1 **Decoding the Language of Microbiomes using Word-Embedding Techniques, and Applications in**

2 **Inflammatory Bowel Disease**

3 Christine A. Tataru[1], Maude M. David[1,2]

4

5 **Short title:** Dimensionality Reduction in Microbiome 16S Analysis

6

7 [1]Department of Microbiology, Oregon State University, Corvallis, OR, 97331

8 [2]Department of Pharmaceutical Sciences, Oregon State University, Corvallis, OR, 97331

9

10 *Corresponding Authors:

11     Christine A. Tataru: tataruc@oregonstate.edu

12     Maude David: maude.david@oregonstate.edu

13

14

15

16

17

18

19

20

21

22

# Abstract

Microbiomes are complex ecological systems that play crucial roles in understanding natural phenomena from human disease to climate change. Especially in human gut microbiome studies, where collecting clinical samples can be arduous, the number of taxa considered in any one study often exceeds the number of samples ten to one hundred-fold. This discrepancy decreases the power of studies to identify meaningful differences between samples, increases the likelihood of false positive results, and subsequently limits reproducibility. Currently, most microbiome survey studies focus on differential abundance testing per taxa in pursuit of specific biomarkers for a given phenotype. This methodology assumes differences in individual species, genera, or families can be used to distinguish between microbial communities and ignores community-level action. In this paper, we propose to shift the analysis paradigm from a focus on taxonomic counts to a focus on comprehensive properties that more completely characterize microbial community members' function and environmental relationships. We learn these properties by applying an embedding algorithm to quantify taxa co-occurrence patterns in over 18,000 samples from the American Gut Project (AGP) microbiome crowdsourcing effort. The resulting set of embeddings transforms human gut microbiome data from thousands of taxa counts to a latent variable landscape of only one hundred "properties", or contextual relationships. We then compare the predictive power of models trained using properties, normalized taxonomic count data, and another commonly used dimensionality reduction method, Principal Component Analysis in categorizing samples from individuals with inflammatory bowel disease (IBD) and healthy controls. We show that predictive models trained using property data are the most accurate, robust, and generalizable, and that property-based models can be trained on one dataset and deployed on another with positive results. Furthermore, we find that these properties can be interpreted in the context of current knowledge; properties correlate significantly with known metabolic pathways, and distances between taxa in "property space" roughly correlate with their phylogenetic distances. Using these properties, we are able to extract known and new bacterial metabolic pathways associated with inflammatory bowel disease across two completely independent studies.

More broadly, this paper explores a reframing of the microbiome analysis mindset, from taxonomic counts to comprehensive community-level properties. By providing a set of pre-trained embeddings, we allow

50 any V4 16S amplicon study to leverage and apply the publicly informed properties presented to increase the

51 statistical power, reproducibility, and generalizability of analysis.

52

# 53 1.Introduction

**54 1.1 Microbial survey studies**

55 Microorganisms are biochemically potent entities that influence the biochemistry of surrounding organisms at

56 all ecological scales. Recent findings suggest that resident microbiomes of the human anatomy influence our

57 bodies and minds in ways we have only just begun to understand. Microbiomes have been implicated in the

58 development of diseases of nearly all types, both acute and chronic, infectious and systemic. The vaginal

59 microbiome has been implicated in preterm birth (1), the skin microbiome in acne (2) and eczema (3), and the

60 gut microbiome in a spectrum of diseases including inflammatory bowel disease (IBD) (4–6,6–9), anxiety (10–

61 12), major depressive disorder (13–15), autism (16–20), and Parkinson's Disease (21–23).

62 To analyze microbiome compositions, current technology sequences various hypervariable regions of

63 the 16S rRNA gene, which acts as an accessible taxonomic tag to measure the abundances of taxa in a

64 community. Studies using this 16S survey technique have reported incredibly diverse collections of microbes in

65 several systems. Multiple individuals studies, along with the American Gut Project (AGP) (24)  and the Human

66 Microbiome Project (25),  have invested colossal effort to document that diversity by creating publicly available

67 reference repositories. Amongst these are repositories of stool-associated microbiota that have furthered our

68 understanding of the role of the microbiome in several diseases, especially inflammatory bowel disease (IBD)

69 (4)

70 Though these and other studies have presented highly relevant findings, 16S microbiome survey

71 studies in general tend to suffer from lack of reproducibility (26,27). Difficulties in reproducibility can be

72 attributed to several technological and analysis-based issues (26,28,29) , including two major problems

73 addressed here. First, due to logistical restrictions, especially in human gut microbiome studies where

74 collecting clinical samples can be arduous, the number of taxa considered in any one study often exceeds the

75 number of samples ten to one hundred-fold. Even the largest microbiome studies only include roughly as many

76 samples as taxa analyzed (24,25) . As the number of samples necessary to present a statistically sound and

3

77  reproducible result increases with the number of variables being considered, individual studies with low

78  sample-to-variable ratios risk being underpowered and reporting false positives, especially when effect sizes

79  are estimated to be small (27,30,31).

80      Second, the most commonly employed analysis techniques assume independence of bacterial species

81  (32–34). In biological contexts, the presence and function of each microbe is deeply dependent on the

82  characteristics of its surrounding neighbors. Differences in microbial function also occur as genes are turned

83  on or off as appropriate for that microbe's environment at any given time. For instance, Belenguer et al. show

84  that *Roseburia* strain A2-183 is unable to use lactate as a carbon source except in the presence of *Bacteroides*

85  *adolecentis* (35). Because of functional dependence, findings of differential abundance or function of a single

86  species must be considered within its wider context of associated species and environmental factors (36).

87  More specifically, predictive models that differentiate between disease and healthy guts based on microbiome

88  composition in one dataset can rarely be successfully applied to samples from the same patient population

89  collected independently (27).

90      Navigating the highly related and very large microbiome space can be done with the help of

91  dimensionality reduction methods. The goal of this project is to create an unbiased method to project

92  taxonomic data into a lower dimensional space that represents taxa properties based on their relationships

93  with each other and their environment. In this context, a property is a pattern that underlies co-occurrences

94  between taxa. The lower dimensional space is learned from public datasets using an embedding algorithm,

95  and allows the integration of patterns from massive datasets into specialized studies to increase reproducibility

96  and statistical power.

97

98  **1.2 Current Methods for Dimensionality Reduction**

99  Currently, most microbiome survey studies focus on differential abundance testing per taxa in pursuit of

100  specific biomarkers for a given phenotype. Often, some form of dimensionality reduction is performed to

101  reduce the data to a manageable size. For example, taxa may be filtered to consider only the common or very

102  rare, however this approach may filter potentially valuable data. In another approach, taxa can be categorized,

103  or binned, by their phylogenetic relationships (e.g. all taxa that share a family are analyzed as one unit)

104  (37,38). Such binning methods may obscure meaningful biological signal, and are also heavily database

4

105     dependent not all microbes are clearly classified by taxonomy. Alternatively, taxa can be clustered based on

106     the similarity of their 16S rRNA gene, which has been used as a proxy for evolutionary relatedness (39).

107     However, in this approach, clustering may hinder comparisons across studies, and may result in biologically

108     unfounded taxonomic units (28). Such taxonomic count-based methodologies, while they have led to

109     interesting and crucial discoveries in stool-associated microbiome surveys, assume that differences in

110     individual species, genera, or families can be used to distinguish between microbial communities and ignore

111     community-level action between and among species.

112          Rather than searching for individual biomarkers, ordination may instead be used to reduce data

113     dimensionality and identify broad patterns in microbiome compositions between samples. Samples, each

114     represented by a vector of taxa, can be projected into a lower dimensional space using a wide array of

115     ordination techniques including principal component analysis (PCA) (40) and multidimensional scaling (41).

116     Broadly used, ordination has played a critical role in associating microbial structure with specific features or

117     phenotypes of interest, but has also proven to be overly sensitive to normalization and study bias (e.g.

118     technological noise, DNA preparation protocol, sequencing error) (42).

119          To adapt ordination to a microbiome-specific technique, Fukuyama et. al integrated phylogenetic

120     information via a Bayesian prior to a standard principal coordinate analysis. In a similar attempt to integrate the

121     concept of distance between 16S gene variants, several authors have proposed to represent each 16S

122     sequence by the set of k-length nucleotide sequences (k-mers) it includes. Woloszynek et. al embed those k-

123     mers to create a vector representation of each sequence, and show that representing taxa as the average of

124     their embedded k-mers results in a meaningful representation of taxa that can be beneficial to exploratory

125     analysis or supervised machine learning (43)

126          Finally, Sankaran et. al model taxa as units drawn from an underlying distribution of latent variables

127     (36). Each sample is modeled as originating from a multinomial across some underlying biological "topics", and

128     taxa counts are modeled as Dirichlet multinomial mixtures across all topics. Under this model, a sample is

129     ultimately represented by its k latent topic distribution instead of by its taxa counts. This method aptly describes

130     samples by assigning topic distributions to them, but does not directly relate taxa to each other.

