

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Evaluating the transcriptional fidelity of cancer models

Da Peng^{1*}, Rachel Gleyzer^{2*}, Wen-Hsin Tai², Pavithra Kumar², Qin Bian², Bradley Issacs²,
Edroaldo Lumertz da Rocha³, Stephanie Cai¹, Kathleen DiNapoli^{4,5}, Patrick Cahan^{1,2,6}

¹Department of Biomedical Engineering, Johns Hopkins University School of Medicine,
Baltimore MD 21205 USA

²Institute for Cell Engineering, Johns Hopkins University School of Medicine,
Baltimore MD 21205 USA

³Department of Biochemistry and Molecular Pharmacology, Boston Children's Hospital and
Harvard Medical School, Boston MA 02115 USA

⁴Department of Cell Biology, Johns Hopkins University School of Medicine, Baltimore, MD
21205 USA

⁵Department of Electrical and Computer Engineering, Johns Hopkins University, Baltimore MD
21218 USA

⁶Department of Molecular Biology and Genetics, Johns Hopkins University School of Medicine,
Baltimore MD 21205 USA

* These authors made equal contributions.

Correspondence to: patrick.cahan@jhmi.edu

Article type: Analysis

Website: http://www.cahanlab.org/resources/cancerCellNet_web

Code: <https://github.com/pcahan1/cancerCellNet>

50 **ABSTRACT**

51
52 Cancer researchers use cell lines, patient derived xenografts, and genetically engineered mice
53 as models to investigate tumor biology and to identify therapies. The generalizability and power
54 of a model derives from the fidelity with which it represents the tumor type of investigation,
55 however, the extent to which this is true is often unclear. The preponderance of models and the
56 ability to readily generate new ones has created a demand for tools that can measure the extent
57 and ways in which cancer models resemble or diverge from native tumors. Here, we present a
58 computational tool, CancerCellNet, that measures the similarity of cancer models to 22 naturally
59 occurring tumor types and 36 subtypes, in a platform and species agnostic manner. We applied
60 this tool to 657 cancer cell lines, 415 patient derived xenografts, and 26 distinct genetically
61 engineered mouse models, documenting the most faithful models, identifying cancers
62 underserved by adequate models, and finding models with annotations that do not match their
63 classification. By comparing models across modalities, we find that genetically engineered mice
64 have higher transcriptional fidelity than patient derived xenografts and cell lines in four out of
65 five tumor types. We have made CancerCellNet available as freely downloadable software and
66 as a web application that can be applied to new cancer models.

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81 INTRODUCTION

82 Models are widely used to investigate cancer biology and to identify potential therapeutics.
83 Popular modeling modalities are cancer cell lines (CCLs)¹, genetically engineered mouse
84 models (GEMMs)², and patient derived xenografts (PDXs)³. These classes of models differ in the
85 types of questions that they are designed to address. CCLs are often used to address cell
86 intrinsic mechanistic questions⁴, GEMMs to chart progression of molecularly defined-disease⁵,
87 and PDXs to explore patient-specific response to therapy in a physiologically relevant context⁶.
88 Models also differ in the extent to which they represent specific aspects of a cancer type⁷.
89 Even with this intra- and inter-class model variation, all models should represent the tumor type
90 or subtype under investigation, and not another type of tumor, and not a non-cancerous tissue.
91 Therefore, cancer-models should be selected not only based on the specific biological question
92 but also based on the similarity of the model to the cancer type under investigation^{8,9}.

93 Various methods have been proposed to determine the similarity of cancer models to
94 their intended subjects. Domcke et al devised a 'suitability score' as a metric of the molecular
95 similarity of CCLs to high grade serous ovarian carcinoma based on a heuristic weighting of
96 copy number alterations, mutation status of several genes that distinguish ovarian cancer
97 subtypes, and hypermutation status¹⁰. Other studies have taken analogous approaches by
98 either focusing on transcriptomic or ensemble molecular profiles (e.g. transcriptomic and copy
99 number alterations) to quantify the similarity of cell lines to tumors¹¹⁻¹³. These studies were
100 tumor-type specific, focusing on CCLs that model, for example, hepatocellular carcinoma or
101 breast cancer. More recently, Yu et al compared the transcriptomes of CCLs to the Cancer
102 Genome Atlas (TCGA) by correlation analysis, resulting in a panel of CCLs recommended as
103 most representative of 22 tumor types¹⁴. While all of these studies have provided valuable
104 information, they leave two major challenges unmet. The first challenge is to determine the
105 fidelity of GEMMs and PDXs and whether there are stark differences between these classes of
106 models and CCLs. The other major unmet challenge is to enable the rapid assessment of new,

107 emerging cancer models. This challenge is especially relevant now as technical barriers to
108 generating models have been substantially lowered^{15,16}, and because each PDX can be
109 considered a distinct entity requiring individual validation¹⁷.

110 To address these challenges, we developed CancerCellNet (CCN), a computational tool
111 that uses transcriptomic data to quantitatively assess the similarity between cancer models and
112 22 naturally occurring tumor types and 36 subtypes in a platform- and species-agnostic manner.
113 Here, we describe CCN's performance, and the results of applying it to assess 657 cancer cell
114 lines, 415 patient derived xenografts, and 26 distinct genetically engineered mouse models.
115 This has allowed us to identify the most faithful models currently available, to document cancers
116 underserved by adequate models, and to find models with inaccurate tumor type annotation.
117 Moreover, because CCN is open-source and easy to use, it can be readily applied to newly
118 generated cancer models as a means to assess their fidelity.

119

120 **RESULTS**

121 **CancerCellNet classifies samples accurately across species and technologies**

122 Previously, we had developed a computational tool using the Random Forest
123 classification method to measure the similarity of engineered cell populations to their *in vivo*
124 counterparts based on transcriptional profiles^{18,19}. More recently, we elaborated on this
125 approach to allow for classification of single cell RNA-Seq data in a manner that allows for
126 cross-platform and cross-species analysis²⁰. Here, we used an analogous approach to
127 quantitatively compare cancer models to naturally occurring patient tumors (**Fig 1A**). In brief, we
128 used TCGA RNA-seq expression data from 22 solid tumor types to train a top-pair multi-class
129 Random forest classifier. We combined training data from Rectal Adenocarcinoma (READ) and
130 Colon Adenocarcinoma (COAD) into one COAD_READ category because READ and COAD
131 are considered to be virtually indistinguishable at a molecular level²¹. We included an 'Unknown'

132 category trained using randomly shuffled gene-pair profiles generated from the training data of
133 22 tumor types to identify query samples that are not reflective of any of the training data.

134 We assessed the performance of this approach by computing the area under the
135 precision recall (AUPR) curves derived by 50 iterations of cross validation (**Fig 1B, Supp Fig**
136 **1A**). In the cross validations, the mean AUPR exceeded 0.95 in most of the tumor types. In
137 addition to achieving high mean AUPRs on held-out TCGA data, we found that CCN also
138 achieved high AUPR (above 0.9) when we applied it to independent testing data from the
139 International Cancer Genome Consortium (ICGC) consisting of RNA-Seq data from 886 tumors
140 across 5 tumor types (**Supp Fig 1B**)²².

141 As one of the central aims of our study is to compare distinct cancer models, including
142 GEMMs, our method needed to be able to classify samples from mouse and human samples
143 equivalently. We used the Top-Pair transform²⁰ to achieve this and we tested the feasibility of
144 this approach by assessing the performance of a normal (i.e non-tumor) cell and tissue classifier
145 trained on human data as applied to mouse samples. Consistent with prior applications²³, we
146 found that the cross-species classifier performed well, achieving mean AUPR of 0.96 when
147 applied to mouse data (**Supp Fig 1C**).

148 To evaluate cancer models at a finer resolution, we also developed an approach to
149 perform tumor subtype classifications (**Supp Fig 1D**). We constructed 11 different cancer
150 subtype classifiers based on the availability of expression or histological subtype
151 information^{21,24-34}. We also included non-cancerous, normal tissues as categories for several
152 subtype classifiers when sufficient data was available: breast invasive carcinoma (BRCA),
153 COAD_READ, head and neck squamous cell carcinoma (HNSC), kidney renal clear cell
154 carcinoma (KIRC) and uterine corpus endometrial carcinoma (UCEC). The 11 subtype
155 classifiers all achieved high overall average AUPRs ranging from 0.75 to 0.99 (**Supp Fig 1E**).

156

157 **Fidelity of cancer cell lines**

158 Having validated the performance of CCN, we then used it to determine the fidelity of
159 CCLs. We mined RNA-Seq expression data of 657 different cell lines across 20 cancer types
160 from the Cancer Cell Line Encyclopedia (CCLE) and applied CCN to them, finding a wide
161 classification range for cell lines of each tumor type (**Fig 2A, Supp Tab 1**). To verify the
162 classification results, we applied CCN to CCLE expression profiles generated through
163 microarray expression profiling³⁵. To ensure that CCN would function on microarray data, we
164 first tested it by applying a CCN classifier created to test microarray data to 720 expression
165 profiles of 12 tumor types. The cross-platform CCN classifier performed well, based on the
166 comparison to study-provided annotation, achieving a mean AUPR of 0.944 (**Supp Fig 2A**).
167 Next, we applied this cross-platform classifier to microarray expression profiles of CCLE (**Supp**
168 **Fig 2B**). From the classification results of 571 cell lines that have both RNA-seq and microarray
169 expression profiles, we found a strong positive association between the classification scores
170 from RNA-seq and those from microarray (**Supp Fig 2C**). This comparison supports the notion
171 that the classification scores for each cell line are not artifacts of profiling methodology.
172 Moreover, this comparison shows that the scores are consistent between the times that the cell
173 lines were first assayed by microarray expression profiling in 2012 and by RNA-Seq in 2019.
174 We also observed high level of correlation between our analysis and the analysis done by Yu et
175 al¹⁴(**Supp Fig 2D**), further validating the robustness of the CCN results.

176 Next, we assessed the extent to which CCN classifications agreed with their nominal
177 tumor type of origin. We annotated cell lines based their CCN score profile as follows. 'Correct'
178 Cell lines with CCN score > 0.3 for the tumor type of origin were annotated 'correct'. Those with
179 CCN scores > 0.3 in the tumor type of origin and at least one other tumor type were annotated
180 as 'mixed'. Cell lines with CCN scores > 0.3 for tumor types other than that of the cell lines origin
181 were annotated as 'other', and those lines that did not received a CCN score > 0.3 for any tumor
182 type were annotated as 'none' (**Fig 2B**). We selected a decision threshold of 0.3 based on the
183 average of the threshold that produced the highest Macro F1 measure, harmonic mean of

184 precision and recall, across 50 cross validations. We found that majority of cell lines originally
185 annotated as Breast invasive carcinoma (BRCA), Cervical squamous cell carcinoma and
186 endocervical adenocarcinoma (CESC), Skin Cutaneous Melanoma (SKCM), Colorectal Cancer
187 (COAD_READ) and Sarcoma (SARC) fell into the 'correct' category (**Fig 2B**). On the other
188 hand, no Esophageal carcinoma (ESCA) or Brain Lower Grade Glioma (LGG) were classified as
189 'correct', demonstrating the need for more transcriptionally faithful cell lines that model those
190 general cancer types.

191 There are several possible explanations for cell lines not receiving a 'correct'
192 classification. One possibility is that the sample was incorrectly labeled in the study from which
193 we harvest the expression data. Consistent with this explanation, we found that colorectal
194 cancer line NCI-H684^{36,37}, a cell line labelled as liver hepatocellular carcinoma (LIHC) by CCLE,
195 was classified strongly as COAD_READ (**Supp Tab 1**). Another possibility to explain low CCN
196 score is that cell lines were derived from subtypes of tumors that are not well-represented in
197 TCGA. To explore this hypothesis, we first performed tumor subtype classification on the CCLE
198 lines from 11 tumor types for which we had trained subtype classifiers (**Supp Tab 2**). We
199 reasoned that if a cell was a good model for a rarer subtype, then it would receive a poor
200 general classification but a high classification for the subtype that it models well. Therefore, we
201 counted the number of lines that fit this pattern. We found that of the 223 lines with no general
202 classification, 54 (24%) were classified as a specific subtype, suggesting that derivation from
203 rare subtypes is not the major contributor to the poor overall fidelity of CCLs.

204 Another potential contributor to low scoring cell lines is intra-tumor stromal and immune
205 cell impurity in the training data. If impurity were a confounder of CCN scoring, then we would
206 expect a strong positive correlation between mean purity and mean CCN classification of CCLs
207 per general tumor type. However, the Pearson correlation coefficient between the mean purity
208 of general tumor type and mean CCN classification scores of CCLs in the corresponding

209 general tumor type was low (0.059), suggesting that tumor purity is not a major contributor to
210 the low CCN scores across CCLE (**Supp Fig 2E**).

211 To more directly assess the impact of intra-tumor heterogeneity in the training data on
212 evaluating cell lines, we constructed a classifier using cell types found in human melanoma and
213 glioblastoma scRNA-seq data^{38, 39}. Previously, we have demonstrated the feasibility of using our
214 classification approach on scRNA-seq data²³. Our scRNA-seq classifier achieved a high
215 average AUPR (0.95) when applied to held-out data (**Supp Fig 3A-B**). Comparing the CCN
216 score from bulk RNA-seq general classifier and scRNA-seq classifier, we observed a high level
217 of correlation (Pearson correlation of 0.89) between the SKCM CCN classification scores and
218 scRNA-seq SKCM malignant CCN classification scores for SKCM cell lines (**Fig 2C, Supp Fig**
219 **3C**). Among the 37 SKCM cell lines that were classified as SKCM in general classification, 36
220 SKCM cell lines were also classified as SKCM malignant cells in scRNA-seq classifier.
221 Interestingly, we also observed a high correlation between the SARC CCN classification score
222 and scRNA-seq cancer associated fibroblast (CAF) CCN classification scores (Pearson
223 correlation of 0.89). Six of the 10 SKCM cell lines that had been classified as SARC by CCN
224 were classified as CAF by the scRNAseq classifier (**Fig 2D, Supp Fig 3C**), which suggests the
225 possibility that these cell lines were derived from CAF or other mesenchymal populations, or
226 that they have acquired a mesenchymal character through their derivation. The high level of
227 agreement between scRNA-seq and bulk RNA-seq classification results shows that
228 heterogeneity in the training data of general CCN classifier has little impact in the classification
229 of SKCM cell lines.

