

1 **Predicting cancer prognosis and drug response from the tumor microbiome**

2 Leandro C. Hermida^{1,2†}, E. Michael Gertz^{1†}, Eytan Ruppin^{1*}.

3

4 ¹ Cancer Data Science Laboratory (CDSL), National Cancer Institute (NCI), National Institutes
5 of Health (NIH), Bethesda, MD, USA.

6 ² Department of Computer Science, University of Maryland, College Park, MD, USA.

7

8 † Equally contributing first authors

9 * Corresponding author (eytan.ruppin@nih.gov)

10 Abstract

11 Tumor gene expression is predictive of patient prognosis in some cancers. However, RNA-
12 seq and whole genome sequencing data contain not only reads from host tumor and normal
13 tissue, but also reads from the tumor microbiome, which can be used to infer the microbial
14 abundances in each tumor. Here, we show that tumor microbial abundances, alone or in
15 combination with tumor gene expression data, can predict cancer prognosis and drug response to
16 some extent – microbial abundances are significantly less predictive of prognosis than gene
17 expression, although remarkably, similarly as predictive of drug response, but in mostly different
18 cancer-drug combinations. Thus, it appears possible to leverage existing sequencing technology,
19 or develop new protocols, to obtain more non-redundant information about prognosis and drug
20 response from RNA-seq and whole genome sequencing experiments than could be obtained from
21 tumor gene expression or genomic data alone.

22 Introduction

23 The Cancer Genome Atlas (TCGA), available from the NCI Genomic Data Commons
24 (GDC)¹, provides RNA-seq and whole genomic sequencing (WGS) data for thousands of cases
25 across dozens of cancer types. RNA-seq data is typically used to measure the expression of
26 human genes, and there is a long history linking tumor gene expression to cancer outcomes²⁻⁸.
27 Milanez-Almeida et al.⁹ recently showed that gene expression from TCGA RNA-seq data could
28 predict overall survival (OS) or progression-free interval (PFI) better than classical clinical
29 prognostic covariates – age at diagnosis, gender, and tumor stage. Importantly, Milanez-Almeida
30 et al. used a data-driven machine learning (ML) based approach which selected features that
31 were predictive of and correlated with prognosis, rather than approaches based on classical
32 statistics or biological knowledge that chose features *a priori*.

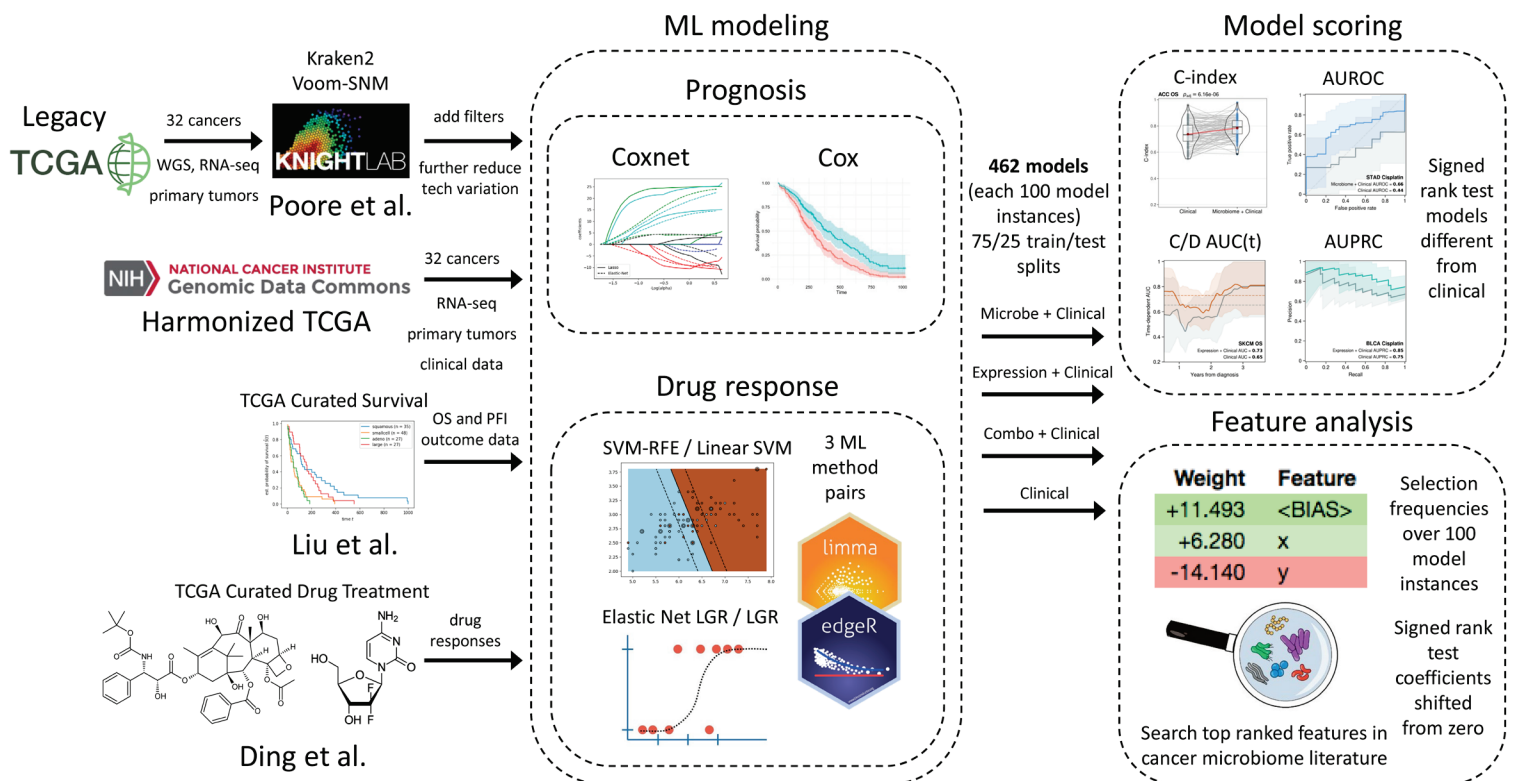
33 Research into the human tumor microbiome has been rapidly expanding, and multiple
34 laboratories have attempted to utilize existing technologies and data to identify microbes and
35 quantify their abundance within human tumors compared to adjacent normal tissue. RNA-seq
36 and WGS data not only contain human sequencing reads, but also reads from the local intratumor
37 microbiome that are typically filtered out from the data when analyzing human gene expression
38 or genomic alterations. Poore et al.¹⁰ recently developed a computational workflow, using two
39 orthogonal microbial detection pipelines, to estimate, decontaminate, normalize, and batch effect
40 correct microbial abundances from human high-throughput sequencing data. They applied this
41 workflow to create a first-of-its-kind comprehensive dataset of pan-cancer tumor microbial
42 abundances derived from WGS or RNA-seq data for the entire TCGA cohort.

43 Our central research questions then were, 1) does a data-driven ML approach reveal that
44 tumor microbial abundances in TCGA data, quantified from these reads, are predictive of cancer
45 prognosis or drug response, 2) what microbial genera are potentially predictive biomarkers of
46 prognosis or drug response, 3) how do these models compare to equivalent models based on
47 tumor gene expression data, and 4) does combining both microbial abundance and gene
48 expression features produce models and select combinations of genes and microbial genera that
49 are more predictive of prognosis or drug response than models from each individual data type?
50 We used the processed microbial abundances directly from the Poore et al. dataset to build
51 predictive models of prognosis and drug response for TCGA. We also used TCGA RNA-seq
52 read counts to build equivalent predictive models for comparison. As a positive control, we also
53 showed that our prognosis ML modeling methods, which differed somewhat from Milanez-
54 Almeida et al.⁹, identified a similar set cancers and outcomes for which gene expression was
55 predictive of prognosis.

56 Results

57 **Tumor microbial abundances are substantially less predictive of prognosis than gene**
58 **expression**

59 An overview of the analytical workflow is presented in **Fig. 1**. It has four major parts, 1) data
60 download and preprocessing, 2) prognosis and drug response ML modeling, 3) model evaluation
61 and scoring, and 4) further feature analysis. A more detailed technical description of the analysis
62 pipeline and computational methods is provided in Methods.

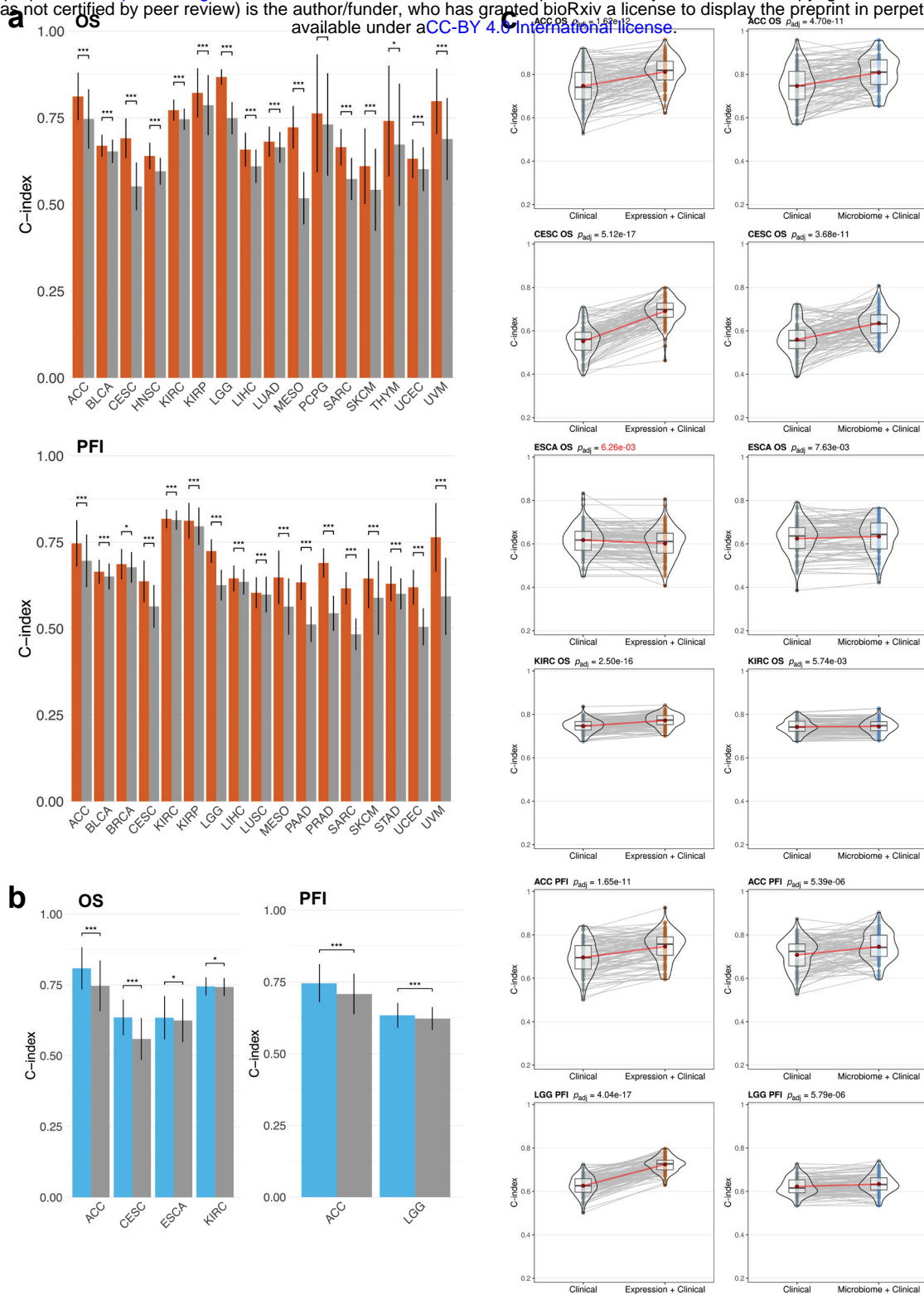


63 **Figure 1. Analysis pipeline overview.** Download and data preprocessing (left) of Poore et al. TCGA primary tumor
 64 Kraken2 Voom-SNM microbial abundances with additional filters to reduce technical variation, NCI Genomic Data
 65 Commons (GDC) harmonized TCGA primary tumor RNA-seq counts and clinical data, Liu et al. TCGA curated
 66 overall survival (OS) and progression-free interval (PFI) outcome data, and Ding et al. TCGA curated drug response
 67 clinical data. Prognosis machine learning (ML) modeling (middle) of microbial abundance, gene expression, and
 68 combined data types with clinical covariates for each cancer using penalized Cox with elastic net penalties (Coxnet)
 69 against matched clinical covariate-only models using standard Cox regression. Drug response classification ML
 70 modeling of the same data types with clinical covariates for each cancer-drug combination using three ML
 71 approaches, 1) SVM-RFE, elastic net logistic regression (LGR), and limma-trend (microbial and combined data
 72 types) or edgeR (gene expression) differential analysis feature scoring and selection with L2 penalized LGR.
 73 Matched clinical covariate-only modeling performed with L2 penalized linear SVM or LGR. ML modeling
 74 generates 100 model instances for each model from 75/25 train/test randomly shuffled and stratified dataset splits.
 75 ML model instance scoring (right top) using concordance index (C-index) and time-dependent cumulative/dynamic
 76 AUC (C/D AUC(t)) for prognosis models and area under receiver-operating characteristic curve (AUROC) and area
 77 under the precision-recall curve (AUPRC) for drug response models. Significance of model performance
 78 improvement over matched clinical covariate-only model determined by signed rank test of C-index or AUROC
 79 scores between each matched model instance for prognosis and drug response models, respectively. Feature analysis
 80 (right bottom) performed using model instance coefficients and selection frequencies. Overall feature importance
 81 ranking and significance determined by signed rank test of model instance feature coefficients shifting from zero
 82 and filtering of top features for selection frequency $\geq 20\%$.

