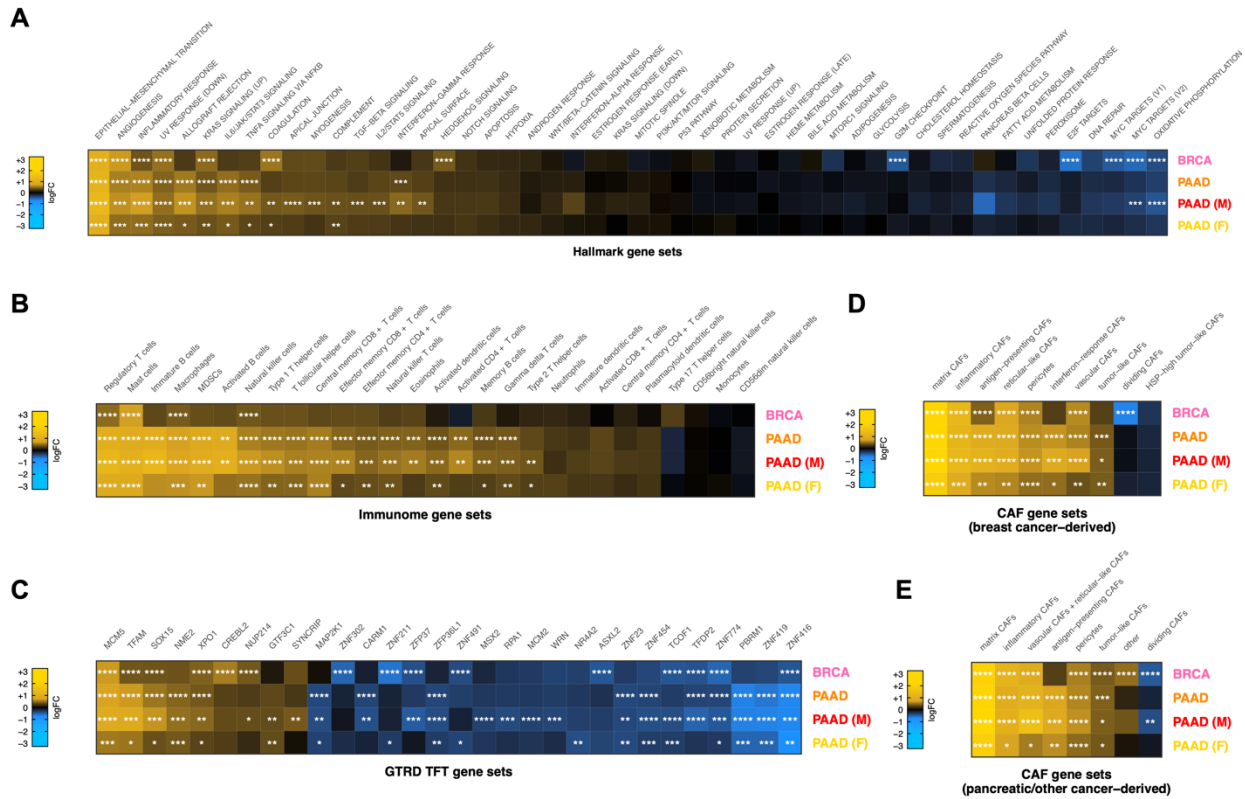
419

## 420 **ColX is associated with regulation of proliferative and mesenchymal cell states**

421 To identify potential impacts on transcription factors (TFs), we assessed representation  
422 of transcription factor target (TFT) gene sets from the Gene Transcription Regulation Database  
423 (GTRD) in each ColX module. Two GTRD-annotated TFT gene sets were significantly enriched  
424 in both BRCA and PAAD ColX modules, corresponding to targets of *FOXH1* and *IGLV5-37*  
425 (Figure 2E). Corroborating the GO enrichment of EMT and related pro-metastatic pathways,  
426 *FOXH1* is an inducer of the TGF- $\beta$ /Nodal/Activin signaling pathway and has been implicated in  
427 proliferation, migration, and invasion of both breast and pancreatic cancer<sup>74,75</sup> (Table S5A). Of  
428 note, the latter TFT gene set appears to have been misattributed to *IGLV5-37* (an  
429 immunoglobulin with no known TF activity) in the GTRD database, and according to the source  
430 publication actually represents targets of the fusion oncoprotein *SS18-SSX*, which alters the  
431 normal regulatory activity of the SWI/SNF (BAF) ATP-dependent chromatin remodeling complex  
432 to drive oncogenesis in synovial sarcoma through induction of *SOX2*<sup>76</sup>. Both of these TFT gene  
433 sets include ColX module genes *ADAMTS6* (a metalloproteinase) and *RUNX1*, which drives  
434 mesenchymal stem cell proliferation and differentiation of myofibroblasts<sup>77</sup>; additionally, the  
435 *SS18-SSX* targets include *LRRC15*, a ColX module gene marking TGF- $\beta$ -driven,  
436 myofibroblastic CAFs which may play an immunoregulatory role in PDAC<sup>78</sup> (Table S5A). The  
437 shared enrichment of these two TFT gene sets in both ColX modules points toward a common  
438 association with fibroblast proliferation and invasion in breast and pancreatic cancer.

439 Among the BRCA ColX module genes, transcription targets of *PRMT5*, *SIX1*, *SUPT16H*,  
440 *PRDM4*, *TFEB*, *NFKBIA*, *ZNF585B*, and *PRDM5* were highly enriched (Figure 2E). These TFs  
441 have been established to regulate numerous pathways in breast and other cancers, ranging  
442 from cell proliferation and tumorigenesis to invasion and immune regulation (Table S5A). In  
443 particular, *SIX1* has been shown to induce EMT and metastasis in breast tumors<sup>79</sup> and  
444 implicated in TGF- $\beta$  regulation of collagen deposition leading to hampered immune infiltration

445 and survival in cancer<sup>80</sup>, while *PRDM5* has been linked to pro-metastatic production of collagen  
 446 in breast cancer<sup>81</sup>. The PAAD CoIX module is enriched for targets of the transcription factor  
 447 *NME2*, which regulates cancer cell proliferation in a variety of contexts through activation of  
 448 *MYC* (Table S5A). Of note, *NME2* is an upstream regulator of *CTSK*, a protease primarily  
 449 expressed in osteoclasts which is involved in ECM degradation and bone remodeling and has  
 450 been found to be expressed in numerous cancers including breast carcinoma<sup>82,83</sup>.



451

452 **Figure 3: Tumors with high ColX module expression exhibit activation of pro-metastatic,**  
453 **immunosuppressive, and myofibroblastic gene signatures. (A–E)** Mean pathway activation (by  
454 QuSAGE) of **(A)** MSigDB hallmark pathway gene sets<sup>54</sup>, **(B)** cancer immunome gene sets<sup>59</sup>, **(C)** top  
455 significant GTRD transcription factor target gene sets<sup>55</sup>, and cancer-associated fibroblast (CAF) gene sets  
456 derived from **(D)** breast tumors and **(E)** pancreas and other tumors<sup>19</sup>, in samples with high ColX module  
457 expression relative to samples with low ColX module expression for each cancer dataset. Gene sets  
458 within each block are ordered by mean pathway activation rank across all 4 datasets. See **Figure S4** for  
459 QuSAGE analysis of all modules and **Table S5B** for reported functions of selected enriched TFTs.  
460 Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*, p.adj < 0.0001.

461 **Increased ColX module expression is associated with activation of pro-metastatic**  
462 **pathways as well as immunosuppressive and myofibroblastic signatures**

463 To investigate the physiological states associated with expression of ColX-related gene  
464 networks, we stratified patient samples into high and low ColX module expression groups and  
465 performed Quantitative Set Analysis for Gene Expression (QuSAGE) to identify differentially-  
466 expressed pathways between these cohorts (see Methods section and **Figure S3** for definition  
467 and derivation of %G.A.M.E. metric). QuSAGE is an efficient alternative to classical Gene Set  
468 Enrichment Analysis (GSEA) which accounts for inter-gene correlations and provides a more  
469 robust quantification of differential pathway expression by generating a complete probability  
470 density to describe the activity of a particular gene set of interest<sup>61–63</sup>. We ran QuSAGE on ColX  
471 module-stratified samples for the MSigDB hallmark pathways<sup>54</sup>, a collection of previously-  
472 characterized immune signature pathways<sup>59</sup>, the MSigDB GTRD TFT gene sets<sup>54,55</sup>, and a  
473 collection of CAF genesets curated from human breast and pancreatic cancer datasets<sup>19</sup>.

474 Increased activity of hallmark pathways associated with aggressive cancer phenotypes  
475 was significantly associated with ColX module expression (**Figure 3A**). Heightened ColX module  
476 expression was positively associated with activation of genes involved in EMT (2.5-fold average  
477 increase across all cohorts), angiogenesis (1.7-fold average increase), coagulation (1.6-fold  
478 increase in BRCA; also significant in male and female PAAD cohorts), and myogenesis, apical  
479 junction and apical surface (each 1.5-fold increase in male PAAD cohort). Additionally, QuSAGE  
480 revealed increased activity of genes typically decreased in response to UV exposure across all  
481 cohorts (1.5-fold average increase). We also observed significant upregulation in high-ColX

482 module samples of several regulatory pathways involved in tumor progression and immune  
483 response that were not identified through gene enrichment testing, including inflammatory  
484 response and KRAS signaling (all cohorts), Hedgehog signaling (BRCA), allograft rejection, IL-  
485 6/JAK/STAT3 signaling, and TNF- $\alpha$  signaling via NFKB (all PAAD cohorts), IFN- $\gamma$  response  
486 (PAAD and male PAAD cohorts), complement (male and female PAAD cohorts), and IL-  
487 2/STAT5 signaling and TGF- $\beta$  signaling pathways (male PAAD cohort). Several hallmark gene  
488 sets related to proliferation were downregulated in high-ColX module BRCA samples, including  
489 oxidative phosphorylation, MYC targets (V1, V2), E2F targets, and G2/M checkpoint; the high-  
490 ColX module male PAAD cohort also exhibited downregulation of MYC targets (V2) and  
491 oxidative phosphorylation. These findings indicate that tumors with high *COL10A1* expression  
492 exhibit increased activity of numerous pathways related to pro-metastatic and inflammatory  
493 processes, along with a concomitant decrease in activity of homeostatic pathways such as  
494 oxidative metabolism and cell cycle regulation. Thus, activity of ColX module genes highlights  
495 more aggressive cancer phenotypes.

496 ColX module expression was also significantly associated with increased expression of  
497 numerous immune signatures in both BRCA and PAAD, including immunosuppressive and  
498 tumor-promoting features such as regulatory T cells, mast cells, and tumor-associated  
499 macrophages (TAMs) (Figure 3B). In pancreatic cancer specifically, we observed broader  
500 upregulation of numerous immune cell types, but notably no upregulation of the activated CD8<sup>+</sup>  
501 T cell signature, and greater upregulation of immature vs. activated B-cell signatures. Although  
502 the activation of regulatory T cells was not uniquely associated with ColX module expression,  
503 we note that other WGCNA modules tracking positively with regulatory T cell activity were  
504 almost invariably also associated with activated CD8<sup>+</sup> T cell signal as expected in controlled  
505 physiological immune responses, while the ColX modules were not (Figure S4A–D). ColX  
506 module expression was consistently associated with regulatory T cell activity even in the

507 absence of activated CD8<sup>+</sup> T cells. These findings suggest that the ColX modules are strongly  
508 associated with immunosuppressive environments.

509 QuSAGE analysis of GTRD TFT gene sets highlighted numerous transcription factors  
510 whose downstream targets were differentially expressed in tumors with high ColX module  
511 expression, several of which are known to modulate tumor advantage in breast or pancreatic  
512 cancer (Figure 3C; see Table S5B for detailed descriptions of significant TFs). Targets of  
513 *MCM5*, a crucial component of the MCM replicative helicase complex which has been linked to  
514 negative prognosis in breast cancer and implicated as a marker of pancreatic malignancy, were  
515 the most significantly upregulated in both BRCA and PAAD high-ColX module tumors, an effect  
516 which was also observed in the pancreatic gender-specific analysis. Expression of *NME2*  
517 targets was also increased significantly in high-ColX module PAAD tumors; this protein is a well-  
518 established activator of *MYC* and while competing evidence suggests that it may have pro- or  
519 anti-tumor effects in diverse contexts, it has been identified as upregulated in metastatic PDAC  
520 samples compared to primary tumors by snRNA-Seq<sup>84</sup>. *XPO1* and *NUP214*, two associated  
521 proteins whose regulatory targets are mutually upregulated in high-ColX module tumors, have  
522 collectively been described as drivers of breast cancer and markers of poor survival in  
523 pancreatic cancer. *NUP214* is of particular interest as its varied fusion products have been  
524 implicated as drivers of breast cancer (via fusion with *NOTCH*)<sup>85</sup> and leukemogenesis (via  
525 aberrant activation of HOX genes)<sup>86</sup>; similarly, *HOXC8* co-modularizes with *COL10A1* in BRCA  
526 (Table S2A) and *Hoxa3* has been shown to be upregulated alongside *Col10a1* and implicated in  
527 OA progression in hypertrophic mouse chondrocytes<sup>87</sup>. Targets of *GTF3C1* and *SYNCRIP*, both  
528 of which are upregulated in metastatic pancreatic tumors compared to primary tumors, were  
529 also significantly increased in expression in high-ColX module gender-specific PAAD cohorts.

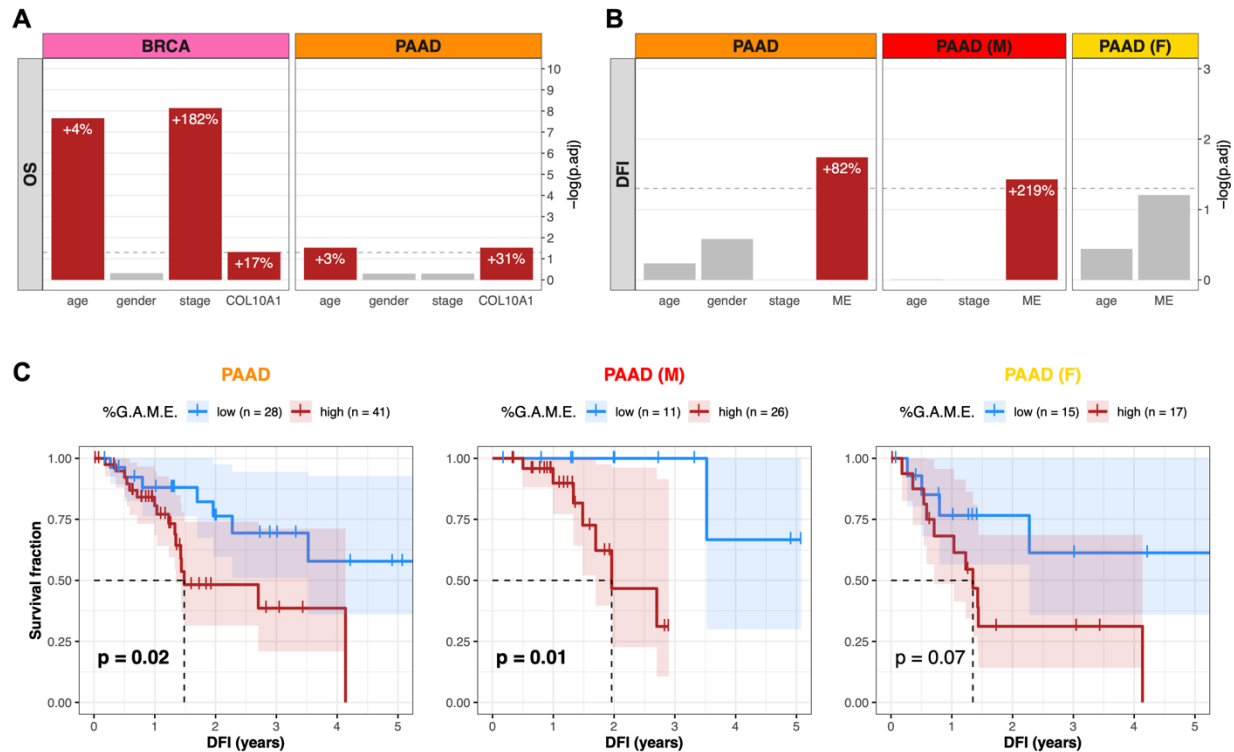
530 Conversely, targets of multiple key genes crucial to the maintenance of various tumor-  
531 suppressive protein complexes were significantly downregulated in high-ColX module samples  
532 (Figure 3C), including *ASXL2* in BRCA and *PBRM1* in PAAD (irrespective of gender). Significant

533 downregulation in high-ColX module tumors was also observed for targets of numerous zinc-  
534 finger TFs (*ZNF302*, *ZNF211*, *ZFP37*, *ZNF23*, *ZNF454*, *ZNF774*, *ZNF419*, and *ZNF416*)  
535 associated with *TRIM28*, a pro-tumorigenic driver of EMT which stabilizes *TWIST1* to promote  
536 cancer cell invasion and migration<sup>88</sup>. While individual coexpression of each significant zinc-  
537 finger TF with *TRIM28* and *TWIST1* varied across cohorts, on average they were negatively  
538 correlated with expression of both pro-metastatic genes in BRCA and PAAD, potentially  
539 corroborating the increased activity of EMT drivers in high-ColX module tumors.

540 ColX module expression was additionally associated with increased activity of multiple  
541 CAF gene sets. We analyzed gene signatures from diverse CAF populations recently  
542 characterized by Cords et al. using scRNA-Seq data from breast tumors (~14,000 CAFs from 14  
543 patients with breast invasive carcinoma) and pancreatic tumors (~5,700 CAFs from 4 cancer  
544 types, of which 30% were sourced from patients with pancreatic ductal adenocarcinoma)<sup>19</sup>. As  
545 these gene signatures have been shown to be highly conserved across cancer types, we  
546 performed matched and cross-comparisons between breast- and pancreatic-derived CAF  
547 signatures and our breast and pancreatic cancer data, stratified by ColX module expression. We  
548 found that gene sets associated with numerous breast-derived myofibroblastic CAF populations  
549 (matrix, inflammatory, and antigen-presenting) were significantly upregulated in high-ColX  
550 module samples compared to low-ColX module samples in both BRCA and PAAD cohorts, as  
551 were genes associated with vascular CAFs, which have been shown to facilitate angiogenesis  
552 and tumor vascularization (Figure 3D). High ColX module expression in the PAAD cohorts was  
553 uniformly associated with increased interferon-response and tumor-like CAF signatures as well.  
554 Matrix CAFs (mCAFs) play a role in ECM remodeling, migration, TGF- $\beta$ -driven myofibroblastic  
555 activation, and EMT; numerous mCAF markers are present in the ColX modules, including  
556 *COMP*, *MMP11*, *POSTN*, *COL1A1*, *COL1A2*, *LRRRC15*, and the pro-myofibroblastic markers  
557 *FAP* and *PDPN* (Table S2A). Various other CAF-specific genes are co-modularized with  
558 *COL10A1*, including inflammatory CAF markers *CXCL12* (BRCA) and *CXCL14* (female PAAD),



559 vascular CAF markers *ACTA2* (all ColX modules) and *NOTCH3* (PAAD and male PAAD), and  
560 tumor-like CAF markers *PDPN* (BRCA, PAAD, and male PAAD) and *TMEM158* (male PAAD).  
561 Additionally, the dividing CAF signature was significantly decreased in high-ColX module BRCA  
562 tumors, corroborating the decrease in homeostatic cell cycle control indicated by the hallmark  
563 pathway analysis of the BRCA cohort (Figure 3D). Similar enrichment of pancreatic-derived  
564 myofibroblastic CAF signatures in high-ColX module tumors was observed in both BRCA and  
565 PAAD cohorts, concomitant with increased signal associated with tumor-like CAFs and  
566 decreased signal associated with dividing CAFs (Figure 3E). These observations further support  
567 the theme of ColX being a marker of CAF heterogeneity in stromal environments across diverse  
568 cancers.



569

570 **Figure 4: COL10A1 and ColX module expression stratify breast and pancreatic cancer cohorts by**  
571 **survival outcomes. (A and B)** BH-adjusted significance values for multivariate Cox proportional hazards  
572 models conditioning either **(A)** overall survival (OS) on age, gender, binarized tumor stage, and COL10A1  
573 gene expression, or **(B)** disease-free interval (DFI) on age, gender, binarized tumor stage, and ColX  
574 module eigengene expression (ME). Dotted lines correspond to  $p_{adj} = 0.05$ . Per-unit contributions of  
575 each significant variable (red bars) to the overall hazard risk in each panel is indicated by white  
576 percentages (for COL10A1 and ColX ME expression, 1 standard deviation = 1 unit). Note that the Cox  
577 model for the female pancreatic cancer cohort was not conditioned on binarized tumor stage because all  
578 samples were classified into the same group. See **Figure S5** for full survival analysis results for all  
579 WGCNA-inferred modules. **(C)** DFI Kaplan-Meier survival curves for the pancreatic cancer cohorts based  
580 on %G.A.M.E. groupings. Dotted lines correspond to median “survival” (i.e., time to recurrence). Shaded  
581 regions represent the 95% confidence intervals for each group. Log-rank test p-values are shown for  
582 each panel.

