



Published in final edited form as:

*Nat Methods*. 2019 December ; 16(12): 1306–1314. doi:10.1038/s41592-019-0616-3.

## Learning accurate representations of microbe-metabolite interactions

**James T. Morton**<sup>1,2</sup>, **Alexander A. Aksenov**<sup>3,4</sup>, **Louis Felix Nothias**<sup>3,4</sup>, **James R. Foulds**<sup>5</sup>, **Robert A. Quinn**<sup>6</sup>, **Michelle H. Badri**<sup>7</sup>, **Tami L. Swenson**<sup>8</sup>, **Marc W. Van Goethem**<sup>8</sup>, **Trent R. Northen**<sup>8,9</sup>, **Yoshiki Vazquez-Baeza**<sup>10,11</sup>, **Mingxun Wang**<sup>3,4</sup>, **Nicholas A. Bokulich**<sup>12,13</sup>, **Aaron Watters**<sup>14</sup>, **Se Jin Song**<sup>1,11</sup>, **Richard Bonneau**<sup>7,14,15,16</sup>, **Pieter C. Dorrestein**<sup>3,4</sup>, **Rob Knight**<sup>1,2,17,11</sup>

<sup>1</sup>Department of Pediatrics, University of California, San Diego, La Jolla, CA, USA

<sup>2</sup>Department of Computer Science & Engineering, University of California, San Diego, La Jolla, CA, USA

<sup>3</sup>Collaborative Mass Spectrometry Innovation Center, University of California San Diego, La Jolla, CA, USA

<sup>4</sup>Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA

<sup>5</sup>Department of Information Systems, University of Maryland Baltimore County, Baltimore, MD, USA

<sup>6</sup>Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI, USA

<sup>7</sup>Department of Biology, New York University, New York, 10012 NY, USA

<sup>8</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Rd, Berkeley, CA, 94720, USA

<sup>9</sup>DOE Joint Genome Institute, 2800 Mitchell Dr., Walnut Creek, CA, 94598, USA

<sup>10</sup>Jacobs School of Engineering, University of California San Diego, La Jolla, CA, USA

<sup>11</sup>Center for Microbiome Innovation, University of California, San Diego, La Jolla, CA, USA

<sup>12</sup>The Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA

<sup>13</sup>Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:[http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

### <sup>6</sup>Author contributions

J.T.M wrote the mmvec algorithm, conducted the benchmarks and ran all of the analyses. A.A.A. and L.F.N. preprocessed and annotated the metabolomics data. A.A.A. provided insights in the high fat diet study. J.R.F. provided insights behind word2vec and topic modeling. M.H.B. benchmarked SPIEC-EASI. R.A.Q. provided insights behind the cystic fibrosis study and simulations. Y.V.B. provided insights behind the interpretation of the IBD analysis. M.W. developed the GNPS workflow for mmvec. N.A.B developed the heatmap visualizations. A.W developed the network visualizations. A.W developed the network visualizations. T.L.S. M.W.V.G and T.N. provided insights behind the biocrust soils experiment. All authors were involved with writing the manuscript.

### <sup>7</sup>Competing interests

Mingxun Wang is the founder of Ometa Labs LLC. None of the remaining authors have any competing interests.

<sup>14</sup>Flatiron Institute, Simons Foundation, New York, 10010 NY, USA

<sup>15</sup>Computer Science Department, Courant Institute, New York, 10012 NY, USA

<sup>16</sup>Center For Data Science, NYU, New York, NY 10008, USA

<sup>17</sup>Department of Bioengineering University of California, San Diego, La Jolla, CA, USA

## Abstract

Integrating multi-omics datasets is critical for microbiome research, but multiple statistical challenges can confound traditional correlation techniques. We solve this problem by using neural networks to estimate the conditional probability that each molecule is present given the presence of each specific microbe. We show with known environmental (desert biological soil crust wetting) and clinical (cystic fibrosis lung) examples, our ability to recover microbe-metabolite relationships, and demonstrate how the method can discover relationships between microbially-produced metabolites and inflammatory bowel disease.

## 2 Introduction

Knowledge gained by integrating complementary “-omics” data with a multi-omics approach will lead to improved diagnostics, automated drug discovery, and optimized culturing conditions for uncharacterized microbes [1]. Previous work have been able to predict metabolite abundance profiles from microbe abundance profiles [2, 3]. However, because conventional correlation techniques have unacceptably high false discovery rates, finding meaningful relationships between genes within complex microbiomes and their products in the metabolome is challenging.

Although there has been a widespread effort to develop multi-omics approaches, several conceptual challenges limit techniques that integrate disparate “omics” data in general, including linking the microbial sequencing and untargeted mass spectrometry. Therefore, new approaches are needed that can handle disparate data types [4]. Relative abundances of thousands of microbes and metabolites can be measured using sequencing and mass spectrometry, resulting in the generation of very high dimensional microbiome and metabolomics datasets. Quantifying microbe-metabolite interactions from these abundances requires estimating a distribution across all possible microbe-metabolite interactions.

Techniques such as Canonical Correspondence Analysis (CCA) and Partial Least Squares (PLS) approximate this joint distribution using a low dimensional representations [5, 6, 7]. Network models have been shown to improve classification accuracy using multiple datasets [8]. Factor models have been proposed to incorporate multiple datasets for biomarker analysis [9]. Despite of the wide application of these methods, they are notoriously difficult to interpret [10, 11, 12] and it remains unclear whether these models can obtain individual microbe-metabolite interactions.

Pearson and Spearman correlations assume independence between interactions, simplifying the estimation procedure by reducing it to a combination of independent two dimensional problems. However, many studies have shown that these methods are not statistically valid

for compositional data, a fact first recognized by Pearson in 1895 and followed up in numerous studies [13, 14, 15, 16, 17]. This problem is further complicated because both microbiome [17] and mass spectrometry [18, 19, 20, 21] datasets are also compositional, meaning that the absolute abundances are not measured, which can confound statistical inference. For example, in untargeted mass spectrometry experiments, the set of molecules detected and their relative abundance vary depending on the extraction protocol and analytic methods used, which leads to only a partial snapshot of the metabolome. Moreover, measuring the total mass of molecules extracted is often not performed in large scale metabolomics efforts, due to the highly laborious nature of that step.

To understand how issues associated with compositional data impact inference on microbe-metabolite interactions, consider the example in Figure S1. There are two microbes and two metabolites in Figure S1a. All are increasing exponentially at different rates and are highly correlated with each other. If proportions are estimated from the absolute abundances via sampling, the information about the total microbe population size and the total metabolite abundance is lost, and the correlations between the microbes and the metabolites disappear. False positives can also appear as shown in Figure S1b, microbe and metabolite interactions that have no apparent correlation structure may appear to be correlated when investigating the proportions. These issues alone can give rise to overwhelming false positives and false negatives, making Pearson and Spearman in some scenarios comparable to random coin flips. Experimental validation currently takes large laboratories multiple years to perform [22], often requiring time-consuming manual examinations of erroneous correlations.

There are other compositional techniques such as SparCC[13] and proportionality[23] that are scale-invariant when analyzing a single dataset, but lose scale-invariance when analyzing multiomics datasets. This was shown in the context of identifying microbe-fungal interactions [24], which provided motivation to extend SPIEC-EASI [14] to handle multiomics datasets. We show that this approach does not work for microbe-metabolite interactions because of differences of measurement units between sequencing and mass spectrometry measurements (Supplementary materials). An alternative approach is to consider co-occurrence probabilities instead of correlations. Here, co-occurrence probabilities refer to the conditional probability of observing a metabolite given that a microbe was observed, thereby allowing us to identify the most likely microbe-metabolite interactions. To do this, we propose “mmvec”, (microbe-metabolite vectors), to learn these co-occurrence probabilities between microbes and metabolites (Figure 1). Due to its scale-robustness properties, the microbial-metabolite relationships learned by mmvec are consistent between the absolute and relative abundances. The microbe-metabolite interactions can be ranked [25] and visualized through standard dimensionality reduction interfaces, enabling interpretable findings. The computations behind mmvec can take advantage of modern GPU architectures using Tensorflow [26], enabling scalable inference on large multiomics datasets. Furthermore, we provide evidence in two benchmarks and four case studies that mmvec outperforms existing statistical methods.

### 3 Results and Discussion

We performed benchmarks comparing mmvec to Pearson, Spearman, SPIEC-EASI, SparCC and Proportionality [23] using a cystic fibrosis biofilm simulation. We then show that mmvec can resolve contradictory cyanobacteria-metabolite relationships in a desert soil biocrust wetting study. We also demonstrate recovery of known associations of *P. aeruginosa*-produced metabolites observed in cystic fibrosis [27]. Finally, we explore the relationships of microbiota and metabolic changes in mice fed a high fat diet [28] and inflammatory bowel disease [29], showing how this approach can be used to determine microbial origin of novel molecules even in extremely complex real-life biological systems with limited knowledge of existing associations.

#### 3.1 Simulation benchmarks

To compare mmvec performance to Pearson, Spearman, Proportionality, SparCC and SPIEC-EASI correlations, we used data from existing studies in which the relationships between microbes and metabolites were the central focus of investigation. One such study simulated spatial-temporal dynamics in a microbial biofilm [27]. The original study tested the hypothesis that the cystic fibrosis (CF) microbiome community within human lungs can be manipulated by altering its chemical environment. Changes in pH and oxygen saturation suppress the principal pathogen, *P. aeruginosa*, without using antibiotics, by promoting the growth of a community of fermenters that out-compete the pathogen. The simplicity of this system allowed the high-level ecological patterns to be modelled. In the original simulations, the interactions between two microbes (fermenters denoted by  $\Theta_f$  and *P. aeruginosa* denoted by  $\Theta_p$ ) and multiple molecules were modeled using Monod kinetics and diffusion processes [27] (Figure 2a).

We simulated the measurement process for microbial DNA sequencing and untargeted mass spectrometry for metabolites as discussed in the Online Methods, providing ground truth information on their interactions. The model simulates interactions between *P. aeruginosa* and the fermenters, and their interactions with the environment. It also simulates known interactions between microbes and molecules, such as sugar consumption by fermenters and ammonia production by the pathogen. For example, the fermenters are positively associated with sugars and ammonium concentration, and negatively associated with inhibitor concentration; *P. aeruginosa* is positively associated with amino acids and pH.

Therefore, we can test whether the top  $K$  metabolites associated with each microbe by each tool includes the correct microbe-metabolite interactions. Figure 2c shows specificity and sensitivity for each tools as a function of  $K$ . In these simulations, random chance outperformed all of the tools except for mmvec and SPIEC-EASI, with mmvec performing the best. As shown in Figure 2d and Figure S2, mmvec is the only method robust to scale deviations amongst the methods tested. This is critical for maintaining consistency between absolute and relative abundances, which can otherwise lead to inflated false positives and false negatives [16].

### 3.2 Soil biocrust wetting event case study

Many studies produce inconsistent results that can be resolved with improved data analysis, especially in environmental and clinical settings. To test whether mmvec can resolve unexplained discrepancies in microbe-metabolite interactions across studies, we applied it to a study of biocrust wetting [30]. In this study, laboratory-based exometabolite patterns observed with bacterial isolates were reproduced in the environment. Specifically, in this work authors identified metabolites that were consumed and released by multiple biocrust isolates including *Microcoleus vaginatus* and two *Bacillus* strains [31], and compared these patterns with closely-related environmental taxa and metabolites observed in situ [30].

While almost 70% of the examined microbe-metabolite relationships following the wetting event were validated [30], some contradicted the microbe-metabolite relationships observed in cultures [31]. These contradictions stemmed from Spearman correlations between *M. vaginatus* abundances and the observed metabolite abundances, but were resolved by mmvec (Figure 3a).