83 We built OS and PFI gene expression ML models of 32 TCGA tumor types (see
84 **Supplementary Data 1** for cohort information) using the Coxnet¹¹ algorithm, which jointly
85 selects the most predictive subset of features via cross-validation (CV) while simultaneously
86 being able to control for prognostic clinical covariates. In our models, we included and
87 controlled for the clinical covariates age at diagnosis, gender, and tumor stage. For comparison,
88 we also built standard Cox regression models based on the clinical covariates alone. We
89 evaluated the predictive performance of our models using Harrell's concordance index (C-
90 index), which is a metric of survival model predictive accuracy. Each model analysis generated
91 100 model instances and C-index scores from randomly shuffled train-test CV splits on the data.
92 We found 33 OS and PFI models for 21 tumor types that had a mean C-index score ≥ 0.6 and
93 significantly outperformed their corresponding clinical covariate-only models (**Fig. 2a & c,**
94 **Supplementary Figs. 1a, 2a**). Our models were predictive of prognosis in 11 of the same 13
95 tumor types that were reported by Milanez-Almeida et al.⁹ (**Supplementary Table 1**). We did
96 not analyze one tumor type that Milanez-Almeida did, acute myeloid leukemia (LAML), because
97 Poore et al. excluded it from their analysis. Among the cancers and outcomes that Milanez-
98 Almeida et al. analyzed, our methodology produced predictive models for four additional tumor
99 types: breast cancer (BRCA), cervical squamous cell carcinoma (CESC), sarcoma (SARC), and
100 uterine corpus endometrial carcinoma (UCEC), as well as quite a few predictive models for
101 additional cancers and outcomes that were not analyzed in their study (**Supplementary Table 1**).
102 We also evaluated prognosis model performance by calculating the time-dependent,
103 cumulative/dynamic area under the curve ($AUC^{C/D}(t)$)^{12,13}, which is an extension of the area
104 under the receiver-operating characteristic curve (AUROC) for continuous outcomes and can
105 provide a more detailed resolution picture of predictive performance throughout the test outcome
106 time range compared to the C-index score. Although 33 of our OS and PFI gene expression
107 models had a statistically significant C-index score improvement compared to clinical covariates

108 alone, only 22 of these models showed an improvement in $AUC^{C/D}(t)$, where the improvement in
109 mean $AUC^{C/D}(t)$ over the entire test time range after diagnosis was ≥ 0.025 (**Supplementary**
110 **Figs. 1b, 2b**).

111 We applied Coxnet¹¹ using the same methodology to build prognosis models using the
112 microbial abundance estimates provided by Poore et al.¹⁰. We found six microbial abundance
113 models that had a mean C-index score ≥ 0.6 and significantly outperformed their corresponding
114 clinical covariate-only models (**Fig. 2b & c, Supplementary Fig. 3a**). We found that in only two
115 of the six models, microbial abundances outperformed clinical covariates alone in terms of
116 $AUC^{C/D}(t)$, where the improvement in mean $AUC^{C/D}(t)$ over the entire test time range after
117 diagnosis was ≥ 0.025 (**Supplementary Fig. 3b**). In adrenocortical carcinoma (ACC), microbial
118 features predicted OS significantly better than clinical prognostic covariates starting at
119 approximately 6 years after diagnosis. In CESC, microbial abundances predicted OS better than
120 clinical covariates from approximately 6 months to 10 years after diagnosis. Overall, we found
121 that tumor microbial abundances from Poore et al. were only marginally predictive of prognosis
122 across the TCGA cohort, and that gene expression was a significantly more powerful predictor of
123 prognosis (**Fig. 2, Supplementary Figs. 1-3**).



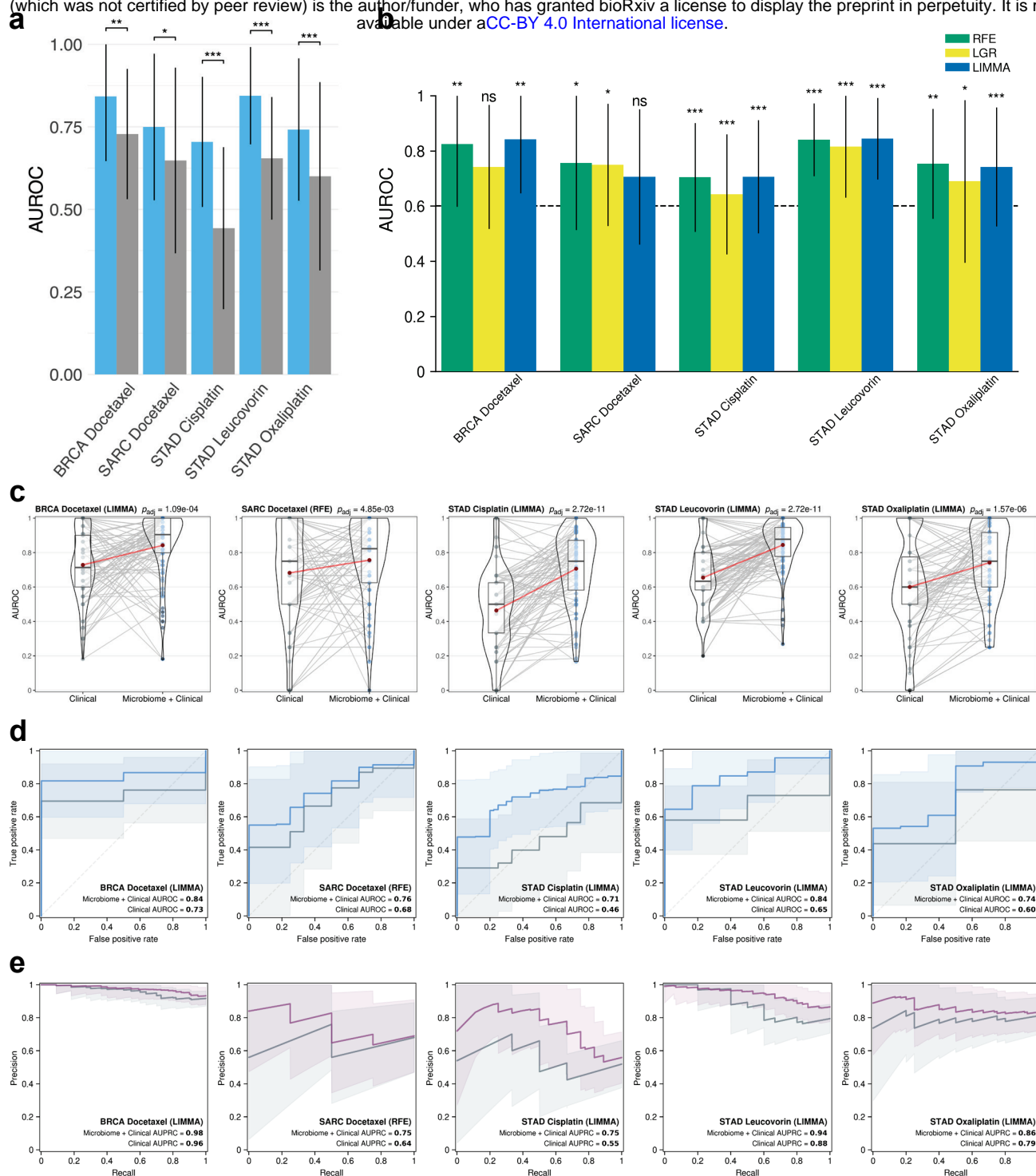
124 **Figure 2. Performance of gene expression and microbial abundance prognosis prediction models where**
 125 **features add predictive power to clinical covariates.** Mean C-index scores for (a) gene expression with clinical
 126 covariate models (orange) and (b) microbial abundance with clinical covariate models (blue) vs clinical covariate-
 127 only models (grey). Error bars denote standard deviations. Significance: * ≤ 0.01 , ** ≤ 0.001 , *** ≤ 0.0001 . (c) C-
 128 index score violin density plots for the six models where microbial abundance with clinical covariate features
 129 outperform clinical covariate-only models. Corresponding gene expression models shown for comparison. Lines
 130 connecting points (light grey) represent score pairs from same train-test split on the data. Mean C-index scores
 131 shown as red dots with red lines connecting the means. Significance for the prediction improvement over clinical
 132 covariate-only models was calculated using a two-sided Wilcoxon signed-rank test and adjusted for multiple testing
 133 using the Benjamini-Hochberg method with adjusted p-values shown at top. Adjusted p-values colored in red signify
 134 difference where clinical covariate-only model is better.

135 **Tumor microbial abundances are predictive of chemotherapy drug response in some**
136 **cancers and in mostly different cancer-drug combinations than gene expression**

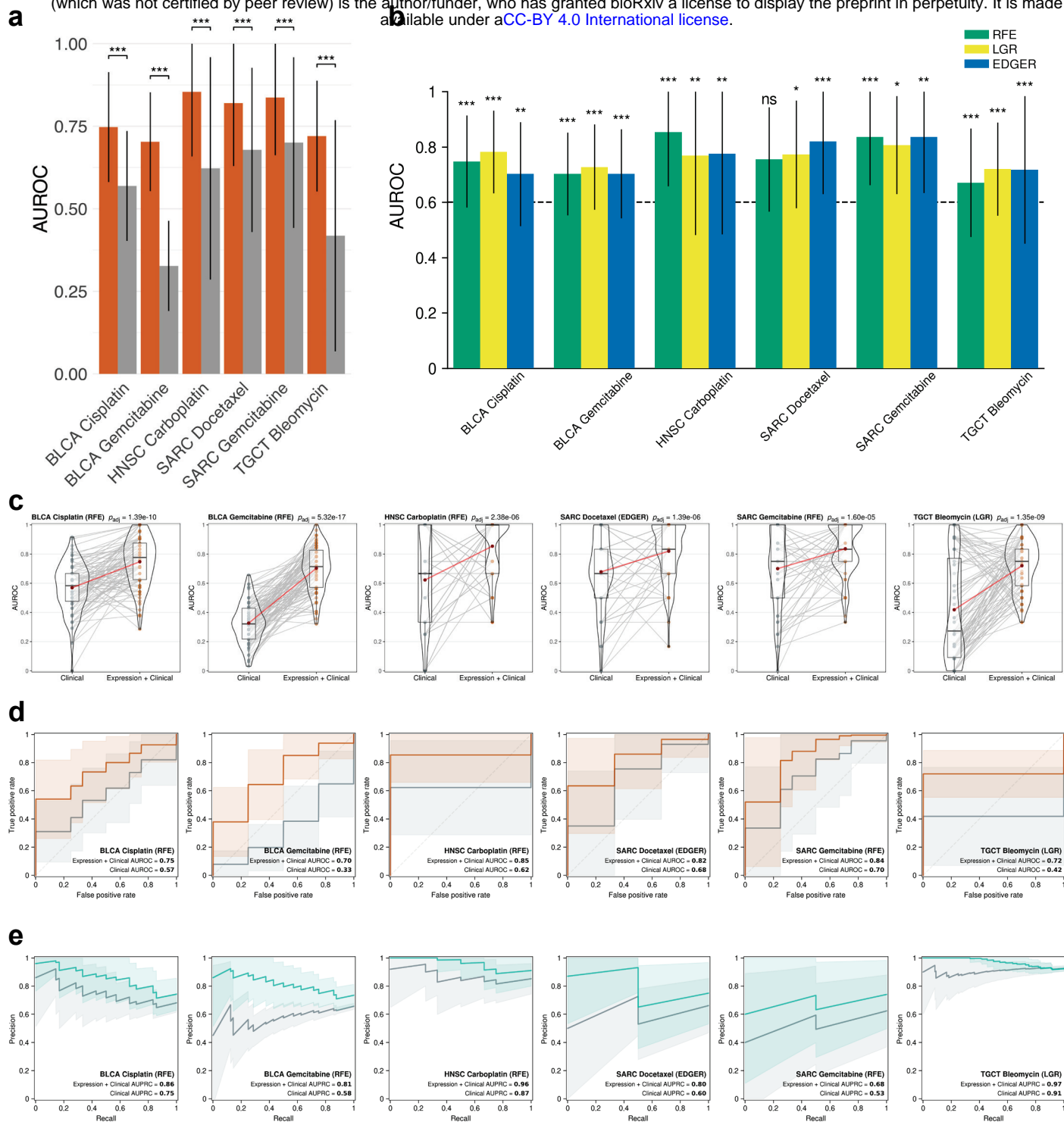
137 We next asked whether tumor microbial abundances from pre-treatment biopsies could
138 predict drug response better than the clinical covariates age at diagnosis, gender, and tumor stage
139 alone. TCGA drug response clinical data were obtained from Ding et al.¹⁴ as described in
140 Methods. Cases with complete response (CR) or partial response (PR) were labeled as
141 responders and those with stable disease (SD) or progressive disease (PD) as non-responders.
142 Thirty TCGA cancer-drug combinations met our minimum dataset size thresholds (see
143 **Supplementary Data 1** for cohort information). We built drug response models using three
144 different ML methods: 1) a variant of the linear support vector machine recursive feature
145 elimination (SVM-RFE) algorithm¹⁵ that we developed, 2) logistic regression (LGR) with elastic
146 net¹⁶ (L1 + L2) penalties and embedded feature selection, and 3) logistic regression with an L2
147 penalty and limma¹⁷ (for microbial abundance and combined data type datasets) or edgeR^{18,19}
148 (for RNA-seq count datasets) differential abundance/expression feature scoring and wrapper
149 selection methods (see Methods for details). All three ML methods unconditionally included the
150 clinical covariates – age at diagnosis, gender, and tumor stage – in the model (bypassing feature
151 selection) while selecting the most predictive subset of microbial abundance or gene expression
152 features. For comparison, we built standard linear SVM or LGR models using the clinical
153 covariates alone. We evaluated the predictive performance of drug response models using
154 AUROC. Each analysis generated 100 model instances, AUROC, and area under the precision-
155 recall curve (AUPRC) scores from randomly shuffled train-test CV splits on the data.