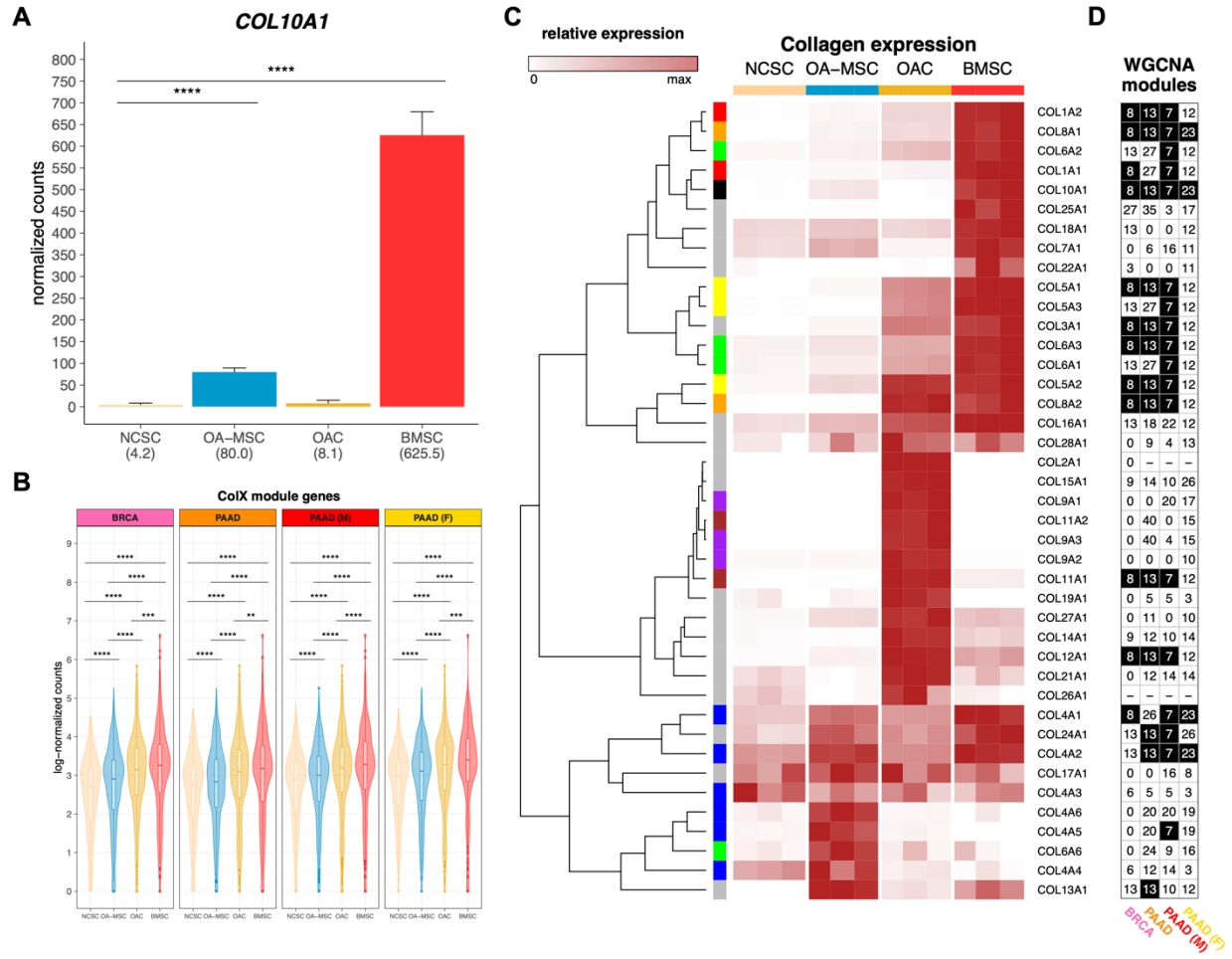
### 583 ColX gene and module expression levels are prognostic of differential survival outcomes

584 Given the relevance of the stromal microenvironment to cancer progression and patient  
585 outcomes, we assessed the predictive value of ColX with regard to breast and pancreatic  
586 survival metrics. Using a multivariate Cox proportional hazards model to assess overall survival  
587 risk, increased COL10A1 expression was found to be significantly associated with negative  
588 prognosis (+17% risk per standard deviation in breast cancer,  $p_{adj} = 0.048$ ; +31% risk per

589 standard deviation in pancreatic cancer,  $p_{\text{adj}} = 0.030$ ) (Figure 4A). These hazard contributions  
590 were conferred in addition to the significant effects of age (+4% risk per year in breast cancer,  
591  $p_{\text{adj}} = 2.2 \times 10^{-8}$ ; +3% risk per year in pancreatic cancer,  $p_{\text{adj}} = 0.030$ ) and advanced tumor  
592 stage (+182% risk for stages 3-4 compared to stages 1-2 in breast cancer,  $p_{\text{adj}} = 7.3 \times 10^{-9}$ ).

593 To determine whether ColX-associated gene modules preserved the prognostic value of  
594 ColX itself, we performed a similar analysis using the ColX module eigengene (ME) in place of  
595 *COL10A1* expression. Increased ColX module expression was significantly associated with  
596 shortened disease-free interval (DFI) in pancreatic cancer (+82% risk per standard deviation,  
597  $p_{\text{adj}} = 0.018$ ) (Figure 4B). This effect was also significant in the male pancreatic cancer  
598 subcohort (+219% risk per standard deviation,  $p_{\text{adj}} = 0.037$ ), but was not preserved in the  
599 female subcohort ( $p_{\text{adj}} = 0.062$ ). The ColX ME did not significantly impact DFI in breast cancer,  
600 possibly due to the greater availability and efficacy of curative therapies for breast cancer  
601 relative to pancreatic cancer. However, the fact that the ColX ME was one of only two modules  
602 in the pancreatic cancer cohort (as well as in the male subset) whose expression tracked  
603 significantly with increased DFI risk suggests that the ColX signature is uniquely associated with  
604 likelihood of recurrence of advanced cancer (Figures S5B, S5C).

605 We then investigated whether the rescaled metric of ColX module expression  
606 (%G.A.M.E.) could be used to effectively stratify pancreatic cancer patients based on DFI.  
607 Patients with high %G.A.M.E. had significantly worse DFI prognosis compared to those with low  
608 %G.A.M.E. ( $p = 0.02$ ) by Kaplan-Meier analysis. Interestingly, this effect was preserved in the  
609 male subcohort ( $p = 0.01$ ), but attenuated in the female subcohort ( $p = 0.07$ ) (Figure 4C). These  
610 findings recapitulated the multivariate hazard risk analysis, corroborating the prognostic value of  
611 ColX expression at both the gene and module level.



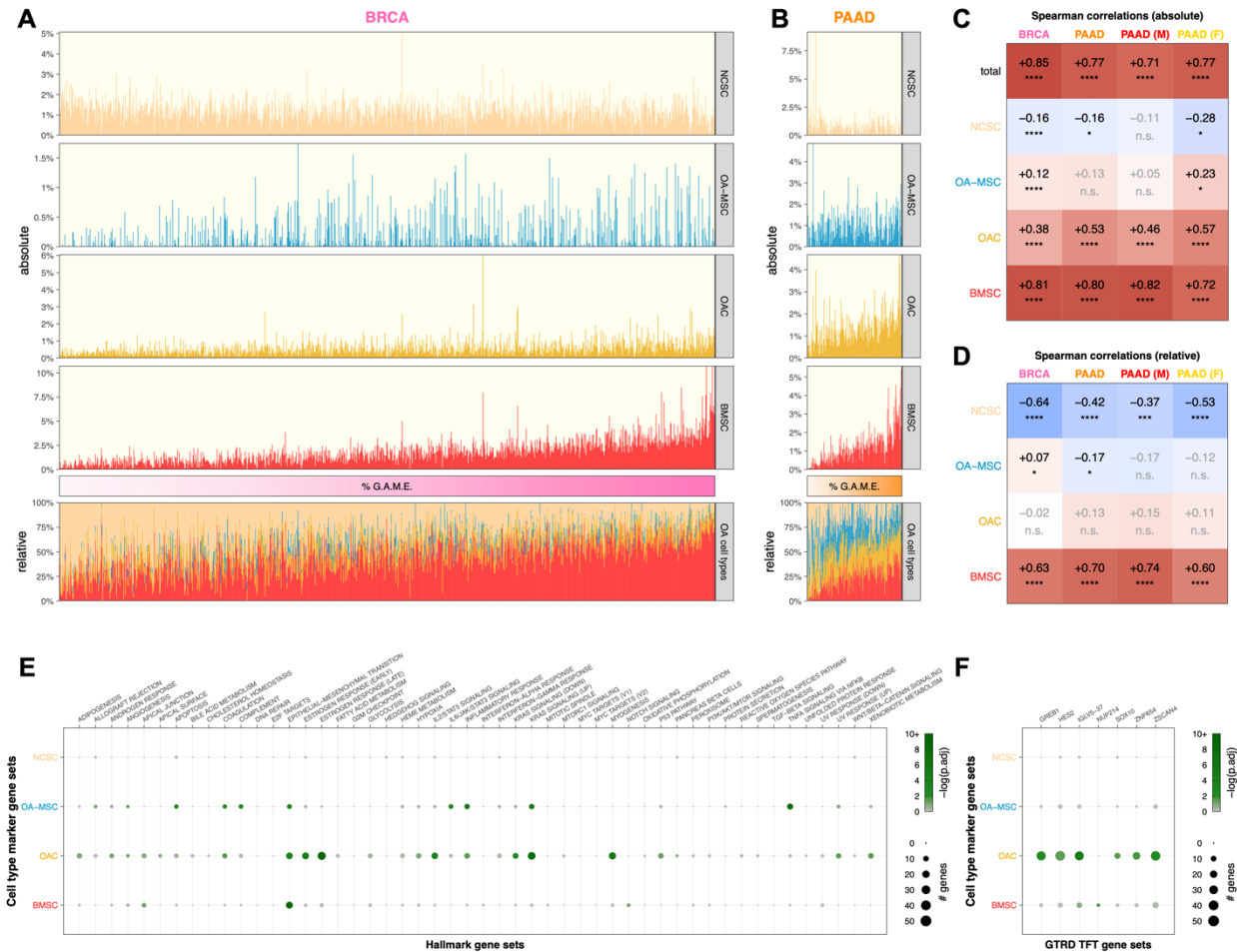
612

613 **Figure 5: Bone marrow and cartilage cells are differentiated by collagen expression clusters. (A)**  
614 Normalized expression of *COL10A1* across bone marrow and articular cartilage cell types ( $n = 3$  each).  
615 Mean normalized expressions are indicated in parentheses. Significance values were computed using  
616 DESeq2, BH-adjusted across all genes analyzed. **(B)** Log-normalized expression of ColX module genes  
617 in each cell type ( $n = 3$  each). Normalized gene-wise expressions were averaged prior to log-  
618 transformation with pseudocount of 1. Significance values were computed using the paired Wilcoxon  
619 signed-rank test on log-normalized counts. **(C)** Relative expression of collagen genes across cell samples  
620 ( $n = 3$  each). Normalized expression values are scaled by rows. *COL10A1* is indicated by a black box in  
621 the left column; genes contributing to the same parent collagen are indicated by same-color boxes; and  
622 genes which are the sole contributor to their parent collagen are indicated by gray boxes. Note that  
623 *COL6A4P1*, *COL6A5*, *COL20A1*, and *COL23A1* were filtered out as “low-expression” genes. See [Figure](#)  
624 [S7A](#) for normalized (unscaled) expression of all collagen genes across cell types. **(D)** WGCNA module  
625 assignments for all collagens in each cancer dataset. See [Table S2B](#) for WGCNA module assignments  
626 for all genes. Black cells indicate the ColX module for each column. “-” indicates low-expression genes  
627 which were filtered out; genes labeled “0” were not assigned to any module by WGCNA. See [Figure S7B](#)  
628 for normalized expression of all collagen genes across TCGA cohorts.

### 629 **ColX module expression correlates with increased activity of bone marrow stroma and** 630 **osteoarthritic cartilage signatures**

631 Identifying connections between the roles of ColX in cancer and non-cancer cells may  
632 provide novel insights into common factors at play in these tissues. In its canonical contexts,  
633 *COL10A1* is a marker for hypertrophic chondrocytes and bone marrow-derived mesenchymal  
634 stem cells and contributes to cellular senescence, ECM degradation, and angiogenic invasion  
635 during OA<sup>6,9,11,89–92</sup>. Therefore, we examined the bone marrow and OA character of breast and  
636 pancreatic tumors based on their levels of ColX module expression, using RNA-Seq data from  
637 four types of cartilage or bone marrow stromal cells for comparison<sup>10</sup>. While *COL10A1* was not  
638 strongly expressed by normal cartilage stromal cells (NCSCs), it was significantly upregulated in  
639 both bone marrow stromal cells (BMSCs; over 100-fold increase,  $p_{\text{adj}} = 1.1 \times 10^{-31}$ ) and OA  
640 mesenchymal stromal cells (OA-MSCs; 18.6-fold increase,  $p_{\text{adj}} = 3.4 \times 10^{-11}$ ) ([Figure 5A](#)).  
641 Interestingly, the module-wide expression of *COL10A1*-associated genes in each cancer was  
642 significantly increased in BMSCs, followed by two OA cartilage cell types, OA chondrocytes  
643 (OACs) and OA-MSCs ([Figure 5B](#)); additionally, cell type-specific markers for BMSCs and  
644 OACs were found to be highly enriched for genes in the ColX modules ([Figure S6A–D](#)). These  
645 enrichments were significant by Fisher’s exact test when compared across all WGCNA modules

646 for both BRCA (p.adj =  $1.6 \times 10^{-15}$  for BMSC and p.adj =  $4.7 \times 10^{-3}$  for OAC) and PAAD (p.adj =  
647  $1.7 \times 10^{-11}$  for BMSC and p.adj = 0.02 for OAC) (Figure S6E). Comparative analysis of the  
648 expression of various collagens across cartilage and bone marrow cell types revealed that the  
649 collagen gene signatures of BMSCs were closely aligned with the ColX modules in BRCA and  
650 PAAD, in particular male but not female PAAD (Figures 5C, 5D). Numerous BMSC-specific  
651 collagens also co-modularized with *COL10A1* in tumors, including *COL8A1* and *COL8A2*,  
652 *COL1A1* and *COL1A2*, *COL4A1* and *COL4A2*, and several genes contributing to collagen types  
653 III, V and VI. These latter three collagens have all been implicated in tissue repair, wound  
654 healing, expression profiles of diverse CAF subtypes, and regulation of both BMSC and OAC  
655 physiology<sup>6,10,93</sup>. Like *COL10A1*, collagen type IV is a network-forming collagen; its  $\alpha112$  triple  
656 helical form is ubiquitously expressed in basement membranes and is present in both normal  
657 and OA articular cartilage, where it is expressed by chondrocytes<sup>6,94</sup>. *COL8A1* is highly  
658 expressed in OA tissues and, like *COL10A1*, is a myofibroblastic-specific marker in cancer<sup>93,95</sup>.  
659 These results suggest that, in addition to *COL10A1*, BRCA, PAAD and male PAAD tumors  
660 especially share common ECM gene signatures with BMSCs.



661

662 **Figure 6: Tumoral proportions of osteoarthritic cell types trend with ColIX module expression and**  
 663 **cancer-associated pathway activity. (A and B)** Sample-wise absolute (*top*) and relative (*bottom*)  
 664 proportions of 4 bone marrow and cartilage cell types inferred by CIBERSORTx for TCGA **(A)** breast and  
 665 **(B)** pancreatic cancer cohorts. Color bars (*middle*) indicate relative cancer-specific ColIX module  
 666 expression; samples were ordered by increasing %G.A.M.E. See Methods section and **Figure S3** for  
 667 details on %G.A.M.E. metric. **(C and D)** Spearman correlations between relative cancer-specific ColIX  
 668 module expression and **(C)** absolute or **(D)** relative OA cell type proportions inferred by CIBERSORTx.  
 669 Raw p-values were Bonferroni-corrected across rows. Significance values: \*,  $p < 0.05$ ; \*\*,  $p < 0.01$ ; \*\*\*,  $p$   
 670  $< 0.001$ ; \*\*\*\*,  $p < 0.0001$ ; n.s., not significant. NCSC, normal cartilage stromal cells; OA-MSc,  
 671 osteoarthritis mesenchymal stromal cells; OAC, osteoarthritis chondrocytes; BMSc, bone marrow stromal  
 672 cells. **(E and F)** Bubble plots of **(E)** MSigDB hallmark pathway gene sets<sup>54</sup> and **(F)** top significant GTRD  
 673 transcription factor target gene sets<sup>55</sup> in cell type-specific marker gene sets. Of note, the TFT gene set  
 674 attributed to *IGLV5-37* in the GTRD database actually represents targets of the fusion oncoprotein SS18-  
 675 SSX, as described in the text. See Methods section for definition of bone marrow and cartilage cell type-  
 676 specific marker genes. Significance values: \*,  $p_{\text{adj}} < 0.05$ ; \*\*,  $p_{\text{adj}} < 0.01$ ; \*\*\*,  $p_{\text{adj}} < 0.001$ ; \*\*\*\*,  $p_{\text{adj}} <$   
 677  $0.0001$ .

678

Cell type proportion inference using CIBERTSORTx revealed that expression of ColX

679

modules was positively associated with signatures attributable to BMScs, OACs, and OA-MSCs

680 across both breast and pancreatic tumors (Figures 6A, 6B). ColX module expression was  
681 significantly and positively correlated with CIBERSORTx-inferred absolute proportions of  
682 BMSCs as well as, to a lesser extent, two OA cell types (OACs and OA-MSCs) (Figure 6C).  
683 Conversely, ColX module expression was negatively correlated with the absolute inferred  
684 proportion of NCSCs in breast cancer; correlations of similar direction and magnitude were  
685 observed in pancreatic cancer, with some loss of significance likely attributable to decreased  
686 power due to the smaller cohort size (Figure 6C). The relative inferred proportion of NCSCs  
687 among the 4 cell types also decreased significantly with increased ColX module expression in  
688 both cancers, concomitant with a significant increase in the relative inferred proportion of  
689 BMSCs (Figure 6D). Thus, ColX emerged as a biomarker for bone marrow-derived cells in  
690 BRCA and PAAD tumors.

691

#### 692 **ColX highlights similar EMT markers between OA disease states and tumors**

693 EMT comprises a range of markers depending on tissue settings and surrounding  
694 phenotypes, and contributes to diverse processes including gastrulation, fibroblastic  
695 differentiation, and tumor metastasis<sup>96</sup>. Examination of cell type-specific markers for each ColX-  
696 expressing non-cancer cell type (Table S6) revealed that EMT hallmark pathway genes were  
697 significantly overrepresented in BMSCs, OACs, and OA-MSCs (Figure 6E). Notably, several  
698 EMT-associated genes were identified as both OA cartilage cell type markers and ColX module  
699 genes. *COL11A1*, *COMP*, and *MATN3*, which were present in both breast and pancreatic ColX  
700 modules (Table S2A; Figures S2E, S2F), are markers of chondrogenesis expressed distinctly by  
701 OACs (Table S6); additionally, *COMP* is a potent anti-apoptotic factor which may contribute to  
702 survival of both chondrocytes and neoplastic cells. The BRCA ColX module also includes the  
703 OAC-specific markers *LUM* (also present in the PAAD ColX module) and its related genes *DCN*  
704 and *ECM2*, which are collectively involved in collagen fibril organization and epithelial cell  
705 migration (Table S2A). Genes coding for various CXC chemokines, inflammatory cytokines, and

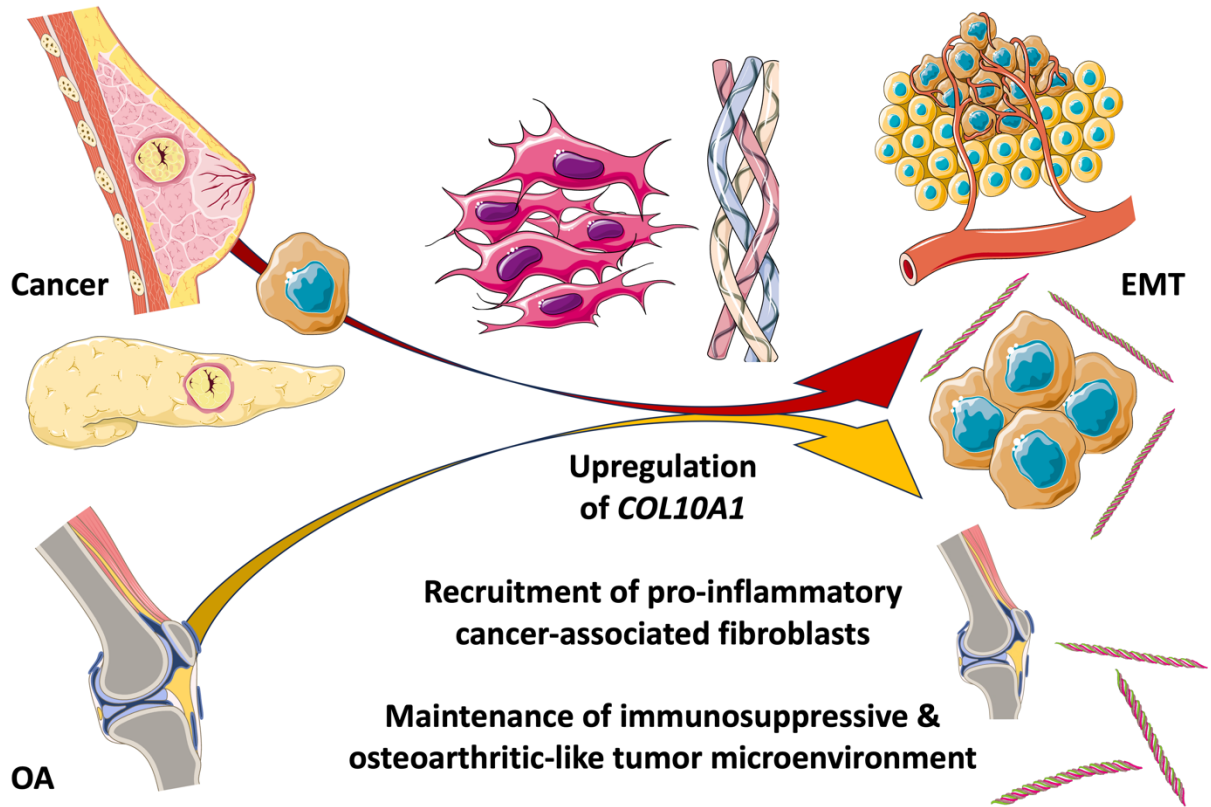


706 matrix metalloproteinases were identified as EMT-associated markers in OA-MSCs (*CXCL1*,  
707 *CXCL6*, *CXCL8*, *IL6*, *MMP1*, *PTX3*, *TNFAIP3*) (Table S6); while these were not specifically co-  
708 modularized with *COL10A1* in tumors, several related genes in both the CXC and MMP families  
709 are present in the ColX modules (pro-metastatic *CXCL12* in BRCA, pro-angiogenic *MMP2* in  
710 both BRCA and PAAD, and related MMPs 2, 3, and 14 in male PAAD) (Table S2A; Figures  
711 S2E, S2F).

