1    **Title**

2    Themetagenomics: Exploring Thematic Structure and Predicted Functionality of 16s rRNA

3    Amplicon Data

4

5    **Authors**

6    Stephen Woloszynek (sw424@drexel.edu) [1] [corresponding author]

7    Joshua Chang Mell (joshua.mell@drexelmed.edu) [2]

8    Zhengqiao Zhao (zz347@drexel.edu) [1]

9    Gideon Simpson (grs53@drexel.edu) [3]

10    Michael P. O'Connor (mike.oconnor@drexel.edu) [4]

11    Gail L. Rosen (gailr@coe.drexel.edu) [1]

12

13    **Affiliations**

14    [1] Department of Electrical and Computer Engineering, Drexel University, Philadelphia, PA,

15    United States of America.

16    [2] Department of Microbiology and Immunology, Drexel University College of Medicine,

17    Philadelphia, PA, United States of America.

18    [3] Department of Mathematics, Drexel University, Philadelphia, PA, United States of America

1

19  [4] Department of Biodiversity, Earth, and Environmental Science, Drexel University,

20  Philadelphia, PA, United States of America

21

22  **Abstract**

23

24  Analysis of microbiome data involves identifying co-occurring groups of taxa associated with

25  sample features of interest (e.g., disease state). Elucidating such relations is often difficult as

26  microbiome data are compositional, sparse, and have high dimensionality. Also, the

27  configuration of co-occurring taxa may represent overlapping subcommunities that contribute

28  to sample characteristics such as host status. Preserving the configuration of co-occurring

29  microbes rather than detecting specific indicator species is more likely to facilitate biologically

30  meaningful interpretations. Additionally, analyses that use taxonomic relative abundances to

31  predict the abundances of different gene functions aggregate predicted functional profiles

32  across taxa. This precludes straightforward identification of predicted functional components

33  associated with subsets of co-occurring taxa. We provide an approach to explore co-occurring

34  taxa using "topics" generated via a topic model and link these topics to specific sample features

35  (e.g., disease state). Rather than inferring predicted functional content based on overall

36  taxonomic relative abundances, we instead focus on inference of functional content within

37  topics, which we parse by estimating interactions between topics and pathways through a

38  multilevel, fully Bayesian regression model. We apply our methods to three publicly available

39  16S amplicon sequencing datasets: an inflammatory bowel disease dataset from Gevers et al., an

40    oral cancer dataset from Schmidt et al., and a time-series dataset from David et al. Using our

41    topic model approach to uncover latent structure in 16S rRNA amplicon surveys, investigators

42    can (1) capture groups of co-occurring taxa termed topics; (2) uncover within-topic functional

43    potential; (3) link taxa co-occurrence, gene function, and environmental/host features; and (4)

44    explore the way in which sets of co-occurring taxa behave and evolve over time. These methods

45    have been implemented in a freely available R package:

46    https://github.com/EESI/themetagenomics.

47

48    **Introduction**

49    High-throughput sequencing now permits for the analysis of multiple large datasets on the

50    microbiome and diseases of interest.  Historically, researchers have sought to reduce the

51    dimensionality of the data and/or perform feature selection to identify species (or other taxa) of

52    interest that are correlated with sample/community-level attributes (which we will refer to as

53    "phenotypic" attributes or "phenotypes") like host health status. Unfortunately, these

54    phenotype-associated species may co-occur with the same or different proportions across

55    samples within the same phenotype. Capturing these configurations is of interest to us, as we

56    contend it is more informative than merely finding specific taxa [1,2].

57    Nevertheless, obtaining meaningful configurations or subsets of taxa is often a daunting task.

58    These high-dimensional microbiome datasets include categorical and numeric features

59    associated with each sample.  These, in turn, may be linked to a set of taxonomic abundances

60    that are derived from clustering similar sequencing reads.  Typically, taxonomic markers, such

3

61    as variable regions of the 16S rRNA gene common to all prokaryotes, are used to perform the

62    clustering based on a fixed degree of sequence similarity among reads. Such clusters are termed

63    Operational Taxonomic Units (OTUs), and each OTU is usually assigned to some level of

64    taxonomy, such as a genus. Identifying OTUs correlating with specific sample features (e.g.,

65    body site, disease presence, diet, age) can be done via unsupervised exploratory methods [3].

66    Unfortunately, complexities inherent to taxonomic abundance data hinders many of these

67    methods. These complexities include vastly more OTUs relative to the number of available

68    samples [4], substantial sparsity in the OTU counts (absence of organisms in most samples), and

69    differences in sampling depth among samples. The sampling depth issue then requires

70    normalization, introducing additional challenges.  In particular, the normalization transforms

71    the abundances into relative abundances within each sample (compositional data)  [5,6].

72    Common approaches (e.g., differential abundance analysis  [3,7,8] and regularized regression

73    [9,10]) associate indicator taxa with sample information, leading to overly simplified biological

74    interpretations.

75    From an ecological perspective, co-occurring OTUs may represent related subcommunities of

76    taxa, which consist of OTUs that are common to (or overlap with) each sample.  This overlap is

77    due to taxa that covary with host or environmental factors; thus, identifying important

78    subcommunities (groups of taxa) and configurations of taxa (the grouping and ratios/relative

79    abundances of co-occurring taxa) may allow for a more biologically meaningful interpretation

80    than identifying indicator OTUs, because identifying subcommunities preserves the groupings

81    and abundances of taxa [2,11–13].  Developing techniques for identifying subcommunities is a

82    fundamental goal of this work.

83   Methods that predict functional profiles from 16S rRNA survey data usually report the overall

84   function of a sample and do not provide granularity on how each subcommunity provides

85   specific functions (Fig 1). Standard methods that predict function from 16S  rRNA survey data

86   include PICRUSt, Tax4fun, Piphillin, and SINAPS [14–17].  These simulate gene abundances

87   from the OTU relative abundance profile by assigning pre-existing gene ontologies, based on

88   whole genome sequences, to the OTUs.  The simulation is trivial for known microbes, but for

89   novel OTUs, gene content is interpolated through its neighbors' genes.  These are determined

90   via an unsupervised phylogenetic tree reconstruction. However, after the gene abundance

91   profiles are simulated for an entire sample, a user cannot view which functional content

92   associates with which taxa, nor how subcommunities contribute to function.

93

94   Fig 1. (Thematic Approach) Given a 16S rRNA gene abundance table, a topic model is used to

95   uncover the thematic structure of the data in the form of two latent distributions: the samples-

96   over-topics frequencies and the topics-over-OTUs frequencies. The samples-over-topics

97   frequencies are regressed against sample features of interest to identify the strength of a topic-

98   covariate relationship to rank topics (top). The topics-over-OTUs frequencies are used in a gene

99   function prediction (FP) algorithm to predict gene content. Important functional categories are

100  identified via a fully Bayesian multilevel negative binomial (NBR) regression model (middle).

101  The topics-over-OTUs distribution is hierarchically clustered to infer relationships between

102  clusters of co-occurring OTUs and topics (bottom). The result is the ability to identify key topics

103  that associate clusters of bacteria and their associated functional content to sample information

104    of interest. (Alternative Approach). A common alternative approach currently used in the

105    literature involves independently (1) characterizing the taxonomic configuration and (2)

106    predicting the functional configuration of the OTU abundance table. Gene function prediction is

107    performed on the full OTU abundance table, followed by a differential abundance analysis to

108    infer differences in specific genes between sample features of interest (top). The OTU table is

109    normalized to overcome library size inconsistencies and then analyzed via two methods: (1) an

110    elastic net (EN) to find sparse sets of OTUs that are predictive for the sample feature of interest

111    (middle) and (2) a multivariate (MV) analysis to identify relationships between beta diversity

112    and the sample feature of interest (bottom). The result are three analyses that summarize the

113    entire OTU relative abundance table, unlike the thematic approach, which characterizes co-

114    occurring sets of OTUs (configurations) in three ways.

115

116    We consequently have developed `themetagenomics`, a novel pipeline for analyzing 16S

117    rRNA amplicon surveys that (1) identifies subcommunities associated with specific sample

118    features and (2) uncovers functional profiles that further characterize these subcommunities.

119    We use a topic model approach to uncover subcommunity structure by estimating taxonomic

120    co-occurrence. Topic models are dimensionality reduction techniques that have had

121    considerable use in natural language processing to represent, as topics, co-occurrence

122    relationships between words from a corpus of documents. They have more recently shown

123    promise as a method for exploring taxonomic abundance data [2,18], where topics act as low-

124    dimensional representations of co-occurring sets of taxa given a set of samples, i.e., far fewer

125    topics than OTUs (Table 1). Unlike other dimensional reduction techniques common to

126    microbiome data analysis (e.g., principal coordinate analysis), topic models provide a new set of

127    features (topics) that should be familiar to microbiome researchers in that they have a form

128    similar to relative abundances: each sample is represented as a vector of frequencies across

129    topics and each topic is represented as a vector of frequencies across taxa. Lower dimensional

130    features that are also familiar may ease their interpretation.

131

132    Table 1. Relationship of Terms

| Topic Model | Pipeline | Description |
|---|---|---|
| Document | Sample | Collection of reads from subject $m$ at time $t$ |
| Topic | Topic | Collection of co-occurring taxa, subcommunity |
| Word | OTU, Gene, Taxa | Features from taxonomic abundance table or predicted functional content |
| Document-Level Covariate | Sample information, Sample class | Sample-level variable of interest – e.g., disease presence, diet, rainfall, time |
| $\theta$ | Samples-Over-Topics Distribution | Vector of topic frequencies in a given sample; probability of a topic occurring in a given sample |
| $\beta$ | Topics-Over-OTUs Distribution | Vector of OTU frequencies in a given topic; probability of an OTU occurring in a given topic |

133

8

134

135     Our pipeline aims to concisely summarize high-dimensional data in the form of OTU

136     abundances as low-dimensional sets of co-occurring taxa (topics) with their corresponding

137     predicted functional potential. When additional high-dimensional data is available (e.g.,

138     predicted gene function abundances), interpretability becomes increasingly difficult. Although

139     topic models have been applied to microbiome data because of their interpretable features, no

140     work has been done to leverage their interpretability to link low-dimensional representations of

141     OTU and predicted gene function abundances. In addition, little research addresses ways to

142     fully leverage the latent features topic models extract from microbiome data. For example,

143     correlated topic models [19] not only capture taxonomic co-occurrence but also topic co-

144     occurrence, such that the frequency of two topics, with different sets of co-occurring taxa,

145     occurring in any given sample, may be positively correlated. This is the basis of our novel

146     approach to exploit the correlation structure of topics across samples to resolve long-term

147     temporal behavior of subcommunities (represented as topics) in microbiome time-series

148     datasets.