156 We found five microbial abundance cancer-drug combinations that had a mean AUROC
157 score ≥ 0.6 and performed better than clinical covariates alone in at least two out of three ML
158 methods (**Fig. 3**). Three of these cancer-drug combinations involved stomach adenocarcinoma

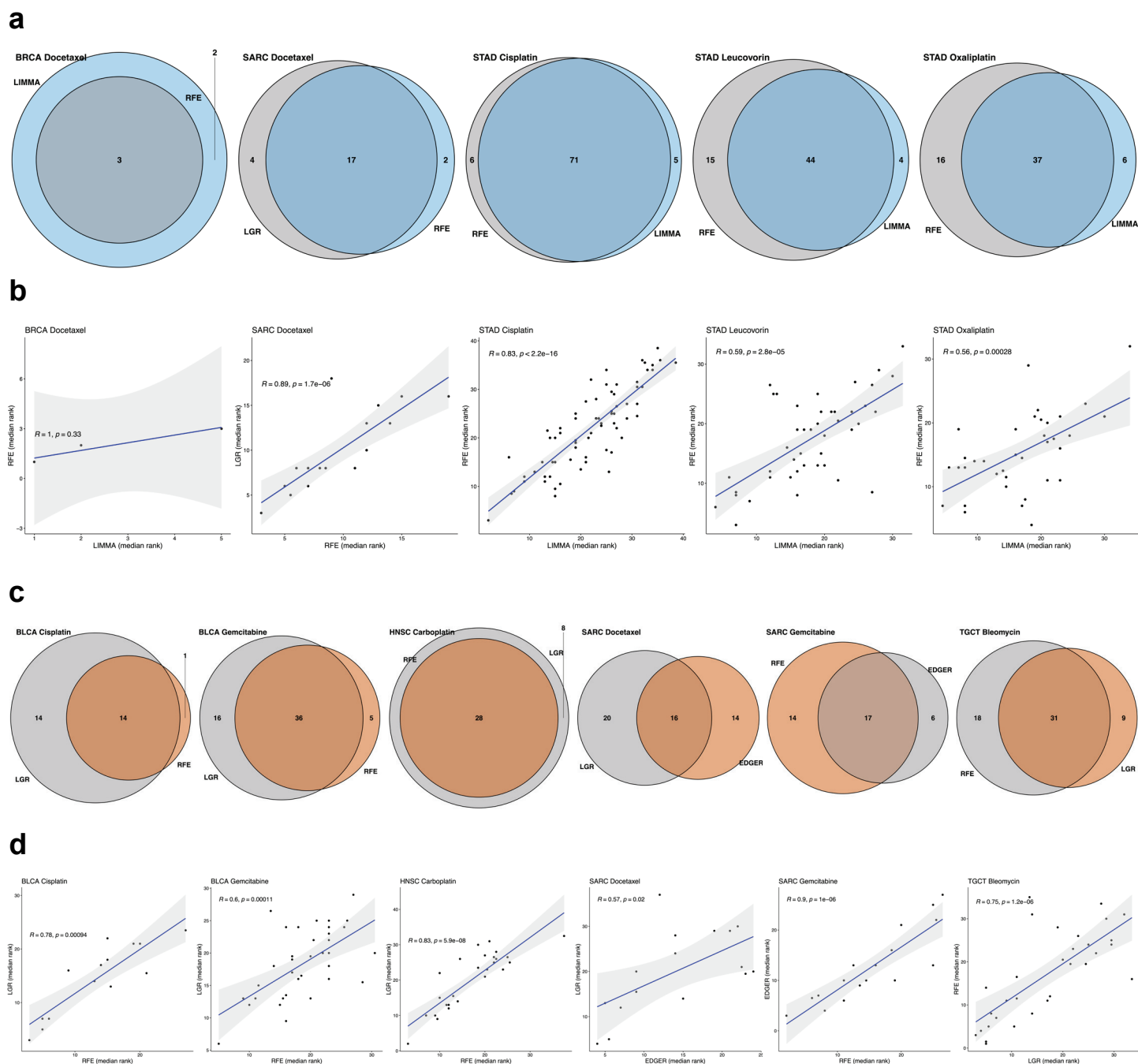
159 (STAD). We performed the same drug response modeling using TCGA gene expression data and
160 here we found six cancer-drug response combinations that had a mean AUROC score ≥ 0.6 and
161 significantly outperformed their corresponding clinical covariate-only models in at least two out
162 of three ML methods (**Fig. 4**). Only one cancer-drug combination, SARC docetaxel, overlapped
163 between the microbial abundance and gene expression drug response model results, suggesting
164 that tumor microbial abundances have independent predictive power. Even though one of our
165 thresholds for a significant drug response model was a mean AUROC score ≥ 0.6 , the 11 total
166 significant models that we found from both data types each had a mean AUROC > 0.7 . We also
167 found there was considerable overlap in the selected microbial abundance and gene expression
168 features reported by each ML method (**Fig. 5a & c**) and frequently found a significant
169 correlation between the feature importance rankings reported by each ML method when
170 comparing the two most significant methods in each cancer-drug combination (**Fig. 5b & d**).
171 These results suggest that our significant drug response models and their inferred important
172 features are not the result a specific ML modeling methodology. Overall, our results support the
173 notion that the tumor microbiome may contain information that is predictive of drug response in
174 some cancers, consistent with recent reports^{20,21}.



175 **Figure 3. Performance of microbial abundance drug response prediction models in the five cancer-drug**
 176 **combinations where models performed better than clinical covariates alone. (a)** Mean AUROC scores for
 177 **microbial abundance with clinical covariate models (blue) vs clinical covariate-only models (grey) and (b)** mean
 178 **AUROC scores for each ML method. In both (a) and (b) error bars denote standard deviations. Significance: * ≤**
 179 **0.01, ** ≤ 0.001, *** ≤ 0.0001. (c)** Violin density plots of AUROC scores for microbial abundance with clinical
 180 **covariate models vs clinical covariate-only models. Lines connecting points (light grey) represent score pairs from**
 181 **same train-test split on the data. Mean AUROC scores are shown as red dots connected by red lines. (d)** Mean ROC
 182 **curves (blue) and (e) precision-recall (PR) curves (purple) for microbial abundance with clinical covariate models vs**
 183 **clinical covariate-only models (grey). Mean AUROC and AUPRC scores shown in legends and shaded areas denote**
 184 **standard deviations. Significance for the prediction improvement over clinical covariate-only models was calculated**
 185 **using a two-sided Wilcoxon signed-rank test and adjusted for multiple testing using the Benjamini-Hochberg**
 186 **method with adjusted p-values shown at top of violin plots in (c). In (c-e) results for the modeling method that had**
 187 **the most significant Wilcoxon signed-rank test are shown.**



188 **Figure 4. Performance of gene expression drug response prediction models in the six cancer-drug**
 189 **combinations where models performed better than clinical covariates alone. (a) Mean AUROC scores for gene**
 190 **expression with clinical covariate models (orange) vs clinical covariate-only models (grey) and (b) mean AUROC**
 191 **scores for each ML method. In both (a) and (b) error bars denote standard deviations. Significance: * ≤ 0.01 , ** \leq**
 192 **0.001, *** ≤ 0.0001 . (c) Violin density plots of AUROC scores for gene expression with clinical covariate models**
 193 **vs clinical covariate-only models. Lines connecting points (light grey) represent score pairs from same train-test**
 194 **split on the data. Mean AUROC scores are shown as red dots connected by red lines. (d) Mean ROC (orange) and**
 195 **(e) precision-recall (PR) curves (green) for gene expression with clinical covariate models vs clinical covariate-only**
 196 **models (grey). Mean AUROC and AUPRC scores shown in legends and shaded areas denote standard deviations.**
 197 **Significance for the prediction improvement over clinical covariate-only models was calculated using a two-sided**
 198 **Wilcoxon signed-rank test and adjusted for multiple testing using the Benjamini-Hochberg method with adjusted p-**
 199 **values shown at top of violin plots in (c). In (c-e) results for the modeling method that had the most significant**
 200 **Wilcoxon signed-rank test are shown.**



201 **Figure 5. Comparison of drug response model top-ranked selected features by each ML method.** For each drug
 202 response model, we selected the two best ML methods by significance for the prediction improvement over their
 203 respective clinical covariate-only model. **(a, c)** Venn diagrams for microbial abundance **(a)** or gene expression **(c)**
 204 models comparing the number of features individually selected by each ML method, and the intersection of the two
 205 ML methods. **(b, d)** Spearman rank correlation plots for microbial abundance **(b)** or gene expression **(d)** models
 206 showing that the median rank of features (among the 100 model instances in which the feature was selected)
 207 often correlated between the two most significant ML methods. The best method is shown on the x-axis, the second best
 208 on the y-axis.

209 **Combining tumor microbial abundance and gene expression features adds a modest** 210 **predictive improvement in some cancers**

211 Finally, we investigated if models built from combining microbial abundance and gene
212 expression features would result in an improvement in predictive power over their corresponding
213 single data type models. Combining data types resulted in a modest predictive improvement in
214 only three prognosis models: SARC OS, STAD PFI, and thymoma (THYM) OS
215 (**Supplementary Fig. 4a**). Although this improvement was not statistically significant in terms
216 of C-index score, the $AUC^{C/D}(t)$ metric showed a clear improvement in prognostic predictive
217 power for these models, where the improvement in mean $AUC^{C/D}(t)$ over the entire time range
218 after diagnosis was ≥ 0.025 compared to their respective single data type models. We also found
219 five combined data type drug response models which performed significantly better than clinical
220 covariates alone, although none of these models reached statistical significance when compared
221 to their respective single data type models in terms of improvement in AUROC score, but one of
222 these models, for BLCA cisplatin, did show an improvement in AUROC ≥ 0.025 compared to its
223 corresponding single data type models (**Supplementary Fig. 4b-c**).

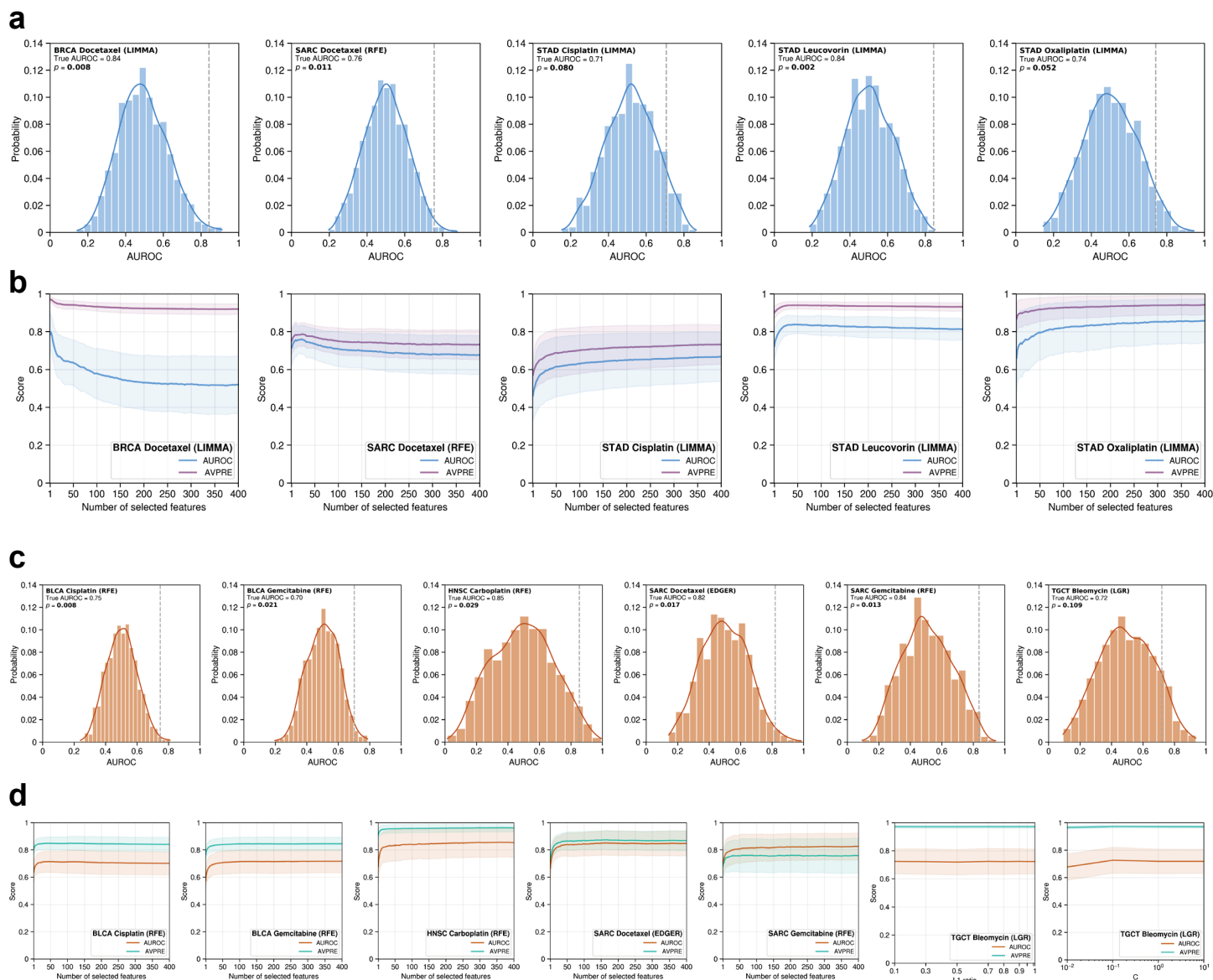
224 **Evaluating the robustness of drug response models**

225 Some of the TCGA drug response cohorts used in our study were of limited size and this
226 could have an impact on the robustness of our analysis (see **Supplementary Data 1** for cohort
227 size information). To study this issue further, we evaluated the significance of model scores
228 using a class label permutation test. We shuffled dataset class labels 1000 times and each time
229 ran the outer CV procedure on the permuted dataset, where for each CV iteration we fit a model
230 instance and calculated an AUROC score. We then calculated a p-value from the fraction of
231 permuted scores that were greater than or equal to the true score. Three of the five microbial
232 abundance drug response models that were reported above to have perform significantly better

233 than clinical covariates alone had a permutation test p-value < 0.05 and the remaining two, for
234 stomach adenocarcinoma (STAD) cisplatin and oxaliplatin, had p-values < 0.08 (**Fig. 6a**).
235 Permutation test scores and significance for microbial abundance models were similar regardless
236 of the modeling method used (**Supplementary Fig. 5a**). Five of the six gene expression drug
237 response models that performed significantly better than clinical covariates alone had a
238 permutation test p-value < 0.05 (**Fig. 6c**). Again here, permutation test scores and significance
239 were similar regardless of the modeling method used (**Supplementary Fig. 6a**). The testicular
240 germ cell tumor (TGCT) bleomycin gene expression model did not quite reach significance,
241 though it is worth mentioning that for the edgeR feature selection and L2 logistic regression
242 modeling method it was close ($p = 0.077$).

243 We further evaluated the robustness of our significant drug response models by examining
244 the effect that the number of selected features had on model performance. During the
245 hyperparameter grid search and tuning that occurred in the nested inner CV during each model
246 instance fitting, scores for every combination of hyperparameter setting and inner CV
247 train/validation fold were saved (see Methods for full details). We plotted how these scores were
248 affected by the hyperparameters that controlled feature selection. Our decision to conservatively
249 limit the feature selection search space in our drug response models to a maximum of 400 best
250 scoring features, to reduce model complexity and the possibility of overfitting, appeared
251 sufficient, as scores for our significant models reached a maximum or leveled off well within this
252 search range (**Fig. 6b & d**). In the five microbial abundance models, predictive power was driven
253 by a small number of features in three models, where selecting more features did not contribute
254 to additional predictive power or it added noise (**Fig. 6b**). Even in the remaining two models,
255 most of the predictive power was driven by the top 50 to 100 features. In the six gene expression
256 models, this finding was even more stark, where all the predictive power was achieved by a

257 small number of features in each model (**Fig. 6d**). In all the significant models from both data
258 types, the variance in scores was not significantly affected by the number of selected features and
259 feature-to-sample ratio within our chosen hyperparameter search range. As with the permutation
260 test results, we found the effect that the number of selected features had on model performance
261 was similar regardless of the feature selection or modeling method used (**Supplementary Fig.**
262 **5b, 6b**). In summary, these two comprehensive analyses suggest that the significant cancer-drug
263 response combinations found in this study and the most important features inferred from their
264 models represent a potentially real and robust biological signal.



265 **Figure 6. Evaluation of drug response model robustness.** Model significance and robustness was further
 266 evaluated using a class label permutation test and examination of the effect feature selection had on model
 267 performance. Results for the modeling method which had the most significant Wilcoxon signed-rank test are shown.
 268 **(a, c)** Permutation test result histograms and significance for microbial abundance **(a)** or gene expression **(c)**
 269 showing the distribution of permutation mean AUROC scores. True mean AUROC score shown as dotted vertical
 270 grey line and kernel density estimate shown as a curve over the histogram. **(b, d)** Curves showing the effect that
 271 model hyperparameters which control the number of selected features had on mean AUROC and average precision
 272 (AVPRE) scores during hyperparameter grid search across all 100 model instances for microbial abundance **(b)** or
 273 gene expression **(d)** models. Shaded areas denote standard deviations.

