# BIRDMAn: A Bayesian differential abundance framework that enables robust inference of host-microbe associations

Gibraan Rahman[1,2], James T. Morton[3], Cameron Martino[1,2], Gregory D. Sepich-Poore[4], Celeste Allaband[1], Caitlin Guccione[1,2,5], Yang Chen[1,6,7], Daniel Hakim[1,2], Mehrbod Estaki[8], Rob Knight[1,9]

[1]Department of Pediatrics, University of California San Diego, La Jolla, CA, USA
[2]Bioinformatics and Systems Biology Program, University of California San Diego, La Jolla, CA, USA
[3]Biostatistics & Bioinformatics Branch, Eunice Kennedy Shriver National Institute of Child Health and Human Development, National Institutes of Health, Bethesda, MD, USA
[4]Micronoma, San Diego, CA, USA
[5]Division of Biomedical Informatics, Department of Medicine, University of California San Diego, La Jolla, CA
[6]Department of Dermatology, University of California San Diego, La Jolla, CA, USA
[7]Biomedical Sciences Graduate Program, University of California San Diego, La Jolla, CA
[8]Department of Physiology & Pharmacology, University of Calgary, Calgary, Canada
[9]Department of Computer Science and Engineering, University of California, San Diego, La Jolla, California, USA

# Abstract

Quantifying the differential abundance (DA) of specific taxa among experimental groups in microbiome studies is challenging due to data characteristics (e.g., compositionality, sparsity) and specific study designs (e.g., repeated measures, meta-analysis, cross-over). Here we present BIRDMAn (**B**ayesian **I**nferential **R**egression for **D**ifferential **M**icrobiome **An**alysis), a flexible DA method that can account for microbiome data characteristics and diverse experimental designs. Simulations show that BIRDMAn models are robust to uneven sequencing depth and provide a >20-fold improvement in statistical power over existing methods. We then use BIRDMAn to identify antibiotic-mediated perturbations undetected by other DA methods due to subject-level heterogeneity. Finally, we demonstrate how BIRDMAn can construct state-of-the-art cancer-type classifiers using The Cancer Genome Atlas (TCGA) dataset, with substantial accuracy improvements over random forests and existing DA tools across multiple sequencing centers. Collectively, BIRDMAn extracts more informative biological signals while accounting for study-specific experimental conditions than existing approaches.

# Main

Advances in sequencing technology and computational methods have enabled researchers to experimentally characterize microbiomes across wide ranges of biological conditions, including psychiatric diseases[1,2], cancer[3,4], and COVID-19[5,6]. However, as the understanding of microbial effects on human health and disease has increased, the experimental questions, hypotheses, and concomitant statistics have grown in complexity, with study designs now commonly involving longitudinal analyses[7–9], experimental interventions[10–12], and meta-analyses[7]. Although such approaches can provide mechanistic insights into the microbiome's effect(s) on the host, their conclusions are often limited by the ability to perform valid statistical analyses that are sufficiently flexible to account for the added experimental complexity.

One common but critical challenge in these contexts is when population-level heterogeneity (such as subject-to-subject variation) is confounded by technical variability. For example, samples originating from the same sequencing center will tend to be more similar to each other than those sequenced from different centers[13]. The confounding factors that may explain these differences make it difficult to determine consistent microbial biomarkers associated with biological variables or conditions of interest[8]—an effect compounded by other microbiome data difficulties, such as high sparsity, high-dimensionality, and compositionality. Moreover, statistical tools that can properly assess and account for strong structural effects while still indicating which microbes truly vary between biological conditions are limited to date[15].

Making matters more difficult, disagreement exists about how to benchmark differential abundance (DA) tools and methods. Previous efforts have commonly focused on comparing the results of hypothesis testing while accounting for the multiplicity of features through false-discovery-rate (FDR) correction[15–17]. Studies have demonstrated that tools designed for differential abundance often report contradictory results with different microbial abundances among biologically distinct sampling groups[19].

Addressing these challenges requires a more robust statistical framework for benchmarking differential abundance methods and would benefit from flexible DA modeling approaches. Thus, we developed BIRDMAn (**B**ayesian **I**nferential **R**egression for **D**ifferential **M**icrobiome **An**alysis), a flexible computational framework for hierarchical Bayesian modeling of microbiome data that simultaneously accounts for its high sparsity, high-dimensionality, and compositionality.

The Bayesian approach to statistical modeling provides unique advantages compared to frequentist solutions, such as the inclusion of prior information, uncertainty estimation of parameters, native hierarchical modeling, and edge case smoothing (e.g., estimating log fold changes when a feature is only present in one group). Implemented within the Stan programming language (commonly used for designing probabilistic models), BIRDMAn flexibly enables parameter estimation of all biological variables and non-biological covariates. These advantages allow us to demonstrate how explicitly modeling population-level effects in probabilistic BIRDMAn models increases the amount of true biological signal recovered compared to existing tools on both simulated and real-world datasets. Moreover, the BIRDMAn

78  workflow significantly lowers the barrier of entry for differential abundance methods
79  development and implementation. Additionally, to address reproducibility issues of prior DA tool
80  benchmarking, we present a novel approach that employs techniques from compositional data
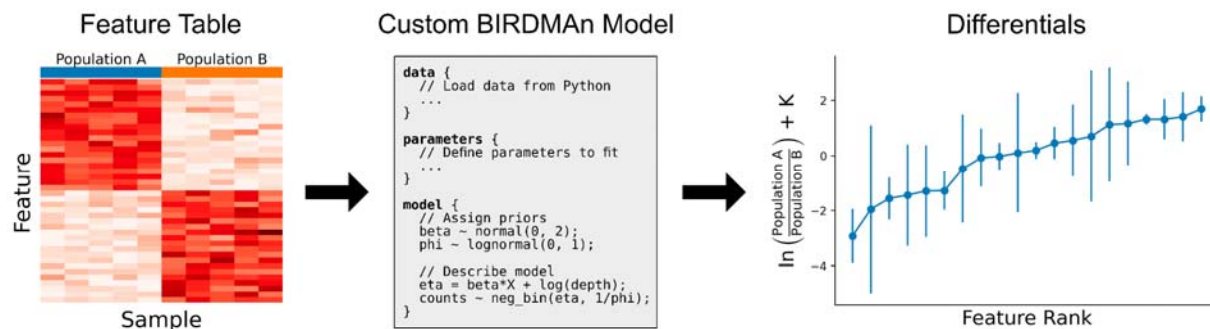81  analysis, making the comparison of tools more interpretable and statistically valid.
82



84  **Fig 1**: *Overview of BIRDMAn workflow for customizable differential abundance analysis. A table*
85  *of counts by features is modeled using Bayesian probabilistic programming, resulting in credible*
86  *intervals of the estimated parameter posterior distributions. The statistical model can be*
87  *customized using the Stan probabilistic programming language and fit using the BIRDMAn*
88  *Python interface.*

# Results

90  BIRDMAn is implemented as a Python interface to the Stan probabilistic programming
91  language, which utilizes Hamiltonian Monte Carlo sampling, one of the state-of-the-art
92  approaches for Bayesian uncertainty estimation[20]. Users can employ pre-configured model
93  designs or flexibly customize inputs to account for their specific experimental design and
94  biological questions; BIRDMAn then fits and processes these models (Fig 1). The results of
95  these analyses are the posterior distributions of the defined parameters of interest, such as log-
96  fold changes and their uncertainty given the data (see Methods).
97

98  To showcase the statistical properties of BIRDMAn models, we first leverage simulations to
99  evaluate the accuracy of estimating differential uncertainty in the context of realistic biological
100 scenarios. Then, we apply BIRDMAn models on real-world data, demonstrating superiority for
101 resolving subject-level heterogeneity in an antibiotics experiment, as well as alleviating
102 sequencing center-specific effects in a cancer genomics dataset, each while capturing
103 biologically-informative signals.

## Simulations demonstrate BIRDMAn model accuracy and precision

105 A common difficulty in benchmarking differential abundance methods is the lack of ground truth.
106 We typically do not know which microbial taxa are truly increasing or decreasing across
107 experimental conditions. To gain insights into the robustness of BIRDMAn models, we
108 performed a data-driven simulation of a case-control microbiome dataset with one binary
109 covariate, large batch effects (10 features, 10 batches, and 300 samples), data overdispersion,