149     Our approach at linking taxonomic composition to predicted functional content (obtained via

150     methods that leverage preexisting gene ontologies) within topics is unique.  We apply a recently

151     developed structural topic model (STM) [20] to a novel domain (16S rRNA amplicon surveys),

152     where each topic represents a cluster of co-occurring OTUs and each OTU can occur in multiple

153     topics with varying frequency. Functional content is then predicted within-topic, allowing the

154     topics to act as low-dimensional taxonomic and functional summaries of the input data. The

9

155    topics are then linked to sample-information that reflects host or environment status. Topics-of-

156    interest (e.g., those that contain differentially-enriched functional profiles) can easily be

157    identified in our pipeline via a fully Bayesian multilevel regression model. We also apply our

158    approach to empirical time-series data where we characterized events in terms of sets of

159    correlated topics to explore how the taxonomic configurations evolved over time.

160    Our pipeline has been implemented in the R package `themetagenomics`:

161    https://github.com/EESI/themetagenomics.

162

163    **Results and Discussion**

164

165    Here we explore the use of `themetagenomics` on publicly available datasets studying Crohn's

166    disease microbiota (Gevers et al. [21]), oral cancer microbiota (Schmidt et al. [22]), and the

167    variation of microbiota as a function of time (David et al. [23]).  With the larger Gevers et al.

168    Crohn's dataset, we validate the ability of `themetagenomics` to capture microbial profile

169    "signatures" (configurations of taxa which are groups with specific ratios/relative abundances

170    of co-occurring taxa). We show that (1) topics generalize well to test data not initially seen by

171    the model (generalizable topics are topics robust to overfitting, such that they avoid fitting noise

172    and thus can capture important signals representative of true taxonomic co-occurrence profiles),

173    and (2) topics capture distinct microbial signatures found in the original OTU relative

174    abundance data.

175    After validating the configuration of taxa within-topic (by assessing classification performance

176    to evaluate topic generalizability and OTU co-occurrence to evaluate topic quality) and the

177    configuration of predicted gene functions within-topic (via a permutation test using

178    metagenomic data), we assess the biological relevance of our low-dimensional summaries

179    (topics).  We then apply our complete pipeline to Gevers et al. to link a topic's functional

180    content, taxonomic co-occurrence, and sample information (clinical diagnosis of Crohn's disease

181    (CD)), and we compare these results to those obtained by the original authors. We compare our

182    results to those obtained by DESeq2 and an alternative topic-model based microbiome analysis

183    tool, BioMiCo [2]. We validate the functional prediction of our pipeline with the oral cancer

184    Schmidt et al. dataset by showing the low-dimensional topic profiles identified by

185    `themetagenomics` are also present in complementary metagenomic shotgun (MGS) sequence

186    data. We lastly implement our approach on time-series gut microbiome data from David et al.

187    We interpret the results in terms of topics and posterior uncertainty and compare our findings

188    to those obtained by a HC approach, as well as the results reported by David et al.

189

190    **Topic Modeling Feasibility and Generalizability**

191    We assess (1) if topics correlate to sample phenotypes (e.g., disease state) and (2) whether those

192    topics generalize well – that is, can the learned topics predict phenotypes from new data. Using

193    a random forest classifier, we compared the classification performance between two different

194    sets of predictors: (1) frequencies of topics-across-samples, $\theta$, from the STM, and (2) OTU

195    relative abundances across samples generated from QIIME [24]. For this analysis, we focused on

11

196     the Crohn's disease study from Gevers et al. given its large sample size (555 terminal ileum

197     samples).

198     To assess generalizability, we used a training/testing approach. We randomly selected 80% of

199     samples as our training set; the remaining 20% were set aside for testing (Table S1). Class labels

200     were binary, with positive (CD+) and negative (CD-) clinical diagnoses acting as the positive

201     and negative classes, respectively. For classifying CD diagnosis, we hypothesized that using

202     topics as predictors would outperform using relative abundances of OTUs, since the relative

203     abundance-based predictors are sparser, whereas topic modeling performs dimensionality

204     reduction, resulting in a relatively smaller set of topics that are less sparse relative to OTUs.

205     There was little difference between the topic model with at least 25 topics and the OTU table to

206     train the classifier (S1 Fig, Table S2). During testing, however, using topics as features

207     outperformed relative abundances, particularly in the F1 score, with relative abundances

208     achieving 80.8% and at least 25 topics achieving greater than 82.1% (Table S3). Using OTU

209     relative abundances as predictive features resulted in a larger proportion of false negatives,

210     which was likely due to its reliance on few, relatively rare taxa. Topics, on the other hand, are

211     less reliant on rare taxa because dimensionality reduction generates less sparse features (S2

212     appendix).

213     **Correlation Between Topics and Phenotype**

214     To identify topics of interest that were strongly associated with phenotype, we again

215     implemented `themetagenomics` on the Crohn's disease dataset, using the same binary

216     indicator for CD diagnosis as above. We then performed posterior inference. The primary

217    output of the topic model, as with any Bayesian analysis, is a posterior distribution of quantities

218    that estimate latent variables-of-interest (e.g., the frequencies of topics, $\theta$, in a particular sample)

219    given the observed data (e.g., OTU abundances). Posterior inference involves sampling these

220    latent variables-of-interest from the posterior distribution of the fitted topic model to calculate

221    expected means and assess uncertainty in those expectations.

222    With the posterior distribution, we identified topics-of-interest based on their "topic-sample-

223    effects" – the regression coefficients that represent differences in topic frequencies between CD+

224    and CD- samples. We performed permutation tests to ensure that detected topic-sample-effects

225    were not spurious (S2 appendix). For the model with 25 topics (K25), we performed 25

226    permutations, where we randomly permuted class label assignments (CD+, CD-), refit the topic

227    model, and calculated the mean regression coefficient for each topic. Of the 25 topics, 8 topics

228    had 95% uncertainty intervals for the effect size (differences between CD+ and CD-) that did not

229    span 0 (S2 Fig). We consider these "high-ranking-topics." Topics T15, T12, T2, and T14 had

230    estimates greater than 0 (implying robust associations with CD+), whereas topics T11, T25, T13,

231    and T19 had estimates less than 0 (implying robust associations with CD-). Increasing the

232    number of fitted topics gave similar results; for K75, 14 topics did not span 0 (S3 Fig).

233    We next tested how well a topic model (fit with the binary CD encoding) could capture the

234    severity of disease using the Pediatric Crohn's Disease Activity Index (PCDAI) associated with

235    CD+ that increases as CD severity increases (CD- samples were set to PCDAI=0). The frequency

236    of a sample containing a particular topic given its PCDAI is shown in Fig 2A for models K25

237    and K75. Topics are color-coded based on their association with CD, which is estimated using

13

238    their topic-sample-effects (yellow and violet represent topics most and least associated with CD,

239    respectively). Each overlapping line represents one of 25 replicate simulations. Both panels

240    demonstrate that as PCDAI increases, the thematic profile shifts from one dominated by a single

241    CD- associated topic (T8) to a set of CD+ topics (T12, T15, T45). The transition occurs at

242    approximately PCDAI=35. Because the K25 model had greater separation of high probability

243    topics, it will be the focus for the remainder of analyses involving Gevers et al. data.

244

245    Fig 2A. The relationship between topic frequency within a sample and that sample's Crohn's

246    Disease (CD) severity (PCDAI score) for the 25-topic STM. Each line represents the frequency of

247    a topic as a function of sample PCDAI score. High frequency topics are labeled. Violet and

248    yellow color-coded trajectories designate CD- and CD+ associated topics, respectively. Posterior

249    sampling was performed across 25 replicates, with each line plotted to represent the

250    distribution of the topic frequency trajectories. Fig 2B. Trajectories for the 75-topic model. Fig

251    2C. The relative abundance of OTUs in the (input) OTU relative abundance table for

252    "noteworthy" OTUs from high-ranking-topics. The left and right panels show the relative

253    abundance of these OTUs in each CD- and CD+ sample, respectively. Noteworthy OTUs are

254    defined as high-frequency OTUs, sampled from the posterior distribution, that concentrate into

255    high-ranking-topics (yellow=CD+ topic group, violet=CD- topic group, green=unassociated

256    topic group). The horizontal line marks a subset of samples that contain a large proportion of

257    the OTU profile associated with CD+ high-ranking-topics.

258

259     From the posterior topics-over-OTUs distribution (β) for the K25 model, we identified OTUs

260     highly associated with CD, that is, OTUs with high frequency in high-ranking-topics (CD+

261     associated topics T19, T13, T25, T11; CD- associated topics T14, T2, T12, T15) in more than 99%

262     of posterior samples (arbitrary threshold). We categorized these OTUs as CD+ associated OTUs,

263     CD- associated OTUs, and unassociated OTUs. Fig 2C shows the relative abundances of the 3

264     groups for each sample in the QIIME-generated OTU abundance table. Of CD+ samples (right

265     of vertical black bar), approximately 25% were characterized by a greater proportion of CD+

266     associated OTUs relative to CD- (marked by the horizontal black bar). The ratio of CD-

267     associated OTUs to unassociated OTUs had a similar distribution among CD+ and CD- samples,

268     suggesting that the OTU profile from CD+ high-ranking-topics is specific for the CD+ disease

269     status. Lastly, when we regressed PCDAI against the relative abundances of the CD+ associated

270     OTU profile, we found a significant positive relationship ($\beta$=0.057, p=0.01, 100 permutations),

271     albeit explanatory for only a small portion of the variation ($R^2$=8.64%), suggesting that presence

272     of this OTU profile may be weakly indicative of severe cases of CD (S2 appendix).

273     **Comparison to BioMiCo**. We compared our approach's performance to BioMiCo, a topic model

274     that identifies meaningful sets of "assemblages" (analogous to topics – i.e., sets of cooccurring

275     taxa) by directly incorporating sample- or environmental level features (labels) during the

276     training procedure. It is fully supervised and assumes that a sample is comprised of a mixture

277     of communities that share sample- or environmental level features. These communities are

278     described by a set of high probability assemblages which are in turn described by a set of high

279     probability taxa.

15

280   We fit BioMiCo using 25 and 50 assemblages and compared its ability to distinguish CD from

281   control using held-out testing data (same train/test splits as described previously) and then

282   compared these results to the prediction performance of the STM. Testing performance was

283   similar between the two approaches (Table S3, S6). The balanced accuracy was highest for the

284   25-topic STM model, but the STM's performance varied as a function of topic number. F1 score,

285   however, was much worse for BioMiCo due to its low precision.