712 Significant overlap in EMT pathway genes was observed between BMSC-specific  
713 markers and the ColX modules; 14 BMSC markers co-modularized with *COL10A1* in BRCA ( $p =$   
714  $1.8 \times 10^{-5}$  by Fisher's exact test), 10 of which also co-modularized in PAAD ( $p = 2.8 \times 10^{-3}$  by  
715 Fisher's exact test) (Figure S6F). Several of these overlapping EMT pathway genes are also  
716 highly expressed in matrix CAFs (*COL1A2*, *LRRC15*, *POSTN*); other overlapping genes have  
717 been shown to contribute to vascular remodeling (*ACTA2*, *CTHRC1*) and myofibroblastic  
718 motility (*CALD1* and *TPM1*, both of which are significantly upregulated in metastatic pancreatic  
719 cancer compared to primary tumors<sup>84</sup>) (Table S2A; Table S6; Figures S2E, S2F). In particular,  
720 *FN1*, *VCAN*, and *LRRC15* are markers of BMSCs which also co-modularized with *COL10A1* in  
721 cancer; all three genes are involved in cell migration, and *FN1* contributes directly to osteoblast  
722 mineralization as well as metastasis. Together, these results suggest a shared activation of  
723 common EMT-related pathways in cancer, bone marrow, and pathological OA contexts,  
724 highlighting the role of *COL10A1* as a marker of both BMSCs and pro-metastatic tumors. The  
725 ColX modules thus highlight well-characterized inflammatory OA-like disease states which  
726 contribute to a pro-metastatic tumor microenvironment in both breast and pancreatic cancer.

727 Several TFs whose targets were found to be differentially expressed in high-ColX  
728 module vs. low-ColX module tumors have also been reported to play roles in OA disease states  
729 (Figure 3C). Upregulated TFT sets included targets of *MCM5*, whose expression is increased in  
730 rat chondrocytes following inflammatory stimulation by IL-1 $\beta$ ; *TFAM*, which promotes  
731 mitochondrial biogenesis in OA chondrocytes; and *NME2*, which has been shown to be

732 upregulated in early OA (Table S5B). Downregulated TFT sets included targets of the  
733 associated proteins *RPA1*, whose expression is decreased in OA patients, and *WRN*, which is  
734 mutated in Werner's syndrome and confers early-onset OA risk; as well as *PBRM1*, which has  
735 been implicated in GWAS of OA and regulation of BMP signaling and osteogenic fate  
736 determination (Table S5B). Cell type-specific markers for both OACs and OA-MSCs were also  
737 significantly enriched for various inflammatory/immune (e.g., interleukin/STAT protein signaling)  
738 and proliferative (e.g., KRAS signaling) pathways (Figure 6E). Additionally, OAC- and BMSC-  
739 specific markers were found to be enriched for targets of the TFs *SS18-SSX* (whose targets  
740 were misattributed in the GTRD database to the immunoglobulin *IgLV5-37* as noted above),  
741 and *NUP214*, respectively (Figure 6F). These findings corroborate the similar enrichment of  
742 *SS18-SSX* targets in both the breast and pancreatic ColX modules (Figure 2E), suggesting that  
743 common pathways involving dysregulated cell proliferation and tissue remodeling are involved  
744 in both cancer and OA (Table S5A), as well as the differential activation of *NUP214* targets in  
745 high-ColX module breast and male pancreatic tumor samples (Figure 3C), which additionally  
746 suggest a shared pathogenic mechanism linked to aberrant transcription regulation (Table S5B).  
747 Many TFs thus appear to be perturbed similarly in both high-ColX module tumors and OA  
748 development, suggesting similar dysregulation of transcriptional programs across these disease  
749 states.



750

751 **Figure 7: Pathological expression of COL10A1 fosters immunosuppressive, fibroblastic**  
752 **microenvironments in cancer, bone marrow, and cartilage.** Graphical abstract summarizing the  
753 findings presented in this study. In brief, a specifically expressed collagen, COL10A1, connects the ECM  
754 and tissue microenvironments across cancer and bone. The pathological contributions of ColX and its  
755 associated gene networks are especially prominent in breast and pancreatic tumors, and mimic the  
756 development of a similar inflammatory and fibroblast-dominated environment seen in bone marrow and  
757 cartilage changes in OA. A central outcome of this shared impact is the epithelial-to-mesenchymal  
758 transition (EMT), which contributes to disease progression in both contexts. *All images were sourced from*  
759 *Bioicons* (<https://bioicons.com>) and are licensed for public use by Servier (<https://smart.servier.com>)  
760 under CC-BY 3.0 (<https://creativecommons.org/licenses/by/3.0>) or by DBCLS  
761 (<https://togotv.dbcls.jp/en/pics.html>) under CC-BY 4.0 (<https://creativecommons.org/licenses/by/4.0>).

762

## 763 **DISCUSSION**

764 The pathological contributions of collagen type X (COL10A1, ColX) and its associated  
765 gene networks are especially notable in breast and pancreatic tumors, where they contribute to  
766 fostering an inflammatory and immunosuppressive microenvironment that contributes to  
767 aggressive tumorigenesis, metastasis, and poor clinical prognosis. A significant component of  
768 this pathological role appears to be related to ColX's identity as a marker for bone marrow-

769 derived cells and cancer-associated fibroblasts (CAFs), both of which are strongly associated  
770 with extracellular matrix (ECM) remodeling, the epithelial-to-mesenchymal transition (EMT), and  
771 progression of disease in both cancer and OA. The findings of this study, summarized in a  
772 graphical abstract for visualization ([Figure 7](#)), suggest that there is a close relationship between  
773 these two ColX-expressing stromal cell populations; one in breast and pancreatic cancer and  
774 the other in bone marrow and OA cartilage.

775 ColX is normally expressed only during skeletal development at the cartilage-to-bone  
776 transition and in bone marrow during adulthood. However, ColX is also highly expressed in  
777 certain disease contexts including OA cartilage and a variety of solid tumors, especially breast  
778 and pancreatic cancer. Here, we demonstrated that ColX expression is aberrantly elevated in  
779 solid tumors from TCGA; namely, breast invasive carcinoma (BRCA) and pancreatic  
780 adenocarcinoma (PAAD). From RNA-Seq analysis, we found that ColX is also highly expressed  
781 by bone marrow stromal cells (BMSCs) and senescent mesenchymal stromal cells (OA-MSCs)  
782 in OA. Recent work demonstrates that chromatin accessibility of *COL10A1* is highest in  
783 hypertrophic chondrocytes and increases with progression of chondrocyte maturation in skeletal  
784 development<sup>97</sup>. Indeed, in BRCA tumors we observed that *COL10A1* expression is  
785 anticorrelated with the somitic mesoderm marker *TCF15* (Spearman's  $\rho = -0.16$ ,  $p = 4.5 \times 10^{-8}$ )  
786 and positively correlated with the sclerotome marker *PAX9* (Spearman's  $\rho = +0.31$ ,  $p = 2.0 \times 10^{-$   
787 <sup>26</sup>), supporting its contribution to the development of paraxial mesoderm and ossification in  
788 cancer. ColX's patterns of co-expression and variability across breast and pancreatic tumors  
789 suggest that it may be associated with more aggressive disease as indicated by activation of  
790 EMT and angiogenic pathways as well as heightened relapse and survival risk.

791 There are compelling reasons to extrapolate from the roles that *COL10A1* plays in OA to  
792 putative functions within the tumor microenvironment. ColX is known to contribute to calcified  
793 zones of articular cartilage when expressed by hypertrophic chondrocytes in OA, a phenotype  
794 similar to mechanical stiffness conferred by the stroma of many solid tumors<sup>6,98</sup>. The EMT

795 process is also enriched during the progression of OA; we have previously described this  
796 “chondrocyte-mesenchymal transition” (CMT) and note that it mimics the loss of tissue-specific  
797 gene expression and inflammatory increase in invasive capacity which is the hallmark of the  
798 metastatic process<sup>10</sup>. Invasion and aggression of breast tumors are known to be correlated with  
799 mechanical stress and stiffened phenotypes<sup>99,100</sup>, and pancreatic tumors have similarly been  
800 associated with such stress owing to their significant stromal fraction<sup>101,102</sup>. The shared nature of  
801 these phenotypes highlights the role of the ECM and stroma in tumor malignancy, and suggests  
802 that ColX may contribute in a similar pathological manner to both OA and cancer progression.

803 Numerous biological pathways and transcriptional networks implicated in pro-OA and  
804 pro-metastatic disease states are significantly enriched within or activated in conjunction with  
805 increased expression of the ColX modules characterized here. The hallmark pathway most  
806 significantly enriched in each ColX module was the epithelial-to-mesenchymal transition (EMT),  
807 which was also found to be differentially activated in high-ColX module expression samples  
808 compared to low-ColX module expression samples. EMT is intimately involved in the  
809 intercellular remodeling processes of embryonic development, fibrosis, and wound healing, and  
810 its activity in cancer is broadly considered crucial to the invasion of ECM and neighboring  
811 tissues which initiates metastasis in epithelial malignancies<sup>103</sup>. Upregulation of EMT was found  
812 to coincide with increased activity of angiogenic pathways which facilitate tumor cell  
813 intravasation and motility into the circulatory system, as well as key immunosuppressive cellular  
814 signatures (regulatory T cells and TAMs) which permit tumor development. Activity of both of  
815 these immune cell types have been associated with poor outcomes in breast and pancreatic  
816 cancer<sup>104–107</sup>. These results suggest that heightened ColX module activity signifies aggressive  
817 tumor states marked by impaired immune response and greater potential to metastasize.

818 ColX modules also exhibited significant overrepresentation and/or differential activity of  
819 transcription factor (TF) networks crucial to numerous developmental pathways known to play a  
820 role in cancer aggressiveness and metastasis (Notch, Wnt, and Hedgehog signaling), including

821 *FOXH1*, *SOX15*, and several *PRDM* gene family members<sup>108–113</sup>. GO enrichment of each ColX  
822 module revealed multiple genes implicated in ossification, such as *RUNX2* and *TGFB3*, both of  
823 which have been shown to play a role in OA, influence developmental pathways such as Wnt  
824 signaling, and contribute to EMT in breast and pancreatic cancer<sup>114–116</sup>. The male PAAD ColX  
825 module was additionally enriched for GO terms relating to EMT and Wnt signaling, both of which  
826 exert pro-metastatic effects that impair treatment responses<sup>117</sup>. Of particular interest were  
827 enriched targets of two oncogenic fusion proteins, *SS18-SSX* (misattributed in the GTRD  
828 database as *IGLV5-37* as noted above) and *SET-NUP214*, whose regulatory targets were  
829 enriched or differentially activated in association with both ColX module expression and  
830 signatures for osteogenesis and pathological OA cell types, suggesting common transcriptional  
831 dysregulation between these disease states.

832 Together, these results demonstrate that the genes in these cancer-specific ColX  
833 modules are involved in a diverse array of biological functions which may exacerbate the  
834 pathophysiology of both cancer and OA. Numerous genes implicated in both OA pathology and  
835 metastasis co-modularize with *COL10A1* in tumors, including *COMP*, *CXCL12*, and *FN1*, further  
836 suggesting common EMT-mediated dysregulation across these disease states. ColX module  
837 expression thus appears to highlight breast and pancreatic tumor states marked by activation of  
838 developmental and pro-OA pathways which may contribute to heightened metastatic potential  
839 and poorer clinical outcomes.

840 Survival analysis of the two TCGA cohorts reveals the prognostic value of ColX and its  
841 associated modules in both breast and pancreatic cancer. Increased expression of stromal ColX  
842 has previously been shown to correlate with worse survival outcomes in breast cancer cohorts  
843 with diverse tumor mutational burdens<sup>13,118</sup>, and *COL10A1* has also been linked to negative  
844 prognosis in pancreatic cancer<sup>93</sup>. Here, we corroborated that increased ColX expression confers  
845 significantly increased overall survival risk in both breast and pancreatic cancer. High ColX  
846 module expression is also significantly associated with decreased DFI in pancreatic cancer, an

847 effect which is stronger in the male PAAD cohort compared to the female PAAD cohort. We  
848 observed that DFI events actually occurred more frequently in the high-ColX module group in  
849 the female PAAD cohort (65% event rate) compared to the male PAAD cohort (27% event rate),  
850 but the groupwise Kaplan-Meier analysis for the female PAAD cohort was ultimately not  
851 significant, likely due to being underpowered (when the analysis was rerun with all female  
852 samples artificially duplicated, the difference between high- and low-ColX module expression  
853 groups became significant at an  $\alpha$ -level of 0.05). Analysis of a larger sample cohort would help  
854 clarify whether this is truly a gender-specific effect or merely the result of the cohort size  
855 analyzed here. Of note, a prior meta-analysis of multiple studies comprising 1,000 pancreatic  
856 cancer patients in total found that increased EMT is vital to the process of tumor budding, which  
857 confers significantly worse outcomes in terms of both overall survival and disease-free  
858 survival<sup>109</sup>. Thus, our results corroborate earlier findings regarding the predictive value of ColX  
859 expression and suggest the clinical utility of the ColX modules for risk stratification and  
860 prognostication.

861 Although *COL10A1* has been reported to be expressed by pancreatic cancer cells  
862 directly<sup>73</sup>, we observed that ColX is most strongly detected in the stromal region of PAAD  
863 tumors (Figure 1B), suggesting that it more likely serves as a marker for infiltration of fibroblasts  
864 and subsequent induction of an immunosuppressive microenvironment. Additionally, stromal  
865 markers such as *ACTA2* are co-modularized with *COL10A1* in all 4 ColX modules. These  
866 observations are consistent with *COL10A1*'s previously-characterized role as a marker of  
867 myofibroblastic TGF- $\beta$ -driven fibroblasts in pancreatic cancer, based on single-cell analysis of  
868 human PDAC samples<sup>78,119</sup>. Recent work by Thorlacius-Ussing et al. has further highlighted  
869 *COL10A1* (as well as *COL8A1*, *COL11A1*, and *COL12A1*, all of which are present in the ColX  
870 modules) as a marker of myCAFs in PDAC as well as in other cancer types<sup>93</sup>. In tumors, the  
871 ECM is primarily produced by fibroblasts and activated myofibroblasts, the latter of which are  
872 additionally responsible for driving fibrosis in response to inflammatory stimuli such as TGF- $\beta$

873 signaling from immune cells<sup>1</sup>. While CAFs comprise a heterogeneous collection of cells with a  
874 diverse array of functions, as a whole they exhibit numerous shared features across breast,  
875 pancreatic and other cancers, notably plasticity of their developmental pathways as well as  
876 purported regulatory roles in the functioning of NK, T, and other immune cells<sup>120</sup>. In breast and  
877 pancreatic tumors with high ColX module expression, we observed significant differential  
878 activation of numerous CAF types including matrix CAFs, inflammatory CAFs, antigen-  
879 presenting CAFs, vascular CAFs, and tumor-like CAFs, coupled with a decreased signal  
880 associated with dividing CAFs especially in breast cancer. This activation of CAF-specific gene  
881 signatures was corroborated by co-modularization of key CAF type markers with *COL10A1*  
882 across the breast and pancreatic ColX modules, which highlight CAF heterogeneity as well as  
883 common features of CAF-mediated aggression in both cancer types. CAFs have been shown to  
884 exert a variety of tumor-protective effects through direct interactions with cancer cells, ranging  
885 from promoting cancer stemness and preventing T cell recognition of cancer cells to facilitating  
886 invasion of the basement membrane through deposition of collagen “migratory tracks” and  
887 integrin-mediated interactions with fibronectin<sup>1</sup>. Although CAF populations are often  
888 characterized by their dominant role, considerable diversity of function has also been reported  
889 for several specific subtypes<sup>19</sup>. For example, although matrix CAFs are primarily responsible for  
890 producing and remodeling the ECM, they also appear to be capable of producing pro-  
891 inflammatory cytokines and chemokines to facilitate adhesion and migration. Similarly, tumor-  
892 like CAFs typically mimic tumor expression patterns and interact directly with tumor cells to  
893 promote stemness, chemoresistance and immunosuppression, but may also produce MMPs  
894 and other matrix proteins which contribute to remodeling. Our observed activation of antigen-  
895 presenting CAFs in pancreatic cancer is of particular interest, as prior studies in pancreatic  
896 cancer have shown that this CAF subtype interacts significantly with tumor-infiltrating T cells  
897 and TAMs through MHC II expression which induces naive CD4<sup>+</sup> T cells to differentiate into  
898 regulatory T cells, thereby contributing to immune evasion<sup>120,121</sup>. Regulatory T cells and TAMs



899 were among the most strongly activated immune cell types in high-ColX module pancreatic (as  
900 well as breast) tumors in our analysis, an effect which links high ColX module expression to  
901 both CAF activity and immunosuppression in aggressive tumors. Of note, the activation of  
902 regulatory T cells in the absence of activated CD8<sup>+</sup> T cell signatures is an effect which is almost  
903 uniquely associated with the ColX modules (as opposed to other WGCNA modules), suggesting  
904 predominance of a pathological immunosuppressive environment which lacks a strong cell-  
905 mediated immune component. Additionally, recent work has shown that increased CAF density  
906 is associated with an inflammatory, pro-EMT environment in PDAC<sup>122</sup>. *COL10A1* is thus a  
907 valuable biomarker for CAF-mediated ECM remodeling, immunosuppression, inflammation, and  
908 induction of invasion in breast and pancreatic tumors; future work will further elucidate its  
909 complex roles across diverse CAF and cancer types.

910         Activated CAFs have been shown to derive from various stromal origins including bone  
911 marrow<sup>1</sup>. The upregulation of the TGF- $\beta$ -induced *TGFBI* in BMSCs (Table S6) supports the  
912 hypothesis that ColX may originate from “activated,” dysregulated bone marrow-derived  
913 mesenchymal cells which infiltrate breast and pancreatic tumors. There is evidence for this  
914 effect in mouse models<sup>123</sup>, but effective markers are still being sought for human tumors. ColX is  
915 the strongest potential marker so far due to its specific expression in bone marrow during  
916 adulthood. We found that increased ColX module expression correlated significantly with the  
917 BMSC signature, as evidenced by CIBERSORTx analysis and co-modularization of *COL10A1*  
918 with BMSC-associated collagens type I, IV, VIII, XI, and XII. The roles that such bone-derived  
919 cells play within the tumor microenvironment is an area of active study. BMSCs are intricately  
920 involved in regulation of osteoblastic differentiation and the vascular microenvironment, and OA-  
921 MSCs, which are also ColX-positive, drive the chronic fibrosis and inflammation which  
922 characterizes the disease state<sup>10,124</sup>. In the context of solid tumors, high *COL10A1* expression  
923 may therefore indicate the presence of bone marrow-derived cells which contribute to an  
924 inflammatory, pro-EMT microenvironment. Coupled with the enrichment for gene ontologies

925 related to cartilage/skeletal development and ossification, this suggests that increased ColX  
926 module expression may signify development of ECM pathology within the tumor mesenchyme.  
927 A compelling piece of evidence for this common developmental feature is the fact that the Wnt  
928 family gene *WNT2* emerged as a ColX module gene in both BRCA and PAAD (Table S2A) as  
929 well as a specific marker of BMSCs (Table S6); the Wnt signaling cascade is heavily implicated  
930 in joint degeneration and OA pathogenesis<sup>125</sup>, and *WNT2* has been established as a highly-  
931 expressed gene in breast cancer as well as a pro-metastatic activator in pancreatic cancer<sup>126,127</sup>.  
932 The interplay between various CAF subtypes and bone marrow cells is further supported by the  
933 co-modularization of *COL10A1* with key markers such as *CXCL12*, which is a marker for both  
934 inflammatory CAFs as well as hypertrophic chondrocytes which give rise to osteoblasts and  
935 recruit vasculature during skeletal development<sup>19,128</sup>. Corroborated by the capacity of NCSCs to  
936 transform into diverse pathological OA cell types through a “senescence-associated cell  
937 transition and interaction” (SACTAI)<sup>124</sup> which mimics EMT-associated changes occurring in  
938 breast and pancreatic cancer, these results further support the idea that advanced tumors may  
939 contain significant subpopulations of cells with OA character. Additionally, the mutual  
940 enrichment of specific genes implicated in EMT across breast and pancreatic cancer and  
941 pathological OA cell types suggests that similar dysregulatory mechanisms are crucial to  
942 development of both disease states. Spatial transcriptomic studies would provide greater insight  
943 to the specific tumoral regions exhibiting *COL10A1* expression and help to further validate this  
944 hypothesis. Thus, differential activity of the ColX module reveals a more complete picture of the  
945 contribution of *COL10A1* to malignant CAF and pro-metastatic activity in advanced tumors.

946 Although our findings suggest a multifaceted role for *COL10A1* in the maintenance and  
947 remodeling of the ECM with involvement of CAFs, immune cells, and other key players in the  
948 tumor microenvironment, they are necessarily limited by the descriptive nature of this study. In  
949 particular, while we have highlighted numerous pathological mechanisms which appear to  
950 correlate and trend directionally with ColX module expression, direct experimentation will be

951 necessary to robustly validate these relationships in an in vitro or in vivo setting. Survival  
952 analyses (in particular for the gender-specific PAAD cohorts) were restricted by the number of  
953 TCGA samples with complete data for each outcome of interest, and the ColX module-based  
954 survival impacts presented here would benefit from validation in a larger dataset. Nevertheless,  
955 the inflammatory, immunosuppressive, and pro-metastatic pathways we have identified here  
956 supplement the current knowledge regarding *COL10A1* and its significance to solid tumors, and  
957 suggest multiple avenues of exploration to further characterize its importance in breast,  
958 pancreatic, and other cancers.

959

## 960 **CONCLUSIONS**

961 In this study, we have demonstrated numerous links of biological and clinical interest  
962 between the stromal ECM in breast and pancreatic cancer as well as bone marrow and OA  
963 cartilage, highlighted by shared expression of *COL10A1* and its associated gene networks  
964 which contribute to development of an inflammatory, immunosuppressive, and CAF-dominated  
965 microenvironment to facilitate EMT and metastasis. *COL10A1* is an important and specifically  
966 expressed collagen whose role in the progression of solid tumors is an area of active study; our  
967 results suggest that it holds substantial value as a regulator and biomarker of aggressive tumor  
968 phenotypes with implications for ECM-targeted therapies and clinical outcomes. Identification of  
969 tumors which exhibit high expression of *COL10A1* and its associated genes may reveal the  
970 presence of more aggressive pathological microenvironments with heightened EMT and  
971 metastatic potential. These findings may enable more effective risk assessment and treatment  
972 of patients with breast and pancreatic cancer.