110  and known differentials associated with case status (see Methods) (Fig 2a). We then used
111  BIRDMAn to estimate the model parameters for each feature and compared the Bayesian
112  posterior estimates with the true value, finding that BIRDMAn models recovered the ground
113  truth differentials with high accuracy and precision (Fig 2b) while outperforming other tools in
114  terms of root mean square error (RMSE) (Fig 2c). This highlights how BIRDMAn model
115  customization permits more accurate estimations of differentials.
116
117  One advantage of Bayesian models is that they can leverage posterior estimates to summarize
118  the uncertainty of these differentials, taking into account the sample size and the sequencing
119  depth. As expected, we show that when BIRDMAn models are fitted on larger sample sizes, the
120  uncertainty decreases, highlighting how incorporating more data, and avoiding rarefaction,
121  enables a more accurate estimation of the differentials (Fig 2d). Furthermore, we show that
122  decreasing the sequencing depth also increases the uncertainty, highlighting how rarefaction
123  could degrade parameter estimates' precisions in BIRDMAn models (Fig 2e). Since BIRDMAn
124  can handle variable sequencing depths, there is no need to perform rarefaction before model
125  fitting, which is desirable when analyzing microbiome datasets[21].
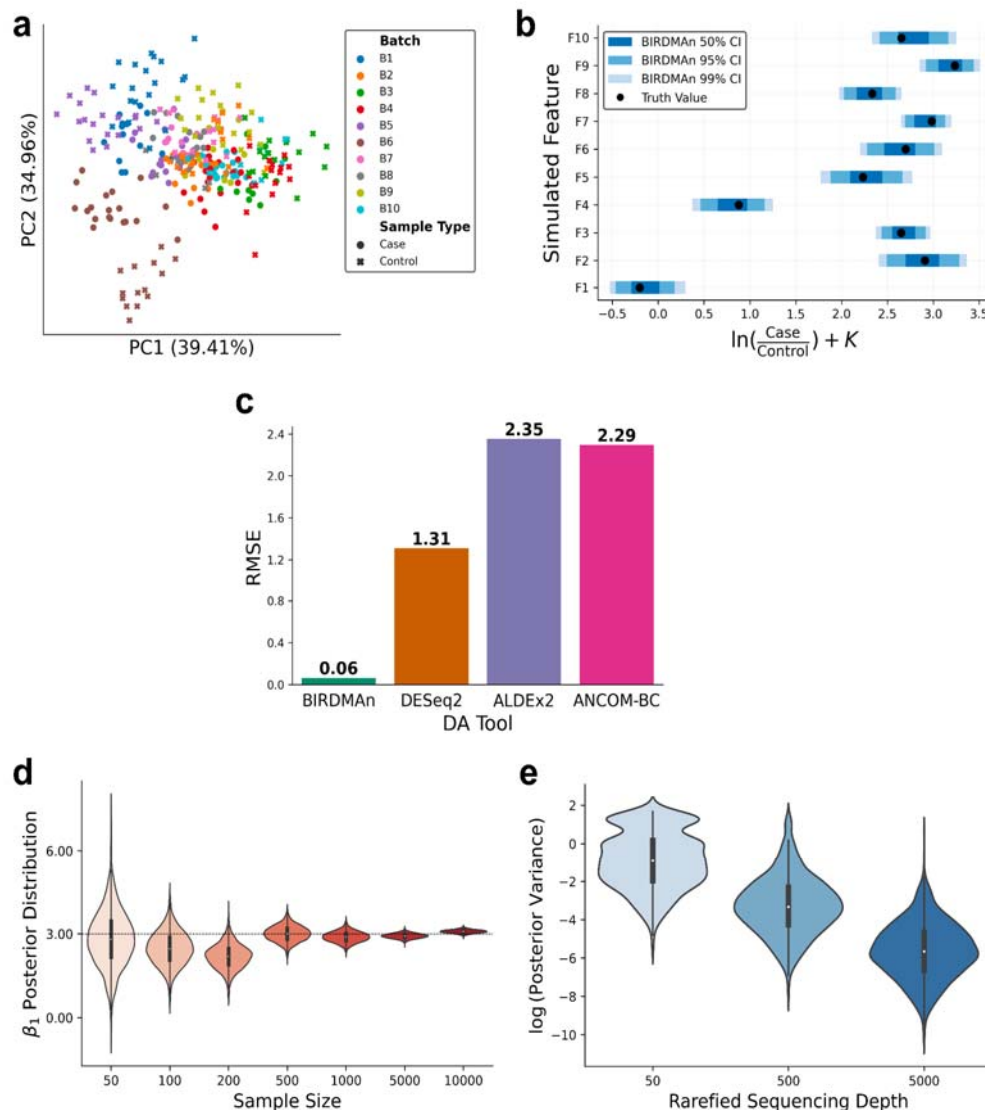126

**Fig 2**: (a) Robust Aitchison principal components plot of the simulated data, showing the large separation by batch effect. Simulations of 10 batches (B1 to B10) of microbiome results, each containing 10 features (F1 to F10), where each feature has a true differential abundance between cases and controls that is the same for each batch, and also a random per batch bias. (b) Recovery of the true simulated log ratio between cases and controls for each feature (black dots), with credible intervals on average centered on the true log ratio (blue bars). (c) Superior performance of BIRDMAn over other differential abundance methods in minimizing the RMSE of the difference between the estimated mean posterior log ratio between cases and controls, revealing a >20-fold improvement in RMSE over the nearest competitor, DESeq2. (d) Estimated distributions of log-fold changes from Bayesian analysis tighten as the number of samples increases. Dashed line represents the true simulated value for each simulation. (e) Rarefaction simulation performed using multinomial count generative models (1000 features) at three

141 *different sequencing depths shows that the variance of the posterior distribution decreases as*
142 *depth increases.*

## BIRDMAn models capture biological signals missed by other methods during dual-course longitudinal antibiotics

145
146 Another challenge for DA methods is to compare multiple samples from the same subject
147 longitudinally (repeated measures) since concomitant host-specific variation can obscure
148 phenotypically-associated microbial changes. Methods designed for longitudinal data[22–26]
149 cannot easily account for modeling perturbations and struggle with scaling to high dimensions.
150 To demonstrate the use of BIRDMAn on repeated measure study designs, we evaluated a
151 published longitudinal study of two courses of the antibiotic ciprofloxacin (Cp) (3 subjects, 7
152 timepoints)[27]. Notably, this study originally concluded that inter-subject variability drove the
153 response to antibiotics by examining beta-diversities, which do not account for auto-correlation
154 effects of repeated measures[28] (Fig 3a). Other studies have also highlighted the importance of
155 properly accounting for the microbial community composition prior to antibiotics when assessing
156 varying responses[29,30], which requires accurate temporal modeling.

157
158 Given BIRDMAn's flexibility, we constructed a customized DA model that leverages Linear
159 Mixed Effects models, accounting for repeated measurements from subjects while computing
160 temporal differences (see Methods). This model design then enabled the exploration of common
161 microbial community changes associated with antibiotic perturbation, which the originally
162 published methods could not identify. With the computed log-fold changes over time (Supp Fig
163 1a), we investigated how consistent antibiotic induced shifts were across subjects. For each
164 temporal difference, we took the top and bottom 40 OTUs to calculate sample log-ratios, which
165 were used to predict antibiotics intake[31]. From these log-ratios, we observed strong, statistically
166 significant temporal shifts associated with each successive time interval (Supp Fig 1b).

167
168 To determine if existing tools could have identified these timepoint-specific perturbations, we
169 also developed a multinomial logistic regression classifier based on the BIRDMAn results to
170 predict the corresponding time interval. We then compared our prediction performances against
171 classifiers built using ALDEx2[32], ANCOM-BC[33], and DESeq2[34] results on the same samples, as
172 well as a classifier built on the center log-ratio transformed table (see Methods). Remarkably,
173 BIRDMAn-informed classifiers were able to accurately differentiate between the different
174 treatment groups (accuracy > 0.65) (Supp Fig 1c) and showed substantially better prediction
175 accuracy compared to all other methods (Fig 3b). We also verified that this superior
176 performance held across varying numbers of OTUs used in log-ratio calculation (Supp Fig 1d).
177 Ultimately, these findings show how BIRDMAn can identify clear-cut biological changes that
178 were missed or obscured by other approaches, highlighting its ability to confirm expected
179 biological hypotheses.

180
181 We used the sample log-ratios associated with the First and Second Cp applications and plotted
182 the dynamics over time (Fig 3c, d). Accordingly, we plotted the corresponding derivative log-fold