All metabolites released from the *M. vaginatus* isolate have higher conditional probabilities than the average metabolite following biocrust wetting, and are among the top 80 co-occurring metabolites with *M. vaginatus* (of 485 molecules total). This result supports the original finding that *M. vaginatus* actually releases these molecules after the wetting event. In contrast, Spearman labels 7 of 13 of these molecules with a negative correlation, indicating that these molecules were consumed by *M. vaginatus* rather than released, as originally stated in [30]. When the annotation detection rates amongst different statistical methodologies, mmvec has a substantially higher true positive rate as shown in Figure 3b.

The conflicting results between mmvec and Spearman could be explained by the growing microbial biomass and shift in available resources after wetting (Figure 3 c, d). Total biomass is expected to increase, because *M. vaginatus* releases metabolites that enable the growth of many other microbes. Because DNA sequencing can only measure proportions, the growth in other microbes could cause the proportions of *M. vaginatus* to decrease, leading to a misleading anti-correlation with 4-guanidinobutanoate (Figure 3d). However, it is not possible to infer whether *M. vaginatus* is decreasing in abundance [25] or 4-guanidinobutanoate is increasing in abundance.

The change in the total biomass and the total available resources could explain the contradiction between the Spearman correlations and the isolate results. *M. vaginatus* likely grows at a slower rate relative to other microbes that benefit from the metabolite release. Because mmvec does not rely on knowledge of the total biomass or normalize to relative abundance, these contradictions are avoided.

### 3.3 Cystic Fibrosis case study

To further validate if mmvec can detect known microbe-metabolite interactions in a biological setting, we re-analyzed a study on lung mucus microbiome of patients with cystic fibrosis [27, 32]. Cystic fibrosis has been shown to be dominated by two major groups of microbes, anaerobes and pathogens that occupy unique niches, and their interactions are defined by the environment. Anaerobes dominate in low oxygen and low pH environments,

while pathogens, in particular *P. aeruginosa*, dominate in the opposite conditions [27]. Mmvec clearly separates anaerobes and pathogens (Figure 4a), with known anaerobic microbes (*Veillonella*, *Fusobacterium*, *Prevotella* and *Streptococcus*) on the left, and notable pathogens, such as *P. aeruginosa*, on the right.

*P. aeruginosa* is known to produce small-molecule virulence factors [33]. In the original study, based on annotations from GNPS[34], the bacterium was found to produce six molecules: 4-hydroxy-2-heptylquinoline (HHQ), Pyocyanin (PYO), Phenazine-1-carboxylic acid (PCA), 2-nonyl-4-hydroxy-quinoline (NHQ), 2-heptyl-3,4-dihydroxyquinoline (PQS, *Pseudomonas* quinolone signal) and Pyochelin [27]. As shown in Figure 4a, mmvec identifies these molecules with a high co-occurrence probability with *P. aeruginosa*. Mmvec also identifies a cluster of rhamnolipids likely produced by *P. aeruginosa*. Rhamnolipids are well characterized and are an important virulence factors for *P. aeruginosa*, contributing to biofilm development, motility on surfaces and antagonistic interactions with host inflammatory cells [35, 36]. These rhamnolipids were not identified in the original study [27]. The annotations for these compounds have been established using GNPS [34].

There is a negative correlation between the first principal component learned from mmvec and the metabolites log-fold change across the oxygen gradient (Figure 4b) (Pearson  $r = -0.59$ ,  $p$ -value  $1.8 \times 10^{-44}$ ,  $n = 442$  molecules), which is consistent with the findings in the original work. No such correlation between the oxygen gradient and the first microbial principal component was found by Pearson ( $r = 0.11$ ,  $p = 0.16$ ,  $n = 138$  microbes). There exist two notable microbes on opposing ends of the first microbial principal component: *P. aeruginosa*, a known pathogen, and *Streptococcus*, a known anaerobe. The top 100 metabolites that are specific to *P. aeruginosa* and *Streptococcus* are shown to have drastically different profiles in samples where *P. aeruginosa* and *Streptococcus* were the most abundant species (Figure 4d,e) (logratio t-test=6.51,  $p = 4.4 \times 10^{-8}$ ,  $n = 49$  samples). This provides evidence that in the context of this study, the metabolomic profiles can be largely influenced by the most abundant microbes, a notion that has important implications for understanding CF etiology. To further support this, the learned metabolite conditional probabilities for *P. aeruginosa* can be used to predict the metabolite proportions in the 41 samples where *P. aeruginosa* is the most abundant taxa. The predicted *P. aeruginosa* metabolite profiles alone can explain 10% of the metabolite variation in these samples ( $r = 0.319$ ,  $p = 1.18 \times 10^{-11}$ ,  $n = 442$  molecules).

Of 14 quinolone molecules known to be produced by *P. aeruginosa*, Pearson correlation detected 9 with  $p < 0.05$  without FDR correction, and only 5 with FDR correction. For example, Pyocyanin, does not appear related to *P. aeruginosa* by the raw proportions ( $r = 0.158$ , FDR-corrected  $p$ -value=0.089, rank=96,  $n = 172$  samples), but is ranked 34th most associated with *P. aeruginosa* by mmvec (Figure S3c), consistent with culturing experiments that demonstrate that *P. aeruginosa* produces this molecule [37]. 18 rhamnolipids are among the top 25 metabolites most associated with *P. aeruginosa* by mmvec, and have higher ranks with mmvec than with Pearson correlation (Figure S3b).

### 3.4 Effects of high fat diet in murine model case study

We then tested whether mmvec could determine the microbial origin of specific molecules in a complex biological system. We recently discovered a new kind of bile acid, where cholate is conjugated to amino acids other than glycine and taurine [38]. These molecules increased in abundance with high-fat diet in humans. We determined that these molecules are microbially-made since they were present in specific pathogen free, but not in germ free mice. We therefore set out to identify candidate producers. We were able to confirm that one of these bile acids, cholate phenylalanine amidate, was associated with high-fat diet in well-controlled study that investigated the development of non-alcoholic fatty liver disease (NAFLD), cirrhosis, and hepatocarcinoma (HCC) in a mouse model [28]. When re-analyzing these datasets for differential abundances via multinomial regression, the strong association of the novel bile acid with HFD became immediately apparent. The use of mmvec showed distinct associated groups of microbes and HFD (Figure 5a) and a clear stratification of the mass spectrometry data according to diet (Figure 5b). Several *Clostridium spp.* correlated with the cholate phenylalanine conjugate. Indeed, we showed that *Clostridium spp.* were found to produce this bile acid [38]. This result demonstrates mmvec's ability to streamline the discovery of microbes that produce specific molecules of interest *in vivo*.

### 3.5 Microbe-metabolite interactions in Inflammatory Bowel Disease

Finally, microbe-metabolite interactions were investigated for samples of IBD patients generated under the integrative Human Microbiome Project [29]. The role of the microbiome in IBD is acknowledged, but still poorly understood. The original study uncovered shifts in metabolomic and microbial profiles associated with the IBD. In particular, levels of carnitines and bile acids were shown to be affected [29]. Using mmvec we confirmed the core findings in the previous study, such as the co-occurrence between *R. hominis* and multiple carnitines, including previously noted C20, which have anti-inflammatory properties (Figure 6a) [29]. We also found high correlation of *Klebsiella spp.* with IBD status and that it co-occurs with high probability with several bile acids (Figure 6b). Although *Klebsiella* itself does not produce these compounds, some pathogens (including *Klebsiella*) are known to be resistant to bile acids [39]. Excessive production of some bile acids and bile acid malabsorption can lead to overabundance of bile acids, which is a hallmark of IBD [40], although the exact mechanisms remain unknown. The ability of *Klebsiella* to thrive in concentrated bile acid environments is consistent with the high co-occurrence probabilities shown in Figure 6b. We also noted that three *Klebsiella* species are the top drivers of the IBD-associated molecules (Figure 6c). It is important to delineate different reasons for co-occurrence. Unlike *Klebsiella*, *Clostridium* species are known for bile acid manipulation, including production of bile that can germinate *Clostridium difficile* spores or that have anti-microbial properties [41, 42].

Therefore, it is possible that in case of *Clostridia*, the existing co-occurrences (Figure 6b) are due to actual biosynthesis of the metabolites by the microbial species indicated rather than ability to withstand them.

In addition to recapitulating reported findings, mmvec also yielded previously undetected relationships. The major microbe that was found to be associated with healthy patients is *Propionibacteriaceae*, which was not detected in Price et al 2019 (Figure 6cd). This relationship is corroborated by other published studies. In one study, it has been shown that some members of the *Propionibacterium* genus produce 1,4-Dihydroxy-2-naphthoic acid (DHNA), a growth stimulator for bacteria such as *Bifidobacterium* that are thought to reduce the symptoms of IBD [43]. Also, in a survey of *in vivo* vs. *in vitro* bacterial activity, *Propionibacterium freudenreichii* was shown to play an immunomodulatory role in the context of an ulcerative colitis mice model [44]. In another study it was shown that *Propionibacterium freudenreichii* is a viable core component in an anti-inflammatory probiotic fermented dairy product [45]. The members of this family have been considered beneficial for intestinal immunoregulation; *Propionibacteriaceae* have been observed to be enriched in human breast milk and have been shown to restore Th17 differentiation [46]. Thus, it appears that the existing knowledge supports the statistically-inferred interaction uncovered by mmvec, but not identified in the published analysis of the dataset.

## 4 Conclusion

In both simulation benchmarks and annotated dataset, mmvec shows promise for inferring microbe-metabolite interactions from multiomics datasets. Our benchmarks suggest that mmvec outperforms all existing tools that aim to infer interactions between paired microbe-metabolite abundance datasets, both in simulations and in experimental data. In the biocrust wetting experiment, mmvec resolved conflicting findings between the *in vitro* validated *M. vaginatus* released metabolites and the sequencing/mass spectrometry analysis of environmental samples. In the cystic fibrosis study, mmvec can reliably identify all of the experimentally determined *P. aeruginosa*-produced molecules of interest. We show in the example of bile acid production that mmvec enables exploratory analysis in complex biological systems and streamlined discovery of the microbial origin of specific metabolites. Finally, mmvec was able to identify the strongest microbial contributions to the metabolite abundances in the IBD study, where one of those microbes was missed in the original study.

In light of these findings, the current methodology still has limitations. It remains unclear how to access statistical significance of an interaction using co-occurrence probabilities. Similarly, confidence intervals for the strength of each microbe-metabolite interaction can not yet be calculated. Furthermore, more theoretical work will be required to handle continuous-valued inputs.

The concepts outlined here should generalize beyond microbe-metabolite interactions to handle other paired multi-omic data types, provided that the input dataset is made up of counts (as in metagenomics, transcriptomics, etc.). With the exponential growth of multiomics datasets, there is much potential to use these methods to reveal microbial metabolism, including for microbes that are not cultivable in the laboratory. Approaches utilizing co-occurrence probabilities have the potential to enable more targeted experimental assays, accelerating the discovery of microbe-metabolite interactions, paving the way towards new ecosystems engineering approaches in clinical, environmental and industrial applications.



## 10 Methods

### 10.1 Mmvec neural network architecture

The development of our proposed neural network was inspired by applications in natural language processing. The underlying model can also be referred to as a bi-loglinear multinomial regression. Our mmvec model posits an assumed generative process for the data, which leads to an inference algorithm to recover the model's parameters from multi-omics data. The model's assumed generative model for metabolite  $v$ , microbe  $\mu$  and sample  $k$  given as follows.