274 **Feature analysis reveals a wide range of predictive microbial genera**

275 To learn more about the most predictive features, we determined the top microbial genera
276 and top genes (**Supplementary Data 2**) selected by each of the significantly predictive microbial
277 abundance and gene expression models, respectively, according to their selection frequency and
278 model coefficients across the 100 model instances from each analysis. There were 428 distinct
279 microbial genera appearing in at least one prognosis or drug response model. Of these 428
280 genera, 160 were individually significantly predictive of prognosis or drug response by a
281 Wilcoxon test, indicating that the other genera were significantly predictive in combination. The
282 median number of genera selected per model was 52, with a minimum of 3 (BRCA docetaxel)
283 and a maximum of 78 (STAD cisplatin). Of the 428 genera, 95 were selected in more than one
284 model and only 13 were selected in more than two models. This is consistent with the
285 observation of Nejman et al.²² that the tumor microbiome is tumor type specific. The predictive
286 genera we found span all non-eukaryotic domains of life, in total encompassing 365 bacterial, 17
287 archaeal, and 46 viral genera (**Supplementary Data 2**).

288 **Discussion**

289 In summary, we find that the microbial abundance estimates generated by Poore et al.¹⁰ are
290 predictive of cancer patient prognosis and response to chemotherapy in a subset of tumor types,
291 survival outcomes, and treatments. Machine learning methods, such as those applied in this
292 study, are not able to infer causality, but only inform on the positive or negative predictive
293 associations covariates have with the response variable. The potential causal role that those
294 covariates may play in determining patient prognosis or drug response can only be ascertained
295 via dedicated mechanistic studies. Overall, in terms of the number of significant models, based
296 on their cross-validated C-index or AUROC scores and improvement over clinical covariates

297 alone, the tumor microbiome is considerably less predictive than the tumor human transcriptome
298 at predicting patient prognosis, but notably, performs similarly to gene expression at predicting
299 chemotherapy response and in mostly different cancer-drug combinations. Our investigation
300 motivates future studies investigating the role of the tumor microbiome in predicting the
301 response to targeted therapies and immunotherapies.

302 There are also some limitations to our current study. As we described previously, some
303 TCGA drug response cohorts were of limited size or had relatively few responder or non-
304 responder cases within these cohorts and this could have an impact on the interpretability of the
305 results. Vabalas et al.²³ conducted a literature review of ML algorithm validation of high-
306 dimensional biological data models with limited sample size and performed their own
307 independent simulation analyses evaluating different techniques. They found that, consistent
308 with previous literature, nested CV was the optimal validation method and gives unbiased
309 performance estimates regardless of sample size. They also found that performing feature
310 selection and other model development steps (e.g., normalization, outlier removal) fully within
311 the inner nested CV is essential to avoid overfitting and to produce unbiased results, and that
312 hyperparameter tuning should ideally also be performed in nested fashion. Finally, they found
313 that performing an adequate number of CV folds was important to reduce bias. Our analyses
314 have followed their observations and recommendations, employing them at every level of model
315 development and evaluation, including additional techniques not reviewed in their work (see
316 Methods for full details).

317 There are further limitations to this study inherited from limitations in the data, originally
318 raised by Poore et al.¹⁰ First, the study was retrospective, using existing data from the TCGA. As
319 such, it did not involve any specific protocols to capture microbial reads or to control for
320 contamination. Second, decontamination of such retrospective data is a highly involved and

321 dataset-specific process, which they made great effort to validate. Poore et al. conclude from this
322 validation that the retrospective study of TCGA was successful, and that similar retrospective
323 studies would be valuable. A third point, which they touch on briefly, is that the protocols that
324 were used have limitations with respect to capturing microbial reads and cannot distinguish if the
325 source of microbial reads is intracellular or extracellular, or alive or dead when the sample was
326 taken. Poore et al. suggest, correctly we believe, that additional protocols need to be developed
327 for prospective studies.

328 Accepting the limitations of the study, we observed certain trends. Proteobacteria and
329 Firmicutes were the most frequent phyla identified as predictive features (**Supplementary Data**
330 **2**), followed by Actinobacteria and Bacteroidetes. Among viruses, Herpesvirales were the most
331 frequent. More microbial genera were negatively predictive of drug response or prognosis than
332 those positively predictive (negative for 306/537 features; two-sided binomial test p-value =
333 0.0014). Firmicutes reversed this trend, being more often positively predictive (positive for 49/82
334 features, two-sided Fisher's exact test p-value = 0.0036; **Supplementary Table 2**).

335 Further examining the predictive features of our significant models and their cancer types, we
336 found that several genera of Firmicutes were predictive of OS in CESC, including genera of
337 *Lactobacillales* were found to be negatively predictive of survival. We also found that the genus
338 *Chlamydia* had an even stronger negatively predictive association with OS in CESC. Notably,
339 though CESC is known to often arise from HPV infection, the presence of other microbial
340 species, in particular the genera *Chlamydia* and *Lactobacillales*, have been reported in the
341 literature to be associated with the risk of developing CESC^{24,25}.

342 Our prognosis analysis results were different than two recent reports^{26,27} which found some
343 intratumor microbes that were potentially correlated with prognosis in three TCGA cancers. We

344 did not find that the Poore et al. tumor microbial abundances estimated from TCGA were
345 predictive of OS or PFI in these three cancers using our data-driven, regularized ML
346 computational approach. A few important possible reasons for this difference in results are that
347 different source data and methods were used to perform prognosis analysis compared to these
348 studies. Gnanasekar et al.²⁶ analyzed the THCA cohort by tumor subtype, they used harmonized
349 and normalized GDC TCGA data instead of legacy TCGA followed by normalization and batch
350 effect correction as in Poore et al., they only used RNA-seq data instead of WGS and RNA-seq
351 data, they applied different methods for extraction of microbial reads and decontamination, and
352 finally they did not perform any direct analysis of correlation of their derived microbial
353 abundances with survival outcomes. Dohlman et al.²⁷ analyzed colorectal cancers (colon
354 (COAD) and rectum (READ) adenocarcinomas) also using harmonized and normalized GDC
355 TCGA data, they used WGS and whole exome sequencing (WXS) data instead of WGS and
356 RNA-seq, they also used different methods for extraction and decontamination of microbial
357 reads, and finally they also applied classical univariate statistics on their entire data to infer
358 correlation with overall survival (OS). While we believe the use of harmonized GDC TCGA data
359 is superior to legacy TCGA, Poore et al. applied robust computational methods to remove
360 technical variation from legacy TCGA data and validated that their approach was effective. We
361 also applied additional filters of TCGA samples to further remove technical variation. We also
362 believe that, in general, applying classical univariate statistics on the entire data to find
363 correlations has the potential to overfit the specific dataset and it does not consider the
364 multivariate nature of high-dimensional biological data like intratumor microbial abundances. A
365 data-centric, multivariate, and regularized ML approach focused on fitting models on training
366 data and evaluating on unseen test data has potential to generalize better and discover whether
367 features are potentially predictive of and correlated with the response variable, such as survival
368 outcomes or drug response.

369 Looking at our drug response model results, in STAD, tumor microbial abundances were
370 predictive of response to three different drugs: cisplatin, leucovorin, and oxaliplatin. The genus
371 *Helicobacter* was a quantified microbial abundance feature in the Poore et al. dataset although
372 notably, even though it is well established that patients infected with *H. pylori* have an increased
373 risk of developing gastric cancer²⁸, *Helicobacter* was not identified as a predictive feature of
374 drug response in our STAD models. This finding is in line with recent research indicating
375 reduced microbial diversity, decreased abundance of *H. pylori*, and enrichment of other mostly
376 commensal bacterial genera in gastric carcinoma²⁹. Instead, in STAD we found that known
377 opportunistic bacteria *Cedecea* and *Sphingobacterium* were both strongly negatively predictive
378 of leucovorin response, *Sphingobacterium* was strongly negatively predictive of cisplatin
379 response, and the opportunistic bacteria *Rouxiiella* was strongly negative predictive of oxaliplatin
380 response. *Cedecea* and *Sphingobacterium* have been implicated in bacteremia in
381 immunocompromised individuals in rare cases, including cancer^{30,31,32,33}. As dysbiosis is
382 frequent in stomach cancer^{34,35}, and considering the mechanism of action of leucovorin, it may
383 be of interest to study whether organisms from these two genera may sequester or prevent the
384 bacterial production of folinic acid³⁶.

385 We found three microbial genera whose abundances were strongly associated with breast
386 cancer response to doxetaxel. Indeed, the involvement of the microbiome in breast cancer
387 (BRCA)^{22,37} has recently received considerable attention. In BRCA, we found that the genus
388 containing Epstein-Barr virus (EBV) was negatively associated with response to docetaxel,
389 which is concordant with previous findings that EBV is associated with chemoresistance to
390 docetaxel in gastric cancer³⁸. Interestingly, Cyanobacteria were predictive features in several
391 cancers in our study and we identified a genus of Cyanobacteria as predictive of response to
392 docetaxel in BRCA. Notably, the presence of Cyanobacteria in BRCA was recently confirmed

393 by Nejman et al.²² by 16S-rRNA sequencing. While the genus we identified, *Raphidiopsis*, a
394 planktonic Cyanobacteria that produces toxins harmful to human health and found in freshwater,
395 is possibly a taxonomic identification error in the original microbial abundance estimates, our
396 findings may point to a related genus under the recently discovered clade Melainabacteria of
397 Cyanobacteria³⁹, which is present in humans. Though Melainabacteria are difficult to culture, we
398 believe that confirmation of the relationship between BRCA to response to docetaxel and
399 Melainabacteria should be tested, and a first step would be to confirm our computationally
400 derived findings in a dedicated 16S-rRNA analysis.

401 Interestingly, in sarcoma (SARC), among the most predictive microbial features we found
402 the genus *Lactococcus* to be positively associated with response to docetaxel. *Lactococcus*
403 contains species that can sometimes cause opportunistic infections in humans, as *Lactococcus*
404 are similar to *Streptococcus* and formerly belonged to that genus. The result that this genus was
405 positively associated with response in our model initially appeared counterintuitive, although
406 while the use of therapeutic bacteria as antitumor agents has not been an extensively studied
407 field, there have been some limited findings in the literature that suggest the use of
408 bacteriotherapy as anticancer agents⁴⁰. Historically, the intentional use of the toxins of various
409 *Streptococcus* species showing significant antitumor activity in SARC has been
410 documented^{41,42,43}. One possible testable explanation for some microbes being strongly positive
411 predictive of docetaxel response in our model is that they might produce some extracellular
412 products or toxins that could work as an adjuvant to the chemotherapy.

413 In summary, while these findings and others reported in this study are computationally
414 derived associations, we believe that they can serve as leads for further experimental studies of
415 the role of microbial species in modulating patients survival and drug response, potentially by
416 metabolizing drug levels in the tumor microenvironment as suggested above, or by altering the

417 immune response, either by changing the levels of specific immunometabolites or by having the
418 tumors present specific bacterial antigens⁴⁴.

419 **Methods**

420 **Data retrieval and processing**

421 Normalized and batch effect corrected microbial abundance data for 32 TCGA tumor types
422 were downloaded from the online data repository referenced in Poore et al.¹⁰
423 (ftp://ftp.microbio.me/pub/cancer_microbiome_analysis). Specifically, the “Kraken-TCGA-
424 Voom-SNM-Plate-Center-Filtering-Data.csv” microbial abundance data file and adjoining
425 “Metadata-TCGA-Kraken-17625-Samples.csv” metadata file were used as the starting input for
426 further data processing.

427 We first filtered the data for primary tumor samples (TCGA “Primary Tumor” or “Additional
428 - New Primary” sample types). Poore et al. generated microbial abundances from all the
429 available WGS and RNA-seq data in legacy TCGA (after some quality filters), which frequently
430 contained replicate WGS and RNA-seq data for the same case and sample type. It was common
431 in legacy TCGA to increase WGS sequencing coverage by performing an additional sequencing
432 run from the same sample and these secondary runs typically had a much lower number of reads
433 and coverage compared to their corresponding primary sequencing runs. When comparing the
434 normalized and batch effect corrected read counts between these WGS runs, we found that
435 microbial abundance data which came from lower coverage secondary runs could be
436 substantially different from abundances derived from the larger primary sequencing runs.
437 Therefore, we excluded microbial abundance data which came from secondary runs. In addition,
438 legacy TCGA commonly contained data for the same samples analyzed using different
439 computational pipeline versions. We excluded replicate microbial abundance data from older

440 TCGA analysis pipeline versions if a replicate from a newer version existed. After the above
441 filters, the Poore et al. data went from 17,625 samples and 10,183 unique cases to 12,111
442 samples and 9,812 unique cases (comprising of 1,944 WGS samples from 1,904 unique cases
443 and 10,167 RNA-seq samples from 9,745 unique cases).

444 TCGA gender, age at diagnosis, and tumor stage demographic and clinical data and as well
445 as primary tumor RNA-seq read count data for the 32 TCGA tumor types included in our study
446 were obtained from the NCI Genomic Data Commons (GDC Data Release v29.0) using the R
447 Bioconductor package GenomicDataCommons. TCGA GENCODE v22 gene annotations were
448 obtained from the GDC data portal and Ensembl Gene v98 using the R package rtracklayer and
449 R Bioconductor packages AnnotationHub and ensemblDb. The downloaded GDC primary tumor
450 cohort with RNA-seq read count data comprised of 9,735 samples from 9,680 unique cases.
451 There were 68 cases at the GDC which had missing age of diagnosis but existing values in the
452 Poore et al. data and we chose not to exclude these data and used the Poore et al. age of diagnosis
453 values for these cases. TCGA curated survival phenotypic data⁴⁵ were obtained from UCSC
454 Xena. Cases which had both missing overall survival (OS) and progressive-free interval (PFI)
455 outcome data were excluded from survival modeling.

456 TCGA curated drug response clinical data were compiled from Ding et al.¹⁴ Our drug
457 response models used the following binary classification targets: complete response (CR) and
458 partial response (PR) were labeled as responders and stable disease (SD) and progressive disease
459 (PD) as non-responders. All TCGA samples with drug response phenotypic data were from pre-
460 treatment biopsies. Due to the limited cancer-drug combination cohort sizes in TCGA, we
461 modeled each drug individually, even if a patient received multiple drugs concurrently. If the
462 same drug was given at multiple timepoints to a patient, we only considered their first drug
463 response. We considered cancer-drug combinations that contained a minimum of 18 cases and at

464 least 4 cases per response binary class, except for STAD oxaliplatin, where we allowed a
465 minimum of 14 cases so that the gene expression dataset could be included. In total, we analyzed
466 30 cancer-drug combinations which had paired microbial abundance and gene expression data
467 that met the above thresholds. Combined feature microbial abundance and gene expression
468 datasets were created by joining data from each individual dataset which had matching TCGA
469 sample UUIDs. For some TCGA cases, data existed from multiple different aliquots per sample
470 or multiple technical runs per aliquot, therefore in these cases all combinations were joined at the
471 sample UUID level. Cross-validation sampling probability weights as well as model and scoring
472 sample weights were applied to account and adjust for any imbalance caused by the process.