973

## 974 **LIST OF ABBREVIATIONS**

- 975 • BH: Benjamini-Hochberg  
976 • BMSCs: bone marrow stromal cells

- 977 ● BRCA: breast invasive carcinoma
- 978 ● CAFs: cancer-associated fibroblasts
- 979 ● ColX: collagen type X (*COL10A1*)
- 980 ● CMT: chondrocyte-mesenchymal transition
- 981 ● DFI: disease-free interval
- 982 ● ECM: extracellular matrix
- 983 ● EMT: epithelial-to-mesenchymal transition
- 984 ● %G.A.M.E.: percentage of Genes Above-Median Expression
- 985 ● GEO: Gene Expression Omnibus
- 986 ● GO: gene ontology
- 987 ● GSEA: Gene Set Enrichment Analysis
- 988 ● GTRD: Gene Transcription Regulation Database
- 989 ● IHC: immunohistochemistry
- 990 ● ME: module eigengene
- 991 ● NCSCs: normal cartilage stromal cells
- 992 ● OA: osteoarthritis, osteoarthritic
- 993 ● OA-MSCs: OA mesenchymal stromal cells
- 994 ● OACs: OA chondrocytes
- 995 ● OS: overall survival
- 996 ● PAAD: pancreatic adenocarcinoma
- 997 ● PACA-AU: Pancreatic Cancer Australian
- 998 ● PDAC: pancreatic ductal adenocarcinoma
- 999 ● QuSAGE: Quantitative Set Analysis for Gene Expression
- 1000 ● SACTAI: senescence-associated cell transition and interaction
- 1001 ● TAMs: tumor-associated macrophages
- 1002 ● TCGA: The Cancer Genome Atlas

- 1003 • TF: transcription factor
- 1004 • TFTs: transcription factor targets
- 1005 • WGCNA: weighted gene co-expression network analysis

1006

## 1007 **DECLARATIONS**

### 1008 **Ethics approval and consent to participate**

1009           Use of patient material was approved by the Lifespan institutional review board approval  
1010 (IRB #1070389–9). All procedures were performed in accordance with the relevant guidelines  
1011 and regulations.

1012

### 1013 **Consent for publication**

1014           Patient consent for sample publication was obtained appropriately.

1015

### 1016 **Availability of data and materials**

1017           The datasets supporting the conclusions of this article are available in the NCI Genomic  
1018 Data Commons (<https://gdc.cancer.gov/about-data/publications/pancanatlas>), in the Broad  
1019 Institute GDAC FireBrowse portal (<http://firebrowse.org/>), and in the GEO accession  
1020 GSE176199 (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE176199>).

1021

### 1022 **Competing interests**

1023           ASB and MBR are co-inventors of the following patent: Brodsky, Alexander S. and  
1024 Wang, Yihong and Resnick, Murray. 2017. Collagens as markers for breast cancer treatment.  
1025 US Patent US09784743B2, filed Jun 20, 2016, and issued Oct 10, 2017. The other authors  
1026 declare that they have no competing interests.

1027

1028 **Funding**

1029 This study was supported by R01AG080141 and P30GM122732.

1030

1031 **Authors' contributions**

1032 EHHFY and AFY contributed equally. ASB and QC conceived and designed the study.

1033 WL generated the cartilage RNA-Seq data. DY performed IHC. EYW and AS interpreted IHC.

1034 MBR contributed to interpretation of all *COL10A1* and tumor features. EHHFY and AFY

1035 acquired, analyzed and interpreted the data, and drafted the manuscript. ASB and QC edited

1036 the manuscript with input from all authors.

1037

1038 **Acknowledgements**

1039 The results described here are in part based upon data generated by the TCGA

1040 Research Network (<https://www.cancer.gov/tcga>). We thank the patients and their families for

1041 their participation in the individual TCGA projects.

1042

1043 **REFERENCES**

1044 1. Winkler, J., Abisoye-Ogunniyan, A., Metcalf, K. J. & Werb, Z. Concepts of extracellular  
1045 matrix remodelling in tumour progression and metastasis. *Nat. Commun.* **11**, 5120 (2020).

1046 2. Naba, A. *et al.* The extracellular matrix: Tools and insights for the “omics” era. *Matrix Biol.*  
1047 **49**, 10–24 (2016).

1048 3. Cox, T. R. & Eler, J. T. Remodeling and homeostasis of the extracellular matrix:  
1049 implications for fibrotic diseases and cancer. *Dis. Model. Mech.* **4**, 165–178 (2011).

1050 4. Guo, K. S. & Brodsky, A. S. Tumor collagens predict genetic features and patient  
1051 outcomes. *NPJ Genomic Med.* **8**, 15 (2023).

1052 5. Kendall, R. T. & Feghali-Bostwick, C. A. Fibroblasts in fibrosis: novel roles and mediators.  
1053 *Front. Pharmacol.* **5**, 123 (2014).

- 1054 6. Luo, Y. *et al.* The minor collagens in articular cartilage. *Protein Cell* **8**, 560–572 (2017).
- 1055 7. Von Der Mark, K. *et al.* Type x collagen synthesis in human osteoarthritic cartilage.  
1056 indication of chondrocyte hypertrophy. *Arthritis Rheum.* **35**, 806–811 (1992).
- 1057 8. Jayasuriya, C. T. *et al.* Molecular characterization of mesenchymal stem cells in human  
1058 osteoarthritis cartilage reveals contribution to the OA phenotype. *Sci. Rep.* **8**, 7044 (2018).
- 1059 9. Hoyland, J. A. *et al.* Distribution of type X collagen mRNA in normal and osteoarthritic  
1060 human cartilage. *Bone Miner.* **15**, 151–163 (1991).
- 1061 10. Liu, W. *et al.* Senescent Tissue-Resident Mesenchymal Stromal Cells Are an Internal  
1062 Source of Inflammation in Human Osteoarthritic Cartilage. *Front. Cell Dev. Biol.* **9**, 725071  
1063 (2021).
- 1064 11. He, Y. *et al.* Type X collagen levels are elevated in serum from human osteoarthritis  
1065 patients and associated with biomarkers of cartilage degradation and inflammation. *BMC*  
1066 *Musculoskelet. Disord.* **15**, 309 (2014).
- 1067 12. Guilak, F., Nims, R., Dicks, A., Wu, C.-L. & Meulenbelt, I. Osteoarthritis as a disease of the  
1068 cartilage pericellular matrix. *Matrix Biol.* **71–72**, 40–50 (2018).
- 1069 13. Brodsky, A. S. *et al.* Identification of stromal ColXα1 and tumor-infiltrating lymphocytes as  
1070 putative predictive markers of neoadjuvant therapy in estrogen receptor-positive/HER2-  
1071 positive breast cancer. *BMC Cancer* **16**, 274 (2016).
- 1072 14. Wang, Y. *et al.* ColXα1 is a stromal component that colocalizes with elastin in the breast  
1073 tumor extracellular matrix. *J. Pathol. Clin. Res.* **5**, 40–52 (2018).
- 1074 15. Zhao, C. L. *et al.* Stromal ColXα1 expression correlates with tumor-infiltrating lymphocytes  
1075 and predicts adjuvant therapy outcome in ER-positive/HER2-positive breast cancer. *BMC*  
1076 *Cancer* **19**, 1036 (2019).
- 1077 16. Hingorani, S. R. Epithelial and stromal co-evolution and complicity in pancreatic cancer.  
1078 *Nat. Rev. Cancer* **23**, 57–77 (2023).
- 1079 17. Olivares, O. *et al.* Collagen-derived proline promotes pancreatic ductal adenocarcinoma

- 1080 cell survival under nutrient limited conditions. *Nat. Commun.* **8**, 16031 (2017).
- 1081 18. Kawai, H. *et al.* Characterization and potential roles of bone marrow-derived stromal cells  
1082 in cancer development and metastasis. *Int. J. Med. Sci.* **15**, 1406–1414 (2018).
- 1083 19. Cords, L. *et al.* Cancer-associated fibroblast classification in single-cell and spatial  
1084 proteomics data. *Nat. Commun.* **14**, 4294 (2023).
- 1085 20. Raz, Y. *et al.* Bone marrow-derived fibroblasts are a functionally distinct stromal cell  
1086 population in breast cancer. *J. Exp. Med.* **215**, 3075–3093 (2018).
- 1087 21. Sweeney, E., Roberts, D., Corbo, T. & Jacenko, O. Congenic mice confirm that collagen X  
1088 is required for proper hematopoietic development. *PLoS One* **5**, e9518 (2010).
- 1089 22. R Core Team. R: A Language and Environment for Statistical Computing. R Foundation  
1090 for Statistical Computing (2020).
- 1091 23. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York  
1092 (2016).
- 1093 24. Warnes, G. R. *et al.* *gplots: Various R Programming Tools for Plotting Data*. (2022).
- 1094 25. Larsson, J. & Gustafsson, P. A Case Study in Fitting Area-Proportional Euler Diagrams  
1095 with Ellipses using *eulerr*. *Proc. Int. Workshop Set Vis. Reason.* **2116**, 84–91 (2018).
- 1096 26. Larsson, J. *eulerr: Area-Proportional Euler and Venn Diagrams with Ellipses*. (2022).
- 1097 27. Grossman, R. L. *et al.* Toward a Shared Vision for Cancer Genomic Data. *N. Engl. J. Med.*  
1098 **375**, 1109–1112 (2016).
- 1099 28. Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.*  
1100 **45**, 1113–1120 (2013).
- 1101 29. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene  
1102 expression and hybridization array data repository. *Nucleic Acids Res.* **30**, 207–210  
1103 (2002).
- 1104 30. Jiang, J. *et al.* Tumour-Infiltrating Immune Cell-Based Subtyping and Signature Gene  
1105 Analysis in Breast Cancer Based on Gene Expression Profiles. *J. Cancer* **11**, 1568–1583



- 1106 (2020).
- 1107 31. Loi, S. *et al.* Definition of clinically distinct molecular subtypes in estrogen receptor-positive  
1108 breast carcinomas through genomic grade. *J. Clin. Oncol. Off. J. Am. Soc. Clin. Oncol.* **25**,  
1109 1239–1246 (2007).
- 1110 32. Loi, S. *et al.* Predicting prognosis using molecular profiling in estrogen receptor-positive  
1111 breast cancer treated with tamoxifen. *BMC Genomics* **9**, 239 (2008).
- 1112 33. Desmedt, C. *et al.* Multifactorial approach to predicting resistance to anthracyclines. *J.*  
1113 *Clin. Oncol. Off. J. Am. Soc. Clin. Oncol.* **29**, 1578–1586 (2011).
- 1114 34. Sircoulomb, F. *et al.* Genome profiling of ERBB2-amplified breast cancers. *BMC Cancer*  
1115 **10**, 539 (2010).
- 1116 35. Li, Y. *et al.* Amplification of LAPT4B and YWHAZ contributes to chemotherapy  
1117 resistance and recurrence of breast cancer. *Nat. Med.* **16**, 214–218 (2010).
- 1118 36. Kao, K.-J., Chang, K.-M., Hsu, H.-C. & Huang, A. T. Correlation of microarray-based  
1119 breast cancer molecular subtypes and clinical outcomes: implications for treatment  
1120 optimization. *BMC Cancer* **11**, 143 (2011).
- 1121 37. Dedeurwaerder, S. *et al.* DNA methylation profiling reveals a predominant immune  
1122 component in breast cancers. *EMBO Mol. Med.* **3**, 726–741 (2011).
- 1123 38. Sabatier, R. *et al.* A gene expression signature identifies two prognostic subgroups of  
1124 basal breast cancer. *Breast Cancer Res. Treat.* **126**, 407–420 (2011).
- 1125 39. Sabatier, R. *et al.* Down-regulation of ECRG4, a candidate tumor suppressor gene, in  
1126 human breast cancer. *PLoS One* **6**, e27656 (2011).
- 1127 40. Clarke, C. *et al.* Correlating transcriptional networks to breast cancer survival: a large-  
1128 scale coexpression analysis. *Carcinogenesis* **34**, 2300–2308 (2013).
- 1129 41. Huang, C.-C. *et al.* Concurrent gene signatures for han chinese breast cancers. *PLoS One*  
1130 **8**, e76421 (2013).
- 1131 42. Sonnenblick, A. *et al.* Integrative proteomic and gene expression analysis identify potential

- 1132 biomarkers for adjuvant trastuzumab resistance: analysis from the Fin-her phase III  
1133 randomized trial. *Oncotarget* **6**, 30306–30316 (2015).
- 1134 43. Zhang, J. *et al.* The International Cancer Genome Consortium Data Portal. *Nat.*  
1135 *Biotechnol.* **37**, 367–369 (2019).
- 1136 44. Carvalho, B. S. & Irizarry, R. A. A framework for oligonucleotide microarray preprocessing.  
1137 *Bioinformatics* **26**, 2363–2367 (2010).
- 1138 45. Carlson, M. hgu133plus2.db: Affymetrix Human Genome U133 Plus 2.0 Array annotation  
1139 data (chip hgu133plus2). (2016).
- 1140 46. Pagès, H., Carlson, M., Falcon, S. & Li, N. AnnotationDbi: Manipulation of SQLite-based  
1141 annotations in Bioconductor. (2020).
- 1142 47. Wickham, H. *et al.* Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).
- 1143 48. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network  
1144 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 1145 49. Langfelder, P. & Horvath, S. Fast R Functions for Robust Correlations and Hierarchical  
1146 Clustering. *J. Stat. Softw.* **46**, 1–17 (2012).
- 1147 50. Leek, J. *et al.* sva: Surrogate Variable Analysis. (2020).
- 1148 51. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion  
1149 for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 1150 52. Uhlén, M. *et al.* Proteomics. Tissue-based map of the human proteome. *Science* **347**,  
1151 1260419 (2015).
- 1152 53. Kamburov, A., Stelzl, U., Lehrach, H. & Herwig, R. The ConsensusPathDB interaction  
1153 database: 2013 update. *Nucleic Acids Res.* **41**, D793–D800 (2013).
- 1154 54. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set  
1155 collection. *Cell Syst.* **1**, 417–425 (2015).
- 1156 55. Kolmykov, S. *et al.* GTRD: an integrated view of transcription regulation. *Nucleic Acids*  
1157 *Res.* **49**, D104–D111 (2021).

- 1158 56. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and  
1159 Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300  
1160 (1995).
- 1161 57. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased  
1162 coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic  
1163 Acids Res.* **47**, D607–D613 (2019).
- 1164 58. Langfelder, P., Luo, R., Oldham, M. C. & Horvath, S. Is My Network Module Preserved  
1165 and Reproducible? *PLOS Comput. Biol.* **7**, e1001057 (2011).
- 1166 59. Charoentong, P. *et al.* Pan-cancer Immunogenomic Analyses Reveal Genotype-  
1167 Immunophenotype Relationships and Predictors of Response to Checkpoint Blockade.  
1168 *Cell Rep.* **18**, 248–262 (2017).
- 1169 60. Rabosky, D. L. *et al.* BAMMtools: an R package for the analysis of evolutionary dynamics  
1170 on phylogenetic trees. *Methods Ecol. Evol.* **5**, 701–707 (2014).
- 1171 61. Yaari, G., Bolen, C. R., Thakar, J. & Kleinstein, S. H. Quantitative set analysis for gene  
1172 expression: a method to quantify gene set differential expression including gene-gene  
1173 correlations. *Nucleic Acids Res.* **41**, e170 (2013).
- 1174 62. Turner, J. A., Bolen, C. R. & Blankenship, D. M. Quantitative gene set analysis  
1175 generalized for repeated measures, confounder adjustment, and continuous covariates.  
1176 *BMC Bioinformatics* **16**, 272 (2015).
- 1177 63. Meng, H., Yaari, G., Bolen, C. R., Avey, S. & Kleinstein, S. H. Gene set meta-analysis with  
1178 Quantitative Set Analysis for Gene Expression (QuSAGE). *PLOS Comput. Biol.* **15**,  
1179 e1006899 (2019).
- 1180 64. Therneau, T. M. & Grambsch, P. M. *Modeling Survival Data: Extending the Cox Model.*  
1181 (Springer, New York, NY, 2000). doi:10.1007/978-1-4757-3294-8.
- 1182 65. Therneau, T. A Package for Survival Analysis in R. (2022).
- 1183 66. Kassambara, A., Kosinski, M. & Biecek, P. survminer: Drawing Survival Curves using

- 1184 'ggplot2'. (2021).
- 1185 67. Newman, A. M. *et al.* Determining cell type abundance and expression from bulk tissues  
1186 with digital cytometry. *Nat. Biotechnol.* **37**, 773–782 (2019).
- 1187 68. Naba, A., Clauser, K. R., Lamar, J. M., Carr, S. A. & Hynes, R. O. Extracellular matrix  
1188 signatures of human mammary carcinoma identify novel metastasis promoters. *eLife* **3**,  
1189 e01308 (2014).
- 1190 69. Tian, C. *et al.* Cancer-cell-derived matrisome proteins promote metastasis in pancreatic  
1191 ductal adenocarcinoma. *Cancer Res.* **80**, 1461–1474 (2020).
- 1192 70. Tian, C. *et al.* Proteomic analyses of ECM during pancreatic ductal adenocarcinoma  
1193 progression reveal different contributions by tumor and stromal cells. *Proc. Natl. Acad. Sci.*  
1194 *U. S. A.* **116**, 19609–19618 (2019).
- 1195 71. Naba, A. Ten Years of Extracellular Matrix Proteomics: Accomplishments, Challenges,  
1196 and Future Perspectives. *Mol. Cell. Proteomics* **22**, 100528 (2023).
- 1197 72. Zhang, K. *et al.* The collagen receptor discoidin domain receptor 2 stabilizes SNAIL1 to  
1198 facilitate breast cancer metastasis. *Nat. Cell Biol.* **15**, 677–687 (2013).
- 1199 73. Wen, Z. *et al.* COL10A1-DDR2 axis promotes the progression of pancreatic cancer by  
1200 regulating MEK/ERK signal transduction. *Front. Oncol.* **12**, (2022).
- 1201 74. Liu, Y., Zhang, L., Meng, Y. & Huang, L. Benzyl isothiocyanate inhibits breast cancer cell  
1202 tumorigenesis via repression of the FoxH1-Mediated Wnt/ $\beta$ -catenin pathway. *Int. J. Clin.*  
1203 *Exp. Med.* **8**, 17601–17611 (2015).
- 1204 75. Katoh, M. & Katoh, M. Integrative genomic analyses of CXCR4: Transcriptional regulation  
1205 of CXCR4 based on TGF $\beta$ , Nodal, Activin signaling and POU5F1, FOXA2, FOXC2,  
1206 FOXH1, SOX17, and GFI1 transcription factors. *Int. J. Oncol.* **36**, 415–420 (2010).
- 1207 76. McBride, M. J. *et al.* The SS18-SSX Fusion Oncoprotein Hijacks BAF Complex Targeting  
1208 and Function to Drive Synovial Sarcoma. *Cancer Cell* **33**, 1128-1141.e7 (2018).
- 1209 77. Kim, W. *et al.* RUNX1 is essential for mesenchymal stem cell proliferation and

- 1210 myofibroblast differentiation. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 16389–16394 (2014).
- 1211 78. Dominguez, C. X. *et al.* Single-Cell RNA Sequencing Reveals Stromal Evolution into  
1212 LRRC15+ Myofibroblasts as a Determinant of Patient Response to Cancer  
1213 Immunotherapy. *Cancer Discov.* **10**, 232–253 (2020).
- 1214 79. Micalizzi, D. S. *et al.* The Six1 homeoprotein induces human mammary carcinoma cells to  
1215 undergo epithelial-mesenchymal transition and metastasis in mice through increasing  
1216 TGF- $\beta$  signaling. *J. Clin. Invest.* **119**, 2678–2690 (2009).
- 1217 80. Liu, W. *et al.* Homeoprotein SIX1 compromises antitumor immunity through TGF- $\beta$ -  
1218 mediated regulation of collagens. *Cell. Mol. Immunol.* **18**, 2660–2672 (2021).
- 1219 81. Couto, J. P. *et al.* Nicotinamide N-methyltransferase sustains a core epigenetic program  
1220 that promotes metastatic colonization in breast cancer. *EMBO J.* **42**, e112559 (2023).
- 1221 82. Littlewood-Evans, A. J. *et al.* The osteoclast-associated protease cathepsin K is expressed  
1222 in human breast carcinoma. *Cancer Res.* **57**, 5386–5390 (1997).
- 1223 83. Wu, N. *et al.* Cathepsin K regulates the tumor growth and metastasis by IL-17/CTSK/EMT  
1224 axis and mediates M2 macrophage polarization in castration-resistant prostate cancer.  
1225 *Cell Death Dis.* **13**, 1–9 (2022).
- 1226 84. Terekhanova, N. V. *et al.* Epigenetic regulation during cancer transitions across 11 tumour  
1227 types. *Nature* **623**, 432–441 (2023).
- 1228 85. Zhao, J. *et al.* RWCFusion: identifying phenotype-specific cancer driver gene fusions  
1229 based on fusion pair random walk scoring method. *Oncotarget* **7**, 61054–61068 (2016).
- 1230 86. Oka, M. *et al.* Chromatin-bound CRM1 recruits SET-Nup214 and NPM1c onto HOX  
1231 clusters causing aberrant HOX expression in leukemia cells. *eLife* **8**, e46667 (2019).
- 1232 87. Chen, J. *et al.* DLX5 promotes Col10a1 expression and chondrocyte hypertrophy and is  
1233 involved in osteoarthritis progression. *Genes Dis.* **10**, 2097–2108 (2023).
- 1234 88. Czerwińska, P., Mazurek, S. & Wiznerowicz, M. The complexity of TRIM28 contribution to  
1235 cancer. *J. Biomed. Sci.* **24**, 63 (2017).

