1 **TITLE**

2 Uncovering thematic structure to link co-occurring taxa and predicted functional content in 16S
3 rRNA marker gene surveys

4

5 **WORKING TITLE**

6 Uncovering thematic structure in 16S rRNA marker gene surveys

7

8 **AUTHORS**

9 Stephen Woloszynek (sw424@drexel.edu) [1]

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

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

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

13 Gail L. Rosen (gailr@coe.drexel.edu) [1] [corresponding author]

14

15 **AFFILIATIONS**

16 [1] Department of Electrical and Computer Engineering, Drexel University, Philadelphia, PA,
17 United States of America.

18 [2] Department of Microbiology and Immunology, Drexel University College of Medicine,
19 Philadelphia, PA, United States of America.

20 [3] Department of Mathematics, Drexel University, Philadelphia, PA, United States of America

21 [4] Department of Biodiversity, Earth, and Environmental Science, Drexel University,
22 Philadelphia, PA, United States of America

23

24

25

26

27

28

29

30

31 **ABSTRACT**

32

33 **Background:** Analysis of microbiome data involves identifying co-occurring groups of taxa
34 associated with sample features of interest (*e.g.*, disease state). But elucidating key associations
35 is often difficult since microbiome data are compositional, high dimensional, and sparse. Also,
36 the configuration of co-occurring taxa may represent overlapping subcommunities that
37 contribute to, for example, host status. Preserving the configuration of co-occurring microbes
38 rather than detecting specific indicator species is more likely to facilitate biologically
39 meaningful interpretations. In addition, analyses that utilize both taxonomic and predicted
40 functional abundances typically independently characterize the taxonomic and functional
41 profiles before linking them to sample information. This prevents investigators from identifying
42 the specific functional components associate with which subsets of co-occurring taxa.

43 **Results:** We provide an approach to explore co-occurring taxa using "topics" generated via a
44 topic model and then link these topics to specific sample classes (*e.g.*, diseased versus healthy).
45 Rather than inferring predicted functional content independently from taxonomic abundances,
46 we instead focus on inference of functional content within topics, which we parse by estimating
47 pathway-topic interactions through a multilevel, fully Bayesian regression model. We apply our
48 methods to two large publically available 16S amplicon sequencing datasets: an inflammatory
49 bowel disease (IBD) dataset from Gevers *et al.* and data from the American Gut (AG) project.
50 When applied to the Gevers *et al.* IBD study, we demonstrate that a topic highly associated with
51 Crohn's disease (CD) diagnosis is (1) dominated by a cluster of bacteria known to be linked
52 with CD and (2) uniquely enriched for a subset of lipopolysaccharide (LPS) synthesis genes. In
53 the AG data, our approach found that individuals with plant-based diets were enriched with
54 Lachnospiraceae, *Roseburia*, *Blautia*, and *Ruminococcaceae*, as well as fluorobenzoate degradation
55 pathways, whereas pathways involved in LPS biosynthesis were depleted.

56 **Conclusions:** We introduce an approach for uncovering latent thematic structure in the context
57 of sample features for 16S rRNA surveys. Using our topic-model approach, investigators can (1)
58 capture groups of co-occurring taxa termed topics, (2) uncover within-topic functional potential,
59 and (3) identify gene sets that may guide future inquiry. These methods have been
60 implemented in a freely available R package https://github.com/EESI/themetagenomics.

61

62 **KEYWORDS**

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

65

66

67

68    **BACKGROUND**

69

70    With the decreasing cost of high-throughput sequencing, large datasets are becoming
71    increasingly available, particularly microbiome datasets rich in sample data. These data include
72    categorical and numeric features associated with each sample, which, in turn, may be linked to
73    a set of taxonomic abundances that are derived from clustering sequencing reads. Typically,
74    taxonomic marker genes, such as a portion of the 16S rRNA gene common to all bacteria, are
75    used to perform the clustering based on a fixed degree of sequence similarity among reads.
76    These clusters are termed Operational Taxonomic Units (OTUs), and each OTU corresponds to
77    a taxonomic level, such as a genus.

78    Analysis of OTU abundances often involves identifying OTUs associated with specific sample
79    features (*e.g.*, body site, disease presence, diet, age) via unsupervised exploratory methods such
80    as principal component analysis, correspondence analysis, multidimensional scaling, and
81    hierarchical clustering. Analyses may also include statistical inference strategies aimed at
82    identifying differentially abundant OTUs and differences in alpha and beta diversity.
83    Nevertheless, model building is hindered by the complexity inherent to abundance data of this
84    type, which have a disproportionate number of OTUs relative to samples [1], a substantial
85    degree of sparsity, and are typically strictly non-negative and constrained to sum to 1, *i.e.*,
86    compositional  [2, 3].

87    From an ecological perspective, co-occurring OTUs may represent related, overlapping groups
88    of taxa that correlate with, for example, host status. Identifying important groups of co-
89    occurring taxa, which we will refer to as "subcommunities," facilitates a more biologically
90    meaningful interpretation than identifying single indicator OTUs. This is because identifying
91    subcommunities preserves the natural groupings of taxa when making inferences with respect
92    to important sample features [4–7].

93    Still, suitable approaches are lacking that uncover the potential mechanistic relationships
94    between subcommunities and sample features, namely how functional content within a
95    subcommunity could correlate with sample features, such as genes belonging to a particular
96    metabolic pathway correlating with disease status. That is, few methods successfully integrate
97    subcommunity and sample features with functional profiles specific to these subcommunities.

98    In 16S rRNA survey approaches, analyses using taxonomic and functional abundance
99    information typically involve independently inferring how the taxonomic and functional
100   profiles of the samples associate with host features (Figure 1). The functional profiles are
101   predicted based on taxonomic abundances with methods that use preexisting gene annotations.
102   Examples of these methods include PICRUSt, Tax4fun, and Piphillin [8–10]. Because the
103   inference of taxonomic and functional abundances occurs in two stages, it remains difficult to
104   identify which functional content associates with which subsets of co-occurring taxa.

3

105  In the context of 16S rRNA surveys, we had two objectives: (1) implement a model framework
106  that identifies subcommunities associated with specific sample features, and (2) uncover
107  functional properties that further characterize these subcommunities.

108  To satisfy our first objective, we employed a topic model approach. Topic models have had
109  considerable use in natural language processing, but have also shown promise as a method for
110  exploring taxonomic abundance data. Knights *et al.* [11] used latent Dirichlet allocation (LDA) to
111  infer the relative contributions of an unknown number of source environments to a set of
112  indoor samples. Shafiei *et al.* [4], alternatively, took a supervised approach where they first
113  trained their model on sets of co-occurring OTUs to learn how they correlate with sample
114  classes of interest. They were then able to predict the class of new samples given the trained
115  model.

116  Our approach uses a structural topic model (STM) [12], which generalizes previously described
117  topic models such as LDA, the correlated topic model [13], and the Dirichlet-Multinomial
118  regression topic model [14]. Like the Dirichlet-Multinomial regression topic model, the STM
119  permits the use of sample features to help inform the frequency of topics occurring in a given
120  sample. LDA, on the other hand, can only incorporate sample information if done in a two-
121  stage process – first performing topic extraction, and then identifying linear relationships
122  between the topic assignments and sample features [15]. A two-stage approach limits the
123  breadth of sample information one can use, typically forcing the user to use only a single
124  sample feature [12]. It also prevents propagating uncertainty throughout the model. Similar to
125  the correlated topic model, the STM's Logistic-Normal (LN) distribution defines the frequency
126  of topics for a given sample and permits correlation between topics.

127  We use the STM to uncover a thematic representation of 16S rRNA survey abundance data and
128  jointly measure its relationship with sample features (Figure 1). The aim is to cluster co-
129  occurring OTUs into overlapping "topics," where a given OTU can occur in multiple topics,
130  albeit with varying frequency. The model also estimates the frequency of each topic in each
131  sample based on that sample's OTU abundances. Thus, the STM infers two latent, unobserved
132  distributions from a table of OTU abundances: a samples-over-topics distribution (the
133  frequency of each topic in each sample) and a topics-over-OTUs distribution (the frequency of
134  each OTU in each topic). In addition, by utilizing sample features as covariates, the STM can
135  also determine whether particular features increase or decrease the frequency of a topic
136  occurring in a given sample, providing a means to identify topics associated with sample
137  features.

138  Our second objective is to exploit the estimated topics-over-OTUs distribution, which dictates
139  the taxonomic composition of each topic and therefore should capture meaningful
140  subcommunities. We can therefore infer the functional content of individual topics using tools
141  such as PICRUSt and a database of gene annotations.  By identifying topics of interest based on
142  their relationship to specific sample features and then inferring within-topic functional profiles,
143  we can infer the specific functional content within subcommunities associated with sample
144  features.

145    We apply our approach on two large 16S rRNA survey datasets: an inflammatory bowel disease
146    (IBD) dataset from Gevers *et al.* [16] and data from the American Gut (AG) project. After
147    confirming the generalizability of extracted topics, we identified distinct taxonomic
148    subcommunities that, in the case of the Gevers *et al.* dataset, were consistent with published
149    results. These subcommunities were composed of distinct functional profiles. Also, our
150    approach provided gene sets specific to topics of interest that may warrant further exploration.

151    These methods have been implemented in a freely available R package themetagenomics:
152    https://github.com/EESI/themetagenomics. In a companion paper, we performed simulations to
153    further validate using topic models for 16S rRNA survey data and to determine a suitable
154    normalization strategy (Woloszynek, S., Zhao, Z., Simpson, G., Mell, J. C., and Rosen, *in prep*).