230 In contrast, we observed a weaker correlation between GBM CCN classification scores
231 and scRNA-seq GBM neoplastic CCN classification scores (Pearson correlation of 0.58) for
232 GBM cell lines (**Fig 2E, Supp Fig 3D**). Of the 32 GBM lines that were not classified as GBM
233 with CCN, 26 were classified as GBM neoplastic cells with the scRNAseq classifier. Among the
234 27 GBM lines that were classified as SARC with CCN, 15 cell lines were classified as CAF (**Fig**

235 **2F**), and 10 of 15 lines were classified as both GBM neoplastic and CAF in the scRNA-seq
236 classifier. Similar to the situation with SKCM lines that classify as CAF, this result is consistent
237 with the possibility that some GBM lines classified as SARC by CCN could be derived from
238 mesenchymal subtypes exhibiting both strong mesenchymal signatures and glioblastoma
239 signatures or that they have acquired a mesenchymal character through their derivation⁴⁰. The
240 lower level of agreement between scRNA-seq and bulk RNA-seq classification results for GBM
241 models suggests that the heterogeneity of glioblastomas⁴¹ can impact the classification of GBM
242 cell lines, and that the use of scRNA-seq classifier can resolve this deficiency.

243 Next, we explored the subtype classification of CCLs from three general tumor types in
244 more depth, focusing first on UCEC. The histologically defined subtypes of UCEC, endometrioid
245 and serous, differ in prevalence, molecular properties, prognosis, and treatment. For instance,
246 the endometrioid subtype, which accounts for approximately 80% of uterine cancers, retains
247 estrogen receptor and progesterone receptor status and is responsive towards progestin
248 therapy^{42,43}. Serous, a more aggressive subtype, is characterized by the loss of estrogen and
249 progesterone receptor and is not responsive to progestin therapy^{42,43}. CCN classified the
250 majority of the UCEC cell lines as serous except for JHUEM-1 which is classified as
251 endometrioid (**Fig 3A**). The preponderance CCLE lines of serous versus endometrioid character
252 may be due to properties of serous cancer cells that promote their *in vitro* propagation, such as
253 upregulation in cell adhesion⁴⁴. Some of our subtype classification results are consistent with
254 prior observations. For example, HEC-1A, HEC-1B, and KLE were previously characterized as
255 type II endometrial⁴⁵. On the other hand, our subtype classification results contradict prior
256 observations in at least one case. For instance, the Ishikawa cell line was derived from type I
257 endometrial cancer (endometrioid histological subtype)^{45,46}, however CCN classified a derivative
258 of this line, Ishikawa 02 ER-, as serous. The high serous CCN score could result from a shift in
259 phenotype of the line concomitant with its loss of estrogen receptor (ER) as this is a

260 distinguishing feature of type II endometrial cancer (serous histological subtype)⁴². Taken
261 together, these results indicate a need for more endometroid-like CCLs.

262 Next, we examined the subtype classification of Lung Squamous Cell Carcinoma
263 (LUSC) cell lines (**Fig 3C**). We found that of the 22 lines unclassified or misclassified in the
264 general classifier, 6 (27%) were classified as a subtype. Among the LUSC cell lines that were
265 classified as unknown in the general classifier and classified with a subtype, several cell lines
266 had general classification scores modestly below the threshold. All the LUSC lines with at least
267 one subtype classification had an underlying primitive subtype classification. This is consistent
268 either with the ease of deriving lines from tumors with a primitive character, or with a process by
269 which cell line derivation promotes similarity to more the primitive subtype, which is marked by
270 increased cellular proliferation²⁶. Some of our results are consistent with prior reports that have
271 investigated the resemblance of some lines to LUSC subtypes. For example, HCC-95, classified
272 as classical and primitive subtype, has previously been characterized as classical^{26,47}. Further,
273 LUDLU-1, classified as a primitive subtype, has classification signal in classical subtype which
274 was previously characterized as resembling classical⁴⁷. Likewise, although EPLC-272H was
275 also classified as primitive subtype, it has relatively high CCN score in the basal subtype, which
276 corresponds to its previous characterization as basal⁴⁷. Lung Adenocarcinoma (LUAD) cell lines
277 had classification results similar to LUSC: most lines did not classify as LUAD in the general
278 classifier (66 of 76) (**Fig 3B**). The cell lines that were classified as a subtype were either
279 classified as proximal inflammation, proximal proliferation or a mix of the two. RERF-LC-Ad1
280 had the highest general classification score and the highest proximal inflammation subtype
281 classification score. Taken together, these subtype classification results have revealed an
282 absence of cell lines models for basal, classical, and secretory LUSC, and for the Terminal
283 respiratory unit (TRU) LUAD subtype.

284 Finally, we sought to measure the extent to which cell line transcriptional fidelity related
285 to model use. We used the number of papers in which a model was mentioned, normalized by

286 the number of years since the cell line was documented, as a rough approximation of model
287 usage. To explore this metric, we plotted the normalized citation count versus general
288 classification score, labeling the highest cited and highest classified cell lines from each general
289 tumor type (**Fig 3D**). For most of the general tumor types, the highest cited cell line is not the
290 highest classified cell line except for Hep G2 and ML-1, representing liver hepatocellular
291 carcinoma (LIHC) and thyroid carcinoma (THCA), respectively. On the other hand, the general
292 scores of the highest cited cell lines representing BLCA, SKCM, BRCA, PRAD and
293 COAD_READ fall below the classification threshold of 0.3. Notably, each of these tumor types
294 have other lines with scores exceeding 0.7, which should be considered as more faithful
295 transcriptional models when selecting lines for a study (**Supp Table 1 and**
296 http://www.cahanlab.org/resources/cancerCellNet_results/).

297

298 **Evaluation of patient derived xenografts**

299 Next, we sought to evaluate a more recent class of cancer models: PDX. To do so, we
300 subjected the RNA-Seq expression profiles of 415 PDX models from 13 different types of
301 cancer types generated previously¹⁷ to CCN. Similar to the results of CCLE, the PDXs exhibited
302 a wide range of classification scores (**Fig 4A, Supp Tab 3**). By categorizing the CCN scores of
303 PDX based on the proportion of samples associated with each tumor type that were correctly
304 classified, we found that SARC, SKCM, COAD_READ and BRCA have higher proportion of
305 correctly classified PDX than those of other cancer categories (**Fig 4B**). In contrast to CCLE, we
306 found a higher proportion of correctly classified PDX in Stomach adenocarcinoma (STAD) and
307 KIRC (**Fig 4B**). However, similar to CCLE, no ESCA PDXs were classified as such. This held
308 true when we performed subtype classification on PDX samples: none of the PDX in ESCA
309 were classified as any of the ESCA subtypes (**Supp Tab 4**). UCEC PDXs had both
310 endometrioid subtypes, serous subtypes, and mixed subtypes, which provided a broader
311 representation than in CCLE (**Fig 4C**). Many LUSC PDXs that were classified as a subtype

312 were also classified as Head and Neck squamous cell carcinoma (HNSC) (**Fig 4D**). This could
313 be due to the similarity in expression profiles of basal and classical subtypes of HNSC and
314 LUSC^{26,48}, which is consistent with the observation that these PDXs were also subtyped as
315 basal and classical. No LUSC PDXs were classified as the secretory subtype. While eight of the
316 LUAD PDX samples were classified as the unknown subtype class classification, the remaining
317 six classified as proximal proliferative or proximal inflammatory (**Fig 4E**). Finally, similar to the
318 CCLE, there were no TRU subtypes in the PDX cohort. In summary, we found that while
319 individual PDXs can reach extremely high transcriptional fidelity to both general tumor types and
320 subtypes, many PDXs were not classified as the general tumor type from which they originated.

321

322 **Evaluation of GEMMs**

323 Next, we used CCN to evaluate GEMMs of six general tumor types from nine studies for
324 which expression data was publicly available⁴⁹⁻⁵⁷. As was true for CCLs and PDXs, GEMMs
325 also had a wide range of CCN scores (**Fig 5A, Supp Tab 5**). We next categorized the CCN
326 scores based on the proportion of samples associated with each tumor type that were correctly
327 classified (**Fig 5B**). In contrast to CCLs and PDXs, the GEMM dataset included multiple
328 replicates per model, which allowed us to examine intra-GEMM variability. Both at the level of
329 CCN score and at the level of categorization, GEMMs were highly invariant. For example,
330 replicate of UCEC GEMMs driven by *Prg(cre/+)**Pten(lox/lox)* received almost identical general
331 and subtype classification profiles (**Supp Fig 4, Supp Tab 6**). GEMMs sharing genotypes
332 across studies such as LUAD GEMMs driven by *Kras* mutation and loss of *p53*^{49,55,57} received
333 similar general and subtype classification scores (**Fig. 5A,B,D**). Even GEMMs with mixed
334 classifications received consistent CCN scores. For example, LGG GEMMs, generated by *Nf1*
335 mutations expressed in different neural progenitors in combination with *Pten* deletion⁵⁶,
336 consistently received mixed classification as both LGG and GBM (**Fig 5A**).

337 To explore the extent to which driver genotype impacted subtype classification, we
338 examined two general tumor types in which there were GEMMs with different tumor drivers:
339 LUSC and LUAD. The LUSC GEMMs were generated using loss of Lkb1 and either
340 overexpression of Sox2 (via two distinct mechanisms) or loss of Pten⁵⁵. Although most of the
341 lenti-Sox2-Cre-infected;Lkb1^{fl/fl} and Rosa26LSL-Sox2-IRES-GFP;Lkb1^{fl/fl} samples were
342 classified as unknown, their general LUSC CCN scores were only modestly lower than the
343 decision threshold and consistently throughout (**Fig 5C**). Those two models also classified
344 mostly as secretory subtype of LUSC. The consistency is not surprising given both models
345 overexpress Sox2 and lose Lkb1. Most of the Lkb1^{fl/fl};Pten^{fl/fl} GEMMs received unknown general
346 classifications with general LUSC CCN scores substantially lower than those of lenti-Sox2-Cre-
347 infected;Lkb1^{fl/fl} samples and Rosa26LSL-Sox2-IRES-GFP;Lkb1^{fl/fl} samples. Moreover, our
348 subtype classification indicated that this GEMM was mostly classified as unknown, in contrast to
349 prior reports suggesting that it is most similar to a basal subtype⁵⁸. The lenti-Sox2-Cre-
350 infected;Lkb1^{fl/fl} samples received high secretory subtype scores, whereas the Rosa26LSL-
351 Sox2-IRES-GFP;Lkb1^{fl/fl} samples were classified as a more balanced mix of secretory and
352 primitive subtypes. None of the three LUSC GEMMs have strong classical or basal sub-type
353 CCN scores.

354 All of the LUAD GEMMs, which were generated using various combinations of activating
355 Kras mutation, loss of Trp53, and loss of Smarca4L^{49,55,57}, were correctly classified (**Fig 5D**).
356 There were no substantial differences in general, or subtype classification across driver
357 genotypes. Notably, the subtypes tended to have CCN scores in mixture of proximal
358 proliferation, proximal inflammation and TRU. Taken together, this analysis suggests that there
359 is a degree of similarity, and perhaps plasticity between the primitive and secretory (but not
360 basal or classical) subtypes of LUSC. On the other hand, while the LUAD GEMMs classify
361 strongly as LUAD, do not have strong particular subtype classification -- a result that does not
362 vary by genotype.

363

364 **Comparison of CCLs, PDXs, and GEMMs**

365 Finally, we sought to estimate the comparative transcriptional fidelity of the three cancer
366 models modalities, limiting our comparison to those five general tumor types for which there
367 were at least two examples per modality: UCEC, Pancreatic adenocarcinoma (PAAD), LUSC,
368 LUAD, and LIHC. We compared the general CCN scores of each model on a per tumor type
369 basis (**Fig 6A**). In the case of GEMMs, we used the mean classification score of all samples
370 with shared genotypes. We found that GEMMs had the highest median general classification
371 scores in four out of the five tumor types. However, some PDXs achieved the highest
372 classification scores. In UCEC, LUAD and LIHC, the maximum classification score of PDXs
373 exceeded 0.75 and were thus comparable to the majority of scores on held out TCGA data,
374 highlighting the potential for PDXs to mirror the transcriptional state of natural tumors (**Fig 6A**).
375 Because the CCN score is based on a moderate number of gene pairs (i.e. 1647) relative to the
376 total number of protein-coding genes, it is possible that a cancer model with a high CCN score
377 might not have a high global similarity to a naturally occurring tumor. Therefore, we also
378 calculated the GRN status, a metric of the extent to which tumor-type specific gene regulatory
379 network is established¹⁸, for all models (**Supp Fig 5**). We observed high level of correlation
380 between the two similarity metrics, which suggests that although CCN classifies on a selected
381 set of genes, its scores are highly correlated with global assessment of transcriptional similarity.

382 We also sought to compare model modalities in terms of the diversity of subtypes that
383 they represent (**Supp Fig 6**). As a reference, we also included in this analysis the overall
384 subtype incidence, as approximated by incidence in TCGA. In models of UCEC, there is a
385 notable difference in endometroid incidence, and the proportion of models classified as
386 endometroid, with only PDX having any representatives (**Fig 6B**). The vast majority of CCLs
387 and all of the GEMM models of PAAD have an unknown subtype classification. However, the
388 majority of PDXs are subtyped as either a mixture of basal and classical, or basal and classical

389 alone. LUSC have proximal inflammation and proximal proliferation subtypes modelled by CCLs
390 and PDX, and TRU subtype modelled by GEMMs exclusively (**Fig 6B**). Likewise, LUAD have
391 basal, classical and primitive subtypes modelled by CCLs and PDXs, and secretory subtype
392 modelled by GEMMs exclusively (**Fig 6B**). Taken together, these results demonstrate the need
393 to carefully select different model systems to more suitably model certain cancer subtypes.