- 1236 89. Eyre, D. R. The collagens of articular cartilage. *Semin. Arthritis Rheum.* **21**, 2–11 (1991).
- 1237 90. Luckman, S. P., Rees, E. & Kwan, A. P. L. Partial characterization of cell-type X collagen  
1238 interactions. *Biochem. J.* **372**, 485–493 (2003).
- 1239 91. Shen, G. The role of type X collagen in facilitating and regulating endochondral ossification  
1240 of articular cartilage. *Orthod. Craniofac. Res.* **8**, 11–17 (2005).
- 1241 92. Fernandes, A. M. *et al.* Similar Properties of Chondrocytes from Osteoarthritis Joints and  
1242 Mesenchymal Stem Cells from Healthy Donors for Tissue Engineering of Articular  
1243 Cartilage. *PLoS One* **8**, e62994 (2013).
- 1244 93. Thorlacius-Ussing, J. *et al.* The collagen landscape in cancer: profiling collagens in tumors  
1245 and in circulation reveals novel markers of cancer-associated fibroblast subtypes. *J.*  
1246 *Pathol.* **n/a**, (2023).
- 1247 94. Alcaide-Ruggiero, L., Molina-Hernández, V., Granados, M. M. & Domínguez, J. M. Main  
1248 and Minor Types of Collagens in the Articular Cartilage: The Role of Collagens in Repair  
1249 Tissue Evaluation in Chondral Defects. *Int. J. Mol. Sci.* **22**, 13329 (2021).
- 1250 95. Todhunter, R. J. *et al.* Gene expression in hip soft tissues in incipient canine hip dysplasia  
1251 and osteoarthritis. *J. Orthop. Res.* **37**, 313–324 (2019).
- 1252 96. Zeisberg, M. & Neilson, E. G. Biomarkers for epithelial-mesenchymal transitions. *J. Clin.*  
1253 *Invest.* **119**, 1429–1437 (2009).
- 1254 97. Tani, S. *et al.* Stem cell-based modeling and single-cell multiomics reveal gene-regulatory  
1255 mechanisms underlying human skeletal development. *Cell Rep.* **42**, (2023).
- 1256 98. Paszek, M. J. *et al.* Tensional homeostasis and the malignant phenotype. *Cancer Cell* **8**,  
1257 241–254 (2005).
- 1258 99. Acerbi, I. *et al.* Human Breast Cancer Invasion and Aggression Correlates with ECM  
1259 Stiffening and Immune Cell Infiltration. *Integr. Biol. Quant. Biosci. Nano Macro* **7**, 1120–  
1260 1134 (2015).
- 1261 100. Pratt, S. J. P., Lee, R. M. & Martin, S. S. The Mechanical Microenvironment in Breast

- 1262 Cancer. *Cancers* **12**, 1452 (2020).
- 1263 101. Hosein, A. N., Brekken, R. A. & Maitra, A. Pancreatic cancer stroma: an update on  
1264 therapeutic targeting strategies. *Nat. Rev. Gastroenterol. Hepatol.* **17**, 487–505 (2020).
- 1265 102. Maneshi, P., Mason, J., Dongre, M. & Öhlund, D. Targeting Tumor-Stromal Interactions in  
1266 Pancreatic Cancer: Impact of Collagens and Mechanical Traits. *Front. Cell Dev. Biol.* **9**,  
1267 787485 (2021).
- 1268 103. Das, S. & Batra, S. K. Pancreatic Cancer Metastasis: Are we being Pre-EMT'ed? *Curr.*  
1269 *Pharm. Des.* **21**, 1249–1255 (2015).
- 1270 104. Plitas, G. & Rudensky, A. Y. Regulatory T Cells in Cancer. *Annu. Rev. Cancer Biol.* **4**,  
1271 459–477 (2020).
- 1272 105. Mota Reyes, C. *et al.* Regulatory T Cells in Pancreatic Cancer: Of Mice and Men. *Cancers*  
1273 **14**, 4582 (2022).
- 1274 106. Qiu, S.-Q. *et al.* Tumor-associated macrophages in breast cancer: Innocent bystander or  
1275 important player? *Cancer Treat. Rev.* **70**, 178–189 (2018).
- 1276 107. Yang, S., Liu, Q. & Liao, Q. Tumor-Associated Macrophages in Pancreatic Ductal  
1277 Adenocarcinoma: Origin, Polarization, Function, and Reprogramming. *Front. Cell Dev.*  
1278 *Biol.* **8**, 607209 (2021).
- 1279 108. Zhang, Z. & Xu, Y. FZD7 accelerates hepatic metastases in pancreatic cancer by  
1280 strengthening EMT and stemness associated with TGF- $\beta$ /SMAD3 signaling. *Mol. Med.*  
1281 *Camb. Mass* **28**, 82 (2022).
- 1282 109. Palamaris, K., Felekouras, E. & Sakellariou, S. Epithelial to Mesenchymal Transition: Key  
1283 Regulator of Pancreatic Ductal Adenocarcinoma Progression and Chemoresistance.  
1284 *Cancers* **13**, 5532 (2021).
- 1285 110. Gaijanigo, N., Melisi, D. & Carbone, C. EMT and Treatment Resistance in Pancreatic  
1286 Cancer. *Cancers* **9**, 122 (2017).
- 1287 111. Riobo-Del Galdo, N. A., Lara Montero, Á. & Wertheimer, E. V. Role of Hedgehog Signaling

- 1288 in Breast Cancer: Pathogenesis and Therapeutics. *Cells* **8**, 375 (2019).
- 1289 112. Kumar, V. *et al.* The Role of Notch, Hedgehog, and Wnt Signaling Pathways in the  
1290 Resistance of Tumors to Anticancer Therapies. *Front. Cell Dev. Biol.* **9**, (2021).
- 1291 113. Moradi, A. *et al.* The cross-regulation between SOX15 and Wnt signaling pathway. *J. Cell.*  
1292 *Physiol.* **232**, 3221–3225 (2017).
- 1293 114. Chuang, L. S. H. & Ito, Y. The Multiple Interactions of RUNX with the Hippo–YAP  
1294 Pathway. *Cells* **10**, 2925 (2021).
- 1295 115. Wysokinski, D., Blasiak, J. & Pawlowska, E. Role of RUNX2 in Breast Carcinogenesis. *Int.*  
1296 *J. Mol. Sci.* **16**, 20969–20993 (2015).
- 1297 116. Kayed, H. *et al.* Regulation and functional role of the Runt-related transcription factor-2 in  
1298 pancreatic cancer. *Br. J. Cancer* **97**, 1106–1115 (2007).
- 1299 117. Ram Makena, M. *et al.* Wnt/ $\beta$ -Catenin Signaling: The Culprit in Pancreatic Carcinogenesis  
1300 and Therapeutic Resistance. *Int. J. Mol. Sci.* **20**, 4242 (2019).
- 1301 118. Zhang, M., Chen, H., Wang, M., Bai, F. & Wu, K. Bioinformatics analysis of prognostic  
1302 significance of COL10A1 in breast cancer. *Biosci. Rep.* **40**, BSR20193286 (2020).
- 1303 119. Peng, J. *et al.* Single-cell RNA-seq highlights intra-tumoral heterogeneity and malignant  
1304 progression in pancreatic ductal adenocarcinoma. *Cell Res.* **29**, 725–738 (2019).
- 1305 120. Luo, H. *et al.* Pan-cancer single-cell analysis reveals the heterogeneity and plasticity of  
1306 cancer-associated fibroblasts in the tumor microenvironment. *Nat. Commun.* **13**, 6619  
1307 (2022).
- 1308 121. Huang, H. *et al.* Mesothelial cell-derived antigen-presenting cancer-associated fibroblasts  
1309 induce expansion of regulatory T cells in pancreatic cancer. *Cancer Cell* **40**, 656-673.e7  
1310 (2022).
- 1311 122. Guinn, S. *et al.* Transfer Learning Reveals Cancer-Associated Fibroblasts Are Associated  
1312 with Epithelial–Mesenchymal Transition and Inflammation in Cancer Cells in Pancreatic  
1313 Ductal Adenocarcinoma. *Cancer Res.* OF1–OF17 (2024) doi:10.1158/0008-5472.CAN-23-



- 1314 1660.
- 1315 123. Quante, M. *et al.* Bone marrow-derived myofibroblasts contribute to the mesenchymal  
1316 stem cell niche and promote tumor growth. *Cancer Cell* **19**, 257–272 (2011).
- 1317 124. Liu, Y., Schwam, J. & Chen, Q. Senescence-Associated Cell Transition and Interaction  
1318 (SACTAI): A Proposed Mechanism for Tissue Aging, Repair, and Degeneration. *Cells* **11**,  
1319 1089 (2022).
- 1320 125. Sayegh, E. T. *et al.* Inhibition of Wnt pathway activity as a treatment approach for human  
1321 osteoarthritis: a systematic review. *J. Cartil. Jt. Preserv.* **2**, 100069 (2022).
- 1322 126. Xu, X., Zhang, M., Xu, F. & Jiang, S. Wnt signaling in breast cancer: biological  
1323 mechanisms, challenges and opportunities. *Mol. Cancer* **19**, 165 (2020).
- 1324 127. Jiang, H. *et al.* Activation of the Wnt pathway through Wnt2 promotes metastasis in  
1325 pancreatic cancer. *Am. J. Cancer Res.* **4**, 537–544 (2014).
- 1326 128. Long, J. T. *et al.* Hypertrophic chondrocytes serve as a reservoir for marrow-associated  
1327 skeletal stem and progenitor cells, osteoblasts, and adipocytes during skeletal  
1328 development. *eLife* **11**, e76932 (2022).

1329

## 1330 **FIGURES**

1331 **Figure 1: COL10A1 is highly expressed in breast and pancreatic tumors.**

1332 **(A)** Expression of COL10A1 across all TCGA sample cohorts. **(B)** Representative 400x ColX $\alpha$ 1  
1333 immunohistochemistry staining in pancreatic tumors. **(C)** Outline of study samples and  
1334 analyses. See Methods and Results sections for details.

1335

1336 **Figure 2: ColX modules are enriched for ECM genes, pro-metastatic pathways, and**  
1337 **developmental/regulatory transcription factor targets.**

1338 **(A and B)** Overlap of genes within ColX WGCNA modules from TCGA **(A)** breast and  
1339 pancreatic cancer and **(B)** gender-segregated pancreatic cancer datasets. See **Table S2A** for

1340 lists of ColX module genes, **Table S2B** for full list of WGCNA modules for each dataset, and  
1341 **Table S3** for gene ontology and Reactome pathway enrichment analysis of cancer-specific and  
1342 overlapping ColX-associated genes. **(C and D)** Enrichment of **(C)** human *in silico* matrixome  
1343 gene sets<sup>2</sup> and **(D)** MSigDB hallmark pathway gene sets<sup>54</sup> within ColX modules. Gene sets  
1344 within each block are ordered by mean significance rank (by Fisher's exact test) across all 4  
1345 modules. See **Figure S2A–D** and **Table S4** for hallmark pathway enrichment analysis of all  
1346 WGCNA-inferred modules for each dataset. Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01;  
1347 \*\*\*, p.adj < 0.001; \*\*\*\*, p.adj < 0.0001. **(E)** Enrichment of Gene Transcription Regulation  
1348 Database (GTRD) transcription factor (TF) targets<sup>55</sup> within breast and pancreatic cancer ColX  
1349 modules. Of note, the TFT gene set attributed to *IGLV5-37* in the GTRD database actually  
1350 represents targets of the fusion oncoprotein *SS18-SSX*, as described in the text. Dotted lines  
1351 correspond to p.adj = 0.10. Green labels/points indicate TFs whose targets are enriched in both  
1352 breast and pancreatic cancer ColX modules. See **Table S5A** for reported TF functions and lists  
1353 of overlapping TFTs. Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*,  
1354 p.adj < 0.0001.

1355

1356 **Figure 3: Tumors with high ColX module expression exhibit activation of pro-metastatic,**  
1357 **immunosuppressive, and myofibroblastic gene signatures.**

1358 **(A–E)** Mean pathway activation (by QuSAGE) of **(A)** MSigDB hallmark pathway gene sets<sup>54</sup>, **(B)**  
1359 cancer immunome gene sets<sup>59</sup>, **(C)** top significant GTRD transcription factor target gene sets<sup>55</sup>,  
1360 and cancer-associated fibroblast (CAF) gene sets derived from **(D)** breast tumors and **(E)**  
1361 pancreas and other tumors<sup>19</sup>, in samples with high ColX module expression relative to samples  
1362 with low ColX module expression for each cancer dataset. Gene sets within each block are  
1363 ordered by mean pathway activation rank across all 4 datasets. See **Figure S4** for QuSAGE  
1364 analysis of all modules and **Table S5B** for reported functions of selected enriched TFTs.  
1365 Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*, p.adj < 0.0001.

1366

1367 **Figure 4: COL10A1 and ColX module expression stratify breast and pancreatic cancer**  
1368 **cohorts by survival outcomes.**

1369 **(A and B)** BH-adjusted significance values for multivariate Cox proportional hazards models  
1370 conditioning either **(A)** overall survival (OS) on age, gender, binarized tumor stage, and  
1371 COL10A1 gene expression, or **(B)** disease-free interval (DFI) on age, gender, binarized tumor  
1372 stage, and ColX module eigengene expression (ME). Dotted lines correspond to  $p_{adj} = 0.05$ .  
1373 Per-unit contributions of each significant variable (red bars) to the overall hazard risk in each  
1374 panel is indicated by white percentages (for COL10A1 and ColX ME expression, 1 standard  
1375 deviation = 1 unit). Note that the Cox model for the female pancreatic cancer cohort was not  
1376 conditioned on binarized tumor stage because all samples were classified into the same group.  
1377 See **Figure S5** for full survival analysis results for all WGCNA-inferred modules. **(C)** DFI Kaplan-  
1378 Meier survival curves for the pancreatic cancer cohorts based on %G.A.M.E. groupings. Dotted  
1379 lines correspond to median “survival” (i.e., time to recurrence). Shaded regions represent the  
1380 95% confidence intervals for each group. Log-rank test p-values are shown for each panel.

1381

1382 **Figure 5: Bone marrow and cartilage cells are differentiated by collagen expression**  
1383 **clusters.**

1384 **(A)** Normalized expression of COL10A1 across bone marrow and articular cartilage cell types ( $n$   
1385 = 3 each). Mean normalized expressions are indicated in parentheses. Significance values were  
1386 computed using DESeq2, BH-adjusted across all genes analyzed. **(B)** Log-normalized  
1387 expression of ColX module genes in each cell type ( $n = 3$  each). Normalized gene-wise  
1388 expressions were averaged prior to log-transformation with pseudocount of 1. Significance  
1389 values were computed using the paired Wilcoxon signed-rank test on log-normalized counts. **(C)**  
1390 Relative expression of collagen genes across cell samples ( $n = 3$  each). Normalized expression  
1391 values are scaled by rows. COL10A1 is indicated by a black box in the left column; genes

1392 contributing to the same parent collagen are indicated by same-color boxes; and genes which  
1393 are the sole contributor to their parent collagen are indicated by gray boxes. Note that  
1394 *COL6A4P1*, *COL6A5*, *COL20A1*, and *COL23A1* were filtered out as “low-expression” genes.  
1395 See **Figure S7A** for normalized (unscaled) expression of all collagen genes across cell types.  
1396 **(D)** WGCNA module assignments for all collagens in each cancer dataset. See **Table S2B** for  
1397 WGCNA module assignments for all genes. Black cells indicate the ColX module for each  
1398 column. “-” indicates low-expression genes which were filtered out; genes labeled “0” were not  
1399 assigned to any module by WGCNA. See **Figure S7B** for normalized expression of all collagen  
1400 genes across TCGA cohorts.

1401

1402 **Figure 6: Tumoral proportions of osteoarthritic cell types trend with ColX module**  
1403 **expression and cancer-associated pathway activity.**

1404 **(A and B)** Sample-wise absolute (*top*) and relative (*bottom*) proportions of 4 bone marrow and  
1405 cartilage cell types inferred by CIBERSORTx for TCGA **(A)** breast and **(B)** pancreatic cancer  
1406 cohorts. Color bars (*middle*) indicate relative cancer-specific ColX module expression; samples  
1407 were ordered by increasing %G.A.M.E. See Methods section and **Figure S3** for details on  
1408 %G.A.M.E. metric. **(C and D)** Spearman correlations between relative cancer-specific ColX  
1409 module expression and **(C)** absolute or **(D)** relative OA cell type proportions inferred by  
1410 CIBERSORTx. Raw p-values were Bonferroni-corrected across rows. Significance values: \*,  $p <$   
1411 0.05; \*\*,  $p < 0.01$ ; \*\*\*,  $p < 0.001$ ; \*\*\*\*,  $p < 0.0001$ ; n.s., not significant. NCSC, normal cartilage  
1412 stromal cells; OA-MSc, osteoarthritis mesenchymal stromal cells; OAC, osteoarthritis  
1413 chondrocytes; BMSC, bone marrow stromal cells. **(E and F)** Bubble plots of **(E)** MSigDB  
1414 hallmark pathway gene sets<sup>54</sup> and **(F)** top significant GTRD transcription factor target gene  
1415 sets<sup>55</sup> in cell type-specific marker gene sets. Of note, the TFT gene set attributed to *IGLV5-37* in  
1416 the GTRD database actually represents targets of the fusion oncoprotein *SS18-SSX*, as  
1417 described in the text. See Methods section for definition of bone marrow and cartilage cell type-

1418 specific marker genes. Significance values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001;  
1419 \*\*\*\*, p.adj < 0.0001.

1420

1421 **Figure 7: Pathological expression of *COL10A1* fosters immunosuppressive, fibroblastic**  
1422 **microenvironments in cancer, bone marrow, and cartilage.**

1423 Graphical abstract summarizing the findings presented in this study. In brief, a specifically  
1424 expressed collagen, *COL10A1*, connects the ECM and tissue microenvironments across cancer  
1425 and bone. The pathological contributions of ColX and its associated gene networks are  
1426 especially prominent in breast and pancreatic tumors, and mimic the development of a similar  
1427 inflammatory and fibroblast-dominated environment seen in bone marrow and cartilage changes  
1428 in OA. A central outcome of this shared impact is the epithelial-to-mesenchymal transition  
1429 (EMT), which contributes to disease progression in both contexts. *All images were sourced from*  
1430 *Bioicons (<https://bioicons.com>) and are licensed for public use by Servier*  
1431 *(<https://smart.servier.com>) under CC-BY 3.0 (<https://creativecommons.org/licenses/by/3.0>) or by*  
1432 *DBCLS (<https://togotv.dbcls.jp/en/pics.html>) under CC-BY 4.0*  
1433 *(<https://creativecommons.org/licenses/by/4.0>).*

1434

### 1435 **SUPPLEMENTAL FIGURES**

1436 **Figure S1: TCGA ColX modules are preserved in cancer microarray datasets.**

1437 **(A and B)** Module preservation scores ( $Z_{summary}$ ) for all TCGA RNA-Seq-derived **(A)** breast and  
1438 **(B)** pancreatic cancer WGCNA modules in comparably-sized microarray tumor datasets. ColX  
1439 modules are indicated by **(A)** pink or **(B)** orange bars. Dotted lines indicate “high preservation”  
1440 threshold of  $Z_{summary} = 10$  as defined by the authors of WGCNA.

1441

1442 **Figure S2: TCGA WGCNA modules are enriched for numerous hallmark pathways.**

1443 **(A–D)** Bubble plots of MSigDB hallmark pathway gene sets<sup>54</sup> enrichment in WGCNA modules  
1444 from **(A)** breast cancer, **(B)** pancreatic cancer, **(C)** male pancreatic cancer, and **(D)** female  
1445 pancreatic cancer cohorts. ColX modules for each dataset are indicated by bolded labels. See  
1446 **Table S4** for significance values and **Figure 2D** for ColX module-specific enrichment results. **(E**  
1447 **and F)** Overlap of EMT hallmark pathway genes within ColX WGCNA modules from TCGA **(E)**  
1448 breast and pancreatic cancer and **(F)** gender-segregated pancreatic cancer datasets. Genes  
1449 comprising each sector are listed alphabetically.

1450

1451 **Figure S3: The G.A.M.E. metric effectively proxies ColX module expression and improves**  
1452 **sample stratification.**

1453 **(A–D)** Relationship between ColX module eigengene (ME) expression and proportion of “Genes  
1454 Above Median Expression” (G.A.M.E.) for **(A)** breast cancer, **(B)** pancreatic cancer, **(C)** male  
1455 pancreatic cancer, and **(D)** female pancreatic cancer cohorts. Colored curves represent  
1456 densities of ME (*right*) and G.A.M.E. (*top*) variables for each panel. Blue dotted lines indicate  
1457 Jenks natural breakpoints defining 3 clusters (“low”, “medium”, and “high”). Spearman  
1458 correlations ( $\rho$ ) are shown for each panel.

1459

1460 **Figure S4: TCGA WGCNA module expression tracks differential activity of immune,**  
1461 **transcription factor, and cancer-associated fibroblast signatures.**

1462 **(A–P)** Mean pathway activation (by QuSAGE) of various gene sets in samples with high module  
1463 expression relative to samples with low module expression for each WGCNA module and TCGA  
1464 cohort. **(A–D)** Results for cancer immunome genesets in **(A)** breast cancer, **(B)** pancreatic  
1465 cancer, **(C)** male pancreatic cancer, and **(D)** female pancreatic cancer cohorts. Gene sets within  
1466 each block are ordered by significance in respective ColX module; see **Figure 3B** for focused  
1467 comparison. **(E–H)** Results for transcription factor target lists in **(E)** breast cancer, **(F)** pancreatic

1468 cancer, **(G)** male pancreatic cancer, and **(H)** female pancreatic cancer cohorts. Gene sets within  
1469 each block are ordered by effect size in respective ColX module; see **Figure 3C** for focused  
1470 comparison. **(I–P)** Results for cancer-associated fibroblast (CAF) gene sets derived from **(I–L)**  
1471 breast tumor scRNA-Seq data and **(M–P)** pancreas/other tumor scRNA-Seq data, in **(I and M)**  
1472 breast cancer, **(J and N)** pancreatic cancer, **(K and O)** male pancreatic cancer, and **(L and P)**  
1473 female pancreatic cancer cohorts. Gene sets within each block are ordered by significance in  
1474 respective ColX module; see **Figures 3D–E** for focused comparison. See Methods section for  
1475 sample grouping process (note the same procedure for calculating G.A.M.E. and assigning  
1476 “low” and “high” labels was applied to each WGCNA module respectively). Significance values  
1477 were corrected across each module individually (by rows). Significance values: \*, p.adj < 0.05;  
1478 \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*, p.adj < 0.0001.

1479

1480 **Figure S5: TCGA WGCNA module expression correlates with variable survival risk.**

1481 **(A–D)** BH-adjusted significance values for multivariate Cox proportional hazards models  
1482 conditioning overall survival (OS), disease-specific survival (DSS), disease-free interval (DFI), or  
1483 progression-free interval (PFI) on age, gender, binarized tumor stage, and module eigengene  
1484 (ME) expression for WGCNA modules from **(A)** breast cancer, **(B)** pancreatic cancer, **(C)** male  
1485 pancreatic cancer, and **(D)** female pancreatic cancer cohorts. ColX modules for each dataset  
1486 are indicated by bolded labels (see **Figure 4B** for focused ColX module results). Significance  
1487 values: \*, p.adj < 0.05; \*\*, p.adj < 0.01; \*\*\*, p.adj < 0.001; \*\*\*\*, p.adj < 0.0001.

