

UC San Diego

UC San Diego Previously Published Works

Title

Learning representations of microbe–metabolite interactions

Permalink

<https://escholarship.org/uc/item/3wd3p1tm>

Journal

Nature Methods, 16(12)

ISSN

1548-7091

Authors

Morton, James T

Aksenov, Alexander A

Nothias, Louis Felix

et al.

Publication Date

2019-12-01

DOI

10.1038/s41592-019-0616-3

Supplemental Material

<https://escholarship.org/uc/item/3wd3p1tm#supplemental>

Peer reviewed

# 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>1</sup> Mingxun Wang,<sup>3,4</sup> Aaron Watters,<sup>10</sup> Se Jin Song,<sup>1,11</sup> Richard Bonneau,<sup>7,10,12,13</sup> Pieter C. Dorrestein,<sup>3,4</sup> and Rob Knight<sup>1,2,14,11</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>*Flatiron Institute, Simons Foundation, New York, 10010 NY, USA*

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

<sup>12</sup>*Computer Science Department, Courant Institute, New York, 10012 NY, USA*

<sup>13</sup>*Center For Data Science, NYU, New York, NY 10008, USA*

<sup>14</sup>*Department of Bioengineering University of California, San Diego, La Jolla, CA, USA*

## I. 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.

## II. 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]. 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 [2]. Relative abundances of thousands of microbes and metabolites can be measured using sequencing and mass spectrometry, which can result in the generation of very high dimensional microbiome and metabolomics datasets. Quantifying microbe-metabolite interactions 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 [3–5]. Network models have been shown to improve classification accuracy using multiple datasets [6]. Factor models have been proposed to incorporate multiple datasets for biomarker analysis [7]. Despite of the wide application of these methods, they are notoriously difficult to interpret [8–10] and it remains unclear whether these models can obtain individual microbe-metabolite interactions.

57 Pearson and Spearman correlations assume independence between interactions, simplify-  
58 ing the estimation procedure by reducing it to a combination of independent two dimensional  
59 problems. However, many studies have shown that these methods are not statistically valid  
60 for compositional data, a fact first recognized by Pearson in 1895 and followed up in nu-  
61 merous studies [11–15]. This problem is further complicated because both microbiome [15]  
62 and mass spectrometry [16–19] datasets are also compositional, meaning that the absolute  
63 abundances are not measured, which can confound statistical inference. For example, in  
64 untargeted mass spectrometry experiments, the set of molecules detected and their relative  
65 abundance vary depending on the extraction protocol and analytic methods used, which  
66 leads to only a partial snapshot of the metabolome. Moreover, measuring the total mass of  
67 molecules extracted is often not performed in large scale metabolomics efforts, due to the  
68 highly laborious nature of that step.

69 To understand how issues associated with compositional data impact inference on  
70 microbe-metabolite interactions, consider the example in Figure S1. There are two mi-  
71 crobes and two metabolites in Figure S1a. All are increasing exponentially at different rates  
72 and are highly correlated with each other. If proportions are estimated from the absolute  
73 abundances via sampling, the information about the total microbe population size and  
74 the total metabolite abundance is lost, and the correlations between the microbes and the  
75 metabolites disappear. False positives can also appear as shown in Figure S1b, microbe and  
76 metabolite interactions that have no apparent correlation structure may appear to be corre-  
77 lated when investigating the proportions. These issues alone can give rise to overwhelming  
78 false positives and false negatives, making Pearson and Spearman in some scenarios com-  
79 parable to random coin flips. Experimental validation currently takes large laboratories  
80 multiple years to perform [20], often requiring time-consuming manual examinations of  
81 erroneous correlations.

82 There are other compositional techniques such as SparCC[11] and proportionality[21] that  
83 are scale-invariant when analyzing a single dataset, but lose scale-invariance when analyzing  
84 multiomics datasets. This was shown in the context of identifying microbe-fungal inter-  
85 actions [22], which provided motivation to extend SPIEC-EASI [12] to handle multiomics  
86 datasets. We show that this approach does not work for microbe-metabolite interactions  
87 because of differences of measurement units between sequencing and mass spectrometry mea-  
88 surements (Supplementary materials). An alternative approach is to consider co-occurrence

89 probabilities instead of correlations. Here, co-occurrence probabilities refer to the condi-  
 90 tional probability of observing a metabolite given that a microbe was observed, thereby  
 91 allowing us to identify the most likely microbe-metabolite interactions. To do this, we  
 92 propose “mmvec”, (microbe-metabolite vectors), to learn these co-occurrence probabilities  
 93 between microbes and metabolites. Due to its scale-robustness properties, the microbial-  
 94 metabolite relationships learned by mmvec are consistent between the absolute and relative  
 95 abundances. The microbe-metabolite interactions can be ranked [23] and visualized through  
 96 standard dimensionality reduction interfaces, enabling interpretable findings. The compu-  
 97 tations behind mmvec can take advantage of modern GPU architectures using Tensorflow  
 98 [24], enabling scalable inference on large multiomics datasets. Furthermore, we provide evi-  
 99 dence in two benchmarks and four case studies that mmvec outperforms existing statistical  
 100 methods.

### 101 III. RESULTS AND DISCUSSION

102 We performed benchmarks comparing mmvec to Pearson, Spearman, SPIEC-EASI,  
 103 SparCC and Proportionality [21] using a cystic fibrosis biofilm simulation. We then show  
 104 that mmvec can resolve contradictory cyanobacteria-metabolite relationships in a desert  
 105 soil biocrust wetting study. We also demonstrate recovery of known associations of *P.*  
 106 *aeruginosa*-produced metabolites observed in cystic fibrosis [25]. Finally, we explore the  
 107 relationships of microbiota and metabolic changes in mice fed a high fat diet [26] and inflam-  
 108 matory bowel disease [27], showing how this approach can be used to determine microbial  
 109 origin of novel molecules even in extremely complex real-life biological systems with limited  
 110 knowledge of existing associations.

#### 111 A. Simulation benchmarks

112 To compare mmvec performance to Pearson, Spearman, Proportionality, SparCC and  
 113 SPIEC-EASI correlations, we used data from existing studies in which the relationships  
 114 between microbes and metabolites were the central focus of investigation. One such study  
 115 simulated spatial-temporal dynamics in a microbial biofilm [25]. The original study tested  
 116 the hypothesis that the cystic fibrosis (CF) microbiome community within human lungs can

117 be manipulated by altering its chemical environment. Changes in pH and oxygen saturation  
 118 suppress the principal pathogen, *P. aeruginosa*, without using antibiotics, by promoting the  
 119 growth of a community of fermenters that out-compete the pathogen. The simplicity of this  
 120 system allowed the high-level ecological patterns to be modelled. In the original simulations,  
 121 the interactions between two microbes (fermenters denoted by  $\theta_f$  and *P. aeruginosa* denoted  
 122 by  $\theta_p$ ) and multiple molecules were modeled using Monod kinetics and diffusion processes[25]  
 123 (Figure 2a).

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

132 Therefore, we can test whether the top  $K$  metabolites associated with each microbe by  
 133 each tool includes the correct microbe-metabolite interactions. Figure 2c shows specificity  
 134 and sensitivity for each tools as a function of  $K$ . In these simulations, random chance  
 135 outperformed all of the tools except for mmvec and SPIEC-EASI, with mmvec performing  
 136 the best. As shown in Figure 2d and Figure S2, mmvec is the only method robust to scale  
 137 deviations amongst the methods tested. This is critical for maintaining consistency between  
 138 absolute and relative abundances, which can otherwise lead to inflated false positives and  
 139 false negatives [14].

## 140 B. Soil biocrust wetting event case study

141 Many studies produce inconsistent results that can be resolved with improved data anal-  
 142 ysis, especially in environmental and clinical settings. To test whether mmvec can resolve  
 143 unexplained discrepancies in microbe-metabolite interactions across studies, we applied it  
 144 to a study of biocrust wetting [28]. In this study, laboratory-based exometabolite patterns  
 145 observed with bacterial isolates were reproduced in the environment. Specifically, in this  
 146 work authors identified metabolites that were consumed and released by multiple biocrust

147 isolates including *Microcoleus vaginatus* and two *Bacillus* strains [29], and compared these  
148 patterns with closely-related environmental taxa and metabolites observed in situ [28].

149 While almost 70% of the examined microbe-metabolite relationships following the wetting  
150 event were validated [28], some contradicted the microbe-metabolite relationships observed  
151 in cultures [29]. These contradictions stemmed from Spearman correlations between *M.*  
152 *vaginatus* abundances and the observed metabolite abundances, but were resolved by mmvec  
153 (Figure 3a).