First generate microbe vector  $\mathbf{u}_\mu$  for microbe  $\mu \in \{1, \dots, N\}$  and metabolite vectors  $\mathbf{v}_v$  for metabolite  $v \in \{1, \dots, M\}$ ,

$$\mathbf{u}_\mu \sim \mathcal{N}(\mathbf{0}, \sigma_u I) \quad \mathbf{v}_v \sim \mathcal{N}(\mathbf{0}, \sigma_v I),$$

These vectors are length  $p$ , corresponding to the number of latent vectors dimensions. Each of these vectors are drawn from a normal prior centered around zero and a diagonal covariance matrix with variances  $\sigma_u$  and  $\sigma_v$ , namely to serve regularization purposes and avoid overfitting. For a given microbial sample  $x_k$ , the model's generative process draws a single microbe from a single draw from the categorical distribution

$$\mu \sim \text{Categorical}(\mathbf{x}_k).$$

That microbe  $\mu$  can be used to index  $U$  in order to generate conditional probabilities  $\mathbf{q}_\mu$

$$p(v|\mu) = \frac{\exp(\mathbf{v}_v \cdot \mathbf{u}_\mu + v_{v0} + u_{\mu0})}{\sum_j \exp(\mathbf{v}_j \cdot \mathbf{u}_\mu + v_{j0} + u_{\mu0})},$$

$$\mathbf{q}_\mu = [p(v_1|\mu), \dots, p(v_M|\mu)]$$

Here,  $v_{j0} + u_{\mu0}$  are row and column biases, which are required to accurately estimate the conditional probabilities. The above transformation is the softmax transform [47] to compute probabilities from real-valued quantities. This transformation is also known as the inverse clr transform [48], which enforces scale invariance as shown in the simulations. In the mmvec model's generative process, these conditional probabilities generate the metabolite abundances  $\mathbf{y}_k$  for a given sample  $k$  through a multinomial distribution.

$$\mathbf{y}_k \sim \text{Multinomial}(n, \mathbf{q}_\mu),$$

where  $n$  is the total metabolite abundances across sample  $k$ . It is important to note that metabolite abundances themselves are not counts, but rather a continuous representation of

molecule counts. We make the simplifying assumption that these continuous valued abundances can be approximated by Multinomial count models.

This model bears resemblance to how word2vec estimates word probabilities conditioned on a single particular word [49]. There are a couple of major differences to be considered. First, in the original application of word2vec, a skipgram was proposed. Skipgrams [49] have been designed to account for the sequential nature of text. There is no such sequential nature with microbiome or metabolite samples, the only ordering information that is known is the sample membership. As a result, the skipgrams can be replaced using multinomial sampling, where a single microbe is randomly sampled from a microbiome sample at each gradient descent step.

Second, in the original word2vec application a single input/output word pair were evaluated at each gradient descent step, which is required to incorporate the contextual information of words within sentences. In the application of multiomics, this is unnecessarily complicated, since there is no such contextual with regards to microbes and metabolites. Instead, all of the metabolite abundances can be simultaneously evaluated for each gradient descent step, ultimately speeding up computations. Specifically, these metabolite abundances are simultaneously considered in order to estimate the conditional probabilities  $q_k$  for the given microbial count  $u_{jk}$ . From these conditional probabilities, the metabolite abundances  $y_k$  are generated from a Multinomial distribution. This process is repeated across all of the microbial reads. To show that  $p(y|\mu)$  truly approximates the probability of observing a metabolite given a microbe, we first need to make the simplifying assumption that the conditional distribution of a metabolite given the presence of a single microbe also follows a multinomial distribution as follows

$$p(Y = y | X_{\mu} = 1) = \text{Multinomial}(y | q_{\mu})$$

Where  $y$  is the vector of observed metabolites,  $Y$  is the random variable modeling metabolite abundances,  $X$  is a random variable modeling microbe abundances,  $x$  is a vector of observed microbes and  $\mu$  is a single microbe. Given these modeling assumptions, we can parameterize the conditional Multinomial distributions with embedding vectors as described above. This estimation procedure can be reformulated as a matrix factorization, where the conditional probability matrix is decomposed into two weight matrices  $\mathbf{U}$  and  $\mathbf{V}$ , which are comprised of microbe-metabolite vectors as follows

$$\mathbf{U} = [\mathbf{0}, \mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_N]^T \quad \mathbf{V} = [\mathbf{v}_0, \mathbf{0}, \mathbf{v}_1, \dots, \mathbf{v}_M].$$

Here  $\mathbf{U} \in \mathbb{R}^{N \times p}$  and  $\mathbf{V} \in \mathbb{R}^{(M-1) \times p}$  represents the corresponding embeddings for  $N$  microbes and  $M$  metabolites. The number dimensions  $p$  for both  $\mathbf{U}$  and  $\mathbf{V}$  as well as the priors are specified by the user, but can also be evaluated during cross-validation. The biases  $\mathbf{u}_0$  and  $\mathbf{v}_0$  are critical for estimating accurate co-occurrence probabilities, as suggested by similar methodologies used in recommender systems [50]. The  $\mathbf{U}$  and  $\mathbf{V}$  matrices are estimated through maximum a posteriori (MAP) estimation using ADAM [51] with the following log-posterior

$$\mathcal{L} = \mathcal{L}_Y + \mathcal{L}_U + \mathcal{L}_V$$

$$\mathcal{L}_U = \sum_{\mu} \sum_{q=1}^p \mathcal{N}(U_{\mu,q} | \theta_{\mu}, \sigma_u)$$

$$\mathcal{L}_V = \sum_{\nu} \sum_{q=1}^p \mathcal{N}(V_{\nu,q} | \theta_{\nu}, \sigma_{\nu})$$

$$\mathcal{L}_Y = \sum_k \sum_{r \in x_k} \text{Multinomial}(y_k | q_{\mu})$$

Within a single iteration of stochastic gradient descent a single microbial sequence  $i$  is randomly drawn and compared to a complete set of metabolite abundances  $y_i$  for that given sample. If there are a total of  $R$  microbial reads across all of the microbial samples, there will be  $R$  iterations for a complete epoch over the microbial dataset. This means that the running time of this training process is  $O(RM)$  for a single epoch. Cross validation can be performed by holding out samples measuring the predictive power by looking at the sum of squares errors. Predictions can be made as follows

$$SSE = \sum_{k,i} \left( y_k - m_k \cdot \text{softmax}(\mathbf{V}\mathbf{U}_{u_{ik}, \cdot}) \right)^2$$

Where the predictive metabolite abundances are compared to the holdout abundances  $y_k$  across all microbial reads  $i$  in the holdout samples  $k$ .  $m_k$  denotes the total metabolite abundances in sample  $k$

## 10.2 Microbe-metabolite vectors in simplicial coordinates

Here, we will provide some insights behind the underlying geometry behind this neural network. Doing so will provide intuition behind the algebraic operations commonly applied in the context of word2vec, suggesting the possibility of performing similar tasks in the context of microbe-metabolite interactions. Furthermore, this will motivate the use of the Aitchison distance to quantify microbe-microbe and metabolite-metabolite interactions. Finally we will make a connection to topic modeling, providing another means to potentially interpret the latent dimensions in the model. The connection between the softmax and the inverse clr transform suggests that the inputs to this transform can be represented in clr coordinates. The softmax function and its corresponding inverse, the clr transform, is given as follows

$$\text{softmax}(x) = \left[ \frac{e^{x_1}}{\sum_i e^{x_i}}, \dots, \frac{e^{x_1}}{\sum_i e^{x_i}} \right]$$

$$\text{clr}(z) = \left[ \log \frac{z_1}{g(z)}, \dots, \log \frac{z_D}{g(z)} \right]$$

Since biases are incorporated into the mmvec model, by construction  $Q=UV^T$  is both row centered and column centered, meaning that the sum of rows are zero and the sum of the columns are zero. Given this the following holds

**Theorem:** If  $Q=UV$  and  $\mathbf{1}_N Q = \mathbf{0}$  and  $Q \mathbf{1}_M = \mathbf{0}$  then  $U \mathbf{1}_p = \mathbf{0}$  and  $V \mathbf{1}_p = \mathbf{0}$

Suppose that there exists another solution  $Q = UV^{*T}$  where  $V = V - \mathbf{1}_M \lambda_v^T$  and  $\lambda_v \in R^p$ . Then

$$Q = U(V - \mathbf{1}_M \lambda_v^T).$$

Given that the rows of  $Q$  sum to 0, then

$$U(V - \mathbf{1}_M \lambda_v^T)^T \mathbf{1}_M = 0$$

$$U \lambda_v^T M = 0.$$

This means that only the trivial solution  $\lambda_v = 0$  exists, therefore the rows of  $V$  do sum to 0.

Using the same reasoning above, suppose that there exists another solution  $Q = U^* V^T$  where  $U^* = U - \mathbf{1}_N \lambda_u^T$  and  $\lambda_u \in R^p$ . Then

$$Q = (U - \mathbf{1}_N \lambda_u^T) V^T.$$

Given that the columns of  $Q$  sum to 0, then

$$\mathbf{1}_N^T (U - \mathbf{1}_N \lambda_u^T) V^T = 0$$

$$N \lambda_u^T V = 0.$$

This means that only the trivial solution  $\lambda_u = 0$  exists, therefore the rows of  $U$  do sum to 0.

Therefore the rows of both  $\mathbf{U}$  and  $\mathbf{V}$  must sum to zero if  $\mathbf{U}$  and  $\mathbf{V}$  are non-trivial.

As noted in previous compositional data analysis work, the sum of the components within a vector in clr coordinates is zero. Given that the row vectors within  $U$  and  $V$  both sum to zero, that suggests that each of these vectors are also in clr coordinates. This means the following properties are satisfied

**Topic proportions**—Since the  $U$  and  $V$  row vectors are in clr coordinates, that implies that these row vectors can be directly converted to  $p$ -dimensional proportions, yielding a similar interpretation to topics used in models such as LDA [52, 53].

**Linearity**—Vectors in clr coordinates are known to satisfy linearity, namely

$$\text{clr}(\alpha x + y) = \alpha \text{clr}(x) + \text{clr}(y)$$

for  $\alpha \in \mathbb{R}$ ,  $x \in \mathcal{S}^p$  and  $y \in \mathcal{S}^p$ . This linearity property was leveraged in word2vec models to perform analogy reasoning. Since both microbes and metabolites are in clr coordinates, it should be possible to categorize microbe-microbe and metabolite-metabolite interactions.

**Isometry**—The clr transform is distance preserving, meaning that the Aitchison distance on proportions is equivalent to the Euclidean distance on clr vectors. This provides motivation for using Euclidean distances to compute microbe-microbe and metabolite-metabolite similarities.

### 10.3 Visualization through biplots

Visualization techniques from compositional data analysis can aid with interpretation [54, 55].  $\mathbf{U}$  and  $\mathbf{V}$  can be visualized as factors within a biplot to visualize the microbe-metabolite embeddings on a single plot. The first two latent dimensions of  $\mathbf{U}$  represent microbial coordinates on a 2D scatter plot and the first two latent dimensions of  $\mathbf{V}$  represent metabolite coordinates on a 2D scatter plot. Typically the coordinate from the  $\mathbf{V}$  matrix are plotted as arrows from the origin in order to identify features that explain the variance in  $\mathbf{U}$ . However, in our case studies, there are typically many more metabolites than microbes - so we opt to visualize the metabolites as points and microbes as arrows for a simpler visualization. As suggested by the above theorem, the distance between points approximates the Aitchison distance between metabolites, and the distance between arrow tips approximates the Aitchison distance between microbes. As suggested in [56], the Aitchison distance is also equivalent to the variance of the log ratios, suggesting that microbe-microbe and metabolite-metabolite distances could also be interpreted as a measure of proportionality [23].