183  changes computed from BIRDMAn (Fig 3e, f) and see that our trajectories match between the
184  sample log-ratios and the estimated log-fold changes, indicating that our model was able to
185  successfully capture the overall signal independent of subject.
186
187  The antibiotic used in the original work, Cp, is known to primarily target (though not exclusively)
188  gram negative bacteria[35,36]. We thus hypothesized that the differential abundance results should
189  reflect the longitudinal dynamics of gram negative bacterial abundance. In the top and bottom
190  40 most changed taxa after FirstCp, 17.5% of the numerator taxa were gram negative, whereas
191  27.5% of the denominator were gram negative (Supp Fig 2e). Given the Cp antibiotic
192  mechanism, it is likely that gram negative taxa in the denominator decreased which caused the
193  increased log-ratio[37,38] (Figure 2c). We see that there is a sharp decrease in this log-ratio at
194  FirstWPC, which could be attributed to gut homeostasis[37,38]. However, we see a weaker pattern
195  in the top/bottom 40 microbes after SecondCp, where 2.5% of the numerator taxa were gram
196  negative and 10% of the denominator taxa were gram negative. In contrast to the FirstCp, the
197  microbes most affected by SecondCp quickly returned to their original abundances.
198  Furthermore, we see that the microbes most altered by FirstCp were not affected by SecondCp.
199  Altogether this hints at newly acquired antimicrobial resistant genes after the application of
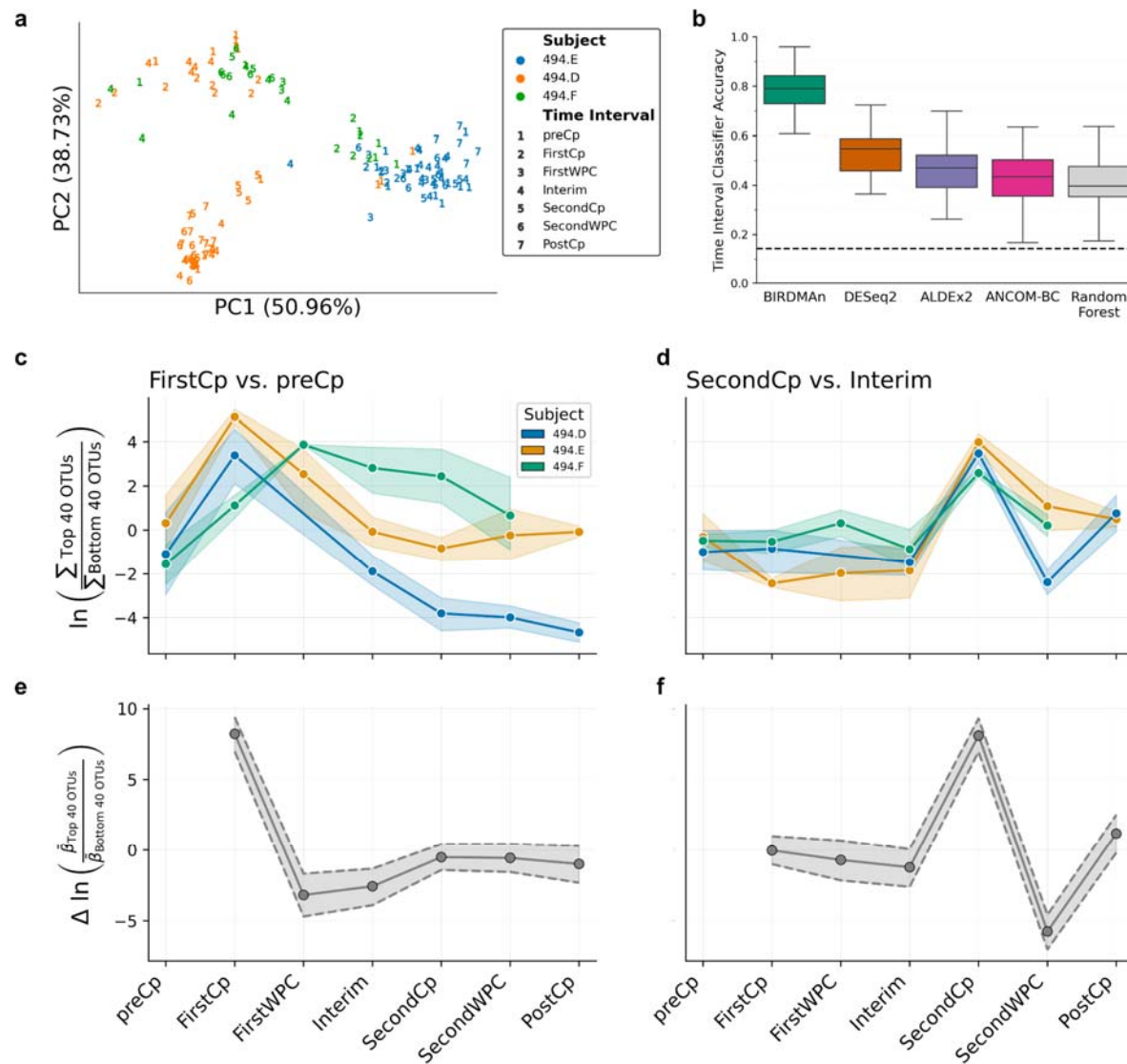200  FirstCp.

**Fig 3**: *(a) Robust Aitchison principal components plot of full dataset shows samples cluster primarily by host subject. (b) Balanced accuracy of multinomial classification of time point by tool. Differential abundance classifiers were constructed using logistic regression with the log-ratios of the top 40 and bottom 40 OTUs associated with each timepoint as predictors. Repeated k-fold cross-validation was performed with 5 splits and 10 repeats. The mean classifier error is at least twice as great with all other differential abundance tools as with BIRDMAn. Dashed line represents random guessing performance among the seven timepoints. (c, d) Dynamics of sample log-ratios of (c) first Cp course and (d) second Cp course colored by subject. (e, f) Dynamics of BIRDMAn-estimated log-fold changes associated with (e) FirstCp effect with preCp as reference and (f) SecondCp effect with Interim as reference. Shaded intervals represent the 90% credible interval of the estimated posterior distributions.*

## BIRDMAn models mitigate batch effects in cancer microbiome data

214
215 To investigate how generalizable BIRDMAn models are with respect to population
216 heterogeneity, we conducted a meta-analysis using cancer microbiome data derived from The
217 Cancer Genome Atlas (TCGA). This dataset is known to have large structural batch effects[4],
218 where the samples were processed at multiple centers across North America, resulting in an
219 artificial separation of cancer microbiomes by sequencing center if not otherwise accounted for
220 (Fig 4a, Supp Fig 2a)[4,39]. These effects can make it difficult to determine microbial biomarkers
221 associated with tumors rather than artifacts of technical variation, but correcting for this could
222 enable downstream host-microbial cancer analyses. We thus tested how well BIRDMAn models
223 could extract biological signals from this dataset while accounting for technical batch effects
224 modeled as random effects. We additionally modeled each microbial feature's abundance using
225 this approach to determine the specificity of these microbes for each cancer type (see Methods
226 and Code).
227
228 Since cancer types are known to have distinct microbiomes[4,40], we first confirmed that BIRDMAn
229 models could extract cancer type-specific differences despite the technical variation observed in
230 this study. From our log-ratio classification benchmarks, we observe that our custom BIRDMAn
231 model can detect a substantially stronger differential signature between the cancer types
232 compared to ALDEx2, ANCOM-BC, DESeq2, and Random Forests (Fig 4b; note the axis log-
233 scaling) after controlling for the batch effects due to the sequencing center (Supp Fig 2c).
234
235 To determine the generalizability of our results, we then constructed a leave-one-center-out
236 cross-validation benchmark using logistic regression on the BIRDMAn-computed log-ratios.
237 Four cancer types with at least three represented data submitting centers (head and neck
238 cancer [HNSC], bladder cancer [BLCA], thyroid cancer [THCA], and cervical cancer [CESC])
239 were included in this benchmark. The receiver operating characteristic (ROC) curves
240 demonstrated strong classification performance (Fig 4c), indicating that BIRDMAn captures
241 generalizable microbial signals across multiple sequencing centers. Generalizability can be a
242 major challenge in microbiome studies[3], where classifiers become overfitted for individual
243 cohorts. We observe this with other DA tools (ALDEx2, DESeq2, ANCOM-BC) and even
244 Random Forests (Supp Fig 2d), where most tools struggle to achieve an area under the ROC
245 curves (AUROC) of >0.8. BIRDMAn is competitive with these tools, achieving an AUROC >0.9
246 in HNSC, BLCA, and CESC cancers while achieving the highest predictive accuracy in BLCA
247 and CESC cancers. The high classifier accuracy leaving out each individual center
248 demonstrates that no one center's data strongly affects the classifier accuracy, with the
249 exception of BI for THCA.
250
251 To investigate the heterogeneity across different cancer types, we next computed Kendall
252 correlations of BIRDMAn-estimated microbial log-fold changes across all pairs of cancer types.
253 This analysis revealed similarities among cancer types that we would expect, including strong
254 similarities between kidney cancer subtypes (KIRC, KICH, KIRP), lung cancer subtypes (LUAD,
255 LUSC), and gastrointestinal (GI) cancers (COAD, ESCA, HNSC, STAD), Additionally, the
256 BIRDMAn-informed data suggested some novel associations, such as the similarity between
257 kidney cancers and liver cancer (LIHC). When clustering the individual microbes' differentials