473 **ML modeling**

474 Machine learning (ML) models were built using the scikit-learn⁴⁶ and scikit-survival
475 libraries^{47,48,49}. Custom extensions to scikit-learn and scikit-survival were developed to add new
476 methods and functionalities required by this project. Survival models were built using Coxnet –
477 regularized Cox regression with elastic net penalties¹¹. Coxnet models controlled for gender, age
478 at diagnosis, and tumor stage clinical prognostic covariates by including them as unpenalized
479 features in the model (Coxnet penalty factor = 0). Drug response classification models were built
480 using three different ML methods: 1) a variant of the linear support vector machine recursive
481 feature elimination (SVM-RFE) algorithm¹⁵ that we developed with a number of additional
482 features and better performance than the scikit-learn built-in version, 2) logistic regression
483 (LGR) with elastic net¹⁶ (L1 + L2) penalties and embedded feature selection, and 3) LGR with
484 an L2 penalty and limma¹⁷ (for tumor microbial and combination datasets) or edgeR^{18,19} (for
485 RNA-seq count datasets) differential abundance/expression feature scoring inside a k-best
486 wrapper feature selection method around the learning algorithm. Limma differential abundance
487 analysis was run inside the ML pipeline with default parameters except for fitting an intensity-

488 dependent trend to the prior variances and running a robust empirical Bayes procedure (eBayes
489 function parameters trend = TRUE and robust = TRUE). edgeR differential expression analysis
490 was run inside the ML pipeline with default parameters except for enabling robust estimation of
491 the negative binomial dispersion (calcDispersions function robust = TRUE) and robust
492 estimation of the prior quasi-likelihood (QL) dispersion (glmQLFit function robust = TRUE).
493 Both limma and edgeR methods scored and ranked features by differential abundance/expression
494 p-value.

495 All three drug response ML methods unconditionally included the same three clinical
496 covariates in the model as in the prognosis models by having them bypass feature selection in the
497 ML pipeline, though in drug response models, clinical covariates were modeled as L2 penalized
498 features. In SVM-RFE, clinical covariate features bypassed recursive feature elimination but
499 were always included at each RFE recursive feature elimination model fitting step as well as
500 final model refitting. To the best of our knowledge, no available comprehensive ML library in
501 python or R currently provides an elastic net LGR algorithm with the functionality to specify
502 features that can bypass embedded feature selection and be modeled with an L2 penalty (setting
503 the R glmnet penalty factor, for example, does not provide this functionality as it is not a penalty
504 factor per regularization term but a factor applied to the sum of both L1 and L2 terms). In order
505 to develop this functionality for our study, our elastic net LGR model pipeline was designed as a
506 two-level LGR, 1) an elastic net LGR and embedded feature selection on only microbial
507 abundance or gene expression features with clinical covariates bypassing this step, followed by
508 2) an L2 penalized LGR on features selected by the elastic net LGR step and the clinical
509 covariates. We know this design does not likely produce the exact same model settings and
510 results of a single-level elastic net LGR algorithm with the functionality we needed, if such an
511 implementation it existed, though we tested every drug response model through an ML pipeline

512 with elastic net LGR and no clinical feature selection bypass and found that model predictive
513 performance, feature coefficients and signs, and feature importance rankings were similar to our
514 two-level ML pipeline setup.

515 Gender was one-hot encoded and tumor stage ordinal encoded by major stage. In the final
516 cohort included in our prognosis and drug response models, 3363 out of 9708 tumor microbial
517 abundance cases (34.64%) and 3244 out of 9484 gene expression cases (34.21%) had tumor
518 stage “not reported” or BRCA stage “X”. Since missing tumor stage metadata is so prevalent in
519 TCGA, we took the approach of including these in our study and modeled missing tumor stage
520 with as neutral an ordinal encoding as possible. Looking at the distribution of reported major
521 tumor stages in our cohort, we determined that encoding missing data as an ordinal between
522 tumor stage II and III was as close to the middle of the distribution of stages in TCGA as we
523 could possibly achieve with ordinal encoding.

524 All prognosis and drug response models included the previously described feature selection
525 as well as normalization and transformation steps integrated into the ML modeling pipeline using
526 an extended version of the scikit-learn Pipeline framework. Each cancer, data type, and survival
527 or drug response target type combination was modeled individually using a nested cross-
528 validation (CV) strategy to perform model selection and evaluation on held-out test data.
529 Training data splits always underwent feature selection, normalization, and transformation
530 through the ML pipeline independently from held-out test or validation data splits before
531 learning. Models built using gene expression read count data included edgeR low count filtering,
532 weighted trimmed mean of M-values (TMM) normalization, and log counts per million (CPM)
533 transformation steps within the ML pipeline. These were developed and integrated into our
534 scikit-learn-based framework via R and rpy2. All models also included standardization of
535 features within the ML pipeline just before learning. During prediction, held-out test or

536 validation data were feature selected, normalized, and transformed through the ML pipeline
537 using the parameters learned from the training data at each pipeline step before model prediction
538 and scoring. Hyperparameter search and optimization of all model pipeline steps was performed
539 in nested fashion within the inner nested CV. All cross-validation iterators kept replicate sample
540 data per case grouped together such that data would only reside in either the train or test split
541 during each CV iteration.

542 Survival models used a stratified and randomly shuffled outer CV with 75% train and 25%
543 test split sizes that was repeated 100 times. The CV procedure stratified the splits on event status.
544 Each training set from the outer CV was used to perform hyperparameter tuning and model
545 selection by optimizing Harrell's concordance index (C-index) over a stratified, randomly
546 shuffled, 4-fold inner CV on the training set repeated 5 times. A few cancer datasets contained
547 fewer than four uncensored cases which required reducing the number of inner CV folds for
548 these models such that at least one case per fold was uncensored. The data derived from Poore et
549 al. often included more than one sample per case, and an unequal number of samples between
550 cases, therefore requiring either ML model sample weighting or CV random sampling per case.
551 The Coxnet implementation in scikit-survival does not currently support sample weights,
552 therefore our custom outer CV iterator randomly sampled one replicate sample per case during
553 each iteration, using a sampling procedure with probability weights that balanced the probability
554 that a replicate WGS- or RNA-seq-based sample was selected during each CV iteration. Model
555 selection grid search was performed on the following hyperparameters: elastic net penalty L1
556 ratios 0.1, 0.3, 0.5, 0.7, 0.8, 0.9, 0.95, 0.99, and 1, and for each L1 ratio a default alpha path of
557 100 alphas using an alpha min ratio of 10^{-2} . Alpha is the constant multiplier of the penalty terms
558 in the Coxnet objective function. Optimal alpha and L1 ratio settings were determined via inner
559 CV and a model with these settings was then refit on the entire outer CV train data split. Model

560 performance was evaluated in both inner and outer CV on each held-out validation or test data
561 split, respectively, by generating model test predicted risk scores and using these scores to
562 directly calculate C-index scores. We also evaluated and compared model predictive
563 performance for each test data split survival time period by calculating time-dependent
564 cumulative/dynamic AUCs^{12,13}.

565 Drug response models used a stratified, randomly shuffled, 4-fold outer CV that was repeated
566 25 times (i.e., 100 model instances). Each training set from the outer CV was used to perform
567 hyperparameter tuning and model selection by optimizing the area under the receiver-operating
568 characteristic curve (AUROC) over a stratified, randomly shuffled, 3-fold inner CV repeated 5
569 times. Case replicate sample weights were provided to SVM-RFE and LGR learning algorithms
570 and all model selection and evaluation scoring methods. Class weights were provided to SVM-
571 RFE and LGR learning algorithms to adjust for any class imbalance. Model selection grid search
572 was performed on the following hyperparameters: L2 penalized SVM and LGR C regularization
573 parameter from a range of 10^{-5} to 10^3 , elastic net LGR L1 ratios of 0.1, 0.3, 0.5, 0.7, 0.8, 0.9,
574 0.95, 0.99 and 1, elastic net LGR C regularization parameter from a range of 10^{-2} to 10^3
575 (microbial abundance) or from 10^{-2} to 10^1 (gene expression and combined data type), and finally
576 RFE, elastic net LGR, and limma and edgeR feature scorer k-best feature selection search range
577 from 1 to 400 top scoring microbial abundance, gene expression, or combined data type features.
578 SVM-RFE models performed a feature elimination procedure of the one worst feature per
579 recursive step for microbial abundance models (which started with 1287 features in the Poore et
580 al. data) and 5% of worst remaining features per recursive step until 1300 features were reached
581 followed by the one worst feature per recursive step for gene expression (starting with 60,483
582 features in GENCODE v22) and combined data type models (starting with 61,770 features).
583 Optimized hyperparameter settings were determined via inner CV and a model with the

584 optimized settings was then refit on the entire outer CV train data split. Model performance was
585 evaluated in both inner and outer CV on each held-out validation or test data split, respectively,
586 by AUROC, average precision (AVPRE) or area under precision-recall curve (AUPRC), and
587 balanced accuracy (BCR). AUROC was used to evaluate and select the best model and
588 optimized hyperparameter settings from the grid search.

589 Gender, age at diagnosis, and tumor stage clinical covariate-only survival models were built
590 using standard unpenalized Cox regression. Clinical covariate-only drug response models were
591 built using L2 penalized linear SVM or LGR. Models included standardization of features as part
592 of the ML pipeline. Models were trained and tested using the same outer CV iterators and
593 train/test data splits as their corresponding microbial abundance, gene expression, or
594 combination data type models. To test whether a Coxnet, SVM-RFE, or LGR microbial
595 abundance or gene expression model was significantly better than their corresponding Cox,
596 linear SVM, or LGR clinical covariate-only model, respectively, a two-sided Wilcoxon signed-
597 rank test was performed between the 100 pairs of C-index or AUROC scores between both
598 models. All raw p-values generated from the signed-rank test across survival or drug response
599 analyses from the same data type were adjusted for multiple testing using the Benjamini-
600 Hochberg (BH) procedure to control the false discovery rate (FDR), and a threshold $FDR \leq 0.01$
601 was used to determine statistical significance. To test whether a combined data type model was
602 significantly better than its corresponding microbial abundance or gene expression model, a two-
603 sided Dunn test was performed between all three groups of data type model scores. Each Dunn
604 test raw p-value was adjusted for multiple testing using the Benjamini-Hochberg (BH) procedure
605 to control the false discovery rate (FDR), and a threshold $FDR \leq 0.05$ was used to determine
606 statistical significance.

607 Permutation tests were performed by shuffling dataset class labels 1000 times and each time
608 running the outer CV procedure on the permuted dataset, where for each CV iteration we fit a
609 model instance and calculated an AUROC score, totaling 100,000 fits and scores for each model.
610 Permutation mean AUROC scores were compared to the true mean AUROC score for the model
611 and a p-value was calculated from the fraction of permutation mean scores that were greater than
612 or equal to the true mean score. A p-value ≤ 0.05 was used to determine statistical significance.
613 The Freedman-Draconis rule was used in permutation test histogram plots to compute the bin
614 width. Analysis of the effect of number of selected features on model performance was
615 performed via the hyperparameter grid search and tuning that occurred in the nested inner CV
616 during each model instance fitting, where scores for every combination of hyperparameter
617 setting and inner CV train/validation fold were saved for all model instances and used for
618 plotting.

619 **Microbial abundance model feature analysis**

620 For each analysis, 100 prognosis or drug response model instances were generated from the
621 outer CV procedure. Each model instance selected a subset of features that performed best
622 during CV and the model algorithm learned coefficients (or weights) for each feature. To select
623 microbial genera for downstream investigation from the feature results across all these model
624 instances, we proceeded as follows. First, we applied a two-sided Wilcoxon signed-rank test that
625 the mean feature coefficient rank generated by the model is shifted away from zero, and thus that
626 the genus is identifiably positively or negatively associated with survival or drug response. For
627 all Wilcoxon tests, we used the package `coin`⁵⁰, which allows exact calculation of p-values.
628 Coefficients were ignored when a genus was assigned a zero coefficient or absent from a model.
629 Second, within each model, all coefficients, ignoring the results of the Wilcoxon test, were
630 ranked by absolute magnitude. We then kept genera that were among the top 50 features in at

631 least 20% of the models and for which the Holm-adjusted, two-sided Wilcoxon signed-rank test
632 p-value was ≤ 0.01 . Having a Coxnet feature coefficient equal to zero or feature being absent
633 from an SVM-RFE or LGR model was not strong enough evidence that the genus has no effect,
634 but rather that one or more features with stronger effect were chosen. Thus, we ignored genera
635 with a zero coefficient or absent from a model when computing mean coefficient weight and
636 Wilcoxon statistics on the means.

637 For the drug response models, where three ML methods were tested, we noted the features
638 selected by individual models and the median rank the feature attained in the instances in which
639 it appeared, but further filtered the features to account for the consensus between ML models.
640 We kept features selected in any two ML model methods that individually met our criteria for
641 inclusion, ignoring features in ML models that did not meet these criteria. We then computed the
642 Spearman correlation between the median ranks attained by the features.

643 For each selected microbial feature, we tested whether it was a significantly univariate
644 feature of survival or drug response. This is a strictly different question than whether the
645 coefficient of a feature has consistent sign – sign may be consistent when used in combination
646 with other features, but the feature may not be individually predictive. For drug response models,
647 we divided individuals into responders and non-responders, and for survival data we divided
648 individuals whose survival time was greater or less than the censored median, ignoring those
649 who were lost to follow up before median time. For each cancer-test type pair, we applied a two-
650 sided Wilcoxon rank-sum test. We applied a Benjamini-Hochberg multiple hypothesis correction
651 for each cancer-test type pair and report the false discovery rate in **Supplementary Data 2**.

652 We analyzed the distribution of features, selected by the rules described above, that had
653 positive or negative signs for their mean coefficient. We used a two-sided binomial test to show

654 that selected features had significantly more negative the positive mean coefficients. We used a
655 two-sided Fisher's exact test to determine if selected genera belonging to Firmicutes had a
656 statistically significant difference in the breakdown between positive and negative mean
657 coefficients than selected features as a whole.

658 Data availability

659 All results generated from this work are available under <https://doi.org/10.5281/zenodo.5221525>.

660 Code availability

661 All code and data used to produce this work are available under

662 <https://github.com/ruppinlab/tcga-microbiome-prediction>.