154 All metabolites released from the *M. vaginatus* isolate have higher conditional proba-  
155 bilities than the average metabolite following biocrust wetting, and are among the top 80  
156 co-occurring metabolites with *M. vaginatus* (of 485 molecules total). This result supports  
157 the original finding that *M. vaginatus* actually releases these molecules after the wetting  
158 event. In contrast, Spearman labels 7 of 13 of these molecules with a negative correlation,  
159 indicating that these molecules were consumed by *M. vaginatus* rather than released, as  
160 originally stated in [28]. When the annotation detection rates amongst different statistical  
161 methodologies, mmvec has a substantially higher true positive rate as shown in Figure 3b.

162 The conflicting results between mmvec and Spearman could be explained by the growing  
163 microbial biomass and shift in available resources after wetting (Figure 3 c, d). Total biomass  
164 is expected to increase, because *M. vaginatus* releases metabolites that enable the growth of  
165 many other microbes. Because DNA sequencing can only measure proportions, the growth  
166 in other microbes could cause the proportions of *M. vaginatus* to decrease, leading to a mis-  
167 leading anti-correlation with 4-guanidinobutanoate (Figure 3d). However, it is not possible  
168 to infer whether *M. vaginatus* is decreasing in abundance [23] or 4-guanidinobutanoate is  
169 increasing in abundance.

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

### C. Cystic Fibrosis case study

175

176 To further validate if mmvec can detect known microbe-metabolite interactions in a bi-  
177 ological setting, we re-analyzed a study on lung mucus microbiome of patients with cystic  
178 fibrosis [25, 30]. Cystic fibrosis has been shown to be dominated by two major groups  
179 of microbes, anaerobes and pathogens that occupy unique niches, and their interactions  
180 are defined by the environment. Anaerobes dominate in low oxygen and low pH environ-  
181 ments, while pathogens, in particular *P. aeruginosa*, dominate in the opposite conditions  
182 [25]. Mmvec clearly separates anaerobes and pathogens (Figure 4a), with known anaerobic  
183 microbes (*Veillonella*, *Fusobacterium*, *Prevotella* and *Streptococcus*) on the left, and notable  
184 pathogens, such as *P. aeruginosa*, on the right.

185 *P. aeruginosa* is known to produce small-molecule virulence factors [31]. In the origi-  
186 nal study, based on annotations from GNPS[32], the bacterium was found to produce six  
187 molecules: 4-hydroxy-2-heptylquinoline (HHQ), Pyocyanin (PYO), Phenazine-1-carboxylic  
188 acid (PCA), 2-nonyl-4-hydroxy-quinoline (NHQ), 2-heptyl-3,4-dihydroxyquinoline (PQS,  
189 *Pseudomonas* quinolone signal) and Pyochelin [25]. As shown in Figure 4a, mmvec identi-  
190 fies these molecules with a high co-occurrence probability with *P. aeruginosa*. Mmvec also  
191 identifies a cluster of rhamnolipids likely produced by *P. aeruginosa*. Rhamnolipids are  
192 well characterized and are an important virulence factors for *P. aeruginosa*, contributing to  
193 biofilm development, motility on surfaces and antagonistic interactions with host inflamma-  
194 tory cells [33, 34]. These rhamnolipids were not identified in the original study [25]. The  
195 annotations for these compounds have been established using GNPS [32].

196 There is a negative correlation between the first principal component learned from mmvec  
197 and the metabolites log-fold change across the oxygen gradient (Figure 4b) (Pearson  $r=-$   
198  $0.59$ ,  $p\text{-value } 1.8 \times 10^{-44}$ ), which is consistent with the findings in the original work. No  
199 such correlation between the oxygen gradient and the first microbial principal component  
200 was found by Pearson ( $r=0.01$ ,  $p=0.89$ ). There exist two notable microbes on opposing  
201 ends of the first microbial principal component: *P. aeruginosa*, a known pathogen, and  
202 *Streptococcus*, a known anaerobe. The top 100 metabolites that are specific to *P. aerug-*  
203 *inosa* and *Streptococcus* are shown to have drastically different profiles in samples where  
204 *P. aeruginosa* and *Streptococcus* were the most abundant species (Figure 4d,e) (logratio  
205  $t\text{-test}=6.51$ ,  $p=4.4 \times 10^{-8}$ ). This provides evidence that in the context of this study, the



206 metabolomic profiles can be largely influenced by the most abundant microbes, a notion  
207 that has important implications for understanding CF etiology. To further support this, the  
208 learned metabolite conditional probabilities for *P. aeruginosa* can be used to predict the  
209 metabolite proportions in the 41 samples where *P. aeruginosa* is the most abundant taxa.  
210 The predicted *P. aeruginosa* metabolite profiles alone can explain 10% of the metabolite  
211 variation in these samples ( $r=0.319$ ,  $p=1.18 \times 10^{-11}$ ).

212 Of 14 quinolone molecules known to be produced by *P. aeruginosa*, Pearson correla-  
213 tion detected 9 with  $p < 0.05$  without FDR correction, and only 5 with FDR correction.  
214 For example, Pyocyanin, does not appear related to *P. aeruginosa* by the raw proportions  
215 ( $r=0.158$ , FDR-corrected  $p$ value= $0.089$ , rank= $96$ ), but is ranked 34th most associated with  
216 *P. aeruginosa* by mmvec (Figure S3c), consistent with culturing experiments that demon-  
217 strate that *P. aeruginosa* produces this molecule [35]. 18 rhamnolipids are among the top  
218 25 metabolites most associated with *P. aeruginosa* by mmvec, and have higher ranks with  
219 mmvec than with Pearson correlation (Figure S3b).

#### 220 D. Effects of high fat diet in murine model case study

221 We then tested whether mmvec could determine the microbial origin of specific molecules  
222 in a complex biological system. We recently discovered a new kind of bile acid, where  
223 cholate is conjugated to amino acids other than glycine and taurine [36]. These molecules  
224 increased in abundance with high-fat diet in humans. We determined that these molecules  
225 are microbially-made since they were present in specific pathogen free, but not in germ free  
226 mice. We therefore set out to identify candidate producers. We were able to confirm that  
227 one of these bile acids, cholate phenylalanine amidate, was associated with high-fat diet in  
228 well-controlled study that investigated the development of non-alcoholic fatty liver disease  
229 (NAFLD), cirrhosis, and hepatocarcinoma (HCC) in a mouse model [26]. When re-analyzing  
230 these datasets for differential abundances via multinomial regression, the strong association  
231 of the novel bile acid with HFD became immediately apparent. The use of mmvec showed  
232 distinct associated groups of microbes and HFD (Figure 5a) and a clear stratification of the  
233 mass spectrometry data according to diet (Figure 5b). Several *Clostridium spp.* correlated  
234 with the cholate phenylalanine conjugate. Indeed, we showed that *Clostridium spp.* were  
235 found to produce this bile acid [36]. This result demonstrates mmvec’s ability to streamline

236 the discovery of microbes that produce specific molecules of interest *in vivo*.

### 237 **E. Microbe-metabolite interactions in Inflammatory Bowel Disease**

238 Finally, microbe-metabolite interactions were investigated for samples of IBD patients  
239 generated under the integrative Human Microbiome Project [27]. The role of the microbiome  
240 in IBD is acknowledged, but still poorly understood. The original study uncovered shifts in  
241 metabolomic and microbial profiles associated with the IBD. In particular, levels of carnitines  
242 and bile acids were shown to be affected [27]. Using mmvec we confirmed the core findings in  
243 the previous study, such as the co-occurrence between *R. hominis* and multiple carnitines,  
244 including previously noted C20, which have anti-inflammatory properties (Figure 6a) [27].  
245 We also found high correlation of *Klebsiella spp.* with IBD status and that it co-occurs  
246 with high probability with several bile acids (Figure 6b). Although *Klebsiella* itself does not  
247 produce these compounds, some pathogens (including *Klebsiella*) are known to be resistant  
248 to bile acids [37]. Excessive production of some bile acids and bile acid malabsorption  
249 can lead to overabundance of bile acids, which is a hallmark of IBD [38], although the  
250 exact mechanisms remain unknown. The ability of *Klebsiella* to thrive in concentrated bile  
251 acid environments is consistent with the high co-occurrence probabilities shown in Figure  
252 6b. We also noted that three *Klebsiella* species are the top drivers of the IBD- associated  
253 molecules (Figure 6c). It is important to delineate different reasons for co-occurrence. Unlike  
254 *Klebsiella*, *Clostridium* species are known for bile acid manipulation, including production  
255 of bile that can germinate *Clostridium difficile* spores or that have anti-microbial properties  
256 [39, 40].

