

1 **TITLE**

2 Exploring thematic structure in 16S rRNA marker gene surveys

3

4 **AUTHORS**

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

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

7 Gideon Simpson (simpson@math.drexel.edu) [3]

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

9

10 **AFFILIATIONS**

11 [1] Department of Electrical and Computer Engineering, Drexel University, Philadelphia, PA,
12 United States of America.

13 [2] Department of Microbiology and Immunology, Drexel University College of Medicine,
14 Philadelphia, PA, United States of America.

15 [3] Department of Mathematics, Drexel University, Philadelphia, PA, United States of America

16

17 **KEYWORDS**

18 American Gut, Bayesian, Crohn's Disease, Diet, Inflammatory Bowel Disease, KEGG,
19 Metagenomics, Microbiome, PICRUSt, Topic Model

20 **ABSTRACT**

21

22 **Background:** Analysis of microbiome data involves identifying co-occurring taxa associated
23 with a specific set of sample attributes (e.g., disease presence) but is often hindered by the data
24 being compositional, high dimensional, and sparse. Also, the configuration of co-occurring taxa
25 may represent overlapping subcommunities that contribute to host status. Preserving the
26 configuration of co-occurring microbes is superior to detecting indicator species since this
27 approach is more likely to represent underlying microbiome mechanisms and thus facilitate
28 more biologically meaningful interpretations. Moreover, analysis which simultaneously utilizes
29 taxonomic and functional abundances typically requires independent characterization of
30 taxonomic and functional profiles before linking them to sample information. However, this
31 limits investigators from identifying which specific functional components associate with which
32 subsets of co-occurring taxa.

33 **Results:** We provide a pipeline to explore co-occurring taxa using topics generated via a topic
34 model approach and then link these topics to specific sample classes. Also, rather than inferring
35 predicted functional content independently from taxonomic information, we instead focus on
36 within-topic functional content, which we parse via estimating pathway-topic interactions
37 through a multilevel fully Bayesian regression model. We apply our methods to two large 16S
38 amplicon sequencing datasets: an inflammatory bowel disease (IBD) dataset from Gevers et al.
39 and data from the American Gut (AG) project. When applied to the Gevers et al. IBD study, we
40 determine that a topic highly associated with Crohn's disease (CD) diagnosis is (1) dominated
41 by a cluster of bacteria known to be linked with CD and (2) uniquely enriched for a subset of
42 lipopolysaccharide (LPS) synthesis genes. In the AG data, our approach found that individuals
43 with plant-based diets were enriched with Lachnospiraceae, *Roseburia*, *Blautia*, and
44 *Ruminococcaceae*, as well as fluorobenzoate degradation pathways, whereas pathways involved
45 in LPS biosynthesis were depleted.

46 **Conclusions:** We therefore introduce an approach for uncovering latent thematic structure in
47 the context of host state for 16S rRNA surveys. Using our topic-model approach, investigators
48 can (1) capture sets of co-occurring taxa, (2) uncover their functional potential, and (3) identify
49 gene sets that may help guide future inquiry. These methods have been implemented in a freely
50 available R package <https://github.com/EESI/themetagenomics>.

51 **LIST OF ABBREVIATIONS**

52

53 AG, American gut

54 CD, Crohn's disease

55 CV, cross validation

56 IBD, inflammatory bowel disease

57 LDA, latent Dirichlet allocation

58 LFC, log-fold change

59 LPS, lipopolysaccharide

60 OTU, operational taxonomic unit

61 PPD, posterior predictive distribution

62 RF, random forest

63 STM, structural topic model

64 BACKGROUND

65

66 With the decreasing cost of high-throughput sequencing, large datasets are becoming
67 increasingly available, particularly microbiome datasets rich in sample data. These data consist
68 of categorical and numeric information associated with each sample, which in turn are linked to
69 a set of taxonomic abundances that are derived from clustering sequencing reads. Such clusters
70 are based on taxonomic marker genes – typically a portion of the 16S rRNA gene that meet a
71 fixed degree of sequence similarity, termed Operational Taxonomic Units (OTUs). Analysis of
72 these data often involves identifying co-occurring groups of taxa associated with specific
73 sample features via unsupervised exploratory methods such as principal component analysis,
74 correspondence analysis, multidimensional scaling, and hierarchical clustering, in addition to
75 statistical inference strategies aimed at identifying differentially abundant taxa and differences
76 in alpha and beta diversity. Nevertheless, model building is hindered by the complexity
77 inherent to these data, which have a disproportionate number of samples relative to features
78 (Knights et al., 2011), a substantial degree of sparsity, and are typically strictly positive and
79 constrained to sum to 1, i.e., compositional (Gilbert et al., 2016; Li, 2015).

80 From an ecological perspective, the configuration of these co-occurring microbiota may
81 represent related, overlapping sets of subcommunities consisting of taxa that correlate with, for
82 example, host phenotype. Identifying subcommunities that contribute to host status as opposed
83 to single indicator species facilitates a more biologically meaningful interpretation by
84 preserving the natural configuration of co-occurring bacteria when making inferences with
85 respect to host phenotype (Shafiei et al., 2015). Recent work has attempted to explore such
86 relationships (Jiang et al., 2012; Ning & Beiko, 2015; Ren et al., 2016; Shafiei et al., 2015).

87 Still, suitable approaches to uncover these relationships in the context of functional information
88 is deficient – that is, few methods successfully integrate subcommunity-host phenotype with
89 functional profiles specific to these subcommunities. In both metagenomic and 16S rRNA
90 surveys, analyses utilizing taxonomic and functional abundance information typically involve
91 independently characterizing the taxonomic and functional profiles of the samples and
92 subsequently associating these profiles with host information. In the former, this is done
93 directly by analyzing taxonomic and functional sequence information, whereas the latter
94 requires predicting functional profiles from 16S rRNA survey information using methods such
95 as PICRUSt, Tax4fun, and Piphillin (Aßhauer et al., 2015; Iwai et al., 2016; Langille et al., 2013).
96 In either approach, despite obtaining information regarding which taxa co-occur and whether
97 specific taxa or functional categories are associated with sample data, the investigator remains
98 limited from identifying which specific functional components associate with which subsets of
99 co-occurring taxa.

100 In the context of 16S rRNA gene sequencing data, our objective is therefore two-fold: to
101 implement a modeling framework that can (1) capture sets of co-occurring taxa associated with

102 specific sample data and (2) uncover the functional potential that further characterizes the
103 configuration of these subcommunities.

104 For our first objective, we will employ a topic model approach. Topic models have had
105 considerable use in natural language processing, but have also been explored as a method for
106 exploring genomic count data. Knights et al. (Knights et al., 2013) utilized latent Dirichlet
107 allocation (LDA) to infer the relative contributions of an unknown number of source
108 environments to a set of indoor samples. Shafiei et al. (Shafiei et al., 2015), alternatively, took a
109 supervised approach where they first trained their model on sets of co-occurring OTUs to learn
110 how they correlate with sample classes of interest. They were then able to predict the class of
111 new samples given the trained model.

112 Our approach utilizes a structural topic model (STM) (Roberts et al., 2014), which generalizes
113 previously described topic models such as LDA, the correlated topic model (Blei & Lafferty,
114 2007), and Dirichlet-Multinomial regression topic model (Mimno & McCallum, 2012). Like the
115 Dirichlet-Multinomial regression topic model, the STM permits the influence of sample data on
116 the distribution of samples-over-topics; LDA, on the other hand, can only incorporate sample
117 information if done so in a two-stage process – first performing topic extraction, and then
118 identifying linear relationships between the topic assignments and sample information (Blei et
119 al., 2003; Roberts et al., 2014). A two-stage approach limits the breadth of sample information
120 one can use, typically forcing the user to use only a single vector of covariate information, and
121 moreover prevents propagating uncertainty throughout the model. Similar to the correlated
122 topic model, the STM's logistic normal distribution defines the prevalence of topics across
123 samples and permits correlation between topics.

124 With the STM, we will uncover a thematic representation of 16S rRNA survey abundance data
125 and jointly measure its relationship with sample information (figure 1.1). Two latent
126 distributions will be estimated: a samples-over-topics and a topics-over-OTU distribution,
127 which represent the probability of a topic occurring in a sample and the probability of an OTU
128 occurring in a topic, respectively (figure 1.2). By utilizing sample information (figure 1.3), we
129 will then be able to determine whether particular sample covariates increase or decrease the
130 probability of a given topic occurring in a set of samples (figure 1.4).

131 Our second objective will exploit the estimated topics-over-OTUs distribution. These posterior
132 probabilities dictate the taxonomic composition of the topics and therefore should capture
133 meaningful co-occurrences. Moreover, these probabilities resemble relative abundances of
134 samples across taxa. We can therefore infer the functional potential of these topics using tools
135 such as PICRUST, allowing us to predict topic-specific gene composition, using a database of
136 reference genomes (figure 1.5). Then, by identifying topics of interest based on their
137 relationship to sample covariates, we can subsequently link this predicted within-topic
138 functional profile to both within-topic taxonomic abundances, as well as the specific samples
139 that have high probability of containing these topics.

140 It should be noted how this approach differs from the naïve approach where taxonomic and
141 functional profiles are independently estimated and then jointly interpreted. A naïve approach
142 will successfully identify taxonomic abundances that associate with covariate information, and
143 the same for (predicted) functional abundances, but the result lacks the ability to infer which
144 sets of functions are directly linked to specific sets of taxa. The ability to uncover such
145 information provides context as to why specific co-occurrences are present.