663 References

- 664 1. Grossman, R. L. et al. Toward a Shared Vision for Cancer Genomic Data. *New England*
665 *Journal of Medicine* **375**, 1109–1112 (2016).
- 666 2. Ahluwalia, P., Kolhe, R. & Gahlay, G. K. The clinical relevance of gene expression based
667 prognostic signatures in colorectal cancer. *Biochimica et Biophysica Acta (BBA) - Reviews*
668 *on Cancer* **1875**, 188513 (2021).
- 669 3. Brodsky, A. S. et al. Expression profiling of primary and metastatic ovarian tumors reveals
670 differences indicative of aggressive disease. *PLoS One* **9**, e94476 (2014).
- 671 4. Liu, Y. et al. Pan-cancer analysis of clinical significance and associated molecular features of
672 glycolysis. *Bioengineered* **12**, 4233–4246 (2021).
- 673 5. Selfors, L. M., Stover, D. G., Harris, I. S., Brugge, J. S. & Coloff, J. L. Identification of
674 cancer genes that are independent of dominant proliferation and lineage programs. *Proc Natl*
675 *Acad Sci U S A* **114**, E11276–E11284 (2017).
- 676 6. Shimoni, Y. Association between expression of random gene sets and survival is evident in
677 multiple cancer types and may be explained by sub-classification. *PLoS Comput Biol* **14**,
678 e1006026 (2018).
- 679 7. Shukla, S. et al. Development of an RNA-Seq Based Prognostic Signature in Lung
680 Adenocarcinoma. *J Natl Cancer Inst* **109**, (2017).

- 681 8. Venet, D., Dumont, J. E. & Detours, V. Most random gene expression signatures are
682 significantly associated with breast cancer outcome. *PLoS Comput Biol* **7**, e1002240 (2011).
- 683 9. Milanez-Almeida, P., Martins, A. J., Germain, R. N. & Tsang, J. S. Cancer prognosis with
684 shallow tumor RNA sequencing. *Nature Medicine* **26**, 188–192 (2020).
- 685 10. Poore, G. D. *et al.* Microbiome analyses of blood and tissues suggest cancer diagnostic
686 approach. *Nature* **579**, 567–574 (2020).
- 687 11. Simon, N., Friedman, J., Hastie, T., & Tibshirani, R. Regularization Paths for Cox’s
688 Proportional Hazards Model via Coordinate Descent. *Journal of Statistical Software* **39** (5),
689 1-13 (2011).
- 690 12. Hung, H. & Chiang, C.T. Estimation methods for time-dependent AUC models with survival
691 data. *Canadian Journal of Statistics* **38** (1), 8–26 (2010).
- 692 13. Lambert, J. & Chevret, S. Summary measure of discrimination in survival models based on
693 cumulative/dynamic time-dependent ROC curves. *Statistical methods in medical research* **25**
694 (5), 2088–2102 (2016).
- 695 14. Ding Z *et al.* Evaluating the molecule-based prediction of clinical drug responses in cancer.
696 *Bioinformatics* **32**, (19): 2891-5 (2016).
- 697 15. Guyon, I., Weston, J., Barnhill, S. & Vapnik, V. Gene Selection for Cancer Classification
698 using Support Vector Machines. *Machine Learning* **46**, 389–422 (2002).
- 699 16. Zou H. & Hastie T. Regularization and variable selection via the elastic net. *J R Statist Soc B*
700 **67** (2), 301–320 (2005).
- 701 17. Ritchie, M.E. *et al.* *limma* powers differential expression analyses for RNA-sequencing and
702 microarray studies. *Nucleic Acids Res* **43** (7), e47 (2015).
- 703 18. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for
704 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140
705 (2010).
- 706 19. McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor
707 RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* **40**, 4288–4297
708 (2012).
- 709 20. Geller, L. T. *et al.* Potential role of intratumor bacteria in mediating tumor resistance to the
710 chemotherapeutic drug gemcitabine. *Science* **357**, 1156–1160 (2017).
- 711 21. Pushalkar, S. *et al.* The Pancreatic Cancer Microbiome Promotes Oncogenesis by Induction
712 of Innate and Adaptive Immune Suppression. *Cancer Discov* **8**, 403–416 (2018).
- 713 22. Nejman, D. *et al.* The human tumor microbiome is composed of tumor type-specific
714 intracellular bacteria. *Science* **368**, 973–980 (2020).
- 715 23. Vabalas A. *et al.* Machine learning algorithm validation with a limited sample size. *PLoS*
716 *One* **14** (11): e0224365 (2019).

- 717 24. Lin, D. *et al.* Microbiome factors in HPV-driven carcinogenesis and cancers. *PLoS Pathog*
718 **16** (6), e1008524 (2020).
- 719 25. Zhu, H. *et al.* Chlamydia Trachomatis Infection-Associated Risk of Cervical Cancer: A
720 Meta-Analysis. *Medicine (Baltimore)* **95** (13): e3077 (2016).
- 721 26. Gnanasekar A. *et al.* The intratumor microbiome predicts prognosis across gender and
722 subtypes in papillary thyroid carcinoma. *Comput Struct Biotechnol J* **19**, 1986-1997 (2021).
- 723 27. Dohlman *et al.* The cancer microbiome atlas: a pan-cancer comparative analysis to
724 distinguish tissue-resident microbiota from contaminants. *Cell Host Microbe* **10** (2), 281-
725 298.e5 (2021).
- 726 28. Parsonnet, J. *et al.* Helicobacter pylori infection and the risk of gastric carcinoma. *N Engl J*
727 *Med* **325**, 1127–1131 (1991).
- 728 29. Ferreira R.M. *et al.* Gastric microbial community profiling reveals a dysbiotic cancer-
729 associated microbiota. *Gut* **67** (2), 226-236 (2018).
- 730 30. Abate, G., Qureshi, S. & Mazumder, S. A. *Cedecea davisae* bacteremia in a neutropenic
731 patient with acute myeloid leukemia. *J Infect* **63**, 83–85 (2011).
- 732 31. Akinosoglou, K. *et al.* Bacteraemia due to *Cedecea davisae* in a patient with sigmoid colon
733 cancer: a case report and brief review of the literature. *Diagn Microbiol Infect Dis* **74**, 303–
734 306 (2012).
- 735 32. Koh, Y. R. *et al.* The first Korean case of *Sphingobacterium spiritivorum* bacteremia in a
736 patient with acute myeloid leukemia. *Ann Lab Med* **33**, 283–287 (2013).
- 737 33. Wu, P. *et al.* Profiling the Urinary Microbiota in Male Patients with Bladder Cancer in
738 China. *Front Cell Infect Microbiol* **8**, 167 (2018).
- 739 34. Coker, O. O. *et al.* Mucosal microbiome dysbiosis in gastric carcinogenesis. *Gut* **67**, 1024–
740 1032 (2018).
- 741 35. Castaño-Rodríguez N. *et al.* Dysbiosis of the microbiome in gastric carcinogenesis. *Sci Rep* **7**
742 (1), 15957 (2015).
- 743 36. Ogowang, S. *et al.* Bacterial Conversion of Folinic Acid Is Required for Antifolate Resistance.
744 *Journal of Biological Chemistry* **286**, 15377–15390 (2011).
- 745 37. Eslami-S, Z., Majidzadeh-A, K., Halvaei, S., Babapirali, F. & Esmaili, R. Microbiome and
746 Breast Cancer: New Role for an Ancient Population. *Front Oncol* **10**, (2020).
- 747 38. Shin, H. J., Kim, D. N. & Lee, S. K. Association between Epstein-Barr virus infection and
748 chemoresistance to docetaxel in gastric carcinoma. *Mol. Cells* **32**, 173–179 (2011).
- 749 39. Di Rienzi, S. C. *et al.* The human gut and groundwater harbor non-photosynthetic bacteria
750 belonging to a new candidate phylum sibling to Cyanobacteria. *eLife* **2**, e01102 (2013).

- 751 40. Sedighi M. *et al.* Therapeutic bacteria to combat cancer; current advances, challenges, and
752 opportunities. *Cancer Med* **8** (6), 3167-3181 (2019).
- 753 41. Coley WB. Late results of the treatment of inoperable sarcoma by the mixed toxins of
754 erysipelas and bacillus prodigiosus. *Trans Southern Surg Gynecol Ass* **18**, 197 (1906).
- 755 42. Fehleisen F. Ueber die Züchtung der Erysipelkokken auf künstlichem Nährboden und ihre
756 Übertragbarkeit auf den Menschen. *Dtsch Med Wochenschr* **8**, 553-554 (1882).
- 757 43. Busch W. Aus der Sitzung der medicinischen Section vom 13 November 1867. *Berl Klin*
758 *Wochenschr* **5**, 137 (1868).
- 759 44. Kalaora S. *et al.* Identification of bacteria-derived HLA-bound peptides in melanoma. *Nature*
760 **592** (7852), 138-143 (2021).
- 761 45. Liu, J. *et al.* An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-Quality
762 Survival Outcome Analytics. *Cell* **173**, 400-416.e11 (2018).
- 763 46. Pedregosa *et al.* Scikit-learn: Machine Learning in Python. *J Mach Learn Res* **12**, 2825-2830
764 (2011).
- 765 47. Pölsterl, S., Navab, N. & Katouzian, A. Fast Training of Support Vector Machines for
766 Survival Analysis. in *Machine Learning and Knowledge Discovery in Databases* (eds.
767 Appice, A. *et al.*) 243–259 (Springer International Publishing, 2015).
- 768 48. Pölsterl, S., Navab, N., & Katouzian, A., An Efficient Training Algorithm for Kernel
769 Survival Support Vector Machines. *4th Workshop on Machine Learning in Life Sciences*, 23
770 September 2016, Riva del Garda, Italy.
- 771 49. Pölsterl, S. *et al.* Heterogeneous ensembles for predicting survival of metastatic, castrate-
772 resistant prostate cancer patients. *F1000Res* **5**, 2676 (2017).
- 773 50. Hothorn, T., Hornik, K., Wiel, M. A. van de & Zeileis, A. Implementing a Class of
774 Permutation Tests: The coin Package. *Journal of Statistical Software* **28**, 1–23 (2008).

775 Acknowledgements

776 The results shown here are in whole or part based upon data generated by the TCGA Research
777 Network (<https://www.cancer.gov/tcga>). This research was supported by the Intramural Research
778 Program of the National Institutes of Health, National Cancer Institute and by NIH grant
779 1ZIABC011803-03. The authors would like to personally thank Christopher Buck from the NCI,
780 Pedro Milanez-Almeida and John Tsang from NIH NIAID, and Alejandro Schäffer, Welles
781 Robinson, Fiorella Schischlik, Sanju Sinha, and Sanna Madan from NCI CDSL for their

782 assistance in this project. This study utilized the high-performance computational capabilities of
783 the HPC Biowulf Linux cluster at the National Institutes of Health, Bethesda, MD
784 (<https://hpc.nih.gov>). The authors would like to thank Richard Lehr, Tim Miller, Wolfgang
785 Resch, and Steve Fellini for their assistance in running this analysis on the NIH HPC Biowulf
786 cluster. The authors would also like to thank Joel Nothman, Andreas Mueller, and Adrin Jalali
787 from the scikit-learn core development team and scikit-survival author Sebastian Pölsterl for
788 their assistance with developing extensions to their libraries. Figure 1 embedded image credits
789 (all CC 3.0, BSD 3-Clause, or Apache 2.0 license): scikit-learn.org, [scikit-](https://scikit-survival.readthedocs.io)
790 [survival.readthedocs.io](https://scikit-survival.readthedocs.io), eli5.readthedocs.io, gdc.cancer.gov,
791 github.com/Bioconductor/BiocStickers, [Fiorella Schischlik](https://www.fiorella.com), wikipedia.org.

792 **Author information**

793 **Affiliations**

794 **Cancer Data Science Laboratory (CDSL), National Cancer Institute (NCI), National**
795 **Institutes of Health (NIH), Bethesda, MD, USA**

796 Leandro Cruz Hermida, E. Michael Gertz & Eytan Ruppín

797

798 **Department of Computer Science, University of Maryland, College Park, MD, USA**

799 Leandro Cruz Hermida

800 **Contributions**

801 L.C.H., E.M.G., and E.R. designed the study. L.C.H. and E.M.G. performed all computational

802 analyses and results interpretation. L.C.H., E.M.G., and E.R. wrote the paper.

803 **Corresponding authors**

804 Correspondence to Eytan Ruppín.

805 **Ethics declaration**

806 **Competing interests**

807 All authors declare that they have no competing interests.

808 **Supplementary Tables**

809 **Supplementary Table 1.** Prognosis model results comparison to Milanez-Almeida et al.