1488

1489 **Figure S6: Osteoarthritis cell type-specific markers are enriched for ColX module genes.**

1490 **(A–D)** Bubble plots of OA cell type-specific marker gene enrichment in WGCNA modules from  
1491 **(A)** breast cancer, **(B)** pancreatic cancer, **(C)** male pancreatic cancer, and **(D)** female pancreatic  
1492 cancer cohorts. ColX modules for each dataset are indicated by bolded labels. **(E)** Overlap of  
1493 genes within breast and pancreatic cancer ColX modules and OA cell type-specific gene sets.

1494 **(F)** Overlap of EMT pathway genes within breast and pancreatic cancer ColX modules and OA  
1495 cell type-specific gene sets. Unlabeled sectors represent 0 gene overlap. NCSC is not shown as  
1496 no NCSC-specific markers overlap with EMT pathway genes. NCSC, normal cartilage stromal  
1497 cells; OA-MSc, osteoarthritis mesenchymal stromal cells; OAC, osteoarthritis chondrocytes;  
1498 BMSC, bone marrow stromal cells.

1499

1500 **Figure S7: Collagen gene expression varies across bone and cartilage cell types and**  
1501 **TCGA cohorts.**

1502 **(A)** Normalized expression of all collagen genes expressed at nontrivial levels in bone marrow  
1503 and cartilage cell types. Note that *COL6A4P1*, *COL6A5*, *COL20A1*, and *COL23A1* were filtered  
1504 out as “low-expression” genes and are omitted here. NCSC, normal cartilage stromal cells; OA-  
1505 MSC, osteoarthritis mesenchymal stromal cells; OAC, osteoarthritis chondrocytes; BMSC, bone  
1506 marrow stromal cells. **(B)** Normalized expression of all collagen genes expressed at nontrivial  
1507 levels in TCGA cohorts. Note that *COL6A4P1*, *COL6A5*, *COL20A1*, and *COL26A1* were filtered  
1508 out as “low-expression” genes and are omitted here; additionally, *COL2A1* was only nontrivially  
1509 expressed in BRCA.

1510

## 1511 **SUPPLEMENTAL TABLES**

1512 **Table S1: Overview of cancer datasets used in this study.**

1513 **(A)** Breast and pancreatic cancer datasets used in this study. **(B)** Citations of breast cancer  
1514 GEO accessions used in this study.

1515

1516 **Table S2: Gene modules characterized in this study.**

1517 **(A)** WGCNA-generated *COL10A1* (ColX) modules from TCGA breast and pancreatic cancer  
1518 datasets, as well as gender-specific pancreatic cancer datasets. **(B)** WGCNA module  
1519 assignments for all expressed genes in each dataset. Modules are numbered in descending



1520 order of size, beginning with module 1. “-” indicates genes which were filtered out based on low  
1521 expression; genes labeled with a 0 were not assigned to any module by WGCNA.

1522

1523 **Table S3: Gene ontology and Reactome pathway enrichment analysis of ColX modules.**

1524 **(A–E)** Gene ontology pathway enrichment for **(A)** breast cancer, **(B)** pancreatic cancer, **(C)**  
1525 overlap between breast and pancreatic cancer, **(D)** male pancreatic cancer, and **(E)** female  
1526 pancreatic cancer ColX modules. **(F–J)** Reactome pathway enrichment for **(F)** breast cancer,  
1527 **(G)** pancreatic cancer, **(H)** overlap between breast and pancreatic cancer, **(I)** male pancreatic  
1528 cancer, and **(J)** female pancreatic cancer ColX modules. All significantly enriched Reactome  
1529 pathways/GO terms ( $q < 0.05$ ) are shown.

1530

1531 **Table S4: Hallmark pathway enrichment analysis of WGCNA modules.**

1532 **(A–D)** Hallmark pathway enrichment for **(A)** breast cancer, **(B)** pancreatic cancer, **(C)** male  
1533 pancreatic cancer, and **(D)** female pancreatic cancer ColX modules. All p-values shown were  
1534 BH-corrected across all 50 hallmark pathways for each module. ColX modules for each dataset  
1535 (#8, #13, #7, and #23, respectively) are shown in the first column for clarity. Related to **Figure**  
1536 **S2A–D**.

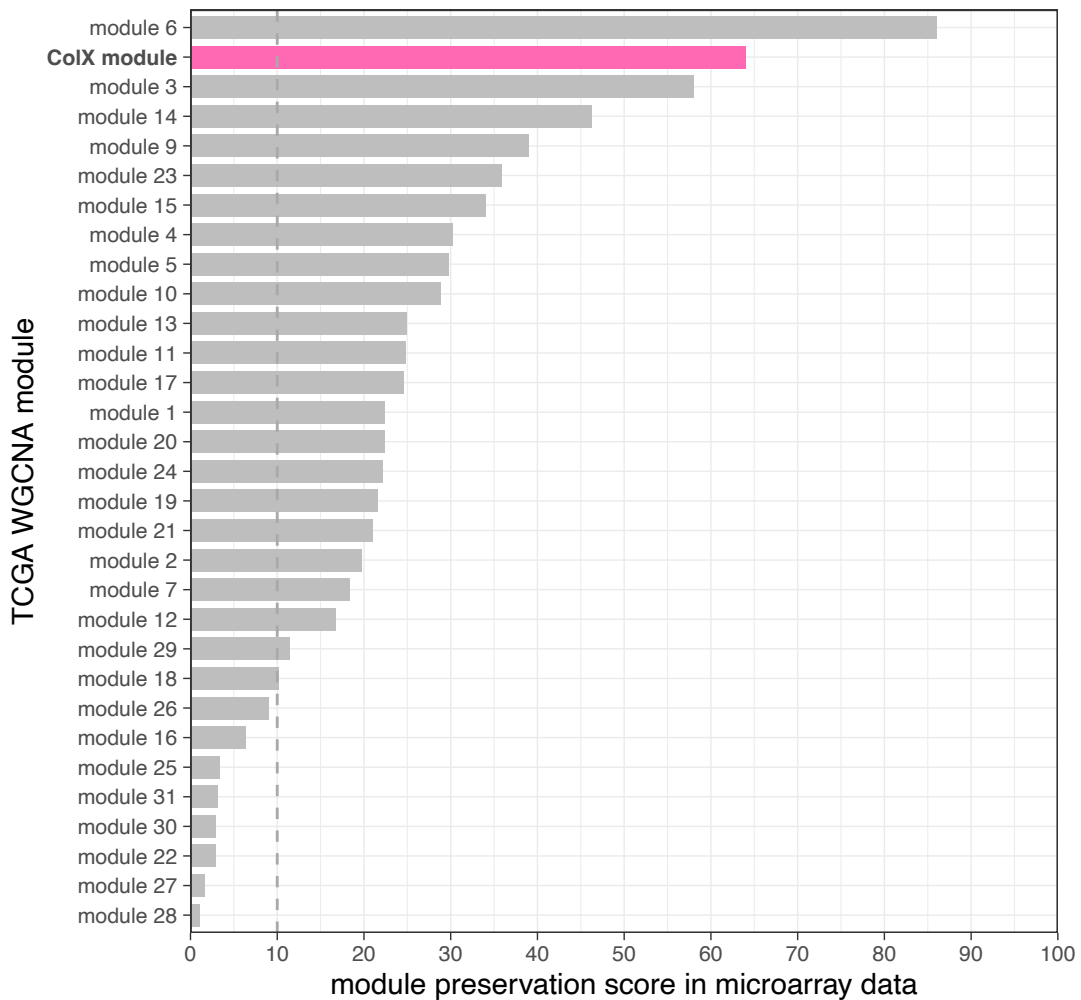
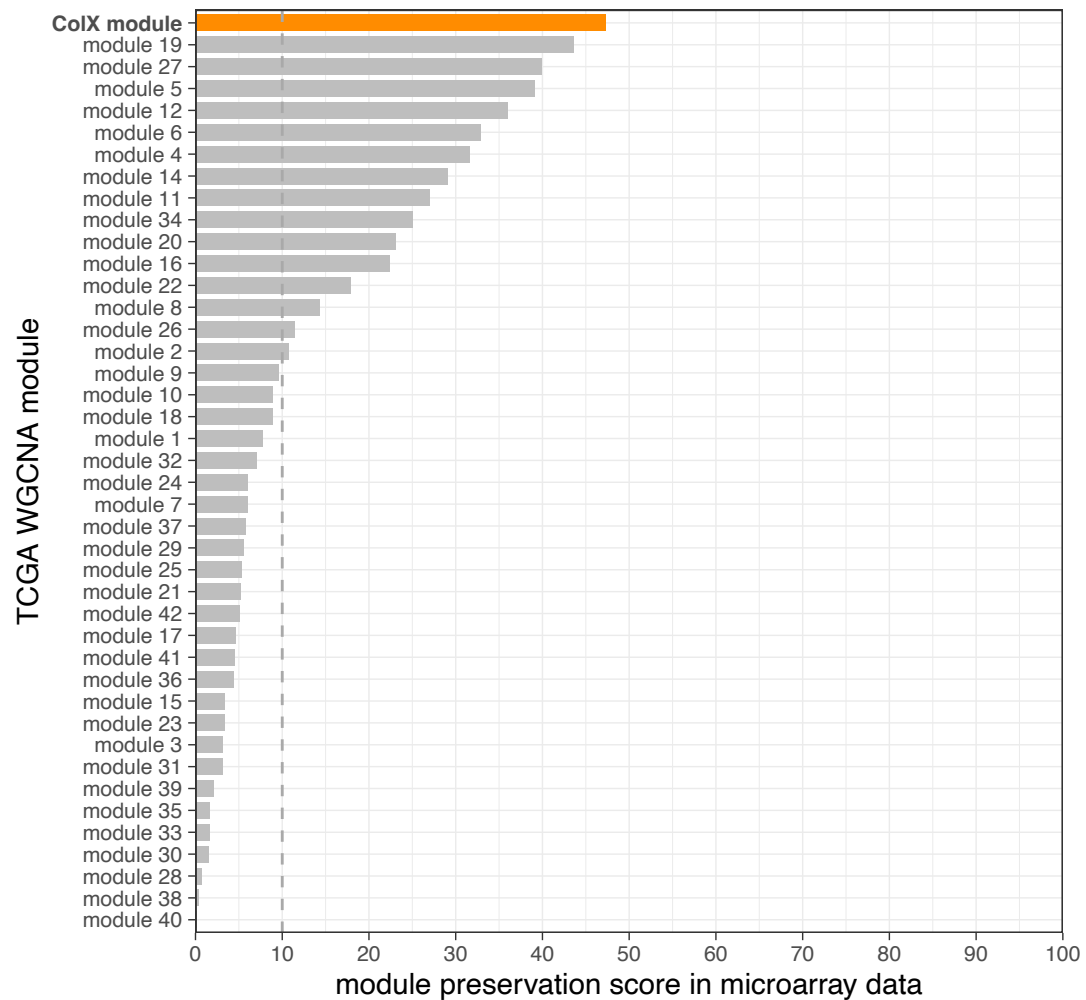
1537

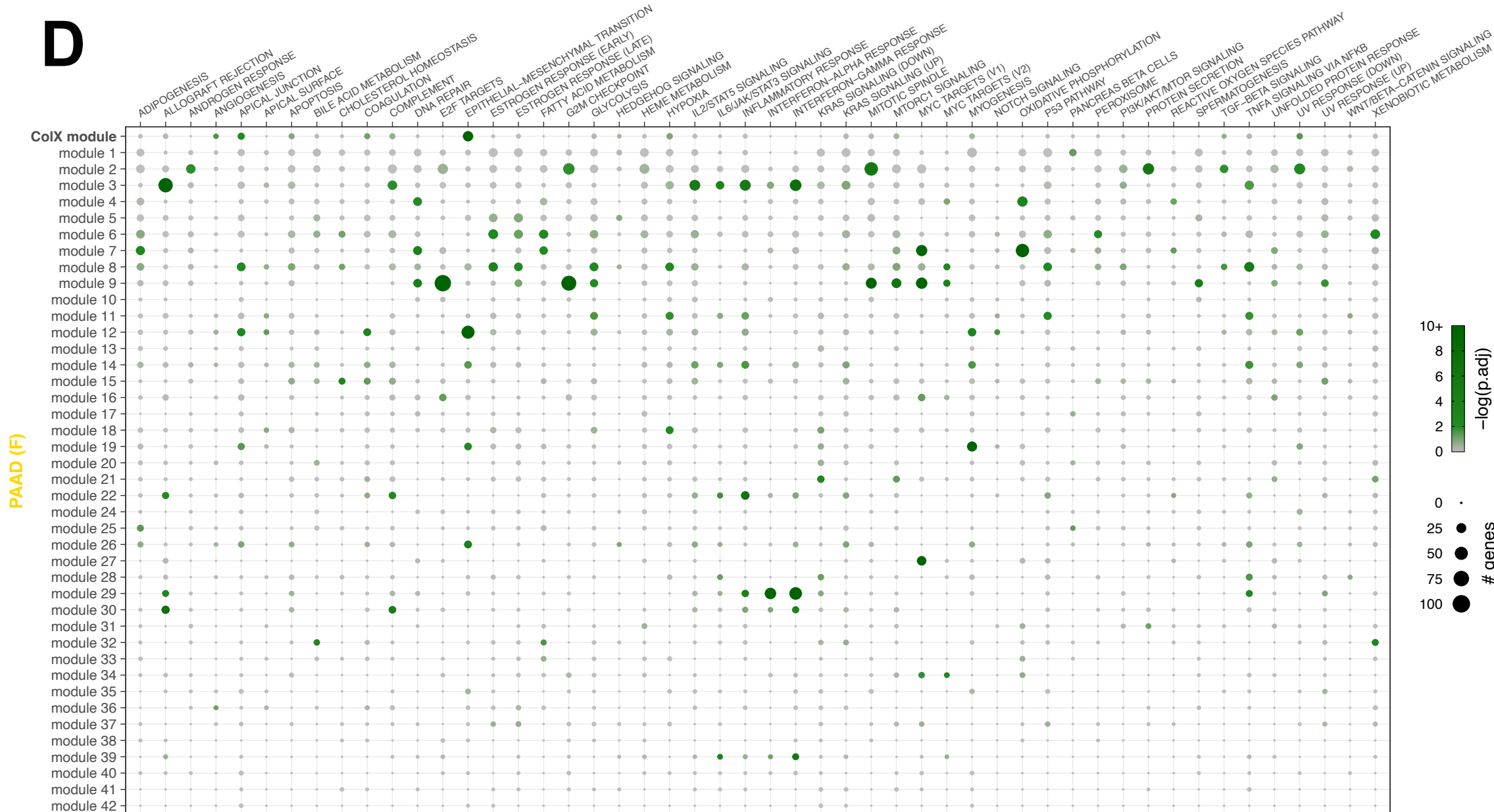
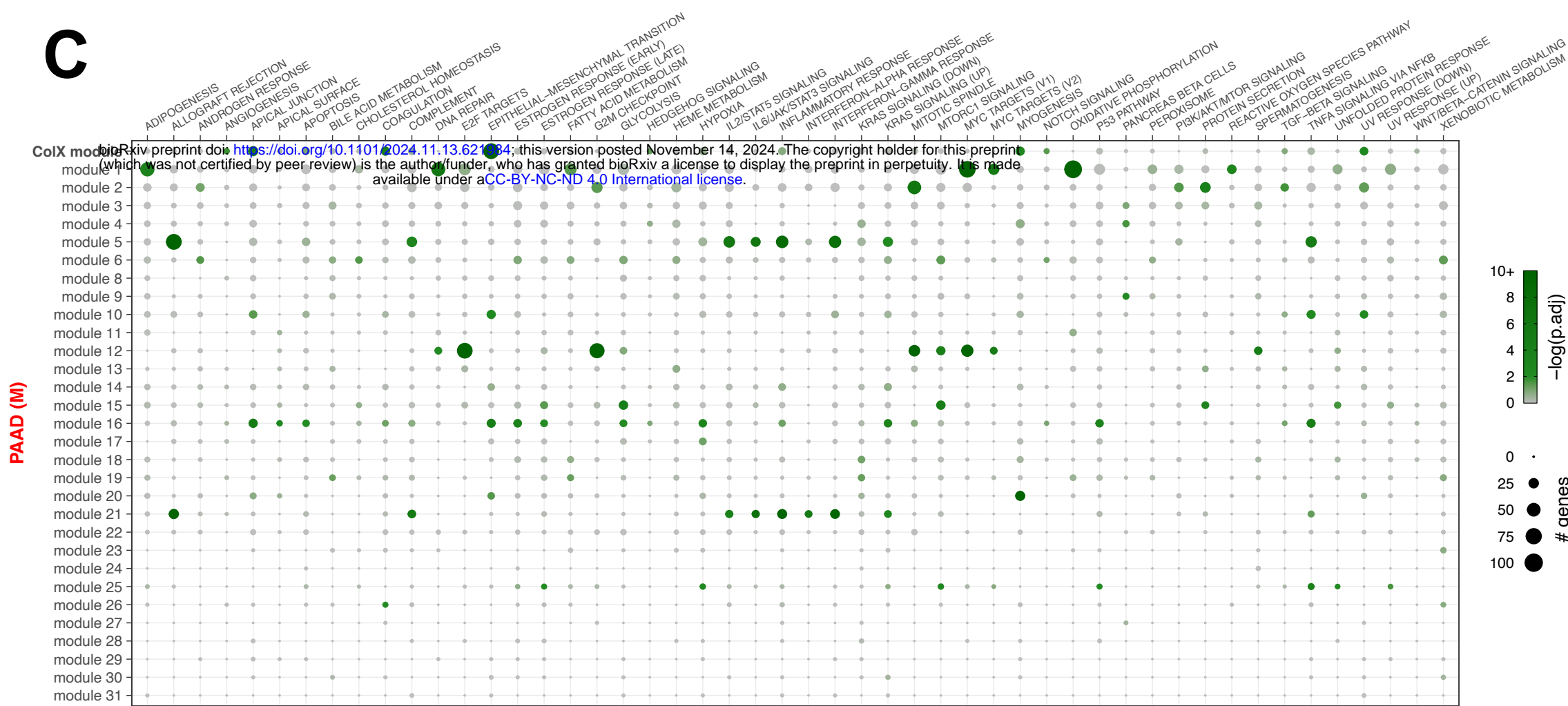
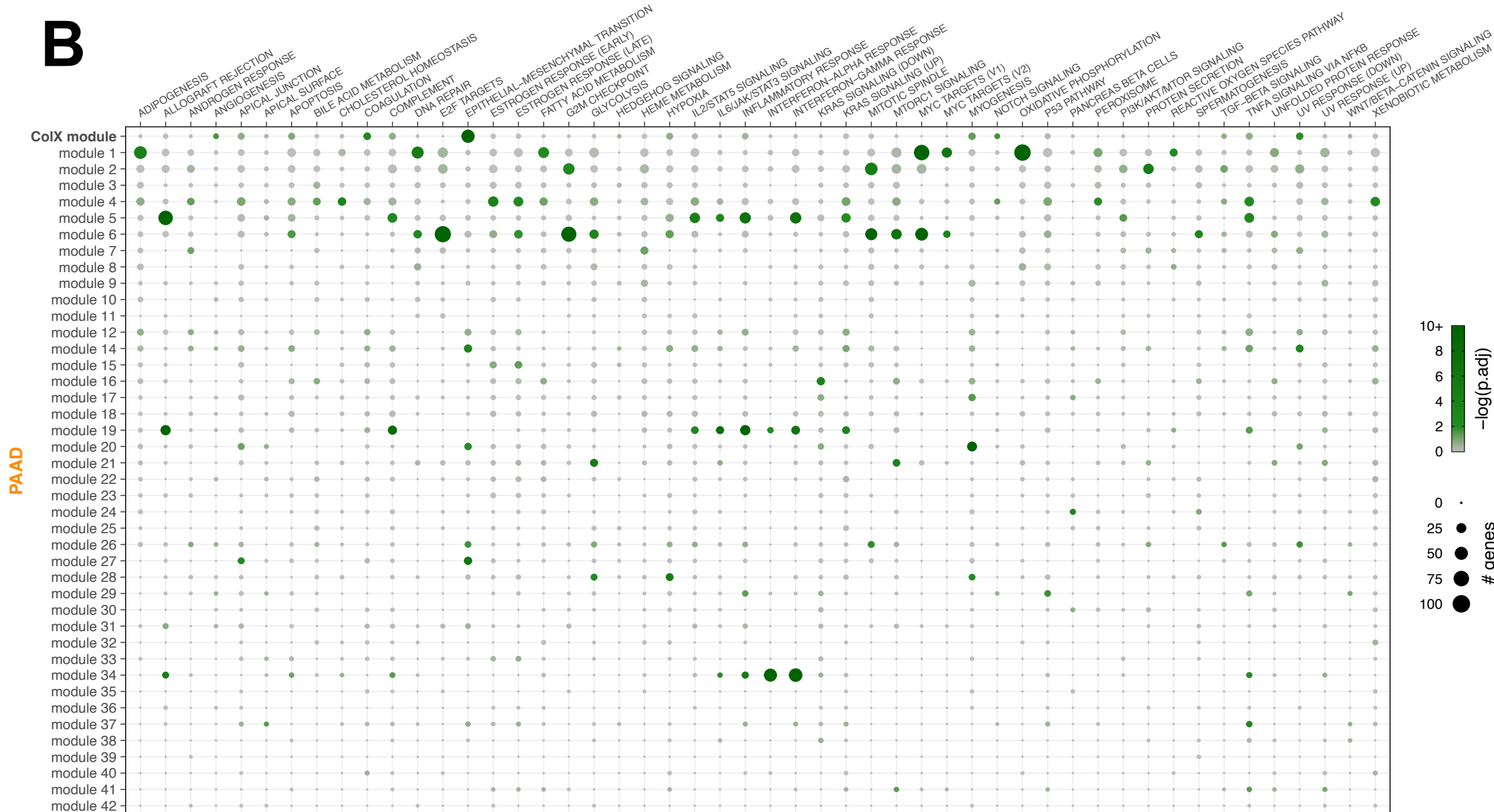
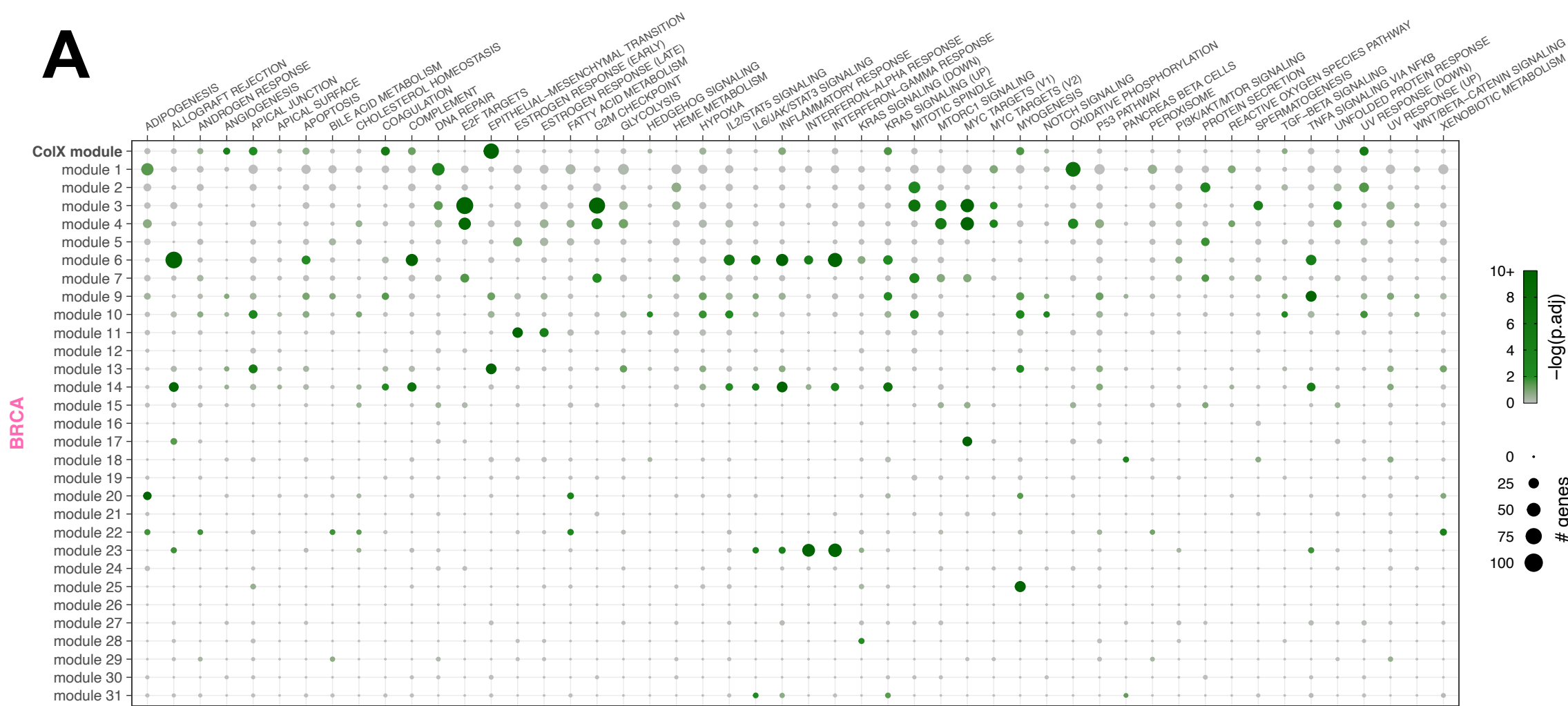
1538 **Table S5: TFs implicated in ColX modules.**