### 10.4 Benchmarks

The simulated data was based on a cystic fibrosis biofilm model derived in Quinn et al [27] shown in Figure S12 in the paper. The biofilm model was built to explain how fermenters and *P. aeruginosa* responded to different concentrations of sugars, amino acids, pH, oxygen and antibiotics across the Winogradsky column. These models solved for differential

equations integrating Monod kinetics and diffusion processes and was run in Matlab using the code provided at [https://github.com/zhangzhongxun/WinCF\\_model\\_Code](https://github.com/zhangzhongxun/WinCF_model_Code)

From this simulation, we only focus 2 microbes and 5 compounds. The two microbes are *P. aeruginosa* ( $\Theta_p$ ) and fermenters ( $\Theta_f$ ). The five compounds (SG), acids (F), ammonium (P), amino acids (SA) and inhibition molecules (I). In order to simulate a high dimensional dataset, each microbial taxon was split into 50 different subtaxa and each compound was split into 50 molecular subclasses. The partitioning procedure is given as follows

$$\mathbf{p}_i \sim \mathcal{N}(\mathbf{0}, \sigma_o \mathbf{I}) \quad \mathbf{q}_i \sim \mathcal{N}(\mathbf{0}, \sigma_c \mathbf{I})$$

$$\mathbf{o}_{ij} = \kappa_{ij} \text{ilr}^{-1}(\mathbf{p}_i) \quad \mathbf{c}_{ik} = \eta_{ik} \text{ilr}^{-1}(\mathbf{q}_i),$$

where  $\mathbf{p}_i$  is a vector proportions representing how the subtaxa corresponding to  $j$  will be distributed in sample  $i$ .  $\kappa_{ij}$  represents the absolute abundance of taxon  $j$  in sample  $i$ .  $\mathbf{o}_{ij}$  represents a vector of the absolute abundances for all of the subtaxa corresponding to taxon  $j$ . These are the absolute abundances that are used for comparison in Figure 2.

Here we use the  $\text{ilr}^{-1}$  transform to generate proportions from a multivariate normal distribution. Here the multivariate normal distribution is centered around zero, and the covariance matrix  $\sigma_o \mathbf{I}$  has only a constant diagonal structure with a tunable parameter  $\sigma_o$  specifying the variability of the partitioning procedure. Larger values of  $\sigma_o$  will cause the allocations of the microbes to be increasingly uneven.

The partitioning procedure is identical for the metabolites.  $\mathbf{q}_i$  is a vector proportions representing how the subcompounds corresponding to  $k$  will be distributed in sample  $i$ .  $\eta_{ik}$  represents the absolute abundance of compound  $k$  in sample  $i$ .  $\mathbf{c}_{ik}$  represents a vector of the absolute abundances for all of the subtaxa corresponding to compound  $k$ . The multivariate normal distribution used to generate the proportions is centered around zero. The covariance matrix  $\sigma_c \mathbf{I}$  has only a constant diagonal structure with a tunable parameter  $\sigma_c$  specifying the variability of the partitioning procedure. Larger values of  $\sigma_c$  will cause the allocations of the metabolites to be increasingly uneven.

Once the subtaxa and subcompounds absolute abundances have been simulated, the microbial relative counts and metabolite abundances are simulated. The sampling procedure is performed as follows

$$\xi_i \sim \mathcal{LN}(n, \tau_o) \quad \omega_i \sim \mathcal{LN}(m, \tau_c)$$

$$x_i \sim \mathcal{PLN}(\xi_i, \mathbf{C}(\mathbf{o}_i), \varepsilon_o) \quad y_i \sim \mathcal{LN}(\omega_i, \mathbf{C}(\mathbf{c}_i), \varepsilon_c).$$

The total sequencing depths and total intensities for sample  $i$  are drawn from Lognormal distributions with means parameterized by  $n$  and  $m$  and overdispersion parameters  $\tau_o$  and  $\tau_c$ . We chose to use the lognormal distribution for three reasons. First, the lognormal distribution models overdispersion. Second, the lognormal distribution has a simpler interpretation than other overdispersed distributions such as the negative binomial, since the parameters can be directly interpreted as a normal distribution and consequentially has a compositional interpretation due to its connection to the ilr transform. Finally, the lognormal distribution is commonly used for modeling in the ecological literature in the context of studying species populations in Niche theory and Neutral theory, leading to a natural biological interpretation.

Once the total sequencing depth and the total intensities are sampled, the microbial sequencing counts and metabolite abundances are then sampled. A Poisson lognormal distribution is used to generate the microbial counts from the microbial proportions  $C(\mathbf{o}_i)$  scaled by the sequencing depth  $\zeta_i$ . The counts are sampled with error  $\epsilon_o$ . A Lognormal distribution is used to generate the metabolite abundances from metabolite proportions  $C(\mathbf{c}_i)$  scaled by the total intensity  $\omega_i$ . The abundances are sampled with error  $\epsilon_c$ . All of the code used to generate the benchmarks can be found at <https://github.com/knightlab-analyses/multiomic-cooccurrences>.

## 10.5 Software workflows

To facilitate utilization of the mmvec tool, we have developed two different user-facing interfaces. First, we have developed a qiime2 plugin [57], where mmvec can be run using a simple command line interface. This interface is complemented using [26], where users can monitor convergence rates for their models in real-time and evaluate how different parameters will affect their model fit (Figure S4). Second, we have integrated mmvec into the Global Natural Product Social Molecular Networking (GNPS) platform that can be accessed by the public. The online interface through GNPS resolves several usability issues. First, GNPS facilitates import of metabolomics data into qiime2 by pre-processing, importing, and sample renaming. This is performed as part of the standard metabolomics analysis at GNPS (e.g. molecular networking and feature-based molecular networking). Second, since it is possible to both download and re-use outputs of workflows run at GNPS directly, it is straightforward to select the GNPS qza and molecule annotations needed for mmvec. The user will need to upload the accompanying feature and taxonomy data for qiime2 and the analysis will begin. Once the workflow completes, the biplots can be viewed directly in the browser and other outputs (e.g. ranks) are available for download (Figure S5).

The mmvec implementation is written using Tensorflow and can leverage GPUs for computation. The number of gradient descent iterations is specified by the user and model fit diagnostics can be monitored in real time using Tensorboard. The runtime of mmvec across 16 cores can take multiple days until a model convergence reaches convergence. With GPUs, the running time is reduced to a few hours. Using a Tesla GPU, the model can reach convergence within 4 hours on the IBD dataset comprised of 562 microbial taxa, 26,966

metabolite features and 400 samples. However, there is a trade-off of accuracy and running time. More accurate models require smaller learning rates and may take longer to run.

## 10.6 Data Analysis

Due to the overwhelming sparsity in microbiome datasets, some filtering is required in order to infer microbe-metabolite interactions. We chose to filter out microbes that appear in less than 10 samples, since these microbes don't have enough information to infer which metabolites are co-occurring with them. In other words the mmvec model has too many degrees of freedom to perform inference on these microbes. For the cystic fibrosis study, there were 172 samples and after filtering there were 138 unique microbial taxa and 462 metabolite features. For the biocrust soils study, there were 19 samples and after filtering there were 466 unique microbial taxa and 85 metabolite features. For the murine high fat diet study, there were 434 samples and after filtering there were 902 microbes and 11978 metabolites. For the IBD dataset, there were 13920 features in the c18 LCMS dataset, 26966 features in the c8 LCMS dataset and 562 taxa. Cross validation was performed across all studies to evaluate overfitting. In the desert biocrust soils experiment, 1 sample out of 19 samples was randomly chosen to be left out for cross-validation. In all of the other studies, 10 samples were randomly chosen to be left out for cross-validation. All of the analyses can be found under <https://github.com/knightlab-analyses/multiomic-cooccurences>.

## 10.7 Data availability

The cystic fibrosis sequencing and metadata data can be found under <http://qiita.microbio.me>; study id: 10863. The corresponding GNPS analysis can be accessed at <http://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=34d825dbf4e9466e81d809faf814995b>.

The biocrust soils data was retrieved from the supplemental section in Swenson et al [30]. The High fat diet murine model case study 16S rRNA data can be found under <http://qiita.microbio.me>; study id: 10856. The High fat diet murine model case study are publicly available at <https://massive.ucsd.edu/> at MassIVE ID MSV000080918. The GNPS analysis for this study can be accessed at <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=977d85bba47b4e96bf69872b961b8edd>

The IBD data used can be found under <https://ibdmdb.org>.

See Life Sciences Reporting Summary for more details on the experimental design.

## 10.7 Software availability

The software implementing the mmvec algorithm can be found under <https://github.com/biocore/mmvec>.

Differential abundance analyses in the high fat diet study was performed using L2-regularized multinomial regression using software available at <https://github.com/biocore/songbird>

The software used to build the multiomics network can be found at [https://github.com/mortonjt/multiomics\\_network](https://github.com/mortonjt/multiomics_network).



Biplots were generated using Emperor [58].

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

We would like to thank Vera Pawlowsky, Juan Jose Egozcue and Susan Holmes for their insights behind the geometry of this neural network model. In addition, we would also like to thank Nicholas Bokulich for their feedback and contributions on the mmvec software package. T.L.S., M.W.V.G and T.R.N greatly acknowledge funding from the Office Science Early Career Research Program, Office of Biological and Environmental Research, of the U.S. Department of Energy under contract number DE-AC02-05CH11231 to Lawrence Berkeley National Laboratory. This was in part supported by P41GM103484 Center for Computational Mass Spectrometry, Instrument support through NIH S10RR029121 and R03 CA211211 on reuse of metabolomics data. Y. V. B. is funded by the Janssen Human Microbiome Institute through a collaboration with the Center for Microbiome Innovation. J.T.M. was funded by NSF grant GRFP DGE-1144086. R. K. and S. J. S. have been funded by Janssen under grant number 20175015 and the Alfred P. Sloan Foundation under grant number G-2017-9838.

## 8 References

- [1]. Jansson Janet K and Baker Erin S. A multi-omic future for microbiome studies. *Nat Microbiol*, 1(16049):645, 2016.
- [2]. Noecker Cecilia, Eng Alexander, Srinivasan Sujatha, Theriot Casey M, Young Vincent B, Jansson Janet K, Fredricks David N, and Borenstein Elhanan. Metabolic model-based integration of microbiome taxonomic and metabolomic profiles elucidates mechanistic links between ecological and metabolic variation. *MSystems*, 1(1):e00013–15, 2016.
- [3]. Mallick Himel, Franzosa Eric A, McIver Lauren J, Banerjee Soumya, Sirota-Madi Alexandra, Kostic Aleksandar D, Clish Clary B, Vlamakis Hera, Xavier Ramnik J, and Huttenhower Curtis. Predictive metabolomic profiling of microbial communities using amplicon or metagenomic sequences. *Nature communications*, 10(1):3136, 2019.
- [4]. Knight Rob, Vrbanac Alison, Taylor Bryn C, Aksenov Alexander, Callewaert Chris, Debelius Justine, Gonzalez Antonio, Kosciolk Tomasz, McCall Laura-Isobel, McDonald Daniel, et al. Best practices for analysing microbiomes. *Nature Reviews Microbiology*, page 1, 2018.
- [5]. Meng Chen, Zeleznik Oana A, Thallinger Gerhard G, Kuster Bernhard, Gholami Amin M, and Culhane Aedín C. Dimension reduction techniques for the integrative analysis of multi-omics data. *Brief. Bioinform*, 17(4):628–641, 7 2016. [PubMed: 26969681]
- [6]. Le Gall Gwénaëlle, Noor Samah O, Ridgway Karyn, Scovell Louise, Jamieson Crawford, Johnson Ian T, Colquhoun Ian J, Kemsley E Kate, and Narbad Arjan. Metabolomics of fecal extracts detects altered metabolic activity of gut microbiota in ulcerative colitis and irritable bowel syndrome. *Journal of proteome research*, 10(9):4208–4218, 2011. [PubMed: 21761941]
- [7]. Rohart Florian, Gautier Benoit, Singh Amrit, and Le Cao Kim-Anh. mixomics: An r package for ‘omics feature selection and multiple data integration. *PLoS computational biology*, 13(11):e1005752, 2017. [PubMed: 29099853]
- [8]. Wang Bo, Mezlini Aziz M, Demir Feyyaz, Fiume Marc, Tu Zhuowen, Brudno Michael, Haibe-Kains Benjamin, and Goldenberg Anna. Similarity network fusion for aggregating data types on a genomic scale. *Nature methods*, 11(3):333, 2014. [PubMed: 24464287]
- [9]. Argelaguet Ricard, Velten Britta, Arnol Damien, Dietrich Sascha, Zenz Thorsten, Marioni John C, Buettner Florian, Huber Wolfgang, and Stegle Oliver. Multi-omics factor analysis—a framework for unsupervised integration of multi-omics data sets. *Molecular systems biology*, 14(6):e8124, 2018. [PubMed: 29925568]
- [10]. Cajo JF Braak Ter and Verdonschot Piet FM. Canonical correspondence analysis and related multivariate methods in aquatic ecology. *Aquatic sciences*, 57(3):255–289, 1995.