394

395 **DISCUSSION**

396 A major goal in the field of cancer biology is to develop models that mimic naturally occurring
397 tumors with enough fidelity to enable therapeutic discoveries. However, methods to measure
398 the extent to which cancer models resemble or diverge from native tumors are lacking. This is
399 especially problematic now because there are many existing models from which to choose, and
400 it has become easier to generate new models. Here, we present CancerCellNet (CCN), a
401 computational tool that measures the similarity of cancer models to 22 naturally occurring tumor
402 types and 36 subtypes. Because CCN is platform- and species-agnostic, it can be applied
403 across many model modalities, including CCLs, PDXs, and GEMMs, and thus it represents a
404 consistent platform to compare models across modalities. Here, we applied CCN to 657 cancer
405 cell lines, 415 patient derived xenografts, and 26 distinct genetically engineered mouse models.
406 Several lessons emerged from our computational analyses that have implications for the field of
407 cancer biology.

408 First, CancerCellNet indicates that GEMMs are transcriptionally the most faithful models
409 of four out of five general tumor types for which data from all modalities was available. This is
410 consistent with the fact that GEMMs are typically derived by recapitulating well-defined driver
411 mutations of natural tumors, and thus this observation corroborates the importance of genetics
412 in the etiology of cancer⁵⁹. Moreover, in contrast to most PDXs, GEMMs are typically generated
413 in immune replete hosts. Therefore, the higher fidelity of GEMMs may also be a result of the
414 influence of a native immune system on GEMM tumors⁶⁰. Second, PDXs and CCLs have lower

415 scores that are comparable to each other. This is consistent with the observation that PDXs can
416 undergo selective pressures in the host that distort the progression of genomic alterations away
417 from what is observed in natural tumors⁶¹. Furthermore, the observation that a few PDXs have
418 very high classification scores, approaching a level that is indistinguishable from held out TCGA
419 data, suggests that under certain conditions, PDX can almost perfectly mimic natural tumors
420 transcriptionally. It is unclear what are these conditions; it may be that these few PDXs were
421 profiled prior to the acquisition of non-typical genomic alterations. Third, we have found that
422 none of the samples that we evaluated here are transcriptionally adequate models of ESCA,
423 and therefore this tumor type requires further attention to derive new models. Fourth, we found
424 that in several tumor types, GEMMs tend to reflect mixtures of subtypes rather than conforming
425 strongly to single subtypes. The reasons for this are not clear but it is possible that in the cases
426 that we examined the histologically defined subtypes have a degree of plasticity that is
427 exacerbated in the murine host environment. We have made the results of our analyses
428 available online so that researchers can easily explore the performance of selected models or
429 identify the best models for any of the 22 general tumor types and the 36 subtypes presented
430 here.

431 Currently, there are several limitations to our CCN tool, and caveats to our analyses
432 which indicate areas for future work and improvement. First, CCN is based on transcriptomic
433 data but other molecular readouts of tumor state, such as profiles of the proteome⁶²,
434 epigenome⁶³, non-coding RNA-ome⁶³, and genome⁵⁹ would be equally, if not more important, to
435 mimic in a model system. Therefore, it is possible that some models reflect tumor behavior well,
436 and because this behavior is not well predicted by transcriptome alone, these models have
437 lower CCN scores. To both measure the extent that such situations exist, and to correct for
438 them, we plan in the future to incorporate other omic data into CCN so as to make more
439 accurate and integrated model evaluation possible. A second limitation is that in the cross-
440 species analysis, CCN implicitly assumes that homologs are functionally equivalent. The extent

441 to which they are not functionally equivalent determines how confounded the CCN results will
442 be. This possibility seems to be of limited consequence based on the high performance of the
443 normal tissue cross-species classifier and based on the fact that GEMMs have the highest
444 median CCN scores. Finally, the TCGA training data is made up of RNA-Seq from bulk tumor
445 samples, which necessarily includes non-tumor cells, whereas the CCLs are by definition cell
446 lines of tumor origin. Therefore, CCLs theoretically could have artificially low CCN scores due to
447 the presence of non-tumor cells in the training data. This problem appears to be limited as we
448 found no correlation between tumor purity and CCN score in the CCLE samples. However, this
449 problem is related to the question of intra-tumor heterogeneity. We demonstrated the feasibility
450 of using CCN and single cell RNA-seq data to refine the evaluation of cancer cell lines
451 contingent upon availability of scRNA-seq training data. As more sufficient training single cell
452 RNA-Seq data accrues, CCN would be able to not only evaluate models on a per cell type
453 basis, but also based on cellular composition.

454 To ensure that CCN is widely available we have developed a free web application,
455 which performs CCN analysis on user-uploaded data and allows for direct comparison their to
456 the cancer models evaluated here. We have also made the CCN code freely available under an
457 Open Source license and as an easily installed R package, and we are actively supporting its
458 further development. The documentation describes how to analyze model(s) and compare the
459 results to the panel of models that we evaluated here, thereby allowing researchers to
460 immediately compare their models to the broader field in a comprehensive and standard
461 fashion.

462

463 **Online Methods**

464 **Training General CancerCellNet Classifier**

465 To generate training data sets, we downloaded 8991 patient tumor RNA-seq expression
466 count matrix and their corresponding sample table across 22 different tumor types from TCGA

467 using TCGAWorkflowData, TCGAAbiolinks⁶⁴ and SummarizedExperiment⁶⁵ packages. We used
468 all the patient tumor samples for training the general CCN classifier. Later, we found the
469 intersecting genes between TCGA dataset and all the query samples (CCLs, PDXs, GEMMs),
470 and used them as features for the feature engineering and selection process of building the
471 classifier. To train the top pair Random Forest classifier, we used a method similar to our
472 previous method²³. CCN first normalized the training counts matrix by down-sampling the
473 counts to 500,000 counts per cell. To significantly reduce the time and resource of generating
474 gene pairs for all possible genes, CCN then selected 30 up-regulated genes, 30 down-regulated
475 genes and 30 least differentially expressed genes for each of the 23 cancer categories using
476 template matching⁶⁶ as the genes to generate top scoring gene pairs. In short, for each tumor
477 type, CCN defined a template vector that labelled the training tumor samples in cancer type of
478 interest as 1 and all other tumor samples as 0 CCN then calculated the Pearson correlation
479 coefficient between template vector and gene expressions for all genes. The genes with strong
480 match to template as either upregulated or downregulated had large absolute Pearson
481 correlation coefficient. CCN chose the upregulated, downregulated and least differentially
482 expressed genes based on the magnitude of Pearson correlation coefficient.

483 After CCN selected the genes for each cancer type, CCN generated gene pairs among
484 those genes. Gene pair transformation was a method inspired by the top-scoring pair classifier⁶⁷
485 to allow compatibility of classifier with query expression profiles that were collected through
486 different platforms (e.g. microarray query data applied to RNA-seq training data). In brief, the
487 gene pair transformation compares 2 genes within an expression sample and encodes the
488 “gene1_gene2” gene-pair as 1 if the first gene has higher expression than the second gene.
489 Otherwise, gene pair transformation would encode the gene-pair as 0. Using all the gene pair
490 combinations generated through the gene sets across all cancer types, CCN then selected top
491 75 discriminative gene pairs for each category using template matching (with large absolute
492 Pearson correlation coefficient) described above.

493 After the top discriminative gene pairs were selected for each cancer categories, CCN
494 concatenate all the gene pairs into a vector and gene pair transformed the training samples into
495 a binary matrix with all the discriminative gene pairs as row names and all the training samples
496 as column names. Using the binary gene pair matrix, CCN randomly shuffled the binary values
497 across rows then across column generating random profiles that should not resemble training
498 data from any of the cancer categories. CCN then sampled 70 random profiles, annotated them
499 as “Unknown” and appended them to the training gene pair binary matrix as training data for the
500 “Unknown” category.

501 Using gene pair binary training matrix, CCN constructed a multi-class Random Forest
502 classifier of 2000 trees and used stratified sampling of 60 sample size to ensure balance of
503 training data in constructing the decision trees. The specific parameters for the final CCN
504 classifier using the function “broadClass_train” in the package cancerCellNet are in **Supp Tab**
505 **7**. The gene-pairs are in **Supp Tab 8**.

506

507 **Validating General CancerCellNet Classifier**

508 2/3 of patient tumor data from each cancer type were randomly sampled as training data
509 to construct a CCN classifier. After the classifier was built, 35 held-out samples from each
510 cancer categories were sampled and 40 “Unknown” profiles were generated for validation. CCN
511 gene pair transformed the held-out data for assessment based on the top gene-pairs selected to
512 construct the classifier. The process of randomly sample training set from 2/3 of all patient
513 tumor data, train classifier and validate using validation set was repeated 50 times to have a
514 more comprehensive assessment of the classifier. We used precision-recall curve and area
515 under the precision-recall curve (AUPR) as our metric of assessing the classifiers.

516

517 **Classifying Query Data into General Cancer Categories**

518 We downloaded the RNA-seq cancer cell lines expression profiles and sample table
519 from (<https://portals.broadinstitute.org/ccle/data>), and microarray cancer cell lines expression
520 profiles and sample table from Barretina et al ³⁵. We received PDX expression estimates and
521 sample annotations from the authors of Gao et al ¹⁷. We gathered GEMM expression profiles
522 from 9 different studies⁴⁹⁻⁵⁷. To use CCN classifier on GEMM data, the mouse genes from
523 GEMM expression profiles were converted into their human homologs. The query samples were
524 gene pair transformed using gene pairs selected from the training step, and then inputted into
525 CCN classifier for classification. Each query classification profile was labelled as one of the four
526 classification categories: “correct”, “mixed”, “none” and “other” based on classification profiles. If
527 a sample has a CCN score higher than the decision threshold (0.3) in the labelled cancer
528 category, we assign that as “correct”. If a sample has CCN score higher than the decision
529 threshold in labelled cancer category and in other cancer categories, we assign that as “mixed”.
530 If a sample has no CCN score higher than the decision threshold in any cancer category or has
531 the highest CCN score in ‘Unknown’ category, then we assign it as “none”. If a sample has CCN
532 score higher than the decision threshold in a cancer category or categories not including the
533 labelled cancer category, we assign it as “other”. We analyzed and visualized the results using
534 R and R packages pheatmap⁶⁸ and ggplot2⁶⁹.

535

536 **Cross-Species Assessment**

537 To assess the performance of cross-species classification, we downloaded 1003
538 labelled human tissue/cell type and 1993 labelled mouse tissue/cell type RNA-seq expression
539 profiles from Github (<https://github.com/pcahan1/CellNet>). We first converted the mouse genes
540 into human homologous genes. Then we found the intersecting genes between mouse
541 tissue/cell expression profiles and human tissue/cell expression profiles. Limiting the input of
542 RNA-seq profiles to the intersecting genes, we trained a CCN classifier with all the human
543 tissue/cell expression profiles. The parameters used for the function “broadClass_train” in the

544 package cancerCellNet are in **Supp Tab 7**. After the classifier was trained, we randomly
545 sampled 75 samples from each tissue category in mouse tissue/cell data and applied the
546 classifier on those samples to assess performance.

547

548 **Cross-Technology Assessment**

549 To assess the performance of CCN in applications to microarray data, we gathered
550 6,219 patient tumor microarray profiles across 12 different cancer types from more than 100
551 different projects (**Supp Tab 9**). We found the intersecting genes between the microarray
552 profiles and TCGA patient RNA-seq profiles. Limiting the input of RNA-seq profiles to the
553 intersecting genes, we created a CCN classifier with all the TCGA patient profiles using
554 parameters for the function “broadClass_train” listed in **Supp Tab 7**. After the microarray
555 specific classifier was trained, we randomly sampled 60 microarray patient samples from each
556 cancer category, and applied CCN classifier on them as assessment of the cross-technology
557 performance in **Supp Fig 3A**. The same CCN classifier was used to assess microarray CCL
558 samples **Supp Fig 3B**.

559

560 **Training and validating scRNA-seq Classifier**

561 We extracted labelled human melanoma and glioblastoma scRNA-seq expression
562 profiles^{38,39}, and compiled the two datasets excluding 3 cell types T.CD4, T.CD8 and Myeloid
563 due to low number of cells for training. 60 cells from each of the 11 cell types were sampled for
564 training a scRNA-seq classifier. The parameters for training a general scRNA-seq classifier
565 using the function “broadClass_train” are in **Supp Tab 7**. 25 cells from each of the 11 cell types
566 from the held-out data were selected to assess the single cell classifier. Using the PR curve and
567 maximizing Macro F1 measure, we selected the decision threshold of 0.255. We then applied
568 the scRNA-seq classifier on SKCM CCLs and GBM CCLs.

569

570 **Training Subtype CancerCellNet**

571 We found 11 cancer types (BRCA, COAD, ESCA, HNSC, KIRC, LGG, PAAD, UCEC,
572 STAD, LUAD, LUSC) which have meaningful subtypes based on either histology or molecular
573 profile and have sufficient samples to train a subtype classifier with high AUPR. We also
574 included normal tissues samples from BRCA, COAD, HNSC, KIRC, UCEC to create a normal
575 tissue category in the construction of their subtype classifiers. Training samples were either
576 labelled as a cancer subtype for the cancer of interest or as “Unknown” if they belong to other
577 cancer types. Similar to general classifier training, CCN performed gene pair transformation and
578 selected the most discriminate gene pairs for each cancer subtype. In addition to the gene pairs
579 selected to discriminate cancer subtypes, CCN also performed general classification of all
580 training data and appended the classification profiles of training data with gene pair binary
581 matrix as additional features. The reason behind using general classification profile as additional
582 features is that many general cancer types may share similar subtypes, and general
583 classification profile could be important features to discriminate the general cancer type of
584 interest from other cancer types before performing finer subtype classification. The specific
585 parameters used to train individual subtype classifiers using “subClass_train” function of
586 CancerCellNet package can be found in **Supp Tab 7** and the gene pairs are in **Supp Tab 8**.