146 We apply our methods on two large 16S rRNA amplicon sequencing datasets: an inflammatory
147 bowel disease (IBD) dataset from Gevers et al. (Gevers et al., 2014) and data from the American
148 Gut (AG) project. After confirming the generalizability of extracted topics, we identified distinct
149 taxonomic subcommunities that, in the case of the Gevers dataset, were consistent with
150 published results. These subcommunities were in turn composed of distinct predicted
151 functional profiles, and moreover, our approach provided gene-sets specific to topics of interest
152 that may warrant further exploration. In a companion paper, we performed simulations to
153 further validate a topic model approach for 16S survey data and to determine a suitable
154 normalization strategy (Woloszynek et al., 2017). Our simulations suggested that predefined
155 taxonomic subcommunities concentrate with high probability to extracted topics and that no
156 library size normalization is required to maximize power or ability to infer taxonomic structure,
157 thus making a topic model approach a more direct, suitable procedure for inferring the
158 subcommunity configuration. Also, in the context of topic models, while DESeq2 normalization
159 outperforms rarefying, it results in decreased power compared to simply using raw,
160 unnormalized abundances.

161 These methods have been implemented in a freely available R package themetagenomics:
162 <https://github.com/EESI/themetagenomics>.

163 METHODS

164

165 Review of the Structural Topic Model

166 The STM is a Bayesian generative model such that, given a set of M samples, each consisting N
167 OTUs, belonging to a vocabulary of V unique OTU terms, K latent topics (chosen *a priori*) are
168 generated from the data. These topics consist of overlapping sets of co-occurring OTUs,
169 potentially sharing some biological context. The samples-over-topics distribution is given a
170 logistic Normal (LN) prior, which allows for estimation of topic-topic correlations, giving a
171 means to infer co-occurring topics across samples. The topics-over-OTUs prior, on the other
172 hand, estimates the deviation of OTU frequencies from a background distribution that
173 encompasses all samples in the dataset (Eisenstein et al., 2011). Word and topic assignments are
174 both generated via V - and K -multinomial distributions, respectively.

175 The STM is estimated by a semi-collapsed variational expectation maximization procedure.
176 Convergence is reached when a relative change in the variational objective (i.e., the estimated
177 lower bound) falls below a predetermined tolerance.

178

179 Datasets and Preprocessing

180 16S rRNA sequencing data from two human microbiome studies were downloaded from their
181 corresponding repositories. The Gevers et al. dataset (“Gevers”) (PRJNA237362, 03/30/2016) is a
182 multicohort, IBD dataset that includes control, Crohn's disease (CD), and ulcerative colitis
183 samples taken from multiple locations throughout the gastrointestinal tract (Gevers et al., 2014).
184 The AG project (“AG”) (ERP012803, 02/21/2017), on the other hand, is a crowd sourced dataset
185 that includes user-submitted microbiome samples from a variety of body sites and associated
186 subject information provided through questionnaires (<http://americangut.org/>).

187 **Human gut microbiota from an inflammatory bowel disease cohort (Gevers).** Paired-end
188 reads were joined and quality filtered (maximum unacceptable Phred quality score = 32;
189 maximum number of consecutive low quality base calls before read truncation = 3; minimum
190 number of consecutive high quality base calls included per read as a fraction of input read
191 length = 0.75) using QIIME version 1.9.1. Closed-reference OTU picking was performed using
192 SortMeRNA against GreenGenes v13.5 at 97% sequence identity. This was followed by copy
193 number normalization via PICRUST version 1.0.0 (Kembel et al., 2012).

194 We selected only terminal ileum samples. Those with fewer than 1000 total reads were omitted.
195 We subsequently removed OTUs with fewer than 10 total reads across samples and OTUs that
196 lacked a known classification at the Phylum level.

197 **Human gut microbiota from samples differing in terms of diet (AG).** Quality trimming and
198 filtering were performed in the following manner on single-end reads using the fastqFilter

199 command found in the dada2 R package. The first 10 bases were trimmed from each read.
200 Reads were then trimmed to position 135 based on visualizing the quality score of sampled
201 reads as a function of base position. Further truncation occurred at positions with quality scores
202 less than or equal to 2. Any truncated read with total expected errors greater than 2 were
203 removed. A portion of AG samples were affected by bacterial blooming during shipment. These
204 sequences were removed using the protocol provided in the AG documentation (02-
205 filter_sequences_for_blooms.md).

206 OTU picking and copy number normalization were implemented as above. Samples with fewer
207 than 1000 reads, and OTUs with fewer than 10 total reads across samples and lacking any
208 known classification at the Phylum level were discarded. We filtered samples falling into the
209 “baby” age category (thus the minimum age was 3) and retained only fecal samples. Within the
210 diet category, unknown, vegetarian-with-shellfish, and omnivore-without-red-meat diets types
211 were removed. We then merged vegan and vegetarian-without-shellfish into one class,
212 resulting in a binary set of labels: “O” for omnivores and “V” for vegans and vegetarians.

213

214 **Structural Topic Model Fitting** (figure 1.1)

215 Each resulting OTU table consists of sets of raw counts normalized by 16S rRNA copy number.
216 No other normalization was conducted based on the simulation results in Woloszynek et al.
217 (2017). A series of topic models with different parameterizations in terms of topic number ($K \in$
218 15, 25, 50, 75, 100, 150, 250) and sample covariates (e.g., indicators for presence of disease, diet
219 type, etc.) were fit to the OTU tables.

220 We evaluated each model fit for presence of overdispersed residuals. We also conducted
221 permutation tests where the covariate of interest is randomly assigned to a sample, prior to
222 STM fitting. To compare parameterizations between models, we evaluated predictive
223 performance using held-out likelihood estimation (Blei et al., 2003).

224

225 **Assessment of topic generalizability**

226 We performed classification to assess the generalizability of the extracted topics. No sample
227 information was used as covariates. OTU tables were first split into 80/20 training-testing
228 datasets. A topic model was trained to estimate the topics-over-OTUs distribution. We then
229 held this distribution fixed; hence, only the testing set’s samples-over-topics distribution was
230 estimated. For both the training and testing sets, simulated posterior samples from the samples-
231 over-topics distribution were averaged. The resulting posterior topic probabilities in the
232 training set were then used as predictors to classify sample labels, similar to using \bar{Z} in
233 supervised LDA (Blei et al., 2008). The generalization error was then assessed by using the

234 optimal parametrization based on cross validation performance (CV) on the test set topic
235 probabilities. Classification was performed using a random forest (RF).

236 For the RF, parameter tuning to determine the number of variables for each split was
237 accomplished through repeated (10x) 10-fold CV, using up- or down-sampling to overcome
238 class unbalance (for Gevers and AG, respectively). We performed a parameter sweep over the
239 number of randomly selected features, while setting the number of trees fixed at 128. The
240 optimal parameterizations were selected based on maximizing ROC area under the curve
241 (AUC).

242

243 **Assessing Concentration of OTUs as a function of topic number**

244 Comparison of Shannon entropy across topics was performed via ANOVA and Tukey HSD
245 post-hoc analysis. To quantify the relationship between taxonomic abundance and continuous
246 predictors (e.g., PCDAI), we performed negative binomial regression (log link), using total
247 sample coverage as an offset. The family-wise error rate was adjusted via Bonferroni correction.
248 Critical values for hypothesis testing were set at 0.05 unless stated otherwise.

249

250 **Comparison of topic taxonomic profile to a network approach**

251 To further validate the clusters of high probability taxa identified in the topics-over-OTUs
252 distribution, we compared our results to those generated from an OTU-OTU association
253 network on the raw (copy number normalized) OTU tables using SPIEC-EASI's neighborhood
254 selection method (Kurtz et al., 2015).

255

256 **Inferring within-topic functional potential** (figure 1.5)

257 We obtained the topics-over-OTUs distribution for each model fit and mapped the within-topic
258 OTU probabilities to integers ("pseudo-counts") using a constant: $10000 \times \beta$. A large constant
259 was used to prevent low probability OTUs from being set to zero, although their contribution to
260 downstream analysis was likely negligible. Gene prediction was then performed on each topic-
261 OTU pseudo-count table using PICRUSt version 1.0.0 (Langille et al., 2013). Recall that copy
262 number normalization was performed prior to topic model fitting.

263

264 **Identifying topics of interest** (figure 1.3, 1.4)

265 Topics of interest were identified by regressing the sample-specific topic probabilities against
266 their set of sample covariates. We calculated 95% uncertainty intervals using an approximation
267 that accounts for uncertainty in estimation of both the coefficients and the topic probabilities.

268 **Identifying predicted functions that distinguish topics** (figure 1.6)

269 To determine which predicted gene functions best distinguish topics, we utilized the following
270 multilevel negative binomial regression model:

$$271 \theta_{k,c} = \exp[\mu + \beta_k + \beta_c + \beta_{k,c}]$$

$$272 y_{k,c} \sim \text{NB}(\theta_{k,c}, \lambda)$$

273 where μ is the intercept, β_k is the per topic weight, β_c is the per level-3 gene category weight, $\beta_{k,c}$
274 is the weight for a given topic-gene category combination, $y_{k,c}$ is the count for a given topic-gene
275 category combination, and λ is the dispersion parameter. All weights were given normal priors.
276 Convergence was assessed across 4 chains using diagnostic plots to assess mixing and by
277 evaluating the Gelman-Rubin convergence diagnostic (Gelman & Rubin, 1992). To reduce
278 model size, we used genes belonging to only 15 (arbitrary number) level-2 KEGG pathway
279 categories (table S1). For large topic models, we fit only the top 25 topics, ranked in terms of the
280 regression weights that measure the degree of association between sample-over-topic
281 probabilities and our covariate of interest.