- [11]. Witten Daniela M, Tibshirani Robert, and Hastie Trevor. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009. [PubMed: 19377034]
- [12]. Bodein Antoine, Chapleur Olivier, Droit Arnaud, and Lê Cao Kim-Anh. A generic multivariate framework for the integration of microbiome longitudinal studies with other data types. *bioRxiv*, page 585802, 2019.
- [13]. Friedman Jonathan and Alm Eric J. Inferring correlation networks from genomic survey data. *PLoS computational biology*, 8(9):e1002687, 2012. [PubMed: 23028285]
- [14]. Kurtz Zachary D, Müller Christian L, Miraldi Emily R, Littman Dan R, Blaser Martin J, and Bonneau Richard A. Sparse and compositionally robust inference of microbial ecological networks. *PLoS computational biology*, 11(5):e1004226, 2015. [PubMed: 25950956]
- [15]. Weiss Sophie, Van Treuren Will, Lozupone Catherine, Faust Karoline, Friedman Jonathan, Deng Ye, Xia Li Charlie, Xu Zhenjiang Zech, Ursell Luke, Alm Eric J, et al. Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *The ISME journal*, 10(7):1669, 2016. [PubMed: 26905627]
- [16]. Vandeputte Doris, Kathagen Gunter, Kevin D’hoë Sara Vieira-Silva, Mireia Valles-Colomer João Sabino, Wang Jun, Tito Raul Y, De Commer Lindsey, Darzi Youssef, Vermeire Séverine, Falony Gwen, and Raes Jeroen. Quantitative microbiome profiling links gut community variation to microbial load. *Nature*, 551(7681):507–511, 11 2017. [PubMed: 29143816]
- [17]. Gloor Gregory B, Macklaim Jean M, Pawlowsky-Glahn Vera, and Egozcue Juan J. Microbiome datasets are compositional: And this is not optional. *Front. Microbiol*, 8, 2017. [PubMed: 29118739]
- [18]. Tang Keqi, Page Jason S, and Smith Richard D. Charge competition and the linear dynamic range of detection in electrospray ionization mass spectrometry. *J. Am. Soc. Mass Spectrom*, 15(10): 1416–1423, 2004. [PubMed: 15465354]
- [19]. King Richard, Bonfiglio Ryan, Fernandez-Metzler Carmen, Miller-Stein Cynthia, and Olah Timothy. Mechanistic investigation of ionization suppression in electrospray ionization. *J. Am. Soc. Mass Spectrom*, 11(11):942–950, 2000. [PubMed: 11073257]
- [20]. Matuszewski BK, Constanzer ML, and Chavez-Eng CM. Strategies for the assessment of matrix effect in quantitative bioanalytical methods based on hplc-ms/ms. *Anal. Chem*, 75(13):3019–3030, 2003. [PubMed: 12964746]
- [21]. Alžběta Kalivodová Karel Hron, Filzmoser Peter, Najdekr Lukáš, Janešková Hana, and Adam Tomáš. PIs-da for compositional data with application to metabolomics. *Journal of Chemometrics*, 29(1):21–28, 2015.
- [22]. Jansson JK and Baker ES A multi-omic future for microbiome studies. *Nat Microbiol*, 1:16049, 04 2016. [PubMed: 27572648]
- [23]. Lovell David, Pawlowsky-Glahn Vera, José Egozcue Juan, Marguerat Samuel, and Bähler Jürg. Proportionality: a valid alternative to correlation for relative data. *PLoS computational biology*, 11(3):e1004075, 2015. [PubMed: 25775355]
- [24]. Tipton Laura, Müller Christian L, Kurtz Zachary D, Huang Laurence, Kleerup Eric, Morris Alison, Bonneau Richard, and Ghedin Elodie. Fungi stabilize connectivity in the lung and skin microbial ecosystems. *Microbiome*, 6(1):12, 2018. [PubMed: 29335027]
- [25]. Morton James T, Marotz Clarrise, Washburne Alex, Silverman Justin, Zaramela Livia S, Edlund Anna, Zengler Karsten, and Knight Rob. Establishing microbial composition measurement standards with reference frames. *Nature Communications*, 10(1):2719, 2019.
- [26]. Abadi Martín, Barham Paul, Chen Jianmin, Chen Zhifeng, Davis Andy, Dean Jeffrey, Devin Matthieu, Ghemawat Sanjay, Irving Geoffrey, Isard Michael, et al. Tensorflow: A system for large-scale machine learning. In 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16), pages 265–283, 2016.
- [27]. Quinn Robert A, Comstock William, Zhang Tianyu, Morton James T, da Silva Ricardo, Tran Alda, Aksenov Alexander, Nothias Louis-Felix, Wangpraseurt Daniel, Melnik Alexey V, Ackermann Gail, Conrad Douglas, Klapper Isaac, Knight Rob, and Dorrestein Pieter C. Niche partitioning of a pathogenic microbiome driven by chemical gradients. *Sci Adv*, 4(9):eaau1908, 9 2018. [PubMed: 30263961]

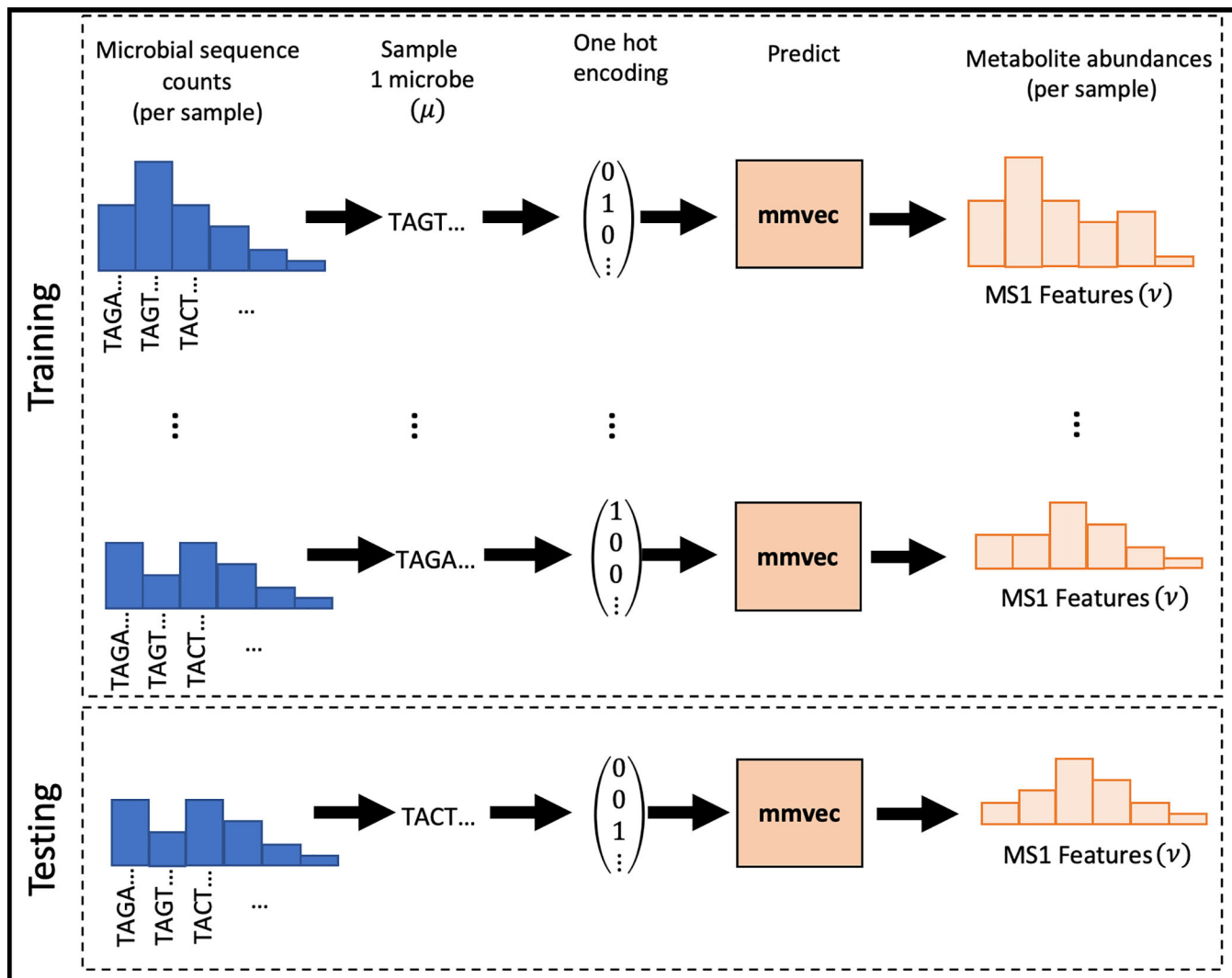
- [28]. Shalpour Shabnam, Lin Xue-Jia, Bastian Ingmar N, Brain John, Burt Alastair D, Aksenov Alexander A, Vrbanac Alison F, Li Weihua, Perkins Andres, Matsutani Takaji, et al. Inflammation-induced iga+ cells dismantle anti-liver cancer immunity. *Nature*, 551(7680):340, 2017. [PubMed: 29144460]
- [29]. Lloyd-Price J, Arze C, Ananthkrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW, Andrews E, Ajami NJ, Bonham KS, Brislawn CJ, Casero D, Courtney H, Gonzalez A, Graeber TG, Hall AB, Lake K, Landers CJ, Mallick H, Plichta DR, Prasad M, Rahnavard G, Sauk J, Shungin D, Vazquez-Baeza Y, White RA, Braun J, Denson LA, Jansson JK, Knight R, Kugathasan S, McGovern DPB, Petrosino JF, Stappenbeck TS, Winter HS, Clish CB, Franzosa EA, Vlamakis H, Xavier RJ, Huttenhower C, Bishai J, Bullock K, Deik A, Dennis C, Kaplan JL, Khalili H, McIver LJ, Moran CJ, Nguyen L, Pierce KA, Schwager R, Sirota-Madi A, Stevens BW, Tan W, Ten Hoeve JJ, Weingart G, Wilson RG, and Yajnik V Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature*, 569(7758):655–662, 5 2019. [PubMed: 31142855]
- [30]. Swenson Tami L, Karaoz Ulas, Swenson Joel M, Bowen Benjamin P, and Northen Trent R. Linking soil biology and chemistry in biological soil crust using isolate exometabolomics. *Nature communications*, 9(1):19, 2018.
- [31]. Baran Richard, Brodie Eoin L, Mayberry-Lewis Jazmine, Hummel Eric, Da Rocha Ulisses Nunes, Chakraborty Romy, Bowen Benjamin P, Karaoz Ulas, Cadillo-Quiroz Hinsby, Garcia-Pichel Ferran, et al. Exometabolite niche partitioning among sympatric soil bacteria. *Nature communications*, 6:8289, 2015.
- [32]. Quinn Robert A, Whiteson Katrine, Lim Yan-Wei, Salamon Peter, Bailey Barbara, Mienardi Simone, Sanchez Savannah E, Blake Don, Conrad Doug, and Rohwer Forest. A winogradsky-based culture system shows an association between microbial fermentation and cystic fibrosis exacerbation. *ISME J*, 9(4):1024–1038, 3 2015. [PubMed: 25514533]
- [33]. Moree Wilna J, Phelan Vanessa V, Wu Cheng-Hsuan, Bandeira Nuno, Cornett Dale S, Duggan Brendan M, and Dorrestein Pieter C. Interkingdom metabolic transformations captured by microbial imaging mass spectrometry. *Proceedings of the National Academy of Sciences*, 109(34):13811–13816, 2012.
- [34]. Wang Mingxun, Carver Jeremy J, Phelan Vanessa V, Sanchez Laura M, Garg Neha, Peng Yao, Nguyen Don Duy, Watrous Jeramie, Kapono Clifford A, Luzzatto-Knaan Tal, Porto Carla, Bouslimani Amina, Melnik Alexey V, Meehan Michael J, Liu Wei-Ting, Crüsemann Max, Boudreau Paul D, Esquenazi Eduardo, Sandoval-Calderón Mario, Kersten Roland D, Pace Laura A, Quinn Robert A, Duncan Katherine R, Hsu Cheng-Chih, Floros Dimitrios J, Gavilan Ronnie G, Kleigrewe Karin, Northen Trent, Dutton Rachel J, Parrot Delphine, Carlson Erin E, Aigle Bertrand, Michelsen Charlotte F, Jelsbak Lars, Sohlenkamp Christian, Pevzner Pavel, Edlund Anna, McLean Jeffrey, Piel Jörn, Murphy Brian T, Gerwick Lena, Liaw Chih-Chuang, Yang Yu-Liang, Humpf Hans-Ulrich, Maansson Maria, Keyzers Robert A, Sims Amy C, Johnson Andrew R, Sidebottom Ashley M, Sedio Brian E, Klitgaard Andreas, Larson Charles B, Cristopher A Boya P, Torres-Mendoza Daniel, Gonzalez David J, Silva Denise B, Marques Lucas M, Demarque Daniel P, Pociute Egle, O'Neill Ellis C, Briand Enora, Helfrich Eric J N, Granatosky Eve A, Glukhov Evgenia, Ryffel Florian, Houson Hailey, Mohimani Hosein, Kharbush Jenan J, Zeng Yi, Vorholt Julia A, Kurita Kenji L, Charusanti Pep, McPhail Kerry L, Nielsen Kristian Fog, Vuong Lisa, Elfeki Maryam, Traxler Matthew F, Engene Niclas, Koyama Nobuhiro, Vining Oliver B, Baric Ralph, Silva Ricardo R, Mascuch Samantha J, Tomasi Sophie, Jenkins Stefan, Macherla Venkat, Hoffman Thomas, Agarwal Vinayak, Williams Philip G, Dai Jingqui, Neupane Ram, Gurr Joshua, Rodríguez Andrés M C, Lamsa Anne, Zhang Chen, Dorrestein Kathleen, Duggan Brendan M, Almaliti Jehad, Allard Pierre-Marie, Phapale Prasad, Nothias Louis-Felix, Alexandrov Theodore, Litaudon Marc, Wolfender Jean-Luc, Kyle Jennifer E, Metz Thomas O, Peryea Tyler, Nguyen Dac-Trung, VanLeer Danielle, Shinn Paul, Jadhav Ajit, Müller Rolf, Waters Katrina M, Shi Wenyuan, Liu Xueting, Zhang Lixin, Knight Rob, Jensen Paul R, Palsson Bernhard O, Pogliano Kit, Linington Roger G, Gutiérrez Marcelino, Lopes Norberto P, Gerwick William H, Moore Bradley S, Dorrestein Pieter C, and Bandeira Nuno. Sharing and community curation of mass spectrometry data with global natural products social molecular networking. *Nat. Biotechnol*, 34(8):828–837, 8 2016. [PubMed: 27504778]