587

588 **Validating Subtype CancerCellNet**

589 Similar to validating general class classifier, we randomly sampled 2/3 of all samples in
590 each cancer subtype as training data and sampled an equal amount across subtypes in the
591 held-out data for assessing subtype classifiers. We repeated the process 20 times for more
592 comprehensive assessment of subtype classifiers.

593

594 **Classifying Query Data into Subtypes**

595 We assigned subtype to query sample if the query sample has CCN score higher than
596 the decision threshold. If a query sample has CCN scores higher than decision threshold, which
597 was chosen through maximizing Macro F1 measure. The table of decision threshold for subtype
598 classifiers are in **Supp Tab 10**. If a query sample with no CCN score higher than decision
599 threshold in any subtype or has the highest CCN score in 'Unknown' category, then we
600 assigned that sample as 'Unknown'. Analysis and visualizations were done in R and
601 ComplexHeatmap package⁷⁰.

602

603 **Tumor Purity Analysis**

604 We used the R package ESTIMATE⁷¹ to calculate the ESTIMATE scores from TCGA
605 tumor expression profiles that we used as training data for CCN classifier. To calculate tumor
606 purity we used the equation described in YoshiHara et al., 2013⁷¹:

$$\text{Tumour purity} = \cos(0.6049872018 + 0.0001467884 \times \text{ESTIMATE score})$$

607

608 **Extracting Citation Counts**

609 We used the R package RISmed⁷² to extract the number of citations for each cell line
610 through query search of "*cell line name*[Text Word] AND cancer[Text Word]" on PubMed. The
611 citation counts were normalized by dividing the citation counts with the number of years since
612 first documented.

$$\text{Normalized citation counts} = \frac{\text{citation counts}}{\# \text{ years since first documented}}$$

613

614 **GRN construction and GRN Status**

615 GRN construction was extended from our previous method¹⁸. 80 samples per cancer
616 type were randomly sampled and normalized through down sampling as training data for the
617 CLR GRN construction algorithm. Cancer type specific GRNs were identified by determining the

618 differentially expressed genes per each cancer type and extracting the subnetwork using those
619 genes.

620 To extend the original GRN status algorithm¹⁸ across different platforms and species, we
621 devised a rank-based GRN status algorithm. Like the original GRN status, rank based GRN
622 status is a metric of assessing the similarity of cancer type specific GRN between training data
623 in the cancer type of interest and query samples. Hence, high GRN status represents high level
624 of establishment or similarity of the cancer specific GRN in the query sample compared to those
625 of the training data. The expression profiles of training data and query data were transformed
626 into rank expression profiles by replacing the expression values with the rank of the expression
627 values within a sample (highest expressed gene would have the highest rank and lowest
628 expressed genes would have a rank of 1). Cancer type specific mean and standard deviation of
629 every gene's rank expression were calculated using training data. The modified Z-score values
630 for genes within cancer type specific GRN were calculated for query sample's rank expression
631 profiles to quantify how dissimilar the expression values of genes in query sample's cancer type
632 specific GRN compared to those of the reference training data:

$$633 \quad Z - score(gene\ i)_{mod} =$$
$$634 \quad \begin{cases} 0, & \text{if } Z - score \text{ is positive and the gene is found to be upregulated} \\ 0, & \text{if } Z - score \text{ is negative and the gene is found to be downregulated} \\ abs(Z - score), & \text{otherwise} \end{cases}$$

635 If a gene in the cancer type specific GRN is found to be upregulated in the specific
636 cancer type relative to other cancer types, then we would consider query sample's gene to be
637 similar if the ranking of the query sample's gene is equal to or greater than the mean ranking of
638 the gene in training sample. As a result of similarity, we assign that gene of a Z-score of 0. The
639 same principle applies to cases where the gene is downregulated in cancer specific subnetwork.

640 GRN status for query sample is calculated as the weighted mean of the (1000 –
641 $Zscore(gene\ i)_{mod}$) across genes in cancer type specific GRN. 1000 is an arbitrary large

642 number, and larger dissimilarity between query's cancer type specific GRN indicate high Z-
643 scores for the GRN genes and low GRN status.

$$RGS = \sum_{i=1}^n (1000 - Zscore(gene\ i)_{mod}) weight_{gene\ i}$$
$$GRN\ Status = \frac{RGS}{\sum_{i=1}^n weight_{gene\ i}}$$

644 The weight of individual genes in the cancer specific network is determined by the importance of
645 the gene in the Random Forest classifier. We later normalize the GRN status in respect to the
646 GRN status of the cancer type of interest and the cancer type with the lowest mean GRN status.

$$Normalized\ GRN\ status = \frac{GRN\ status_{query} - avg(GRN\ status_{min\ cancer})}{avg(GRN\ status_{cancer\ type\ interest})}$$

647 Where “min cancer” represents the cancer type where its training data have the lowest mean
648 GRN status in the cancer type of interest, and $avg(GRN\ status_{min\ cancer})$ represents the
649 average GRN status of cancer type with the lowest average GRN status in the “min cancer”.
650 $avg(GRN\ status_{cancer\ type\ interest})$ represents average GRN status of the cancer type of interest
651 in the training data.

652

653 **Code availability**

654 CancerCellNet code and documentation is available at GitHub:

655 <https://github.com/pcahan1/cancerCellNet>

656

657 **FIGURE LEGENDS**

658 **Fig. 1** CancerCellNet (CCN) workflow and performance. **(A)** Schematic of CCN training (top)
659 and usage (bottom). CCN was designed to assess and compare the expression profiles of
660 cancer models such as CCLs, PDXs, and GEMMs with native patient tumors. First, CCN takes
661 patient tumor expression profiles of 23 different cancer types from TCGA to train a multi-class

662 Random Forest classifier and performs gene-pair transformation on tumor expression profiles.
663 Then CCN selects the most discriminative gene pairs for each cancer type as features. Lastly,
664 CCN trains a multi-class Random Forest classifier using gene-pair transformed training data
665 and feature gene pairs. To use trained classifier, CCN inputs the query samples (e.g.
666 expression profiles from CCLs, PDXs, GEMMs) and generates a classification profile for the
667 query samples. The column names of the classification heatmap represent sample annotation
668 and the row names of the classification heatmap represent different cancer types. Each grid is
669 colored from black to yellow representing the lowest classification score (e.g. 0) to highest
670 classification score (e.g. 1). **(B)** Mean and standard deviation of area under the precision recall
671 curve (AUPR) of classifiers based on 50 iterations of cross-validation: random sampling of
672 training data (2/3 of samples for each cancer category), training CCN classifiers using training
673 data and testing the classifiers on held-out data (1/3 of samples for each cancer category).

674

675 **Fig. 2** Evaluation of cancer cell lines. **(A)** General classification heatmap of CCLs extracted from
676 CCLE. Column annotations of the heatmap represent the labelled cancer category of the CCLs
677 given by CCLE and the row names of the heatmap represent different cancer categories. CCLs'
678 general classification profiles are categorized into 4 categories: correct (red), correct mixed
679 (pink), no classification (light green) and other classification (dark green) based on the decision
680 threshold of 0.3. **(B)** Bar plot represents the proportion of each classification category in CCLs
681 across cancer types ordered from the cancer types with the highest proportion of correct and
682 correct mixed CCLs to lowest proportion. **(C)** Comparison between SKCM general CCN scores
683 from bulk RNA-seq classifier and SKCM malignant CCN scores from scRNA-seq classifier for
684 SKCM CCLs. **(D)** Comparison between SARC general CCN scores from bulk RNA-seq
685 classifier and CAF CCN scores from scRNA-seq classifier for SKCM CCLs. **(E)** Comparison
686 between GBM general CCN scores from bulk RNA-seq classifier and GBM neoplastic CCN
687 scores from scRNA-seq classifier for GBM CCLs. **(F)** Comparison between SARC general CCN

688 scores and CAF CCN scores from scRNA-seq classifier for GBM CCLs. The green lines
689 indicate the decision threshold for scRNA-seq classifier and general classifier.

690

691 **Fig. 3** Subtype classification of CCLs. The heatmap visualizations represent subtype
692 classification of UCEC CCLs **(A)**, LUAD CCLs **(B)** and LUSC CCLs **(C)**. The row names
693 represent CCLs and column names represent cancer subtypes of UCEC, LUSC and LUAD. The
694 bar plots to the right of the subclass heatmaps represent the general classification scores in cell
695 lines' annotated cancer category, and the color strips to the right of the subclass heatmaps
696 represent subclass classification (left) and general classification (right). **(D)** Comparison of
697 normalized citation counts and general CCN classification scores of CCLs. Labelled cell lines
698 either have the highest CCN classification score in their labelled cancer category or highest
699 normalized citation count. Hep G2 and ML-1 have both the highest CCN classification score in
700 their labelled cancer category and the highest normalized citation count. Each citation count
701 was normalized by number of years since first documented.

702

703 **Fig. 4** Evaluation of patient derived xenografts. **(A)** General classification heatmap of PDXs.
704 Column annotations represent annotated cancer type of the PDXs, and row names represent
705 cancer categories. **(B)** Proportion of classification categories in PDXs across cancer types is
706 visualized in the bar plot and ordered from the cancer type with highest proportion of correct and
707 mixed correct classified PDXs to the lowest. Subtype classification heatmaps of UCEC PDXs
708 **(C)**, LUSC PDXs **(D)** and LUAD PDXs **(E)**.

709

710 **Fig. 5** Evaluation of genetically engineered mouse models. **(A)** General classification heatmap
711 of GEMMs. Column annotations represent annotated cancer type of the GEMMs, and row
712 names represent cancer categories. **(B)** Proportion of classification categories in GEMMs
713 across cancer types is visualized in the bar plot and ordered from the cancer type with highest

714 proportion of correct and mixed correct classified GEMMs to the lowest. Subtype classification
715 heatmap of LUSC GEMMs **(C)** and LUAD GEMMs **(D)**.

716

717 **Fig. 6** Comparison of CCLs, PDXs, and GEMMs. **(A)** Box-and-whiskers plot comparing general
718 CCN scores across CCLs, GEMMs, PDXs of five general tumor types. **(B)** Proportion of UCEC
719 (top-left), PAAD (top-right), LUAD (bottom-left) and LUSC (bottom-right) subtypes across cancer
720 model modalities and TCGA patient data. For GEMMs, all classification profiles of replicates
721 with the same genotype from the same study are averaged into one classification profile when
722 calculating the proportion.

723

724 **Supplementary Information**

725 **Supplementary Figure 1** Assessment of CCN general classifier and subtype classifier. **(A)**
726 Mean and range of CCN classifier's PR curves from 50 cross validations. **(B)** AUPR of CCN
727 classifier when applied to independent patient tumor data from ICGC. **(C)** AUPR of CCN human
728 tissue classifier when applied to mouse tissue data. **(D)** The schematic of training a subtype
729 classifier in CCN. CCN uses patient tumor expression profiles from cancer of interest as training
730 data. CCN performs gene-pair transformation and selects the most discriminative gene pairs
731 among the cancer subtypes from training data as features. CCN then applies the general
732 classification on training data and uses the general classification profile as features in addition
733 to gene pairs for training a Random Forest classifier. The weight of the general classification
734 profiles as features can be tuned to maximize AUPR. **(E)** The mean and standard deviation of
735 AUPR for 11 subtype classifiers based on 20 iterations of random sampling of training and held-
736 out data, training subtype classifier using training data, classification of held-out data, and
737 calculation of recall and precision.

738

739 **Supplementary Figure 2** Further validation of CCN and classification results. To validate the
740 cross-platform classification performance of CCN, a new classifier specifically trained to classify
741 microarray data was trained using RNA-seq data from TCGA as training data and intersecting
742 genes between RNA-seq data and microarray data. **(A)** AUPR of CCN classifier when applied to
743 primary microarray testing tumor data. **(B)** Classification heatmap of CCLs using microarray
744 expression data. **(C)** Pearson correlation between CCN scores of CCLs generated from
745 RNA-seq data and microarray data. **(D)** Comparison between CCLs' CCN scores and the
746 median correlation results from Yu et al. **(E)** Comparison of mean tumor purity of training data
747 and mean CCN score of CCLs for each cancer category.

748

749 **Supplementary Figure 3** Single-cell classification of SKCM and GBM cell lines. **(A)**
750 Classification heatmap of held-out scRNA-seq data. **(B)** AUPR of the scRNA-seq classifier
751 when applied to scRNA-seq held-out data. **(C)** Single-cell classification of SKCM CCLs. **(D)**
752 Single-cell classification of GBM CCLs.

753

754 **Supplementary Figure 4** Cancer subtype classification heatmap of UCEC GEMMs.

755

756 **Supplementary Figure 5** Correlation between cancer type subnetwork GRN status and general
757 CCN scores.

758

759

760 **Supplementary Figure 6** Proportion of cancer subtypes in different cancer models and TCGA
761 tumor data across 11 general cancer types.

762

763

764 **Supplementary Table 1** General classification profiles of CCLs.

765

766 **Supplementary Table 2** Subtype classification profiles of CCLs.

767

768 **Supplementary Table 3** General classification profiles of PDXs.

769

770 **Supplementary Table 4** Subtype classification profiles of PDXs.

771

772 **Supplementary Table 5** General classification profiles of GEMMs

773
774
775
776
777
778
779
780
781
782
783
784

Supplementary Table 6 Subtype classification profiles of GEMMs.