131

132

5

**1.3 Current study proposal**

While compelling, the dimensionality reduction methods described above do not consider taxonomic

relationships *within a biological contex*t, or make use of information already available from previous datasets.

By integrating previous studies and subsequently putting 16S rRNA gene into context, our study proposes to

describe inherent properties of a microbial communities in a manner consistent with their functional utility in

their environmental context.

To deduce the above-mentioned properties, we turn to embedding techniques from natural language

processing. The use of natural language methods in microbiome analysis is not new. As noted by Sankaran et.

al (36), there exist some easily drawn parallels between natural language data and microbiome data, namely

that documents are equivalent to biological samples, words to taxa, and topics to microbial neighborhoods.

Just as a book may be defined by the aggregate topics it discusses, a microbial environment may be defined

the neighborhoods or communities it contains.

There is another connection between words and microbes not currently discussed in the literature, and

that is the capacity of both entities to be described by a finite set of discrete, characteristic properties. For

instance, the word 'apple' in English can be defined as an edible, red, non-gendered, crunchy, object. Similarly,

the species *Clostridium difficile* can be defined as a spore-forming, infectious, spindle-shaped bacteria. While it

would be difficult to distinguish between a recipe book and a magazine of food reviews by enumerating

differences in the occurrence of individual words, differentiating the two becomes simple if we select

appropriate properties. While both media use words that have high scores in the property "edibility", the recipe

book also uses words that have a high declarative score, like 'cut', 'wash', and 'prepare', while the food review

uses words that have high descriptive scores, like 'fantastic', 'delectable', or 'abysmal'. Just as the properties of

"declarative" and "descriptive" allow us to differentiate texts more effectively, property-based analysis of

microbiomes allow us to distinguish between two microbial scenarios more easily than individual taxa counts.

Analysis on the level of properties thus provides a more accurate and generalizable representation of the

data's structure.

In this study, the properties mentioned above were learned from patterns in a large microbial dataset

provided by the American Gut Project (AGP). An unsupervised embedding algorithm developed for natural

language processing called GloVe (44) was applied to over 15,000 AGP samples to learn an embedding space

6

161  by quantifying co-occurrence patterns between taxa. The resulting set of embeddings transforms human gut

162  microbiome data from thousands of taxa counts to a property space of only one hundred to seven hundred

163  variables. We quantify the quality of the properties by predicting the Inflammatory Bowel Disease (IBD) status

164  of samples using properties, normalized taxonomic count data, and principal component analysis. We show

165  that predictive random forest models trained using property data are the most accurate, robust, and

166  generalizable, and that property-based models can be trained on one dataset and deployed on an independent

167  one with positive results. Strong correlation between learned properties and annotated metabolic pathways

168  allow us to implicate both known and new metabolic pathways in IBD such as steroid degradation,

169  lipopolysaccharide biosynthesis, and various types of glycan biosynthesis. Lastly, by projecting taxonomic data

170  into property space, the scientific community can integrate patterns from massive public datasets into specific,

171  targeted studies. Analysis in property space means models requires fewer samples to produce robust results,

172  and exploratory studies simultaneously gain increased power and decreased risk of spurious associations.

173  We not only advocate the use of this method, but also propose to shift the analysis paradigm from a

174  focus on taxonomic counts to a focus on comprehensive properties that more completely characterize

175  microbial community members' function and environmental relationships.  The human gut microbiome has the

176  potential to be used as a low-cost environmental barometer for the diagnosis and monitoring of disease, but

177  first we must prioritize model reproducibility and move beyond the concept of the taxonomic unit.

178

179  **Figure 1:** Workflow of data transformation to prediction of host phenotype. First, taxa-taxa co-occurrence

180  (binary) data from the American Gut Project (A) are input into the GloVe embedding algorithm (B) to produce a

181  taxa (Amplicon Sequence Variant or ASV) by property transformation matrix (C). Then, we take the dot product

182  between a sample by taxa table of interest (D) and the transformation matrix (C) to project that table into

183  embedding space (E). This table is used to train a random forest model (F) along with sample associated

184  lifestyle and dietary information (G) to predict the IBD status of the host (H). As points of comparison, random

185  forest models are also built without embedding, after transforming the same sample by taxa table (D) using

186  PCA (I) and normalizing (J).

187

# 2. Results

## 2.1 Model performance

In order to determine the value of the set of embedding produced by GloVe, we tested the performance of classifiers built using GloVe embedded, PCA transformed, and non-embedded normalized count data. We evaluated two main performance metrics in predicting the IBD status of the host: area under the receiver operating curve (AUROC) and area under the precision-recall curve (AUPR). The receiver operating curve plots true positive calls against false positive calls. The higher the AUROC, the more confident you can be that a positive prediction by the classifier is correct. The precision-recall curve plots the precision, how confident you are that a positive call is correct, against recall, what percentage of the positive samples in the dataset were identified. A high AUPR means the classifier is able to identify most of the positive samples without making too many false positive calls. Both curves plot these values over a range of decision thresholds. For both metrics, a value of 1 is a perfect classifier.

## 2.2 Pick optimal number of properties to define a community

We found random forest classifiers trained using GloVe embedded data produce a significantly higher average area under the Receiver Operating Curve (AUROC) across all choices of hyperparameters and number of dimensions (Fig 2) than non-embedded data and PCA-embedded data (p << 0.05, rank sum test). Notably, embedded data consistently produces better results with far fewer features than taxonomic counts. The use of fewer features makes the model less likely to overfit the data and more likely to be reproducible. We run all future tests using 100 properties, as models trained with 100 properties show the most consistently high performance and small variance across all hyperparameter choices.

Figure 2: Transforming ASV tables into GloVe embedding space before training a model produces more accurate host phenotype predictions (IBD vs. healthy control) and makes models more robust to hyperparameter choice. Each point represents a triplet of choices for number of trees, depth of each tree, and weight on a positive prediction of IBD in a random forest model. Each model was trained on the data input type indicated by color (Normalized, non-embedded

8

214  counts is purple, pca embedded data is pink, and GloVe embedded data is blue). Models trained on GloVe embedded

215  data produce higher ROC AUCs with less variance across hyperparameter choice.

216

217

218  **2.3 Models built with embedded data perform better on a held out test set**

219  We then train three separate models on the training portion of the AGP dataset, and test each model on a held

220  out portion of the same dataset that has been used neither for model nor embedding training (Fig 3 panel A).

221  Each model uses a different data input type, GloVe embedded, PCA-transformed, or non-embedded

222  normalized taxa counts, and has hyperparameters optimized using cross-validation over the training set. We

223  see comparable performance between the classifier using GloVe embedded data and the other two methods

224  (Fig 3 panel B). The model with non-embedded data, which uses 26,739 features, has an area under the

225  Receiver Operating Curve (AUROC) of 0.79 and an area under the Precision-Recall curve (AUPR) of 0.46 (Fig

226  3, panel B.1). In contrast, the model using GloVe embedded data, which uses only 113 features, has a higher

227  AUROC of 0.81 but slightly lower AUPR of 0.44 (Fig 3 panel B.2).  A 200-fold decrease in number of features

228  used results in little change in relevant performance metrics. In comparison, the model using PCA-transformed

229  data with 113 features performs only slightly worse, with an AUROC of 0.77 and an AUPR of 0.42 (Fig 3 panel

230  B.3)

231

232  **Fig. 3:** Embeddings trained on American Gut training set, model trained on American Gut training set, model tested on

233  American Gut held out test set (A). Models trained on GloVe embedded data have higher ROC AUC but slightly lower

234  Precision-Recall AUC on a held out test set (B)

235

236  **2.4 Properties are generalizable to independent stool-associated datasets**

237  We find that GloVe embedded data generalizes to a completely independent datasets, and significantly

238  improves performance when fewer than 400 training samples are available. Using data from Halfvarson et. al

239  (8), we train random forest classifiers on gut microbiome data to differentiate between IBD vs. healthy control

240  (Fig. 4A).  Again, we train classifiers using normalized count data, PCA-embedded data, and GloVe embedded

241  data, and optimize over hyperparameters using cross-validation for each model independently. To test the

242  effect of training set size on performance outcomes, we train models using from 50 to 450 samples in the

9

243 training set, and the rest in the test set.  In this dataset, we have 564 samples from 118 patients and 17, 775

244 Amplicon Sequence Variances (ASVs).  We do not include any associated metadata; predictions are made

245 solely based off of the microbiome compositions.