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

260 In addition to recapitulating reported findings, mmvec also yielded previously undetected  
261 relationships. The major microbe that was found to be associated with healthy patients is  
262 *Propionibacteriaceae*, which was not detected in Price et al 2019 (Figure 6cd). This relation-  
263 ship is corroborated by other published studies. In one study, it has been shown that some  
264 members of the *Propionibacterium* genus produce 1,4-Dihydroxy-2-naphthoic acid (DHNA),  
265 a growth stimulator for bacteria such as *Bifidobacterium* that are thought to reduce the

266 symptoms of IBD [41]. Also, in a survey of *in vivo* vs. *in vitro* bacterial activity, *Probion-*  
 267 *ibacterium freudenreichii* was shown to play an immunomodulatory role in the context of  
 268 an ulcerative colitis mice model [42]. In another study it was shown that *Propionibacterium*  
 269 *freudenreichii* is a viable core component in an anti-inflammatory probiotic fermented dairy  
 270 product [43]. The members of this family have been considered beneficial for intestinal im-  
 271 munoregulation; *Propionibacteriaceae* have been observed to be enriched in human breast  
 272 milk and have been shown to restore Th17 differentiation [44]. Thus, it appears that the  
 273 existing knowledge supports the statistically-inferred interaction uncovered by mmvec, but  
 274 not identified in the published analysis of the dataset

#### 275 IV. CONCLUSION

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

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

293 The concepts outlined here should generalize beyond microbe-metabolite interactions to  
 294 handle other paired multi-omic data types, provided that the input dataset is made up of  
 295 counts (as in metagenomics, transcriptomics, etc.). With the exponential growth of multi-

296 omics datasets, there is much potential to use these methods to reveal microbial metabolism,  
297 including for microbes that are not cultivable in the laboratory. Approaches utilizing co-  
298 occurrence probabilities have the potential to enable more targeted experimental assays,  
299 accelerating the discovery of microbe-metabolite interactions, paving the way towards new  
300 ecosystems engineering approaches in clinical, environmental and industrial applications.

## 301 V. ACKNOWLEDGEMENTS

302 We would like to thank Vera Pawlowsky, Juan Jose Egozcue and Susan Holmes for their  
303 insights behind the geometry of this neural network model. T.L.S., M.W.V.G and T.R.N  
304 greatly acknowledge funding from the Office Science Early Career Research Program, Office  
305 of Biological and Environmental Research, of the U.S. Department of Energy under contract  
306 number DE-AC02-05CH11231 to Lawrence Berkeley National Laboratory. This was in part  
307 supported by P41GM103484 Center for Computational Mass Spectrometry, Instrument sup-  
308 port through NIH S10RR029121 and R03 CA211211 on reuse of metabolomics data. Y. V.  
309 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

310

311

## VI. AUTHOR CONTRIBUTIONS

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

320

## VII. COMPETING INTERESTS

321 None of the authors have any competing interests.

322

## VIII. REFERENCES

---

- 323 [1] Janet K Jansson and Erin S Baker. A multi-omic future for microbiome studies. *Nat Microbiol*,  
324 1(16049):645, 2016.
- 325 [2] Rob Knight, Alison Vrbanac, Bryn C Taylor, Alexander Aksenov, Chris Callewaert, Justine  
326 Debelius, Antonio Gonzalez, Tomasz Kosciolk, Laura-Isobel McCall, Daniel McDonald, et al.  
327 Best practices for analysing microbiomes. *Nature Reviews Microbiology*, page 1, 2018.
- 328 [3] Chen Meng, Oana A Zeleznik, Gerhard G Thallinger, Bernhard Kuster, Amin M Gholami, and  
329 Aedín C Culhane. Dimension reduction techniques for the integrative analysis of multi-omics  
330 data. *Brief. Bioinform.*, 17(4):628–641, July 2016.
- 331 [4] Gwénaëlle Le Gall, Samah O Noor, Karyn Ridgway, Louise Scovell, Crawford Jamieson, Ian T  
332 Johnson, Ian J Colquhoun, E Kate Kemsley, and Arjan Narbad. Metabolomics of fecal extracts

- 333 detects altered metabolic activity of gut microbiota in ulcerative colitis and irritable bowel  
334 syndrome. *Journal of proteome research*, 10(9):4208–4218, 2011.
- 335 [5] Florian Rohart, Benoit Gautier, Amrit Singh, and Kim-Anh Le Cao. mixomics: An r pack-  
336 age for ‘omics feature selection and multiple data integration. *PLoS computational biology*,  
337 13(11):e1005752, 2017.
- 338 [6] Bo Wang, Aziz M Mezlini, Feyyaz Demir, Marc Fiume, Zhuowen Tu, Michael Brudno, Ben-  
339 jamin Haibe-Kains, and Anna Goldenberg. Similarity network fusion for aggregating data  
340 types on a genomic scale. *Nature methods*, 11(3):333, 2014.
- 341 [7] Ricard Argelaguet, Britta Velten, Damien Arno, Sascha Dietrich, Thorsten Zenz, John C  
342 Marioni, Florian Buettner, Wolfgang Huber, and Oliver Stegle. Multi-omics factor analysis—a  
343 framework for unsupervised integration of multi-omics data sets. *Molecular systems biology*,  
344 14(6):e8124, 2018.
- 345 [8] Cajo JF Ter Braak and Piet FM Verdonschot. Canonical correspondence analysis and related  
346 multivariate methods in aquatic ecology. *Aquatic sciences*, 57(3):255–289, 1995.
- 347 [9] Daniela M Witten, Robert Tibshirani, and Trevor Hastie. A penalized matrix decompo-  
348 sition, with applications to sparse principal components and canonical correlation analysis.  
349 *Biostatistics*, 10(3):515–534, 2009.
- 350 [10] Antoine Bodein, Olivier Chapleur, Arnaud Droit, and Kim-Anh Lê Cao. A generic multivari-  
351 ate framework for the integration of microbiome longitudinal studies with other data types.  
352 *bioRxiv*, page 585802, 2019.
- 353 [11] Jonathan Friedman and Eric J Alm. Inferring correlation networks from genomic survey data.  
354 *PLoS computational biology*, 8(9):e1002687, 2012.
- 355 [12] Zachary D Kurtz, Christian L Müller, Emily R Miraldi, Dan R Littman, Martin J Blaser,  
356 and Richard A Bonneau. Sparse and compositionally robust inference of microbial ecological  
357 networks. *PLoS computational biology*, 11(5):e1004226, 2015.
- 358 [13] Sophie Weiss, Will Van Treuren, Catherine Lozupone, Karoline Faust, Jonathan Friedman,  
359 Ye Deng, Li Charlie Xia, Zhenjiang Zech Xu, Luke Ursell, Eric J Alm, et al. Correlation  
360 detection strategies in microbial data sets vary widely in sensitivity and precision. *The ISME*  
361 *journal*, 10(7):1669, 2016.
- 362 [14] Doris Vandeputte, Gunter Kathagen, Kevin D’hoë, Sara Vieira-Silva, Mireia Valles-Colomer,  
363 João Sabino, Jun Wang, Raul Y Tito, Lindsey De Commer, Youssef Darzi, Séverine Vermeire,

- 364 Gwen Falony, and Jeroen Raes. Quantitative microbiome profiling links gut community vari-  
365 ation to microbial load. *Nature*, 551(7681):507–511, November 2017.
- 366 [15] Gregory B Gloor, Jean M Macklaim, Vera Pawlowsky-Glahn, and Juan J Egozcue. Microbiome  
367 datasets are compositional: And this is not optional. *Front. Microbiol.*, 8, 2017.
- 368 [16] Keqi Tang, Jason S Page, and Richard D Smith. Charge competition and the linear dynamic  
369 range of detection in electrospray ionization mass spectrometry. *J. Am. Soc. Mass Spectrom.*,  
370 15(10):1416–1423, 2004.
- 371 [17] Richard King, Ryan Bonfiglio, Carmen Fernandez-Metzler, Cynthia Miller-Stein, and Timothy  
372 Olah. Mechanistic investigation of ionization suppression in electrospray ionization. *J. Am.*  
373 *Soc. Mass Spectrom.*, 11(11):942–950, 2000.
- 374 [18] B K Matuszewski, M L Constanzer, and C M Chavez-Eng. Strategies for the assessment of  
375 matrix effect in quantitative bioanalytical methods based on HPLCMS/MS. *Anal. Chem.*,  
376 75(13):3019–3030, 2003.
- 377 [19] Alžběta Kalivodová, Karel Hron, Peter Filzmoser, Lukáš Najdekr, Hana Janečková, and  
378 Tomáš Adam. PLS-DA for compositional data with application to metabolomics. *Journal*  
379 *of Chemometrics*, 29(1):21–28, 2015.
- 380 [20] J. K. Jansson and E. S. Baker. A multi-omic future for microbiome studies. *Nat Microbiol.*,  
381 1:16049, 04 2016.
- 382 [21] David Lovell, Vera Pawlowsky-Glahn, Juan José Egozcue, Samuel Marguerat, and Jürg Bähler.  
383 Proportionality: a valid alternative to correlation for relative data. *PLoS computational biol-*  
384 *ogy*, 11(3):e1004075, 2015.
- 385 [22] Laura Tipton, Christian L Müller, Zachary D Kurtz, Laurence Huang, Eric Kleerup, Alison  
386 Morris, Richard Bonneau, and Elodie Ghedin. Fungi stabilize connectivity in the lung and  
387 skin microbial ecosystems. *Microbiome*, 6(1):12, 2018.
- 388 [23] James T Morton, Clarisse Marotz, Alex Washburne, Justin Silverman, Livia S Zaramela, Anna  
389 Edlund, Karsten Zengler, and Rob Knight. Establishing microbial composition measurement  
390 standards with reference frames. *Nature Communications*, 10(1):2719, 2019.
- 391 [24] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean,  
392 Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A  
393 system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Sys-*  
394 *tems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.