Supplementary Table 7 Specific parameters used for training of all classifiers.

Supplementary Table 8 Gene-pairs selected for final training of CCN general and subtype classifiers.

Supplementary Table 9 Accessions of tumor microarray data used in validation.

Supplementary Table 10 Decision thresholds for subtype classifiers.

785 REFERENCES

- 786 1. Sharma, S. V., Haber, D. A. & Settleman, J. Cell line-based platforms to evaluate
787 the therapeutic efficacy of candidate anticancer agents. *Nat. Rev. Cancer* **10**, 241–
788 253 (2010).
- 789 2. Kersten, K., de Visser, K. E., van Miltenburg, M. H. & Jonkers, J. Genetically
790 engineered mouse models in oncology research and cancer medicine. *EMBO Mol.*
791 *Med.* **9**, 137–153 (2017).
- 792 3. Hidalgo, M. *et al.* Patient-derived xenograft models: an emerging platform for
793 translational cancer research. *Cancer Discov.* **4**, 998–1013 (2014).
- 794 4. Klijn, C. *et al.* A comprehensive transcriptional portrait of human cancer cell lines.
795 *Nat. Biotechnol.* **33**, 306–312 (2015).
- 796 5. Koren, S. *et al.* PIK3CA(H1047R) induces multipotency and multi-lineage mammary
797 tumours. *Nature* **525**, 114–118 (2015).
- 798 6. DeRose, Y. S. *et al.* Tumor grafts derived from women with breast cancer
799 authentically reflect tumor pathology, growth, metastasis and disease outcomes.
800 *Nat. Med.* **17**, 1514–1520 (2011).
- 801 7. Sharpless, N. E. & Depinho, R. A. The mighty mouse: genetically engineered
802 mouse models in cancer drug development. *Nat. Rev. Drug Discov.* **5**, 741–754
803 (2006).
- 804 8. Mouradov, D. *et al.* Colorectal cancer cell lines are representative models of the
805 main molecular subtypes of primary cancer. *Cancer Res.* **74**, 3238–3247 (2014).
- 806 9. Stuckelberger, S. & Drapkin, R. Precious GEMMs: emergence of faithful models for
807 ovarian cancer research. *J. Pathol.* **245**, 129–131 (2018).
- 808 10. Domcke, S., Sinha, R., Levine, D. A., Sander, C. & Schultz, N. Evaluating cell lines
809 as tumour models by comparison of genomic profiles. *Nat. Commun.* **4**, 2126
810 (2013).
- 811 11. Jiang, G. *et al.* Comprehensive comparison of molecular portraits between cell lines
812 and tumors in breast cancer. *BMC Genomics* **17 Suppl 7**, 525 (2016).
- 813 12. Chen, B., Sirota, M., Fan-Minogue, H., Hadley, D. & Butte, A. J. Relating
814 hepatocellular carcinoma tumor samples and cell lines using gene expression data
815 in translational research. *BMC Med. Genomics* **8 Suppl 2**, S5 (2015).
- 816 13. Vincent, K. M., Findlay, S. D. & Postovit, L. M. Assessing breast cancer cell lines as
817 tumour models by comparison of mRNA expression profiles. *Breast Cancer Res.*
818 **17**, 114 (2015).

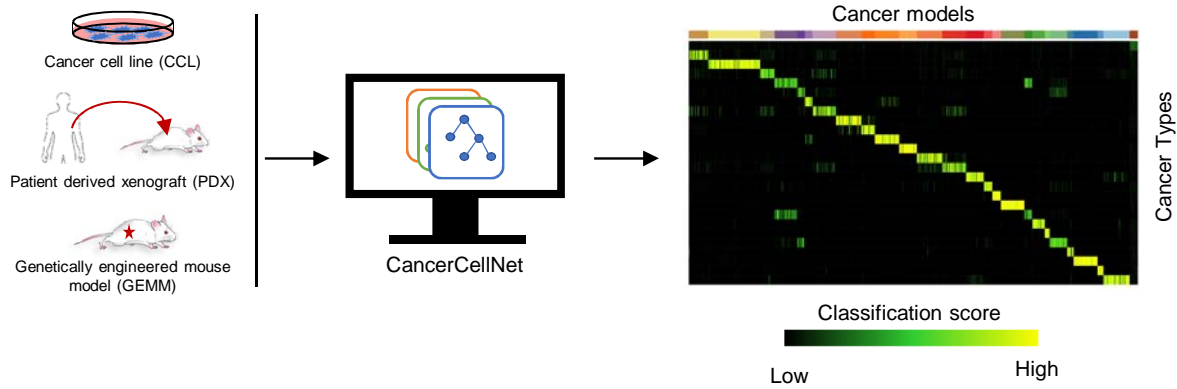
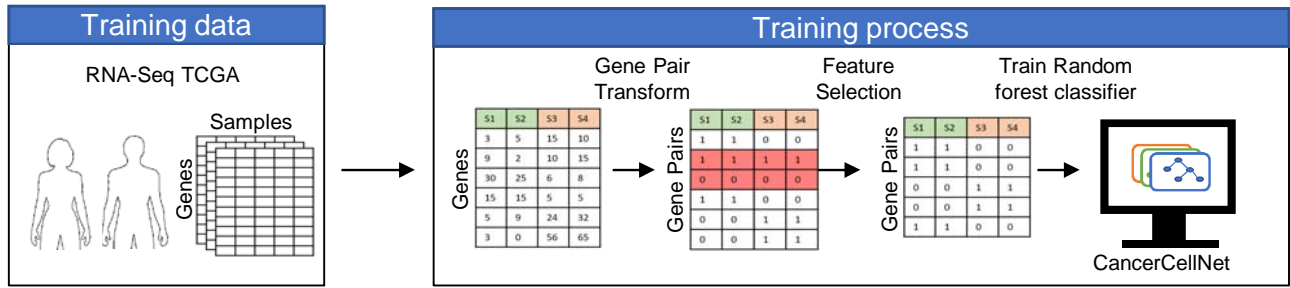
- 819 14. Yu, K. *et al.* Comprehensive transcriptomic analysis of cell lines as models of
820 primary tumors across 22 tumor types. *Nat. Commun.* **10**, 3574 (2019).
- 821 15. Guernet, A. & Grumolato, L. CRISPR/Cas9 editing of the genome for cancer
822 modeling. *Methods* **121-122**, 130–137 (2017).
- 823 16. Gargiulo, G. Next-Generation in vivo Modeling of Human Cancers. *Front. Oncol.* **8**,
824 429 (2018).
- 825 17. Gao, H. *et al.* High-throughput screening using patient-derived tumor xenografts to
826 predict clinical trial drug response. *Nat. Med.* **21**, 1318–1325 (2015).
- 827 18. Cahan, P. *et al.* CellNet: network biology applied to stem cell engineering. *Cell* **158**,
828 903–915 (2014).
- 829 19. Radley, A. H. *et al.* Assessment of engineered cells using CellNet and RNA-seq.
830 *Nat. Protoc.* **12**, 1089–1102 (2017).
- 831 20. Tan, Y. & Cahan, P. SingleCellNet: A Computational Tool to Classify Single Cell
832 RNA-Seq Data Across Platforms and Across Species. *Cell Syst.* **9**, 207–213.e2
833 (2019).
- 834 21. Cancer Genome Atlas Network. Comprehensive molecular characterization of
835 human colon and rectal cancer. *Nature* **487**, 330–337 (2012).
- 836 22. Zhang, J. *et al.* International Cancer Genome Consortium Data Portal--a one-stop
837 shop for cancer genomics data. *Database (Oxford)* **2011**, bar026 (2011).
- 838 23. Tan, Y. & Cahan, P. SingleCellNet: a computational tool to classify single cell RNA-
839 Seq data across platforms and across species. *BioRxiv* (2018). doi:10.1101/508085
- 840 24. Cancer Genome Atlas Network. Comprehensive molecular portraits of human
841 breast tumours. *Nature* **490**, 61–70 (2012).
- 842 25. Parker, J. S. *et al.* Supervised risk predictor of breast cancer based on intrinsic
843 subtypes. *J. Clin. Oncol.* **27**, 1160–1167 (2009).
- 844 26. Wilkerson, M. D. *et al.* Lung squamous cell carcinoma mRNA expression subtypes
845 are reproducible, clinically important, and correspond to normal cell types. *Clin.*
846 *Cancer Res.* **16**, 4864–4875 (2010).
- 847 27. Cancer Genome Atlas Research Network. Electronic address:
848 andrew_aguirre@dfci.harvard.edu & Cancer Genome Atlas Research Network.
849 Integrated genomic characterization of pancreatic ductal adenocarcinoma. *Cancer*
850 *Cell* **32**, 185–203.e13 (2017).
- 851 28. Cancer Genome Atlas Research Network *et al.* Integrated genomic characterization
852 of endometrial carcinoma. *Nature* **497**, 67–73 (2013).
- 853 29. Cancer Genome Atlas Research Network *et al.* Integrated genomic characterization
854 of oesophageal carcinoma. *Nature* **541**, 169–175 (2017).
- 855 30. Cancer Genome Atlas Network. Comprehensive genomic characterization of head
856 and neck squamous cell carcinomas. *Nature* **517**, 576–582 (2015).
- 857 31. Cancer Genome Atlas Research Network. Comprehensive molecular
858 characterization of clear cell renal cell carcinoma. *Nature* **499**, 43–49 (2013).
- 859 32. Verhaak, R. G. W. *et al.* Integrated genomic analysis identifies clinically relevant
860 subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR,
861 and NF1. *Cancer Cell* **17**, 98–110 (2010).
- 862 33. Cancer Genome Atlas Research Network. Comprehensive molecular profiling of
863 lung adenocarcinoma. *Nature* **511**, 543–550 (2014).

- 864 34. Hu, B. *et al.* Gastric cancer: Classification, histology and application of molecular
865 pathology. *J. Gastrointest. Oncol.* **3**, 251–261 (2012).
- 866 35. Barretina, J. *et al.* The Cancer Cell Line Encyclopedia enables predictive modelling
867 of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).
- 868 36. Medico, E. *et al.* The molecular landscape of colorectal cancer cell lines unveils
869 clinically actionable kinase targets. *Nat. Commun.* **6**, 7002 (2015).
- 870 37. Park, J.-G. *et al.* Characteristics of Cell Lines Established from Human Colorectal
871 Carcinoma. *Cancer Res.* (1987).
- 872 38. Jerby-Arnon, L. *et al.* A cancer cell program promotes T cell exclusion and
873 resistance to checkpoint blockade. *Cell* **175**, 984–997.e24 (2018).
- 874 39. Darmanis, S. *et al.* Single-Cell RNA-Seq Analysis of Infiltrating Neoplastic Cells at
875 the Migrating Front of Human Glioblastoma. *Cell Rep.* **21**, 1399–1410 (2017).
- 876 40. Lee, J. *et al.* Tumor stem cells derived from glioblastomas cultured in bFGF and
877 EGF more closely mirror the phenotype and genotype of primary tumors than do
878 serum-cultured cell lines. *Cancer Cell* **9**, 391–403 (2006).
- 879 41. Patel, A. P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in
880 primary glioblastoma. *Science* **344**, 1396–1401 (2014).
- 881 42. Black, J. D., English, D. P., Roque, D. M. & Santin, A. D. Targeted therapy in
882 uterine serous carcinoma: an aggressive variant of endometrial cancer. *Womens*
883 *Health (Lond. Engl.)* **10**, 45–57 (2014).
- 884 43. Yang, S., Thiel, K. W. & Leslie, K. K. Progesterone: the ultimate endometrial tumor
885 suppressor. *Trends Endocrinol. Metab.* **22**, 145–152 (2011).
- 886 44. Huszar, M. *et al.* Up-regulation of L1CAM is linked to loss of hormone receptors and
887 E-cadherin in aggressive subtypes of endometrial carcinomas. *J. Pathol.* **220**, 551–
888 561 (2010).
- 889 45. Kozak, J., Wdowiak, P., Maciejewski, R. & Torres, A. A guide for endometrial
890 cancer cell lines functional assays using the measurements of electronic
891 impedance. *Cytotechnology* **70**, 339–350 (2018).
- 892 46. Korch, C. *et al.* DNA profiling analysis of endometrial and ovarian cell lines reveals
893 misidentification, redundancy and contamination. *Gynecol. Oncol.* **127**, 241–248
894 (2012).
- 895 47. Wu, D. *et al.* Gene-expression data integration to squamous cell lung cancer
896 subtypes reveals drug sensitivity. *Br. J. Cancer* **109**, 1599–1608 (2013).
- 897 48. Walter, V. *et al.* Molecular subtypes in head and neck cancer exhibit distinct
898 patterns of chromosomal gain and loss of canonical cancer genes. *PLoS One* **8**,
899 e56823 (2013).
- 900 49. Adeegbe, D. O. *et al.* BET Bromodomain Inhibition Cooperates with PD-1 Blockade
901 to Facilitate Antitumor Response in Kras-Mutant Non-Small Cell Lung Cancer.
902 *Cancer Immunol Res* **6**, 1234–1245 (2018).
- 903 50. Blaisdell, A. *et al.* Neutrophils oppose uterine epithelial carcinogenesis via
904 debridement of hypoxic tumor cells. *Cancer Cell* **28**, 785–799 (2015).
- 905 51. Fitamant, J. *et al.* YAP inhibition restores hepatocyte differentiation in advanced
906 HCC, leading to tumor regression. *Cell Rep.* **10**, 1692–1707 (2015).
- 907 52. Jia, D. *et al.* Crebbp loss drives small cell lung cancer and increases sensitivity to
908 HDAC inhibition. *Cancer Discov.* **8**, 1422–1437 (2018).