246     It is important to note that the transformation matrix that puts the query dataset into embedding space

247 is trained exclusively on American Gut Project data, and is therefore completely independent of the query

248 dataset. Despite the fact that properties were learned using the American Gut data dataset exclusively, we see

249 better embedding model performance on the independent set from Halfvarson et. al (8) (Fig. 4B). In particular,

250 we see that as the number of training samples becomes smaller, embedding-based models are able to

251 maintain high AUROC (Fig. 4B.1) and AUPR (Fig. 4B.2)  while models based on PCA-transformed data (100

252 features) and  non-dimensionality reduced models (17,775 features) cannot. When large numbers of training

253 samples are available, all methods perform comparably, but only embedding-based models perform well at

254 middling to low (< 400) sample sizes.

255     The patterns learned by the GloVe algorithm from the American Gut data generalize to improve

256 classification performance on an independent dataset. Theoretically, classification accuracy of any host

257 phenotype relating to the gut microbiome could be bolstered by first embedding the input data before model

258 training.

259

260 **Figure. 4:** Embeddings trained on American Gut data, model trained and tested on Halfvarson dataset (A). Transforming

261 microbiome data into GloVe embedding space prior to model training produces more accurate models despite smaller

262 training sample sizes (B).

263

264 **2.5 Models that use properties are generalizable to independent datasets**

265 In the above experiments, all models were trained on the same datasets they were tested on, using cross-

266 validation and a held-out test set. Now, we trained a model on the American Gut data and tested it on the

267 Halfvarson data (Fig 5A). More so than a hold-out test set, this allows us to test the feasibility of deploying a

268 model for diagnosis and monitoring of IBD. Two models were trained, one using normalized taxa counts and

269 the other taxa counts embedded in property space. In this case, only microbiome data and no sample-

270 associated data was included. Hyperparameters that gave the highest F1 score on American Gut data were

10

271 selected, and the trained model was directly applied to the independent dataset without re-tuning

272 hyperparameters or decision thresholds. Both models trained on American Gut taxa count and American Gut

273 embedded data had a precision of 1, meaning that a positive IBD prediction was correct 100% of the time.

274 However, the model trained on taxa counts had a recall of 0.02, meaning that only 2% of the samples from

275 patients with IBD were positively identified. In contrast, the model trained on embedded data recovered 26% of

276 samples from patients with IBD. While the model trained on taxa counts was in no way generalizable to

277 another dataset, the model trained on data in property space was able to make accurate predictions on a

278 completely independent dataset (Fig 5B). This finding demonstrates that in this case, models built from

279 embedded data can generalize to outside data while models built from taxa abundance information cannot.

280

281 **Figure 5:** Models and embeddings trained on American Gut data and tested on Halfvarson data (A). Model trained on

282 properties far outperforms models trained on taxa counts (B).

283

284 **2.6 Distances in embedding space roughly correlate with phylogenetic distance**

285 Taxa close together in embedding space have similar co-occurrence patterns. We expect that phylogenetically

286 close taxa are more likely to fill the same ecological niches than are unrelated taxa. We therefore expect a

287 slight but not extreme correlation between phylogenetic distance and distance in embedding space. Using a

288 Mantel test (45), we do observe a low (coef = 0.12)  but significant (p = 0.001) correlation between the two

289 distance metrics, with more granularity available when comparing taxa in embedding space. This finding

290 demonstrates that co-occurrence patterns as captured by embeddings are a more sensitive distance metric

291 than phylogeny (Fig 6).

292

293 **Figure 6:** The contour plot shows that distances between pairs of taxa in GloVe embedding space roughly correlate with

294 distances between those taxa in phylogenetic space (A). A lighter color signifies a higher density of taxa pairs. There is

295 more granularity along the embedding space axis, implying that related taxa are more easily distinguished from each

296 other in embedding space than they are phylogenetically. A Mantel test shows a low slope but very statistically significant

297 correlation between the two distance metrics (p = 0.001)

298

**2.7 Relationship with Metabolic Capacity:**

We chose to preserve taxa co-occurrence patterns in embedding space because we hypothesize that those

patterns are driven by taxa functionality in an environment. As such, we evaluate the possibility of a connection

between annotated genetic capacity to express metabolic pathways and the properties that make up

embedding space. First, we find each Amplicon Sequence Variant's (ASV) nearest neighbor in the KEGG

database (46) using Piphillian (47), and use the KEGGREST API (48) to determine which pathways are

present in that ASV's genome.  This results in an ASV by pathway table where there are 11,893 ASVs with

near neighbors in the database, and 148 possible metabolic pathways. Then, we identify the significantly

correlated metabolic pathways for each property in embedding space. A permutation test is used to simulate a

null distribution of maximum correlations per embedding property and determine significance. We find that

every property significantly correlates with at least 1 annotated metabolic pathway. Suppl. Table 1 shows each

dimension and its significantly correlated metabolic pathways; each dimension has significant correlation with 3

to 57 pathways. We see that the magnitude of correlations between embedding dimensions and metabolic

pathways are far greater in the GloVe embedding case than in the PCA-transformed case (Fig 7).  Additionally,

none of the correlations between PCA dimensions and metabolic pathways are significant under a permutation

test after multiple hypothesis correction (Suppl. Fig 1). This suggests that the properties learned by the GloVe

algorithm based on co-occurrence patterns between taxa may actually reflect the metabolic capacity of those

taxa.


**Figure 7:** Dimensions in GloVe embedding space correlate with some metabolic pathway annotations (A), but dimensions

in PCA embedding space do not (B). Each column in each heat map represents a metabolic pathway from KEGG (e.g.

ko00983). Each row is a dimension in either GloVe or PCA embedding space.


**2.8 Interpreting the predictive model for IBD with metabolic pathways**

In order to explore the implications of properties and metabolic pathways for IBD, we calculated an association

score (see method section 4.8) between each property and a positive IBD prediction. The full tables of the

most predictive dimensions, their associated metabolic pathways, and their direction of influence on the

prediction can be found in Suppl. Table 2 (AG) and Suppl. Table 3 (Halfvarson). First, we identified those

12

327    properties strongly associated with a positive IBD prediction (association score above 8 in both Halfvarson and

328    AGP datasets). We then selected all metabolic pathways significantly correlated with more than one of those

329    highly relevant properties. In this way, we identified 45 metabolic pathways of interest for IBD (Suppl. Table 4).

330    The pathways fall broadly into 9 main categories according to the KEGG Brite database: steroid metabolism,

331    lipid metabolism, glycan biosynthesis, amino acid metabolism, antibiotic synthesis and resistance, bacterial

332    pathogenic markers, metabolism of terpernoids and polyketides, cell motility and cellular community formation,

333    and xenobiotics biodegradation/other metabolic function.

334

335    **2.9 Explaining the variance in properties**

336    Lastly, we sought to determine how much of the information contained in properties can be recapitulated by

337    looking at the above described annotated metabolic pathways, and how much was unique to each property.

338    For each property, we use a linear regression to predict the property values per taxa from the pathway

339    presence/absence per taxa. We report the $r^2$ statistic per property, and find that metabolic pathways can

340    explain a maximum of 36% of the variance in one property, and a minimum of 11% of another. This means

341    that, while there is a strong correlation between properties and annotated metabolic pathways, most of the

342    information contained in properties are not represented by annotated information (Suppl. Fig 2)

343

# 344   3. Discussion

345    In a data-driven field dominated by small sample sizes and large variable spaces, it is necessary to employ

346    some form of dimensionality reduction. Currently, this is done by filtering by taxa prevalence, clustering based

347    on phylogenetic proximity, or is not done at all. We present here a method to leverage massive public datasets

348    to learn an embedding space that represents the latent properties driving taxonomic abundances. By shifting

349    the paradigm of analysis from taxonomic counts to community-level microbiome properties, we enable more

350    holistic, comprehensive analysis that accounts for taxonomic relationships while simultaneously simplifying the

351    data.

352    We demonstrate that we can learn a fecal microbiome property space that is more apt at predicting the

353    IBD status of the host than non-reduced and PCA-reduced spaces, and remains accurate even at low training

354 sample sizes. We also present a classification model trained on property data that generalizes well between

355 datasets, where models based on taxonomic counts do not.  We lastly define the relationships between

356 embedding space and known metrics used to explore microbiomes like phylogenetic distance and metabolic

357 pathway genetic capacity.