282

283 **Comparison of within-topic pathway profile to OTU-table approach**

284 We compared the profile of predicted functions obtained from the hierarchical negative
285 binomial model to a differential abundance approach. We performed (KEGG) functional
286 prediction via PICRUSt on raw OTU abundances that were copy number normalized. The
287 resulting functional abundances were collapsed into level-3 KEGG pathways. Note that we
288 again restricted our genes to the 15 level-2 KEGG pathways used previously to remain
289 consistent. The resulting level-3 pathway abundances underwent DESeq2 differential
290 abundance analysis followed by Bonferroni correction (McMurdie & Holmes, 2014). Adjusted p-
291 values below 0.1 were deemed significant.

292

293 **Packages utilized**

294 All analysis was done in R version 3.2.3. Topic models, RFs, and NB regression models were fit
295 using *stm* (Roberts, Margaret E., Stewart & Tingley, 2017), *caret* (Kuhn, 2008), and *rstanarm*
296 (Stan Development Team, 2016), respectively. AG filtering was performed using *dada2*
297 (Callahan et al., 2015). SPIEC-EASI was fit using the SPIEC-EASI package (Kurtz et al., 2016).
298 DESeq2 differential abundance analysis was conducted with *phyloseq* (McMurdie & Holmes,
299 2013).

300

301 RESULTS

302

303 We will explore the use of a topic model approach on datasets of gut and fecal microbial
304 community profiles, beginning with the IBD data from Gevers, followed by the dietary data
305 from AG. For each dataset, we show that the topics extracted from the STM generalize well to
306 test set data not initially seen by the model, suggesting that co-occurrence profiles identified by
307 the STM are robust to overfitting. Then, we apply our complete pipeline, where we successfully
308 link within-topic predicted functional profiles to taxonomic subcommunity configurations and
309 host features.

310

311 Thematic Structure of IBD-Associated Microbiota (Gevers)

312 **Dimensionality reduction using topics facilitates classification of CD diagnosis and**
313 **generalizes well to test data.** We aimed to assess whether (1) topics fit in the absence of sample
314 covariates are associated with positive CD diagnosis (CD+), and (2) they generalize to new data
315 – that is, whether they captured meaningful information inherent to the data while ignoring
316 characteristics associated exclusively with the fitted data.

317 The 80/20 training/testing splits for terminal ileum samples from Gevers are shown in table S2.
318 We hypothesized that there would be a drop in performance using OTU relative abundances as
319 features compared to topics, since the former has much higher dimensionality and is sparser.
320 These are both relaxed when using topics, since the size of the feature space is decreased
321 through dimensionality reduction. There was little difference between the two approaches
322 during training CV with at least 25 topics (figure S1, table S3). During testing, however, topics
323 outperformed OTU relative abundances, particularly in terms of F1 score, with scores of 0.808
324 and 0.857 for OTUs and topics (K=25 and K=100), respectively (table S4).

325 As one example, the largest discrepancy in classification performance between OTUs and topics
326 was in terms of their negative predictive value, with the OTU model being correct only half the
327 time (0.517) when predicting the negative class (CD-), whereas the worst performing topic
328 model (K=15) performed slightly better (0.526), and topic models seemingly improved as the
329 number of topics increased: 0.655 (K=25), 0.559 (50), 0.577 (75), 0.682 (100), and 0.643 (150) (table
330 S4). Such a high proportion of false negatives with the OTU model was likely due to its reliance
331 on few, relatively rare taxa (figure S2).

332 As another example, OTU 319708 (Clostridiaceae family) was the fourth most important feature
333 for distinguishing classes. It was over twice as common in CD- training samples. Over 10% of
334 correctly classified CD- samples contained this feature. This was also the case for 10% of
335 misclassified CD+ samples, some of which contained this OTU at a greater proportion than
336 other samples in the training set. A similar scenario can be seen for the OTU with the largest

337 importance score, OTU 186723 (Ruminococcaceae family), which associated predominately with
338 disease presence, and hence its absence in CD+ samples resulted in false negatives.

339 **Concentration of high probability OTUs across topics begins to plateau at 75 topics.** After
340 assessing the generalizability of extracted topics, we implemented our full pipeline using
341 sample covariates, specifically a binary indicator for IBD diagnosis. After fitting the topic model
342 to the OTU abundance data, we aimed to uncover how OTUs concentrate within topics as a
343 function of topic number. We performed an ANOVA to compare Shannon entropy for
344 individual topics across OTUs for topic models of varying sizes, followed by post-hoc multiple
345 comparisons testing using Tukey HSD ($\alpha=0.05$) (figure S2). We found a significant difference in
346 the mean Shannon entropy among the models considered. When we tested for differences
347 between pairwise model combinations, we found that the drop in entropy with increasing topic
348 number diminished, such that differences between models with 75, 100, and 150 topics were not
349 significantly different from one another. This suggests that the probability mass of the topics-
350 over-OTUs distribution concentrates on OTU subsets of similar sizes as topic number increases.
351 Analyzing topics in this way may help guide the user in the selection of topic number.

352 **CD diagnosis was associated with unique thematic and hence taxonomic profiles.** The
353 configuration of topics K25 and K75 are shown in figures S3 and S4, exemplifying how our
354 pipeline represents 16S rRNA abundance data. For the topics shown, their posterior estimates
355 did not span 0, a result that was also present when we performed permutation tests to confirm
356 (figure S5). The panels are ordered in terms of mean effect estimate using the samples-over-
357 topic distribution against sample diagnosis. We consider topics with larger mean effect
358 estimates as “high-ranking topics.” Both panels show that CD- training samples had a topic
359 distribution that differed from CD+ samples. Moreover, a given topic’s association with disease
360 presence in most influenced by the disease burden of its samples, particularly for K25, where
361 CD+ samples with high probability for T19, T13, and T26 (the topics most associated with CD-)
362 tend to have minimal disease burden. We henceforth focus on the K25 model.

363 Focusing on these eight key topics, we identified multiple clusters of bacterial species that
364 disproportionately dominated the top topics associated with CD+ (figure 2, top). For example,
365 T2 contained a cluster dominated by *Enterobacteriaceae* taxa, whereas T12’s cluster contained a
366 mixture of *Fusobacteria* and *Enterobacteriaceae*. The T15 cluster contained *Haemophilus* spp.,
367 *Neisseria*, *Fusobacteria*, and *Streptococcus*, all of which were noted as having a positive correlation
368 with CD+ subjects in Gevers et al, as well as *Aggregatibacter*, a genus reportedly associated with
369 colorectal cancer (Tjalsma et al., 2012).

370

371 Given that T15 contains a cluster of bacteria known for their association with bowel
372 inflammation and this topic occurs disproportionately in subjects with greater disease burden,
373 we asked whether the abundance of these OTUs in CD+ subjects correlated with PCDAI, a
374 clinical measure of CD burden. After performing negative binomial regression (figure 3), we

375 identified significant positive trends as a function of PCDAI for *Aggregatibacter* ($p < 0.0001$),
376 *Erwinia* ($p = 0.0004$), *Fusobacterium* ($p = 0.0001$), and *Haemophilus* ($p = 0.0484$).

377 The topics most associated with CD-, on the other hand, were dominated by taxa belonging to
378 *Lachnospiraceae*, *Roseburia*, *Rubinococcus*, *Blautia*, *Bacteroidetes*, and *Coprococcus*, all of which were
379 noted by Gevers et al. as being negativity associated with CD (figure 2, bottom; figure S6). In
380 addition to these taxa, *Akkermania*, *Dialister*, and *Dorea* contributed to these topics, which is
381 consistent with the findings of Lewis et al. who found a reduction of these taxa in CD+ subjects
382 (Lewis et al., 2015).

383 **Within-topic co-occurrence profiles were confirmed via SPIEC-EASI.** We compared the
384 resulting topics to the correlations obtained via a network approach. The SPIEC-EASI network
385 edges for the clusters of high probability OTUs in our most correlated topics are showed in
386 figure S7. For each of these topic clusters, the majority of taxa were connected by a non-zero
387 edge (table S5). Of the 11 taxa in the T15 cluster, 8 had first order connections (direct
388 connections to other taxa within the cluster, $OTU_c - OTU_c$), whereas 9 had second order
389 connections (indirect connections to other taxa within the cluster via an intermediate OTU not
390 present in the cluster, $OTU_c - OTU_{nc} - OTU_c$). Moreover, the two OTUs connected by largest edge
391 weight, *H. parainfluenzae* and *Haemophilus spp.*, had the largest probabilities of the taxa in the
392 topic cluster, 0.320 and 0.245, respectively. Of these 6 topics, none had more than one OTU with
393 zero connections or fewer than 75% of taxa joined by first order connections. Predictably, the
394 taxa that lacked within-cluster connections received low probability from the topic model, with
395 one exception, *Catenibacterium spp.* in T19. Taken together, this reaffirms that the within-topic
396 co-occurrence profiles are consistent with alternative approaches.

397 **Predicted functional potential of notable topics further described their association with CD.**

398 We sought to further explore the co-occurrence profiles of these topics, thereby exploiting the
399 posterior estimates of the topic model in a way unique compared to other approaches. To do so,
400 we predicted the topic-specific functional content using PICRUST and then performed a fully
401 Bayesian multilevel regression analysis on the abundances of each gene function.