1539 **(A)** GTRD TFTs enriched in BRCA and PAAD ColX modules. Cancer-relevant TF functions and  
1540 overlapping targets are shown where applicable. (Related to **Figure 2E**.) **(B)** Selected QuSAGE-  
1541 significant transcription factors whose targets are differentially activated/inactivated between  
1542 tumors with high and low ColX module expression. Cancer/OA-relevant TF functions are shown  
1543 where applicable. (Related to **Figure 3C**.)

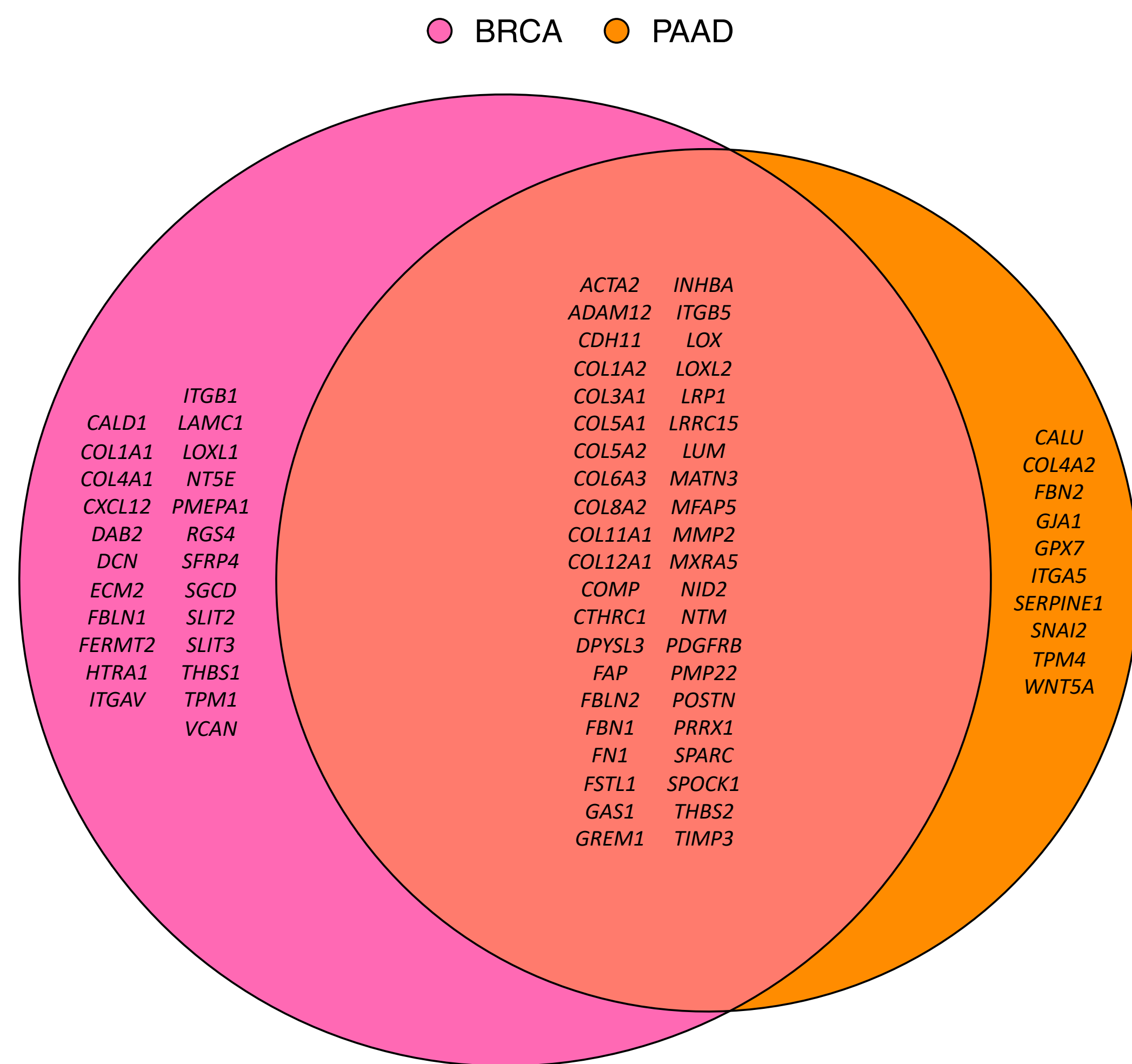
1544

1545 **Table S6: Differential expression of bone marrow and cartilage cell type-specific genes.**  
1546 List of all genes defined as specific to normal cartilage stromal cells/NCSCs, OA mesenchymal  
1547 stromal cells/OA-MSCs, OA chondrocytes/OACs, and bone marrow stromal cells/BMSCs,  
1548 based on OA RNA-Seq data. Differential expression results were computed across all cells and  
1549 statistics were extracted for each pairwise comparison. See Methods section for definition of  
1550 “cell type-specific” genes.

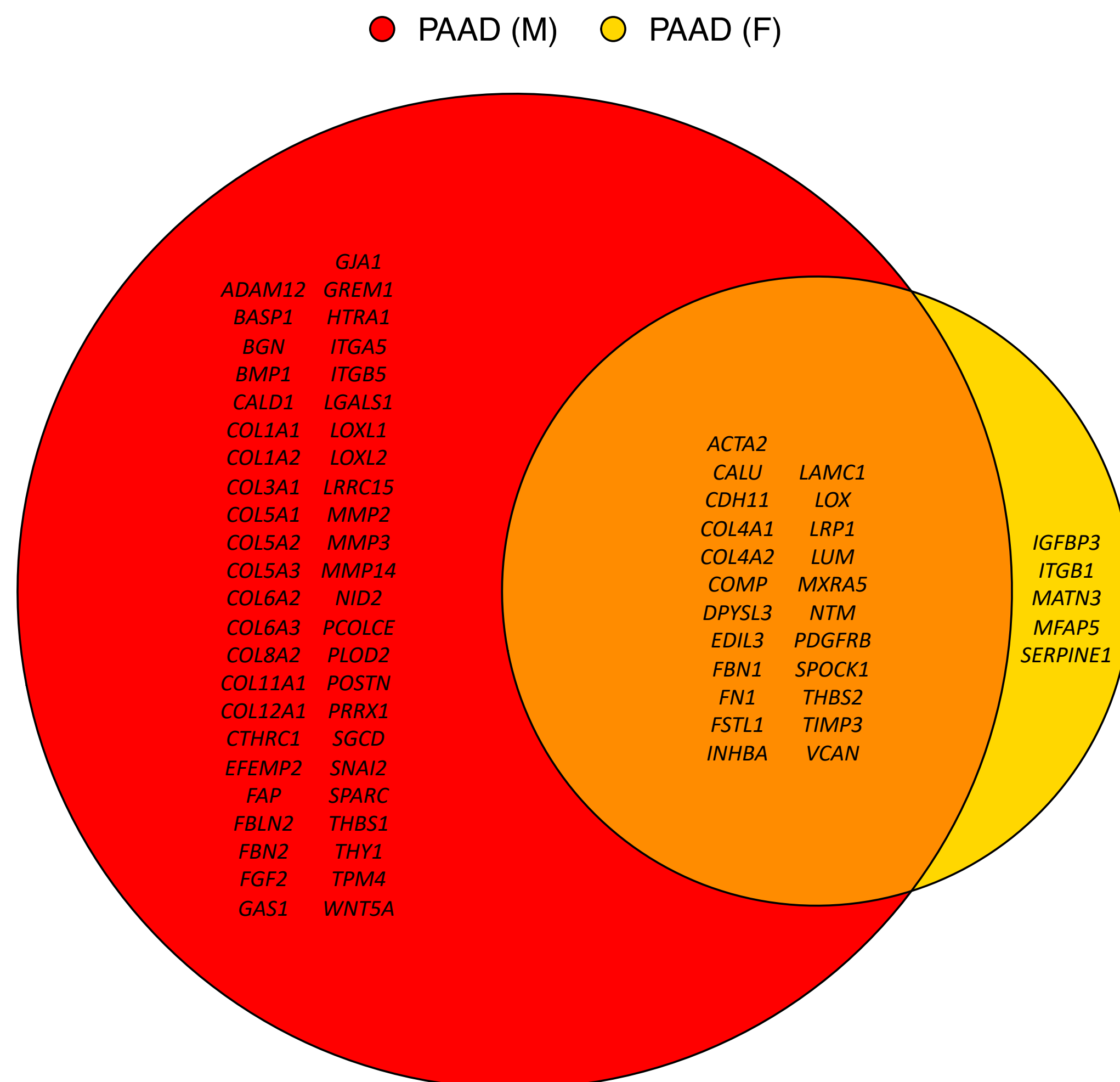
**A****BRCA****B****PAAD**



## E EMT Hallmark Pathway Genes in ColX Modules

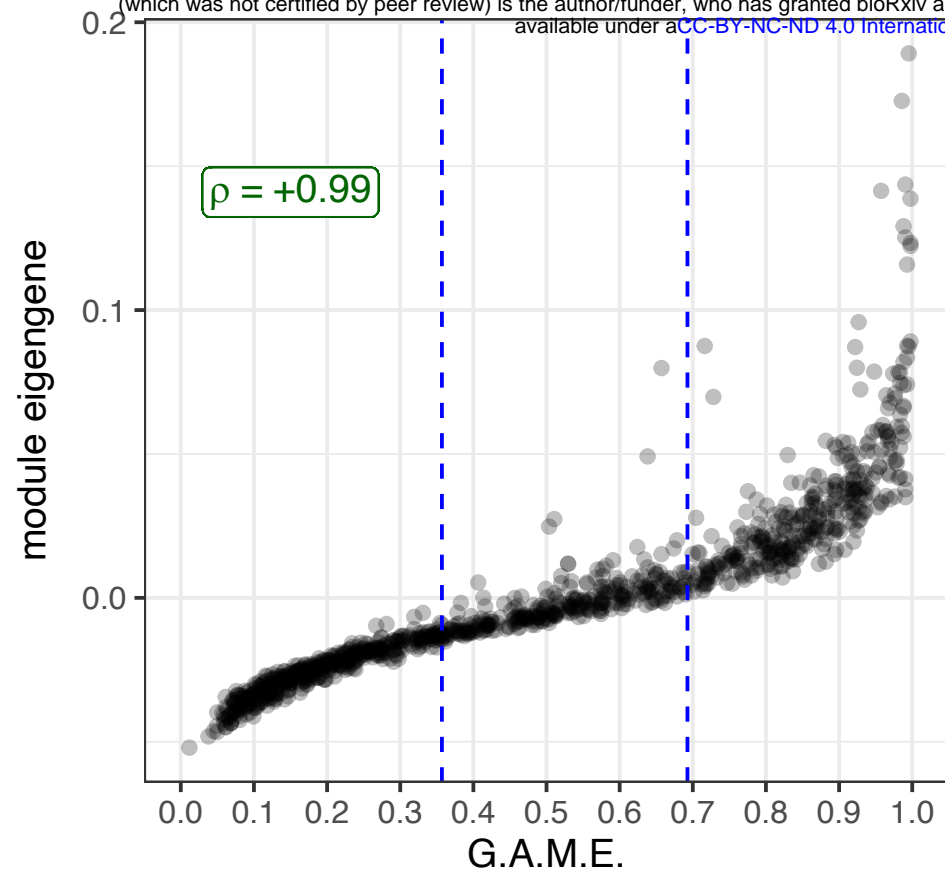
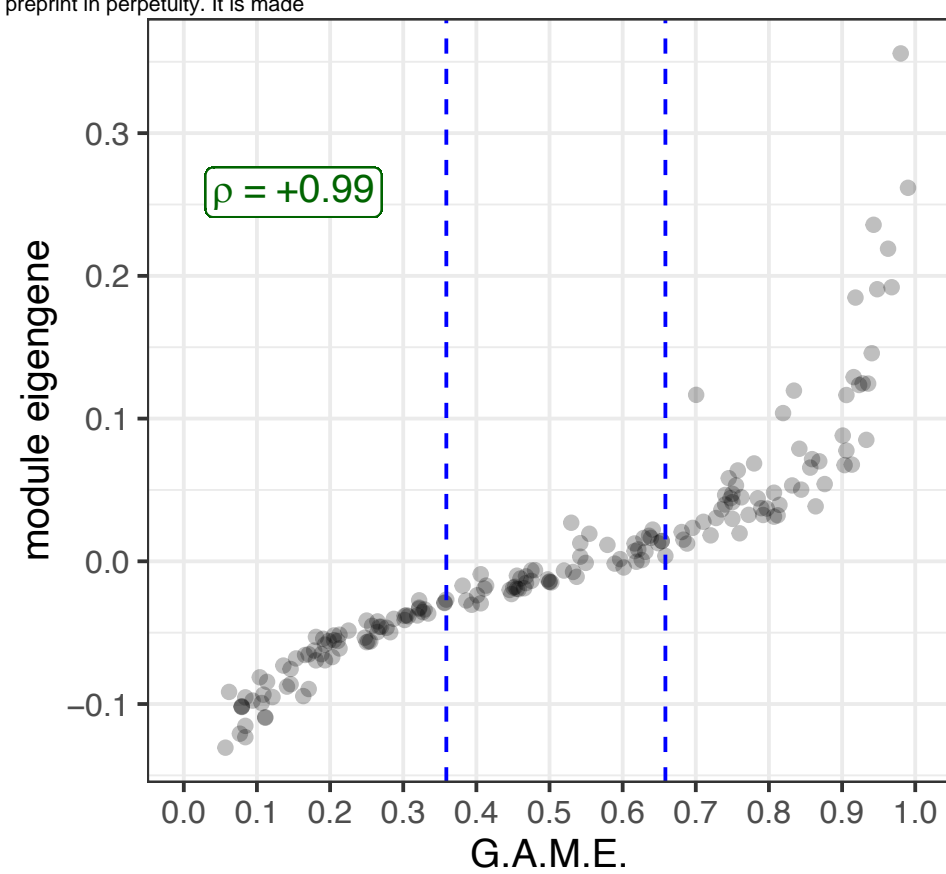
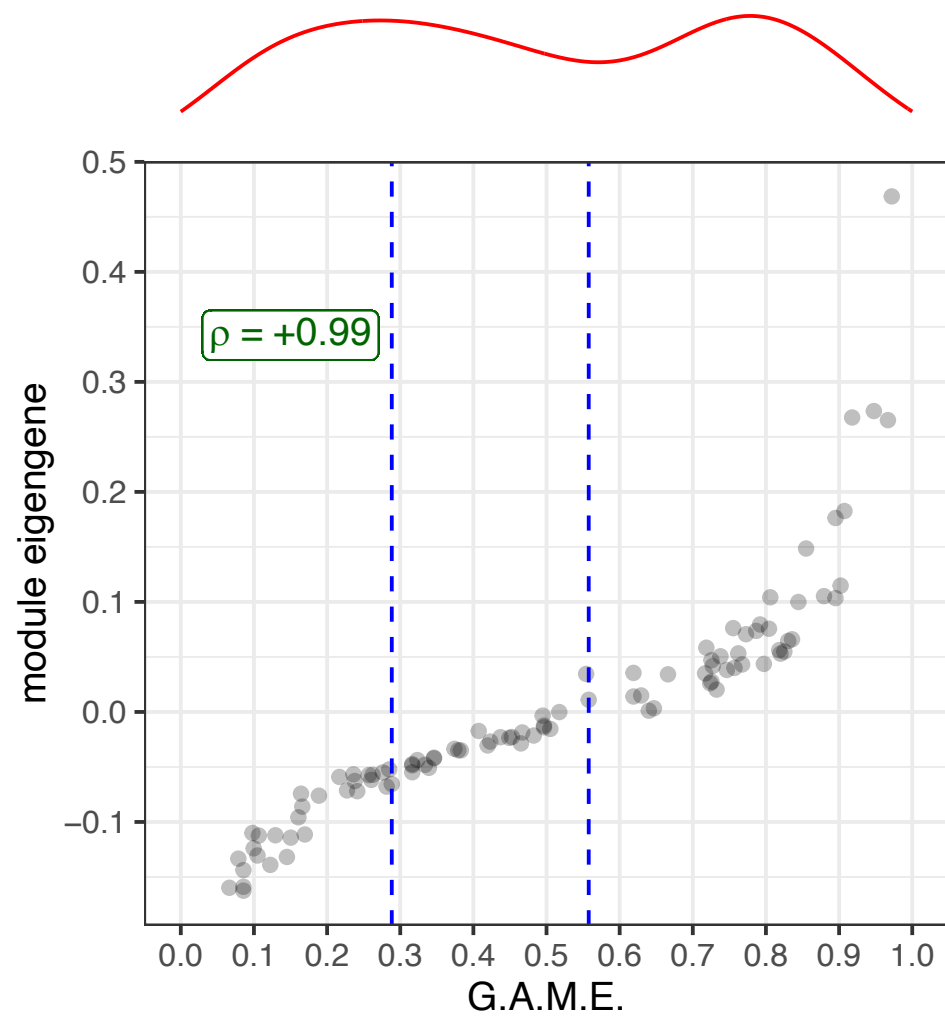
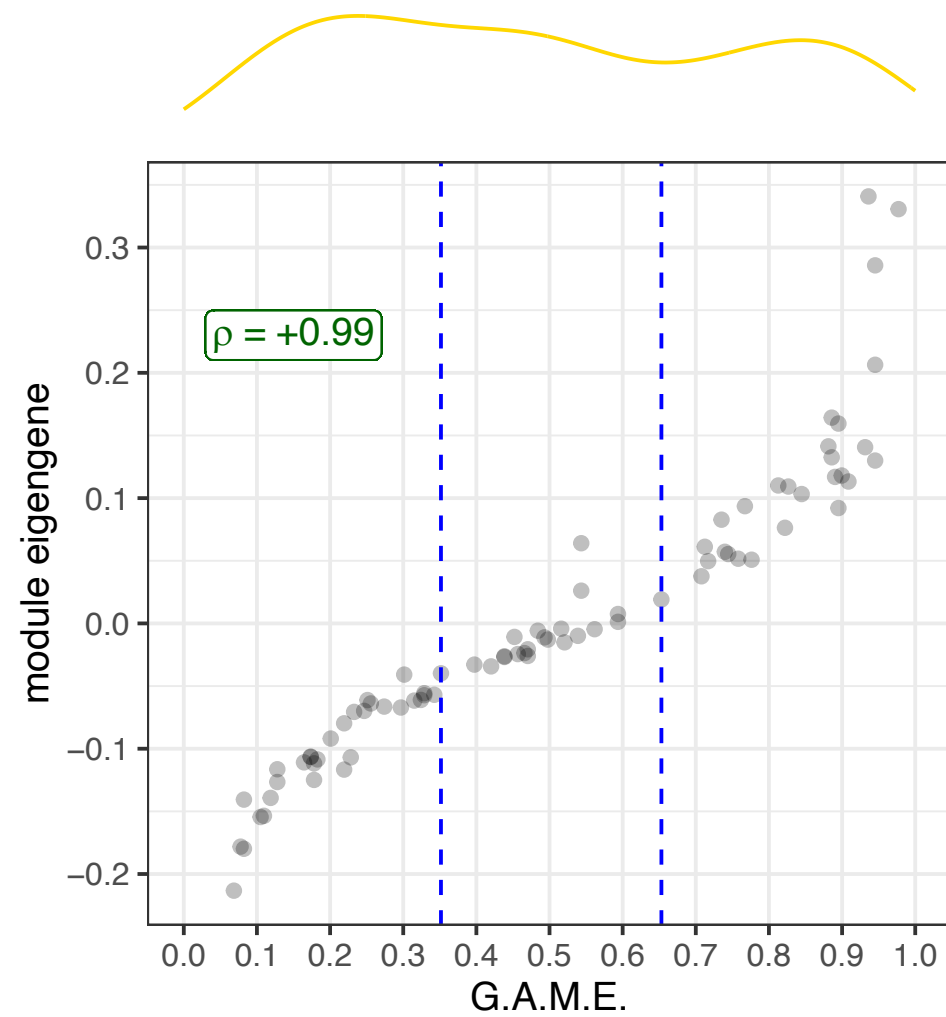