358

359 **3.1 Properties**

360 Embedding is a technique used ubiquitously in machine learning, especially in natural language processing

361 (49–51). Embedding algorithms take discrete units of data (e.g. words or taxa) and embed them into a vector

362 space, preserving proximity between the units based on any metric that can compare two units. In the case of

363 embedding taxa, possible metrics include phylogenetic distance, genome similarity, or morphology: in this

364 paper the chosen metric to determine proximity between units is patterns of co-occurrence. The embedding

365 algorithm used in this paper is GloVe, an algorithm designed for word processing (44). Using this algorithm,

366 two taxa that occur with similar sets of other taxa at similar frequencies should be close in embedding space,

367 and two taxa that are found in the presence of different neighbor sets should be far from each other. To

368 visualize this, we return to the analogy of word analysis. Two words, "apple" and "banana", are close to each

369 other in embedding space because they tend to occur with similar sets of words like "eat", "fruit", "tasty", and

370 "smoothie". Likewise, the words "king" and "marshmallow" tend to occur in different contexts; "king" is most

371 often found in the company of words like "politics", "throne", and "empire" while "marshmallow" is found with

372 words like "toddler", "fluffy", and "scrumptious". Note that there are two ways words may be close in embedding

373 space. First, words may directly co-occur frequently, like the words "apple" and "banana". Instead, words may

374 be synonyms, which do not often co-occur directly with each other, but instead co-occur with similar patterns,

375 like "large" and "huge" both being used to describe giants, mountains, and appetites. Returning to the concept

376 of embedding taxa, we may use embeddings to discover relationships both between taxa that work together

377 directly, and between taxa that are synonymous and likely fill the same niche.

378 Once proximity in embedding space has been established, the data can immediately be used to

379 improve modeling efforts. Subsequently, conceptual properties can be assigned to the learned dimensions by

380 observing which entities have similar values in any given dimension. If "strawberry", "cookies", "cake", and "ice-

14

cream" all have high values in one dimension, and "mud", "medicine", and "brussel sprouts" all have low values

in that same dimension, we may call that dimension the "delicious" property.

We have shown that embedding an ASV table into property space using GloVe integrates patterns from

public data into modeling efforts, producing more accurate diagnostics while decreasing data dimensionality.

Classifiers built after transforming data in this way are more robust, and the same embeddings generalize to

improve the accuracy of classifiers built from completely independent datasets. Properties also allow models

trained on one dataset to be applied to another independent dataset with positive results.

In addition to improving classification accuracy for IBD, the embeddings quantify and simplify the

microbial landscape of gut microbiomes. Rather than considering a microbiome as a collection of bacterial

counts, all of which are mostly independent, we propose to describe a microbiome as a vector of values for the

relevant properties. Consider the example of distinguishing the recipe book from the food magazine; reducing

each into property space allows us to clearly see the differences in declarative and descriptive word usage

rather than counting the number of times the words "spinach", "tomato", and "bowl" were each used. Because

these properties are learned from the data directly, they are much less biased than manually engineered

features. Analysis performed on this latent property space is likely to be much more robust to variations in

datasets, addressing the problem of irreproducibility currently plaguing 16S microbiome studies (26,27).

## 3.2 Biologically driven dimensional reduction

We use unsupervised learning to define an embedding space where taxa proximity represents similarity in co-

occurrence patterns. Unsupervised learning limits the human decision-making bias in property definition, but

also produces unlabeled properties whose interpretation is not immediately obvious. We hypothesize that co-

occurrence patterns are driven by taxa function like metabolism, synthesis of secondary metabolites, and

secretion of antimicrobial products. We show in our pathway analysis that property distributions in fact do

correlate significantly with metabolic pathways. Therefore, the learned property space is likely informed by taxa

function from within a biological context. Some elements of property space may also be informed by other

factors, such as geography or diet commonalities between groups of people, and this should be explored

further.

15

### 3.3 Annotation Independent

While we have explored the associations between embedding properties and the annotated quantities of genetic potential, the power of this embedding technique is that it does not rely on annotations of known taxonomic groupings or full genomes in order to improve prediction accuracy of host phenotype. Because any ASV that has been observed during embedding training can be embedded, it is possible to describe the properties of uncultured and unannotated ASVs, and include this information in a classifier. The transformation into embedding space requires only an ASV table, and uses no sample associated data like lifestyle variables or diagnoses.

### 3.4 Implications for IBD

We were able to identify 9 main categories of KEGG BRITE pathways that were significantly correlated with properties associated with IBD (Suppl. Table 4). Among these pathways, both steroid metabolism and biosynthesis were found to be associated with IBD. Steroids are a well-known and commonly utilized treatment for patients with active Crohn's disease (52). Enrichment in steroid metabolism in the gut microbiome could be reflective of an increase in steroid availability due to treatment.

Several pathways belonging to the rather broad BRITE category of "other metabolic function" have already been well explored and characterized in the literature as related to IBD. Toluene degradation (KEGG pathway 00623) was found to be increased in both Crohn's disease (CD) and Ulcerative Colitis (UC) samples in a microbiome survey meta-analysis (53). Components of the benzoate metabolic pathway, including fluorobenzoate degradation (KEGG pathway 00364), were associated with IBD severity in a treatment-naive cohort with CD (54). Analysis of inflamed gut lining mucosa in patients with IBD also found decreased ascorbic acid content (KEGG pathway 00053)(55) All of these pathways, along with dioxin degradation, inositol phosphate metabolism, and lipoic acid metabolism, were associated with an IBD prediction in our model.

We also found multiple glycan biosynthesis pathways correlated with predictive IBD properties (KEGG pathways 00511, 00514, 00515, 00601). In particular, bacterial glycosphingolipid biosynthesis, a pathway which has anti-inflammatory effects when produced by the host epithelial cells (56), was found to be associated with IBD in our model. We speculate that this may indicate a shortage of glycosphingolipids in the gut environment, exacting positive selection pressure on microbes that can produce their own.

16

437 Lipopolysaccharide biosynthesis and multiple types of O-glycan biosynthesis were also implicated in our

438 model, all of whose association with IBD has been explored, briefly, in the literature (57–59). Given its

439 importance and consistency in our predictive model, this group of pathways may warrant further exploration.

440

441 **3.5 Limitations and future expansion of the work**

442 While embedding Amplicon Sequence Variants (ASVs) affords the benefits to classification and interpretation

443 previously discussed, it relies heavily on the definition of a "biologically meaningful unit" which will then be

444 embedded. For the sake of between-study replicability, we choose to measure the co-occurrence patterns of

445 ASVs (28) as a base unit. It may, however, prove more informative to define a biologically meaningful unit in

446 another way. For example, perhaps ASVs clustered at a 99% threshold more accurately capture meaningful

447 patterns in co-occurrence. We may also consider a variable threshold that is more representative of common

448 ancestry on a phylogenetic tree and aggregate based on clade architecture before embedding.

449 Additionally, the presented set of embeddings was constructed using only the forward reads from the

450 American Gut dataset, as reverse reads were not provided in the EBI database. Future embeddings

451 constructed from full length V4 or other 16S hypervariable regions will likely provide more accuracy and

452 specificity. New embedding transformation matrices would need to be trained for each new biome or segment

453 of 16S gene being explored.

454 In its current form, the algorithm does not make specific considerations for differences in sequencing

455 depth, which affects how many taxa can be observed in a given sample. Future iterations of this method could

456 include weights such that the observed absence of taxa in a sample with a large number of reads is weighted

457 more heavily than the absence of taxa in a sample with fewer reads.

458 While the construction of embeddings is not affected by the inconsistency of self-report data, the

459 accuracy of the classifier may be. In this study, we considered only a self-reported medical professional

460 diagnosis to be accurate, and rejected any self-diagnosis reports. While it is possible that classifier

461 performance would change with the inclusion of more liberal diagnostic criteria, the strict diagnosis definition

462 successfully generalized to an independent dataset, which was not self-reported (8).

463 Properties in embedding space have strong associations with metabolic pathway potentials, but it

464 remains unclear whether they truly represent the expression of those pathways. Future development could

17

465 also consist of integrating multi-omics datasets available in other studies, including the Human Microbiome

466 Project. Wet lab validation of these hypothesized property-metabolic expression associations would verify the

467 ability of GloVe embeddings to predict metabolic expression from 16S data. This would allow for the integration

468 of metabolic data from all observed taxa, not just those few whose full genomes are available in databases.