402 Like Gevers et al., we identified an increase in membrane transport associated with CD+,
403 particularly topics T2 and T12; however, through our approach, we were able to pinpoint the
404 specific topics these functional categories associated with. This, in turn, allowed us to link these
405 categories to specific taxa. For example, the two aforementioned topics were dominated by
406 Enterobacteriaceae (figure S8). Topic T15, on the other hand, contained the cluster of
407 *Haemophilus spp.*, *Neisseria*, and *Fusobacteria* taxa, and despite being most associated with CD+,
408 had a less substantial shift in membrane transport genes, suggesting that this pathogenic cluster
409 contributed less to the shift of those genes.

410 A considerable degree of cell motility genes was found in T19 relative to all other topics, which
411 is consistent with this topic being dominated by mobile bacteria that belong to Lachnospiraceae,
412 *Roseburia*, and Clostridiales. More specifically, this topic was enriched for genes belonging to

413 the following KEGG categories: bacterial motility proteins, bacterial chemotaxis, and flagellar
414 assembly (figure 4). The aforementioned Enterobacteriaceae-enriched topics were also enriched
415 for siderophore and secretion system related genes. Enrichment of two lipopolysaccharide (LPS)
416 synthesis categories were associated with CD+ topics; however, one of these categories was
417 specific for T15 (table S7).

418
419 **Considerably more pathways were deemed significant via a DESeq2 approach on the OTU**
420 **abundance table, hindering interpretability.** We compared our within-topic functional profiles
421 to the profiles obtained by performing PICRUSt on the copy-number normalized OTU
422 abundance table and then performing a DESeq2 differential abundance analysis. Of the 160
423 level-3 KEGG categories, 87 were found significant ($\alpha < .1$) (figure S9) in the DESeq2 approach.
424 Pathways with the largest log-fold change (LFC) associated with CD+ samples included
425 degradation pathways (caprolactam, LFC=0.542; fluorobenzoate, 0.532; geraniol, 0.371; and
426 toluene degradation, 0.371), alphalinolenic acid metabolism (0.641), and electron transfer
427 carriers (0.635). Interestingly, these degradation pathways also demonstrated strong effects
428 between topics; however, they associated with T1, a topic unrelated to disease status. Electron
429 transfer carriers was identified in both approaches, but the topic model approach isolated T12,
430 placing high probability on bacteria also enriched for functions linked to secretion systems, LPS
431 biosynthesis, and motility. The DESeq2 approach also found fewer categories associated with
432 CD- that had large LFC. For example, only 1 category had an LFC less than -0.04, whereas there
433 were 8 greater than 0.04. The categories with the largest LFCs relative to CD- included
434 germination (LFC=-0.450) and sporulation (-0.346). The topic model identified 10 topics with
435 functional profiles significantly enriched or depleted in sporulation genes, three of which were
436 associated with CD- samples. Moreover, multiple topics demonstrated an inverse relationship
437 between sporulation and LPS genes, such that topics that contained taxa enriched in one were
438 depleted in the other.

439

440 **Thematic Structure in Terms of Diet (AG)**

441

442 Despite consisting of far more samples, the AG dataset, split into O and V diet groups from self-
443 reported dietary information, offered a new challenge for our approach, given that there were
444 far more data and features (taxa), as well as severe imbalance between classes. Of the 4864
445 samples that fit into our diet classes, 4527 and only 337 were O and V samples, respectively.
446 This renders comparisons between group means a worse estimate of treatments effects (Gelman
447 & Hill, 2006).

448 **Accounting for potential confounding.** Before applying our pipeline, we aimed to eliminate
449 any potential sources of confounding. Male and female samples were distributed similarly with

450 respect to diet (table S9; figure S10). There was no significant difference in mean age between
451 diet groups ($t=-0.03$, $df=373.93$, $p=0.98$). Sample body mass index was not normally distributed
452 (Shapiro-Wilk: $W=0.86$, $p<0.001$) and was plagued with many mislabeled heights and weights
453 (figure S11). After attempting to remove samples we deemed unreliable, we found a significant
454 mean difference in body mass index between diet groups via a Mann Whitney U test ($p<0.001$).

455 **Classification using topics is less conservative and, for low dimensional models, less**
456 **generalizable.** Unlike Gevers, models with fewer topics ($K < 75$) generalized poorly compared
457 to using OTUs as features, which may be due to AG having nearly 3-times as many unique
458 OTUs, causing too few topics to dampen any meaningful signal (table S11). Interestingly, all
459 parameterizations outperformed the raw data in terms of sensitivity but not specificity (table
460 S11), suggesting that classification using OTU features is more conservative.

461 **Diet was associated with specific taxonomic and predicted functional profiles.** We will
462 henceforth report the results from a 100 topic model fit with dietary prior information. As
463 before, we identified our topics of interest by regressing the samples-over-topics distribution
464 against diet and further validated these results via permutation tests, resulting in 9 topics, 5 of
465 which were associated with the O group, and 4 with the V group (figure S13).

466 Across the 9 topics, members of the family Lachnospiraceae were well represented, which is not
467 surprising given that it typically accounts for over half of bacteria in healthy human fecal
468 samples (Flint, 2012). Within topics, we identified roughly 11 clusters of interest that contained
469 high probability taxa, one of which belonged to T61, the topic most associated with the V group
470 (figure 5). This cluster was dominated by taxa belonging to Lachnospiraceae (11/23), but T61
471 still placed high probability on *Roseburia*, *Blautia*, and *Ruminococcaceae*. Given T61's associated
472 the V diet, this result is consistent with literature associating *Roseburia* and *Ruminococcaceae* with
473 starch and plant polysaccharide metabolism (Flint et al., 2012) and *Roseburia* and *Blautia* with
474 whole grains (Flint et al., 2015; Martínez et al., 2013). Also, consistent with this topic being
475 dominated by Gram positive bacteria, we identified a significant depletion in predicted LPS
476 biosynthesis genes. (figure 6)

477 T12 contained a small yet diverse cluster of bacteria within *Acinetobacter*, a genus often
478 associated with fermented foods and beverages, having high topic probability (Tamang et al.,
479 2016). Quinn et al. (2016), investigating the effect home-fermented foods had on human
480 microbiota, identified enrichment of predicted fluorobenzoate degradation pathways (Quinn et
481 al., 2016). This same pathway, T12, had the largest shift of any predicted pathways within a
482 given topic (figure 6). To further investigate relationship between fluorobenzoate degradation
483 pathways and diet group, we performed a logistic regression (logit link) on all samples aged at
484 least 21y. Diet type ($V=1$) and the z-scored probability of containing T61 were independent
485 variables with alcohol consumption ($n_{no}=837$, $n_{yes}=3692$) as the binary outcome ($yes=1$). Both T61
486 ($\beta_{T61}=1.10$, $z=3.64$, $p<0.001$) and diet ($\beta_{diet}=-0.89$, $z=-6.78$, $p<0.001$) were significant, suggesting a

487 potential relationship with fermented foods (specifically alcohol), *Acinetobacter*, and
488 fluorobenzoate degradation.

489 Finally, while T76 contained bacteria typically associated with a western lifestyle such as
490 *Clostridiales* (Gorvitovskaia et al., 2016), it also placed the most probability mass on the
491 *Faecalibacterium prausnitzii* (figure S14), as well as predicted butyrate production. This is
492 significant because butyrate has not only critical in the fermentation of plant matter (Gill et al.,
493 2006), but reduction of fecal butyrate has been implicated in obesity and a shift toward a less
494 carbohydrate-rich diet (Duncan et al., 2007). Moreover, the remaining bacteria present in this
495 T76 cluster, *Ruminococcus* and *Roseburia*, have been shown to be elevated after fiber consumption
496 (Flint et al., 2015).

497 The topics associated with the O group, on the other hand, had predicted enrichment for LPS
498 and secretion system pathways. A noteworthy cluster in T77 was surprisingly quite similar to
499 the aforementioned cluster in T61. Lachnospiraceae composed the majority of each cluster:
500 47.8% (11/23) of taxa for T61 compared to 20.6% (13/63) for T61. Moreover, the profiles of
501 predicted functional content were analogous for all pathways except carotenoid biosynthesis
502 and porphyrin and chlorophyll metabolism. A notable distinguishing characteristic is the lack
503 of any *Roseburia* in the T77 cluster compared to T61.

504 T20 also was enriched in predicted carotenoid biosynthesis, but the specific genes differed
505 between the two topic clusters (table S12). T77 contained a disproportionate amount of the gene
506 that codes for the enzyme in the final step of the synthesis of bacterial antioxidant
507 staphyloxanthin (Clauditz et al., 2006) (figure S16). T20 was also abundant in genes belonging
508 associated with secretion system function and LPS biosynthesis, and with respect to T20, a
509 relative shift away from a subset of LPS genes key in one specific branch of the LPS pathway.
510 High probability mass was placed on two taxa (order RF32) belonging to the class
511 alphaproteobacteria, which has been identified in a cluster associated with animal based diets
512 (David et al., 2014).