258    (Supp Fig 2b), we also observed that numerous GI-specific microbes differentiated GI cancers
259    from other cancer types.
260
261    When focusing on comparing GI cancers to lung cancers, we found that the resulting BIRDMAn
262    log-fold changes accurately reflected known biology surrounding the niches in which these
263    microbiomes are commonly found. Specifically, *Fusobacterium*[41], *Prevotella*[42], and *Coproccus*[43]
264    are genera commonly found in the GI tract; conversely, *Pseudomonas*[44], *Staphyloccus*[45], and
265    *Sphingobacterium*[46] genera include opportunistic pathogens that are commonly found in lung
266    infections (Fig 4f). We cross-referenced our results against the Tsay *et al.* cohort that utilized
267    16S rRNA sequencing to investigate lung cancer. Out of the 469 genera in the TCGA lung
268    issues, we observed that 39% of these microbes were also observed in the Tsay *et al.* cohort,
269    despite known previous discordant findings comparing 16S rRNA sequencing and whole
270    genome sequencing[47,48]. Furthermore, when we focus on the top 100 microbes that are
271    detected to be associated with lung cancer, 70% of the represented genera were observed in
272    both the TCGA and Tsay *et al.* datasets. Altogether, this shows how BIRDMAn models can
273    provide biologically-informative results while properly accounting for and mitigating strong
274    structural batch effects that currently confound other DA approaches.
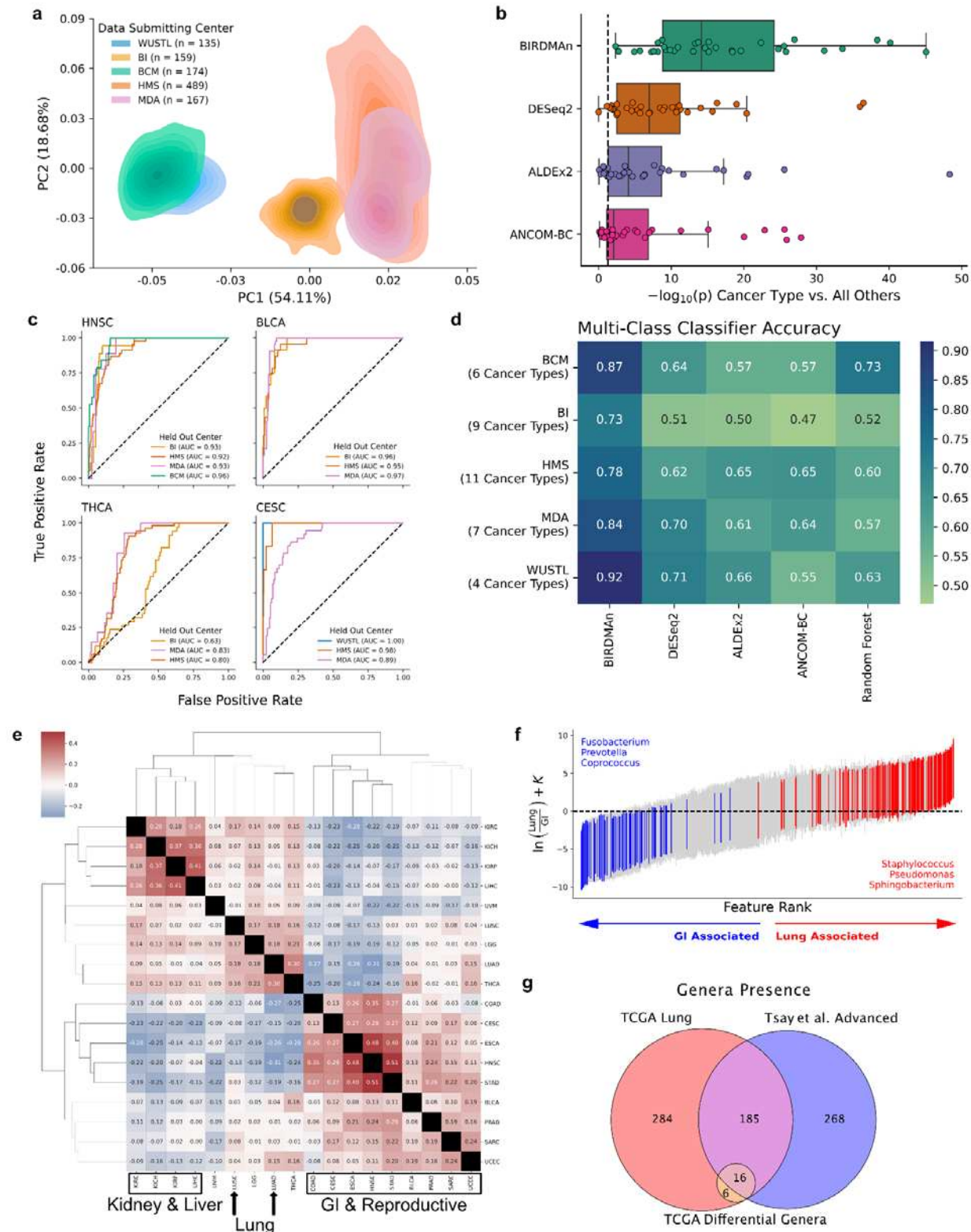275

276
**Fig 4:** *(a) Whole-genome sequenced cancer microbiome data from TCGA shows strong batch*
278 *effects by sequencing center (colored by center; see Supp Fig 2a for per cancer type plots).*
279 *Samples are summarized by the 2D kernel density estimate for each center. (b) T-test p-values*

280  *comparing log-ratios of each cancer type vs. all others within each center. Dashed line*
281  *represents p=0.05. All differential abundance methods show significant differences with log-*
282  *ratios to separate the microbes in each individual cancer type from those found in all other*
283  *cancer types, but BIRDMAn outperforms other methods in highlighting this difference. (c) ROC*
284  *curves for leave-one-center-out cross-validation for four cancer types where at least 3 centers*
285  *sequenced that cancer type (BRCA was not included as it was used as reference). Classifiers*
286  *were built to predict one-vs-rest for that cancer type. BI = Broad Institute of MIT and Harvard;*
287  *BCM = Baylor College of Medicine; HMS = Harvard Medical School; MDA = MD Anderson*
288  *Institute for Applied Cancer Science; WUSTL = Washington University School of Medicine. (d)*
289  *Multinomial (mean) classification accuracy of classifiers to predict cancer type given the log-*
290  *ratios computed from the top and bottom 200 taxa associated with each cancer type. Random*
291  *Forests classifier, which is frequently used in this field but is not based on differential*
292  *abundance, was included as a comparison for this class of methods. Classifications were*
293  *performed within each center to remove batch effects from predictions. BIRDMAn outperforms*
294  *all other methods, including Random Forests, for all tumor types. (e) Clustermap of Kendall tau*
295  *correlation coefficients of pairwise cancer type differentials (breast cancer as reference). (f)*
296  *Comparison of lung-associated genera with GI-associated genera. Highlighted genera are*
297  *known to be associated with either lung or GI microbiome and show strong directionality in the*
298  *BIRDMAn results. (g) Venn diagram of genera present in TCGA lung samples and genera*
299  *present in advanced stage lung cancer from work published by Tsay et al. Additionally, the 22*
300  *genera represented in the top 100 features associated with TCGA lung cancer cancers are*
301  *included. A majority of these genera (16/22) are present in both datasets.*

# Discussion

303  Advances in Bayesian computation have lowered the barriers to developing statistical
304  workflows. To empower microbiome scientists to take advantage of these methods, we
305  developed and implemented a novel approach to differential abundance based on Bayesian
306  hierarchical modeling, with advantages highlighted in simulation benchmarks and real-world
307  datasets. Chiefly, BIRDMAn is designed as a *framework* for researchers to account for the
308  statistical constraints specific to their biological questions. We have demonstrated the benefits
309  of this framework in common biological scenarios involving longitudinal study designs and
310  sequencing center variation — where BIRDMAn can better correct for technical variation than
311  existing methods while identifying biologically-relevant signals. In addition to the ability to
312  construct novel DA models, we presented a robust method for benchmarking and comparing
313  results from different DA tools. In contrast to previous efforts investigating FDR in simulation
314  and reproducibility benchmarks[19,49,50], we show how to construct sample classifiers from the log-
315  fold change estimates, enabling machine learning techniques such as cross-validation on
316  biological datasets.

317

318  Another key challenge of DA benchmarking is the absence of "ground truth," or the true
319  differentials associated with biological conditions, especially in the presence of strong batch
320  effects. Simulations with known parameters for batch and biological effects can address this
321  limitation, and we showed that BIRDMAn models could recover, with high accuracy and

322  precision, these parameters and their uncertainty. Additional simulations on parameter
323  uncertainty further showed decreases with increased sample size and higher sequencing depth,
324  corroborating previous work and traditional statistical knowledge.
325

326  We then investigated two real-world case studies—antibiotics response/recovery and cancer
327  microbiome interactions—demonstrating how BIRDMAn can uncover expected and novel
328  biology. For each dataset, BIRDMAn models were able to account for the inherent effects of
329  center/subject on individual microbial abundances while, when necessary, accounting for
330  complex statistical factors (such as, random intercepts, random slopes, overdispersion). To
331  date, there is no other DA tool that provides a similar and necessary degree of flexible statistical
332  modeling. Our results on the previously published antibiotics dataset revealed the attenuating
333  effect of repeated Cp courses on Gram-negative bacteria, with potential implications for clinical
334  practice using antibiotics. Additionally, BIRDMAn-informed results from the cancer microbiome
335  dataset could be useful in developing novel diagnostic and therapeutic strategies that target or
336  perturb cancer-specific features.
337

338  In light of our findings, there are notable assumptions that need to be considered. Specifically,
339  the choice of prior distributions affects the estimated posterior distributions, especially at low
340  sample sizes. Although priors allow researchers to include their expertise in their modeling
341  procedure, it is often the case that an appropriate prior distribution is unknown, requiring
342  uninformed priors with high uncertainty to be used. However, we note that as more analyses are
343  performed, their results can provide a rationale for picking future priors—a strong advantage of
344  the Bayesian approach over non-Bayesian methods. For our purposes, we defined the same
345  prior distribution for each feature within a dataset, but this can easily be adapted to better model
346  features with their expected parameter range. We also note that the (common) lack of absolute
347  abundance data is a limitation in evaluating differential abundance[51]. Strategies to account for
348  this, such as in Williamson *et al.*[52], could potentially be translated into BIRDMAn models to
349  augment the modeling results.  Furthermore, we model the microbial abundances using the
350  negative binomial approach, which is currently contested as an appropriate model for
351  sequencing count data[53]. Still, an advantage of BIRDMAn is that the likelihood function is not
352  restricted to the negative binomial, and one can exchange it for the Poisson-Lognormal,
353  Multinomial, or any other count distribution[54].
354

355  To summarize, we find that careful statistical consideration during DA analysis enables the
356  identification of microbe-phenotype associations that are missed by existing tools. The flexibility
357  of BIRDMAn can thoroughly account for unwanted confounding factors, such as batch and
358  subject, resulting in higher confidence in reported microbial biomarkers. Moreover, the
359  presented log-ratio benchmarking approach opens up numerous possibilities for testing
360  improved machine learning capabilities on microbiome data. Overall, we posit that BIRDMAn's
361  flexibility and utility will provide impactful statistical results for complex study designs while
362  enabling reproducible science in the microbiome field.
363

364  # Methods

365  ## Performing Bayesian inference with Stan

366  Parameter estimation was performed using Bayesian inference. Our approach utilizes Bayes'
367  Rule where $\theta$ represents the parameter space and $D$ represents our collected data:
368

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

369
370

371  Because the evidence term, $P(D)$, is simply a normalizing constant, we can rewrite Bayes' Rule
372  as follows, substituting terms with their common nomenclature:
373

374  $$\text{Posterior} \sim \text{Likelihood} \cdot \text{Prior}$$

375

376  Thus, our objective with Bayesian inference is to obtain the posterior distribution by modeling
377  the likelihood function of our data as well as our prior knowledge of the parameters. Absent a
378  model formulation involving conjugate priors, we cannot compute the posterior distribution
379  analytically. Instead, we use Stan to draw samples from the posterior distribution using the No-
380  U-Turn Hamiltonian Monte Carlo sampler[20]. A series of Markov chains are initialized and
381  allowed to "warm-up" in their exploration of the parameter posterior distributions. Once the
382  defined number of warm-up iterations has concluded, a set number of samples are drawn from
383  each of the chains. Multiple chains are run to ensure that model convergence occurs.
384

385  We implement Bayesian inference using the CmdStanPy interface in Python, calling the C++
386  Stan toolchain for efficient sampling. The warm-up iterations are discarded by default and the
387  sampling iterations are saved for each Markov chain.