469 It might be possible to use the embedding space to identify taxa that form stable communities together -

470 taxa that are close in embedding space may stabilize each other in culture and *in vivo*. Through mechanisms

471 like cross-feeding, joint nutrient acquisition, and other cooperative behaviors, microbes may form groups that

472 are more versatile and secure than the individual species on their own. Taxa near each other in embedding

473 space, if they are not directly interacting, may have synonymous functions in their respective communities. By

474 clustering and categorizing microbes by their respective roles, we may gain insight into which bacterial

475 populations secure one another's stability. Particularly, the relationship between phylogenetic distance and

476 distance in embedding space may be of interest. Microbes that are very closely related through evolution but

477 have very dissimilar co-occurrence patterns may be particularly predictive of their environment, as they have

478 specialized quickly and efficiently. It may be that different variable regions better capture the co-occurrence

479 patterns of taxa, and so are more representative of taxonomic relationship to the environment.

480 Lastly, the embedding framework can be applied to any system or base unit of interest. It may be

481 particularly illustrative to embed genes from metagenomic datasets instead of taxa. This would allow us to

482 determine mathematical representations of the context of each gene, as well as to glean the robustness and

483 reproducibility benefits from dimensionality reduction for metagenomic data. As always, appropriate

484 benchmarking and exploratory analyses will be necessary to determine the appropriate use cases for this

485 technology.

486

487 **3.6 Conclusion**

488 By integrating patterns from public datasets into individual survey studies, we bring the increased statistical

489 power and generalizability of results of meta-analyses into each independent study. While this work shows the

490 value of an embedding framework for predicting IBD from the gut microbiome, this same framework can be

491 leveraged in any environment with enough data and for any predictive problem of interest. Furthermore, we

18

492  assert that analyses that define microbiomes by their latent properties instead of by their taxa member list are

493  more informative, reproducible, and relevant to the macroscopic world.

494

# 495  4. Materials and Methods

496  Code available at: https://github.com/MaudeDavidLab/embeddings

497

498  **4.1 Embeddings: GloVe algorithm**

499  We used the GloVe algorithm (44) on ASVs to generate embeddings. Briefly, the embedding algorithm (Figure

500  1B) learns taxa representations that maintain patterns in co-occurrence between pairs of taxa, and was used to

501  learn properties of microbial context. In this algorithm, the metric to be preserved is a function of $P\_ik / P\_jk$,

502  the probabilities of co-occurrence of taxa i and j with k. Variables i and j are the taxa being related, and k is a

503  third context taxa. The result from this algorithm is a representation of each taxa in x-dimensional space, where

504  x is chosen by the user. The x-dimensional space is shared across all taxa, and thus each dimension can be

505  interpreted as a property for which each taxa has a value. The number of dimensions, x, is a hyperparameter

506  to be tuned, and results are reported for a range of dimensions: 50, 100, 250, 500, and 750. Embeddings were

507  learned on 85% of the data, which 15% of samples set aside for testing.

508

509  **4.2 Transformation into embedding space**

510  In 16S survey studies, each sample is represented by a vector of its taxa abundances. Thus, we transform

511  samples into embedding space simply by taking the dot product between each sample's taxa vector and a

512  taxa's property vector. This gives an average of property values weighted by taxa abundance. We consider two

513  ASVs the same taxa if they are at least 97% similar and align with an e value less than $10^{-29}$.

514

515  **4.3 PCA transformation**

516  Predictive model performance using embedded data was compared against models trained on data

517  transformed with Principal Coordinate Analysis (PCA). PCA is an ordination technique that projects samples

518  into lower dimensional space while maximizing the variance of the projected data (60).

519

**4.4 Random Forest Predictions:**

521 The value of the embeddings was evaluated by success at predicting host IBD status using a random forest

522 model (60). The model was built using Python sci-kit-learn, and hyperparameters for the depth of tree, number

523 of trees, and weight on a positive prediction were selected using 10-fold cross validation on the training set. A

524 different model with different hyperparameters was built for each data type, normalized taxonomic

525 abundances, PCA embedded abundances, and GloVe embedded abundances. Counts were normalized by

526 applying an inverse hyperbolic sin function. Models also included self-reported sample metadata such as

527 exercise, sex, daily water consumption, probiotic consumption, and dietary habits. Models were evaluated by

528 their performance, namely area under the receiver operating curve, on the held out test set of 15% of samples.

529

**4.5 Correlations with KEGG Pathways**

531 For each ASV, we find its closest match, thresholded at 97% similarity, in the KEGG database using the

532 software Piphillan (47). Each possible metabolic pathway then gets assigned a 0 if it is absent or a 1 if it is

533 present in that nearest neighbor's genome.

534 Limiting the following analysis to include only those taxa that had near neighbors in the database, for

535 each of the properties in embedding space, we find its maximally correlated (absolute value) metabolic

536 pathway. Then, to ascertain whether those correlations were significant, we applied a permutation test (61).

537 We constructed 10,000 null pathway tables by permuting the rows of the original pathway table. We repeated

538 the above procedure, finding the maximally correlated pathway for each of the embedding dimensions in each

539 of the null pathway tables. This results in 10,000 maximum correlation values per embedding dimension, which

540 form a null distribution for each embedding dimension. The significance of the statistic in a permutation test is

541 calculated as the number of times a maximum correlation in a null pathway table was more extreme or equal to

542 the maximum correlation actually observed. Dimensions (columns) in both GloVe transformation matrix and

543 PCA rotation matrix space are normalized to mean 0 and variance 1 to account for differences in scales

544 between the two spaces. We report both the maximally correlated pathway for each property, all of which are

545 significant, and also *all* significantly correlated pathways per property.

546

20

### 4.6 Calculating Phylogenetic Distances

We produced a Multiple Sequence Alignment and subsequently a phylogenetic tree of all ASVs, using Clustal

W2 (62) multiple alignment and phylogeny creation software. The tip-to-tip phylogenetic distances were then

calculated between every pair of taxa using the dendropy python package (63).


### 4.7 Explaining variance of properties with metabolic pathways

For each property, we set up a linear regression where the property values per taxa are the response variable,

and the pathway presence/absence for each taxa are the independent variables. The $r^2$ statistic is reported to

assess the variance in property values explainable by the presence of annotated metabolic pathways.


### 4.8 Importance of properties and pathways in predictive model

In order to calculate the direction of association of a property with disease, we limit each tree in the random

forest to split on 3 variables. We then backtrace; if a higher value of the property led to an IBD prediction, we

add one to the association score between IBD and that variable. Likewise, if a lower value of the property led

to IBD, we subtract one from the association score.

In calculating metabolic pathway importance to the predictive model, we first find all properties that are

consistently associated with health or with IBD. Then, we count the number of times each pathway is

significantly correlated with one of those properties. If a pathway is significantly correlated with more than two

consistently predictive properties, it is considered important in that phenotype.


### 4.9 Dataset

Embeddings were trained using data from the American Gut Project (24). This crowdsourced project provides

16S samples from the United States, United Kingdom, and Australia, along with associated dietary, lifestyle,

and disease diagnosis information. Amplicon Sequences Variants (ASVs) were called using the DADA2

algorithm (64), resulting in 18,750 samples and 335,457 ASVs. Samples with fewer than 5,000 reads and

ASV's not occurring in at least .07% of samples (13 samples) were then discarded, resulting in 18,480 samples

and 26,726 ASVs. Embeddings were trained on a randomly selected 85% of the filtered samples, and the other

15% were set aside for classifier testing.

21

575    Training embeddings does not require labeled data, and so samples could be used irrespective of their

576    available metadata. The machine learning classifier was trained and tested only on samples that had a positive

577    or negative IBD diagnosis, 5018 and 856 samples respectively.  IBD diagnosis was provided in various self-

578    reported options from the American Gut study: "I do not have this condition", "Self-diagnosed", "Diagnosed by a

579    medical professional (doctor, physician assistant)", or "diagnosed by an alternative medicine practitioner". For

580    this study, we considered only samples claiming a medical professional diagnosis to be true.

581    Lastly, in order to test the generalizability of embeddings, we used 16S data on patients with Crohn's

582    Disease (CD) and Ulcerative Colitis (UC) and healthy controls from Halfvarson et. al (8). DADA2 (64)  was

583    again used to call ASVs, samples were discarded if they had fewer than 10,000 reads, and ASVs were not

584    filtered for prevalence. After quality control, 26,251 ASVs remain, 17,775 of which have near neighbor

585    representations in embedding space. The dataset included samples with multiple diagnoses, but for the sake

586    of consistency, we focused on the most common diagnoses of Crohn's disease, Ulcerative Colitis, and healthy

587    control. In total, this left 564 samples from 118 patients, as the dataset contains multiple timepoints for each

588    patient. When models were trained and tested on Halfvarson datasets, timepoints from the same patients were

589    included entirely in the train or test set, so as not to train then test on the samples from the same patient.