513 DISCUSSION

514

515 We have introduced our approach for uncovering latent thematic structure in the context of
516 host state for 16S rRNA surveys. We contend that using a topic model to explore taxonomic and
517 predicted functional structure improves interpretability in its natural ability to uncover the
518 relationship between collections of co-occurring taxa (topics) and samples, topics and
519 individual taxa, as well as topics and host covariates. Also, rather than inferring predicted
520 functional content independently from taxonomic information, we shifted our focus to
521 predicting within-topic functional content, which we parse by estimating pathway-topic
522 interactions using a multilevel fully Bayesian regression model. The result not only provides a
523 means to further explore our topics, it also allows us to link functions to specific clusters which
524 can in turn be linked to sample covariates. This has notable implications in that we are
525 drastically reducing the dimensionality of three sources of information, thus achieving a novel
526 means to interpret these data. Moreover, we can identify gene sets of interest from noteworthy
527 topics. For example, when the pipeline was applied to Gevers, we determined that T15 is (1)
528 associated with CD+ samples; (2) dominated by a cluster of bacteria known to be associated
529 with CD; and (3) uniquely enriched for a subset of LPS synthesis genes. Being able to explore
530 this topic's gene profile demonstrates the utility of this topic model approach. Using this
531 information, one could focus on gene subsets associated with topic specific bacterial clusters
532 that are known disease biomarkers, which in turn may facilitate targeted approaches for
533 manipulating the microbiome.

534 We present our approach at a time when novel means to analyze complex microbiome
535 abundance data is called for. Current methods often link the abundance of a single OTU across
536 samples to some particular sample outcome. These methods routinely identify important
537 subsets of taxa, but ignore OTU co-occurrence. Network methods overcome this concern, but
538 instead fail to do so in the context of sample data and hence are incapable of linking sections of
539 the network with sample subsets of interest. Constrained ordination methods, such as canonical
540 correspondence analysis, do in fact couple inter-community distance with sample information,
541 but the user is limited to specific distance metrics (e.g., Chi-squared) and must follow key
542 assumptions (e.g., the distributions of taxa along environmental gradients are unimodal)
543 (Legendre & Legendre, 1998). Moreover, interpretation of biplots becomes increasingly difficult
544 as more covariates are included, and, unlike our approach, linking key subsets of taxa with
545 corresponding subsets of gene functions is not easily achievable.

546 The ability to make meaningful inferences is further compounded by the fact that microbiome
547 data is often inadequately sampled (justifying some type of normalization procedure),
548 compositional (due to normalization), sparse, and overdispersed. Compositional data restricts
549 the appropriateness of many statistical methods due to the sum constraint placed across
550 samples. SPIEC-EASI provides a robust network approach for overcoming compositional

551 artifacts in an attempt to infer community level interactions. We hence compared our within
552 topic taxonomic clusters to the first and second order interactions identified by SPIEC-EASI, to
553 which we found coherence between the two approaches, suggesting a topic model approach for
554 compositional data is in fact appropriate.

555 Others have explored the use of Dirichlet-Multinomial models, which are well equipped at
556 managing overdispersed count data (Brien & Record, 2016; Holmes et al., 2012; De Valpine &
557 Harmon-Threatt, 2013). The fact that Dirichlet-Multinomial conjugacy is exploited for the
558 topics-over-OTUs component of the topics models described above reflects their suitability for
559 abundance data. We selected the recently developed STM for our workflow because of its
560 ability to not only utilize sample data prior information in the flavor of the Dirichlet-
561 Multinomial topic model, but also its ability to capture topic correlation structure and apply
562 partial pooling over samples or regularization across regression weights.

563 Normalization is also a chief concern when analyzing sequencing abundance data (McMurdie &
564 Holmes, 2014); hence, we found it imperative to determine a suitable approach. In the original
565 LDA paper, the generative process assumed a fixed document length N , but N was considered
566 a simplification and could easily be removed because it is independent of all other components
567 of the model. This allows for the possibility of more realistic document size distributions (Blei et
568 al., 2003). Given this fact, coupled with the ability of the Dirichlet-Multinomial distribution in
569 handling overdispersion, and the results of our simulations, we concluded that raw abundance
570 data could be adequately modeled in our approach (Woloszynek et al., 2017). The variance
571 stabilization through DESeq2, while potentially ideal for large sample sizes with adequate
572 signal, seemed to dampen the ability to identify topic-sample associations. Despite performing
573 well at mapping SCs to topics, the rarefied approach suffered from reduced power when
574 identifying topics with large covariate effects.

575 Finally, there are limitations to our approach. First, the workflow from OTU abundance table
576 through pathway-topic inference scales poorly in terms of computation time for large numbers
577 of topic, which may be more necessary as datasets continue to grow in size. Regularization and
578 sparsity inducing priors help limit the number of important topics; hence, exploring only a
579 subset of topics during the final regression step can offer substantial speed improvements at
580 little cost, but utilizing the complete set of topic information would be ideal. Also, we utilize
581 Hamiltonian MC via Stan. Other posterior inference procedures such as variational inference
582 using software packages such as Edward may provide additional speed enhancements (Brevdo
583 et al., 2017). Second, we are capable of separately estimating the uncertainty in our topic model,
584 the hierarchical regression model, and the functional predictions from PICRUSt, but we
585 currently do not propagate the uncertainty throughout the workflow. Doing so would improve
586 downstream interpretation with better estimation of the topic-sample covariates and pathway-
587 topic effects, which in turn would greatly improve one's confidence with utilizing within-topic
588 gene sets. Third, we do not incorporate phylogenetic branch length information, which could
589 lead to more meaningful topics.

590 **ACKNOWLEDGMENTS**

591

592 We would like to thank Zhengqiao Zhao and Mike O'Connor for their feedback on the
593 manuscript and Zhengqiao in particular for his help in developing the simulation. We would
594 also like to thank the Casey Greene Lab for their information on continuous integration and
595 reproducibility, which was invaluable when developing Themetagenomics.

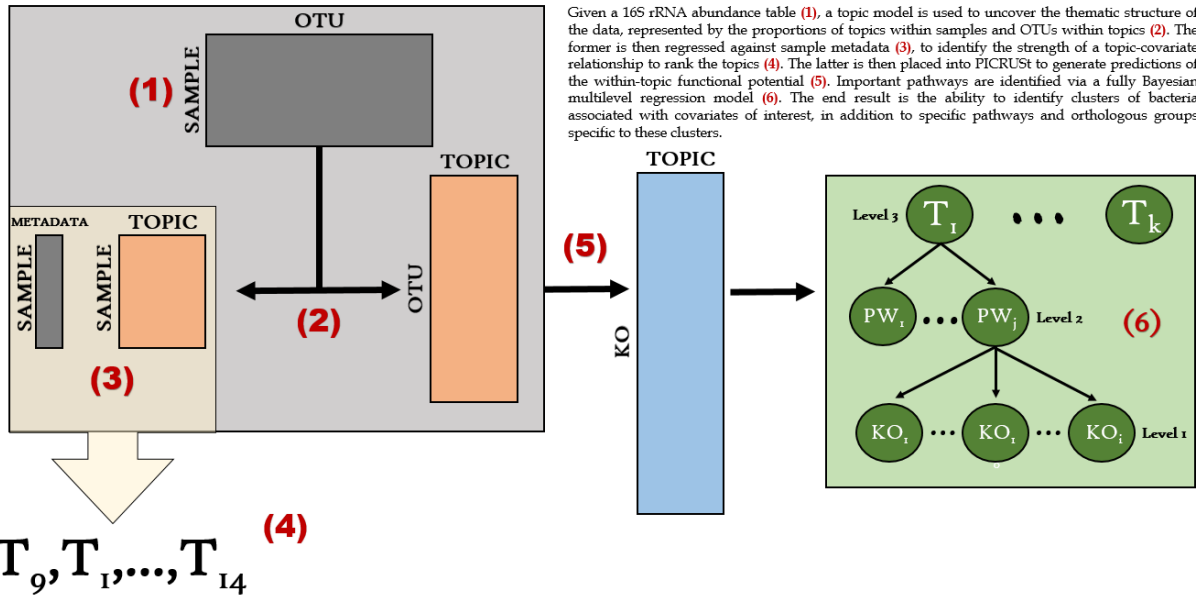
- 596 Aßhauer, K.P., Wemheuer, B., Daniel, R., & Meinicke, P. (2015). Tax4Fun: predicting functional
597 profiles from metagenomic 16S rRNA data: Fig. 1. *Bioinformatics* **31**, 2882–2884.
- 598 Blei, D.M., & Lafferty, J.D. (2007). A correlated topic model of Science. *Ann. Appl. Stat.* **1**, 17–35.
- 599 Blei, D.M., Ng, A.Y., & Jordan, M.I. (2003). Latent Dirichlet Allocation. **3**, 993–1022.
- 600 Blei, D.M., McAuliffe, J.D., & Blei, D.M. (2008). Supervised Topic Models. *Adv. Neural Inf.*
601 *Process. Syst.* **20** **21**, 1–22.
- 602 Brevdo, E., Hoffman, M.D., Murphy, K., Blei, D.M., Tran, D., Saurous, R.A., Hoffman, M.D.,
603 Brevdo, E., Murphy, K., & Blei, D.M. (2017). Deep Probabilistic Programming.
- 604 Brien, J.D.O., & Record, N. (2016). The power and pitfalls of Dirichlet-multinomial mixture
605 models for ecological count data.
- 606 Callahan, B.J., Mcmurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J., & Holmes, S.P. (2015).
607 DADA2 : High resolution sample inference from amplicon data. *bioRxiv* **13**, 0–14.
- 608 Clauditz, A., Resch, A., Wieland, K.P., Peschel, A., & Götz, F. (2006). Staphyloxanthin plays a
609 role in the fitness of *Staphylococcus aureus* and its ability to cope with oxidative stress. *Infect.*
610 *Immun.* **74**, 4950–4953.
- 611 David, L. a, Maurice, C.F., Carmody, R.N., Gootenberg, D.B., Button, J.E., Wolfe, B.E., Ling, A.
612 V, Devlin, a S., Varma, Y., Fischbach, M. a, Biddinger, S.B., Dutton, R.J., & Turnbaugh, P.J.
613 (2014). Diet rapidly and reproducibly alters the human gut microbiome. *Nature* **505**, 559–563.
- 614 Duncan, S.H., Belenguer, a., Holtrop, G., Johnstone, a. M., Flint, H.J., & Lobley, G.E. (2007).
615 Reduced Dietary Intake of Carbohydrates by Obese Subjects Results in Decreased
616 Concentrations of Butyrate and Butyrate-Producing Bacteria in Feces. *Appl. Environ. Microbiol.*
617 **73**, 1073–1078.
- 618 Eisenstein, J., Ahmed, A., & Xing, E.P.E. (2011). Sparse additive generative models of text. *Proc.*
619 *28th Int. Conf. Mach. Learn.* , doi: 10.1.1.206.5167.
- 620 Flint, H.J. (2012). The impact of nutrition on the human microbiome. *Nutr. Rev.* **70**, S10–S13.
- 621 Flint, H.J., Scott, K.P., Duncan, S.H., Louis, P., & Forano, E. (2012). Microbial degradation of
622 complex carbohydrates in the gut. *Gut Microbes* **3**, 289–306.
- 623 Flint, H.J., Graf, D., Cagno, R. Di, Fa, F., Nyman, M., Saarela, M., & Watzl, B. (2015).
624 Contribution of diet to the composition of the human gut microbiota °. *Microb. Ecol.* **1**, 1–11.
- 625 Gelman, A., & Hill, J. (2006). Data analysis using regression and multilevel/hierarchical models
626 (New York, NY: Cambridge University Press).
- 627 Gelman, A., & Rubin, D.B. (1992). Inference from Iterative Simulation Using Multiple
628 Sequences. *Stat. Sci.* **7**, 457–511.
- 629 Gevers, D., Kugathasan, S., Denson, L.A., et al. (2014). The Treatment-Naive Microbiome in
630 New-Onset Crohn’s Disease. *Cell Host Microbe* **15**, 382–392.