286   For the 25-assemblage model, there were roughly four assemblages with high posterior

287   probability for CD samples and low posterior probability for controls. If we focused on the taxa

288   with the top-10 highest posterior probability of belonging to these assemblages, no more than 2

289   taxa were present in the top-10 highest probability taxa in the STM's CD-topics that were most

290   associated with CD, suggesting little correspondence between the composition of assemblages

291   and topics. Alternatively, when focusing on assemblages with high posterior probability for

292   control but not CD, one assemblage had 4 genera in common with the STM's topic 13:

293   *Parabacteroides*, *Bacteroides*, *Ruminoccous*, and *Roseburia*.

294   It is worth noting, however, that the STM and BioMiCo aim to characterize data differently and

295   hence the distribution of taxa within a given topic are expected to be different. Still, both

296   approaches show they similarly generalize to new data. An advantage of `themetagenomics` is

297   that it leverages output inherent to the design of the STM that is not available via BioMiCo,

298   notably topic-topic correlation. Also, the STM is appreciably faster, taking minutes to run on the

299   Gevers data whereas BioMiCo took days. Unlike BioMiCo – as well as the STM which is aimed

300   for more general use – `themetagenomics` delivers a framework that facilitates ease-of-use

301    microbiome analysis using a topic model via an R package with a variety of intuitive functions

302    for preprocessing, analyses, and visualizations. It also provides novel downstream approaches

303    such as time series analysis which leverages the STM's estimation of topic-topic correlation, as

304    well as methods to associate a topic's taxonomic composition to its predicted gene functions.

305

306    **Linking Function to Taxonomy with Topics**

307    We wanted to discern whether the topics would continue to identify meaningful relationships

308    upon introducing another layer of information: predicted function (via abundances of metabolic

309    pathways). Consequently, we applied our full `themetagenomics` pipeline to the Crohn's

310    disease dataset and compared our findings to those of the original authors. To further

311    characterize topics, we applied PICRUSt to the topics-over-OTUs distribution, $\beta$, to predict the

312    functional gene content within topics. The genes were then annotated in terms of their KEGG

313    functional hierarchy designation [25], thereby providing each gene with a metabolic pathway

314    label. We then performed a fully Bayesian multilevel regression analysis on the predicted

315    abundances of each gene to identify strong topic-pathway interactions.

316    Like Gevers et al., we identified an increase in membrane transport associated with CD+

317    subjects' gut microbiome; however, using `themetagenomics`, we were able to pinpoint the

318    specific topics associated with the enrichment of these functional categories, T2 and T12 (Fig

319    3A). We then could link enrichment of membrane transport genes to the taxa that were also

320    enriched in this topic. For example, topics T2 and T12 were dominated by Enterobacteriaceae.

321    These Enterobacteriaceae-enriched topics were also enriched for siderophore and secretion

17

322    system related genes. Like T2 and T12, T15 was highly associated with CD+; however, it was

323    less enriched for membrane transport genes. This suggests that the cluster of bacteria found in

324    T15 (*Haemophilus* spp., *Neisseria*, and *Fusobacteria*) may have contributed less to the shift of

325    transport genes reported by Gevers et al. and instead have distinct pathway associations with

326    CD.

327

328    Fig 3A. Level-3 pathway category-topic interaction regression coefficients from the multiple

329    level negative binomial model. Red asterisks indicate estimated pathway-topic interaction

330    weights that do not span 0 at 80% uncertainty (pathways lacking robust interactions are

331    omitted). Green=large positive coefficients thus enrichment for that pathway in that topic,

332    Violet=large negative coefficients thus depletion for that pathway in that topic. Topics are

333    ordered from CD- associated (left, T19) to CD+ associated (right, T15). High-ranking-topics are

334    delineated by the vertical dotted lines (CD-: T19-T11; CD+: T14-T15). Fig 3B. Volcano plot

335    showing DESeq2 results for differentially abundant predicted level-3 KEGG categories.

336    Functions were predicted using PICRUSt on the copy number normalized OTU abundance

337    table. Blue points represent categories significantly enriched for CD- and red points are

338    categories enriched for CD+, respectively. Gray points are categories with p-values greater than

339    0.1 after Bonferroni correction.

340

341    The strongest topic-pathway interaction was found in T19 for genes encoding bacterial motility

342    proteins. For T19, three motility-related pathways (bacterial motility proteins, bacterial

18

343    chemotaxis, flagellar assembly) had topic-pathway interactions that did not span 0 at 80%

344    uncertainty, suggesting that T19 was more enriched in cell motility genes relative to all other

345    topics. The pathways inferred from T19 are consistent with this taxonomic profile, which

346    consisted of motile bacteria belonging to Lachnospiraceae, Roseburia, and Clostridiales.

347    Enrichment of two lipopolysaccharide (LPS) synthesis categories were associated with CD+

348    topics; however, one of these categories was specific for only T15 (Table S4).

349    **Comparison to DeSeq2.** We compared the topics' functional profiles to the results obtained by

350    performing a DESeq2 differential abundance analysis on functional predictions obtained by

351    applying PICRUSt to the QIIME-generated OTU abundance table. Of the 160 (level-3) KEGG

352    pathway categories, more than half (87) were found significant ($\alpha < 0.1$) in the DESeq2

353    approach, despite using Bonferoni correction (a conservative approach to correct for multiple

354    comparisons), complicating interpretation (Fig 3B). Despite minor differences in specific

355    pathway enrichment between `themetagenomics` and the DESeq2 approach (S2 appendix) the

356    major difference was the greater number of low-uncertainty/significant pathway categories

357    found by DESeq2. While one could reduce the significance level when applying DESeq2 to

358    achieve a smaller subset of significant pathway categories, the choice is arbitrary. Moreover, the

359    predicted functional abundances (via PICRUSt, Tax4fun, etc.) are scaled based on the

360    abundance of taxa from which they were derived. Thus, high taxonomic abundances will often

361    yield high functional abundances. Many of the significant pathway categories identified by

362    DESeq2 may be driven by a small subset of highly abundant taxa. `Themetagenomics`, on the

363    other hand, first groups co-occurring taxa into topics. Because functional prediction is

364    performed within a topic, taxa that are highly abundant in the input OTU abundance table can

19

365   only affect the topics in which they are present at high frequency. Thus, this prevents high

366   abundance taxa associated with a subset of samples (e.g., CD+), and their corresponding

367   predicted pathway abundances, from disproportionately influencing the statistical significance

368   of these pathways.

369

370   **Validating the Functional Predictions of `Themetagenomics` via `Paired MGS Samples`**

371   Using sample-matched (N=12) oral cancer microbiome samples from Schmidt et al. that

372   underwent both 16S rRNA amplicon sequencing and metagenomic shotgun sequencing, we

373   verified enrichment or depletion of predicted functional content (collapsed into metabolic

374   pathway categories) of the `themetagenomics` pipeline. The pipeline processed the 16S rRNA

375   samples and compared the results to metagenome-based gene functional abundance data. Fig

376   4A shows the relative enrichment/depletion of various topic-pathway combinations identified

377   by `themetagenomics`. For example, bacterial motility genes were enriched in topic 25

378   (positive coefficient, shaded green), whereas bacterial motility genes were depleted in topics 3

379   and 9 (negative coefficients, shaded violet).

380

381   Fig 4A. KEGG (level-3) pathway category-topic interaction regression coefficients from the

382   multilevel negative binomial model as a measure of association between pathway and topic.

383   Only pathways present in both the themetagenomics analysis of 16S rRNA data and

384   HUMAnN2 analysis of the metagenomics shotgun sequencing data are shown.

20

385  Green=associated samples with positive cancer diagnosis, Purple=associated with healthy

386  samples. Fig 4B. Pathway category-topic interaction regression coefficients for metagenomic

387  data. Topics were generated based on KOs that belonged to high frequency taxa in the

388  themetagenomics pipeline. Fig 4C. Example topic-pathway heatmaps, similar to Fig 4A and 4B

389  from four of the 100 permuted metagenomic datasets using in the permutation test. Fig 4D.

390  Distribution of root-mean-squared-error (RMSE) scores (between the topic-pathway interaction

391  regression coefficients between themetagenomics and the metagenomic data) from the 100

392  permuted metagenomic datasets. The RMSE score (0.56) for the unpermuted metagenomic

393  dataset is delineated by the red dotted line.

394

395  To compare the results from `themetagenomics` to gene function abundances inferred from

396  metagenomic shotgun sequencing for each topic, we first identified high frequency taxa (those

397  with frequencies greater than 1% in that topic) then identified all reads belonging to these taxa

398  in the metagenomic shotgun data. To identify pathway-topic enrichment/depletion, we then

399  applied a multilevel regression model. The results indicate that the taxa belonging to a topic are

400  associated with an enrichment/depletion of genes present in the shotgun data (Fig 4B). Notably,

401  LPS biosynthesis proteins and porphyrin metabolism pathways were depleted in multiple

402  topics in both sets of results. The relative enrichment/depletion of phosphotransferase system

403  genes was also similar.

404  We performed a permutation test to determine whether the similarities in gene

405  enrichments/depletions between `themetagenomics` and the metagenomic data were spurious.

21

406      We randomly permuted the topic and gene pathway labels in the metagenomic data, refit the

407      multilevel regression model, and then calculated the root mean square error (RMSE) for each

408      topic-pathway interaction regression weight between the `themetagenomics` and permuted

409      metagenomic models. After 100 replicate simulations, the RMSE for the unpermuted

410      metagenomic model was smaller than every permuted metagenomic model ($p < 0.05$) (Fig 4C-

411      D). Therefore, the apparent similarities in the gene enrichment/depletion profiles between

412      `themetagenomics` and the shotgun data were not due to random chance, indicating that

413      using predicted gene enrichment/depletion from 16S rRNA amplicon surveys resulted in

414      similar within-topic predicted functional profiles to those obtained by directly measuring

415      functional content via metagenomic shotgun sequencing.

416

417      **Detection of Events in Subject B from David et al.**

418      The David et al. dataset contains fecal and salivary 16S rRNA gene surveys from two subjects.

419      We focused on fecal samples from subject B. We compared our results to the three profiles

420      described by David et al., which consisted of a pre-food-poisoning profile (days 1-150), food-

421      poisoning profile (151-159), and post-food-poisoning profile (150-318).

422      **The topic model approach identified 3 distinct gut configurations.** In the topic correlation

423      network (Fig 5A), we identified a small subnetwork of three topics (marked by violet bracket)

424      and two large subnetworks that contained 24 and 14 topics each (red and green brackets,

425      respectively). The large subnetworks were connected by a chain of four topics (T9, T24, T2, T37)

426      (blue bracket). We defined the four sets of correlated topics as topic clusters and sampled topic

22

427     frequencies (across samples) and taxa frequencies (across topics) from the topic model's

428     posterior distribution to assess how often topics and taxa occurred within these clusters.