- 395 [25] Robert A Quinn, William Comstock, Tianyu Zhang, James T Morton, Ricardo da Silva,  
396 Alda Tran, Alexander Aksenov, Louis-Felix Nothias, Daniel Wangpraseurt, Alexey V Melnik,  
397 Gail Ackermann, Douglas Conrad, Isaac Klapper, Rob Knight, and Pieter C Dorrestein. Niche  
398 partitioning of a pathogenic microbiome driven by chemical gradients. *Sci Adv*, 4(9):eaau1908,  
399 September 2018.
- 400 [26] Shabnam Shalapour, Xue-Jia Lin, Ingmar N Bastian, John Brain, Alastair D Burt, Alexan-  
401 der A Aksenov, Alison F Vrbanac, Weihua Li, Andres Perkins, Takaji Matsutani, et al.  
402 Inflammation-induced iga+ cells dismantle anti-liver cancer immunity. *Nature*, 551(7680):340,  
403 2017.
- 404 [27] J. Lloyd-Price, C. Arze, A. N. Ananthakrishnan, M. Schirmer, J. Avila-Pacheco, T. W. Poon,  
405 E. Andrews, N. J. Ajami, K. S. Bonham, C. J. Brislawn, D. Casero, H. Courtney, A. Gonza-  
406 lez, T. G. Graeber, A. B. Hall, K. Lake, C. J. Landers, H. Mallick, D. R. Plichta, M. Prasad,  
407 G. Rahnavard, J. Sauk, D. Shungin, Y. Vazquez-Baeza, R. A. White, J. Braun, L. A. Denson,  
408 J. K. Jansson, R. Knight, S. Kugathasan, D. P. B. McGovern, J. F. Petrosino, T. S. Stappen-  
409 beck, H. S. Winter, C. B. Clish, E. A. Franzosa, H. Vlamakis, R. J. Xavier, C. Huttenhower,  
410 J. Bishai, K. Bullock, A. Deik, C. Dennis, J. L. Kaplan, H. Khalili, L. J. McIver, C. J. Moran,  
411 L. Nguyen, K. A. Pierce, R. Schwager, A. Sirota-Madi, B. W. Stevens, W. Tan, J. J. Ten Ho-  
412 eve, G. Weingart, R. G. Wilson, and V. Yajnik. Multi-omics of the gut microbial ecosystem  
413 in inflammatory bowel diseases. *Nature*, 569(7758):655–662, May 2019.
- 414 [28] Tami L Swenson, Ulas Karaoz, Joel M Swenson, Benjamin P Bowen, and Trent R Northen.  
415 Linking soil biology and chemistry in biological soil crust using isolate exometabolomics. *Na-  
416 ture communications*, 9(1):19, 2018.
- 417 [29] Richard Baran, Eoin L Brodie, Jazmine Mayberry-Lewis, Eric Hummel, Ulisses Nunes  
418 Da Rocha, Romy Chakraborty, Benjamin P Bowen, Ulas Karaoz, Hinsby Cadillo-Quiroz,  
419 Ferran Garcia-Pichel, et al. Exometabolite niche partitioning among sympatric soil bacteria.  
420 *Nature communications*, 6:8289, 2015.
- 421 [30] Robert A Quinn, Katrine Whiteson, Yan-Wei Lim, Peter Salamon, Barbara Bailey, Simone  
422 Mienardi, Savannah E Sanchez, Don Blake, Doug Conrad, and Forest Rohwer. A winogradsky-  
423 based culture system shows an association between microbial fermentation and cystic fibrosis  
424 exacerbation. *ISME J.*, 9(4):1024–1038, March 2015.



- 425 [31] Wilna J Moree, Vanessa V Phelan, Cheng-Hsuan Wu, Nuno Bandeira, Dale S Cornett, Bren-  
426 dan M Duggan, and Pieter C Dorrestein. Interkingdom metabolic transformations captured  
427 by microbial imaging mass spectrometry. *Proceedings of the National Academy of Sciences*,  
428 109(34):13811–13816, 2012.
- 429 [32] Mingxun Wang, Jeremy J Carver, Vanessa V Phelan, Laura M Sanchez, Neha Garg, Yao Peng,  
430 Don Duy Nguyen, Jeramie Watrous, Clifford A Kapon, Tal Luzzatto-Knaan, Carla Porto,  
431 Amina Bouslimani, Alexey V Melnik, Michael J Meehan, Wei-Ting Liu, Max Crüsemann,  
432 Paul D Boudreau, Eduardo Esquenazi, Mario Sandoval-Calderón, Roland D Kersten, Laura A  
433 Pace, Robert A Quinn, Katherine R Duncan, Cheng-Chih Hsu, Dimitrios J Floros, Ronnie G  
434 Gavilan, Karin Kleigrewe, Trent Northen, Rachel J Dutton, Delphine Parrot, Erin E Carlson,  
435 Bertrand Aigle, Charlotte F Michelsen, Lars Jelsbak, Christian Sohlenkamp, Pavel Pevzner,  
436 Anna Edlund, Jeffrey McLean, Jörn Piel, Brian T Murphy, Lena Gerwick, Chih-Chuang  
437 Liaw, Yu-Liang Yang, Hans-Ulrich Humpf, Maria Maansson, Robert A Keyzers, Amy C Sims,  
438 Andrew R Johnson, Ashley M Sidebottom, Brian E Sedio, Andreas Klitgaard, Charles B  
439 Larson, Cristopher A Boya P, Daniel Torres-Mendoza, David J Gonzalez, Denise B Silva,  
440 Lucas M Marques, Daniel P Demarque, Egle Pociute, Ellis C O’Neill, Enora Briand, Eric  
441 J N Helfrich, Eve A Granatosky, Evgenia Glukhov, Florian Ryffel, Hailey Houson, Hosein  
442 Mohimani, Jenan J Kharbush, Yi Zeng, Julia A Vorholt, Kenji L Kurita, Pep Charusanti,  
443 Kerry L McPhail, Kristian Fog Nielsen, Lisa Vuong, Maryam Elfeki, Matthew F Traxler,  
444 Niclas Engene, Nobuhiro Koyama, Oliver B Vining, Ralph Baric, Ricardo R Silva, Saman-  
445 tha J Mascuch, Sophie Tomasi, Stefan Jenkins, Venkat Macherla, Thomas Hoffman, Vinayak  
446 Agarwal, Philip G Williams, Jingqui Dai, Ram Neupane, Joshua Gurr, Andrés M C Rodríguez,  
447 Anne Lamsa, Chen Zhang, Kathleen Dorrestein, Brendan M Duggan, Jehad Almaliti, Pierre-  
448 Marie Allard, Prasad Phapale, Louis-Felix Nothias, Theodore Alexandrov, Marc Litaudon,  
449 Jean-Luc Wolfender, Jennifer E Kyle, Thomas O Metz, Tyler Peryea, Dac-Trung Nguyen,  
450 Danielle VanLeer, Paul Shinn, Ajit Jadhav, Rolf Müller, Katrina M Waters, Wenyuan Shi,  
451 Xueting Liu, Lixin Zhang, Rob Knight, Paul R Jensen, Bernhard O Palsson, Kit Pogliano,  
452 Roger G Linington, Marcelino Gutiérrez, Norberto P Lopes, William H Gerwick, Bradley S  
453 Moore, Pieter C Dorrestein, and Nuno Bandeira. Sharing and community curation of mass  
454 spectrometry data with global natural products social molecular networking. *Nat. Biotechnol.*,  
455 34(8):828–837, August 2016.