388  ## Negative binomial model parameterization

389  We fit counts of each microbe in a dataset according to a negative binomial distribution as an
390  approximation of multinomial logistic regression[55]. Due to overdispersion, standard count
391  models such as Poisson are inappropriate for sequencing data[21]. We note that the negative
392  binomial model can be considered an extension to the Poisson model with additional variance
393  components[56].
394

395  The negative binomial models used in this work are described by parameters for both mean and
396  overdispersion. This is in contrast to traditional parameters in negative binomial models
397  described by the probability of success and the number of failures before an instance of a
398  success. The former model, often referred to as the "alternative parameterization," is more
399  amenable to generalized linear modeling through hierarchical models as the mean can be
400  modeled directly.
401

402 The basic format of the alternative parameterization negative binomial model is described below
403 where $n$ corresponds to the count, $\phi$ the overdispersion, and $\mu$ the mean count.
404

$$\mathrm{NB}\left(n \mid \mu, \phi\right) = \binom{n + \phi - 1}{n} \left(\frac{\mu}{\mu + \phi}\right)^{n} \left(\frac{\phi}{\mu + \phi}\right)^{\phi}$$

405
406
407 We use a log-link function, $\mu = \exp\left(\eta\right)$ to model the mean where the log mean count, $\eta$, can be
408 represented by linear terms. To account for variable sequencing depth among samples, we
409 include log sequencing depth as an offset term in our models.

## BIRDMAn framework

411 We developed BIRDMAn as a framework for highly-customizable Bayesian differential
412 abundance modeling. BIRDMAn abstracts much of the Bayesian workflow away for usage with
413 microbiome data. An object-oriented approach allows users to subclass basic models for their
414 custom implementations. BIRDMAn includes, by default, a Negative Binomial model
415 implementation. This can be used without writing any new Stan code or subclassing any
416 BIRDMAn objects.
417
418 BIRDMAn models take BIOM tables[57] as input containing the sample and observation IDs.
419 Sample metadata can be provided as Pandas DataFrames. We provide a method,
420 create_regression, with which users can provide an R-style formula to automatically create the
421 design matrix using the patsy Python package. Another method, specify_model, allows the
422 specification of the desired parameters and dimensions to return. This method is used by
423 create_inference to convert CmdStanPy output to ArviZ[58] InferenceData objects.
424
425 There are two base classes included with BIRDMAn termed the TableModel and the
426 SingleFeatureModel. The TableModel allows fitting an entire dataset at once, while the
427 SingleFeatureModel allows for fitting individual features. The SingleFeatureModel is
428 advantageous as it allows for highly parallelized workflows. Because there are often hundreds
429 or thousands of features in a microbiome dataset, we note that using multiple CPUs to run many
430 features at once is often more efficient than fitting the entire table. We provide a convenience
431 class, ModelIterator, to iterate through the features in a given table. This class also allows for
432 dividing the table into chunks. This allows users to customize the number of features to fit at
433 once depending on their computational resources.

## Simulations

435 All simulations were performed through the fixed_param option in CmdStanPy. Ground-truth
436 parameters were provided into a negative binomial generative model to simulate data from
437 mean and dispersion parameters.
438
439 For the data-driven simulation, we randomly drew values for batch offset, batch dispersion, and
440 base dispersion parameters. These parameters were fed into a model with $\beta_0 = N(-8, 1)$,

441   $\beta_1 = N(2, 1)$. Log sampling depth was simulated from a Poisson-Lognormal distribution with $\lambda$
442   drawn from $N(5000, 0.2)$. We simulated 300 samples comprising 10 total batches with 10 total
443   features.

444

445   For the variable sample size simulations, we simulated feature counts for 500 samples with

446   $\beta_0 = 8$, $\beta_1 = 3$, and $\frac{1}{\phi} = 10$. Log sequencing depths were simulated using a Poisson-
447   Lognormal model with $\lambda$ drawn from $N(50000, 0.5)$ where depth varied.

448

449   To simulate variable rarefaction depth, we first drew ground truth intercept and beta values from
450   $N(-8, 1)$ and $N(2, 1)$ respectively for 1000 features. These values were used to generate
451   counts for 300 samples through the multinomial distribution. We used the multinomial
452   distribution to enforce the same sampling depth for all samples, simulating rarefaction.

## Antibiotics case study

454   16S data was downloaded from Qiita study 494; we used 16S OTUs picked against the
455   GreenGenes_13.8[59] reference database at 97% sequence similarity. OTU picking was
456   performed with SortMeRNA[60] with Qiita default parameter values. Features present in fewer
457   than 10 samples were filtered. We also removed samples with a total sequencing depth less
458   than 1000.

459

460   To account for the longitudinal nature of this design, we used backwards difference encoding
461   such that each time point was compared to the one immediately before it. We implemented the
462   subject identifiers as a random effect with both random intercepts and random slopes. The
463   posterior draws were centered around the mean. Ranking of OTUs by differentials for log-ratio
464   feature selection was done using the posterior means.

465

466   We performed t-tests comparing the log-ratios between groups of samples at different
467   timepoints. The alternative hypothesis was chosen such that samples from the later time point
468   would have higher log-ratios than those from the initial timepoint due to the anticipated effect of
469   Cp on microbial populations.

470

471   We then implemented multinomial logistic regression, random forest classification, and repeated
472   k-fold cross-validation through scikit-learn for our machine learning approach. Because DESeq2
473   supports contrasts natively, we computed the same contrasts as BIRDMAn for parity. With
474   ALDEx2 and ANCOM-BC, we computed the differentials associated with each timepoint using
475   preCp as reference. For the random forest classifier, we used the CLR-transformed feature
476   table (with a pseudocount of 1) entries as the predictors. All models were also provided one-hot-
477   encoded vectors for subject identifiers. Performance was measured using balanced accuracy.
478   For multinomial logistic regression we used the lbfgs solver with 1000 max iterations. For the
479   random forest classifier we used a set random seed and 100 estimators. We used repeated
480   stratified k-fold cross validation with 5 splits and 10 repeats and a random seed. All other
481   parameters not mentioned were set to the scikit-learn defaults.

482

483    Posterior draws for timepoint-contrast differentials were analyzed with (1) FirstCp-associated
484    features with preCp-associated features as reference and (2) SecondCp-associated features
485    with Interim-associated features as reference. In this way, the posterior distribution reflects how
486    each Cp course affects the selected bacterial features over time.

487

488    For determining the Gram status of each OTU, we used the BugBase[61] web interface. We took
489    the set intersection of Gram positive and Gram negative features with the features associated
490    with both FirstCp and SecondCp to determine the Gram status breakdown of both numerator
491    and denominator features.

492 ## TCGA case study

493    The bacterial TCGA tables were obtained from those processed in Narunsky-Haziza et al.[62] and
494    Poore et al.[4] All TCGA sequence data were accessed via the Cancer Genomics Cloud[63] (CGC)
495    as sponsored by SevenBridges (https://cgc.sbgenomics.com) after obtaining data access from
496    the    TCGA    Data    Access    Committee    through    dbGaP    (https://
497    dbgap.ncbi.nlm.nih.gov/aa/wga.cgi?page=login). On Qiita[64], TCGA WGS host-depleted and
498    quality-controlled fastq files were used to generate a metagenomic table by direct genome
499    alignments based on Woltka v0.1.1[65] against the RefSeq[66] release 200 (built as of May 14,
500    2020). The resulting tables can be found on Qiita under study ID 13722, of which we filtered to
501    only analyze the bacteria and then were subsequently decontaminated through decontam[67]
502    (https://github.com/benjjneb/decontam) (version 1.14.0) following the protocol described in
503    Poore et al.[4]

504

505    After initial table generation, we removed samples from data submitting centers with very few
506    samples. We also filtered our data to only include samples from white, African-American, and
507    Asian races. Additionally, we only included samples from patients who were alive at the time of
508    sample procurement and retained only one sample per subject. To filter out lowly prevalent
509    features, we removed features present in fewer than 50 total samples. To remove samples with
510    low sequencing depth, we set a threshold of 500 reads. Finally, we included only cancer types
511    with at least 20 instances in the dataset for statistical power.