- 631 Gilbert, J.A., Quinn, R.A., Debelius, J., Xu, Z.Z., Morton, J., Garg, N., Jansson, J.K., Dorrestein,
632 P.C., & Knight, R. (2016). Microbiome-wide association studies link dynamic microbial consortia
633 to disease. *Nature* **535**, 94–103.
- 634 Gill, S.R., Pop, M., Deboy, R.T., Eckburg, P.B., Turnbaugh, P.J., Samuel, B.S., Gordon, J.I.,
635 Relman, D. a, Fraser-Liggett, C.M., & Nelson, K.E. (2006). Metagenomic analysis of the human
636 distal gut microbiome. *Science* **312**, 1355–1359.
- 637 Gorvitovskaia, A., Holmes, S.P., & Huse, S.M. (2016). Interpreting Prevotella and Bacteroides as
638 biomarkers of diet and lifestyle. *Microbiome* **4**, 15.
- 639 Holmes, I., Harris, K., & Quince, C. (2012). Dirichlet multinomial mixtures: Generative models
640 for microbial metagenomics. *PLoS One* **7**.
- 641 Iwai, S., Weinmaier, T., Schmidt, B.L., Albertson, D.G., Poloso, N.J., Dabbagh, K., & DeSantis,
642 T.Z. (2016). Piphillin: Improved prediction of metagenomic content by direct inference from
643 human microbiomes. *PLoS One* **11**, 1–18.
- 644 Jiang, X., Dushoff, J., Chen, X., & Hu, X. (2012). Identifying enterotype in human microbiome by
645 decomposing probabilistic topics into components. *2012 IEEE Int. Conf. Bioinforma. Biomed.* , doi:
646 10.1109/BIBM.2012.6392720.
- 647 Kembel, S.W., Wu, M., Eisen, J.A., & Green, J.L. (2012). Incorporating 16S Gene Copy Number
648 Information Improves Estimates of Microbial Diversity and Abundance. *PLoS Comput. Biol.* **8**,
649 16–18.
- 650 Knights, D., Costello, E., & Knight, R. (2011). Supervised classification of human microbiota.
651 *FEMS Microbiol. Rev.* **35**, 343–359.
- 652 Knights, D., Kuczynski, J., Charlson, E.S., Zaneveld, J., Mozer, M.C., Collman, R.G., Bushman,
653 F.D., Knight, R., & Kelley, S.T. (2013). Bayesian community-wide culture-independent microbial
654 source tracking. *Nat. Methods* **8**, 761–763.
- 655 Kuhn, M. (2008). Building Predictive Models in R Using the caret Package. *J. Stat. Softw.* **28**, 1–
656 26.
- 657 Kurtz, Z., Mueller, C., Miraldi, E., & Bonneau, R. (2016). SpiecEasi: Sparse Inverse Covariance
658 estimation for Ecological Association and Statistical Inference.
- 659 Kurtz, Z.D., Mueller, C.L., Miraldi, E.R., Littman, D.R., Blaser, M.J., & Bonneau, R.A. (2015).
660 Sparse and Compositionally Robust Inference of Microbial Ecological Networks. *PLoS Comput.*
661 *Biol.* **11**, 1–25.
- 662 Langille, M.G.I., Zaneveld, J., Caporaso, J.G., McDonald, D., Knights, D., Reyes, J. a, Clemente,
663 J.C., Burkpile, D.E., Vega Thurber, R.L., Knight, R., Beiko, R.G., & Huttenhower, C. (2013).
664 Predictive functional profiling of microbial communities using 16S rRNA marker gene
665 sequences. *Nat. Biotechnol.* **31**, 814–821.
- 666 Legendre, P., & Legendre, L. (1998). Numerical Ecology - Second English Edition.

- 667 Lewis, J.D., Chen, E.Z., Baldassano, R.N., et al. (2015). Inflammation, Antibiotics, and Diet as
668 Environmental Stressors of the Gut Microbiome in Pediatric Crohn's Disease. *Cell Host Microbe*
669 **18**, 489–500.
- 670 Li, H. (2015). Microbiome, Metagenomics, and High-Dimensional Compositional Data Analysis.
671 *Annu. Rev. Stat. Its Appl.* **2**, 73–94.
- 672 Martínez, I., Lattimer, J.M., Hubach, K.L., Case, J.A., Yang, J., Weber, C.G., Louk, J.A., Rose, D.J.,
673 Kyureghian, G., Peterson, D.A., Haub, M.D., & Walter, J. (2013). Gut microbiome composition is
674 linked to whole grain-induced immunological improvements. *ISME J.* **7**, 269–280.
- 675 McMurdie, P.J., & Holmes, S. (2013). Phyloseq: An R Package for Reproducible Interactive
676 Analysis and Graphics of Microbiome Census Data. *PLoS One* **8**.
- 677 McMurdie, P.J., & Holmes, S. (2014). Waste Not, Want Not: Why Rarefying Microbiome Data Is
678 Inadmissible. *PLoS Comput. Biol.* **10**.
- 679 Mimno, D., & McCallum, A. (2012). Topic models conditioned on arbitrary features with
680 dirichlet-multinomial regression. *arXiv Prepr. arXiv1206.3278* , doi: 10.1.1.140.6925.
- 681 Ning, J., & Beiko, R.G. (2015). Phylogenetic approaches to microbial community classification.
682 *Microbiome* **3**, 47.
- 683 Quinn, R.A., Navas-Molina, J.A., Hyde, E.R., et al. (2016). From Sample to Multi-Omics
684 Conclusions in under 48 Hours. *mSystems* **1**, e00038-16.
- 685 Ren, B., Bacallado, S., Favaro, S., Holmes, S., & Trippa, L. (2016). Bayesian Nonparametric
686 Ordination for the Analysis of Microbial Communities. *arXiv Prepr. arXiv1601.05156*.
- 687 Roberts, M.E., Stewart, B.M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S.K., Albertson, B.,
688 & Rand, D.G. (2014). Structural topic models for open-ended survey responses. *Am. J. Pol. Sci.*
689 **58**, 1064–1082.
- 690 Roberts, Margaret E., Stewart, B.M., & Tingley, D. (2017). stm: R Package for Structural Topic
691 Models.
- 692 Shafiei, M., Dunn, K.A., Boon, E., MacDonald, S.M., Walsh, D.A., Gu, H., & Bielawski, J.P.
693 (2015). BioMiCo: a supervised Bayesian model for inference of microbial community structure.
694 *Microbiome* **3**, 8.
- 695 Stan Development Team (2016). rstanarm: Bayesian applied regression modeling via Stan.
- 696 Tamang, J.P., Watanabe, K., & Holzapfel, W.H. (2016). Review: Diversity of microorganisms in
697 global fermented foods and beverages. *Front. Microbiol.* **7**.
- 698 Tjalsma, H., Boleij, A., Marchesi, J.R., & Dutilh, B.E. (2012). A bacterial driver-passenger model
699 for colorectal cancer: beyond the usual suspects. *Nat. Rev. Microbiol.* **10**, 575–582.
- 700 De Valpine, P., & Harmon-Threatt, A.N. (2013). General models for resource use or other
701 compositional count data using the Dirichlet-multinomial distribution. *Ecology* **94**, 2678–2687.