- 456 [33] Raina Margaret Maier and G Soberon-Chavez. Pseudomonas aeruginosa rhamnolipids:  
457 biosynthesis and potential applications. *Applied Microbiology and Biotechnology*, 54(5):625–  
458 633, 2000.
- 459 [34] Thammajun L Wood, Ting Gong, Lei Zhu, James Miller, Daniel S Miller, Bei Yin, and  
460 Thomas K Wood. Rhamnolipids from pseudomonas aeruginosa disperse the biofilms of sulfate-  
461 reducing bacteria. *NPJ biofilms and microbiomes*, 4(1):22, 2018.
- 462 [35] Lucy Allen, David H Dockrell, Theresa Pattery, Daniel G Lee, Pierre Cornelis, Paul G  
463 Hellewell, and Moira KB Whyte. Pyocyanin production by pseudomonas aeruginosa induces  
464 neutrophil apoptosis and impairs neutrophil-mediated host defenses in vivo. *The Journal of*  
465 *Immunology*, 174(6):3643–3649, 2005.
- 466 [36] Robert A. Quinn, Alison Vrbanac, Alexey V. Melnik, Kathryn A. Patras, Mitchell Christy,  
467 Andrew T. Nelson, Alexander Aksenov, Anupriya Tripathi, Greg Humphrey, Ricardo da Silva,  
468 Robert Bussell, Taren Thron, Mingxun Wang, Fernando Vargas, Julia M. Gauglitz, Michael J.  
469 Meehan, Orit Poulsen, Brigid S. Boland, John T. Chang, William J. Sandborn, Meerana  
470 Lim, Neha Garg, Julie Lumeng, Barbara I. Kazmierczak, Ruchi Jain, Marie Egan, Kyung E.  
471 Rhee, Gabriel G. Haddad, Dionicio Siegel, Sarkis Mazmanian, Victor Nizet, Rob Knight, and  
472 Pieter C. Dorrestein. Chemical impacts of the microbiome across scales reveal novel conjugated  
473 bile acids. *bioRxiv*, 2019.
- 474 [37] Michelle K. Paczosa and Joan Meccas. Klebsiella pneumoniae: Going on the offense with a  
475 strong defense. *Microbiology and Molecular Biology Reviews*, 80(3):629–661, 2016.
- 476 [38] Elisa Tiratterra, Placido Franco, Emanuele Porru, Konstantinos H Katsanos, Dimitrios K  
477 Christodoulou, and Giulia Roda. Role of bile acids in inflammatory bowel disease. *Annals of*  
478 *gastroenterology*, 31(3):266, 2018.
- 479 [39] Alan F Hofmann and Lars Eckmann. How bile acids confer gut mucosal protection against  
480 bacteria. *Proceedings of the National Academy of Sciences*, 103(12):4333–4334, 2006.
- 481 [40] Máire Begley, Cormac GM Gahan, and Colin Hill. The interaction between bacteria and bile.  
482 *FEMS microbiology reviews*, 29(4):625–651, 2005.
- 483 [41] Y. Okada, Y. Tsuzuki, J. Miyazaki, K. Matsuzaki, R. Hokari, S. Komoto, S. Kato,  
484 A. Kawaguchi, S. Nagao, K. Itoh, T. Watanabe, and S. Miura. Propionibacterium freudenre-  
485 ichii component 1,4-dihydroxy-2-naphthoic acid (DHNA) attenuates dextran sodium sulphate  
486 induced colitis by modulation of bacterial flora and lymphocyte homing. *Gut*, 55(5):681–688,

- 487 May 2006.
- 488 [42] B. Foligne, S. Parayre, R. Cheddani, M. H. Famelart, M. N. Madec, C. Ple, J. Breton,  
489 J. Dewulf, G. Jan, and S. M. Deutsch. Immunomodulation properties of multi-species fer-  
490 mented milks. *Food Microbiol.*, 53(Pt A):60–69, Feb 2016.
- 491 [43] C. Ple, J. Breton, R. Richoux, M. Nurdin, S. M. Deutsch, H. Falentin, C. Herve, V. Chuat,  
492 R. Lemee, E. Maguin, G. Jan, M. Van de Guchte, and B. Foligne. Combining selected im-  
493 munomodulatory *Propionibacterium freudenreichii* and *Lactobacillus delbrueckii* strains: Re-  
494 verse engineering development of an anti-inflammatory cheese. *Mol Nutr Food Res*, 60(4):935–  
495 948, Apr 2016.
- 496 [44] Natacha Colliou, Yong Ge, Bikash Sahay, Minghao Gong, Mojgan Zadeh, Jennifer L Owen,  
497 Josef Neu, William G Farmerie, Francis Alonzo, Ken Liu, et al. Commensal propionibacterium  
498 strain ufl mitigates intestinal inflammation via th17 cell regulation. *The Journal of clinical*  
499 *investigation*, 127(11):3970–3986, 2017.
- 500 [45] Nasser M Nasrabadi. Pattern recognition and machine learning. *Journal of electronic imaging*,  
501 16(4):049901, 2007.
- 502 [46] Vera Pawlowsky-Glahn, Juan José Egozcue, and Raimon Tolosana-Delgado. *Modeling and*  
503 *Analysis of Compositional Data*. John Wiley & Sons, February 2015.
- 504 [47] Mikolov, Tomas and Sutskever, Ilya and Chen, Kai and Corrado, Greg S and Dean, Jeff.  
505 Distributed representations of words and phrases and their compositionality. In *Advances in*  
506 *neural information processing systems*, pages 3111–3119, 2013.
- 507 [48] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recom-  
508 mender systems. *Computer*, (8):30–37, 2009.
- 509 [49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*  
510 *preprint arXiv:1412.6980*, 2014.
- 511 [50] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of*  
512 *machine Learning research*, 3(Jan):993–1022, 2003.
- 513 [51] Kris Sankaran and Susan P Holmes. Latent variable modeling for the microbiome. *arXiv*  
514 *preprint arXiv:1706.04969*, 2017.
- 515 [52] John Aitchison and Michael Greenacre. Biplots of compositional data. *Journal of the Royal*  
516 *Statistical Society: Series C (Applied Statistics)*, 51(4):375–392, 2002.

- 517 [53] John Aitchison, KW Ng, et al. Conditional compositional biplots: theory and application.  
518 2005.
- 519 [54] JA Martín-Fernández, V Pawlowsky-Glahn, JJ Egozcue, and R Tolosona-Delgado. Advances  
520 in principal balances for compositional data. *Mathematical Geosciences*, 50(3):273–298, 2018.
- 521 [55] Michael A Skinnider, Jordan W Squair, and Leonard J Foster. Evaluating measures of asso-  
522 ciation for single-cell transcriptomics. *Nature methods*, 16(5):381, 2019.
- 523 [56] Han Liu, Kathryn Roeder, and Larry Wasserman. Stability approach to regularization se-  
524 lection (stars) for high dimensional graphical models. In *Advances in neural information*  
525 *processing systems*, pages 1432–1440, 2010.
- 526 [57] Evan Bolyen, Jai Ram Rideout, Matthew R Dillon, Nicholas A Bokulich, Christian Abnet,  
527 Gabriel A Al-Ghalith, Harriet Alexander, Eric J Alm, Manimozhiyan Arumugam, Francesco  
528 Asnicar, et al. Qiime 2: Reproducible, interactive, scalable, and extensible microbiome data  
529 science. Technical report, PeerJ Preprints, 2018.
- 530 [58] Yoshiki Vázquez-Baeza, Meg Pirrung, Antonio Gonzalez, and Rob Knight. Emperor: a tool  
531 for visualizing high-throughput microbial community data. *Gigascience*, 2(1):16, 2013.
- 532 [59] Oliver Fiehn, Don Robertson, Jules Griffin, Mariet van der Werf, Basil Nikolau, Norman Mor-  
533 rison, Lloyd W Sumner, Roy Goodacre, Nigel W Hardy, Chris Taylor, et al. The metabolomics  
534 standards initiative (msi). *Metabolomics*, 3(3):175–178, 2007.

## 535 IX. METHODS

### 536 A. Mmvec neural network architecture

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

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

metabolite  $\nu \in \{1, \dots, M\}$ ,

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

544  
545

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 models generative process draws a single microbe from a single draw from the categorical distribution

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

549  
550

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

$$p(\nu|\mu) = \frac{\exp(\mathbf{v}_\nu \cdot \mathbf{u}_\mu + \nu_{\nu 0} + u_{\mu 0})}{\sum_j \exp(\mathbf{v}_j \cdot \mathbf{u}_\mu + \nu_{j 0} + u_{\mu 0})} ,$$

551

$$\mathbf{q}_\mu = [p(\nu_1|\mu), \dots, p(\nu_M|\mu)]$$

552  
553

Here,  $\nu_{j0} + u_{\mu 0}$  are row and column biases, which are required to accurately estimate the conditional probabilities. The above transformation is the softmax transform [45] to compute probabilities from real-valued quantities. This transformation is also known as the inverse clr transform [46], 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.