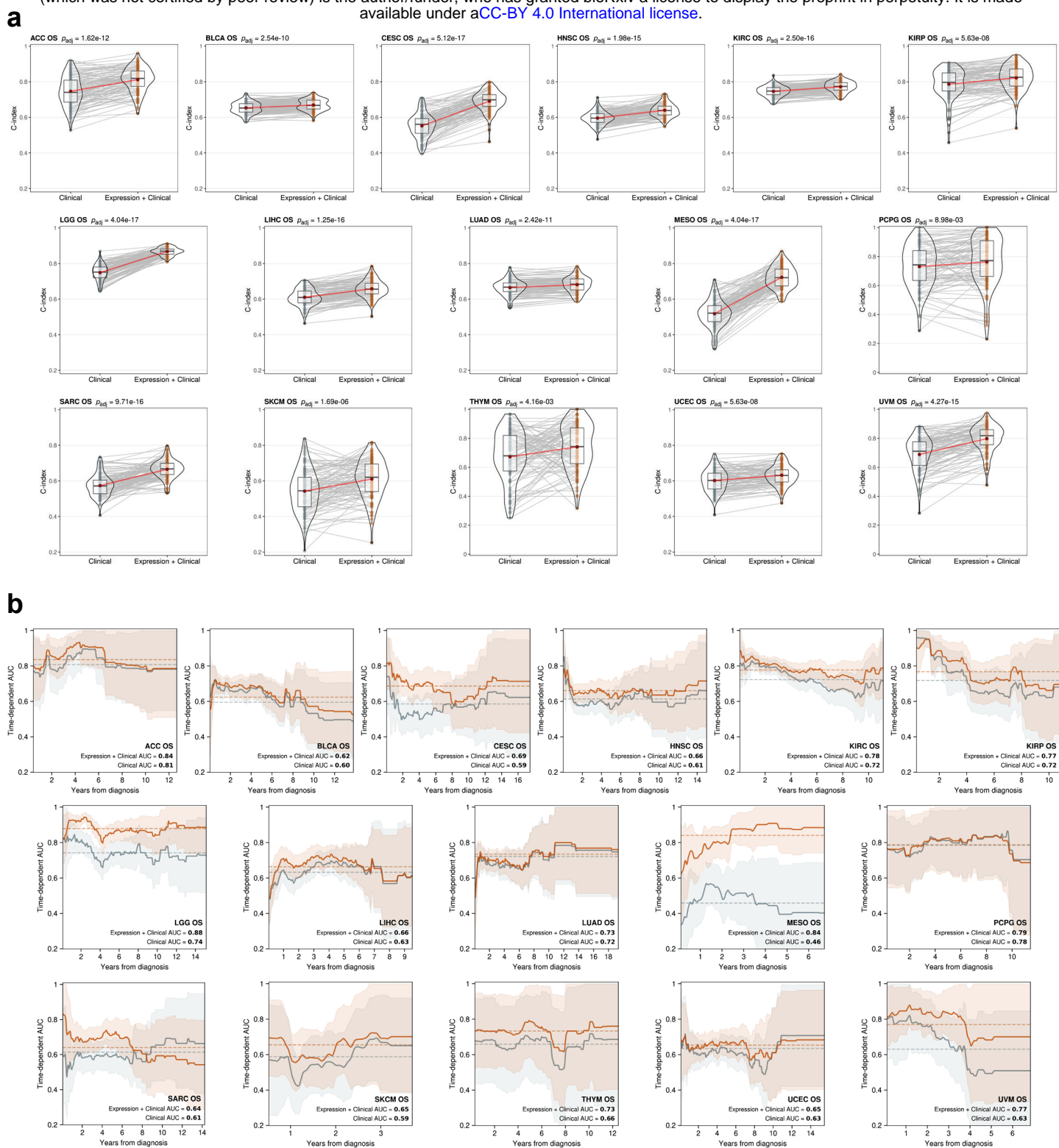
Cancer	Milanez-Almeida		Hermida & Gertz					
	Tested	Hit	Tested	Hit	Common Hit	New Hit	M-A Tested New Hit	Missing M-A Hit
ACC	OS	OS	OS, PFI	OS, PFI	OS	PFI		
BLCA	OS	OS	OS, PFI	OS, PFI	OS	PFI		
BRCA	PFI		OS, PFI	PFI		PFI	PFI	
CESC	OS		OS, PFI	OS, PFI		OS, PFI	OS	
CHOL	<i>EXCL</i>		OS, PFI					
COAD	OS		OS, PFI					
DLBC	<i>EXCL</i>		OS, PFI					
ESCA	OS		OS, PFI					
GBM	OS		OS, PFI					
HNSC	OS	OS	OS, PFI	OS	OS			
KICH	<i>EXCL</i>		OS, PFI					
KIRC	OS	OS	OS, PFI	OS, PFI	OS	PFI		
KIRP	OS	OS	OS, PFI	OS, PFI	OS	PFI		
LAML	OS	OS	<i>EXCL</i>					OS
LGG	PFI	PFI	OS, PFI	OS, PFI	PFI	OS		
LIHC	OS	OS	OS, PFI	OS, PFI	OS	PFI		
LUAD	OS	OS	OS, PFI	OS	OS			
LUSC	OS		OS, PFI	PFI		PFI		
MESO	OS	OS	OS, PFI	OS, PFI	OS	PFI		
OV	OS		OS, PFI					
PAAD	OS	OS	OS, PFI	PFI		PFI		OS

PCPG	<i>EXCL</i>		OS, PFI	OS		OS		
PRAD	PFI	PFI	OS, PFI	PFI	PFI			
READ	PFI		OS, PFI					
SARC	OS		OS, PFI	OS, PFI		OS, PFI	OS	
SKCM	<i>EXCL</i>		OS, PFI	OS, PFI		OS, PFI		
STAD	OS		OS, PFI	PFI		PFI		
TGCT	OS		OS, PFI					
THCA	PFI		OS, PFI					
THYM	PFI		OS, PFI	OS		OS		
UCEC	OS		OS, PFI	OS, PFI		OS, PFI	OS	
UCS	OS		OS, PFI					
UVM	OS	OS	OS, PFI	OS, PFI	OS	PFI		

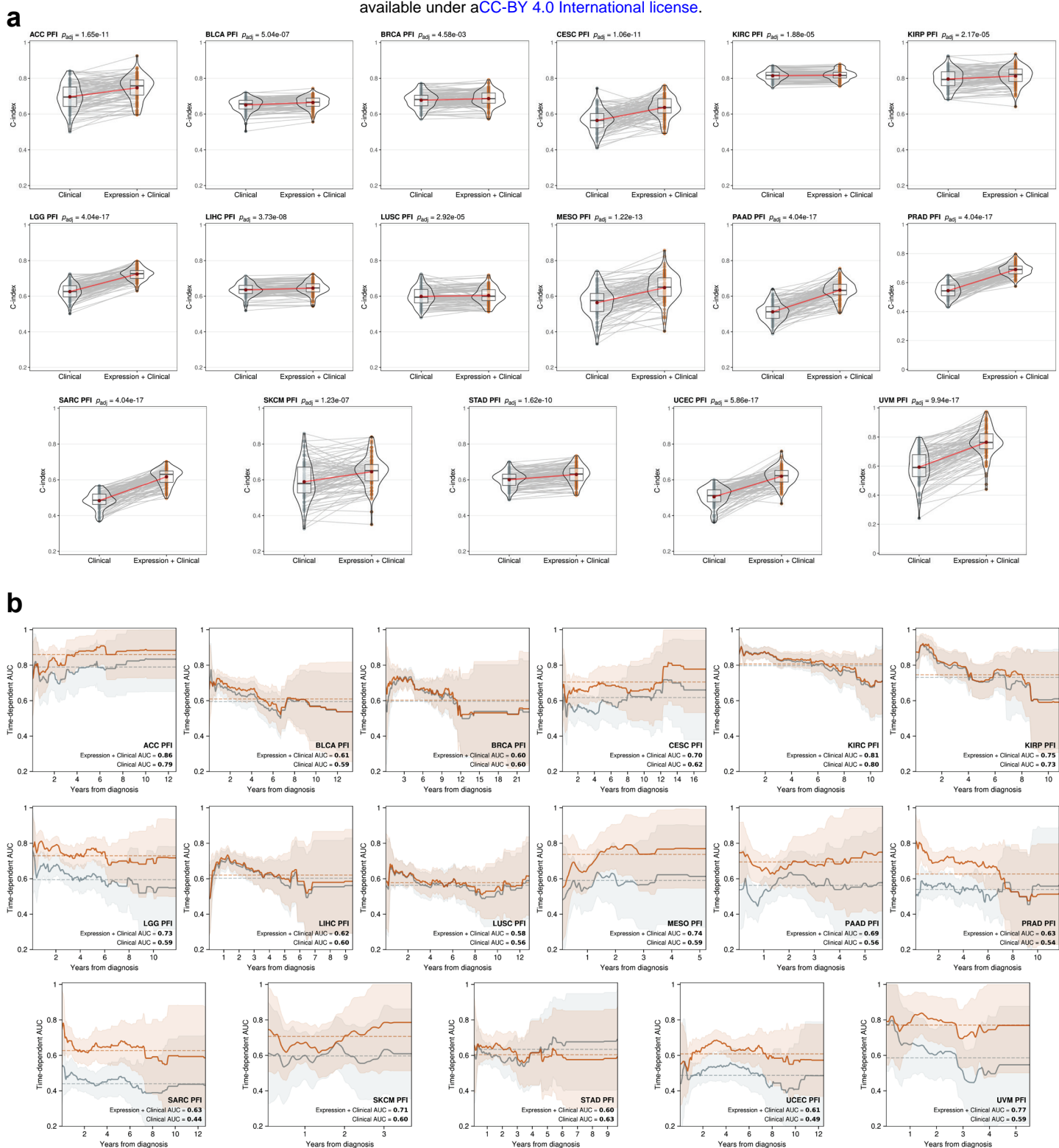
810 **Supplementary Table 2.** Per cancer, the number of times genera from the phylum Firmicutes
811 were found among the selected features, whether positively or negatively associated with drug
812 response or prognosis.

Cancer	Selected Negatively	Selected Positively	Total
ACC	8	13	21
BRCA	1	0	1
CESC	4	7	11
ESCA	0	4	4
KIRC	4	7	11
LGG	1	6	7
SARC	2	2	4
STAD	13	10	23

813 **Supplementary Figures**

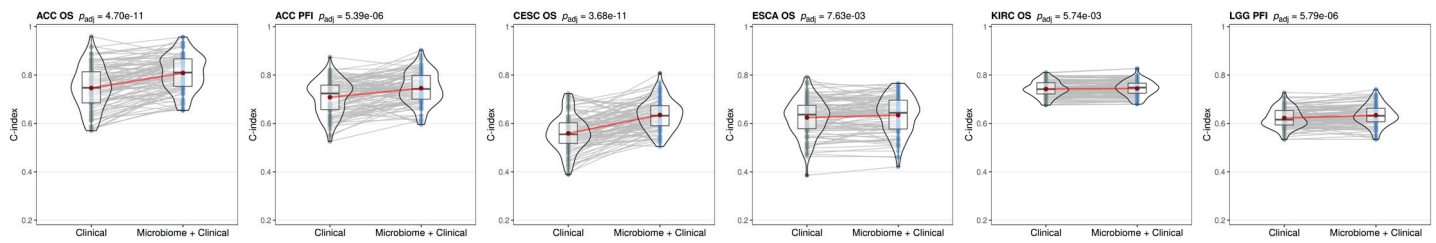


814 **Supplementary Figure 1. Performance of gene expression overall survival (OS) models in the 16 tumor types**
 815 **where gene expression adds predictive power to clinical covariates. (a)** C-index score density distributions for
 816 gene expression with clinical covariate models vs clinical covariate-only models. Lines connecting points (light
 817 grey) represent score pairs from same train-test split on the data. Mean C-index scores and connecting lines shown
 818 in red. Significance for the prediction improvement over clinical covariate-only models was calculated using a two-
 819 sided Wilcoxon signed-rank test and adjusted for multiple testing using the Benjamini-Hochberg method with
 820 adjusted p-values shown at top. **(b)** Time-dependent, cumulative/dynamic AUCs for gene expression with clinical
 821 covariate models (orange) vs clinical covariate-only models (grey) following years after diagnosis. Mean AUCs
 822 across entire test time range after diagnosis shown as a horizontal dotted line and in legends and shaded areas denote
 823 standard deviations.

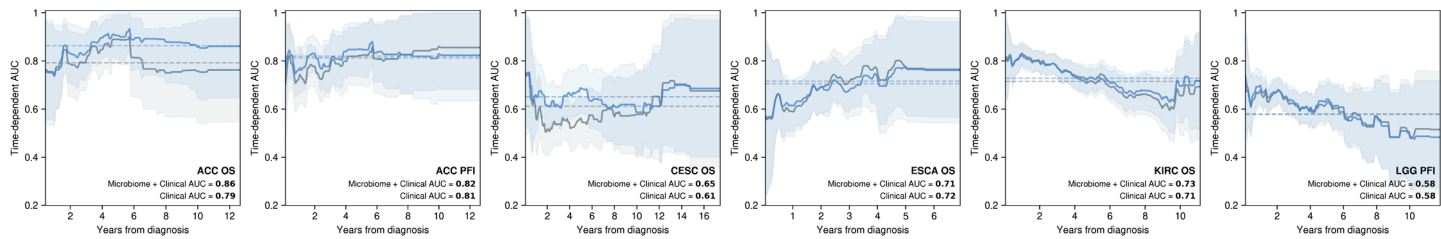


824 **Supplementary Figure 2. Performance of gene expression progression-free interval (PFI) models in the 17**
 825 **tumor types where gene expression adds predictive power to clinical covariates. (a)** C-index score density
 826 distributions for gene expression with clinical covariate models vs clinical covariate-only models. Lines connecting
 827 points (light grey) represent score pairs from same train-test split on the data. Mean C-index scores and connecting
 828 lines shown in red. Significance for the prediction improvement over clinical covariate-only models was calculated
 829 using a two-sided Wilcoxon signed-rank test and adjusted for multiple testing using the Benjamini-Hochberg
 830 method with adjusted p-values shown at top. **(b)** Time-dependent, cumulative/dynamic AUCs for gene expression
 831 with clinical covariate models (orange) vs clinical covariate-only models (grey) following years after diagnosis.
 832 Mean AUCs across entire test time range shown as a horizontal dotted line and in legends and shaded areas denote
 833 standard deviations.

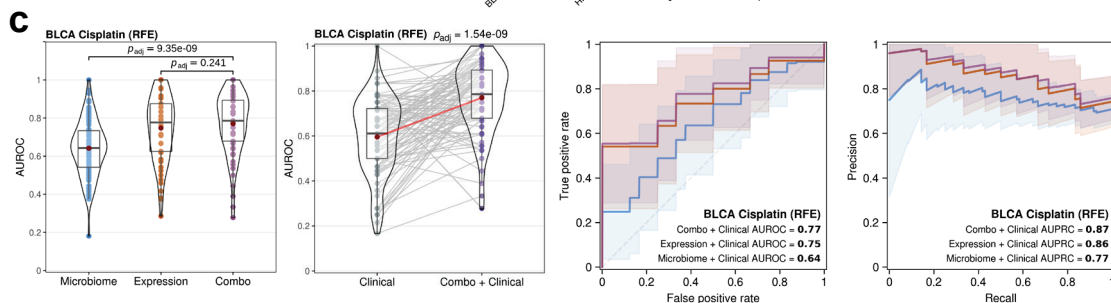
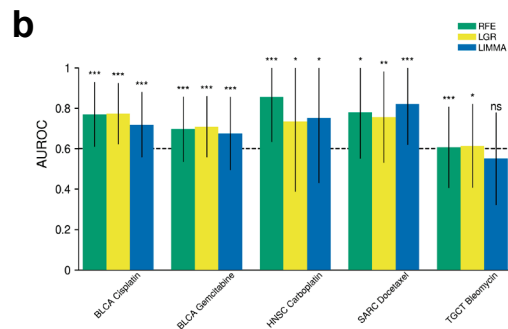
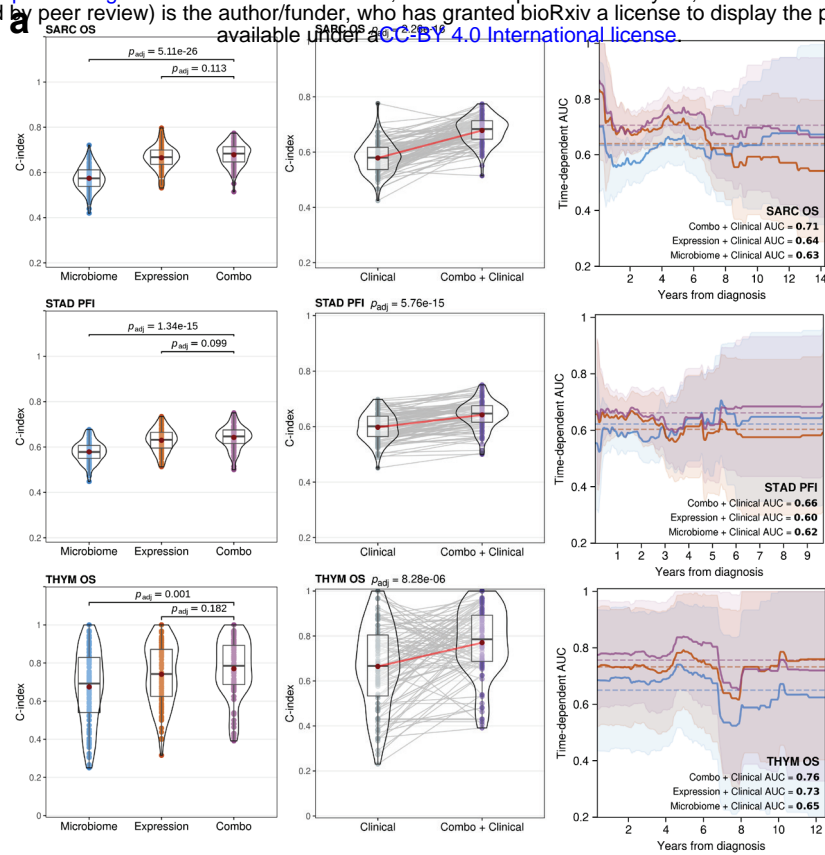
a