512

513    We then built statistical models to model the differential associated with each cancer type.
514    Because TCGA did not include "normal" samples from healthy individuals, we used breast
515    cancer (BRCA) tumor samples as reference. Both race[68] and gender were also included as
516    covariates. Data submitting center was incorporated as a random effect (both random intercepts
517    and random slopes).

518

519    Posterior means were computed for each feature's association with each individual cancer type.
520    For each cancer type, we ranked the differentials and used the top and bottom 200 features
521    associated with that cancer type to compute log-ratios per sample. These log-ratios were used
522    as predictor variables in our machine learning models.

523

524 Because not every cancer type was represented in each center, we performed multi-class
525 classification within centers. For each center, we fit a model to predict cancer type from our log-
526 ratios. This procedure was performed with 5 repeats of stratified 2-fold cross-validation. We
527 repeated this machine learning process for cancer type differentials from DESeq2, ALDEx2, and
528 ANCOM-BC. For comparison, we fit a random forest classifier on the CLR-transformed feature
529 table to predict cancer type as well.

531 The leave-one-center-out models were fit using binomial logistic regression with balanced class
532 weights. For each cancer type, we fit a model on all but one center and used that model to
533 predict cancer type for the held-out center. We also used the same random forest classifier as
534 previously described for comparison.

## Analysis & visualization software

536 Analysis of the results in this work were primarily performed through Python (v3.8.13). Pandas[69]
537 (v1.1.5) and NumPy[70] (v1.22.3) were used for general data analysis. SciPy[71] (v1.7.3) was used
538 for computing statistical tests. For interfacing with multidimensional arrays we used xarray[72]
539 (v0.20.1) and ArviZ[58] (0.12.1). Machine learning models were fit and cross-validated using
540 scikit-learn[73] (v1.0.2). Python figures were generated using seaborn[74] (v0.11.2) and Matplotlib[75]
541 (v3.5.1) as well as Matplotlib-venn (v0.11.7). We used biom-format[57] (2.1.12) and scikit-bio
542 (v0.5.6) for statistical analysis of microbiome data structures.

544 R analysis was performed using the tidyverse[76] packages dplyr (v1.0.9), stringr (v1.4.0), and
545 ggplot2 (v3.3.6). Phylogenetic visualization was performed using treeio[77] (v1.18.0) and ggtree[78]
546 (v3.2.0). BIOM tables were read using the biomformat R package (v1.22.0).

# Code and data availability

548 All data used were downloaded from publicly available Qiita studies. The scripts and Stan
549 models used to analyze these data as well as Jupyter notebooks for the visualizations are
550 available at https://github.com/knightlab-analyses/birdman-analyses-final. The BIRDMAn
551 software package is available at https://github.com/biocore/BIRDMAn and the documentation is
552 available at https://birdman.readthedocs.io/. All analyses in this work were performed using
553 BIRDMAn v0.1.0.

# Acknowledgments

# Author information

564  G.R., J.T.M., and R.K. conceived the idea for the study. G.R. & J.T.M. developed the BIRDMAn
565  software package. G.R., J.T.M., C.G., G.D.S-P., & C.M. contributed to the case study and
566  simulation analysis. C.A., J.T.M., C.M., & R.K. helped to define the scope of the analyses. G.R.
567  & Y.C. contributed to the documentation for BIRDMAn. M.E., Y.C., D.H., & C.M. gave critical
568  feedback on the usage and documentation of the software. All authors helped write and review
569  the manuscript.

# Conflicts of interest

571  G.D.S.-P. and R.K. are inventors on a US patent application (PCT/US2019/059647) submitted
572  by The Regents of the University of California and licensed by Micronoma; that application
573  covers methods of diagnosing and treating cancer using multi-domain microbial biomarkers in
574  blood and cancer tissues. G.D.S.-P. and R.K. are founders of and report stock interest in
575  Micronoma. G.D.S.-P. has filed several additional US patent applications on cancer bacteriome
576  and mycobiome diagnostics that are owned by The Regents of the University of California or
577  Micronoma. R.K. additionally is a member of the scientific advisory board for GenCirq, holds an
578  equity interest in GenCirq, and can receive reimbursements for expenses up to US $5,000 per
579  year.

# References

581  1. Sochocka, M. *et al.* The Gut Microbiome Alterations and Inflammation-Driven Pathogenesis

582     of Alzheimer's Disease—a Critical Review. *Mol. Neurobiol.* **56**, 1841–1851 (2019).

583  2. Fouquier, J. *et al.* The Gut Microbiome in Autism: Study-Site Effects and Longitudinal

584     Analysis of Behavior Change. *mSystems* **6**, e00848-20 (2021).

585  3. Wirbel, J. *et al.* Meta-analysis of fecal metagenomes reveals global microbial signatures that

586     are specific for colorectal cancer. *Nat. Med.* **25**, 679–689 (2019).

587  4. Poore, G. D. *et al.* Microbiome analyses of blood and tissues suggest cancer diagnostic

588     approach. *Nature* **579**, 567–574 (2020).

589  5. Villapol, S. Gastrointestinal symptoms associated with COVID-19: impact on the gut

590    microbiome. *Transl. Res.* **226**, 57–69 (2020).

591    6. Zuo, T. *et al.* Alterations in Fecal Fungal Microbiome of Patients With COVID-19 During Time

592    of Hospitalization until Discharge. *Gastroenterology* **159**, 1302-1310.e5 (2020).

593    7. Poyet, M. *et al.* A library of human gut bacterial isolates paired with longitudinal multiomics

594    data enables mechanistic microbiome research. *Nat. Med.* **25**, 1442–1452 (2019).

595    8. Kostic, A. D. *et al.* The dynamics of the human infant gut microbiome in development and in

596    progression toward type 1 diabetes. *Cell Host Microbe* **17**, 260–273 (2015).

597    9. Proctor, L. M. *et al.* The Integrative Human Microbiome Project. *Nature* **569**, 641–648 (2019).

598    10.    Gopalakrishnan, V. *et al.* Gut microbiome modulates response to anti–PD-1

599    immunotherapy in melanoma patients. *Science* **359**, 97–103 (2018).

600    11.    Spencer, C. N. *et al.* Dietary fiber and probiotics influence the gut microbiome and

601    melanoma immunotherapy response. *Science* **374**, 1632–1640 (2021).

602    12.    Lee, K. A. *et al.* Cross-cohort gut microbiome associations with immune checkpoint

603    inhibitor response in advanced melanoma. *Nat. Med.* **28**, 535–544 (2022).

604    13.    Hiergeist, A., Reischl, U. & Gessner, A. Multicenter quality assessment of 16S ribosomal

605    DNA-sequencing for microbiome analyses reveals high inter-center variability. *Int. J. Med.*

606    *Microbiol.* **306**, 334–342 (2016).

607    14.    Wang, Y. & LêCao, K.-A. Managing batch effects in microbiome data. *Brief. Bioinform.*

608    **21**, 1954–1970 (2020).

609    15.    Chen, W. *et al.* A comparison of methods accounting for batch effects in differential

610    expression analysis of UMI count based single cell RNA sequencing. *Comput. Struct.*

611    *Biotechnol. J.* **18**, 861–873 (2020).

612    16.    Vandeputte, D. *et al.* Quantitative microbiome profiling links gut community variation to

613    microbial load. *Nature* **551**, 507–511 (2017).

614    17.    Kumar, M. S. *et al.* Analysis and correction of compositional bias in sparse sequencing

615    count data. *BMC Genomics* **19**, 799 (2018).

616    18.    Nixon, M. P., Letourneau, J., David, L. A., Mukherjee, S. & Silverman, J. D. A Statistical

617         Analysis of Compositional Surveys. Preprint at https://doi.org/10.48550/arXiv.2201.03616

618         (2022).

619    19.    Nearing, J. T. *et al. Microbiome differential abundance methods produce disturbingly*

620         *different results across 38 datasets.* 2021.05.10.443486

621         https://www.biorxiv.org/content/10.1101/2021.05.10.443486v1 (2021)

622         doi:10.1101/2021.05.10.443486.

623    20.    Hoffman, M. D. & Gelman, A. The No-U-Turn Sampler: Adaptively Setting Path Lengths

624         in Hamiltonian Monte Carlo. Preprint at https://doi.org/10.48550/arXiv.1111.4246 (2011).

625    21.    McMurdie, P. J. & Holmes, S. Waste Not, Want Not: Why Rarefying Microbiome Data Is

626         Inadmissible. *PLOS Comput. Biol.* **10**, e1003531 (2014).

627    22.    Äijö, T., Müller, C. L. & Bonneau, R. Temporal probabilistic modeling of bacterial

628         compositions derived from 16S rRNA sequencing. *Bioinforma. Oxf. Engl.* **34**, 372–380

629         (2018).

630    23.    Silverman, J. D., Durand, H. K., Bloom, R. J., Mukherjee, S. & David, L. A. Dynamic

631         linear models guide design and analysis of microbiota studies within artificial human guts.

632         *Microbiome* **6**, 202 (2018).

633    24.    Joseph, T. A., Shenhav, L., Xavier, J. B., Halperin, E. & Pe'er, I. Compositional Lotka-

634         Volterra describes microbial dynamics in the simplex. *PLOS Comput. Biol.* **16**, e1007917

635         (2020).

636    25.    Joseph, T. A., Pasarkar, A. P. & Pe'er, I. Efficient and Accurate Inference of Mixed

637         Microbial Population Trajectories from Longitudinal Count Data. *Cell Syst.* **10**, 463-469.e6

638         (2020).