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

558  
559

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 [47]. There are a couple of majors differences to be considered. First, in the original application of word2vec, a skipgram was proposed. Skipgrams [47] have been designed to account for the sequential nature of text. There is no such sequential

564

565

566

567

568 nature with microbiome or metabolite samples, the only ordering information that is known  
 569 is the sample membership. As a result, the skipgrams can be replaced using multinomial  
 570 sampling, where a single microbe is randomly sampled from a microbiome sample at each  
 571 gradient descent step.

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

$$585 \quad p(Y = y|X_\mu = 1) = \text{Multinomial}(y|q_\mu)$$

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

$$595 \quad \mathbf{U} = [\mathbf{0}, \mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_N]^T \quad \mathbf{V} = [\mathbf{v}_0, \mathbf{0}, \mathbf{v}_1, \dots, \mathbf{v}_M] .$$

597 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  
 598 and  $M$  metabolites. The number dimensions  $p$  for both  $\mathbf{U}$  and  $\mathbf{V}$  as well as the priors are  
 599 specified by the user, but can also be evaluated during cross-validation. The biases  $\mathbf{u}_0$  and

600  $\mathbf{v}_0$  are critical for estimating accurate co-occurrence probabilities, as suggested by similar  
 601 methodologies used in recommender systems [48]. The  $\mathbf{U}$  and  $\mathbf{V}$  matrices are estimated  
 602 through maximum a posteriori (MAP) estimation using ADAM [49] with the following log-  
 603 posterior

$$\begin{aligned}
 604 \quad \mathcal{L} &= \mathcal{L}_Y + \mathcal{L}_U + \mathcal{L}_V \\
 605 \quad \mathcal{L}_U &= \sum_{\mu} \sum_{\rho=1}^p \mathcal{N}(U_{\mu,\rho} | 0, \sigma_u) \\
 606 \quad \mathcal{L}_V &= \sum_{\nu} \sum_{\rho=1}^p \mathcal{N}(V_{\nu,\rho} | 0, \sigma_v) \\
 607 \quad \mathcal{L}_Y &= \sum_k \sum_{r \in x_k} \text{Multinomial}(\mathbf{y}_k | \mathbf{q}_\mu) . \\
 608
 \end{aligned}$$

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

$$\begin{aligned}
 616 \quad SSE &= \sum_{k,i} (y_k - m_k \cdot \text{softmax}(\mathbf{V}\mathbf{U}_{u_{ik},\cdot}))^2 . \\
 617
 \end{aligned}$$

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

621

## 622 B. Microbe-metabolite vectors in simplicial coordinates

623 Here, we will provide some insights behind the underlying geometry behind this neural  
 624 network. Doing so will provide intuition behind the algebraic operations commonly applied  
 625 in the context of word2vec, suggesting the possibility of performing similar tasks in the  
 626 context of microbe-metabolite interactions. Furthermore, this will motivate the use of the

627 Aitchison distance to quantify microbe-microbe and metabolite-metabolite interactions. Fi-  
 628 nally we will make a connection to topic modeling, providing another means to potentially  
 629 interpret the latent dimensions in the model. The connection between the softmax and the  
 630 inverse clr transform suggests that the inputs to this transform can be represented in clr  
 631 coordinates. The softmax function and its corresponding inverse, the clr transform, is given  
 632 as follows

$$633 \quad \text{softmax}(x) = \left[ \frac{e^{x_1}}{\sum_i e^{x_i}}, \dots, \frac{e^{x_1}}{\sum_i e^{x_i}} \right]$$

$$634 \quad \text{clr}(z) = \left[ \log \frac{z_1}{g(z)}, \dots, \log \frac{z_D}{g(z)} \right]$$

635  
636

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

640 **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  $V \mathbf{1}_p = \mathbf{0}$

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

642 Then

$$643 \quad Q = U(V - \mathbf{1}_M \lambda_v^T).$$

644  
645 Given that the rows of  $Q$  sum to 0, then

$$646 \quad U(V - \mathbf{1}_M \lambda_v^T)^T \mathbf{1}_M = 0$$

$$647 \quad U \lambda_v^T M = 0.$$

648  
649 This means that only the trivial solution  $\lambda_v = 0$  exists, therefore the rows of  $V$  do sum to  
 650 0.

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

$$653 \quad Q = (U - \mathbf{1}_N \lambda_u^T) V^T.$$

654  
655 Given that the columns of  $Q$  sum to 0, then

$$656 \quad \mathbf{1}_N^T (U - \mathbf{1}_N \lambda_u^T) V^T = 0$$

$$657 \quad N \lambda_u^T V = 0.$$

658



659 This means that only the trivial solution  $\lambda_u = 0$  exists, therefore the rows of  $\mathbf{U}$  do sum to  
660 0.

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

661

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

666

#### *Topic proportions*

667

668 Since the  $U$  and  $V$  row vectors are in clr coordinates, that implies that these row vectors  
 669 can be directly converted to  $p$ -dimensional proportions, yielding a similar interpretation to  
 topics used in models such as LDA [50, 51].

670

#### *Linearity*

671

Vectors in clr coordinates are known to satisfy linearity, namely

$$clr(\alpha x + y) = \alpha clr(x) + clr(y)$$

672

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

675

#### *Isometry*

676

677 The clr transform is distance preserving, meaning that the Aitchison distance on propor-  
 678 tions is equivalent to the Euclidean distance on clr vectors. This provides motivation for  
 679 using Euclidean distances to compute microbe-microbe and metabolite-metabolite similari-  
 ties.

680

### **C. Visualization through biplots**

681

682 Visualization techniques from compositional data analysis can aid with interpretation  
 683 [52, 53].  $U$  and  $V$  can be visualized as factors within a biplot to visualize the microbe-  
 metabolite embeddings on a single plot. The first two latent dimensions of  $U$  represent

684 microbial coordinates on a 2D scatter plot and the first two latent dimensions of  $\mathbf{V}$  represent  
 685 metabolite coordinates on a 2D scatter plot. Typically the coordinate from the  $\mathbf{V}$  matrix  
 686 are plotted as arrows from the origin in order to identify features that explain the variance in  
 687  $\mathbf{U}$ . However, in our case studies, there are typically many more metabolites than microbes  
 688 - so we opt to visualize the metabolites as points and microbes as arrows for a simpler  
 689 visualization

690 As suggested by the above theorem, the distance between points approximates the Aitchi-  
 691 son distance between metabolites, and the distance between arrow tips approximates the  
 692 Aitchison distance between microbes. As suggested in [54], the Aitchison distance is also  
 693 equivalent to the variance of the log ratios, suggesting that microbe-microbe and metabolite-  
 694 metabolite distances could also be interpreted as a measure of proportionality [21]

#### 695 D. Benchmarks

696 The simulated data was based on a cystic fibrosis biofilm model derived in Quinn et al [25]  
 697 shown in Figure S12 in the paper. The biofilm model was built to explain how fermenters and  
 698 *P. aeruginosa* responded to different concentrations of sugars, amino acids, pH, oxygen and  
 699 antibiotics across the Winogradsky column. These models solved for differential equations  
 700 integrating Monod kinetics and diffusion processes and was run in Matlab using the code  
 701 provided at [https://github.com/zhangzhongxun/WinCF\\_model\\_Code](https://github.com/zhangzhongxun/WinCF_model_Code)  
 702 From this simulation, we only focus 2 microbes and 5 compounds. The two microbes are  
 703 *P. aeruginosa* ( $\theta_p$ ) and fermenters ( $\theta_f$ ). The five compounds (SG), acids (F), ammonium  
 704 (P), amino acids (SA) and inhibition molecules (I). In order to simulate a high dimensional  
 705 dataset, each microbial taxon was split into 50 different subtaxa and each compound was  
 706 split into 50 molecular subclasses. The partitioning procedure is given as follows

$$\begin{aligned} \mathbf{p}_i &\sim \mathcal{N}(\mathbf{0}, \sigma_o \mathbf{I}) & \mathbf{q}_i &\sim \mathcal{N}(\mathbf{0}, \sigma_c \mathbf{I}) \\ \mathbf{o}_{ij} &= \kappa_{ij} i l r^{-1}(\mathbf{p}_i) & \mathbf{c}_{ik} &= \eta_{ik} i l r^{-1}(\mathbf{q}_i), \end{aligned}$$

707 where  $\mathbf{p}_i$  is a vector proportions representing how the subtaxa corresponding to  $j$  will be  
 708 distributed in sample  $i$ .  $\kappa_{ij}$  represents the absolute abundance of taxon  $j$  in sample  $i$ .  $\mathbf{o}_{ij}$   
 709 represents a vector of the absolute abundances for all of the subtaxa corresponding to taxon

710  $j$ . These are the absolute abundances that are used for comparison in Figure 2.