- 909 53. Kress, T. R. *et al.* Identification of MYC-Dependent Transcriptional Programs in
910 Oncogene-Addicted Liver Tumors. *Cancer Res.* **76**, 3463–3472 (2016).
- 911 54. Li, L. *et al.* GKAP acts as a genetic modulator of NMDAR signaling to govern
912 invasive tumor growth. *Cancer Cell* **33**, 736–751.e5 (2018).
- 913 55. Mollaoglu, G. *et al.* The Lineage-Defining Transcription Factors SOX2 and NKX2-1
914 Determine Lung Cancer Cell Fate and Shape the Tumor Immune
915 Microenvironment. *Immunity* **49**, 764–779.e9 (2018).
- 916 56. Pan, Y. *et al.* Whole tumor RNA-sequencing and deconvolution reveal a clinically-
917 prognostic PTEN/PI3K-regulated glioma transcriptional signature. *Oncotarget* **8**,
918 52474–52487 (2017).
- 919 57. Lissanu Deribe, Y. *et al.* Mutations in the SWI/SNF complex induce a targetable
920 dependence on oxidative phosphorylation in lung cancer. *Nat. Med.* **24**, 1047–1057
921 (2018).
- 922 58. Xu, C. *et al.* Loss of Lkb1 and Pten leads to lung squamous cell carcinoma with
923 elevated PD-L1 expression. *Cancer Cell* **25**, 590–604 (2014).
- 924 59. Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**,
925 719–724 (2009).
- 926 60. Balkwill, F. R., Capasso, M. & Hagemann, T. The tumor microenvironment at a
927 glance. *J. Cell Sci.* **125**, 5591–5596 (2012).
- 928 61. Ben-David, U. *et al.* Patient-derived xenografts undergo mouse-specific tumor
929 evolution. *Nat. Genet.* **49**, 1567–1575 (2017).
- 930 62. Hristova, V. A. & Chan, D. W. Cancer biomarker discovery and translation:
931 proteomics and beyond. *Expert Rev Proteomics* **16**, 93–103 (2019).
- 932 63. Dawson, M. A. & Kouzarides, T. Cancer epigenetics: from mechanism to therapy.
933 *Cell* **150**, 12–27 (2012).
- 934 64. Silva, T. C. *et al.* TCGA Workflow: Analyze cancer genomics and epigenomics data
935 using Bioconductor packages. [version 2; peer review: 1 approved, 2 approved with
936 reservations]. *F1000Res.* **5**, 1542 (2016).
- 937 65. Morgan, M., Obenchain, V., Hester, J. & Pag`es, H. *SummarizedExperiment*:
938 *SummarizedExperiment container.* (2018).
- 939 66. Pavlidis, P. & Noble, W. S. Analysis of strain and regional variation in gene
940 expression in mouse brain. *Genome Biol.* **2**, RESEARCH0042 (2001).
- 941 67. Geman, D., d Avignon, C., Naiman, D. Q. & Winslow, R. L. Classifying gene
942 expression profiles from pairwise mRNA comparisons. *Stat Appl Genet Mol Biol* **3**,
943 Article19 (2004).
- 944 68. Kolde, R. *heatmap: Pretty Heatmaps.* (CRAN, 2019).
- 945 69. Wickham, H. *ggplot2 - Elegant Graphics for Data Analysis* . (Springer-Verlag New
946 York, 2016). doi:10.1007/978-0-387-98141-3
- 947 70. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations
948 in multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).
- 949 71. Yoshihara, K. *et al.* Inferring tumour purity and stromal and immune cell admixture
950 from expression data. *Nat. Commun.* **4**, 2612 (2013).
- 951 72. Kovalchik, S. *RISmed: Download Content from NCBI Databases.* (CRAN.R-project,
952 2017).
- 953