- [35]. Maier Raina Margaret and Soberon-Chavez G. Pseudomonas aeruginosa rhamnolipids: biosynthesis and potential applications. *Applied Microbiology and Biotechnology*, 54(5):625–633, 2000. [PubMed: 11131386]
- [36]. Wood Thammajun L, Gong Ting, Zhu Lei, Miller James, Miller Daniel S, Yin Bei, and Wood Thomas K. Rhamnolipids from pseudomonas aeruginosa disperse the biofilms of sulfate-reducing bacteria. *NPJ biofilms and microbiomes*, 4(1):22, 2018. [PubMed: 30302271]
- [37]. Allen Lucy, Dockrell David H, Pattery Theresa, Lee Daniel G, Cornelis Pierre, Hellewell Paul G, and Whyte Moira KB. Pyocyanin production by pseudomonas aeruginosa induces neutrophil apoptosis and impairs neutrophil-mediated host defenses in vivo. *The Journal of Immunology*, 174(6):3643–3649, 2005. [PubMed: 15749902]
- [38]. Quinn Robert A., Vrbanac Alison, Melnik Alexey V., Patras Kathryn A., Christy Mitchell, Nelson Andrew T., Aksenov Alexander, Tripathi Anupriya, Humphrey Greg, da Silva Ricardo, Bussell Robert, Thron Taren, Wang Mingxun, Vargas Fernando, Gauglitz Julia M., Meehan Michael J., Poulsen Orit, Boland Brigid S., Chang John T., Sandborn William J., Lim Meerana, Garg Neha, Lumeng Julie, Kazmierczak Barbara I., Jain Ruchi, Egan Marie, Rhee Kyung E., Haddad Gabriel G., Siegel Dionicio, Mazmanian Sarkis, Nizet Victor, Knight Rob, and Dorrestein Pieter C. Chemical impacts of the microbiome across scales reveal novel conjugated bile acids. *bioRxiv*, 2019.
- [39]. Paczosa Michelle K. and Meccas Joan. *Klebsiella pneumoniae*: Going on the offense with a strong defense. *Microbiology and Molecular Biology Reviews*, 80(3):629–661, 2016. [PubMed: 27307579]
- [40]. Tiratterra Elisa, Franco Placido, Porru Emanuele, Katsanos Konstantinos H, Christodoulou Dimitrios K, and Roda Giulia. Role of bile acids in inflammatory bowel disease. *Annals of gastroenterology*, 31(3):266, 2018. [PubMed: 29720851]
- [41]. Hofmann Alan F and Eckmann Lars. How bile acids confer gut mucosal protection against bacteria. *Proceedings of the National Academy of Sciences*, 103(12):4333–4334, 2006.
- [42]. Begley Máire, Gahan Cormac GM, and Hill Colin. The interaction between bacteria and bile. *FEMS microbiology reviews*, 29(4):625–651, 2005. [PubMed: 16102595]
- [43]. Okada Y, Tsuzuki Y, Miyazaki J, Matsuzaki K, Hokari R, Komoto S, Kato S, Kawaguchi A, Nagao S, Itoh K, Watanabe T, and Miura S *Propionibacterium freudenreichii* component 1.4-dihydroxy-2-naphthoic acid (DHNA) attenuates dextran sodium sulphate induced colitis by modulation of bacterial flora and lymphocyte homing. *Gut*, 55(5):681–688, 5 2006. [PubMed: 16299037]
- [44]. Foligne B, Parayre S, Cheddani R, Famelart MH, Madec MN, Ple C, Breton J, Dewulf J, Jan G, and Deutsch SM Immunomodulation properties of multi-species fermented milks. *Food Microbiol*, 53(Pt A):60–69, 2 2016. [PubMed: 26611170]
- [45]. Ple C, Breton J, Richoux R, Nurdin M, Deutsch SM, Falentin H, Herve C, Chuat V, Lemee R, Maguin E, Jan G, Van de Guchte M, and Foligne B Combining selected immunomodulatory *Propionibacterium freudenreichii* and *Lactobacillus delbrueckii* strains: Reverse engineering development of an anti-inflammatory cheese. *Mol Nutr Food Res*, 60(4):935–948, 4 2016. [PubMed: 26640113]
- [46]. Colliou Natacha, Ge Yong, Sahay Bikash, Gong Minghao, Zadeh Mojgan, Owen Jennifer L, Neu Josef, Farmerie William G, Alonzo Francis, Liu Ken, et al. Commensal *propionibacterium* strain *ufl* mitigates intestinal inflammation via *th17* cell regulation. *The Journal of clinical investigation*, 127(11):3970–3986, 2017. [PubMed: 28945202]

## 11 Methods only references

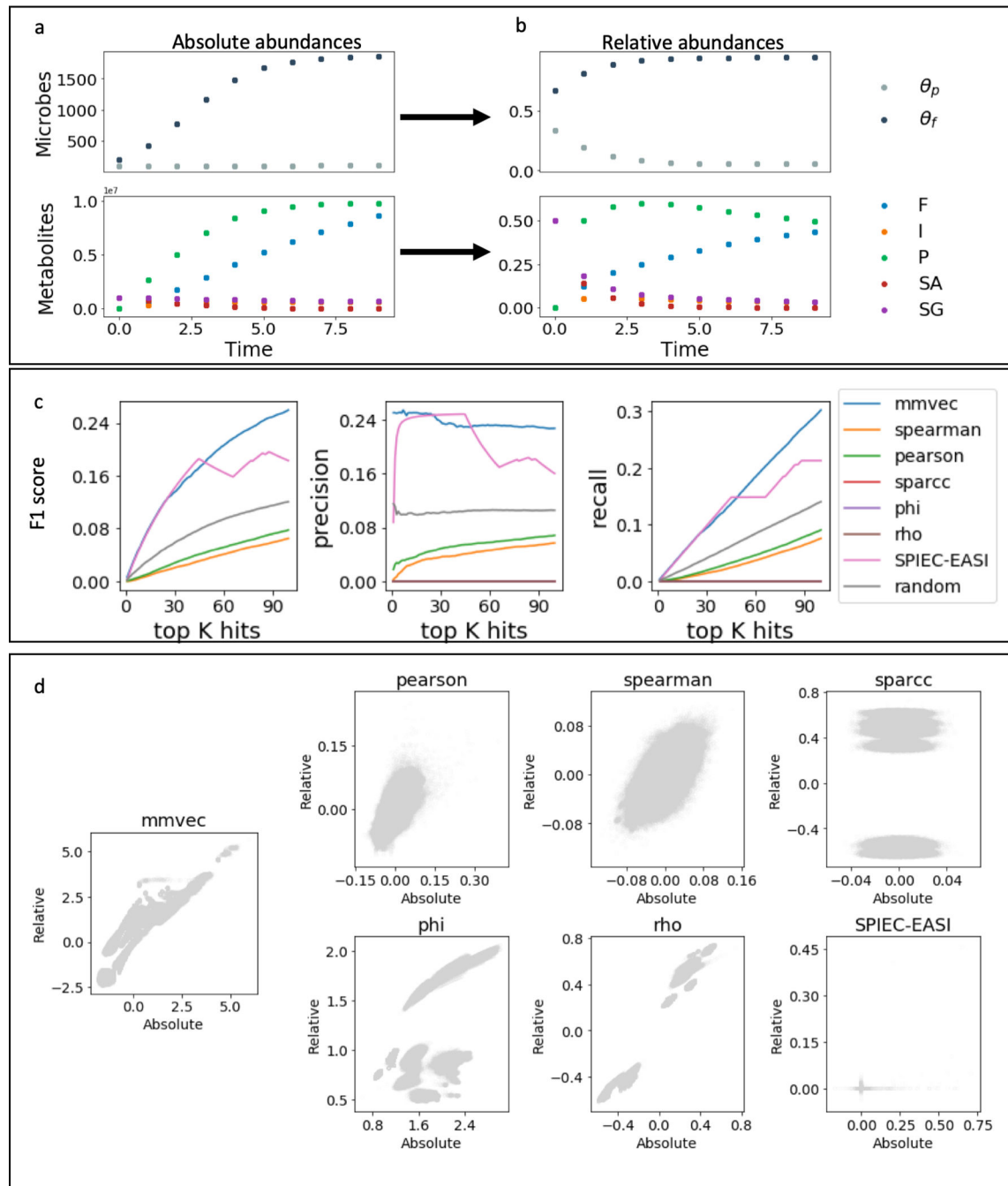
- [47]. Nasrabadi Nasser M. Pattern recognition and machine learning. *Journal of electronic imaging*, 16(4):049901, 2007.
- [48]. Pawlowsky-Glahn Vera, Egozcue Juan José, and Tolosana-Delgado Raimon. *Modeling and Analysis of Compositional Data*. John Wiley & Sons, 2 2015.
- [49]. Mikolov Tomas and Sutskever Ilya and Chen Kai and Corrado Greg S and Dean Jeff. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.