155

156    **METHODS**

157

158    **Review of the Structural Topic Model**

159    The STM is a Bayesian generative model such that, given a set of $M$ samples, each consisting of
160    $N$ OTUs, belonging to a vocabulary of $V$ unique OTU IDs, $K$ (chosen *a priori*) latent topics are
161    assumed to be generated from the data. These topics consist of overlapping groups of co-
162    occurring OTUs. A LN prior is placed on the samples-over-topics distribution. This allows for
163    estimation of topic-topic correlations, giving a means to infer co-occurring topics across
164    samples. The topics-over-OTUs distribution estimates deviation of OTU frequencies from the
165    background OTU distribution [17]. Sparsity inducing priors are placed on the topics-over-OTUs
166    distribution. This ensures a sparse set of estimates, ideal for high dimensional data. Lastly,
167    word and topic assignments are both generated via multinomial distributions with $V$ and $K$
168    classes, respectively. For the relationships between topic model nomenclature and our
169    terminology, see Table 1.

170    The STM is estimated by a semi-collapsed variational expectation maximization procedure [18].
171    Convergence is reached when the relative change in the variational objective (*i.e.,* the estimated
172    lower bound) in successive iterations falls below a predetermined tolerance.

173

174    **Datasets and Preprocessing**

175    16S rRNA survey data from two human microbiome studies were downloaded from their
176    corresponding repositories. The Gevers *et al.* dataset ("Gevers") (PRJNA237362, 03/30/2016) is a
177    multicohort, IBD dataset that includes control, Crohn's disease (CD), and ulcerative colitis
178    samples taken from multiple locations throughout the gastrointestinal tract [16]. The American
179    Gut project ("AG") (ERP012803, 02/21/2017) is a crowd-sourced dataset that includes user-
180    submitted microbiome samples from a variety of body sites and self-reported subject
181    information (http://americangut.org/).

182 **Human gut microbiota from an inflammatory bowel disease cohort (Gevers).** Paired-end
183 reads were joined and quality filtered (maximum unacceptable Phred quality score = 32;
184 maximum number of consecutive low quality base calls before read truncation = 3; minimum
185 number of consecutive high quality base calls included per read as a fraction of input read
186 length = 0.75) using QIIME version 1.9.1. Closed-reference OTU picking was performed using
187 SortMeRNA against GreenGenes v13.5 at 97% sequence identity. This was followed by copy
188 number normalization via PICRUSt version 1.0.0 [19].

189 We selected only terminal ileum samples. Samples with fewer than 1000 total reads were
190 omitted. We removed OTUs with fewer than 10 total reads across samples and OTUs that
191 lacked a known classification at the phylum level.

192 **Human gut microbiota from vegetarian and omnivore subjects (AG).** Quality trimming and
193 filtering were performed in the following manner on single-end reads using the fastqFilter
194 command found in the dada2 R package [20]. The first 10 bases were trimmed from each read.
195 Reads were then trimmed to position 135 based on visualizing the quality score of sampled
196 reads as a function of base position. Further truncation occurred at positions with quality scores
197 less than or equal to 2. Any truncated read with total expected errors greater than 2 were
198 removed. A portion of AG samples were affected by bacterial blooming during shipment. These
199 reads were removed using the protocol provided in the AG documentation (02-
200 filter_sequences_for_blooms.md).

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

208

209 **Fitting Structural Topic Models**

210 Each resulting OTU table consisted of sets of raw counts normalized by 16S rRNA copy
211 number. No other normalization was conducted based on the simulation results in Woloszynek
212 *et al.* (*in prep*). A series of STMs with different parameterizations in terms of topic number (K ∈
213 15, 25, 50, 75, 100, 150, 250) and sample features (*e.g.,* no features, indicators for presence of
214 disease, diet type, etc.) were fit to the OTU tables. STMs were fit via the R package stm [21].

215 We evaluated each model fit for presence of overdispersed residuals. We also conducted
216 permutation tests (permTest in the stm package) where the sample feature of interest is
217 randomly assigned to a sample, prior to STM fitting. To compare parameterizations between
218 models, we evaluated predictive performance using held-out likelihood estimation [15].

219

## Assessing Topic Generalizability

We performed classification to assess the generalizability of the extracted topics. No sample features were used as covariates. OTU abundances were split into 80/20 training-testing datasets. For different number of topics (K ∈ 15, 25, 50, 75, 100, 150), an STM was trained to estimate the topics-over-OTUs distribution. We then held this distribution fixed; hence, only the testing set's samples-over-topics distribution was estimated. For both the training and testing sets, simulated posterior samples from the samples-over-topics distribution were averaged. The resulting posterior topic frequencies in the training set were then used as features to classify sample labels, similar to using $\bar{Z}$ in supervised LDA [22]. Generalization (testing) error was assessed using the optimal parametrization based on cross-validation performance on the test set topic frequencies. Classification was performed using a random forest.

For the random forest, parameter tuning to determine the number of variables for each split was accomplished through repeated (10x) 10-fold cross-validation, using up- or down-sampling to overcome class imbalance (for Gevers and AG, respectively). We performed a parameter sweep over the number of randomly selected OTU features, while setting the number of trees fixed at 128. The optimal parameterizations were selected based on maximizing ROC area under the curve.

The performance of the STMs was compared to the performance using OTUs as features from the starting OTU abundance table. Separately, training and testing set OTU abundances were converted to relative abundances with the following equation: $OTU_{n,m} / \sum_n OTU_{n,m}$. In words, OTU *n* for sample *m* is scaled by the library size of sample *m* (the total abundance of sample *m*). The resulting OTU relative abundance tables were separately z-score normalized. Training cross-validation and testing using a random forest was then performed as above.

## Assessing Concentration of OTUs as a Function of Topic Number

For each STM (K ∈ 15, 25, 50, 75, 100, 150, 250), Shannon entropy was calculated for each topic in the topics-over-OTUs distribution. To compare mean entropy across STMs, we performed an ANOVA, followed by Tukey HSD post-hoc analysis.

## Identifying Within-Topic Clusters of High Frequency OTUs

Using the topics-over-OTUs distribution, we performed hierarchical clustering via Ward's method on Bray-Curtis distances. We refer to high frequency groups of OTUs as "clusters."

## Inferring Within-Topic Functional Potential

We obtained the topics-over-OTUs distribution for each fitted model and mapped the within-topic OTU probabilities to integers ("pseudo-counts") using a constant: $10000 \times \beta$. A large constant was chosen to prevent low frequency OTUs from being set to zero, although their

7

257  contribution to downstream analysis was likely negligible. Gene prediction was performed on
258  each topic-OTU pseudo-count table using PICRUSt version 1.0.0 [8]. (Normalization of 16S copy
259  number was performed prior to topic model fitting using PICRUSt.) Predicted gene content was
260  classified in terms of KEGG orthology (KOs) [23].

261

262  **Identifying Topics of Interest**

263  Topics of interest were identified using the samples-over-topics distribution as a matrix of $K$
264  covariates in a linear regression model. Each column in the samples-over-topics distribution
265  represents the frequency of topic $k$ in sample $m$. The dependent variable chosen depended on
266  the sample feature used as a covariate during model fitting. These include CD presence and
267  PCDAI for Gevers and diet type for AG. We calculated 95% uncertainty intervals using an
268  approximation that accounts for uncertainty in estimation of both the sample covariate
269  coefficients and the topic frequencies. We will refer to these coefficients as "topic-effects."
270  Coefficients whose 95% uncertainty intervals do not span 0 will be referred to as "high ranking
271  topics."

272

273  **Identifying Functional Content that Distinguishes Topics**

274  To determine which predicted functional gene content best distinguished topics, we used the
275  following multilevel negative binomial regression model:

276
$$\theta_{k,c} = \exp[\mu + \beta_k + \beta_c + \beta_{k,c}]$$

277
$$y_{k,c} \sim \mathrm{NB}(\theta_{k,c}, \lambda)$$

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

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

280  for a given topic-function combination, and $\lambda$ is the dispersion parameter. The intercept $\mu$ was
281  given a $\mathrm{Normal}(0, 10)$ prior; all weights were given $\mathrm{Normal}\ 0, 2.5$  priors; and the dispersion
282  parameter $\lambda$ was given a $\mathrm{Cauchy}(0, 5)$ prior.

283  Model inference was performed using Hamiltonian Monte Carlo in the R package rstanarm
284  [24]. Convergence was evaluated across four parallel chains using diagnostic plots to assess
285  mixing and by evaluating the Gelman-Rubin convergence diagnostic [25]. To reduce model size,
286  we used genes belonging to only 15 (arbitrary number) level-2 KEGG pathway categories (Table

8

287  S1). For large topic models, we fit only the top 25 topics, ranked in terms of topic-effects that
288  measure the degree of association between sample-over-topic probabilities and our sample
289  feature of interest.

290

### Assessing Relationships Between Sample Features of Interest and Taxonomic Abundance
291

292  To quantify the relationship between taxonomic abundance and continuous sample features
293  (such as the Pediatric Crohn's Disease Activity Index (PCDAI), a clinical measure of CD
294  severity), we performed negative binomial regression (log-link), using sample library size (sum
295  of OTU abundances across samples) as an offset. The family-wise error rate was adjusted via
296  Bonferroni correction. Critical values for hypothesis testing was set at 0.05.

297

### Comparing Topic Taxonomic Profiles to a Network Approach
298

299  To further validate the clusters of high frequency taxa identified in the topics-over-OTUs
300  distribution, we compared our results to those generated from an OTU-OTU association
301  network on the copy number normalized OTU abundances using SPIEC-EASI's neighborhood
302  selection method (lambda.min.ratio=.01, nlambda=20) [26].

303

### Comparing Within-Topic Functional Profiles to an OTU-Abundance-Based Approach
304

305  We compared the results from the hierarchical negative binomial model to a differential
306  abundance approach. We performed predicted functional content using PICRUSt on copy
307  number normalized OTU abundances. The resulting functional abundances were collapsed into
308  level-3 KEGG pathways. Note that, for consistency, we again restricted our genes to the 15
309  level-2 KEGG pathways used previously. The resulting level-3 pathway abundances underwent
310  DESeq2 differential abundance analysis, which uses negative binomial regression and variance
311  stabilizing transformations to infer the difference log-fold change of OTU abundance [27, 28].
312  The resulting p-values were corrected via the Bonferroni method. Adjusted p-values below 0.1
313  were considered significant.