429

430     Fig 5 Application of the topic model approach to David et al. data. (A) The topic-to-topic

431     correlation graph showing two topic clusters (clusters 1 and 3) connected by a linear chain of

432     topics (cluster 2) that follow the time progression of the taxonomic change due to the food

433     poisoning infection. (B) Distribution of topic assignments as a function of day and cluster

434     (panels), indicating 3 distinct profiles. The interval in which food poisoning symptoms

435     presented (per David et al.) are marked with dotted vertical lines. Gray shading indicated 80%

436     uncertainty intervals. (C) Frequency of cluster assignments as a function of day, indicated day

437     153 marking the shift from profiles 1 to 2 and day 159 marking the shift from profiles 2 to 3. (D)

438     Frequency of taxa assignments given a cluster assignment. Cluster 2 is shown in terms of its

439     topics (9, 24, 2, 37). Topic 20 is also shown (misc. cluster), which lacked any edges in the

440     correlation graph, but marks the initial appearance of *Enterobacteriaceae* on day 153

441     (representing the start of the infection). (E) The probability of the topic assignments given each

442     day for cluster 2. The progression of topics also follows the progression of taxonomic change

443     shown in the correlation graph.

444

445     Fig 5B shows the posterior frequency in which the topic clusters occurred given the day in

446     which the sample was collected (the estimated posterior probability of a cluster occurring on a

447     given day). There were two clear periods of rapid change in cluster frequency, specifically when

23

448     transitioning from cluster 1 to 2 (days 152-154) and clusters 2 to 3 (day 161). Our intervals are

449     similar to the original study's transition points at days 144-145 and 162-163, where the shift

450     from a topic cluster 1 to topic cluster 2 corresponded with subject B's food poisoning diagnosis.

451     The transition between topic clusters 1 and 2 is abrupt and likely occurred around day 153.

452     Taxonomically, this transition is marked by a shift from Bacteroideaceae (posterior

453     frequency=0.338), Lachnospiraceara (0.276), and Rumunococcaceae (0.266) to Enterbacteriaceae

454     (0.246) and Clostridiaceae (0.195) families (Fig 5D). In particular, day 153 was distinctive for

455     topic 20. This rare topic was not correlated with any other topics and hence did not belong to

456     any topic cluster. While its taxonomic profile was quite similar to topic cluster 1, it was

457     distinctly enriched for *Enterobacteriaceaea spp.*, which is consistent with the subject's *Salmonella*

458     diagnosis. Topic 20 likely marks the event of initial exposure to the pathogen.

459     The distribution of topic assignments for topic cluster 2 followed the order in which its topics

460     were positioned in the topic correlation network (the linear chain of topics) (Fig 5E). The start of

461     topic cluster 2, day 155, was dominated by topic 9, characterized by taxa substantially different

462     from topic cluster 1. Bacteria enriched in this topic included *Haemophilus parainfluenzae*,

463     *Clostridium perfringens*, and, notably, *Enterobacteriaceaea spp*. Thus, topic 9 likely represented the

464     disrupted configuration of microbiota due to exposure to *Salmonella*. Enterbacteriaceae spp. and

465     *C. perfringens*, via topic 24, continued to dominate on day 156. Day 157 was best described by

466     topic 2, a topic rich in *Enterobacteriaceae spp.* as well as *Veillonella spp*. It should be noted,

467     however, that our results were more conservative than David et al. in that we confidently

468     estimated that topic cluster 2 lasted roughly 4 days (155 to 158), which is much shorter than the

469     original study's estimate (145 to 162). Our estimated length of illness (153 to 158) was more

24

470     consistent to David et al. (151 to 159), however. At approximately day 159, the taxonomic profile

471     shifted toward cluster 3, which was similar to cluster 1 in terms of Bacteroidaceae (0.369),but

472     enriched in Lachnospiraceae (0.360) and depleted in Rumunoicoccaceae (0.165) (Fig 5D).

473

474     **HC was unable to separate the transition between during- and post-illness periods.** We

475     compared our approach to one using HC. HC cluster 4 contained 360 taxa and corresponded

476     well to the pre-illness period, spanning days 1 to 150. The set of taxa was similar to the taxa

477     identified in topic cluster 1 (S4 Fig). The post-illness period was captured by HC clusters 1 and

478     3, but these clusters failed to completely separate the during- and post-illness periods; they

479     spanned days 151 to 318.

480

481     **Limitations**

482     There are limitations to our approach. First, the topic-pathway inference step currently scales

483     poorly in terms of computation time for large numbers of topics, which may be more important

484     as datasets grow. Regularization and sparsity-inducing priors help limit the number of

485     important topics; hence, exploring only a subset of topics during the final regression step can

486     offer substantial speed improvements at little cost, but utilizing the complete set of topic

487     information would be ideal. Second, we are capable of separately estimating the uncertainty in

488     our topic model, the multilevel regression model, and the functional predictions from PICRUSt,

489     but we currently do not propagate the uncertainty throughout the pipeline. Doing so would

25

490    improve downstream interpretation with better estimation of the uncertainty in topic-sample

491    covariates and topic-pathway interactions, which in turn would greatly improve one's

492    confidence in focusing on within-topic gene sets. Third, we do not incorporate phylogenetic

493    branch length information, which could lead to more meaningful topics.

494

495

496    **III. Conclusion**

497    We present our approach at a time when easily-to-interpret analyses for complex microbiome

498    data are direly needed. Current methods often link the relative abundance of a single OTU to a

499    sample information of interest (e.g., disease state). These methods routinely identify important

500    subsets of taxa but ignore OTU co-occurrence and ratios. Network methods can overcome this

501    concern, but typically don't incorporate phenotypic information within the model;

502    consequently, they are incapable of directly linking sections of the OTU correlation network

503    with sample metadata of interest. Constrained ordination methods, such as canonical

504    correspondence analysis, do in fact couple inter-community distance with sample information,

505    but the user is limited to specific distance metrics (e.g., Chi-squared) and must follow key

506    assumptions (e.g., the distributions of taxa along environmental gradients are unimodal) [26].

507    Moreover, interpretation of biplots becomes increasingly difficult as more covariates are

508    included. While linking key taxa to functional content can be accomplished via sparse canonical

509    correlation analysis [27], this approach is susceptible to many of the interpretability problems

510    found in other ordination approaches, and exploring inferred relationships in the context of

511    taxonomic co-occurrence is not straightforward.

512    The ability to make meaningful inferences using current methods is further limited by the fact

513    that microbiome data is often inadequately sampled (thus justifying some type of normalization

514    procedure), compositional (due to normalization), sparse, and overdispersed. Thus, recent work

515    has explored the use of Dirichlet-Multinomial models, which are well equipped at managing

516    overdispersed count data [28–30]. The fact that Dirichlet-Multinomial conjugacy is exploited in

517    many topics models hints at their applicability for relative abundance data. We selected the

518    recently developed STM for our workflow because of its ability to not only utilize sample data

519    as prior information as in the Dirichlet-Multinomial regression topic model [31], but also

520    capture topic correlation structure and apply partial pooling over samples or regularization

521    across regression weights.

522    Thus, we have proposed an approach for uncovering latent thematic structure in 16S rRNA

523    amplicon data that provides a low-dimensional, biologically interpretable representation of

524    taxonomic and predicted functional content. Rather than inferring functional content

525    independently of taxonomic relative abundances, our approach shifts the focus to investigating

526    within-topic functional content. Unlike other methods, by exploring our topics, we can link

527    categories of functional content to specific clusters of taxa which can in turn be linked to sample

528    features of interest. For example, like Gevers et al., we detected a relationship between

529    membrane transport genes and CD+, but our approach also allowed us to determine which

530    bacteria (OTUs belonging to Enterobacteriaceae) were the prime contributors to the enrichment

27

531    of membrane transport genes. Moreover, the pathogenic set of bacteria reported by Gevers et al.

532    (*Haemophilus* spp., *Neisseria*, and *Fusobacteria*) contributed less to the predicted abundance of

533    membrane transport genes. By independently applying statistical approaches to the OTU and

534    predicted functional content, as is typical, the apparent relationship between membrane

535    transport genes and specific configurations of bacteria would be lost.

536    We have also shown that our approach drastically reduces the dimensionality of two high-

537    dimensional sources of information, taxonomic relative abundances and predicted functional

538    content, increasing the ease in which these data can be interpreted. For Gevers et al., we

539    determined that T15 is (1) associated with CD+ samples; (2) dominated by a cluster of bacteria

540    previously associated with CD; and (3) uniquely enriched for a subset of LPS synthesis genes.

541    With a gene profile from a topic of interest, one could focus on gene subsets associated with

542    topic-specific bacterial clusters that are known disease biomarkers, which in turn may facilitate

543    targeted approaches for future research endeavors.

544    Lastly, our complete pipeline is computationally manageable. Fitting the topic model to a

545    dataset with nearly 5000 samples reached convergence in minutes. Functional prediction via

546    PICRUSt also only takes minutes (using our C++ implementation in `themetagenomics`).

547    Inferring topic-pathway interactions via our multilevel, negative binomial regression approach

548    is comparatively slower, however, taking hours for large datasets. However, this is still

549    manageable. Thus, we offer a viable package that can help researchers discover configurations

550    of taxa and functions that correlate to sample metadata. This is because we implement this

551    model in the probabilistic programming language Stan, which uses Hamiltonian Monte Carlo.

28

552    Maximum likelihood (a much faster alternative) does not provide estimates of uncertainty and

553    generally fails to converge for these data, although the regression weight estimates tend to be

554    quite similar based on our experience.

555

556

557    **Methods**

558

559    **Review of the Structural Topic Model**

560    The STM [20] is a Bayesian generative topic model. It begins with a given a set of M samples,

561    each consisting of N OTUs. These N OTUs are, in turn, elements of a fixed vocabulary of V

562    unique OTU IDs. From this, K (a fixed number chosen a priori) latent topics are assumed to be

563    generated from the data. These topics consist of overlapping sets of co-occurring OTUs. Note

564    that we will describe the STM in the context of the analyses perform herein; for a complete

565    description of the STM, see [20]. The observations include the presence of OTU $w_n$ occurring in

566    sample m and an $M \times P$ matrix of sample-level information such as disease state or age.