Figure 1

A



B

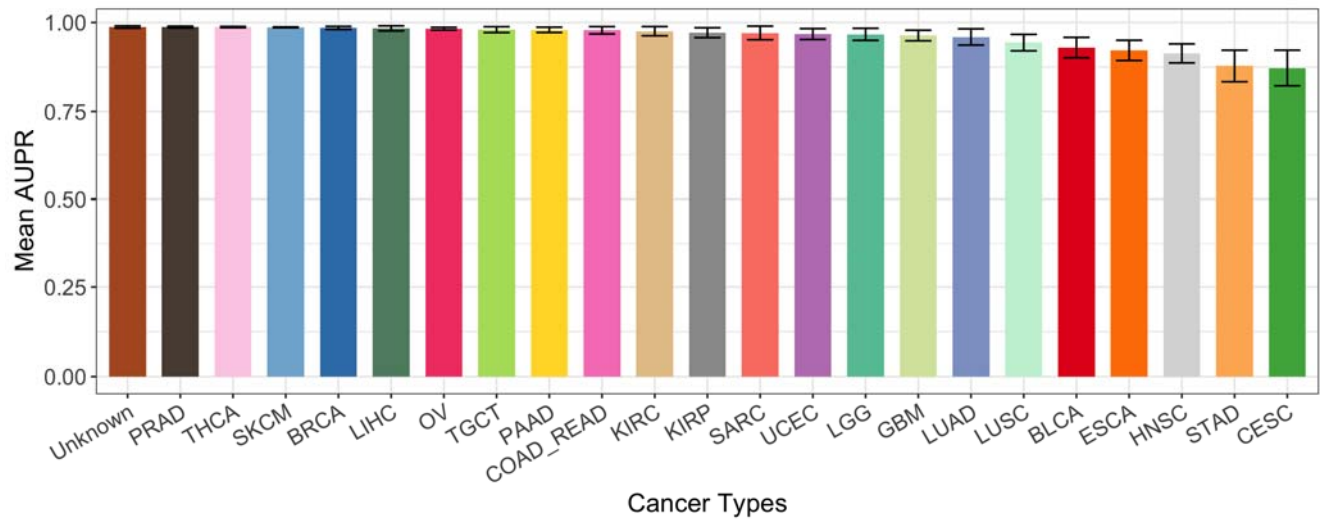


Figure 2

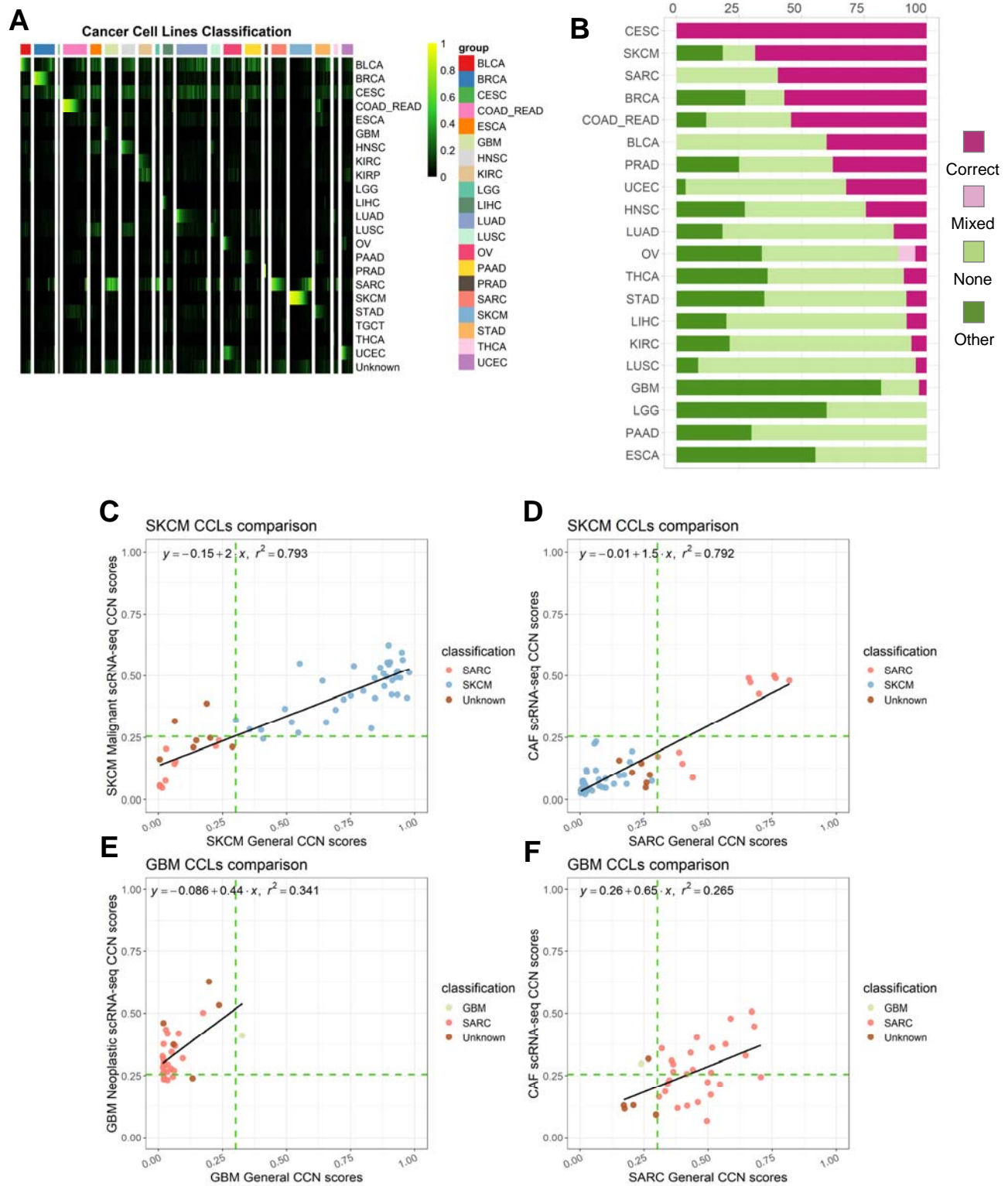


Figure 3

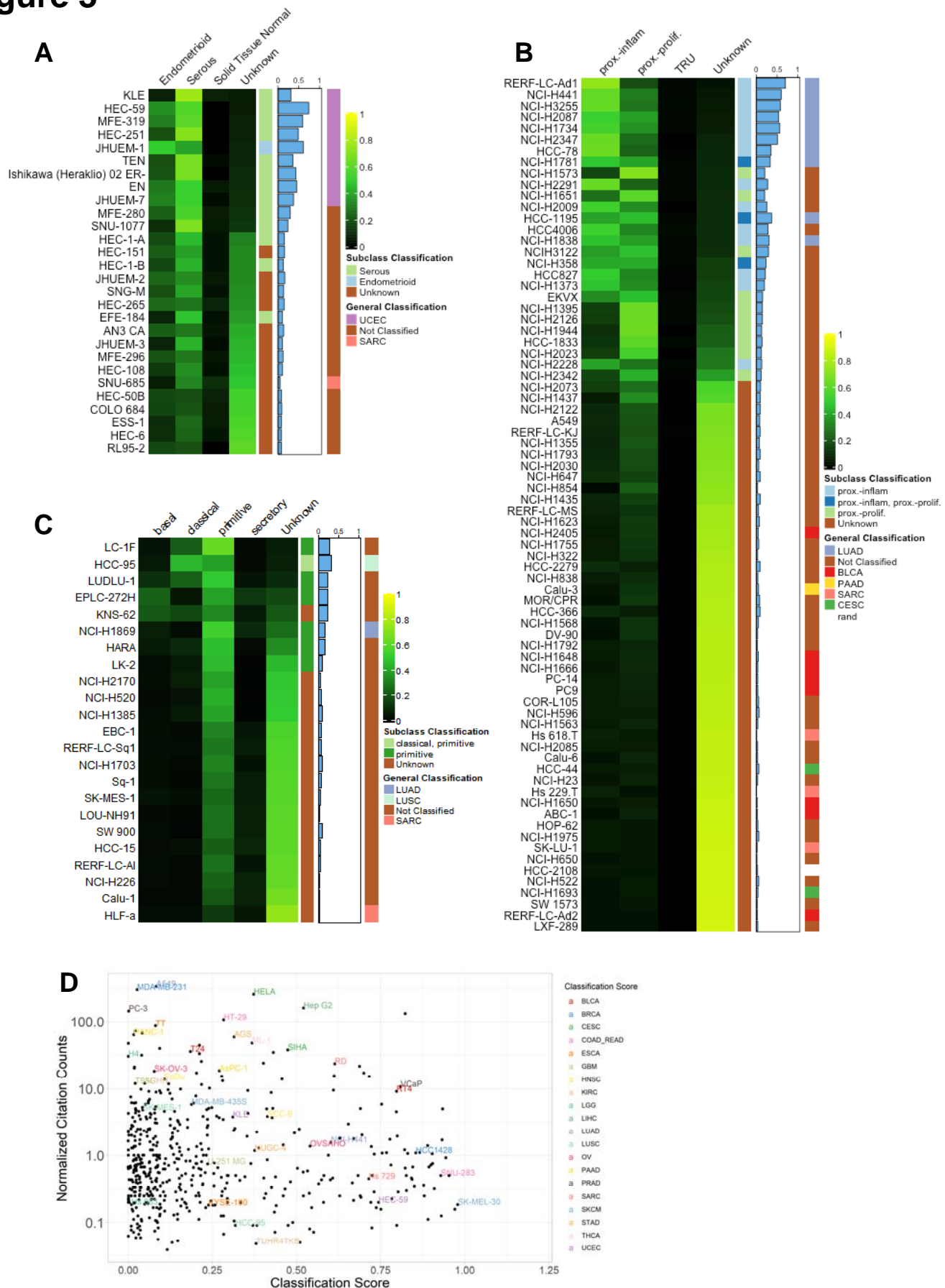


Figure 4

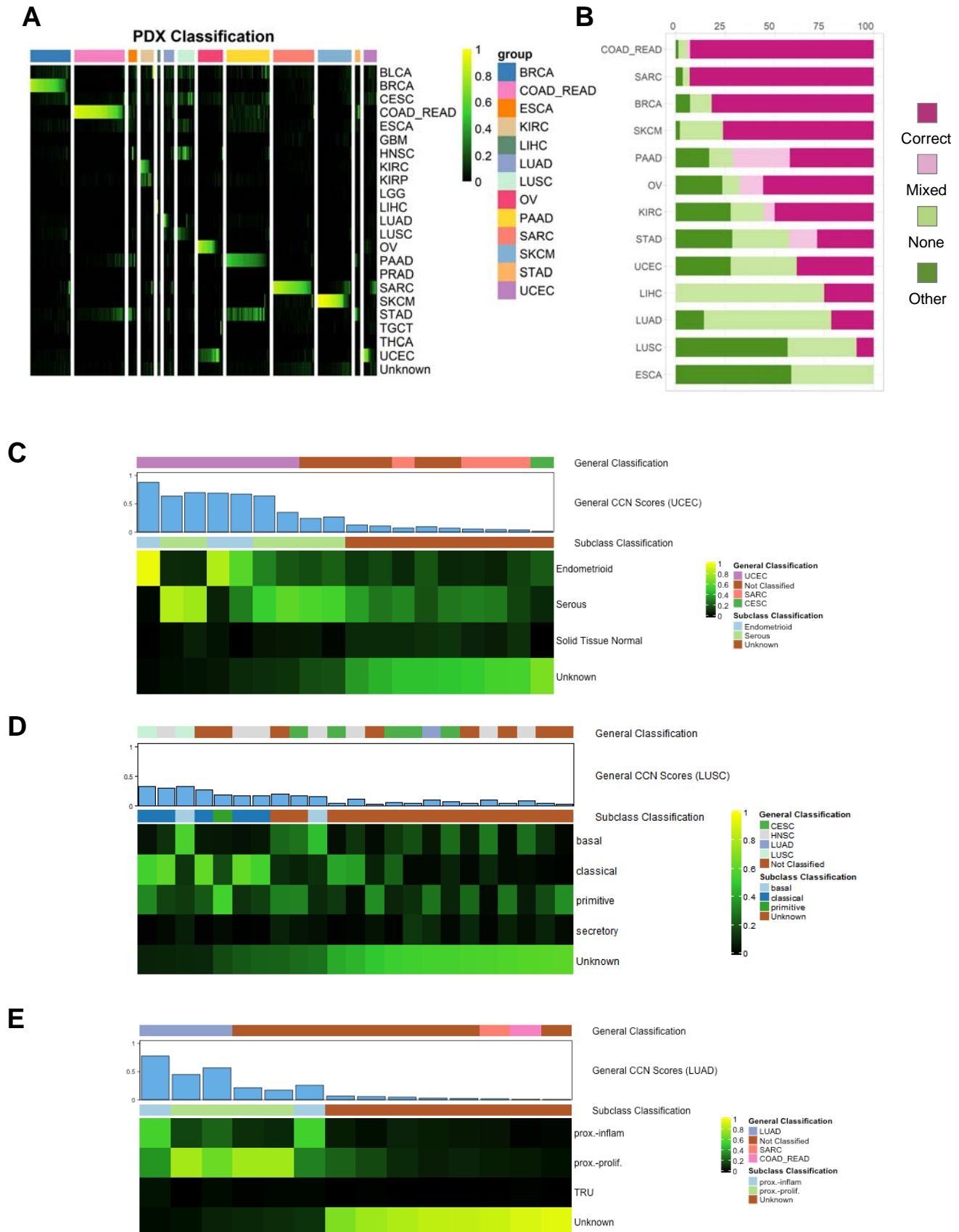


Figure 5

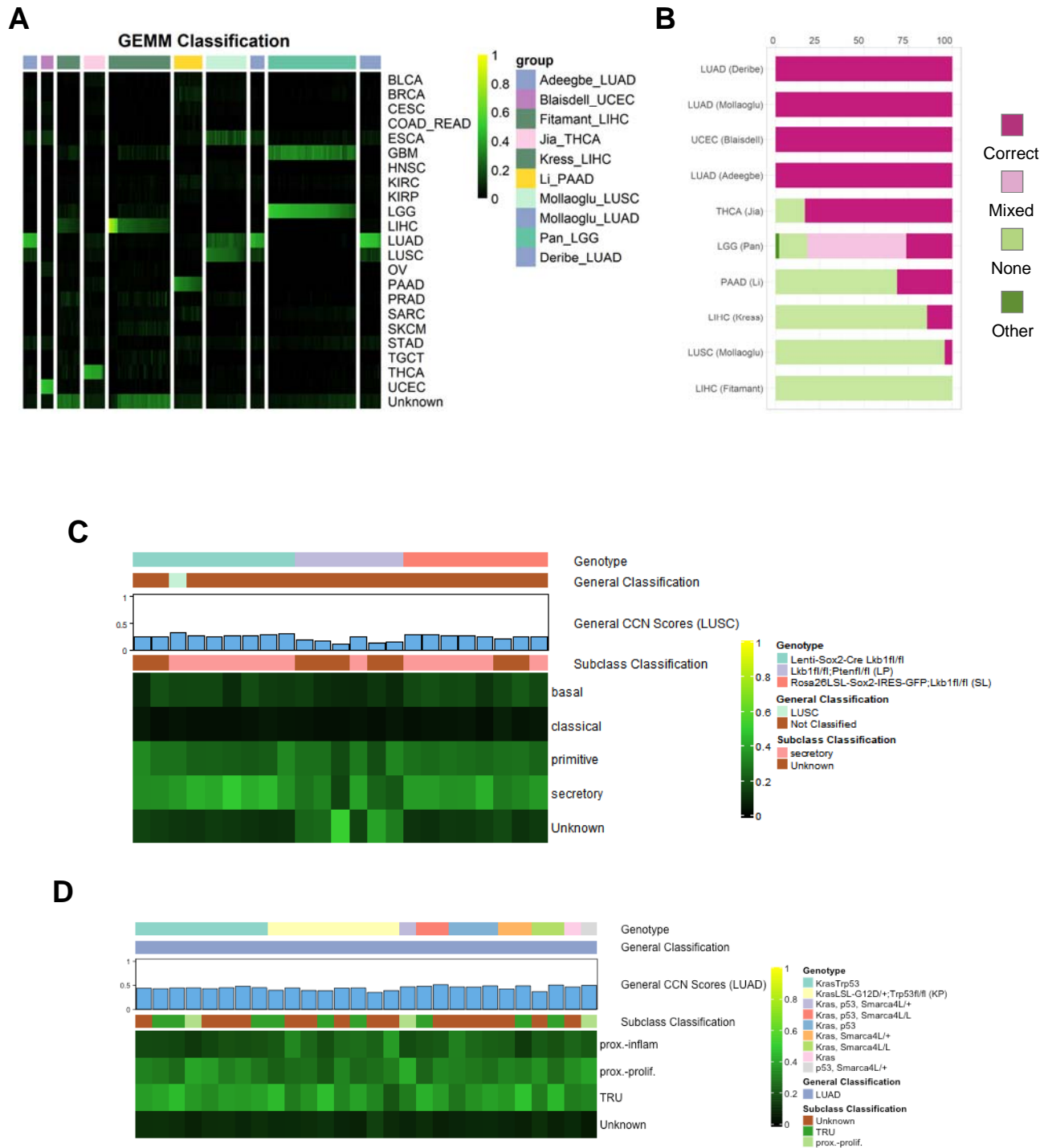
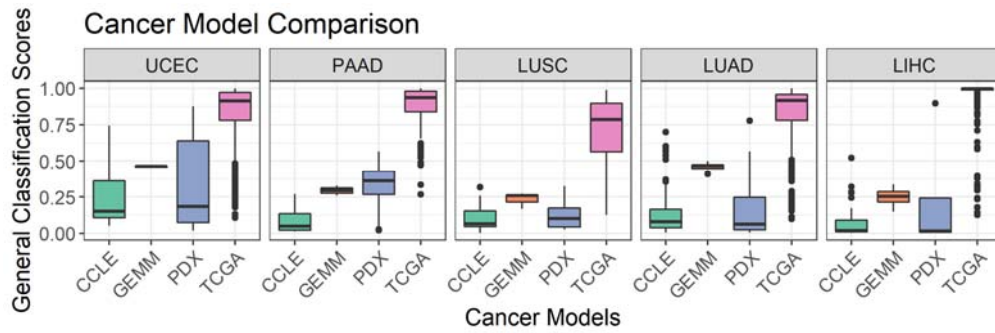
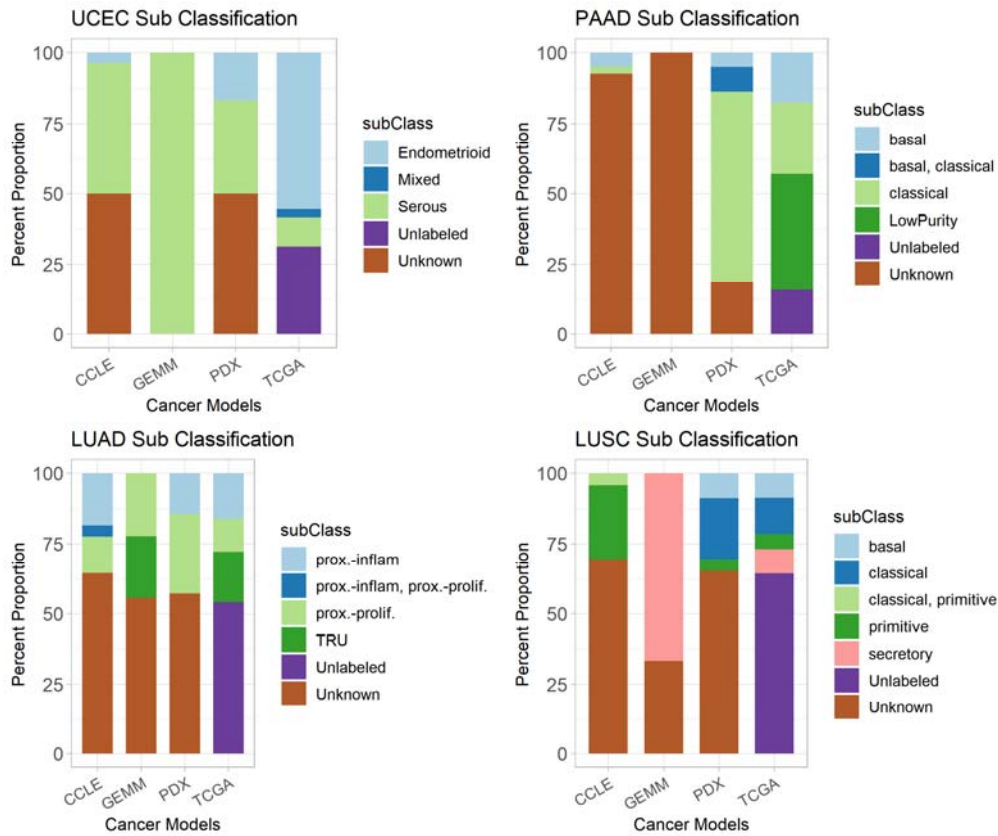