314

### Packages utilized
315

316  All analysis was done in R version 3.2.3. Topic models, random forests, and NB regression
317  models were fit using stm [21], caret [29], and rstanarm [24], respectively. AG filtering was
318  performed using DADA2 [20]. SPIEC-EASI was fit using the SPIEC-EASI package [26]. DESeq2
319  differential abundance analysis was conducted with phyloseq [30]. Shannon entropy was
320  performed with vegan [31].

321

### Implementation
322

9

323 Our approach can be implemented with themetagenomics, an R package that provides a topic
324 model framework for microbiome abundance data, as well as functional prediction for 16S
325 rRNA survey data. Users can choose between our C++ implementation of PICRUSt or Tax4fun
326 functional prediction; therefore, both GreeneGenes and Silva taxonomic annotations are
327 acceptable. Inference of topic-function interactions can be accomplished by maximum
328 likelihood or Hamiltonian Monte Carlo using Rstan, where users can choose between Student-t,
329 Laplace, and Normal prior distributions. The resulting topics, topic-effects, and topic-function
330 interactions can then be explored with a variety of interactive Shiny apps.

331

332 **RESULTS**

333

334 Here we explored the use of our structural topic model (STM) approach on publically available
335 datasets of gut and fecal microbiota, first using the IBD data from Gevers *et al.* [32], and second
336 using the dietary data from AG. For each dataset, we show that the topics extracted from the
337 STM generalize well to test data not initially seen by the model, suggesting that co-occurrence
338 profiles identified by the STM are robust to overfitting. Then, we apply our complete pipeline,
339 where we successfully link, within a given topic, functional content, taxonomic co-occurrence
340 profiles, and sample features of interest.

341

342 **Thematic Structure of IBD-Associated Microbiota (Gevers)**

343 **Dimensionality reduction using topics facilitates classification of CD diagnosis and**
344 **generalizes well to test data.** We sought to assess (1) if topics were associated with positive CD
345 diagnosis (CD+) and (2) whether those topics generalize to new data – that is, did they capture
346 meaningful information inherent to the data while ignoring characteristics associated
347 exclusively with the fitted data. Note that, for this analysis, STM fitting was completely
348 unsupervised, so no sample features were used as covariates.

349 The 80/20 training/testing splits for terminal ilium samples from Gevers are shown in Table S2.
350 We hypothesized that using topics would outperform the relative abundance of OTUs as
351 features for classifying CD diagnosis, since the relative abundance-based features are sparser.
352 Both dimensionality and sparsity are reduced when using topics, since the size of the feature
353 space is decreased through dimensionality reduction. There was little difference between the
354 two approaches during training cross-validation with at least 25 topics (Figure S1, Table S3).
355 During testing, however, topics outperformed OTUs, particularly in F1 score (a measure of
356 performance that considers both sensitivity and positive predictive value), with scores of 0.808
357 for OTUs and at least 0.821 for all models with 25 or more topics (Table S4).

358 The largest discrepancy in classification performance between OTUs and topics was the
359 proportion of true negatives out of total negative classifications (negative predictive value). The
360 OTU model correctly identified CD- subjects only half the time (0.517), whereas the worst

10

361 performing topic model (K=15) performed slightly better (0.526), and topic models improved as
362 the number of topics increased: 0.655 (K=25), 0.559 (50), 0.577 (75), 0.682 (100), and 0.643 (150)
363 (Table S4).

364 The substantially higher proportion of false negatives with the OTU model was likely due to its
365 reliance on few, relatively rare taxa. For example, the random forest importance scores
366 indicated that OTU 319708 (Clostridiaceae family) was the fourth most important feature for
367 distinguishing CD+ from CD-. It was more than twice as common in CD- training samples, and
368 more than 10% of correctly classified CD- samples contained this feature. When OTU 319708
369 was present in CD+ samples, it may have led these samples to be misclassified. Approximately
370 10% of misclassified CD+ samples contained this feature, and some of these samples contained
371 it at a greater proportion than other samples in the training set. A similar scenario can be seen
372 for OTU 186723 (Ruminococcaceae family), which received the highest random forest
373 importance score for classifying CD. It was most common in CD+ samples; hence, its absence in
374 CD+ samples resulted in false negatives.

375 **Concentration of high probability OTUs across topics begins to plateau at 75 topics.** We fit
376 STM to the OTU abundance data and aimed to uncover how specific OTUs concentrate within
377 topics as a function of topic number K (15, 25, 50, 75, 100, 150). We measured concentration
378 using Shannon entropy. We define high quality topics as topics that place high probability on
379 only a few OTUs; thus, high quality topics will have low entropy. A topic that is characterized
380 by a small subset of OTUs is (1) more interpretable as a subcommunity and (2) contains more
381 easily detectable associations with host features of interest.

382 For each K, we calculated the Shannon entropy for each topic, showing decreased entropy with
383 increased topic number. This was supported by a one-way ANOVA ($p<0.0001$, $F_{5,409}=8.327$) and
384 post-hoc pairwise comparisons testing using Tukey HSD ($\alpha=0.05$) (Figure S2). Among pairwise
385 combinations, we found that models with 75 or more topics did not have significantly different
386 Shannon entropies, leading us to focus our attention on topic models with at most 75 topics.

387 **CD diagnosis was associated with distinct thematic profiles and hence distinct**
388 **subcommunity structure.** We implemented our full pipeline using sample features as
389 covariates. Here, we used a binary indicator for CD diagnosis. We identified topics-of-interest
390 based on their "topic-effects" – the regression coefficients estimated when regressing the
391 samples-over-topics distribution against CD presence. We also performed permutation tests to
392 ensure that detected topic-effects were not spurious effects. For model K25, we performed 25
393 permutations and calculated the mean posterior regression coefficient for each topic. Of the 25
394 topics, the 95% uncertainty intervals for 8 topics did not span 0 (Figure S3). We consider these
395 "high ranking topics." Topics T15, T12, T2, and T14 had estimates greater than 0, whereas topics
396 T11, T25, T13, and T19 had estimates less than 0. For K75, 14 topics did not span 0 (Figure S4).

397 The posterior predictive distribution of each sample's topic assignment is shown in Figure 2 for
398 STMs K25 and K75. Each sample is plotted based on its PCDAI, which is 0 for CD- samples and
399 increases as CD severity increases. Topics are ordered in terms of their topic-effects. Both panels
400 demonstrate that as PCDAI increases, the thematic profile changes towards topics with higher

11

401     correlations to CD+. Also, high ranking topics (cyan and red points for CD- and CD+,
402     respectively) concentrate at the extremes, suggesting that the OTU abundance profiles for
403     healthy and high severity CD samples were more precisely captured by the STM. Notably, there
404     was clearer separation between CD- and CD+ associated high ranking topics for the K25 model.
405     The transition point can clearly be seen at approximately PCDAI=35. Based on this result, we
406     will hence focus on the K25 model.

407     Focusing on the high ranking topics for K25 (T19, T13, T25, T11; T14, T2, T12, T15), we
408     identified multiple clusters of bacterial species (via hierarchical clustering) that
409     disproportionately dominated the high ranking topics associated with CD+ (Figure 3A). T2
410     contained a cluster dominated by *Enterobacteriaceae* taxa, whereas T12's cluster contained a
411     mixture of *Fusobacteria* and *Enterobacteriaceae*. The T15 cluster contained *Haemophilus* spp.,
412     *Neisseria*, *Fusobacteria*, and *Streptococcus*, all of which were noted as having a positive correlation
413     with CD+ subjects in Gevers *et al.*, as well as *Aggregatibacter*, a genus reportedly associated with
414     colorectal cancer [33].

415     Given that T15 contains a cluster of bacteria known for their association with bowel
416     inflammation and this topic occurs disproportionately in subjects with greater disease severity,
417     we asked whether the abundances of the OTUs in T15 correlated with PCDAI. After performing
418     negative binomial regression (Figure 4), we identified significant positive trends as a function of
419     PCDAI for *Aggregatibacter* ($p<0.0001$, $\beta=0.089$, $Z=5.285$), *Erwinia* ($p=0.0004$, $\beta=0.103$, $Z=4.116$),
420     *Fusobacterium* ($p=0.0001$, $\beta=0.081$, $Z=6.354$), and *Haemophilus* ($p=0.0484$, $\beta=0.0264$, $Z=2.847$).

421     The high ranking topics for CD-, on the other hand, were dominated by taxa belonging to
422     *Lachnospiraceae*, *Roseburia*, *Rubinococcus*, *Blautia*, *Bacteroidetes*, and *Coprococcus*, all of which were
423     noted by Gevers *et al.* as being negativity associated with CD (Figure 2B). In addition to these
424     taxa, *Akkermania*, *Dialister*, and *Dorea* contributed to these topics, which is consistent with the
425     findings of Lewis *et al.* who found a reduction of these taxa in CD+ subjects [34].

426     **Within-topic co-occurrence profiles were consistent with SPIEC-EASI.** We compared topics to
427     the correlations obtained via a network approach. The edges in the SPIEC-EASI network for the
428     clusters of high probability OTUs in our high ranking topics are shown in Figure S5. For each of
429     these topic clusters, the majority of taxa were connected by a non-zero edge (Table S5). Of the 11
430     taxa in the T15 cluster, 8 had first-order connections (direct connections to other taxa within the
431     cluster, $OTU_c$-$OTU_c$), whereas 9 had second-order connections (indirect connections to other
432     taxa within the cluster via an intermediate OTU not present in the cluster, $OTU_c$-$OTU_{nc}$-$OTU_c$).
433     The two OTUs connected by the largest edge weight, *H. parainfluenzae* and *Haemophilus spp.*, had
434     the highest frequencies in T15, 0.320 and 0.245, respectively. Of topics T15, T12, T2, T19, T13,
435     and T25, none had more than one OTU with zero connections or fewer than 75% of taxa joined
436     by first-order connections. The taxa that lacked within-cluster connections generally had low
437     topic frequencies, with one exception, *Catenibacterium spp.* in T19. Taken together, this reaffirms
438     that the within-topic co-occurrence profiles are consistent with alternative approaches.