567    For our purposes, the posterior distribution of unobserved (latent) parameters given the

568    observed data is given by:

569                    Posterior Distribution:    $p(\theta,\beta,\Sigma,\Gamma,z \mid w,X)$.

570    The generative process is formulated by first specifying the probability

29

571
$$P(\text{Topic } k \text{ occurs in Sample } m) = \theta_{m,k}, \sum_{k=1}^{K} \theta_{m,k} = 1$$

572   and, for each of the samples, is assumed to follow logistic normal distributions,

573
$$\theta \sim LN_{K-1}(\Gamma^T X_m^T, \Sigma)$$

574   where $\Gamma$ is a $P \times (K-1)$ matrix of regression coefficients that estimate the degree of influence a

575   covariate $X_p$ has on $\theta$; and $\Sigma$ is a $K \times K$ covariance matrix. In addition to $\theta$, the probability

576
$$P(\text{OTU } n \text{ occurs in Topic } k) = \beta_{k,n}, \qquad \sum_{n=1}^{N} \beta_{k,n} = 1$$

577   For each topic, $\beta_k$ is assumed to be Dirichlet distributed. Finally, both topic assignments $z_{m,n}$ for

578   each OTU $w_{m,n}$, along with each OTU, obey multinomial distributions,

579
$$z_{m,n} \sim \text{Multinomial}(\theta_m)$$
$$w_{m,n} \sim \text{Multinomial}(\beta, z_{m,n})$$

580       For the relationships between topic model nomenclature and our terminology, see Table

581   1. The posterior distribution is estimated by a semi-collapsed variational expectation

582   maximization procedure. Convergence is reached when the relative change in the variational

583   objective (i.e., the estimated lower bound) in successive iterations falls below a predetermined

584   tolerance.

585

586   **Empirical Datasets**

587   The Gevers et al. dataset (PRJNA237362, 03/30/2016) is a multicohort, IBD dataset that includes

588   16S rRNA amplicon data from control, CD, and ulcerative colitis samples taken from multiple

589   locations throughout the gastrointestinal tract [21]. The Schmidt et al. dataset (PRJEB4953,

590   08/14/2017) consists of human oral microbiota obtained from control subjects and subjects

591   diagnosed with oral cancer. These samples underwent 16S rRNA amplicon sequencing, and a

592   subset (N=12) also underwent metagenomic shotgun sequencing.

593

594   **16S rRNA Amplicon Data Preparation and OTU Picking**

595   Paired-end reads were joined and quality filtered via QIIME v 1.9.1 and dada2 for Gevers et al.

596   and Schmidt et al. data, respectively. Closed-reference OTU picking was performed with QIIME

597   using SortMeRNA against GreenGenes v13.5 at 97% sequence identity. This was followed by

598   copy number normalization via PICRUSt version 1.0.0 [32]. Samples with fewer than 1000 total

599   reads were omitted. OTUs that lacked a known classification at the phylum level were removed.

600   For Gevers et al., we selected only terminal ileum samples and filtered OTUs with fewer than 10

601   total reads across samples, yielding 555 samples over 1500 OTUs. For Schmidt et al., we filtered

602   any OTU with non-zero abundances in fewer than two samples, yielding 81 samples over 1029

603   OTUs.

604

605   **Metagenomic Shotgun Sequence Data Preparation and Functional Genomic Profiling**

606   Low quality reads and human genomic sequences were filtered via KneadData. Functional

607   profiles were then generated using HUMAnN2 with the ChocoPhlAn nucleotide database and

608   UniRef90 protein database. The UniRef90 protein families were collapsed into KEGG

31

609    orthologies (KOs), yielding abundances (copies per million (CPM)) for 12 samples over 36,806

610    KOs.

611

**Structural Topic Model Fitting**

613    The OTU abundance tables consisted of counts normalized by 16S rRNA gene copy number. No

614    other normalization was performed based on the simulation results in [33]. STMs with different

615    parameterizations in terms of topic number (K ∈ 15, 25, 50, 75, 100, 150, 250) and sample

616    features (e.g., no features, indicators for presence of disease, diet type, etc.) were fit to the OTU

617    tables generated from Gevers et al. data via the R package $stm$ [34]. We evaluated each model

618    fit for presence of overdispersed residuals and conducted permutation tests (permTest in the

619    stm package) where the sample feature of interest is randomly assigned to a sample prior to

620    fitting the STM. To compare parameterizations between models, we evaluated predictive

621    performance using held-out likelihood estimation [35].

622

**Assessing Topic Generalizability**

624    We performed classification to assess the generalizability of the extracted topics. No sample

625    information was used as covariates in the logistic normal component of the STM. Samples were

626    split into 80/20 training-testing datasets. For different number of topics (K ∈ 15, 25, 50, 75, 100,

627    150), an STM was trained to estimate the topics-over-OTUs distribution (β). We then held this

628    distribution fixed; hence, only the testing set's samples-over-topics distribution (θ) was

32

629    estimated. For both the training and testing sets, simulated posterior samples from the samples-

630    over-topics distribution (θ) were averaged. The resulting posterior topic frequencies in the

631    training set were then used as features to classify sample labels, similar to using $\overline{Z}$ in supervised

632    LDA [36]. Generalization (testing) error was assessed using the optimal parametrization based

633    on cross-validation performance on the test set topic frequencies. Classification was performed

634    using a random forest classifier, which underwent parameter tuning to determine the number

635    of variables for each split. This was accomplished through repeated (10x) 10-fold cross-

636    validation, using up-sampling to overcome class imbalance. We performed a parameter sweep

637    over the number of randomly selected OTU features, while setting the number of trees fixed at

638    128. The optimal parameterizations were selected based on maximizing ROC area under the

639    curve.

640    The performance of the STMs was compared to the performance using OTUs as features from

641    the original OTU abundance table. Separately, training and testing set OTU abundances were

642    converted to relative abundances with the following equation: $OTU_{n,m}/\Sigma_n OTU_{n,m}$. In words,

643    OTU *n* for sample *m* is scaled by the library size of sample *m* (the total abundance of sample *m*).

644    The resulting OTU relative abundance tables were separately z-score normalized. Training

645    cross-validation and testing using a random forest was then performed as above.

646

647    **Identifying Within-Topic Clusters of High Frequency OTUs**

648    Using the topics-over-OTUs distribution, we performed hierarchical clustering via Ward's

649    method on Bray-Curtis distances. We refer to high frequency groups of OTUs as "clusters."

650

**Inferring Within-Topic Functional Potential**

652    We obtained the topics-over-OTUs distribution (β) for each fitted model and mapped the

653    within-topic OTU probabilities to integers ("pseudo-counts") using a constant: $10000 \times \beta$. A

654    large constant was chosen to prevent low frequency OTUs from being set to zero, although their

655    contribution to downstream analysis was likely negligible. Gene prediction was performed on

656    each topic-OTU pseudo-count table using PICRUSt version 1.0.0 [14]. (Normalization of 16S

657    copy number was performed prior to topic model fitting using PICRUSt.) Predicted gene

658    content was classified in terms of KOs [37].

659

**Identifying Topics of Interest**

661    Topics of interest were identified using the samples-over-topics distribution, where each

662    column represents the frequency of topic $k$ for each sample. Each column was regressed against

663    CD diagnosis. We calculated 95% uncertainty intervals using an approximation that accounts

664    for uncertainty in estimation of both the sample covariate coefficients and the topic frequencies.

665    We refer to these coefficients as "topic-sample-effects." Coefficients whose 95% uncertainty

666    intervals do not span 0 are referred to as "high-ranking-topics."

667

**Validating Within-Topic Co-Occurrence**

669    To determine how well the high-ranking-topics captured co-occurrence in the original OTU

670    relative abundance table, we sampled the top-10 highest frequency taxa in each high-ranking

671    topic's topics-over-OTUs distribution (β). We then normalized the original OTU table using the

672    centered-log-ratio transformation and then evaluated how the high frequency taxa vary as a

673    function of CD diagnosis and PCDAI.

674

675    **Posterior Inference**

676    To determine how well the high-ranking-topics captured the taxonomic profile associated with

677    CD, we performed the following posterior simulation over R=1000 iterations. First, for iteration

678    r, for all samples $m \in M$ (e.g., subject 134), we obtained 100 posterior samples ($i \in \{1,..., 100\}$) of

679    $\theta_m^{(i)}$ from the posterior distribution, $p(\theta,\beta,\Sigma,\Gamma,z \mid w,X)$. For each of these $\theta_m^{(i)}$, we sampled topic

680    assignments $z_{m,n}^{(i)} \sim \text{Multinomial}(\theta_m^{(i)})$, and then OTUs $\hat{w}_{m,n}^{(i)}|z_{m,n}^{(i)} \sim \text{Mulinomial}(z_{m,n}^{(i)},\beta)$.

681    We then recorded whether the topic assignments $z_{m,n}^{(i)}$ belonged to one of the high-ranking-

682    topics and whether they have a positive or negative association with sample covariates of

683    interest, resulting in positive-, negative-, and no-association topic groups. We calculated the

684    frequency $f_n^{(g)}$ in which OTUs $\hat{w}_{m,n}^{(i)}$ were sampled from a given topic group g:

685
$$f_n^{(g)} = \sum_i \sum_{\hat{w}_{m,n}^{(i)}|z_{m,n}^{(i)}} 1[z_{m,n} \in g]$$

686    where $1[\cdot]$ is the indicator function. For each OTU, we calculated which group had the largest

687    sampling frequency:

35

688
$$f^{(g)}_n{}^* = 1\left[f^{(g)}_n = \operatorname*{argmin}_g f_n\right]$$

689 After 1000 iterations, we calculated

690
$$F^{(g)}_n{}^* = \frac{1}{R}\sum_r f^{(g)}_n{}^{*\,(r)}$$

691 For each topic group, we extracted a subset of OTUs that had frequencies above 0.99. In the

692 original relative abundance table, for each sample, we calculated the relative abundance of each

693 group of OTUs.

**Identifying Functional Content that Distinguishes Topics**

695 To determine which predicted functional gene content best distinguished topics, we used the

696 following multilevel negative binomial regression model:

697
$$\theta_{k,c} = \exp\left[\mu + \beta_k + \beta_c + \beta_{k,c}\right]$$

698
$$y_{k,c} \sim \mathrm{NB}(\theta_{k,c}, \lambda)$$

699 where $\mu$ is the intercept, $\beta_k$ is the per topic weight, $\beta_c$ is the per level-3 gene category weight, $\beta_{k,c}$

700 is the interaction weight for a given topic-function (gene category) combination, $y_{k,c}$ is the count

701 for a given topic-function combination, and $\lambda$ is the dispersion parameter. The intercept $\mu$ was

702 given a Normal(0, 10) prior; all weights were given Normal(0, 2.5) priors; and the dispersion

703 parameter $\lambda$ was given a Cauchy(0, 5) prior. Model inference was performed using Hamiltonian

704 Monte Carlo in the R package rstanarm [38]. Convergence was evaluated across four parallel

705 chains using diagnostic plots to assess mixing and by evaluating the Gelman-Rubin

706 convergence diagnostic [39]. To reduce model size, we used genes belonging to only 15

36

707    (arbitrary number) level-2 KEGG pathway categories (Table S5). For large topic models, we fit

708    only the top 25 topics, ranked in terms of topic-sample-effects that measure the degree of

709    association between samples-over-topics probabilities and our sample feature of interest.