702 Woloszynek, S., Zhao, Z., Simpson, G., Mell, J. C., and Rosen, G. Gauging the use of topic
703 models as a means to understand microbiome data structure. *In prep.*

704

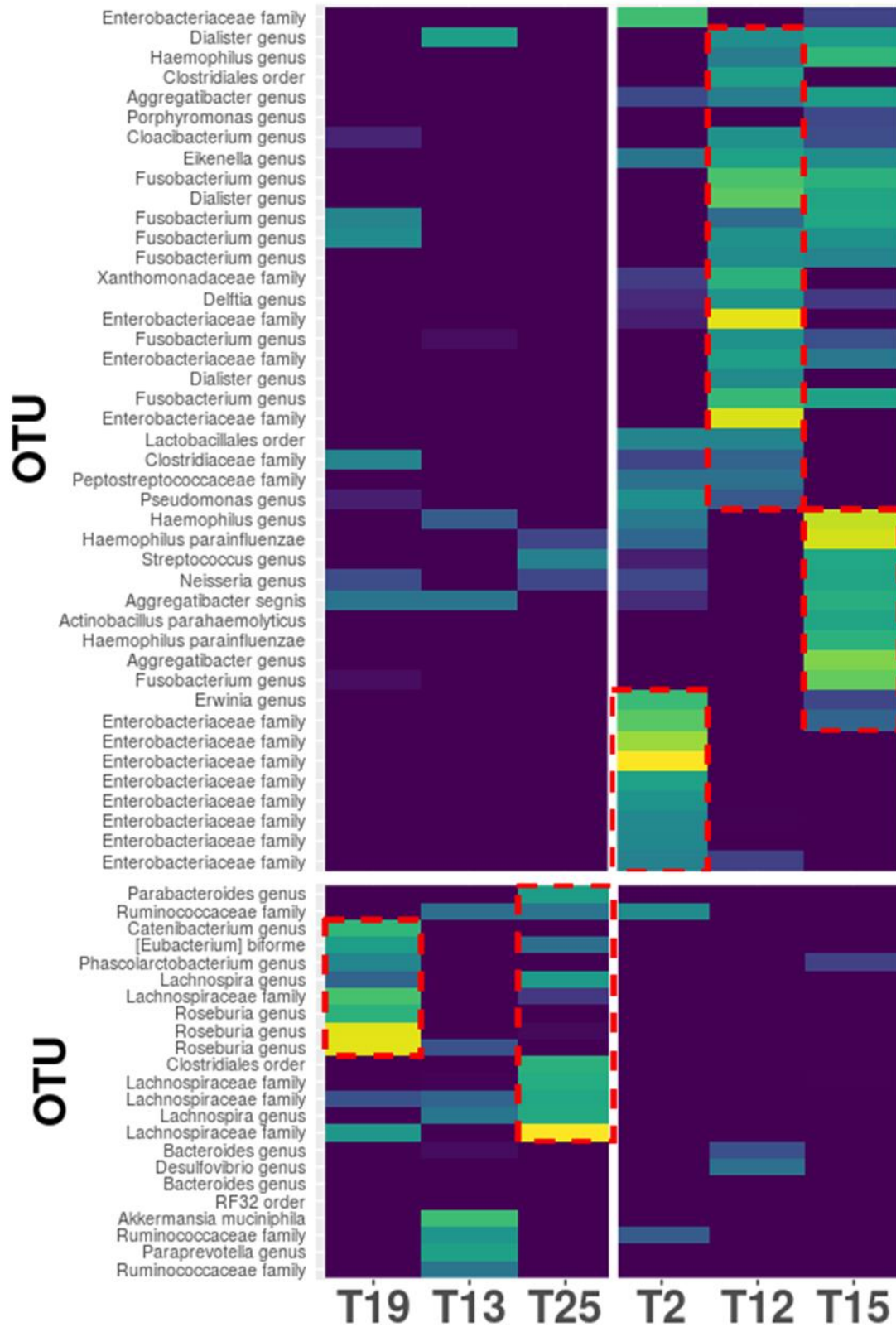


705

706 Figure 1. Given a 16S rRNA abundance table (1), a topic model is used to uncover the thematic structure of the data, represented by
 707 the proportions of topics within samples and OTUs within topics (2). The former is then regressed against sample data (3), to
 708 identify the strength of a topic-covariate relationship to rank the topics (4). The latter is then placed into PICRUSt to generate
 709 predictions of the within-topic functional potential (5). Important pathways are identified via a fully Bayesian multilevel regression
 710 model (6). The end result is the ability to identify clusters of bacteria associated with covariates of interest, in addition to specific
 711 pathways and orthologous groups specific to these clusters.

712

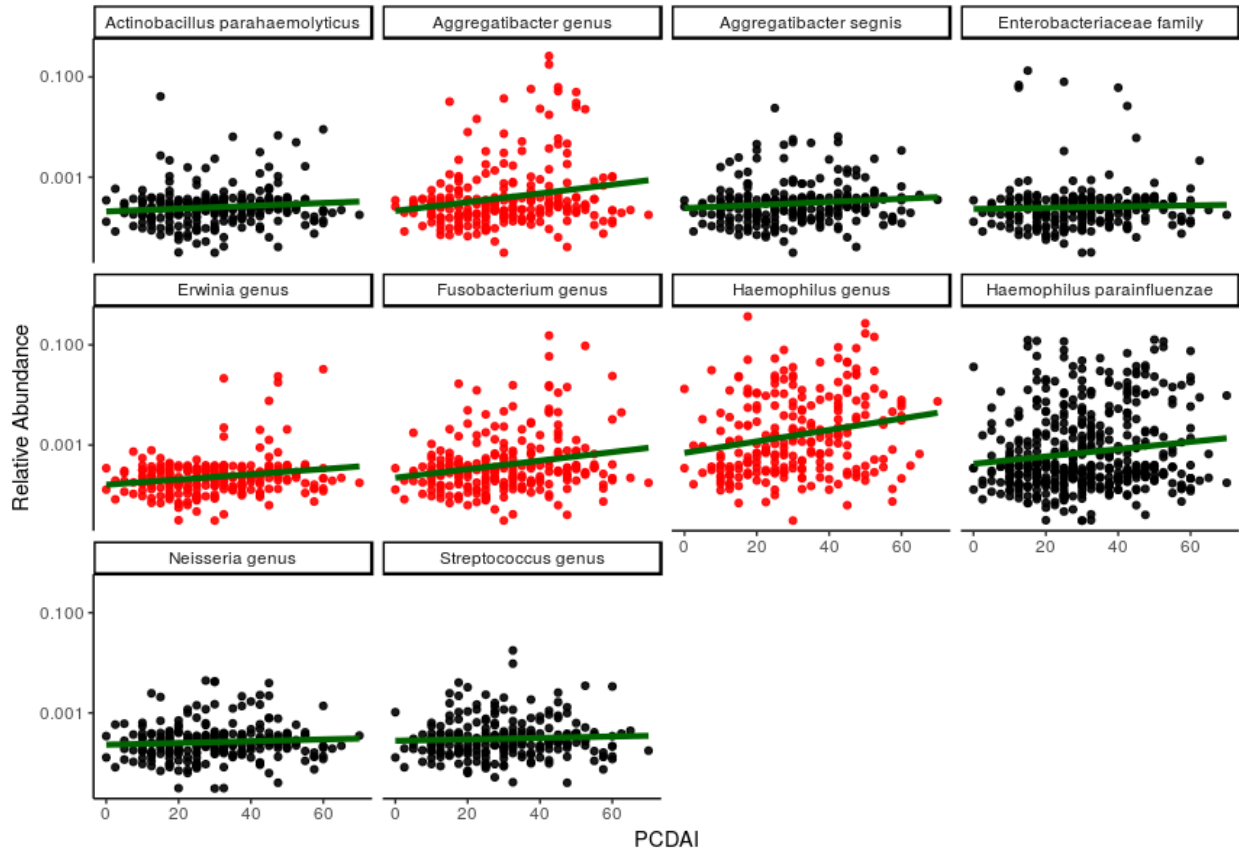
713



714

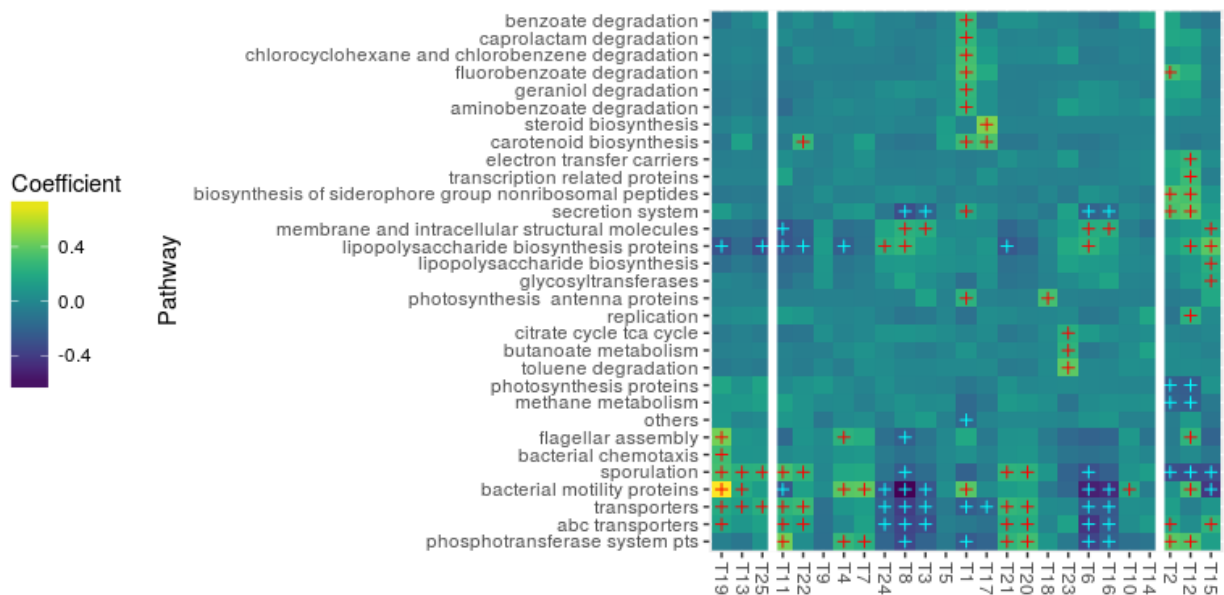
715 Figure 2. Subsections of the heatmap for the Gevers data of the topics over OTU distribution in log space generated by the K25 topic
716 model with covariate prior information. Shown are the top 3 topics associated with CD- and CD+, ordered by mean regression
717 estimate (left to right, respectively, separated by the white line). Clusters of interest are marked with red dotted lines. Clustering
718 was performed via Ward's method on Bray-Curtis distances. Low probabilities ($p < 1 \times 10^{-5}$) are set to 0 to minimize the range of the
719 color gradient to ease visualization. Yellow=high probability, Blue=low probability.