590

591    **4.10 Machine Learning Performance Metrics**

592    We used two main performance metrics: area under the Receiver Operating Curve (AUROC) and area under

593    the Precision-Recall Curve (AUPR). The Receiver Operating Curve plots true positive calls against false

594    positive calls. The higher the AUROC, the more confident you can be that a positive prediction by the classifier

595    is correct. The Precision-Recall Curve plots the precision, how confident you are that a positive call is correct,

596    against recall, how many of the positive samples in the dataset were identified. A high AUPR means the

597    classifier is able to identify most of the positive samples without making too many false positive calls. For both

598    metrics, a value of 1 is a perfect classifier.

599

600    **4.11 Workflow**

601    The workflow is as described in Figure 1:

22

602       First, we learn the embedding space using taxa-taxa co-occurrence data from the American Gut Project

603 (A). The data contains 18,480 samples and 26,726 ASVs. Two taxa are considered co-occurring if they are

604 detected in the same fecal sample. From the patterns of co-occurrence across all samples, the GloVe

605 algorithm produces a transformation matrix, where each Amplicon Sequence Variant (taxa) is represented by a

606 vector in embedding space (B). We call each dimension in embedding space a "property" ($P\_1...P\_k$) as each

607 is a set of numbers used to differentiate taxas' co-occurrence patterns. No metadata is used to create the

608 embeddings; the process is completely unsupervised.

609       To transform the dataset of interest into embedding space (E), we take the dot product between the

610 dataset (D) and the transformation matrix (C). The dot product operation outputs a matrix of samples by

611 properties, where property vectors are calculated as the average of property vectors over all the taxa present

612 in that sample.

613       Lastly, we input the transformed data into a random forest classifier (F), along with 13 sample-

614 associated features like exercise frequency, probiotic consumption, frequency of vegetable intake (G), to train

615 a model that predicts IBD vs. Healthy host status. Samples and their associated features can be found in

616 Supp. Table 5.

617       In total, three random forest classifiers are trained, with the three types of input data: GloVe embedded,

618 PCA transformed, and non-embedded normalized count data. Each classifier was cross-trained on 85% of

619 samples to optimize hyperparameter choices for the number of decision trees, the depth of each tree, and the

620 weight put on a positive classification.

621

622 **4.12 Software and packages**

623 **Python Packages:** Pandas   0.23.4, Numpy 1.16.3, Sklearn 0.20.2, Scipy 1.2.0, Matplotlib 3.0.0, Re    2.2.1,

624 Skbio 0.5.5 **R packages:**  pheatmap_1.0.12, cowplot_0.9.4 , ggplot2_3.2.0 , RColorBrewer_1.1-2,

625 Gtools_3.8.1, Dada2_1.10.1, Rcpp_1.0.1, plyr_1.8.4, stylo_0.6.9, KEGGREST_1.22.0

626

627

628

23

# Acknowledgements

# Supplementary Materials Captions

**Supp Table 1:** Properties matched to all significantly correlated metabolic pathways. Includes KEGG pathway identifier and annotated pathway name.

**Supp Table 2:** Properties listed by importance in differentiating between IBD and healthy control samples in American Gut data using a random forest with a depth of 2. Each property is labeled by its maximally correlated metabolic pathway, and the direction of association it has with disease. The last column reports the cumulative number of trees in a cross-validated random forest that support that association.

**Supp Table 3:** Properties listed by importance in differentiating between IBD and healthy control samples in Halfvarson data using a random forest with a depth of 2. Each property is labeled by its maximally correlated metabolic pathway, and the direction of association it has with disease. The last column reports the cumulative number of trees in a cross-validated random forest that support that association.

**Supp Table 4:** List of pathways significantly correlated with properties strongly associated with IBD.

**Supp Table 5**: Samples and sample associated information converted into numeric quantities for machine learning.

**Supp Figure 1:** Heatmaps showing correlations between dimensions in transformed space and one null annotated metabolic pathway table. We see far fewer and less dramatic correlations between transformed data and metabolic pathways when pathway table has been shuffled.

656 **Supp Figure 2:** Histogram depicting the percent variance of properties explainable by annotated metabolic

657 pathways. Most properties are less than 25% explained by pathways, and no property is more than 36%

658 explained.

659

660 # References

661

662 1. Fettweis JM, Serrano MG, Brooks JP, Edwards DJ, Girerd PH, Parikh HI, et al. The vaginal microbiome

663 and preterm birth. Nat Med. 2019 Jun;25(6):1012–21.

664 2. Xu H, Li H. Acne, the Skin Microbiome, and Antibiotic Treatment. Am J Clin Dermatol. 2019

665 Jun;20(3):335–44.

666 3. Williams MR, Costa SK, Zaramela LS, Khalil S, Todd DA, Winter HL, et al. Quorum sensing between

667 bacterial species on the skin protects against epidermal injury in atopic dermatitis. Sci Transl Med. 2019

668 May 1;11(490):eaat8329.

669 4. Huttenhower C, Kostic AD, Xavier RJ. Inflammatory Bowel Disease as a Model for Translating the

670 Microbiome. Immunity. 2014 Jun 19;40(6):843–54.

671 5. Franzosa EA, Sirota-Madi A, Avila-Pacheco J, Fornelos N, Haiser HJ, Reinker S, et al. Gut microbiome

672 structure and metabolic activity in inflammatory bowel disease. Nat Microbiol. 2019 Feb;4(2):293–305.

673 6. Abbas M, Le T, Bensmail H, Honavar V, EL-Manzalawy Y. Microbiomarkers Discovery in Inflammatory

674 Bowel Diseases using Network-Based Feature Selection. In: Proceedings of the 2018 ACM International

675 Conference on Bioinformatics, Computational Biology, and Health Informatics  - BCB '18 [Internet].

676 Washington, DC, USA: ACM Press; 2018 [cited 2019 Jun 4]. p. 172–7. Available from:

677 http://dl.acm.org/citation.cfm?doid=3233547.3233602

678 7. Abraham C, Medzhitov R. Interactions Between the Host Innate Immune System and Microbes in

679 Inflammatory Bowel Disease. Gastroenterology. 2011 May 1;140(6):1729–37.

680 8. Halfvarson J, Brislawn CJ, Lamendella R, Vázquez-Baeza Y, Walters WA, Bramer LM, et al. Dynamics of

681 the human gut microbiome in inflammatory bowel disease. Nat Microbiol. 2017 Feb 13;2:17004.

682 9. Schirmer M, Franzosa EA, Lloyd-Price J, McIver LJ, Schwager R, Poon TW, et al. Dynamics of

25

683      metatranscription in the inflammatory bowel disease gut microbiome. Nat Microbiol. 2018;3(3):337–46.

684    10.  Bercik P, Verdu EF, Foster JA, Macri J, Potter M, Huang X, et al. Chronic Gastrointestinal Inflammation

685      Induces Anxiety-Like Behavior and Alters Central Nervous System Biochemistry in Mice.

686      Gastroenterology. 2010 Dec 1;139(6):2102-2112.e1.

687    11.  Peirce JM, Alviña K. The role of inflammation and the gut microbiome in depression and anxiety. J

688      Neurosci Res. 2019 May 29;

689    12.  Yang B, Wei J, Ju P, Chen J. Effects of regulating intestinal microbiota on anxiety symptoms: A

690      systematic review. Gen Psychiatry. 2019;32(2):e100056.

691    13.  Cheung SG, Goldenthal AR, Uhlemann A-C, Mann JJ, Miller JM, Sublette ME. Systematic Review of Gut

692      Microbiota and Major Depression. Front Psychiatry [Internet]. 2019 Feb 11 [cited 2019 Aug 5];10.

693      Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6378305/

694    14.  Stower H. Depression linked to the microbiome. Nat Med. 2019;25(3):358.

695    15.  Butler MI, Sandhu K, Cryan JF, Dinan TG. From isoniazid to psychobiotics: the gut microbiome as a new

696      antidepressant target. Br J Hosp Med Lond Engl 2005. 2019 Mar 2;80(3):139–45.

697    16.  Hsiao EY, McBride SW, Hsien S, Sharon G, Hyde ER, McCue T, et al. The microbiota modulates gut

698      physiology and behavioral abnormalities associated with autism. Cell. 2013 Dec 19;155(7):1451–63.

699    17.  Bolte ER. Autism and clostridium tetani. Med Hypotheses. 1998 Aug;51(2):133–44.