639    26.    Shenhav, L. *et al.* Modeling the temporal dynamics of the gut microbial community in

640         adults and infants. *PLOS Comput. Biol.* **15**, e1006960 (2019).

641    27.    Dethlefsen, L. & Relman, D. A. Incomplete recovery and individualized responses of the

642    human distal gut microbiota to repeated antibiotic perturbation. *Proc. Natl. Acad. Sci.* **108**,

643    4554–4561 (2011).

644    28.    Martino, C. *et al.* Context-aware dimensionality reduction deconvolutes gut microbial

645    community dynamics. *Nat. Biotechnol.* **39**, 165–168 (2021).

646    29.    Gibbons, S. M. Keystone taxa indispensable for microbiome recovery. *Nat. Microbiol.* **5**,

647    1067–1068 (2020).

648    30.    Chng, K. R. *et al.* Metagenome-wide association analysis identifies microbial

649    determinants of post-antibiotic ecological recovery in the gut. *Nat. Ecol. Evol.* **4**, 1256–1267

650    (2020).

651    31.    Morton, J. T. *et al.* Establishing microbial composition measurement standards with

652    reference frames. *Nat. Commun.* **10**, 2719 (2019).

653    32.    Fernandes, A. D. *et al.* Unifying the analysis of high-throughput sequencing datasets:

654    characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by

655    compositional data analysis. *Microbiome* **2**, 15 (2014).

656    33.    Lin, H. & Peddada, S. D. Analysis of compositions of microbiomes with bias correction.

657    *Nat. Commun.* **11**, 3514 (2020).

658    34.    Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion

659    for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).

660    35.    Oliphant, C. M. & Green, G. M. Quinolones: A Comprehensive Review. *Am. Fam.*

661    *Physician* **65**, 455–465 (2002).

662    36.    Card, R. M. *et al.* Impact of Ciprofloxacin and Clindamycin Administration on Gram-

663    Negative Bacteria Isolated from Healthy Volunteers and Characterization of the Resistance

664    Genes They Harbor. *Antimicrob. Agents Chemother.* **59**, 4410–4416 (2015).

665    37.    Peterson, C. T., Sharma, V., Elmén, L. & Peterson, S. N. Immune homeostasis,

666    dysbiosis and therapeutic modulation of the gut microbiota. *Clin. Exp. Immunol.* **179**, 363–

667    377 (2015).

668    38.    Ramirez, J. *et al.* Antibiotics as Major Disruptors of Gut Microbiota. *Front. Cell. Infect.*

669        *Microbiol.* **10**, (2020).

670    39.    Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat.*

671        *Genet.* **45**, 1113–1120 (2013).

672    40.    Nejman, D. *et al.* The human tumor microbiome is composed of tumor type-specific

673        intracellular bacteria. *Science* **368**, 973–980 (2020).

674    41.    Bullman, S. *et al.* Analysis of Fusobacterium persistence and antibiotic response in

675        colorectal cancer. *Science* **358**, 1443–1448 (2017).

676    42.    Lo, C.-H. *et al.* Enrichment of Prevotella intermedia in human colorectal cancer and its

677        additive effects with Fusobacterium nucleatum on the malignant transformation of colorectal

678        adenomas. *J. Biomed. Sci.* **29**, 88 (2022).

679    43.    Flemer, B. *et al.* Tumour-associated and non-tumour-associated microbiota in colorectal

680        cancer. *Gut* **66**, 633–643 (2017).

681    44.    Nunley, D. R. *et al.* Allograft Colonization and Infections With Pseudomonas in Cystic

682        Fibrosis Lung Transplant Recipients. *Chest* **113**, 1235–1243 (1998).

683    45.    Laroumagne, S. *et al.* Bronchial colonisation in patients with lung cancer: a prospective

684        study. *Eur. Respir. J.* **42**, 220–229 (2013).

685    46.    Lambiase, A. *et al.* Sphingobacterium respiratory tract infection in patients with cystic

686        fibrosis. *BMC Res. Notes* **2**, 262 (2009).

687    47.    Brumfield, K. D., Huq, A., Colwell, R. R., Olds, J. L. & Leddy, M. B. Microbial resolution

688        of whole genome shotgun and 16S amplicon metagenomic sequencing using publicly

689        available NEON data. *PLOS ONE* **15**, e0228899 (2020).

690    48.    Durazzi, F. *et al.* Comparison between 16S rRNA and shotgun sequencing data for the

691        taxonomic characterization of the gut microbiota. *Sci. Rep.* **11**, 3030 (2021).

692    49.    Hawinkel, S., Mattiello, F., Bijnens, L. & Thas, O. A broken promise: microbiome

693        differential abundance methods do not control the false discovery rate. *Brief. Bioinform.* **20**,

694    210–221 (2019).

695    50.    Yang, L. & Chen, J. A comprehensive evaluation of microbial differential abundance

696    analysis methods: current status and potential solutions. *Microbiome* **10**, 130 (2022).

697    51.    Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome

698    Datasets Are Compositional: And This Is Not Optional. *Front. Microbiol.* **8**, (2017).

699    52.    Williamson, B. D., Hughes, J. P. & Willis, A. D. A multiview model for relative and

700    absolute microbial abundances. *Biometrics* **78**, 1181–1194 (2022).

701    53.    Hawinkel, S., Rayner, J. C. W., Bijnens, L. & Thas, O. Sequence count data are poorly fit

702    by the negative binomial distribution. *PLOS ONE* **15**, e0224909 (2020).

703    54.    Townes, F. W. Review of Probability Distributions for Modeling Count Data.

704    *ArXiv200104343 Stat* (2020).

705    55.    Taddy, M. Distributed multinomial regression. *Ann. Appl. Stat.* **9**, (2015).

706    56.    Lindén, A. & Mäntyniemi, S. Using the negative binomial distribution to model

707    overdispersion in ecological count data. *Ecology* **92**, 1414–1421 (2011).

708    57.    McDonald, D. *et al.* The Biological Observation Matrix (BIOM) format or: how I learned to

709    stop worrying and love the ome-ome. *GigaScience* **1**, 7 (2012).

710    58.    Kumar, R., Carroll, C., Hartikainen, A. & Martin, O. ArviZ a unified library for exploratory

711    analysis of Bayesian models in Python. *J. Open Source Softw.* **4**, 1143 (2019).

712    59.    McDonald, D. *et al.* An improved Greengenes taxonomy with explicit ranks for ecological

713    and evolutionary analyses of bacteria and archaea. *ISME J.* **6**, 610–618 (2012).

714    60.    Kopylova, E., Noé, L. & Touzet, H. SortMeRNA: fast and accurate filtering of ribosomal

715    RNAs in metatranscriptomic data. *Bioinformatics* **28**, 3211–3217 (2012).

716    61.    Ward, T. *et al.* BugBase predicts organism-level microbiome phenotypes. 133462

717    Preprint at https://doi.org/10.1101/133462 (2017).

718    62.    Narunsky-Haziza, L. *et al.* Pan-cancer analyses reveal cancer-type-specific fungal

719    ecologies and bacteriome interactions. *Cell* **185**, 3789-3806.e17 (2022).

720  63.    Lau, J. W. *et al.* The Cancer Genomics Cloud: Collaborative, Reproducible, and

721         Democratized-A New Paradigm in Large-Scale Computational Research. *Cancer Res.* **77**,

722         e3–e6 (2017).

723  64.    Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nat. Methods*

724         **15**, 796–798 (2018).

725  65.    Zhu, Q. *et al.* OGUs enable effective, phylogeny-aware analysis of even shallow

726         metagenome community structures. 2021.04.04.438427 Preprint at

727         https://doi.org/10.1101/2021.04.04.438427 (2021).

728  66.    Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a

729         curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic*

730         *Acids Res.* **35**, D61–D65 (2007).

731  67.    Davis, N. M., Proctor, D. M., Holmes, S. P., Relman, D. A. & Callahan, B. J. Simple

732         statistical identification and removal of contaminant sequences in marker-gene and

733         metagenomics data. *Microbiome* **6**, 226 (2018).

734  68.    Luo, M. *et al.* Race is a key determinant of the human intratumor microbiome. *Cancer*

735         *Cell* **40**, 901–902 (2022).

736  69.    McKinney, W. Data Structures for Statistical Computing in Python. 6 (2010).

737  70.    Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362 (2020).

738  71.    Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.

739         *Nat. Methods* **17**, 261–272 (2020).

740  72.    Hoyer, S. & Hamman, J. xarray: N-D labeled Arrays and Datasets in Python. *J. Open*

741         *Res. Softw.* **5**, 10 (2017).

742  73.    Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**,

743         2825–2830 (2011).