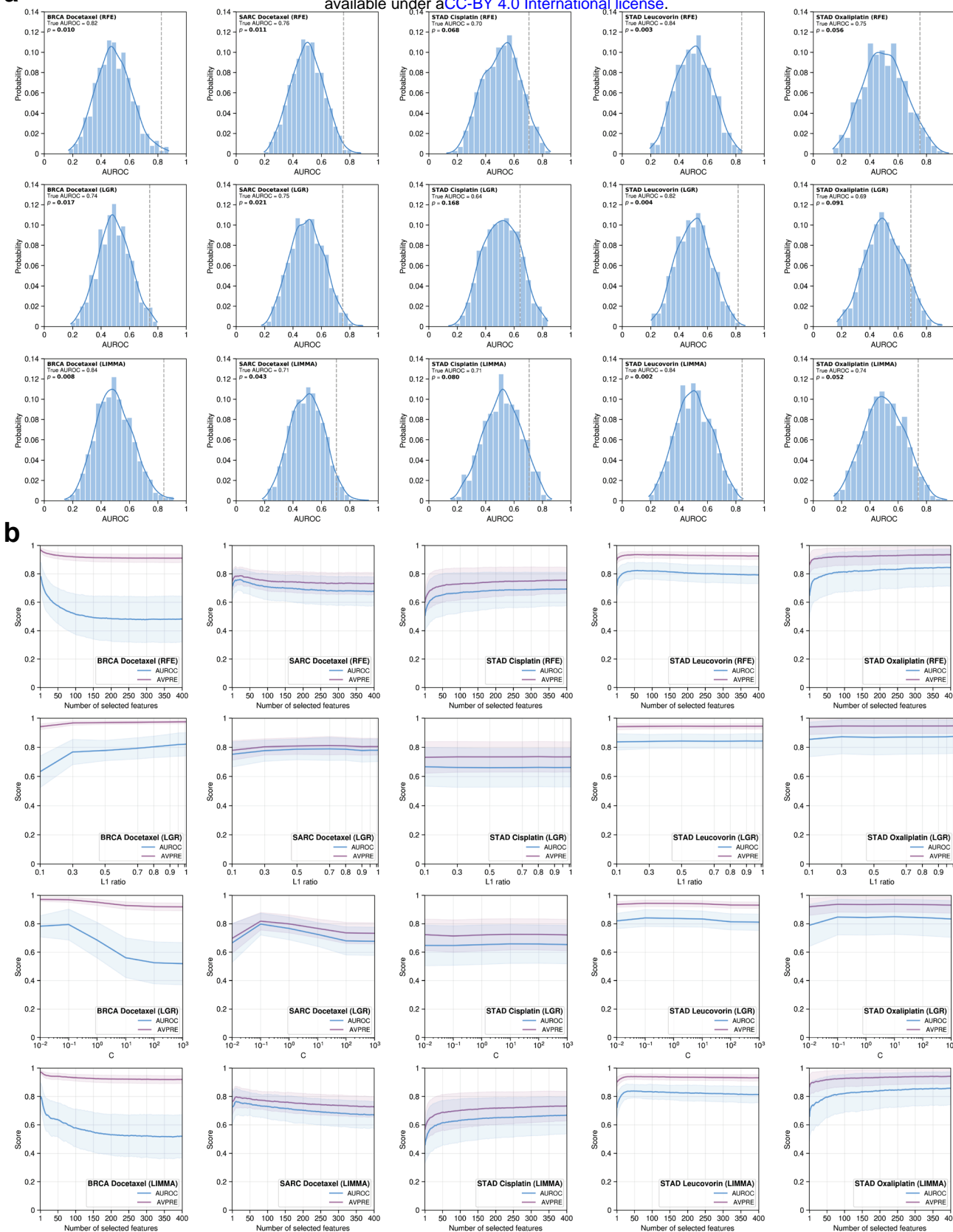
b



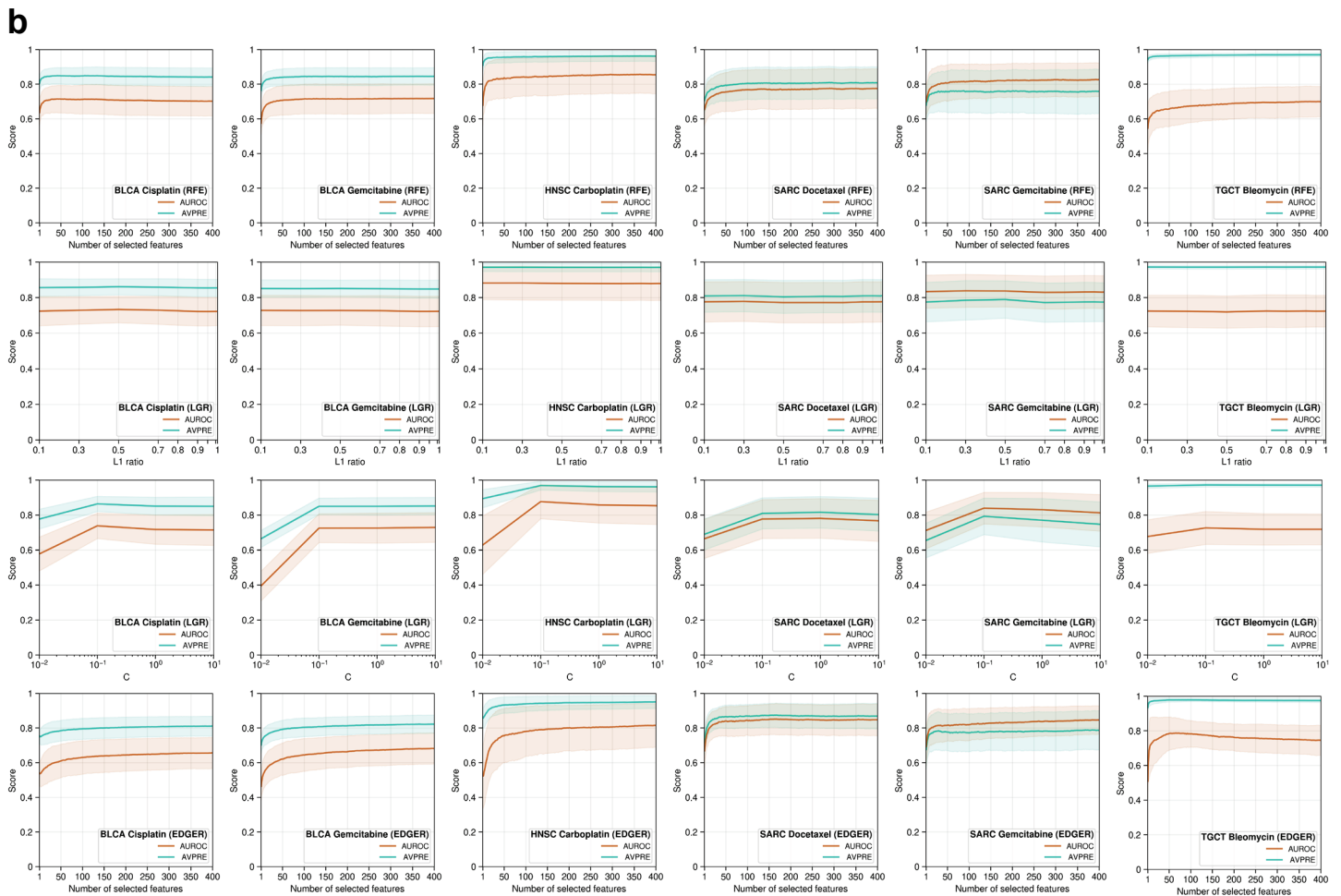
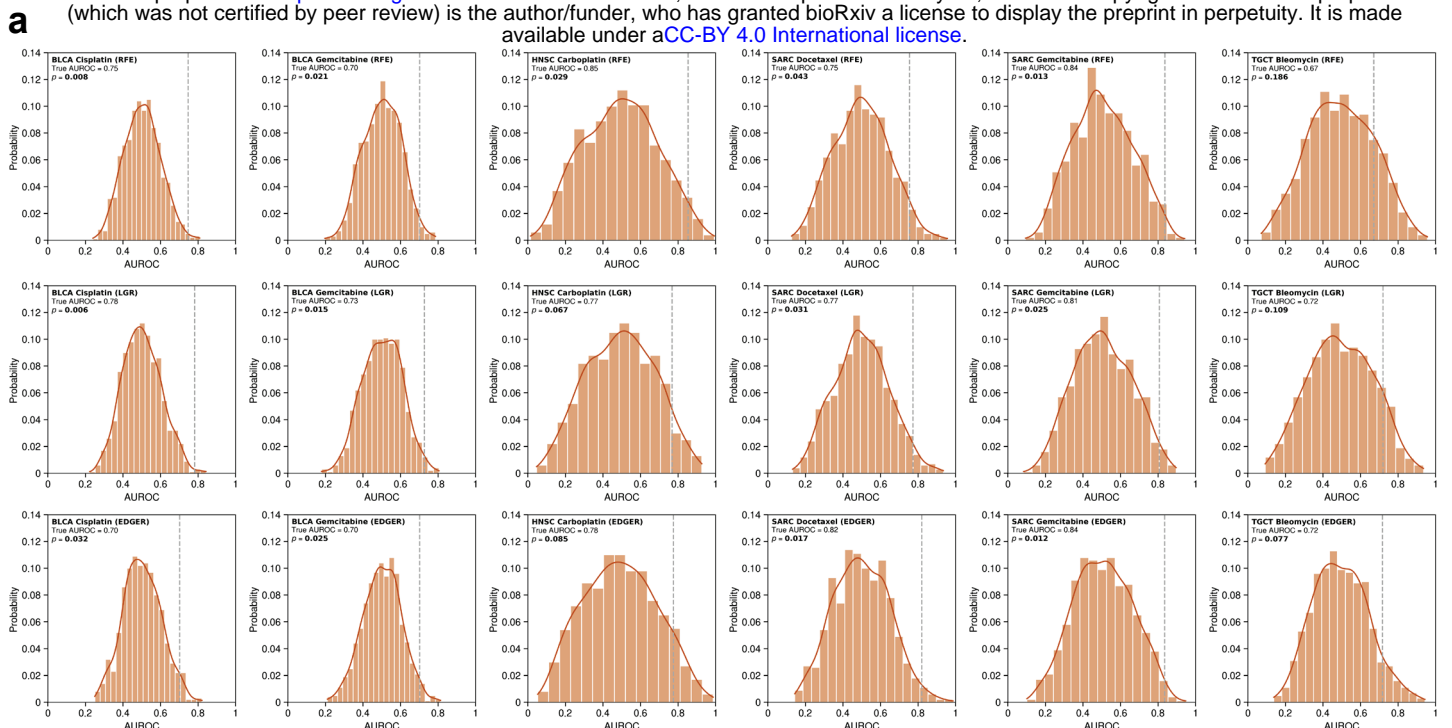
834 **Supplementary Figure 3. Performance of microbial abundance prognosis models in the five tumor types**
 835 **where microbial abundance features add predictive power to clinical covariates. (a)** C-index score density
 836 distributions for microbial abundance with clinical covariate models vs clinical covariate-only models. Lines
 837 connecting points (light grey) represent score pairs from same train-test split on the data. Mean C-index scores and
 838 connecting lines shown in red. **(b)** Time-dependent, cumulative/dynamic AUCs for microbial abundance with
 839 clinical covariate models (blue) vs clinical covariate-only models (grey) following years after diagnosis. Mean
 840 AUCs across entire test time range shown as a horizontal dotted line and in legends and shaded areas denote
 841 standard deviations. Significance was calculated using a two-sided Wilcoxon signed-rank test and adjusted for
 842 multiple testing using the Benjamini-Hochberg method with adjusted p-values shown at top of violin plots in **(a)**.



843 **Supplementary Figure 4. Performance of combined microbial abundance and gene expression models where**
 844 **combining both data types adds to predictive power. (a)** Prognosis model C-index score violin density
 845 distributions (left) between microbial abundance, gene expression, and combined data type models. Mean scores
 846 shown in red. Significance was calculated using a two-sided Dunn test and p-values were adjusted for multiple
 847 testing using the Benjamini-Hochberg method. Model C-index score violin density distributions (middle) for
 848 combined data type with clinical covariate models vs clinical covariate-only models. Time-dependent,
 849 cumulative/dynamic AUCs (right) for combined data type (purple), microbial abundance (blue), and gene expression
 850 (orange) models following years after diagnosis. **(b)** ML method drug response model mean AUROC scores where
 851 combining data types performed better than clinical covariates alone. For the combined data type drug response
 852 model closest to reach significant improvement over respective single data type models **(c)** AUROC score violin
 853 density distributions (left) comparing each data type model, AUROC score violin density distributions (middle left)
 854 comparing combined data type and clinical covariate model to clinical covariate-only model, and mean ROC and PR
 855 curves (right) for each data type model. Mean AUROC and AUPRC scores shown in panel legends and shaded areas
 856 denote standard deviations.



857 **Supplementary Figure 5. Evaluation of microbial abundance drug response model robustness.** (a) Class label
 858 permutation test result histograms and significance showing the distribution of permutation mean AUROC scores.
 859 True mean AUROC score shown as dotted vertical grey line and kernel density estimate shown as a curve over the
 860 histogram. (b) Curves showing the effect that model hyperparameters which control the number of selected features
 861 had on mean AUROC and average precision (AVPRE) scores during hyperparameter grid search across all 100
 862 model instances. Shaded areas denote standard deviations.



863 **Supplementary Figure 6. Evaluation of gene expression drug response model robustness. (a)** Class label
 864 permutation test result histograms and significance showing the distribution of permutation mean AUROC scores.
 865 True mean AUROC score shown as dotted vertical grey line and kernel density estimate shown as a curve over the
 866 histogram. **(b)** Curves showing the effect that model hyperparameters which control the number of selected features
 867 had on mean AUROC and average precision (AVPRE) scores during hyperparameter grid search across all 100
 868 model instances. Shaded areas denote standard deviations.