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

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

724

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

$$\begin{aligned} \zeta_i &\sim \mathcal{LN}(n, \tau_o) & \omega_i &\sim \mathcal{LN}(m, \tau_c) \\ x_i &\sim \mathcal{PLN}(\zeta_i C(\mathbf{o}_i), \epsilon_o) & y_i &\sim \mathcal{LN}(\omega_i C(\mathbf{c}_i), \epsilon_c) . \end{aligned}$$

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

735 Once the total sequencing depth and the total intensities are sampled, the microbial

736 sequencing counts and metabolite abundances are then sampled. A Poisson lognormal dis-  
737 tribution is used to generate the microbial counts from the microbial proportions  $C(\mathbf{o}_i)$   
738 scaled by the sequencing depth  $\zeta_i$ . The counts are sampled with error  $\epsilon_o$ . A Lognormal  
739 distribution is used to generate the metabolite abundances from metabolite proportions  
740  $C(\mathbf{c}_i)$  scaled by the total intensity  $\omega_i$ . The abundances are sampled with error  $\epsilon_c$ . All of  
741 the code used to generate the benchmarks can be found at [https://github.com/knightlab-](https://github.com/knightlab-analyses/multiomic-cooccurrences)  
742 [analyses/multiomic-cooccurrences](https://github.com/knightlab-analyses/multiomic-cooccurrences)

743

## E. Data Analysis

744 Due to the overwhelming sparsity in microbiome datasets, some filtering is required in  
745 order to infer microbe-metabolite interactions. We chose to filter out microbes that appear  
746 in less than 10 samples, since these microbes don't have enough information to infer which  
747 metabolites are co-occurring with them. In other words the mmvec model has too many  
748 degrees of freedom to perform inference on these microbes. For the cystic fibrosis study,  
749 there were 172 samples and after filtering there were 138 unique microbial taxa and 462  
750 metabolite features. For the biocrust soils study, there were 19 samples and after filtering  
751 there were 466 unique microbial taxa and 85 metabolite features. For the murine high fat  
752 diet study, there were 434 samples and after filtering there were 902 microbes and 11978  
753 metabolites. For the IBD dataset, there were 13920 features in the c18 LCMS dataset, 26966  
754 features in the c8 LCMS dataset and 562 taxa. Cross validation was performed across all  
755 studies to evaluate overfitting. In the desert biocrust soils experiment, 1 sample out of 19  
756 samples was randomly chosen to be left out for cross-validation. In all of the other studies,  
757 10 samples were randomly chosen to be left out for cross-validation. All of the analyses can  
758 be found under <https://github.com/knightlab-analyses/multiomic-cooccurrences>.

759

## F. Data availability

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

764 The biocrust soils data was retrieved from the supplemental section in Swenson et al [28].  
765 The High fat diet murine model case study 16S rRNA data can be found under  
766 <http://qiita.microbio.me>; study id: 10856. The High fat diet murine model case study are  
767 publicly available at  
768 <https://massive.ucsd.edu/> at MassIVE ID MSV000080918. The GNPS analysis for this  
769 study can be accessed at  
770 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=977d85bba47b4e96bf69872b961b8edd>  
771 The IBD data used can be found under <https://ibdmdb.org/>

## 772 **G. Software availability**

773 The software implementing the mmvec algorithm can be found under  
774 <https://github.com/biocore/mmvec>  
775 Differential abundance analyses in the high fat diet study was performed using L2-regularized  
776 multinomial regression using software available at <https://github.com/biocore/songbird>  
777 The software used to build the multiomics network can be found at  
778 [https://github.com/mortonjt/multiomics\\_network](https://github.com/mortonjt/multiomics_network)

## 779 **X. SUPPLEMENTAL MATERIALS**

### 780 **A. Challenges of analyzing multiple compositional datasets**

781 One of the challenges involved with inferring microbe-metabolite interactions is resolving  
782 the differences between the absolute abundances in the original environment, and the mea-  
783 sured relative abundances from sequencing and mass spectrometry. In order to guarantee  
784 consistency between absolute and relative abundances, scale invariance must be maintained  
785 [23], otherwise overwhelming false positive or false negative rates can occur (Figure S1).

786 As shown in Figure 2d, most tools are not scale invariant. The reason for the contradiction  
787 is further clarified in Figure S2, from the proportions, it looks like most of the microbes are  
788 decreasing when in fact they are merely increasing with a slower growth rate compared  
789 to the fastest growing microbe. The inability to determine which microbes are actually  
790 increasing or decreasing caused Pearson and Spearman to misannotate the vast majority of

791 the interactions, with the except of the interactions of the fastest growing microbe in this  
792 scenario.

793 It may appear the benchmark in Figure 2d that the proportionality metric rho is scale  
794 invariant in the context of multiomics analysis. However, another benchmark in Figure S2  
795 reveals that rho is not scale invariant. The reason why scale-invariance breaks for phi, rho  
796 and SparCC is because the microbe and metabolite datasets have differing absolute sums.  
797 When analyzing a single dataset all three of these metrics rely on the following quantity to  
798 hold

$$V\left(\log \frac{x_i}{x_j}\right) = V\left(\log \frac{Np_{x_i}}{Np_{x_j}}\right) = V\left(\log \frac{p_{x_i}}{p_{x_j}}\right),$$

799 where  $x_i$  correspond to random variable quantifying absolute abundances of microbe  $i$ ,  $N$   
800 corresponds to a random variable quantifying the total population size of the microbes, and  
801  $p_{x_i}$  correspond to the proportions of microbe  $i$ . Due to the log-ratio, the dependence on  
802 the total population size  $N$  drops out, negating the need to quantify total microbial load.  
803 This is critical for microbiome sequencing applications, since quantifying total microbial  
804 load can be challenging [14, 23]. Furthermore, methods that satisfy scale invariance have  
805 shown shown to be superior to other tools in the context of co-occurrence analysis [55].

However, scale invariance is much harder to enforce when analyzing co-occurrence relationships across multiple datasets. When evaluating the variance of the log ratios across multiple datasets, the scale-invariance relationship is not immediately satisfied

$$V\left(\log \frac{x_i}{y_j}\right) = V\left(\log \frac{Np_{x_i}}{Mp_{y_j}}\right) \neq V\left(\log \frac{p_{x_i}}{p_{y_j}}\right).$$

806 Here,  $y_j$  refers to the absolute abundances of the metabolite  $j$ ,  $M$  refers to the total number  
807 of metabolites in the original environment and  $p_{y_j}$  is the proportion of metabolite  $j$ .

808 This was recognized in Tipton et al [22] and additional modifications were added to  
809 SPIEC-EASI. These modifications explain the superior performance of SPIEC-EASI in the  
810 benchmarks. However, there are two major impediments to the application of SPIEC-EASI,  
811 namely zeros and StARs [56] regularization. SPIEC-EASI still relies on using pseudocounts,  
812 adding bias into the resulting inference. Furthermore, STaRs has been shown to enhance  
813 the interpretation of the SPIEC-EASI results, but STaRs is not a scale-invariant procedure.  
814 Due to this alone, the absolute and relative estimates will not match as shown in Figure 2

815 and Figure S2. This may not be a problem when analyzing multiomics datasets with similar  
816 scales, such as 16S and ITS sequencing datasets. However, these problems will become  
817 exacerbated when analyzing datasets with drastically different scales. Sequencing counts  
818 are usually below 100k reads per samples, where as MS intensities are up to 10e9 intensity  
819 units.

820 In light of the challenges discussed above, there are some scenarios where standard statis-  
821 tical methods will be consistent with the biological reality. As discussed in [23] the differences  
822 between absolute and relative abundances is essentially a constant factor attributed to the  
823 changes in the total biomass. If the total biomass is constant, then traditional statistical  
824 methods will work fine. In the case of the cystic fibrosis dataset [25], microbial communities  
825 were grown in fixed size Winogradsky columns. As a result, the total size of the community  
826 could be constrained due to the limited resources and space. This could explain the consis-  
827 tency between Pearson and mmvec in this particular study (Figure S3). On the other hand,  
828 in the biocrust soils study, the drastic differences between Spearman and mmvec could be  
829 explained by the rapid increase in biomass following the wetting event.