710

711    **Assessing Relationships Between Sample Information of Interest and Taxonomic Relative**

712    **Abundance**

713    To quantify the relationship between taxonomic relative abundance and continuous sample

714    features (such as PCDAI), we performed negative binomial regression (log-link), using sample

715    library size (sum of OTU abundances across samples) as an offset. The family-wise error rate

716    was adjusted via Bonferroni correction. Significance levels for hypothesis testing was set at 0.05.

717

718    **Comparing Within-Topic Functional Profiles to an OTU-Relative-Abundance-Based**

719    **Approach**

720    We compared the results from the hierarchical negative binomial model to a differential

721    abundance approach. We performed predicted functional content using PICRUSt on copy

722    number normalized OTU abundances. The resulting functional abundances were collapsed into

723    level-3 KEGG pathways. Note that, for consistency, we again restricted the KOs to the 15 level-2

724    KEGG pathways used previously. The resulting level-3 pathway abundances underwent

725    DESeq2 differential abundance analysis, which uses negative binomial regression and variance

726    stabilizing transformations to infer the difference log-fold change of OTU relative abundance

37

727    [7,8]. The resulting p-values were corrected via the Bonferroni method. Adjusted p-values

728    below 0.1 were considered significant.

729    **Fitting BioMico**

730    The same training and testing sets were used as described above. Assemblages of 25 and 50

731    were trained with default parameters unless specified: burnin=5000, delay=500 (25 assemblages)

732    or delay=100 (50 assemblages), rarefaction_depth=1000. Parameters were adjusted to decrease

733    training time to less than 3 days. Posterior distributions were evaluated to ensure MCMC

734    convergence.

735

736    **Validating Extracted Functional Profiles using Metagenomic Shotgun Sequencing Data**

737    The themetagenomics pipeline was applied to the Schmidt et al. OTU table: (1) data were

738    normalized for 16S rRNA gene copy number; (2) normalized OTU abundances were fit using a

739    25 topic STM with cancer diagnosis as a binary covariate; (3) within-topic functional content

740    was predicted using PICRUSt; and then (4) topic-pathway effects were detected using the

741    multilevel regression model.

742    For each topic, we identified the high probability OTUs (those with frequencies greater than 1%

743    in that topic), obtained their genus classification, and then subset the metagenomic KO table

744    such that only KOs corresponding to these genera are present. Then, for each level-3 KEGG

745    pathway, we summed the abundances of all remaining KO members. Topic-pathway effects

746    were then detected with the following multilevel regression model:

38

747
$$\theta_{k,c} = \exp\left[\mu + \beta_1 X + \beta_k + \beta_c + \beta_{k,c} + \log Z\right]$$

748
$$y_{k,c} \sim \text{NB}(\theta_{k,c}, \lambda)$$

749 where X is a binary column vector indicating positive cancer diagnosis, $\beta_1$ is the coefficient for

750 cancer diagnosis, and $\log Z$ is an offset accounting for sample library size (sample sum). The

751 remaining parameters are analogous to the model described above.

752 A permutation test was performed to compare the similarity in topic-pathway effects between

753 themetagenomics and the metagenomic model to random sampling. In the metagenomic KO

754 table, topic and pathway labels were randomly permuted. The permuted table was then refit

755 with the regression model described. The root mean squared error was calculated between the

756 topic-pathway regression coefficient $\beta_{k,c}$ for themetagenomics and the metagenomic model:

757
$$RMSE = \sqrt{\frac{\sum_{k,c}\left(\beta_{k,c}^{(theme)} - \beta_{k,c}^{(\text{meta})}\right)^2}{n}}$$

758 This process was repeated over 100 permuted replicates to calculate a null distribution of RMSE

759 scores, which was then compared to the true RMSE between the unpermuted metagenomic KO

760 table and themetagenomics. A p-value ($\alpha$=0.05) was calculated as the proportion of RMSE

761 scores from the 100 permuted metagenomic KO tables that were less than the RMSE score for

762 the unpermuted metagenomic KO table.

763

764 **Exploring Thematic Structure in David et al.**

765

39

766     **Data Preparation and OTU Picking.** The David et al. dataset contains fecal and salivary 16S

767     rRNA surveys from two subjects. The samples were obtained at uneven sampled times from 318

768     days. Data from were downloaded from the European Bioinformatics Institute (EBI) European

769     Nucleotide Archive (ENA) (accession number ERP006059). It consisted of 1.7 million 16S rRNA

770     gene (V4 region) sequencing reads, 100 bp in length. The reads were quality filtered using the

771     fastqFilter command in the dada2 package [40]. Closed reference OTU picking was then

772     performed with QIIME version 1.9.1. using SortMeRNA again GreenGenes v13.5 at 97%

773     sequence identity [24].

774     **Data Preprocessing and STM Fitting.** From the OTU table, we removed any samples with

775     fewer than 1000 total reads, were not of fecal origin, were not from donor B, and did not include

776     sample data for day, donor, and body site. OTUs lacking a known phylum classification or

777     present in fewer than 1% of the remaining samples were removed. The remaining OTUs were

778     normalized in terms of 16S rRNA gene copy number per the table provided by PICRUSt [14].

779     The final OTU table consisted of 1562 OTUs across 189 samples.

780     We fit 7 STMs that varied in terms of topic number $K \in \{15, 25, 50, 75, 105, 155, 250\}$. To infer the

781     relationship between sample data and the samples-over-topics distribution $\theta$, we used two

782     sample covariates: two continuous, integer valued sequences representing the day number in

783     the sequence and the DOW. Given our assumption that fluctuations in microbiota likely varied

784     nonlinearly with respect to day, we used a smoothing spline with 10 degrees of freedom on day

785     and a second-degree polynomial on DOW.

786     **Event detection.** To detect events in subject B, we repeated the approach described for

787     simulation 2 (S2 appendix).

788     **Hierarchical clustering.** We performed HC for comparison. The David et al. data were

789     normalized using the sample geometric mean to correct for library size imbalance. Each feature

790     was then centered and scaled as described for simulation 2. Clustering was performed as

791     detailed for simulation 2. The resulting tree was cut to produce 6 clusters. The choice of 6

792     clusters was based on the three profiles identified by David et al. (days 1-150, 151-159, and 160-

793     318). We included three additional clusters to account for the background taxonomic variation

794     lacking one of the three profiles of interest. Because we are basing our parameter choice on what

795     can be considered the truth, this can be considered a best-case-scenario.

796

797     **Supporting Information**

798

799     S1 supporting figures. Contains supporting figures S1-S4 and tables S1-S6.

800     S2 appendix. Contains additional information regarding the following: (1) simulation 1 which

801     explores different normalization approaches, (2) time series analysis methods for David et al.

802     data including simulation 2; and (3) additional results for Crohn's disease data as well as

803     expansion of results detailed above and comparisons to other approaches such as SPIEC-EASI

804

805     **References**

806

1.    Kurtz ZD, Mueller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA. Sparse and

       Compositionally Robust Inference of Microbial Ecological Networks. PLoS Comput Biol.

       2015;11: 1–25. doi:10.1371/journal.pcbi.1004226

2.    Shafiei M, Dunn KA, Boon E, MacDonald SM, Walsh DA, Gu H, et al. BioMiCo: a

       supervised Bayesian model for inference of microbial community structure. Microbiome.

       2015;3: 8. doi:10.1186/s40168-015-0073-x

3.    Callahan BJ, Sankaran K, Fukuyama JA, McMurdie PJ, Holmes SP. Bioconductor

       workflow for microbiome data analysis: from raw reads to community analyses.

       F1000Research. 2016;5: 1492. doi:10.12688/f1000research.8986.1

4.    Knights D, Costello E, Knight R. Supervised classification of human microbiota. FEMS

       Microbiol Rev. 2011;35: 343–59. doi:10.1111/j.1574-6976.2010.00251.x

5.    Li H. Microbiome, Metagenomics, and High-Dimensional Compositional Data Analysis.

       Annu Rev Stat Its Appl. 2015;2: 73–94. doi:10.1146/annurev-statistics-010814-020351

6.    Gilbert JA, Quinn RA, Debelius J, Xu ZZ, Morton J, Garg N, et al. Microbiome-wide

       association studies link dynamic microbial consortia to disease. Nature. 2016;535: 94–103.

       doi:10.1038/nature18850

7.    McMurdie PJ, Holmes S. Waste Not, Want Not: Why Rarefying Microbiome Data Is

       Inadmissible. PLoS Comput Biol. 2014;10. doi:10.1371/journal.pcbi.1003531

8.    Love MI, Anders S, Huber W. Differential analysis of count data - the DESeq2 package

826       [Internet]. Genome Biology. 2014. doi:110.1186/s13059-014-0550-8

827   9.    Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models

828       via Coordinate Descent. J Stat Softw. 2010;33: 1–22. doi:10.1359/JBMR.0301229

829   10.   Zou H, Hastie T. Regularization and variable selection via the elastic net. J R Stat Soc Ser

830       B Stat Methodol. 2005;67: 301–320. doi:10.1111/j.1467-9868.2005.00503.x

831   11.   Jiang X, Dushoff J, Chen X, Hu X. Identifying enterotype in human microbiome by

832       decomposing probabilistic topics into components. 2012 IEEE Int Conf Bioinforma

833       Biomed. Ieee; 2012; 1–4. doi:10.1109/BIBM.2012.6392720

834   12.   Ning J, Beiko RG. Phylogenetic approaches to microbial community classification.

835       Microbiome. Microbiome; 2015;3: 47. doi:10.1186/s40168-015-0114-5

836   13.   Ren B, Bacallado S, Favaro S, Holmes S, Trippa L. Bayesian Nonparametric Ordination

837       for the Analysis of Microbial Communities. arXiv Prepr arXiv160105156. 2016; 1–25.

838       Available: http://arxiv.org/abs/1601.05156

839   14.   Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes J a, et al.

840       Predictive functional profiling of microbial communities using 16S rRNA marker gene

841       sequences. Nat Biotechnol. Nature Publishing Group; 2013;31: 814–21.