## F EMT Hallmark Pathway Genes in ColX Modules

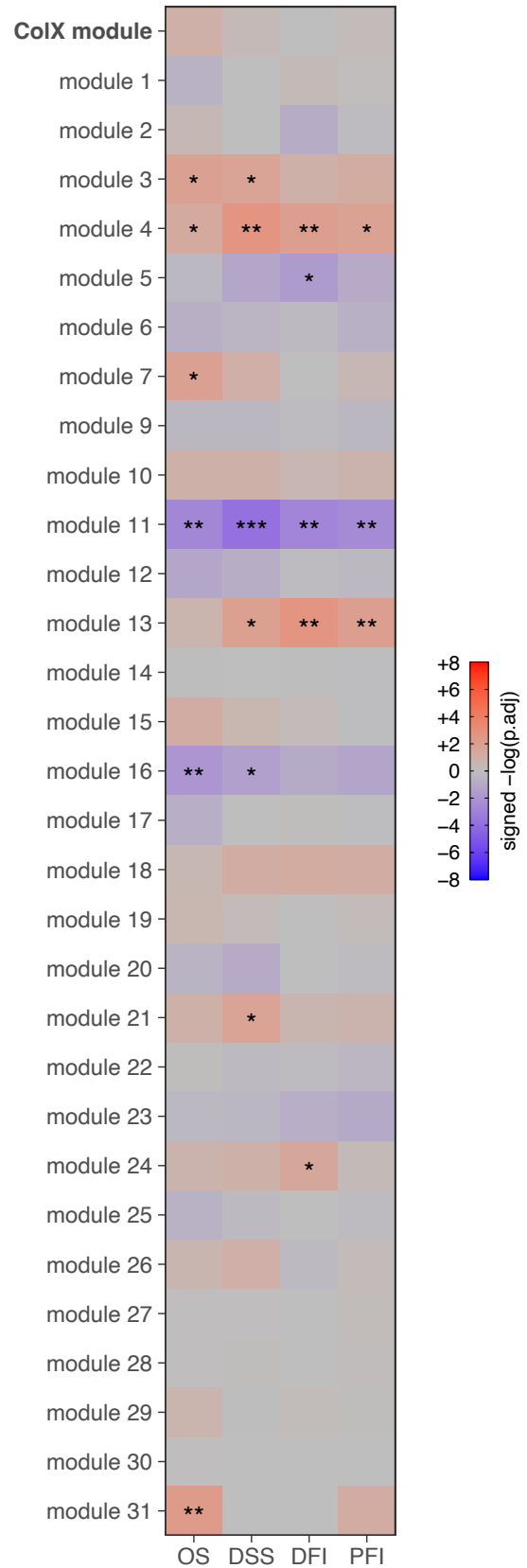
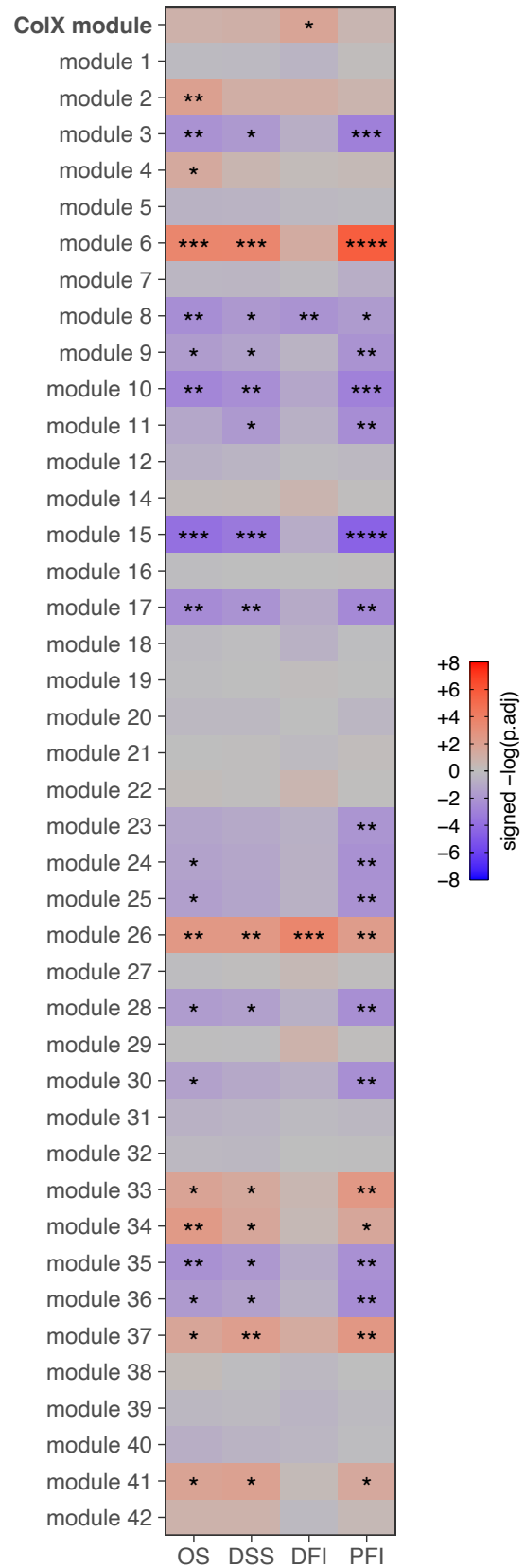
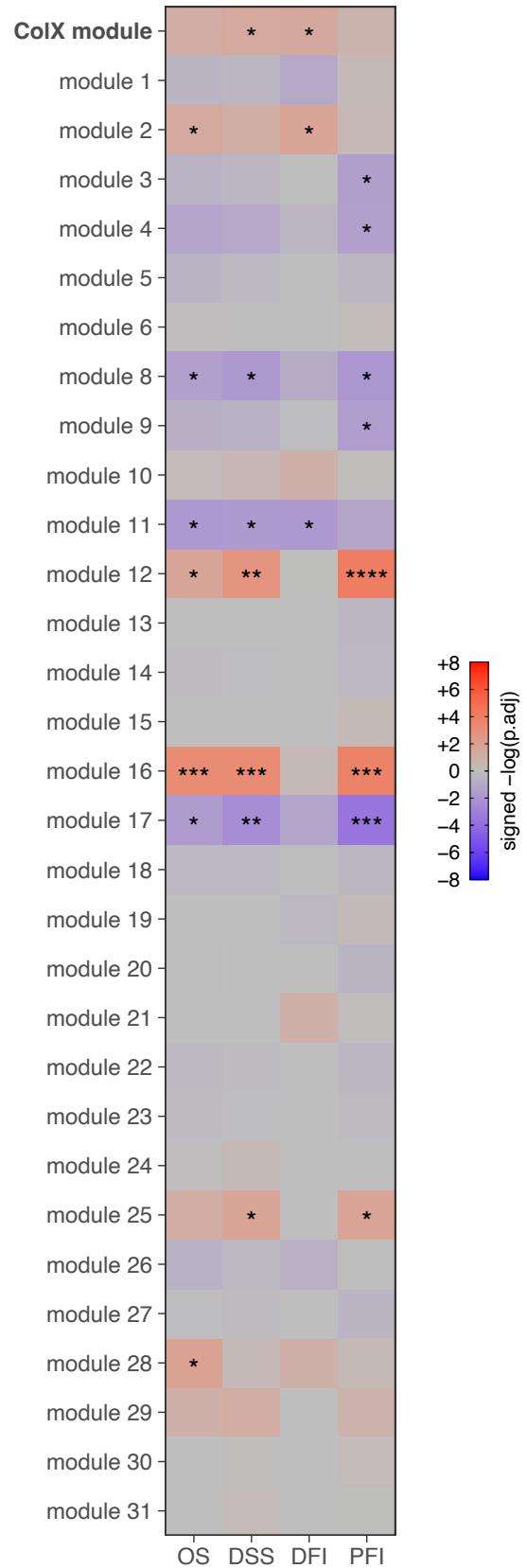
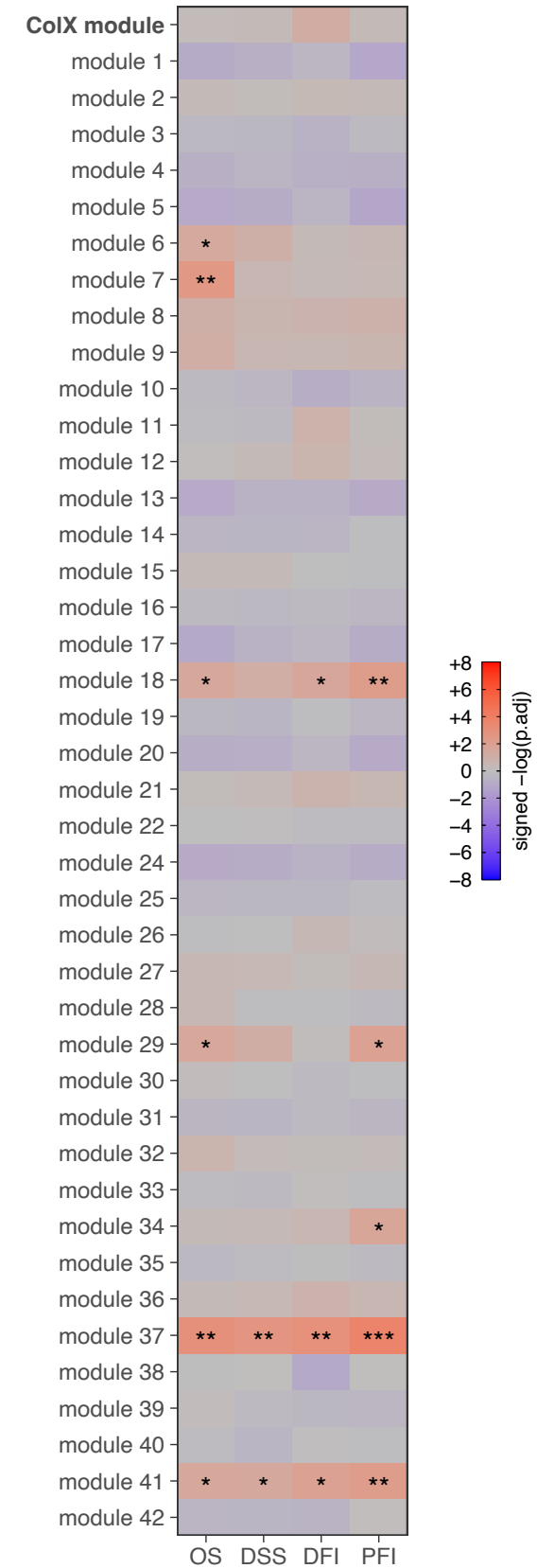


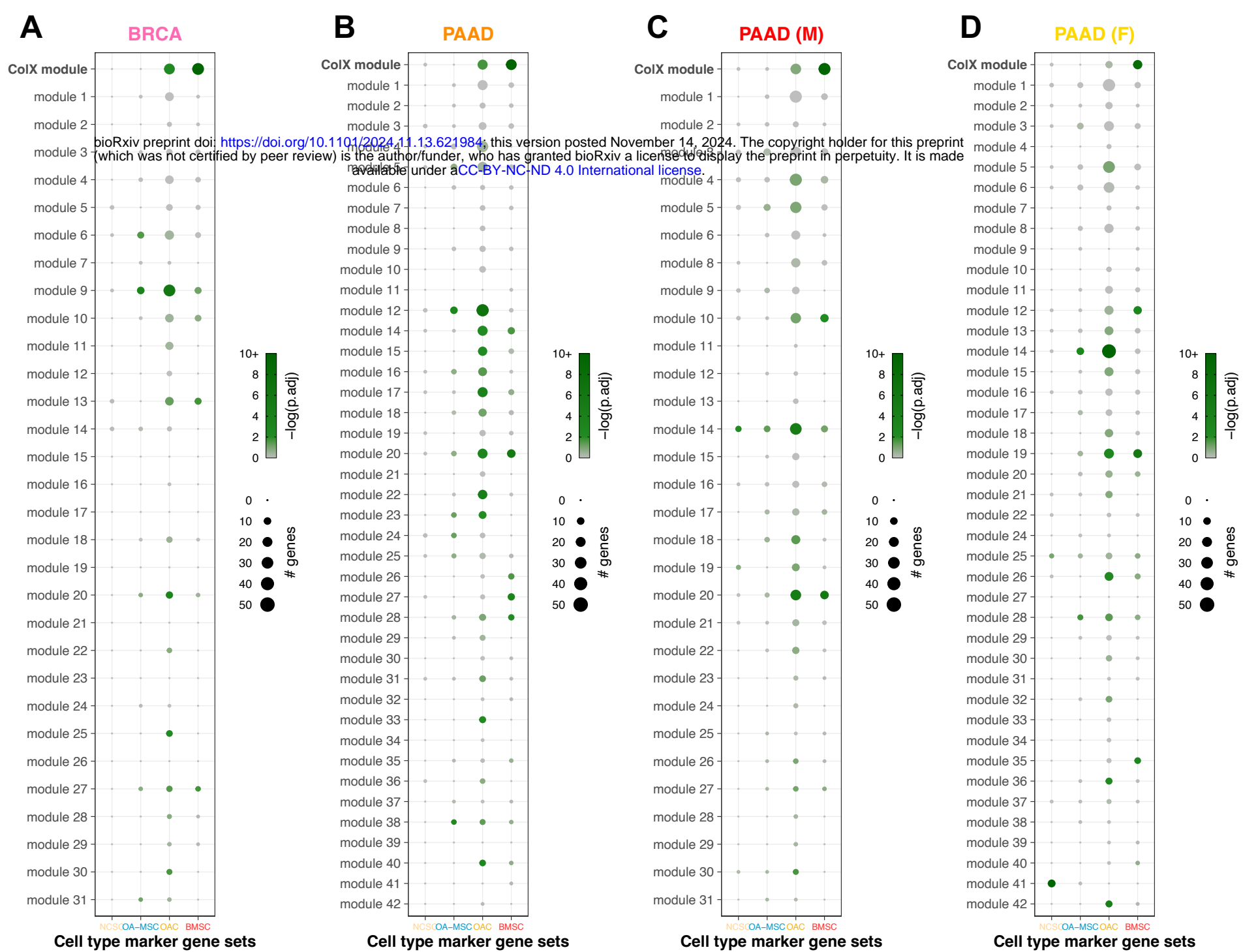
**A****BRCA**

bioRxiv preprint doi: <https://doi.org/10.1101/2024.11.13.621984>; this version posted November 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

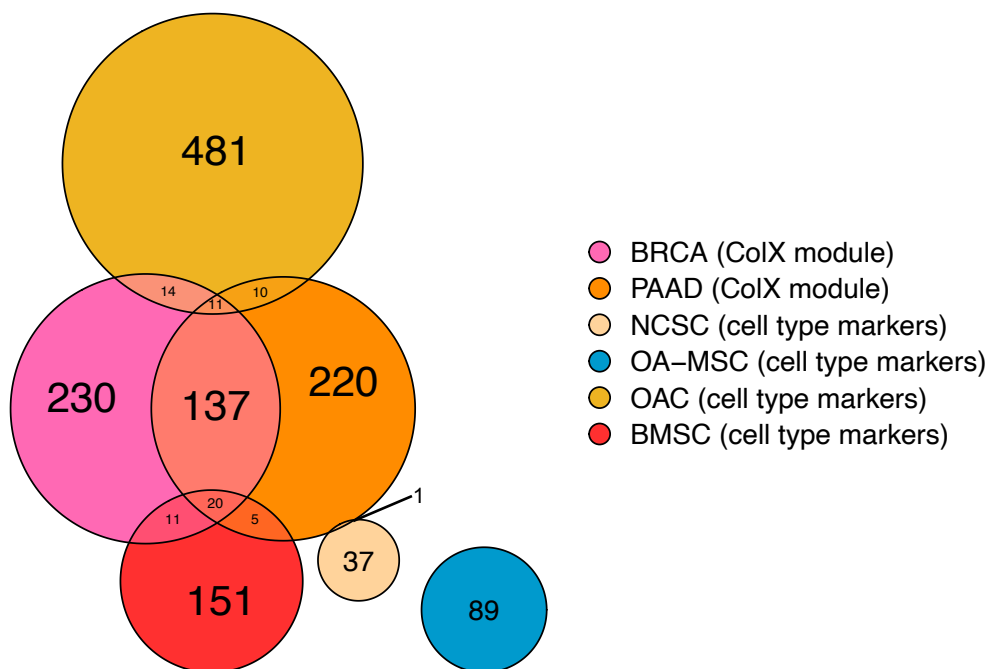
**B****PAAD****C****PAAD (M)****D****PAAD (F)**



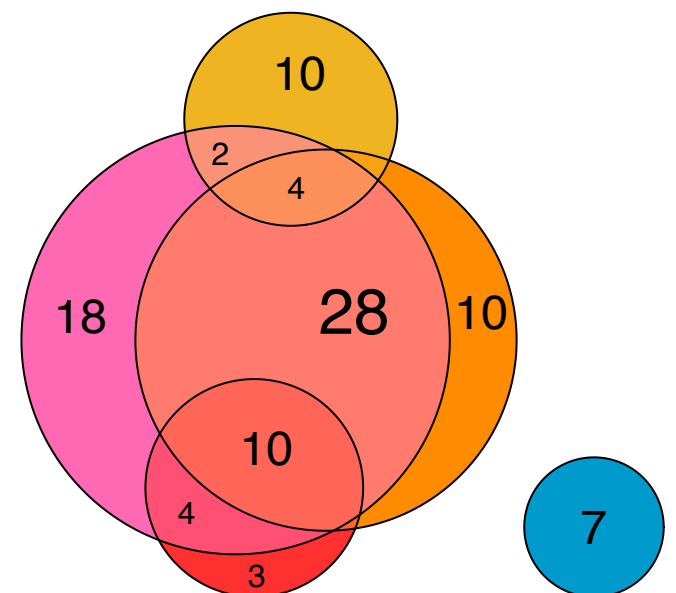
**A****BRCA****B****PAAD****C****PAAD (M)****D****PAAD (F)**



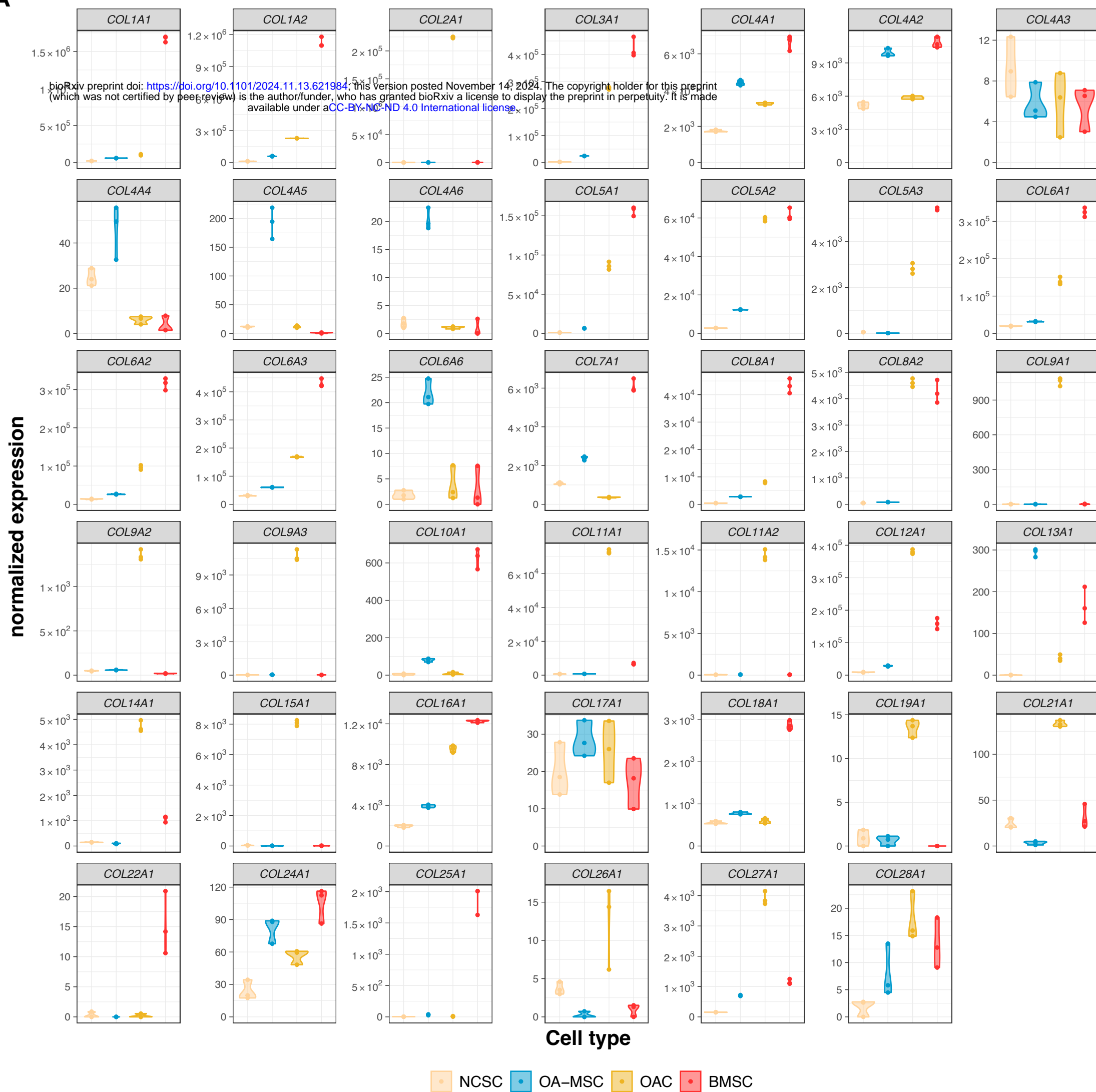
## E Gene set overlaps



## F Gene set overlaps (EMT markers only)





**A****B**