744  74.    Waskom, M. L. seaborn: statistical data visualization. *J. Open Source Softw.* **6**, 3021
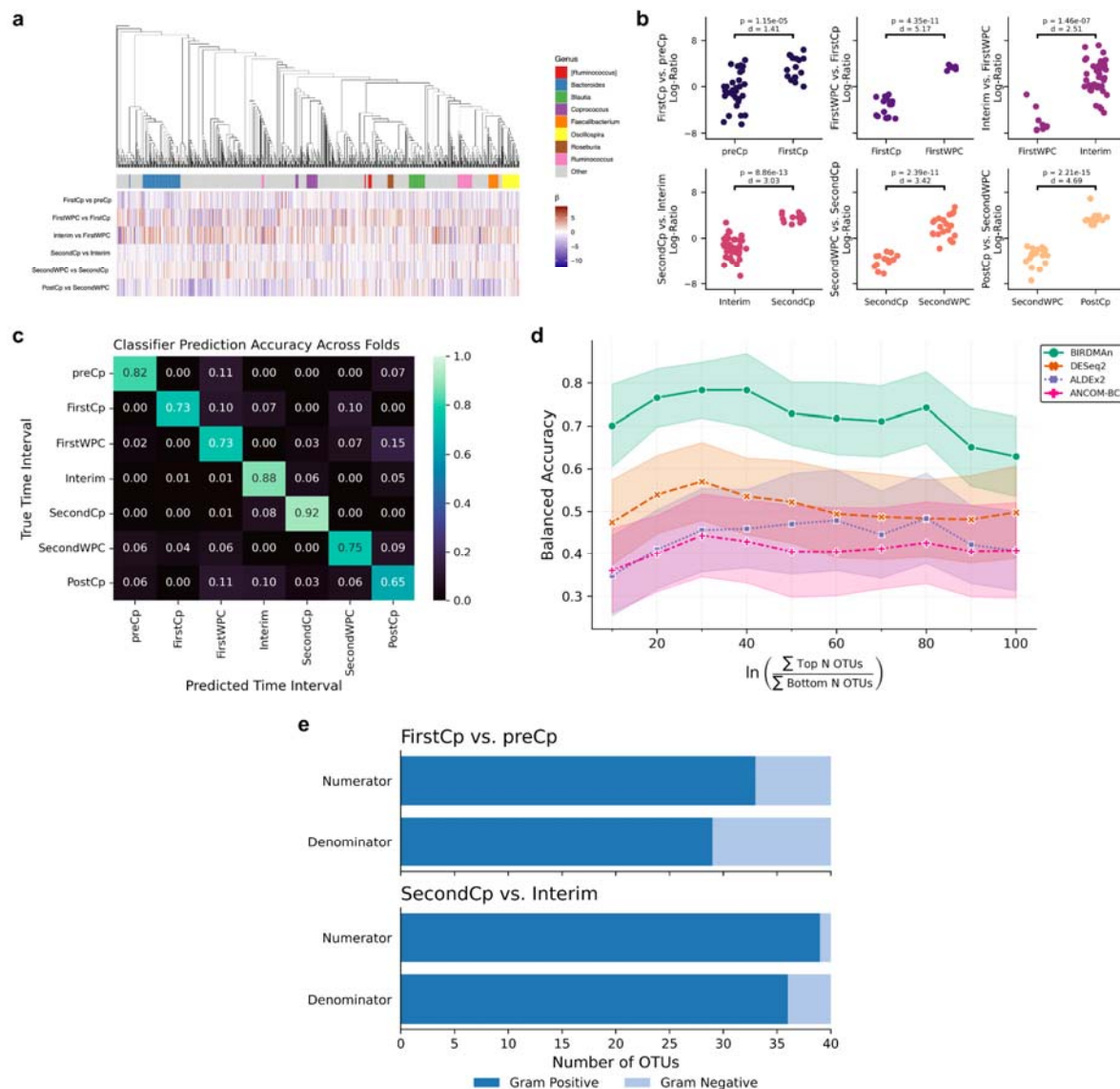
745         (2021).

746    75.    Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90–95

747        (2007).

748    76.    Wickham, H. *et al.* Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).

749    77.    Wang, L.-G. *et al.* Treeio: An R Package for Phylogenetic Tree Input and Output with

750        Richly Annotated and Associated Data. *Mol. Biol. Evol.* **37**, 599–603 (2020).
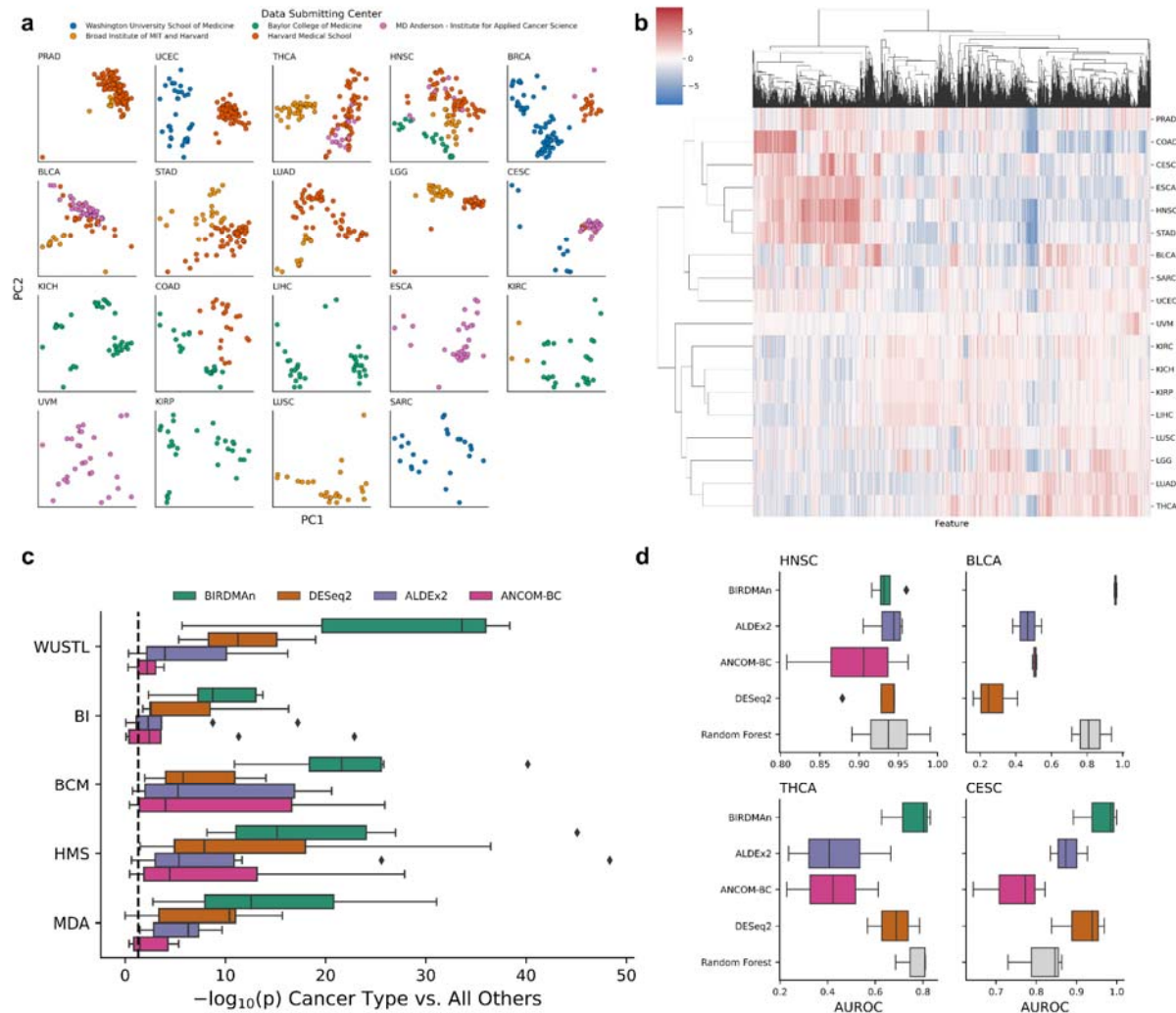
751    78.    Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T.-Y. ggtree: an r package for

752        visualization and annotation of phylogenetic trees with their covariates and other associated

753        data. *Methods Ecol. Evol.* **8**, 28–36 (2017).

754     # Supplementary Figures



755
756     **Supplementary Fig 1**: *(a) Phylogenetic tree of all OTUs with a heatmap of posterior means for*
757     *each time-interval contrast. OTUs assigned to one of the top 8 most abundant genera are*
758     *annotated through the colored strip. (b) When BIRDMAn is used to account for per-subject*
759     *variation, log-ratio comparisons of the top 40 OTUs vs. bottom OTUs are associated with the*
760     *difference between each time point and the next one. For each of these contrasts, the log-ratios*
761     *of the samples between the two time intervals were compared using a one-sided t-test. Plots*
762     *are annotated with p-values. Different taxa contribute to the log ratios for each contrast. (c)*
763     *Overall performance of BIRDMAn classifier on predicting the antibiotics time interval using the*
764     *log-ratios. The classifier prediction accuracies shown are aggregated across folds and repeats*
765     *from repeated k-fold cross-validation. (d) Accuracy of the multinomial classifier by number of*
766     *OTUs used in log-ratio calculations. Points represent mean accuracy across cross-validation*

*767    iterations and shaded areas represent ±1 standard deviation. (e) Distribution of Gram positive*
*768        and Gram negative OTUs associated with FirstCp and SecondCp log-ratios.*

769
770



771
772    **Supplementary Fig 2**: *(a) RPCA projection of the original feature table subset to each*
773    *individual cancer type. Points are colored by data submitting centers, showing that many cancer*
774    *types exhibit strong separation by batch. (b) Posterior means (CLR) of feature differentials*
775    *clustered by cancer type. (c) Log-ratios identified by BIRDMAn separate each tumor type from*
776    *all others when stratified by center. Dashed line represents a t-test p-value at p = 0.05. (d)*
777    *Performance of leave-one-center-out cross-validation logistic regression classifier AUROC of all*
778    *methods.*