720



721

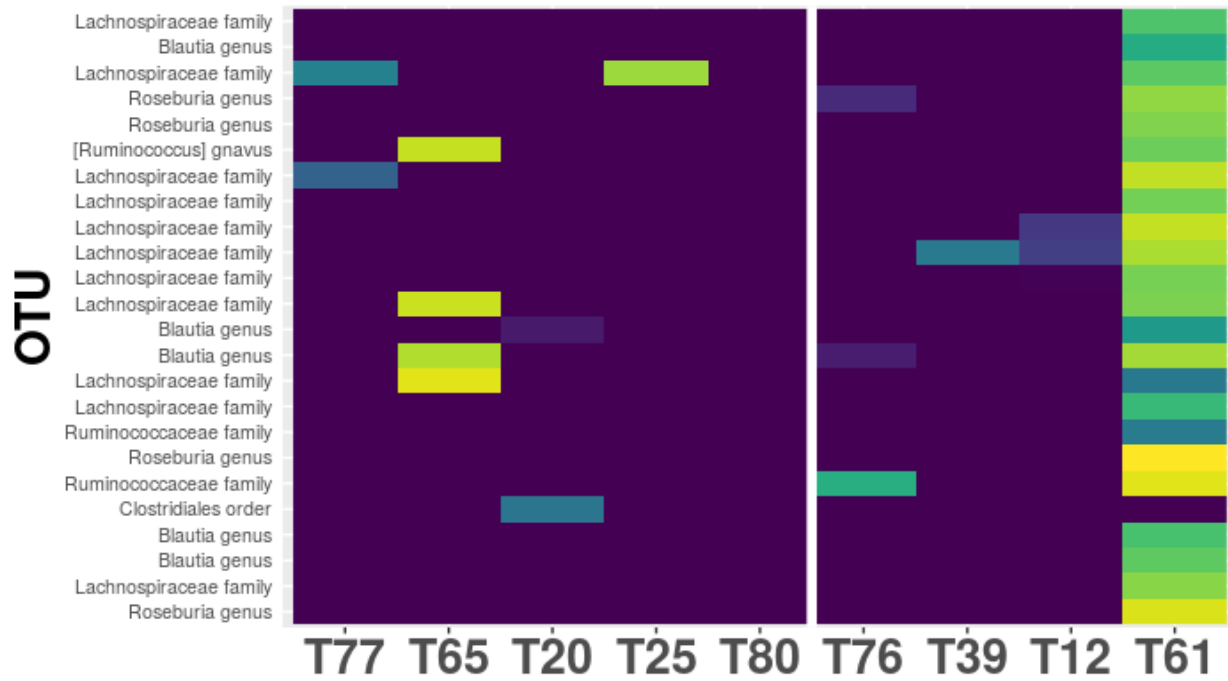
722 Figure 3. Scatterplots of Gevers data for the relative abundance of taxa that compose a high probability cluster in T15 versus PCDAI,
 723 a clinical measure of CD disease burden. Red points reflect significance ($\alpha=0.05$ for negative binomial regression (log linked,
 724 sample coverage offset) with Bonferroni correction.



725

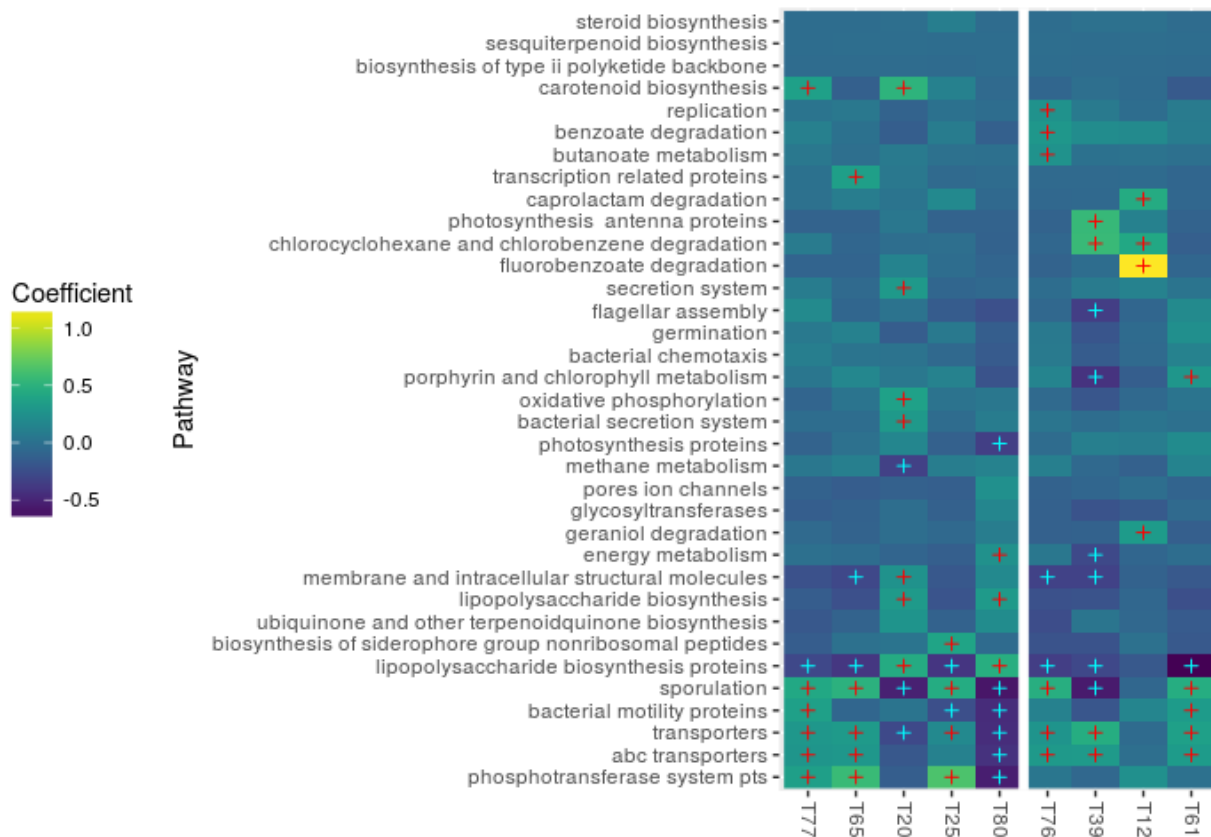
726 Figure 4. Heatmap for Gevers data of the level-3 pathway category-topic interaction regression coefficients from the multiple level
 727 negative binomial model. KEGG information was predicted via PICRUST on the topics over OTU distribution from the K25 topic

728 model with covariate prior information. Topics are ordered based on their mean regression weight when using topic probabilities as
 729 linear predictors for disease presence, where leftmost topics are most associated with CD-, whereas right most topics are most
 730 associated with CD+. Clustering was performed via Ward's method on Bray-Curtis distances. Red and blue crosses indicate weights
 731 or pathway-topic combinations that do not span 0 with 80% uncertainty and are positive or negative, respectively. Only pathways
 732 with at least one such combination are shown.



733
 734 Figure 5. Subsection of the heatmap for AG data for the topics over OTU distribution in log space generated by the K100 topic
 735 model with covariate prior information. Shown are the topics with 95% uncertainty intervals that do not enclose 0 when regressed
 736 against diet type (O=0, V=1), ordered negative to positive by increasing mean regression estimate (left to right), such that T77 is most
 737 associated with O and T61 is most associated with V. The white light signifies a shift from positive to negative means regression

738 estimates. Clustering was performed via Ward's method on Bray-Curtis distances. Low probabilities ($p < 1 \times 10^{-5}$) are set to 0 to
 739 minimize the range of the color gradient to ease visualization. Yellow=high probability, Blue=low probability.



740
 741 Figure 6. Heatmap for AG data of the level 3 pathway category-topic interaction regression coefficients from the multiple level
 742 negative binomial model. KEGG information was predicted via PICRUSt on the topics over OTU distribution from the K100 topic
 743 model with covariate prior information. Only the top 25 topics based on mean regression weight (when using topic probabilities as
 744 linear predictors for disease presence) were chosen for the negative binomial to alleviate computational concerns. Topics are
 745 ordered based on their mean regression weight, where leftmost topics are most associated with O, whereas right most topics are
 746 most associated with V, separated by the white line. Clustering was performed via Ward's method on Bray-Curtis distances. Red
 747 and blue crosses indicate weights for pathway-topic combinations that do not enclose 0 with 80% uncertainty and are positive or
 748 negative, respectively. Only pathways with at least one such combination are shown.

749