700    18.  David MM, Tataru C, Daniels J, Schwartz J, Keating J, Hampton-Marcell J, et al. Crowdsourced study of

701      children with autism and their typically developing siblings identifies differences in taxonomic and

702      predicted function for stool-associated microbes using exact sequence variant analysis. bioRxiv. 2018

703      May 25;319236.

704    19.  Finegold SM, Molitoris D, Song Y, Liu C, Vaisanen M, Bolte E, et al. Gastrointestinal Microflora Studies in

705      Late-Onset Autism. Clin Infect Dis. 2002 Sep;35(s1):S6–16.

706    20.  Sharon G, Cruz NJ, Kang D-W, Gandal MJ, Wang B, Kim Y-M, et al. Human Gut Microbiota from Autism

707      Spectrum Disorder Promote Behavioral Symptoms in Mice. Cell. 2019 May 30;177(6):1600-1618.e17.

708    21.  Dodiya HB, Forsyth CB, Voigt RM, Engen PA, Patel J, Shaikh M, et al. Chronic stress-induced gut

709      dysfunction exacerbates Parkinson's disease phenotype and pathology in a rotenone-induced mouse

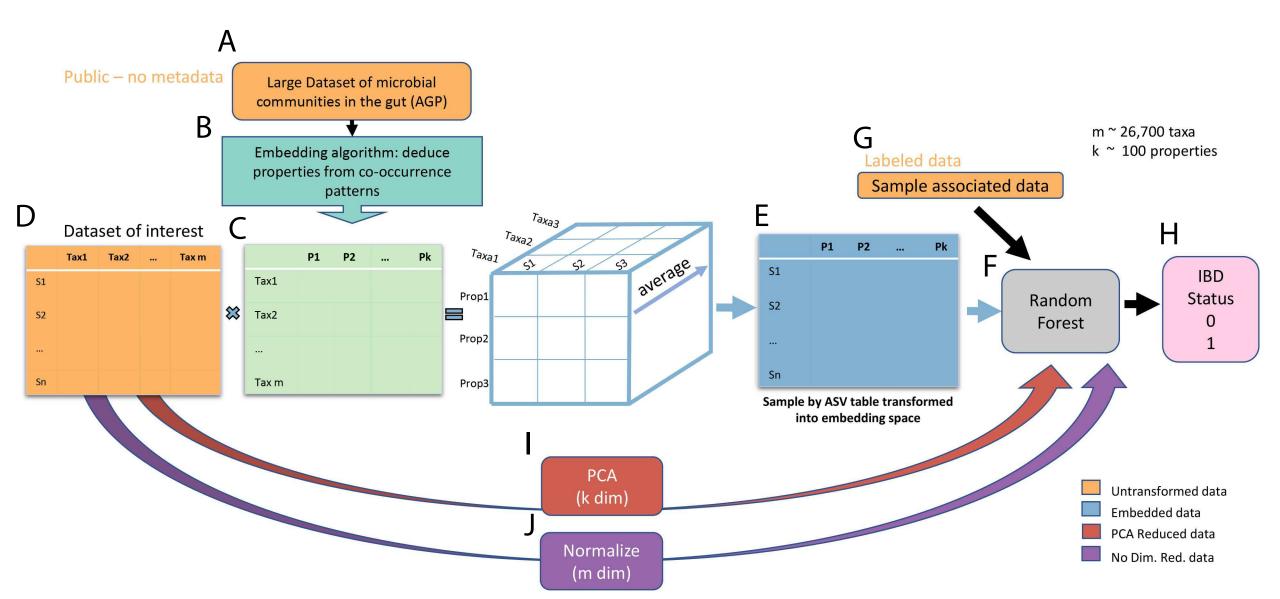710      model of Parkinson's disease. Neurobiol Dis. 2018 Dec 21;

711  22.  Dutta SK, Verma S, Jain V, Surapaneni BK, Vinayek R, Phillips L, et al. Parkinson's Disease: The

712       Emerging Role of Gut Dysbiosis, Antibiotics, Probiotics, and Fecal Microbiota Transplantation. J

713       Neurogastroenterol Motil. 2019 Jul 1;25(3):363–76.

714  23.  Santos SF, de Oliveira HL, Yamada ES, Neves BC, Pereira A. The Gut and Parkinson's Disease-A

715       Bidirectional Pathway. Front Neurol. 2019;10:574.

716  24.  McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, et al. American Gut: an Open

717       Platform for Citizen Science Microbiome Research. mSystems [Internet]. 2018 May 15 [cited 2018 Dec

718       7];3(3). Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5954204/

719  25.  The Human Microbiome Project Consortium, Huttenhower C, Gevers D, Knight R, Abubucker S, Badger

720       JH, et al. Structure, function and diversity of the healthy human microbiome. Nature. 2012

721       Jun;486(7402):207–14.

722  26.  Sinha R, Ahsan H, Blaser M, Caporaso JG, Carmical JR, Chan AT, et al. Next steps in studying the

723       human microbiome and health in prospective studies, Bethesda, MD, May 16–17, 2017. Microbiome

724       [Internet]. 2018 Nov 26 [cited 2019 Aug 21];6. Available from:

725       https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6257978/

726  27.  Schloss PD. Identifying and Overcoming Threats to Reproducibility, Replicability, Robustness, and

727       Generalizability in Microbiome Research. Ravel J, editor. mBio. 2018 Jun 5;9(3):e00525-18,

728       /mbio/9/3/mBio.00525-18.atom.

729  28.  Callahan BJ, McMurdie PJ, Holmes SP. Exact sequence variants should replace operational taxonomic

730       units in marker-gene data analysis. ISME J. 2017 Dec;11(12):2639–43.

731  29.  Poussin C, Sierro N, Boué S, Battey J, Scotti E, Belcastro V, et al. Interrogating the microbiome:

732       experimental and computational considerations in support of study reproducibility. Drug Discov Today.

733       2018 Sep 1;23(9):1644–57.

734  30.  Ioannidis JPA. Why Most Published Research Findings Are False. PLOS Med. 2005 Aug 30;2(8):e124.

735  31.  Fan J, Guo S, Hao N. Variance estimation using refitted cross-validation in ultrahigh dimensional

736       regression: Variance Estimation using Refitted Cross-validation. J R Stat Soc Ser B Stat Methodol. 2012

737       Jan;74(1):37–65.

738  32.  Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with
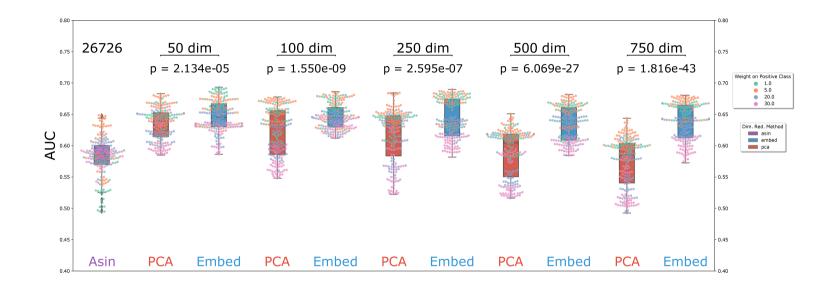
739     DESeq2. Genome Biol. 2014 Dec 5;15(12):550.

740  33. Paulson JN, Stine OC, Bravo HC, Pop M. Differential abundance analysis for microbial marker-gene

741     surveys. Nat Methods. 2013 Dec;10(12):1200–2.

742  34. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD. Analysis of composition of

743     microbiomes: a novel method for studying microbial composition. Microb Ecol Health Dis. 2015;26:27663.

744  35. Belenguer A, Duncan SH, Calder AG, Holtrop G, Louis P, Lobley GE, et al. Two routes of metabolic

745     cross-feeding between Bifidobacterium adolescentis and butyrate-producing anaerobes from the human

746     gut. Appl Environ Microbiol. 2006 May;72(5):3593–9.

747  36. Sankaran K, Holmes SP. Latent Variable Modeling for the Microbiome. ArXiv170604969 Stat [Internet].

748     2017 Jun 15 [cited 2019 May 15]; Available from: http://arxiv.org/abs/1706.04969

749  37. McMurdie PJ, Holmes S. phyloseq: An R Package for Reproducible Interactive Analysis and Graphics of

750     Microbiome Census Data. PLOS ONE. 2013 Apr 22;8(4):e61217.

751  38. Brooks AW, Priya S, Blekhman R, Bordenstein SR. Gut microbiota diversity across ethnicities in the

752     United States. PLOS Biol. 2018 Dec 4;16(12):e2006842.

753  39. Washburne AD, Morton JT, Sanders J, McDonald D, Zhu Q, Oliverio AM, et al. Methods for phylogenetic

754     analysis of microbiome data. Nat Microbiol. 2018 Jun;3(6):652–61.