439     **Predicted functional potential of notable topics further elucidated their association with CD.**
440     We sought to further explore the functional content of topics, thereby exploiting the posterior

441 estimates of the STM in a way other approaches have not. To do so, we predicted the functional
442 content within topics using PICRUSt, which infers the metabolic potential of a microbial
443 community by matching taxonomic classifications made from 16S rRNA gene sequences with a
444 closely related reference genome annotation. We then performed a fully Bayesian multilevel
445 regression analysis on the predicted abundances of each gene to identify topic-function
446 interactions. This identified the bacterial taxa (irrespective of taxonomy) that drive the
447 functional associations.

448 Like Gevers *et al.*, we identified an increase in membrane transport associated with CD+
449 subjects' gut microbiome; however, through our approach, we were able to pinpoint the specific
450 topics (*i.e.*, subcommunities) associated with the enrichment of these functional categories,
451 topics T2 and T12 (Figure 2C). We then could link enrichment of membrane transport genes to
452 the taxa that were also enriched in this topic. For example, topics T2 and T12 were dominated
453 by Enterobacteriaceae. The Enterobacteriaceae-enriched topics (T2, T12) were also enriched for
454 siderophore and secretion system related genes. Like T2 and T12, T15 was highly associated
455 with CD+; however, it was less enriched for membrane transport genes. This suggests that the
456 cluster of bacteria found in T15 (*Haemophilus* spp., *Neisseria*, and *Fusobacteria*) may have
457 contributed less to the shift of transport genes reported by Gevers *et al.* and instead have
458 distinct functional associations with CD.

459 The largest topic-function interaction effect was found in T19 for genes encoding bacterial
460 motility proteins. For T19, three motility-related KEGG categories (bacterial motility proteins,
461 bacterial chemotaxis, flagellar assembly) had topic-function interaction effects that did not span
462 0 at 80% uncertainty, suggesting that T19 was more enriched in cell motility genes relative to all
463 other topics. The gene functions inferred for T19 are consistent with this taxonomic profile,
464 consisting of motile bacteria belonging to Lachnospiraceae, Roseburia, and Clostridiales.
465 Enrichment of two lipopolysaccharide (LPS) synthesis categories were associated with CD+
466 topics; however, one of these categories was specific for only T15 (Table S7).

467 **More functional categories were significant via a DESeq2 approach on the OTU abundance**
468 **table.** We compared our within-topic functional profiles to the results obtained by performing
469 PICRUSt on the copy-number normalized OTU abundance table and then performing a DESeq2
470 differential abundance analysis. Of the 160 level-3 KEGG categories, more than half (87) were
471 found significant ($\alpha < 0.1$), in the DESeq2 approach, leading to a difficulty in making meaning
472 out of the data (Figure 3D). Pathways with the largest log-fold change (LFC) associated with
473 CD+ samples included degradation pathways (caprolactam, LFC=0.542; fluorobenzoate, 0.532;
474 geraniol, 0.371; and toluene degradation, 0.371), alphalinolenic acid metabolism (0.641), and
475 electron transfer carriers (0.635).

476 Interestingly, the degradation pathways associated with CD+ also demonstrated strong topic-
477 function interaction effects; however, they associated most strongly with T1, a topic unrelated to
478 CD presence. Predicted electron transfer carrier genes were identified by both approaches, but
479 the topic model approach isolated the effect to T12, placing high probability on bacteria that are
480 also enriched for functions linked to secretion systems, LPS biosynthesis, and motility.

481      The topic-free DESeq2 approach also found fewer categories associated with CD- that had large
482      LFC. For example, only one category had a LFC less than -0.4, whereas there were 8 greater than
483      0.-4 (enrichment in CD+). The categories with the largest LFCs relative to CD- included
484      germination (LFC=-0.450) and sporulation (-0.346). Similarly, the topic model identified 10
485      topics with functional profiles significantly enriched or depleted in sporulation genes, three of
486      which were associated with CD- samples. Multiple topics demonstrated an inverse relationship
487      between sporulation and LPS genes, such that topics that contained taxa enriched in one were
488      depleted in the other.

489

490      **Thematic Structure of Diet-Associated Microbiota (AG)**

491

492      Despite consisting of far more samples, the AG dataset, split into omnivore (O) and vegetarian
493      (V) diet groups from self-reported dietary information, offered a new challenge for our
494      approach, given that there were far more data and taxonomic features (OTUs), as well as severe
495      imbalance between classes. Of the 4864 samples that fit into our diet groups, 4527 were
496      identified as O and only 337 as V samples. This disparity of group sizes limits the utility of
497      comparing group means, particularly as an estimate of covariate effects [35].

498      **BMI is a potential confounder.** Before applying our pipeline, we considered potential sources
499      of confounding. Male and female samples were distributed similarly with respect to diet (Table
500      S9; Figure S6). There was no significant difference in mean age between diet groups (t=-0.03,
501      df=373.93, p=0.98). Sample body mass index (BMI) was not normally distributed (Shapiro-Wilk:
502      W=0.86, p<0.001) and was plagued with many mislabeled heights and weights (Figure S7). After
503      attempting to remove samples we deemed unreliable (age < 17y, height < 1.4m, height > 2.2m,
504      weight > 200kg), we still found a significant mean difference in BMI between diet groups via a
505      Mann Whitney U test (p<0.001). Despite this difference in BMI, we removed no additional
506      samples since we were concerned about further worsening the imbalance between diet groups.

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

513      **Diet was associated with specific taxonomic and functional profiles.** We will report the results
514      from a 100 topic STM, fit with the binary diet information as a covariate. As before, we
515      identified our topics-of-interest by regressing the samples-over-topics distribution against diet
516      and further validated these results via permutation tests, resulting in 9 high ranking topics, 5 of
517      which were associated with the O group, and 4 with the V group (Figure S9).

518      Across the 9 topics, members of the family Lachnospiraceae were well represented, which is not
519      surprising given that it typically accounts for over half of bacteria in healthy human fecal

14

520  samples [36]. Within the topics, we identified roughly 11 clusters of interest that contained high
521  frequency taxa.

522  The topic most associated with the V group T61, was dominated by taxa belonging to
523  Lachnospiraceae, but was also enriched for *Roseburia*, *Blautia*, and *Ruminococcaceae* (Figure 5A).
524  T61's association with the V diet is consistent with literature linking *Roseburia* and
525  *Ruminococcaceae* to starch and plant polysaccharide metabolism [37] and *Roseburia* and *Blautia* to
526  whole grains [38, 39]. Also, consistent with this topic being dominated by Gram positive
527  bacteria, we identified a significant depletion in LPS biosynthesis genes. (Figure 5B).

528  T12 contained a small yet diverse cluster of bacteria within *Acinetobacter*, a genus often
529  associated with fermented foods and beverages [40]. Quinn *et al.* (2016), investigating the effect
530  home-fermented foods had on human microbiota, identified enrichment of fluorobenzoate
531  degradation pathways [41]. In our results, we found that the fluorobenzoate degradation
532  pathway for T12 had the largest shift of any predicted pathway within a given topic (Figure 6).
533  To further investigate the relationship between fluorobenzoate degradation pathways and diet,
534  we performed a logistic regression (logit link) on all samples with subject ages at least 21y. Diet
535  type and the z-scored frequency of T61 were independent variables with alcohol consumption
536  ($n_{no}$=837, $n_{yes}$=3692) as the binary outcome. Both T61 (Z=3.64, $\beta_{T61}$=1.10, p<0.001) and diet (Z=-
537  6.78, $\beta_{diet}$=-0.89, p<0.001) were significant, suggesting a potential relationship with fermented
538  foods (specifically alcohol), *Acinetobacter*, and fluorobenzoate degradation.

539  Finally, T76 contained bacteria typically associated with a western lifestyle such as *Clostridiales*
540  [42]. It was also enriched for *Faecalibacterium prausnitzii*, as well as butyrate production. This is
541  significant because butyrate is not only critical in the fermentation of plant matter [43], but
542  reduction of fecal butyrate has been implicated in obesity and a shift toward a less
543  carbohydrate-rich diet [44]. The remaining bacteria present in this T76 cluster, *Ruminiococus* and
544  *Roseburia*, have been shown to be elevated after fiber consumption [38].

545  The topics associated with the O group were enriched for LPS and secretion system pathways.
546  A noteworthy cluster in T77 was surprisingly quite similar to the cluster in T61.
547  Lachnospiraceae composed the majority of each cluster: 47.8% (11/23) of taxa for T61 compared
548  to 20.6% (13/63) for T61. The functional profiles were analogous for all pathways except
549  carotenoid biosynthesis and porphyrin and chlorophyll metabolism. A notable distinguishing
550  characteristic is the lack of any *Roseburia* in the T77 cluster in contradistinction to T61.

551  Lastly, we assessed how specific parts of certain metabolic pathways were enriched or depleted
552  within a given topic. We found that T77 was enriched for a specific gene that codes for the
553  enzyme in the final step of the synthesis of bacterial antioxidant staphyloxanthin, a known
554  *Staphylococcus aureus* virulence factor [45] (Figure 5D). T20 was abundant in LPS biosynthesis
555  gene, but depleted in a subset of LPS genes key in a specific branch of the LPS pathway (Figure
556  5E).

557

558  **DISCUSSION**

559

560    We have proposed an approach for uncovering latent thematic structure in 16S rRNA survey
561    data that simultaneously explores taxonomic and predicted functional content. Rather than
562    inferring functional content independently from taxonomic abundances, our approach shifts the
563    focus to investigating within-topic functional content. Unlike other methods, by exploring our
564    topics, we can link categories of functional content to specific clusters of taxa which can in turn
565    be linked to sample features of interest. For example, like Gevers *et al.*, we detected a
566    relationship between membrane transport genes and CD+, but our approach also allowed us to
567    determine which bacteria (OTUs belonging to Enterobacteriaceae) were the prime contributors
568    to the enrichment of membrane transport genes. Moreover, the pathogenic set of bacteria
569    reported by Gevers *et al.* (*Haemophilus* spp., *Neisseria*, and *Fusobacteria*) contributed less to the
570    abundance of membrane transport genes. By applying statistical approaches on the dataset as a
571    whole, as is typical, the apparent relationship between membrane transport genes and specific
572    clusters of bacteria would be lost.