842       doi:10.1038/nbt.2676

843   15.   Aßhauer KP, Wemheuer B, Daniel R, Meinicke P. Tax4Fun: predicting functional profiles

844       from metagenomic 16S rRNA data. Bioinformatics. 2015;31: 2882–2884.

845       doi:10.1093/bioinformatics/btv287

846  16.  Iwai S, Weinmaier T, Schmidt BL, Albertson DG, Poloso NJ, Dabbagh K, et al. Piphillin:

847       Improved prediction of metagenomic content by direct inference from human

848       microbiomes. PLoS One. 2016;11: 1–18. doi:10.1371/journal.pone.0166104

849  17.  Edgar RC. SINAPS: Prediction of microbial traits from marker gene sequences. bioRxiv.

850       2017; doi:10.1101/124156

851  18.  Knights D, Kuczynski J, Charlson ES, Zaneveld J, Mozer MC, Collman RG, et al. Bayesian

852       community-wide culture-independent microbial source tracking. Nat Methods. 2013;8:

853       761–763. doi:10.1038/nmeth.1650.Bayesian

854  19.  Blei DM, Lafferty JD. A correlated topic model of Science. Ann Appl Stat. 2007;1: 17–35.

855       doi:10.1214/07-AOAS136

856  20.  Roberts ME, Stewart BM, Tingley D, Lucas C, Leder-Luis J, Gadarian SK, et al. Structural

857       topic models for open-ended survey responses. Am J Pol Sci. 2014;58: 1064–1082.

858       doi:10.1111/ajps.12103

859  21.  Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The

860       Treatment-Naive Microbiome in New-Onset Crohn's Disease. Cell Host Microbe. 2014;15:

861       382–392. doi:10.1016/j.chom.2014.02.005

862  22.  Schmidt BL, Kuczynski J, Bhattacharya A, Huey B, Corby PM, Queiroz EL, et al. Changes

863       in abundance of oral microbiota associated with oral cancer. PLoS One. 2014;9: e98741.

864       doi:10.1371/journal.pone.0098741

865  23.  David LA, Materna AC, Friedman J, Baptista MIC, Blackburn MC, Perrotta A, et al. Host

866          lifestyle affects human microbiota on daily timescales. Genome Biol. 2016;17: 117.

867          doi:10.1186/s13059-016-0988-y

868   24.   Caporaso J, Kuczynski J, Stombaugh J, Bittinger K, Bushman. QIIME allows analysis of

869          high-throughput community sequencing data. Nat Methods. 2012;7: 335–336.

870          doi:doi:10.1038/nmeth.f.303

871   25.   Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M. KEGG: Kyoto encyclopedia

872          of genes and genomes. Nucleic Acids Research. 1999. pp. 29–34. doi:10.1093/nar/27.1.29

873   26.   Legendre P, Legendre L. Numerical Ecology - Second English Edition. Developments in

874          Environmental Modelling. 1998. doi:10.1017/CBO9781107415324.004

875   27.   Hardoon DDR, Shawe-Taylor J. Sparse canonical correlation analysis. Mach Learn.

876          2011;10: 1–15. doi:10.1007/s10994-010-5222-7

877   28.   De Valpine P, Harmon-Threatt AN. General models for resource use or other

878          compositional count data using the Dirichlet-multinomial distribution. Ecology. 2013;94:

879          2678–2687. doi:10.1890/12-0416.1

880   29.   Holmes I, Harris K, Quince C. Dirichlet multinomial mixtures: Generative models for

881          microbial metagenomics. PLoS One. 2012;7. doi:10.1371/journal.pone.0030126

882   30.   Brien JDO, Record N. The power and pitfalls of Dirichlet-multinomial mixture models for

883          ecological count data. bioRxiv. 2016; 1–22. doi:10.1101/045468

884   31.   Mimno D, McCallum A. Topic models conditioned on arbitrary features with dirichlet-

885          multinomial regression. arXiv Prepr arXiv12063278. 2012; doi:10.1.1.140.6925

886    32.    Kembel SW, Wu M, Eisen JA, Green JL. Incorporating 16S Gene Copy Number

887           Information Improves Estimates of Microbial Diversity and Abundance. PLoS Comput

888           Biol. 2012;8: 16–18. doi:10.1371/journal.pcbi.1002743

889    33.    Woloszynek S, Zhao Z, Simpson G, O'Connor MP, Mell JC, Rosen GL. Evaluating a topic

890           model approach for parsing microbiome data structure. bioRxiv. 2017; 1–36.

891           doi:10.1101/176412

892    34.    Roberts, Margaret E., Stewart BM, Tingley D. stm: R Package for Structural Topic Models

893           [Internet]. 2017. Available: http://www.structuraltopicmodel.com.

894    35.    Blei DM, Ng AY, Jordan MI. Latent Dirichlet Allocation. 2003;3: 993–1022.

895    36.    Blei DM, McAuliffe JD, Blei DM. Supervised Topic Models. Adv Neural Inf Process Syst

896           20. 2008;21: 1–22. doi:10.1002/asmb.540

897    37.    Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. KEGG for integration and

898           interpretation of large-scale molecular data sets. Nucleic Acids Res. 2012;40.

899           doi:10.1093/nar/gkr988

900    38.    Stan Development Team. rstanarm: Bayesian applied regression modeling via Stan

901           [Internet]. 2016. Available: http://mc-stan.org/

902    39.    Gelman A, Rubin DB. Inference from Iterative Simulation Using Multiple Sequences. Stat

903           Sci. 1992;7: 457–511. doi:10.1214/ss/1177011136

904    40.    Callahan BJ, Mcmurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2 : High