- [50]. Koren Yehuda, Bell Robert, and Volinsky Chris. Matrix factorization techniques for recommender systems. *Computer*, (8):30–37, 2009.
- [51]. Kingma Diederik P and Ba Jimmy. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [52]. Blei David M, Ng Andrew Y, and Jordan Michael I. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- [53]. Sankaran Kris and Holmes Susan P. Latent variable modeling for the microbiome. arXiv preprint arXiv:1706.04969, 2017.
- [54]. Aitchison John and Greenacre Michael. Biplots of compositional data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 51(4):375–392, 2002.
- [55]. Aitchison John, Ng KW, et al. Conditional compositional biplots: theory and application. 2005.
- [56]. Martín-Fernández JA, Pawłowsky-Glahn V, Egozcue JJ, and Tolosona-Delgado R. Advances in principal balances for compositional data. *Mathematical Geosciences*, 50(3):273–298, 2018.
- [57]. Bolyen Evan, Rideout Jai Ram, Dillon Matthew R, Bokulich Nicholas A, Abnet Christian, Al-Ghalith Gabriel A, Alexander Harriet, Alm Eric J, Arumugam Manimozhiyan, Asnicar Francesco, et al. Qiime 2: Reproducible, interactive, scalable, and extensible microbiome data science. Technical report, PeerJ Preprints, 2018.
- [58]. Yoshiki Vázquez-Baeza Meg Pirrung, Gonzalez Antonio, and Knight Rob. Emperor: a tool for visualizing high-throughput microbial community data. *Gigascience*, 2(1):16, 2013. [PubMed: 24280061]

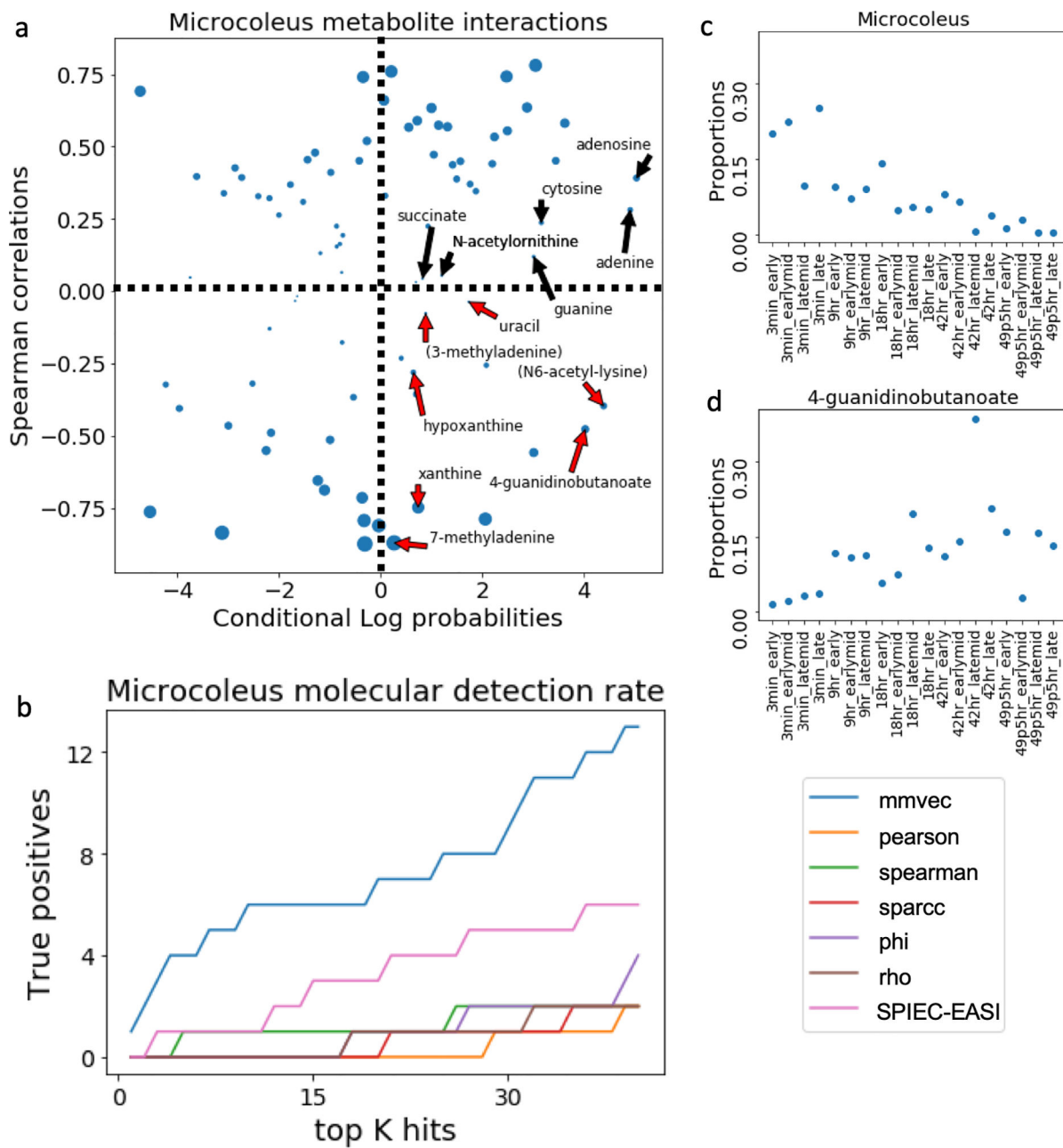


**Figure 1:**

Input data types and mmvec neural network architecture. (a) The neural network architecture where the input layer represents one-hot encodings of  $N$  microbes and the output layer represents the proportions of  $M$  metabolites.  $U$  corresponds to microbial vectors and  $V$  corresponds to metabolite vectors. (b) The pipeline for training mmvec. The objective behind mmvec is to predict metabolite abundances ( $y$ ) given a single input microbe sequence ( $x$ ), also known as a one-hot encoding. This training procedure will estimate conditional probabilities of observing a metabolite given the input microbe sequence. Cross-validation can be performed on hold-out samples to access overfitting.

**Figure 2:**

Simulation benchmarks. (a) Absolute abundances of microbes and metabolites simulated from differential equations derived in [27] for a specific spatial point. (b) Proportions of the abundances shown in (a). (c) F1 score, precision and recall curves comparing mmvec to Pearson, Spearman, SparCC, SPIEC-EASI, and proportionality metrics phi and rho across the top 100 metabolites for each microbe. (d) comparisons of coefficients learned from absolute abundances and relative abundances all of the benchmarked methods.

**Figure 3:**

*M. vaginatus* released metabolites after the biocrust wetting event. (a) Comparison of *M. vaginatus* metabolite interactions estimated from Spearman and mmvec from (n=19 samples). All of the experimentally validated *M. vaginatus* released metabolites are labeled. All metabolites with contradicting findings between the wetting experiment and the *in vitro* experimental results are highlighted in red. Points are resized according to the  $-10 \log(p\text{-value})$  obtained from Spearman correlation. Dashlines mark the cutoff for a Spearman correlation of zero, and the conditional log probabilities of zero. Here a zero log conditional probability represents the conditional probability of the average metabolite because all probabilities here are mean centered. (b) Benchmarks comparing the detection rate of the



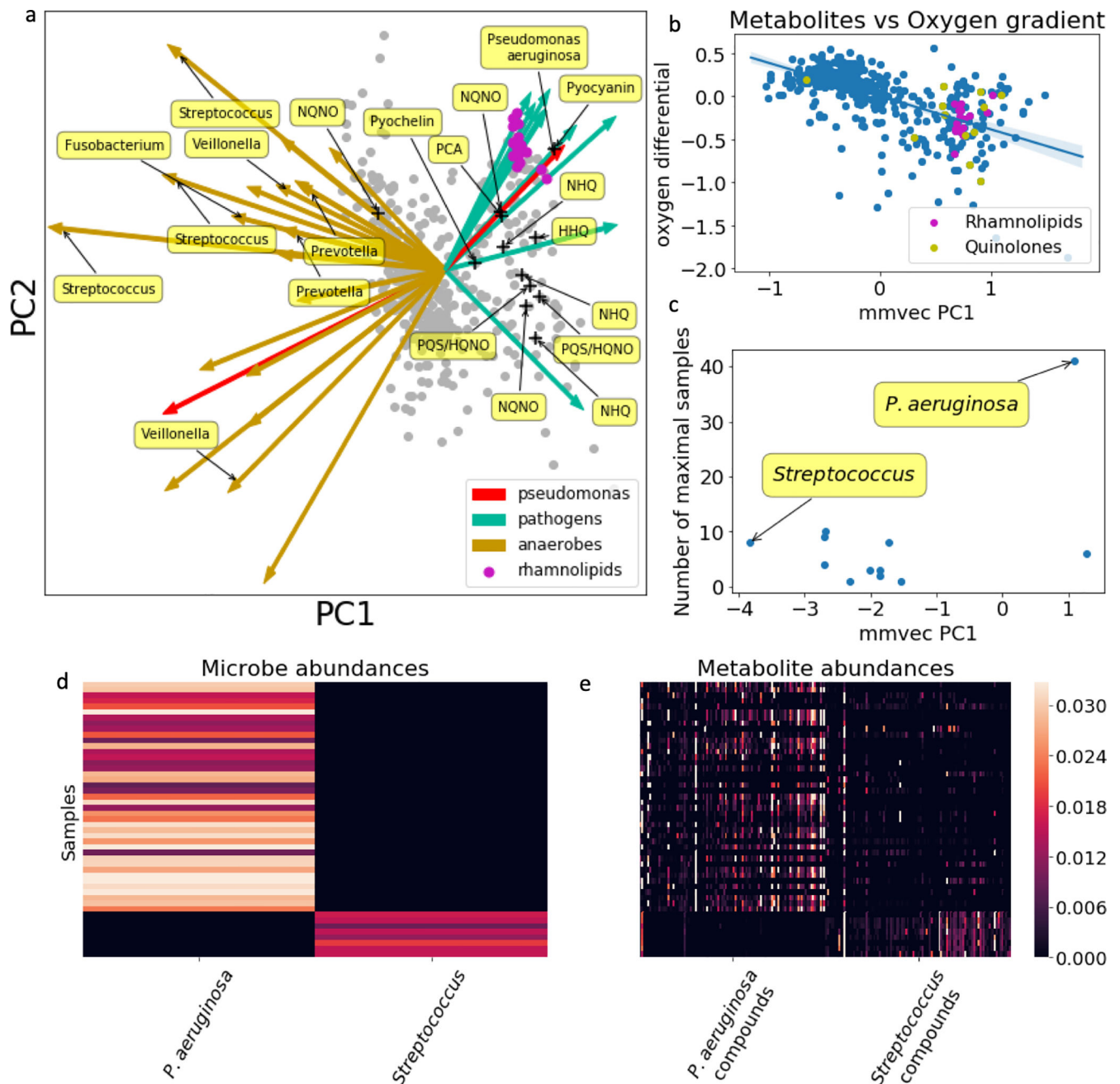
experimentally validated molecules across different statistical methodologies. (c) *M. vaginatus* proportions and (d) 4-guanidinobutanoate proportions following a wetting event.