573    We have also shown that our approach drastically reduces the dimensionality of two high-
574    dimensional sources of information, taxonomic abundances and functional content, increasing
575    the ease in which these data can be interpreted.  For instance, we can identify gene sets of
576    interest from noteworthy topics. For the Gevers *et al.* dataset, we determined that T15 is (1)
577    associated with CD+ samples; (2) dominated by a cluster of bacteria previously associated with
578    CD; and (3) uniquely enriched for a subset of LPS synthesis genes. With a gene profile from a
579    topic of interest, one could focus on gene subsets associated with topic-specific bacterial clusters
580    that are known disease biomarkers, which in turn may facilitate targeted approaches for future
581    research endeavors.

582    Lastly, our complete pipeline is computationally manageable. Fitting the STM to nearly 5000
583    samples from the AG dataset reaches convergence in minutes. Functional prediction via
584    PICRUSt also only takes minutes (using our C++ implementation in themetagenomics).
585    Inferring topic-function interaction effects via our multilevel, negative binomial regression
586    approach is comparatively slower, however, taking hours for large datasets (*e.g.,* AG). This is
587    because we implement this model in the probabilistic programming language Stan, which uses
588    Hamiltonian Monte Carlo. Maximum likelihood (a much faster alternative) generally fails to
589    converge for these data, although the regression weight estimates tend to be quite similar based
590    on our experience.

591    We present our approach at a time when novel means to analyze complex microbiome
592    abundance data is called for. Current methods often link the abundance of a single OTU across
593    samples to a sample feature of interest. These methods routinely identify important subsets of
594    taxa, but ignore OTU co-occurrence. Network methods overcome this concern, but instead fail
595    to do so while including sample information within the model. Consequently, they are
596    incapable of directly linking sections of the OTU correlation network with sample features of
597    interest. Constrained ordination methods, such as canonical correspondence analysis, do in fact
598    couple inter-community distance with sample information, but the user is limited to specific
599    distance metrics (*e.g.,* Chi-squared) and must follow key assumptions (*e.g.,* the distributions of

600  taxa along environmental gradients are unimodal) [46]. Moreover, interpretation of biplots
601  becomes increasingly difficult as more covariates are included, and, unlike our approach,
602  linking key groups of taxa with functional content is not straightforward.

603  The ability to make meaningful inferences is further limited by the fact that microbiome data is
604  often inadequately sampled (thus justifying some type of normalization procedure),
605  compositional (due to normalization), sparse, and overdispersed. Compositional data restricts
606  the appropriateness of many statistical methods due to the sum constraint placed across
607  samples. SPIEC-EASI provides a robust network approach for overcoming compositional
608  artifacts in an attempt to infer community level interactions. We compared our within-topic
609  taxonomic profiles to the first and second order connections identified by SPIEC-EASI, to which
610  we found coherence between the two approaches, suggesting a topic model approach for
611  compositional data is in fact appropriate.

612  Others have explored the use of Dirichlet-Multinomial models, which are well equipped at
613  managing overdispersed count data [47–49]. The fact that Dirichlet-Multinomial conjugacy is
614  exploited for the topics-over-OTUs component of many topics models reflects their suitability
615  for abundance data. We selected the recently developed STM for our workflow because of its
616  ability to not only utilize sample data prior information in the flavor of the Dirichlet-
617  Multinomial regression topic model, but also its ability to capture topic correlation structure
618  and apply partial pooling over samples or regularization across regression weights.

619  Because normalization is also a chief concern when analyzing sequencing abundance data [27],
620  we found it imperative to determine a suitable approach. In the original LDA paper, the
621  generative process assumed a fixed document length $N$, but $N$ was considered a simplification
622  and could easily be removed because it is independent of all other components of the model.
623  This allows for the possibility of more realistic document size distributions [15]. Giving this fact,
624  in a companion paper, we performed simulations to determine a suitable normalization strategy
625  (Woloszynek *et al.*, *in prep*). Our simulations suggested that predefined subcommunities
626  concentrate with high probability to extracted topics and that no library size normalization was
627  required to maximize power or the ability to infer taxonomic structure. Thus, a topic model
628  approach provides a direct, suitable procedure for inferring the subcommunity configuration.
629  The variance stabilization through DESeq2, while potentially ideal for large sample sizes with
630  adequate signal, seemed to dampen the ability to identify topic-sample associations. Despite
631  performing well at mapping subcommunities to topics, the rarefied approach suffered from
632  reduced power when identifying topics with large covariate effects. Taken together, we
633  concluded that raw abundance data could be adequately modeled in our approach.

634  There are limitations to our approach. First, the topic-function inference step currently scales
635  poorly in terms of computation time for large numbers of topics, which may be more important
636  as datasets continue to grow in size. Regularization and sparsity-inducing priors help limit the
637  number of important topics; hence, exploring only a subset of topics during the final regression
638  step can offer substantial speed improvements at little cost, but utilizing the complete set of
639  topic information would be ideal. Also, we used Hamiltonian Monte Carlo via Stan. As of now,
640  Stan lacks within-chain parallelization, although that may change in upcoming releases (see

17

641 Stan MPI prototype). Alternatively, other posterior inference procedures such as variational
642 inference using software packages such as Edward may provide additional speed
643 enhancements [50]. Second, we are capable of separately estimating the uncertainty in our topic
644 model, the hierarchical regression model, and the functional predictions from PICRUSt, but we
645 currently do not propagate the uncertainty throughout the pipeline. Doing so would improve
646 downstream interpretation with better estimation of the topic-sample covariates and pathway-
647 topic effects, which in turn would greatly improve one's confidence with using within-topic
648 gene sets. Third, we do not incorporate phylogenetic branch length information, which could
649 lead to more meaningful topics.

650 **LIST OF ABBREVIATIONS**

651

652 AG, American gut

653 BMI, body mass index

654 CD, Crohn's disease

655 IBD, inflammatory bowel disease

656 LDA, latent Dirichlet allocation

657 LFC, log-fold change

658 LN, logistic Normal

659 LPS, lipopolysaccharide

660 OTU, operational taxonomic unit

661 PCDAI, pediatric Crohn's disease activity index

662 PPD, posterior predictive distribution

663 STM, structural topic model

664

665

666

667

668

669

670

671

18

672

673

674

675

676

677

678 **AVAILABILITY OF DATA AND MATERIAL**

679 The datasets generated and/or analyzed during the current study are available in the tm_pipline
680 github repository, https://github.com/sw1/tm_pipeline

681

682 **FUNDING**

684

685 **AUTHORS' CONTRIBUTIONS**

686 SW chose and preprocessed the datasets; developed the pipeline, modeling approach, and
687 statistical tests; designed the topic-function NBR model; and fit all models and performed all
688 tests. JM and SW conceived of the topic-function prediction approach. SW and GS designed
689 simulations for normalization and determined the approach to assess posterior uncertainty. SW,
690 JM, GR, and MO assisted in interpretation and identifying ways to further validate the pipeline.
691 GR supervised the project. All authors contributed to the writing of the manuscript. All authors
692 read and approved the final manuscript.

693

694 **ACKNOWLEDGEMENTS**

699

700

701

702

703

704

705

706

707

708

709

710     **REFERENCES**

711

712     1. Knights D, Costello E, Knight R. Supervised classification of human microbiota. FEMS
713     Microbiol Rev. 2011;35:343–59. doi:10.1111/j.1574-6976.2010.00251.x.

714     2. Li H. Microbiome, Metagenomics, and High-Dimensional Compositional Data Analysis.
715     Annu Rev Stat Its Appl. 2015;2:73–94. doi:10.1146/annurev-statistics-010814-020351.

716     3. Gilbert JA, Quinn RA, Debelius J, Xu ZZ, Morton J, Garg N, et al. Microbiome-wide
717     association studies link dynamic microbial consortia to disease. Nature. 2016;535:94–103.
718     doi:10.1038/nature18850.

719     4. Shafiei M, Dunn KA, Boon E, MacDonald SM, Walsh DA, Gu H, et al. BioMiCo: a supervised
720     Bayesian model for inference of microbial community structure. Microbiome. 2015;3:8.
721     doi:10.1186/s40168-015-0073-x.

722     5. Jiang X, Dushoff J, Chen X, Hu X. Identifying enterotype in human microbiome by
723     decomposing probabilistic topics into components. 2012 IEEE Int Conf Bioinforma Biomed.
724     2012;:1–4. doi:10.1109/BIBM.2012.6392720.

725     6. Ning J, Beiko RG. Phylogenetic approaches to microbial community classification.
726     Microbiome. 2015;3:47. doi:10.1186/s40168-015-0114-5.

727     7. Ren B, Bacallado S, Favaro S, Holmes S, Trippa L. Bayesian Nonparametric Ordination for the
728     Analysis of Microbial Communities. arXiv Prepr arXiv160105156. 2016;:1–25.
729     http://arxiv.org/abs/1601.05156.

730     8. Langille MGI, Zaneveld J, Caporaso JG, McDonald D, Knights D, Reyes J a, et al. Predictive
731     functional profiling of microbial communities using 16S rRNA marker gene sequences. Nat
732     Biotechnol. 2013;31:814–21. doi:10.1038/nbt.2676.

733     9. Aßhauer KP, Wemheuer B, Daniel R, Meinicke P. Tax4Fun: predicting functional profiles
734     from metagenomic 16S rRNA data. Bioinformatics. 2015;31:2882–4.
735     doi:10.1093/bioinformatics/btv287.

736     10. Iwai S, Weinmaier T, Schmidt BL, Albertson DG, Poloso NJ, Dabbagh K, et al. Piphillin:
737     Improved prediction of metagenomic content by direct inference from human microbiomes.