905           resolution sample inference from amplicon data. bioRxiv. 2015;13: 0–14.
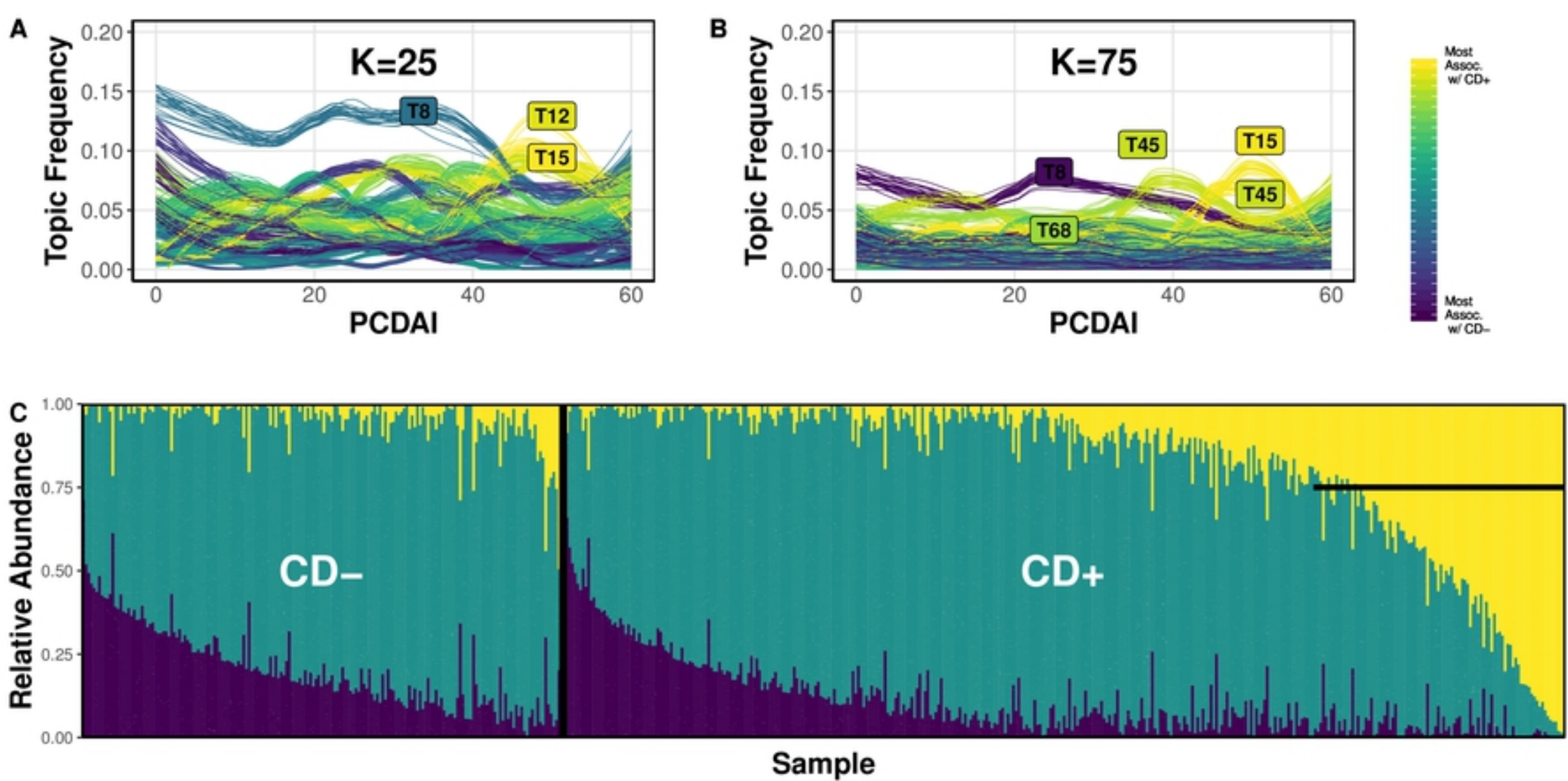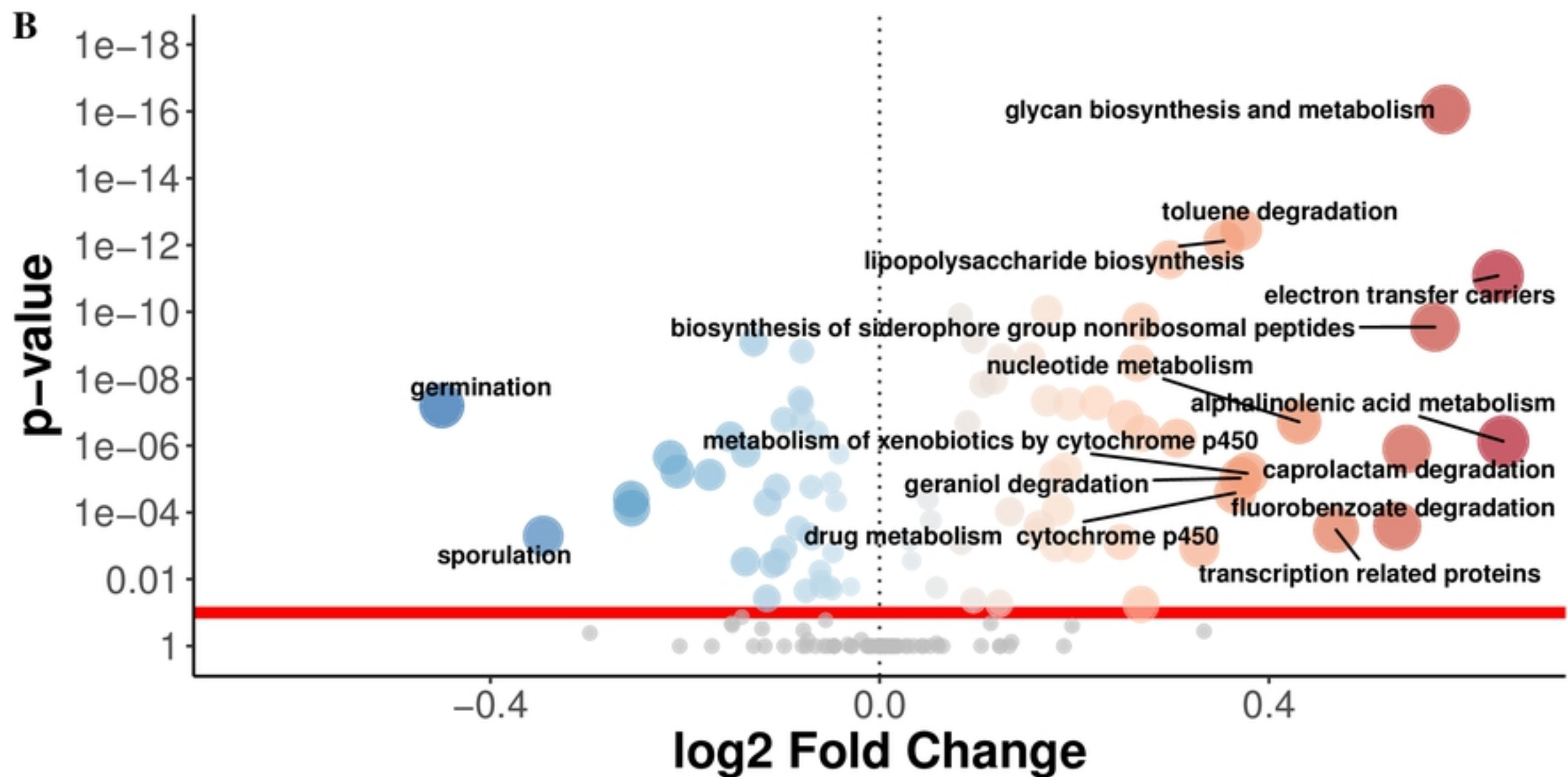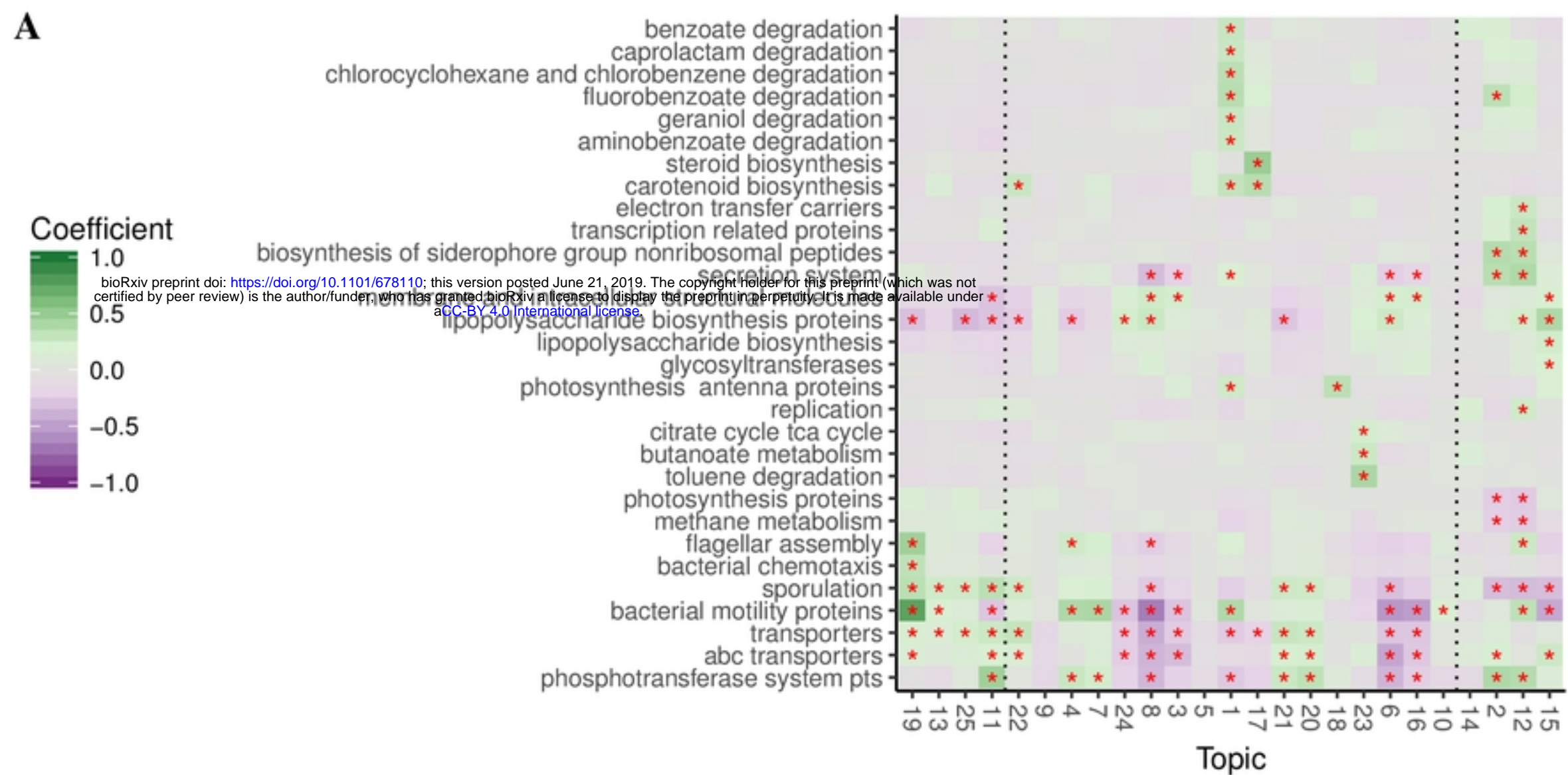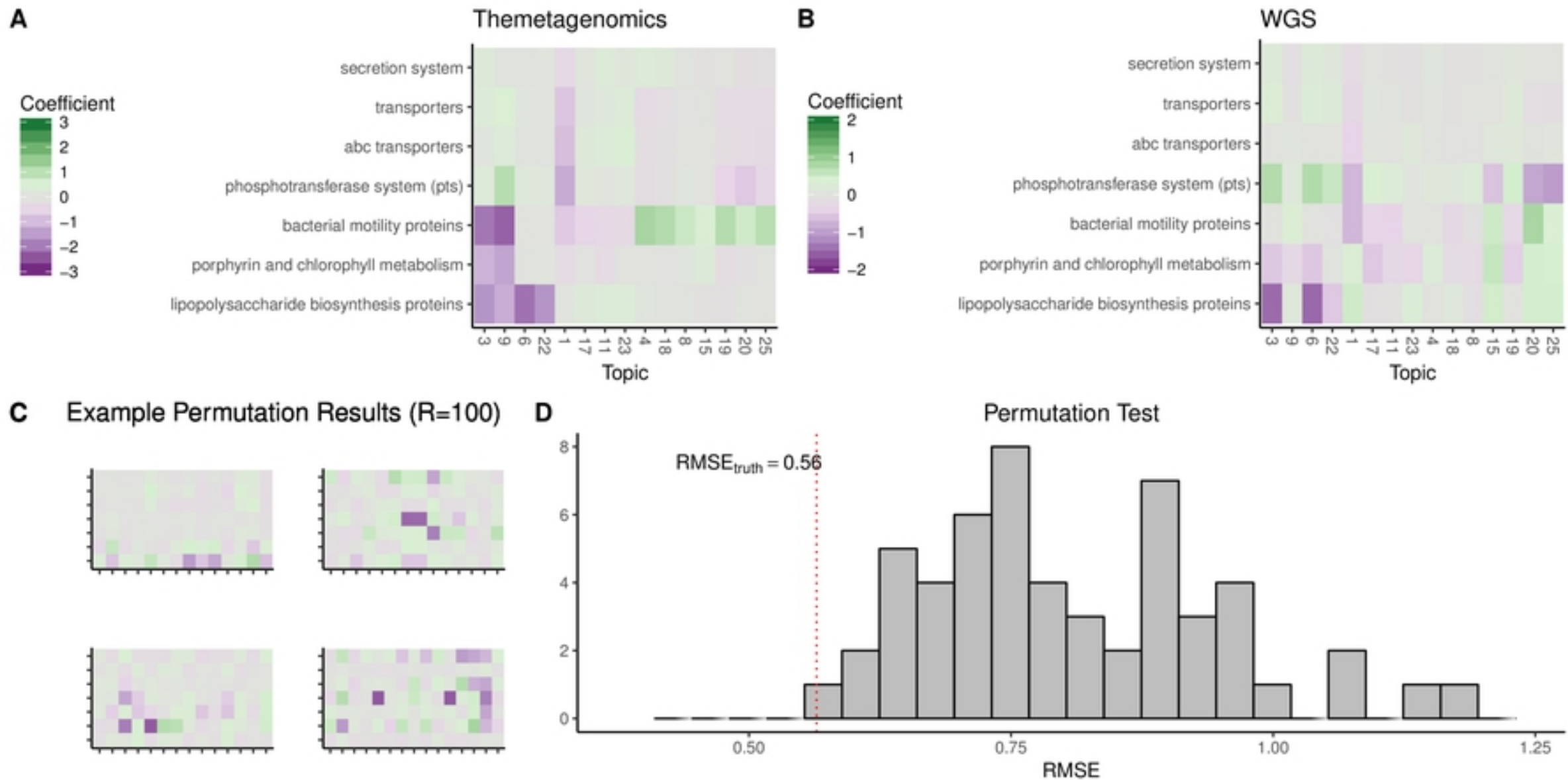
46

906  doi:10.1101/024034

907

Fig 2

Fig 3

**A** Themetagenomics

**B** WGS

**C** Example Permutation Results (R=100)

**D** Permutation Test

RMSE$_{truth}$ = 0.56

Fig 4

Fig 5

Fig 1