830

## B. Software workflows

831 To facilitate utilization of the mmvec tool, we have developed two different user facing  
832 interfaces. First, we have developed a qiime2 plugin [57], where mmvec can be run using a  
833 simple command line interface. This interface is complemented using [24], where users can  
834 monitor convergence rates for their models in real-time and evaluate how different parameters  
835 will affect their model fit (Figure S4). Second, we have integrated mmvec into the Global  
836 Natural Product Social Molecular Networking (GNPS) platform that can be accessed by  
837 the public. The online interface through GNPS resolves several usability issues. First,  
838 GNPS facilitates import of metabolomics data into qiime2 by pre-processing, importing,  
839 and sample renaming, This is performed as part of the standard metabolomics analysis at  
840 GNPS (e.g. molecular networking and feature-based molecular networking). Second, since  
841 it is possible to both download and re-use outputs of workflows run at GNPS directly, it is  
842 straightforward to select the GNPS qza and molecule annotations needed for mmvec. The  
843 user will need to upload the accompanying feature and taxonomy data for qiime2 and the  
844 analysis will be begin. Once the workflow completes, the biplots can be viewed directly in



845 the browser and other outputs (e.g. ranks) are available for download (Figure S5).

846 The mmvec implementation is written using Tensorflow and can leverage GPUs for com-  
847 putation. The number of gradient descent iterations is specified by the user and model fit  
848 diagnostics can be monitored in real time using Tensorboard. The runtime of mmvec across  
849 16 cores can take multiple days until a model convergence reaches convergence. With GPUs,  
850 the running time is reduced to a few hours. Using a Telsa GPU, the model can reach conver-  
851 gence within 4 hours on the IBD dataset comprised of 562 microbial taxa, 26,966 metabolite  
852 features and 400 samples. However, there is a trade-off of accuracy and running time. More  
853 accurate models require smaller learning rates and may take longer to run.

854

## **XI. FIGURE LEGENDS**

Figure legends are below

855 Figure 1: Input data types and mmvec neural network architecture. (a) The neural net-  
856 work architecture where the input layer represents one-hot encodings of  $N$  microbes and  
857 the output layer represents the proportions of  $M$  metabolites.  $U$  corresponds to microbial  
858 vectors and  $V$  corresponds to metabolite vectors. (b) The pipeline for training mmvec.  
859 The objective behind mmvec is to predict metabolite abundances ( $y$ ) given a single input  
860 microbe sequence ( $x$ ), also known as a one-hot encoding. This training procedure will esti-  
861 mate conditional probabilities of observing a metabolite given the input microbe sequence.  
Cross-validation can be performed on hold-out samples to assess overfitting.

862 Figure 2: Simulation benchmarks. (a) Absolute abundances of microbes and metabolites  
863 simulated from differential equations derived in [25] for a specific spatial point. (b) Propor-  
864 tions of the abundances shown in (a). (c) F1 score, precision and recall curves comparing  
865 mmvec to Pearson, Spearman, SparCC, SPIEC-EASI, and proportionality metrics phi and  
866 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.

867 Figure 3: *M. vaginatus* released metabolites after the biocrust wetting event. (a) Compar-  
868 ison of *M. vaginatus* metabolite interactions estimated from Spearman and mmvec from  
869 (n=19 samples). All of the experimentally validated *M. vaginatus* released metabolites are  
870 labeled. All metabolites with contradicting findings between the wetting experiment and  
871 the *in vitro* experimental results are highlighted in red. Points are resized according to  
872 the  $10 \log(\text{p-value})$  obtained from Spearman correlation. Dashlines mark the cutoff for a  
873 Spearman correlation of zero, and the conditional log probabilities of zero. Here a zero  
874 log conditional probability represents the conditional probability of the average metabolite  
875 because all probabilities here are mean centered. (b) Benchmarks comparing the detection  
876 rate of the experimentally validated molecules across different statistical methodologies. (c)  
877 *M. vaginatus* proportions and (d) 4-guanidinobutanoate proportions following a wetting  
event.

878 Figure 4: Investigation of *P.aeruginosa*-associated molecules. (a) Biplot drawn from the  
879 mmvec conditional probabilities estimated for the cystic fibrosis dataset [25]. Arrows rep-  
880 resent microbes and dots represent metabolites. The x and y axes represent principal  
881 components from the SVD of the microbe-metabolite conditional probabilities estimated  
882 from mmvec (n=138 samples). Distances between points quantify co-occurrence strength be-  
883 tween metabolites, with small distances indicating metabolites that have a high probability  
884 of co-occurring with high probability. Distances between arrow tips quantify co-occurrence  
885 strength between microbes. The directionality of the arrows can be used to pinpoint which  
886 microbes can explain the metabolite co-occurrence patterns. Arrows highlighted in green  
887 correspond to putative cystic fibrosis pathogens and yellow arrows highlight known anaer-  
888 obes. Only known molecules produced by *P. aeruginosa* are labeled. (b) Scatter plot of  
889 molecules with respect to the oxygen gradient differential and the first principal component  
890 learned from mmvec (n=442 molecules) with linear regression model and 95% confidence  
891 interval for regression estimate. (c) The first principal component vs the number of samples  
892 where the taxa was the most abundant taxa in that sample . (d) Heatmap of *P. aeruginosa*  
893 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*.

894 Figure 5: Microbe/metabolite co-occurrences across study of HCC progression in the con-  
895 text of innate immunity in a mouse model [26]. (a) Visualization of microbial co-occurrence  
896 patterns, where distances between points approximates the Aitchison distance between  
897 microbes, which quantities microbial occurrences. Small distances are indicative of mi-  
898 crobes with high probability of co-occurring together. Microbes are colored according to  
899 their association with HFD, which was estimated using differential abundance analysis  
900 via multinomial regression. (b) Emperor [58] biplot of microbe-metabolite interactions,  
901 with metabolites colored according to their association with HFD. HFD association was  
902 estimated through differential abundance analysis via multinomial regression. Distances be-  
903 tween points approximate Aitchison distances between metabolites and distances between  
904 arrow tips approximate Aitchison distances between microbes. Several *Clostridium spp.*  
905 appear to co-occur with the new bile acid molecule cholate phenylalanine amidate, also  
referred to as Phe conjugated cholic acid.

906 Figure 6: Microbe-metabolite interactions of the human microbiome in association with  
907 IBD samples [27]. (a) Heatmap visualization of the inferred conditional probabilities for  
908 various bile acids given the presence of *Klebsiella*, *Roseburia* and *Clostridium bolteae*. (b)  
909 Heatmap visualization of the inferred conditional probabilities for the carnitines given the  
910 presence of *Klebsiella*, *Roseburia*, and *Clostridium bolteae*. (c) Multiomics biplot of the  
911 microbe-metabolite interactions learned from metagenomics profiles and C18 negative ion  
912 mode LC-MS. Microbes (arrows) and metabolites (spheres) are colored according to their  
913 differentials estimated from multinomial regression. *Klebsiella spp.* appears to be strongly  
914 associated with IBD, while *Propionibacterium spp.* has strong negative association. (d)  
915 Network of the top 300 edges where only the edges that contain *Klebsiella* and *Propioni-*  
*bacteriaceae* are visualized.

916 Figure S1: Description of the compositionality issue. (a) An illustration of how false neg-  
917 atives can occur - in the absolute abundance data, there is a strong Pearson correlation  
918 between the microbes and the metabolites. These correlations disappear when considering  
919 the corresponding proportions. (b) An illustration of how false positives can occur - in the  
920 absolute abundance data, there is no correlation between the dark green molecule and the  
921 dark blue microbe. However, the proportions of the same dataset show that there is a very  
strong correlation between the dark blue and the dark green molecule.



922 Figure S2: Illustration of how excessive misannotation rates can occur. (a) Absolute abun-  
923 dances and relative abundances of microbes/metabolites observed in an environment over  
924 time, with each microbe/metabolite colored according to its rate of increase / decrease. (b)  
925 A scale-invariance comparison of statistical methodologies. Points are colored by the cor-  
926 responding microbes in the interactions; triangle markers represent increasing metabolites  
927 and decreasing metabolites. Mmvec is the only method that remains consistent between the  
absolute and relative abundances.

928 Figure S3: Comparison of Pearson and mmvec on Cystic Fibrosis study. (a) Estimates  
929 of *P. aeruginosa* associated molecules between Pearson and the conditional probabilities  
930 calculated from the mmvec applied to the cystic fibrosis study dataset. The annotations  
931 correspond to level 2 or 3 of the metabolomics standards initiative [59] and may correspond  
932 to different isomeric species (n=462 molecules). (b) Ranks of Pearson coefficients and condi-  
933 tional probabilities from the mmvec for the Rhamnolipids (n=462 molecules). (c) Pyochelin  
proportions vs *P. aeruginosa* proportions.

934 Figure S4 : Negative log likelihood and prediction accuracy of mmvec. Tensorboard visu-  
935 alization of training error and cross-validation error of mmvec on the IBG dataset. Five  
different runs with differing initialization conditions are shown.

936 Figure S5: GNPS [32] job output. An example of job on the GNPS website with the job description and the downloadable output files from mmvec.