738   PLoS One. 2016;11:1–18.

739   11. Knights D, Kuczynski J, Charlson ES, Zaneveld J, Mozer MC, Collman RG, et al. Bayesian
740   community-wide culture-independent microbial source tracking. Nat Methods. 2013;8:761–3.

741   12. Roberts ME, Stewart BM, Tingley D, Lucas C, Leder-Luis J, Gadarian SK, et al. Structural
742   topic models for open-ended survey responses. Am J Pol Sci. 2014;58:1064–82.

743   13. Blei DM, Lafferty JD. A correlated topic model of Science. Ann Appl Stat. 2007;1:17–35.
744   doi:10.1214/07-AOAS136.

745   14. Mimno D, McCallum A. Topic models conditioned on arbitrary features with dirichlet-
746   multinomial regression. arXiv Prepr arXiv12063278. 2012. doi:10.1.1.140.6925.

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

748   16. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The
749   Treatment-Naive Microbiome in New-Onset Crohn's Disease. Cell Host Microbe. 2014;15:382–
750   92. doi:10.1016/j.chom.2014.02.005.

751   17. Eisenstein J, Ahmed A, Xing EPE. Sparse additive generative models of text. Proc 28th Int
752   Conf Mach Learn. 2011;:1041–8. doi:10.1.1.206.5167.

753   18. Roberts ME, Stewart BM. A model of text for experimentation in the social sciences. Work
754   Pap. 2015.

755   19. Kembel SW, Wu M, Eisen JA, Green JL. Incorporating 16S Gene Copy Number Information
756   Improves Estimates of Microbial Diversity and Abundance. PLoS Comput Biol. 2012;8:16–8.

757   20. Callahan BJ, Mcmurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP. DADA2 : High
758   resolution sample inference from amplicon data. bioRxiv. 2015;13:0–14. doi:10.1101/024034.

759   21. Roberts, Margaret E., Stewart BM, Tingley D. stm: R Package for Structural Topic Models.
760   2017. http://www.structuraltopicmodel.com.

761   22. Blei DM, McAuliffe JD, Blei DM. Supervised Topic Models. Adv Neural Inf Process Syst 20.
762   2008;21:1–22. doi:10.1002/asmb.540.

763   23. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. KEGG for integration and
764   interpretation of large-scale molecular data sets. Nucleic Acids Res. 2012;40.

765   24. Stan Development Team. rstanarm: Bayesian applied regression modeling via Stan. 2016.
766   http://mc-stan.org/.

767   25. Gelman A, Rubin DB. Inference from Iterative Simulation Using Multiple Sequences. Stat
768   Sci. 1992;7:457–511.

769   26. Kurtz Z, Mueller C, Miraldi E, Bonneau R. SpiecEasi: Sparse InversE Covariance estimation
770   for Ecological Association and Statistical Inference. 2016. https://github.com/zdk123/SpiecEasi.

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

772     Inadmissible. PLoS Comput Biol. 2014;10.

773     28. Love MI, Anders S, Huber W. Differential analysis of count data - the DESeq2 package. 2014.
774     doi:110.1186/s13059-014-0550-8.

775     29. Kuhn M. Building Predictive Models in R Using the caret Package. J Stat Softw. 2008;28:1–26.
776     doi:10.1053/j.sodo.2009.03.002.

777     30. McMurdie PJ, Holmes S. Phyloseq: An R Package for Reproducible Interactive Analysis and
778     Graphics of Microbiome Census Data. PLoS One. 2013;8.

779     31. Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'Hara RB, et al. vegan:
780     Community Ecology Package. R package version 2.3-1. 2015;:264.
781     doi:10.4135/9781412971874.n145.

782     32. Gevers D, Kugathasan S, Denson LA, Vázquez-Baeza Y, Van Treuren W, Ren B, et al. The
783     treatment-naive microbiome in new-onset Crohn's disease. Cell Host Microbe. 2014;15:382–92.
784     doi:10.1016/j.chom.2014.02.005.

785     33. Tjalsma H, Boleij A, Marchesi JR, Dutilh BE. A bacterial driver-passenger model for
786     colorectal cancer: beyond the usual suspects. Nat Rev Microbiol. 2012;10:575–82.
787     doi:10.1038/nrmicro2819.

788     34. Lewis JD, Chen EZ, Baldassano RN, Otley AR, Griffiths AM, Lee D, et al. Inflammation,
789     Antibiotics, and Diet as Environmental Stressors of the Gut Microbiome in Pediatric Crohn's
790     Disease. Cell Host Microbe. 2015;18:489–500. doi:10.1016/j.chom.2015.09.008.

791     35. Gelman A, Hill J. Data analysis using regression and multilevel/hierarchical models. New
792     York, NY: Cambridge University Press; 2006.

793     36. Flint HJ. The impact of nutrition on the human microbiome. Nutr Rev. 2012;70:S10–3.
794     doi:10.1111/j.1753-4887.2012.00499.x.

795     37. Flint HJ, Scott KP, Duncan SH, Louis P, Forano E. Microbial degradation of complex
796     carbohydrates in the gut. Gut microbes. 2012;3:289–306.

797     38. Flint HJ, Graf D, Cagno R Di, Fa F, Nyman M, Saarela M, et al. Contribution of diet to the
798     composition of the human gut microbiota. Microb Ecol. 2015;1:1–11.

799     39. Martínez I, Lattimer JM, Hubach KL, Case JA, Yang J, Weber CG, et al. Gut microbiome
800     composition is linked to whole grain-induced immunological improvements. ISME J.
801     2013;7:269–80. doi:10.1038/ismej.2012.104.

802     40. Tamang JP, Watanabe K, Holzapfel WH. Review: Diversity of microorganisms in global
803     fermented foods and beverages. Front Microbiol. 2016;7 MAR.

804     41. Quinn RA, Navas-Molina JA, Hyde ER, Song SJ, Vázquez-Baeza Y, Humphrey G, et al. From
805     Sample to Multi-Omics Conclusions in under 48 Hours. mSystems. 2016;1:e00038-16.
806     doi:10.1128/mSystems.00038-16.

807    42. Gorvitovskaia A, Holmes SP, Huse SM. Interpreting Prevotella and Bacteroides as
808    biomarkers of diet and lifestyle. Microbiome. 2016;4:15. doi:10.1186/s40168-016-0160-7.

809    43. Gill SR, Pop M, Deboy RT, Eckburg PB, Turnbaugh PJ, Samuel BS, et al. Metagenomic
810    analysis of the human distal gut microbiome. Science. 2006;312:1355–9.
811    doi:10.1126/science.1124234.

812    44. Duncan SH, Belenguer a., Holtrop G, Johnstone a. M, Flint HJ, Lobley GE. Reduced Dietary
813    Intake of Carbohydrates by Obese Subjects Results in Decreased Concentrations of Butyrate and
814    Butyrate-Producing Bacteria in Feces. Appl Environ Microbiol. 2007;73:1073–8.
815    doi:10.1128/AEM.02340-06.

816    45. Clauditz A, Resch A, Wieland KP, Peschel A, Götz F. Staphyloxanthin plays a role in the
817    fitness of Staphylococcus aureus and its ability to cope with oxidative stress. Infect Immun.
818    2006;74:4950–3.

819    46. Legendre P, Legendre L. Numerical Ecology - Second English Edition. 1998.

820    47. De Valpine P, Harmon-Threatt AN. General models for resource use or other compositional
821    count data using the Dirichlet-multinomial distribution. Ecology. 2013;94:2678–87.

822    48. Holmes I, Harris K, Quince C. Dirichlet multinomial mixtures: Generative models for
823    microbial metagenomics. PLoS One. 2012;7.

824    49. Brien JDO, Record N. The power and pitfalls of Dirichlet-multinomial mixture models for
825    ecological count data. bioRxiv. 2016;:1–22.

826    50. Brevdo E, Hoffman MD, Murphy K, Blei DM, Tran D, Saurous RA, et al. Deep Probabilistic
827    Programming. 2017; Nips:1–17.