Author Manuscript

Author Manuscript

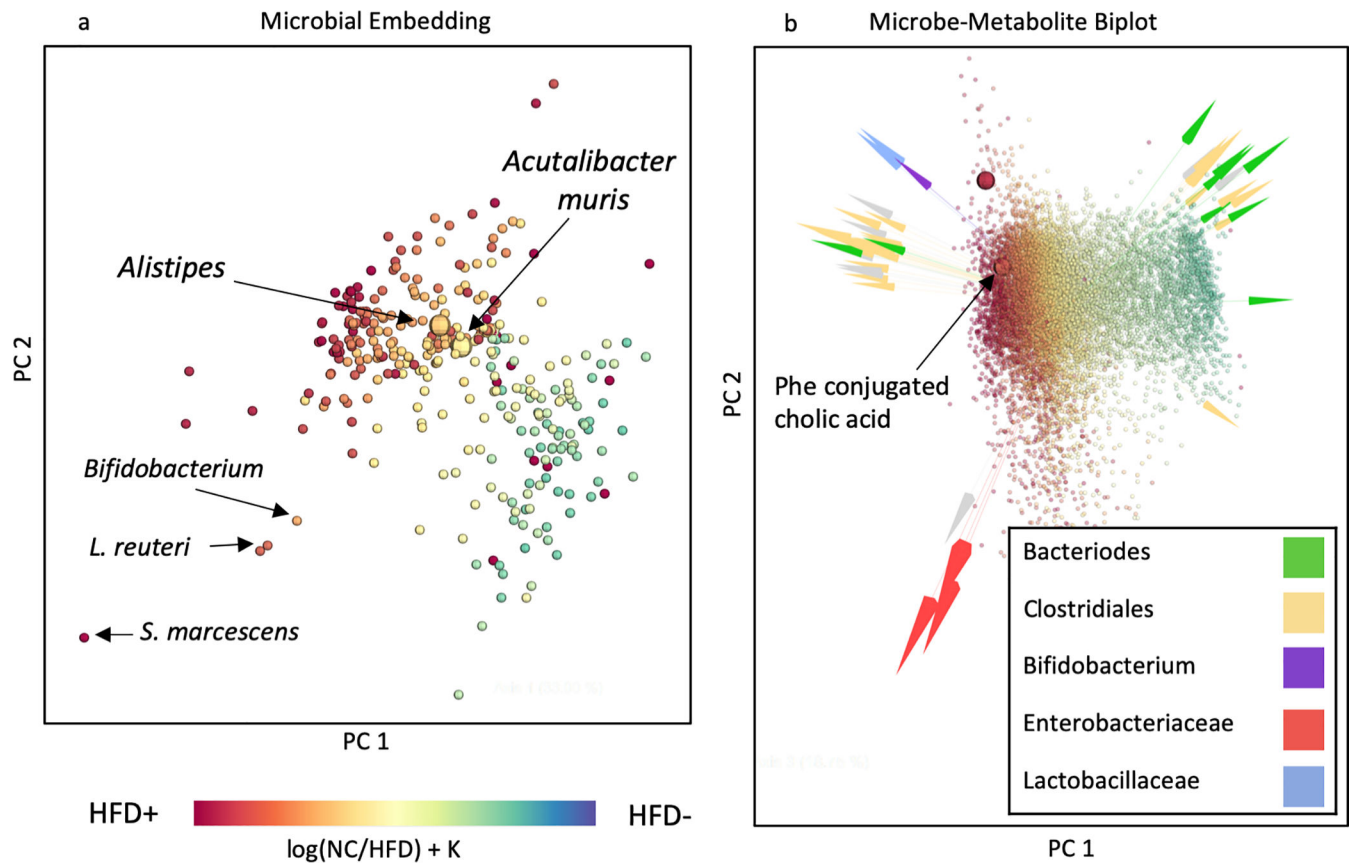
Author Manuscript

Author Manuscript

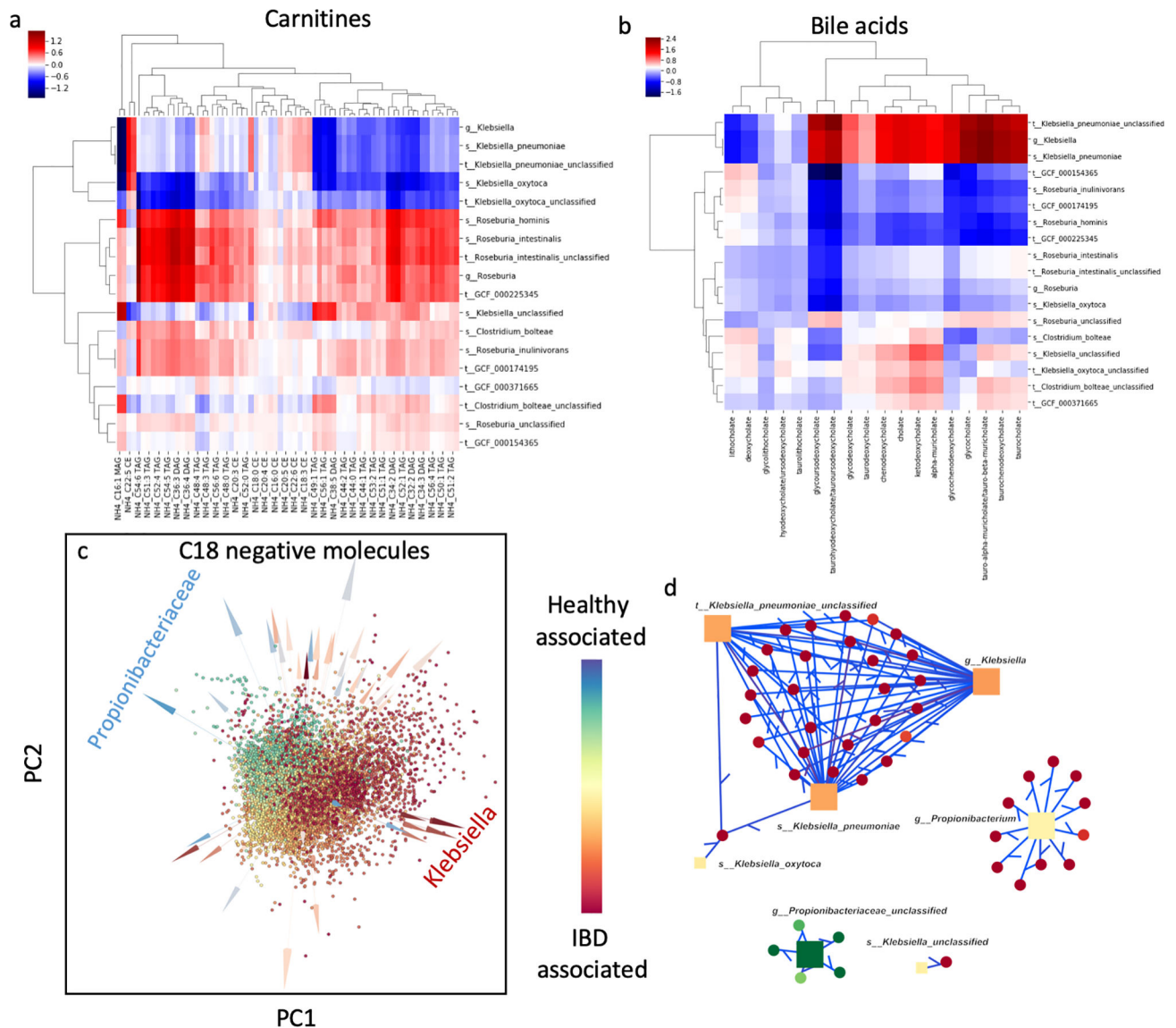


**Figure 4:** Investigation of *P.aeruginosa*-associated molecules. (a) Biplot drawn from the mmvec conditional probabilities estimated for the cystic fibrosis dataset [27]. Arrows represent microbes and dots represent metabolites. The x and y axes represent principal components from the SVD of the microbe-metabolite conditional probabilities estimated from mmvec ( $n=138$  samples). Distances between points quantify co-occurrence strength between metabolites, with small distances indicating metabolites that have a high probability of co-occurring with high probability. Distances between arrow tips quantify co-occurrence strength between microbes. The directionality of the arrows can be used to pinpoint which

microbes can explain the metabolite co-occurrence patterns. Arrows highlighted in green correspond to putative cystic fibrosis pathogens and yellow arrows highlight known anaerobes. Only known molecules produced by *P. aeruginosa* are labeled. (b) Scatter plot of molecules with respect to the oxygen gradient differential and the first principal component learned from mmvec (n=442 molecules) with linear regression model and 95% confidence interval for regression estimate. (c) The first principal component vs the number of samples where the taxa was the most abundant taxa in that sample. (d) Heatmap of *P. aeruginosa* and *Streptococcus* abundances in samples where they are the most abundant species. (e) Heatmap of the top 100 molecules that co-occur with *P. aeruginosa* and *Streptococcus*.



**Figure 5:** Microbe/metabolite co-occurrences across study of HCC progression in the context of innate immunity in a mouse model [28]. (a) Visualization of microbial co-occurrence patterns, where distances between points approximates the Aitchison distance between microbes, which quantities microbial occurrences. Small distances are indicative of microbes with high probability of co-occurring together. Microbes are colored according to their association with HFD, which was estimated using differential abundance analysis via multinomial regression. (b) Emperor [59] biplot of microbe-metabolite interactions, with metabolites colored according to their association with HFD. HFD association was estimated through differential abundance analysis via multinomial regression. Distances between points approximate Aitchison distances between metabolites and distances between arrow tips approximate Aitchison distances between microbes. Several *Clostridium spp.* appear to co-occur with the new bile acid molecule cholate phenylalanine amidate, also referred to as Phe conjugated cholic acid.



**Figure 6:** Microbe-metabolite interactions of the human microbiome in association with IBD samples [29]. (a) Heatmap visualization of the inferred conditional probabilities for various bile acids given the presence of *Klebsiella*, *Roseburia* and *Clostridium bolteae*. (b) Heatmap visualization of the inferred conditional probabilities for the carnitines given the presence of *Klebsiella*, *Roseburia*, and *Clostridium bolteae*. (c) Multiomics biplot of the microbe-metabolite interactions learned from metagenomics profiles and C18 negative ion mode LC-MS. Microbes (arrows) and metabolites (spheres) are colored according to their differentials estimated from multinomial regression. *Klebsiella spp.* appears to be strongly associated with IBD, while *Propionibacterium spp.* has strong negative association. (d) Network of the top 300 edges where only the edges that contain *Klebsiella* and *Propionibacteriaceae* are visualized.