Figure 6

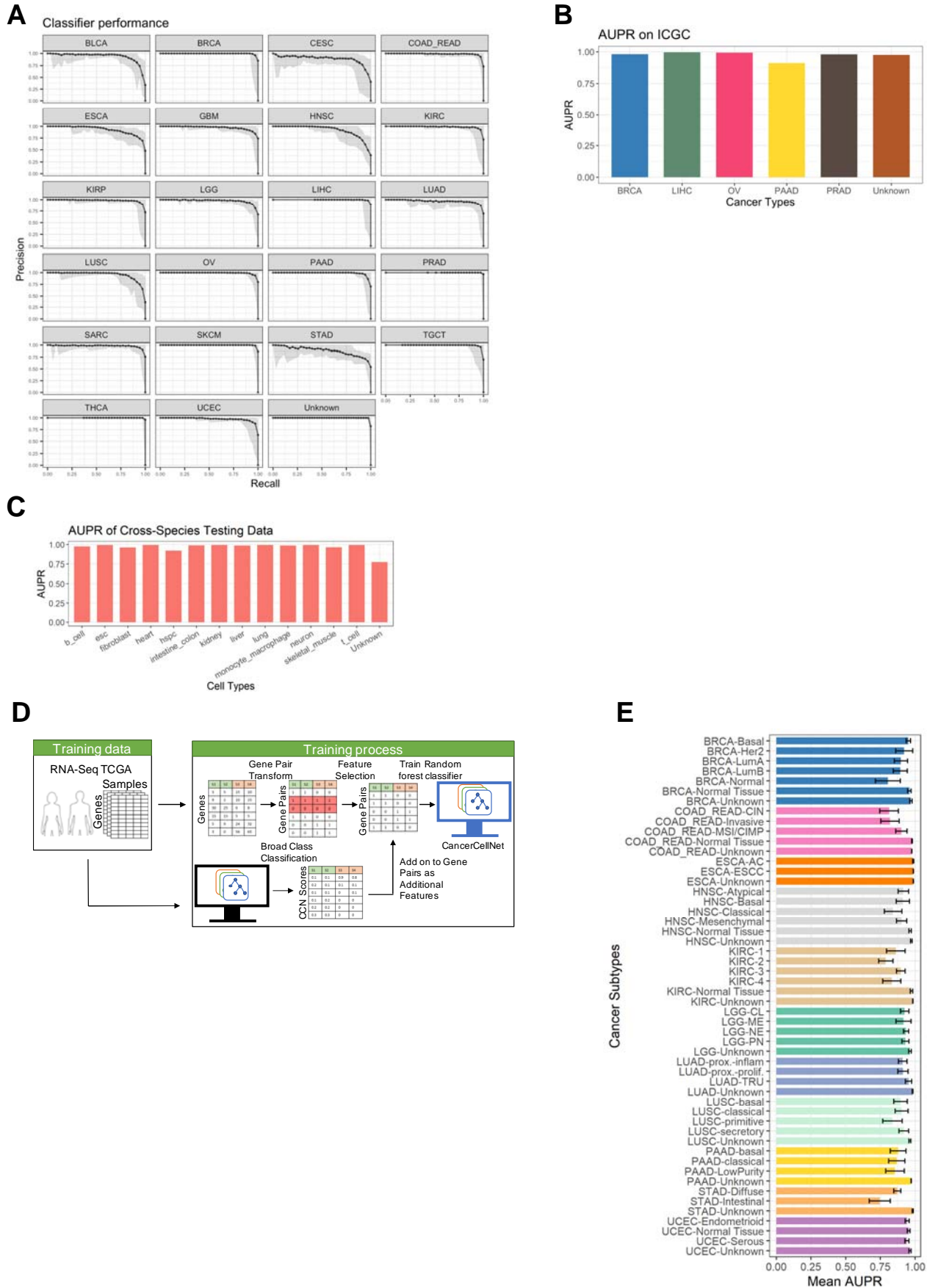
A



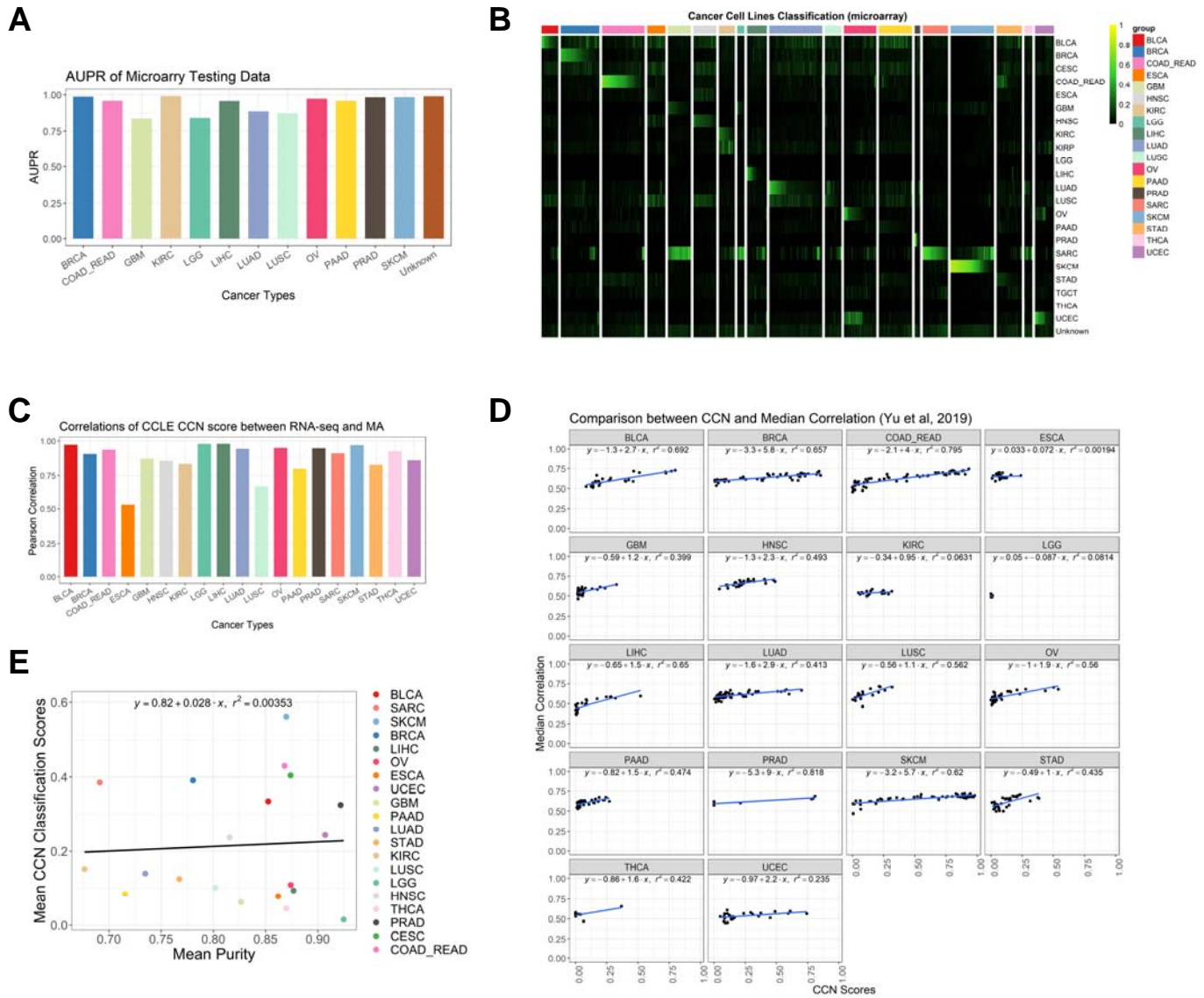
B



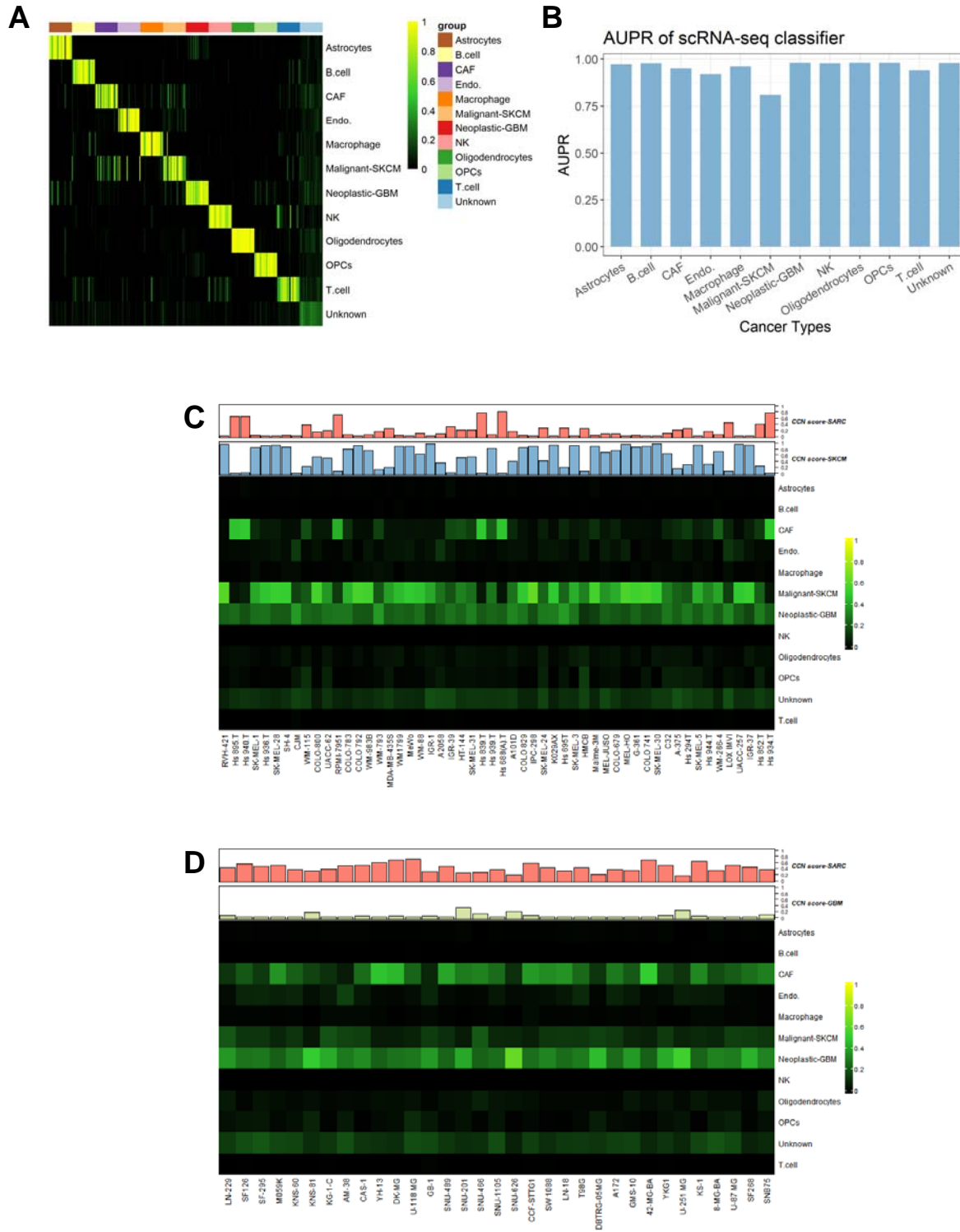
Supplemental Figure 1



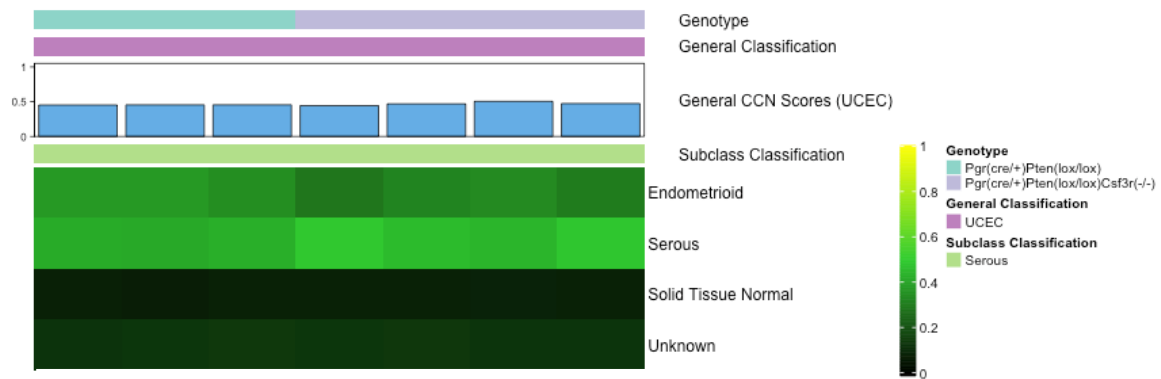
Supplemental Figure 2



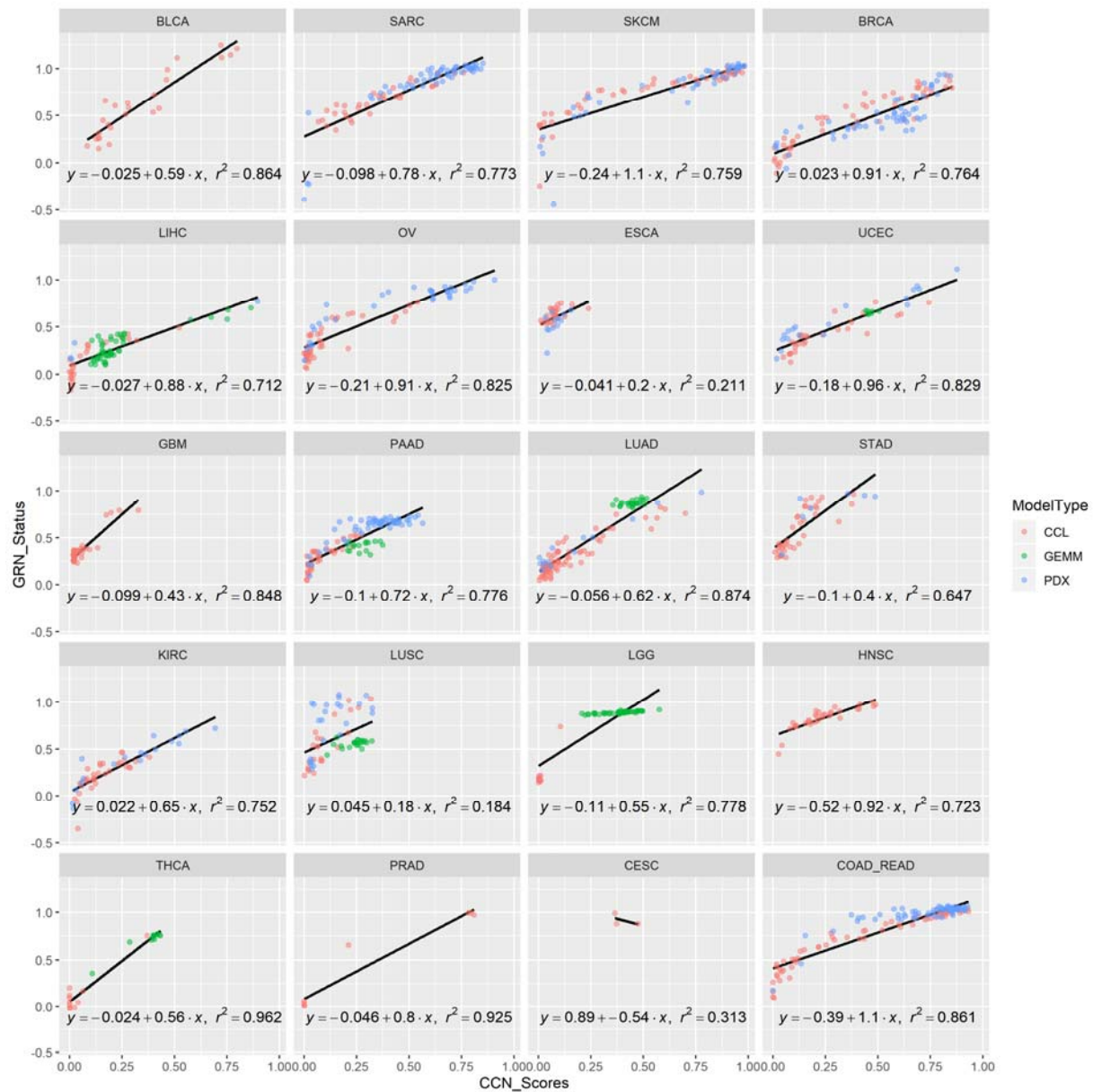
Supplemental Figure 3



Supplemental Figure 4



Supplemental Figure 5



Supplemental Figure 6