828
829
830
831
832
833
834
835
836
837
838
839

840

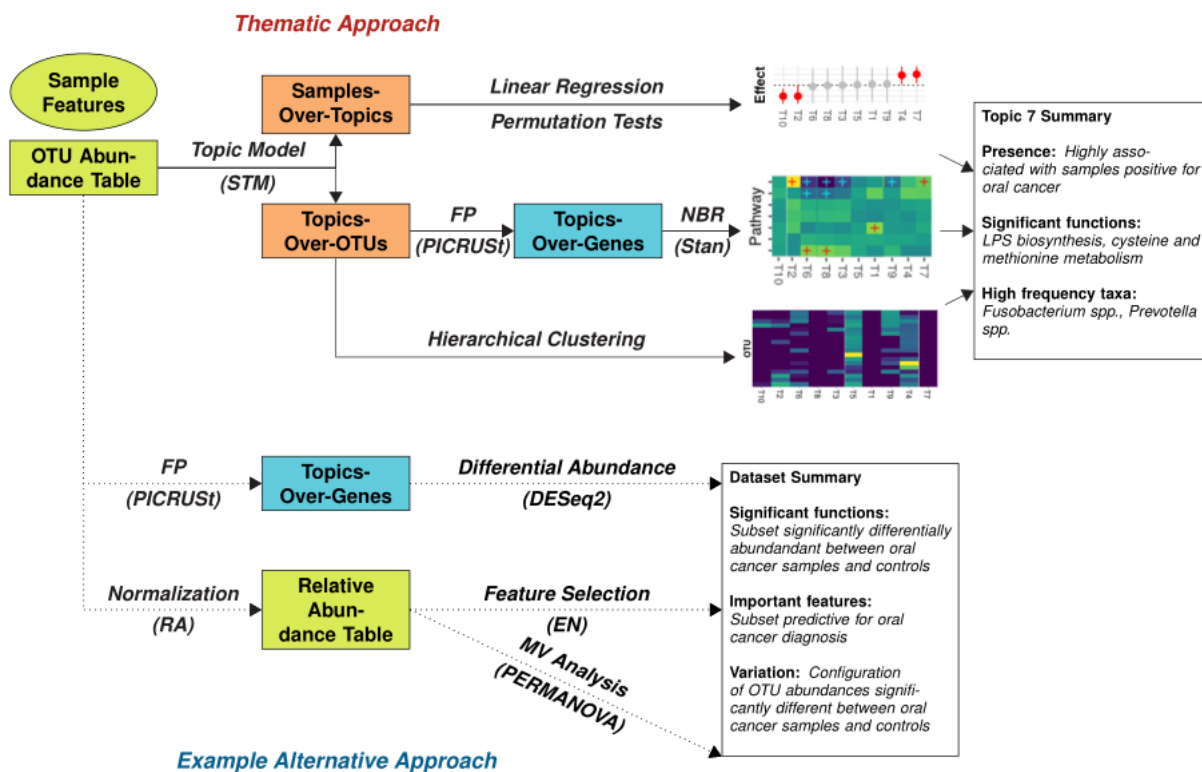841

842

843

844

845

846 **FIGURES**



847

848

849 Figure 1. (Thematic Approach) Given a 16S rRNA abundance table, a topic model is used to
850 uncover the thematic structure of the data in the form of two latent distributions: the samples-
851 over-topics frequencies and the topics-over-OTUs frequencies. The samples-over-topics
852 frequencies are regressed against sample features of interest to identify the strength of a topic-
853 covariate relationship to rank topics (top). The topics-over-OTUs frequencies are used in a gene
854 function prediction (FP) algorithm to predict gene content. Important functional categories are
855 identified via a fully Bayesian multilevel negative binomial (NBR) regression model (middle).
856 The topics-over-OTUs distribution is also hierarchically clustered to infer relationships between
857 clusters of co-occurring OTUs and topics (bottom). The end result is the ability to identify key
858 topics that associate clusters of bacteria, sample features of interest, and functional content.

859 (Alternative Approach). An example alternative approach involves independently
860 characterizing the taxonomic configuration and the predicted functional configuration of the
861 OTU abundance table. Gene function prediction is performed on the full OTU abundance table,
862 followed by a differential abundance analysis to infer differences in specific genes between
863 sample features of interest (top). The OTU table is normalized to overcome library size
864 inconsistencies and then analyzed via two methods: (1) an elastic net (EN) to find sparse sets of
865 OTUs that are predictive for the sample feature of interest (middle) and (2) a multivariate (MV)
866 analysis to identify relationships between beta diversity and the sample feature of interest
867 (bottom). The end result are three analyses that summarize the data as a whole, unlike the
868 thematic approach, which characterizes co-occurring sets of OTUs in three ways.
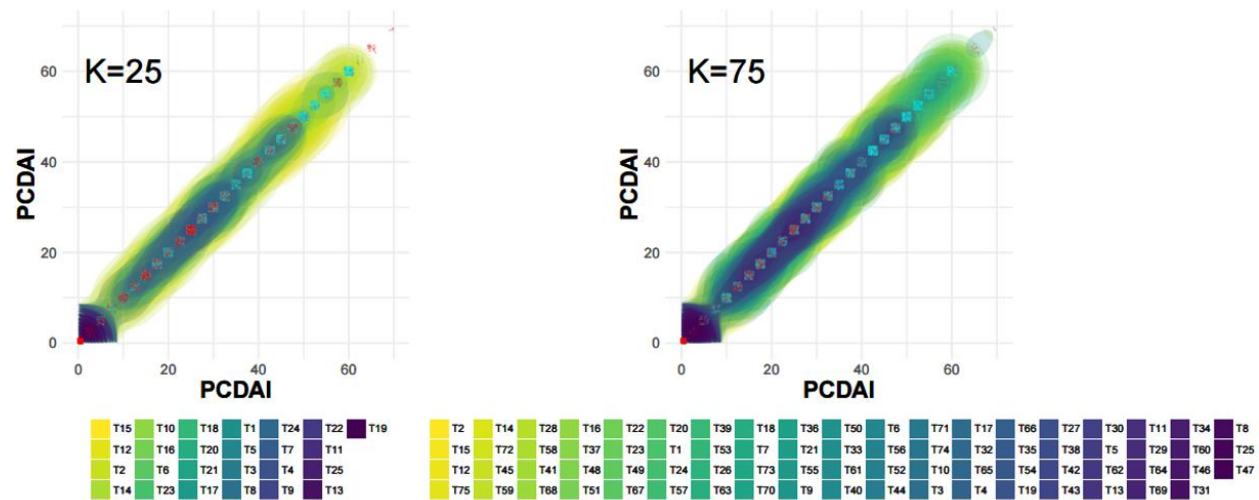


869

870 Figure 2. The posterior predictive distribution (PPD) of topic assignments for each sample in
871 Gevers with reported PCDAI. Shown are the PPDs for K25 and K75. CD- samples were set to
872 PCDAI=0. Topics are ordered based on their topic-effect, the regression weights estimated by
873 regressing the samples-over-topics distribution against CD presence. Thus, left-most topics are
874 most associated with CD- samples, whereas right-most topics are most associated with CD+
875 samples. The cyan (CD-) and red (CD+) points indicate if a topic, drawn from the PPD for a
876 given PCDAI value, were high ranking topics (topics that did not span 0 at 95% uncertainty for
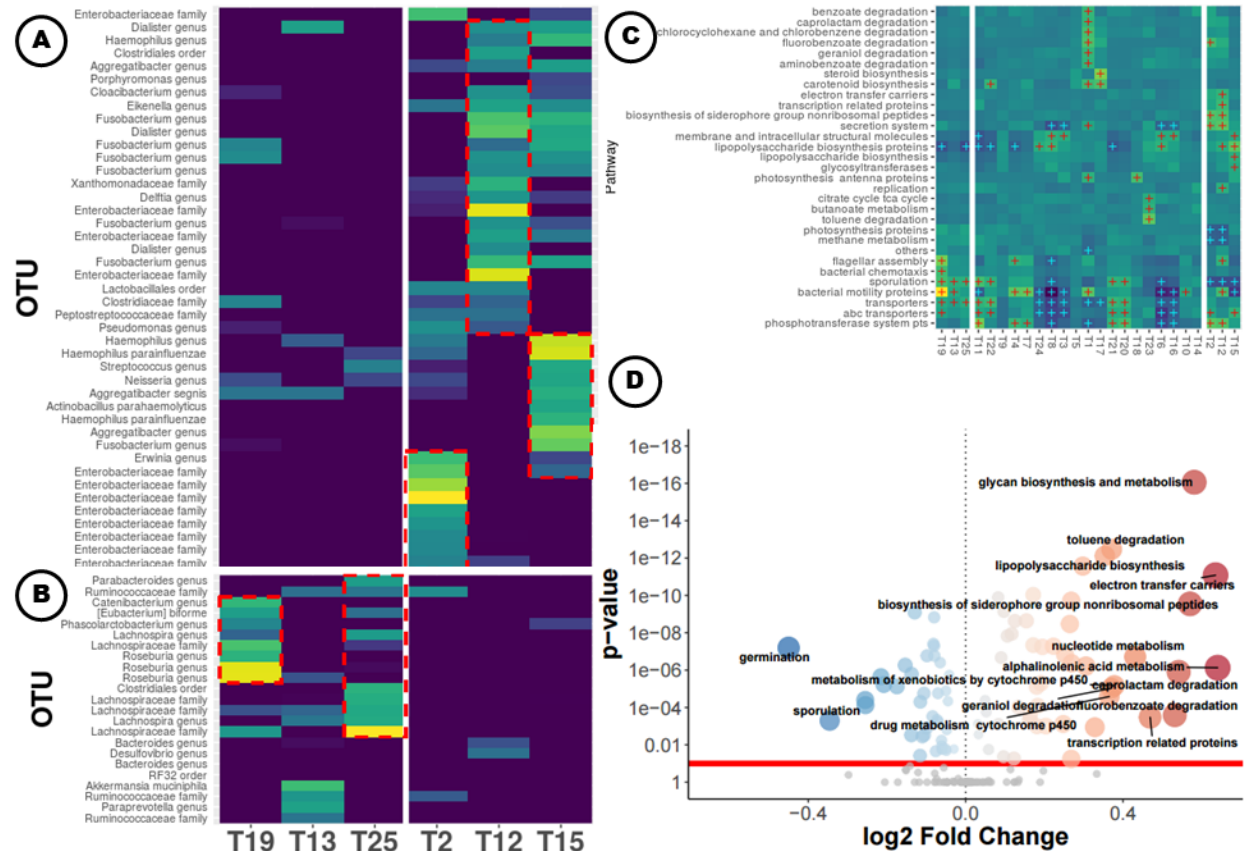877 permutation tests).

878

879    Figure 3A,B. Subsections of the heatmap for the Gevers of the topics-over-OTUs distribution

880    (K25) in log space. Shown are the top 3 topics associated with CD- and CD+, ordered by topic-

881    effects (left to right, respectively, separated by the white line). Clusters of interest are marked

882    with red dotted lines. Clustering was performed via Ward's method on Bray-Curtis distances.

883    Low probabilities ($p < 1 \times 10^{-5}$) are set to 0 to minimize the range of the color gradient to ease

884    visualization. Yellow=high probability, Blue=low probability. Figure 3C. Level-3 pathway

885    category-topic interaction regression coefficients from the multiple level negative binomial

886    model. KEGG information was predicted via PICRUSt on the topics-over-OTUs distribution.

887    Clustering was performed via Ward's method on Bray-Curtis distances. Red and blue crosses

888    indicate estimated pathway-topic interaction weights that do not span 0 at 80% uncertainty and

889    are positive or negative, respectively. Only pathways with at least one such combination are

890    shown. Yellow=large positive weight estimate, Blue=large negative weight estimate. Figure 3D.