755  40. Wakita Y, Shimomura Y, Kitada Y, Yamamoto H, Ohashi Y, Matsumoto M. Taxonomic classification for

756     microbiome analysis, which correlates well with the metabolite milieu of the gut. BMC Microbiol. 2018 Nov

757     16;18(1):188.

758  41. Larsen P, Dai Y. Metabolome of human gut microbiome is predictive of host dysbiosis. GigaScience.

759     2015;4:42.

760  42. Panek M, Paljetak HČ, Barešić A, Perić M, Matijašić M, Lojkić I, et al. Methodology challenges in studying

761     human gut microbiota – effects of collection, storage, DNA extraction and next generation sequencing

762     technologies. Sci Rep. 2018 Mar 23;8(1):1–13.

763  43. Woloszynek S, Zhao Z, Chen J, Rosen GL. 16S rRNA sequence embeddings: Meaningful numeric

764     feature representations of nucleotide sequences that are convenient for downstream analyses. PLOS

765     Comput Biol. 2019 Feb 26;15(2):e1006721.

766  44. Pennington J, Socher R, Manning C. Glove: Global Vectors for Word Representation. In: Proceedings of
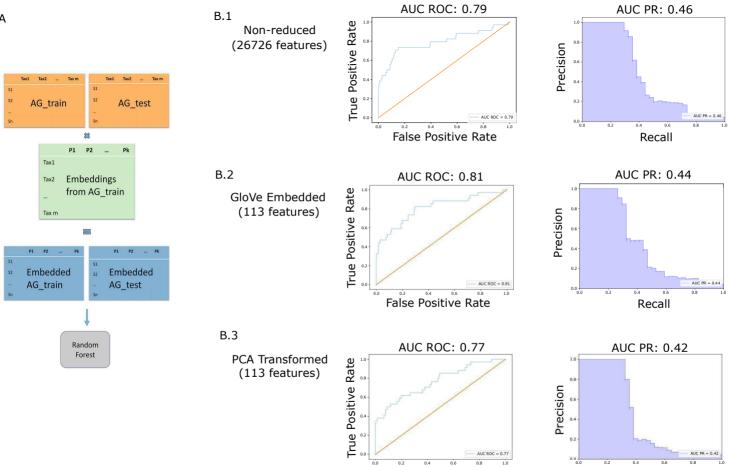
767    the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP) [Internet]. Doha,

768    Qatar: Association for Computational Linguistics; 2014 [cited 2019 Aug 6]. p. 1532–43. Available from:

769    http://aclweb.org/anthology/D14-1162

770    45.    Mantel N. The Detection of Disease Clustering and a Generalized Regression Approach. Cancer Res.

771    1967 Feb 1;27(2 Part 1):209–20.

772    46.    Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic Acids Res. 2000 Jan

773    1;28(1):27–30.

774    47.    Iwai S, Weinmaier T, Schmidt BL, Albertson DG, Poloso NJ, Dabbagh K, et al. Piphillin: Improved

775    Prediction of Metagenomic Content by Direct Inference from Human Microbiomes. PLOS ONE. 2016 Nov

776    7;11(11):e0166104.

777    48.    Tenenbaum D. KEGGREST: Client-side REST access to KEGG. 2018.

778    49.    Tang D, Wei F, Yang N, Zhou M, Liu T, Qin B. Learning Sentiment-Specific Word Embedding for Twitter

779    Sentiment Classification. In: Proceedings of the 52nd Annual Meeting of the Association for

780    Computational Linguistics (Volume 1: Long Papers) [Internet]. Baltimore, Maryland: Association for

781    Computational Linguistics; 2014 [cited 2019 Aug 22]. p. 1555–1565. Available from:

782    https://www.aclweb.org/anthology/P14-1146

783    50.    Wang Y, Liu S, Afzal N, Rastegar-Mojarad M, Wang L, Shen F, et al. A comparison of word embeddings

784    for the biomedical natural language processing. J Biomed Inform. 2018 Nov 1;87:12–20.

785    51.    Zou WY, Socher R, Cer D, Manning CD. Bilingual Word Embeddings for Phrase-Based Machine

786    Translation. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language

787    Processing [Internet]. Seattle, Washington, USA: Association for Computational Linguistics; 2013 [cited

788    2019 Aug 22]. p. 1393–1398. Available from: https://www.aclweb.org/anthology/D13-1141

789    52.    Dubois-Camacho K, Ottum PA, Franco-Muñoz D, De la Fuente M, Torres-Riquelme A, Díaz-Jiménez D,

790    et al. Glucocorticosteroid therapy in inflammatory bowel diseases: From clinical practice to molecular

791    biology. World J Gastroenterol. 2017 Sep 28;23(36):6628–38.

792    53.    Zhou Y, Xu ZZ, He Y, Yang Y, Liu L, Lin Q, et al. Gut Microbiota Offers Universal Biomarkers across

793    Ethnicity in Inflammatory Bowel Disease Diagnosis and Infliximab Response Prediction. mSystems

794    [Internet]. 2018 Jan 30 [cited 2019 Aug 20];3(1). Available from:

795      https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5790872/

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

797      naïve microbiome in new-onset Crohn's disease. Cell Host Microbe. 2014 Mar 12;15(3):382–92.

798  55.  Buffinton GD, Doe WF. Altered ascorbic acid status in the mucosa from inflammatory bowel disease

799      patients. Free Radic Res. 1995 Feb;22(2):131–43.

800  56.  Abdel Hadi L, Di Vito C, Riboni L. Fostering Inflammatory Bowel Disease: Sphingolipid Strategies to Join

801      Forces. Mediators Inflamm [Internet]. 2016 [cited 2019 Aug 20];2016. Available from:

802      https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4736332/

803  57.  Miyahara K, Nouso K, Saito S, Hiraoka S, Harada K, Takahashi S, et al. Serum Glycan Markers for

804      Evaluation of Disease Activity and Prediction of Clinical Course in Patients with Ulcerative Colitis. PLOS

805      ONE. 2013 Oct 7;8(10):e74861.

806  58.  Larsson JMH, Karlsson H, Crespo JG, Johansson MEV, Eklund L, Sjövall H, et al. Altered O-glycosylation

807      profile of MUC2 mucin occurs in active ulcerative colitis and is associated with increased inflammation.

808      Inflamm Bowel Dis. 2011 Nov 1;17(11):2299–307.

809  59.  Caradonna L, Amati L, Magrone T, Pellegrino NM, Jirillo E, Caccavo D. Enteric bacteria,

810      lipopolysaccharides and related cytokines in inflammatory bowel disease: biological and clinical

811      significance. J Endotoxin Res. 2000;6(3):205–14.

812  60.  Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine

813      Learning in Python. J Mach Learn Res. 2011;12:2825–2830.

814  61.  Pitman EJG. Significance Tests Which May be Applied to Samples From any Populations. Suppl J R Stat

815      Soc. 1937;4(1):119–30.

816  62.  Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, et al. The EMBL-EBI search and

817      sequence analysis tools APIs in 2019. Nucleic Acids Res. 2019 Jul;47(W1):W636–41.

818  63.  Sukumaran J, Holder MT. DendroPy: a Python library for phylogenetic computing. Bioinformatics. 2010

819      Jun 15;26(12):1569–71.

820  64.  Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution

821      sample inference from Illumina amplicon data. Nat Methods. 2016;13(7):581–3.
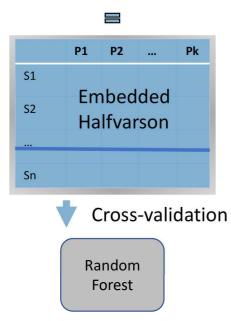
822

823

824

825

826

827

A

B.1 Non-reduced (26726 features)

AUC ROC: 0.79

AUC PR: 0.46

B.2 GloVe Embedded (113 features)

AUC ROC: 0.81

AUC PR: 0.44

B.3 PCA Transformed (113 features)

AUC ROC: 0.77

AUC PR: 0.42

A

B.1

B.2

A

B

| | Normalized Taxa Counts | GloVe embedded counts |
|---|---|---|
| Accuracy | 0.11 | 0.33 |
| Precision | 1.0 | 1.0 |
| Recall | 0.02 | 0.26 |
| F1 | 0.04 | 0.41 |
| F2 | 0.02 | 0.30 |

GloVe Embedded                          PCA Embedded

Dimensions (100)

Metabolic Pathways                        Metabolic Pathways

Dimensions (100)

GloVe Embedded

PCA Embedded

Dimensions (100)

Metabolic Pathways

Metabolic Pathways