891    Volcano plot showing DESeq2 results for differentially abundant predicted level-3 KEGG

892    categories. Functions were predicted using PICRUSt on the copy number normalized OTU

893    abundance table. Blue and red points represent categories significantly enriched for CD- and

894    CD+, respectively. Gray points are categories with p-values greater than 0.1 after Bonferroni
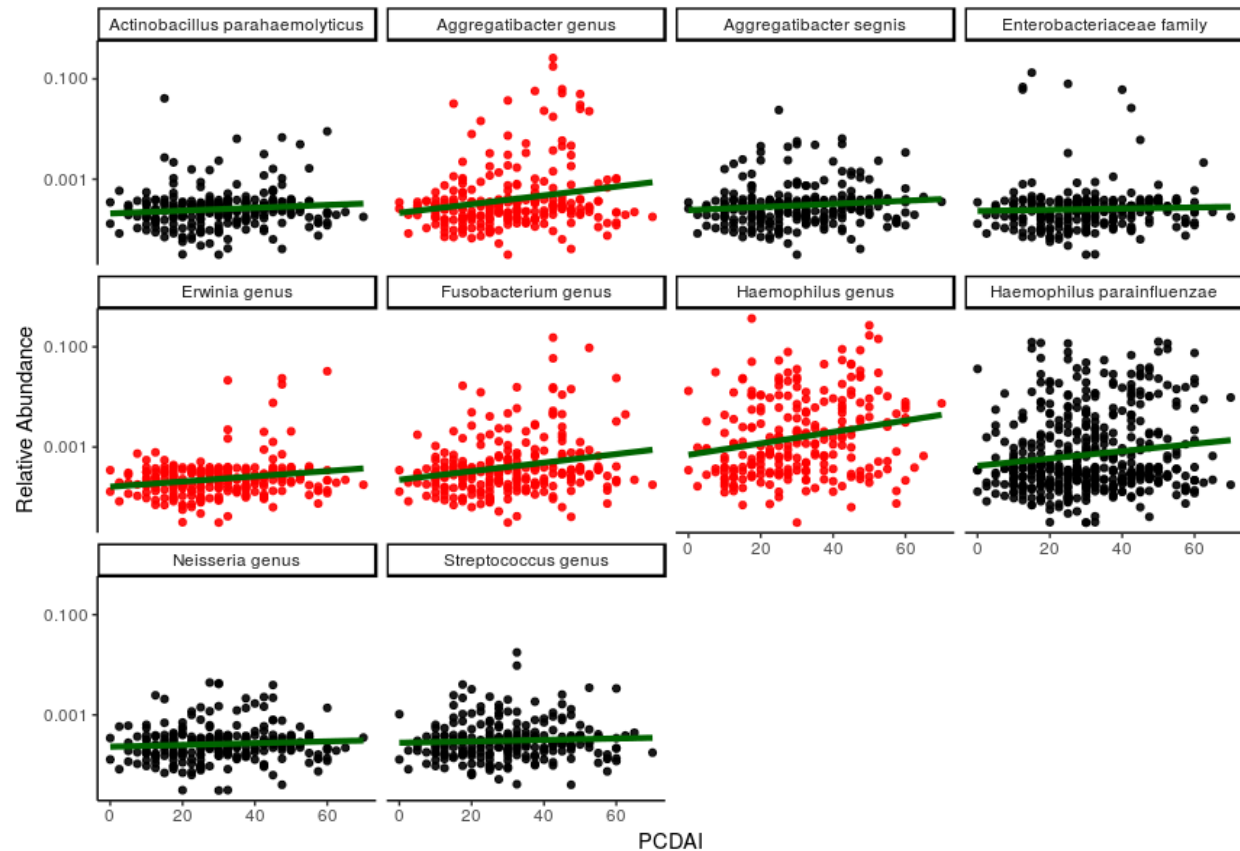
895    correction.

896

897    Figure 4. Scatterplots of Gevers data for the relative abundance of taxa that compose a high
898    probability cluster in T15 versus PCDAI, a clinical measure of CD disease burden. Red points
899    reflect significance (alpha=.05) for negative binomial regression (log linked, sample coverage
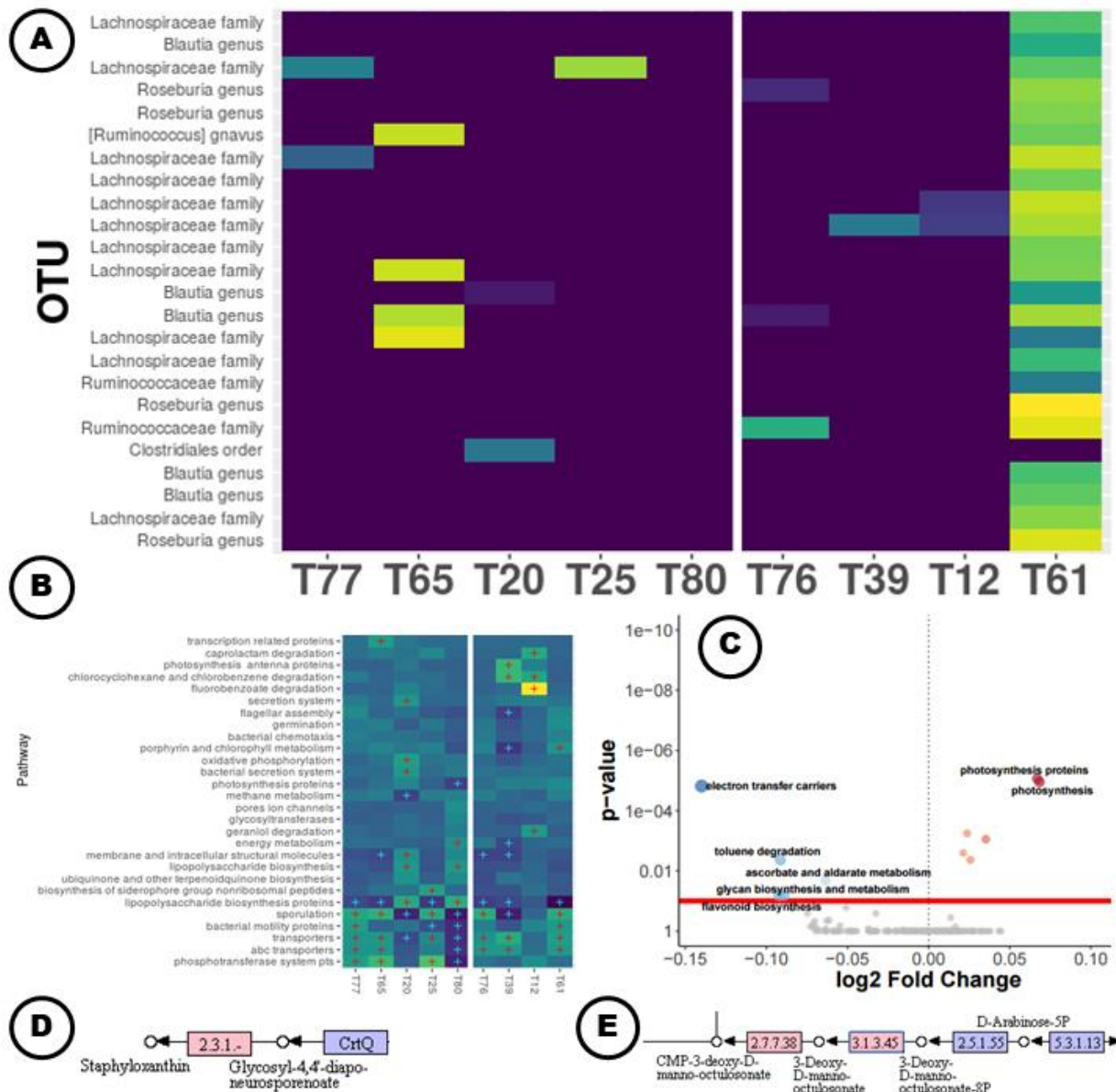900    offset) with Bonferroni correction.

Figure 5A. Subsection of the heatmap for AG for the topics-over-OTUs distribution (K100) in log space. Shown are the topics with 95% uncertainty intervals that do not span 0 when regressed against diet type, ordered by increasing mean regression estimate (left to right). T77 is most associated with O. T61 is most associated with V. White line signifies a shift from positive to negative mean regression estimates. Clustering was performed via Ward's method on Bray-Curtis distances. Low probabilities ($p < 1 \times 10^{-5}$) are set to 0 to ease visualization. Yellow=high probability, Blue=low probability. Figure 5B. Level-3 topic-function interaction weight estimates from the multiple level negative binomial model. KEGG information was predicted via PICRUSt on the topics-over-OTUs distribution. Only the top 25 topics based on mean regression weight were chosen for the negative binomial to alleviate computational concerns. Clustering was performed via Ward's method on Bray-Curtis distances. Red and blue crosses indicate weights for pathway-topic combinations that do not span 0 at 80% uncertainty and are positive

914    or negative, respectively. Only pathways with at least one such combination are shown.

915    Yellow=large positive weight estimate, Blue=large negative weight estimate. Figure 5C. Volcano

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

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

918    table. Blue and red points represent categories significantly enriched for O and V, respectively.

919    Gray points are categories with p-values greater than 0.1 after Bonferroni correction. Figure 5D.

920    Glycosyl-4,4'-diaponeurosporenoate acyltransferase step (red) in carotenoid biosynthesis

921    pathway. This gene is enriched in T77 relative to T20. Figure 5E. Lipopolysaccharide

922    biosynthesis pathway where genes with abundances greater than 50 for T20 are colored red.

923

924    **TABLES**

925

926    Table 1. Relationship of Terms

| Topic Model | Pipeline | Description |
|---|---|---|
| Document | Sample | Collection of reads from subject $m$ on time $t$ |
| Topic | Topic | Collection of co-occurring taxa, subcommunity |
| Word | OTU, Gene, Taxa | Features from taxonomic abundance table or predicted functional content |
| Covariate | Sample feature, Sample class | Sample-level variable of interest – e.g., disease presence, diet, rainfall, time |

927