

PAPER

# Variational inference for microbiome survey data with application to global ocean data

Aditya Mishra,<sup>1,6,\*</sup> Jesse McNichol,<sup>2</sup> Jed Fuhrman,<sup>2</sup> David Blei<sup>1,3</sup>  
and Christian L. Müller<sup>1,4,5</sup>

<sup>1</sup>Center for Computational Mathematics, Flatiron Institute, New York, NY, USA, <sup>2</sup>Department of Biological Sciences, University of Southern California, LA, USA, <sup>3</sup>Department of Statistics and Computer Science, Columbia University, NY, USA, <sup>4</sup>Computational Health Center, Helmholtz Zentrum München, Munich, Germany, <sup>5</sup>Department of Statistics, LMU München, Munich, Germany and <sup>6</sup>Department of Statistics, University of Georgia, Athens, GA, USA

\*Corresponding author. [aditya.mishra@uga.edu](mailto:aditya.mishra@uga.edu)

## Abstract

Linking sequence-derived microbial taxa abundances to host (patho-)physiology or habitat characteristics in a reproducible and interpretable manner has remained a formidable challenge for the analysis of microbiome survey data. Here, we introduce a flexible probabilistic modeling framework, VI-MIDAS (Variational Inference for Microbiome survey Data analysis), that enables *joint* estimation of context-dependent drivers and broad patterns of associations of microbial taxa abundances from microbiome survey data. VI-MIDAS comprises mechanisms for direct coupling of taxon abundances with covariates and taxa-specific latent coupling which can incorporate spatio-temporal information *and* taxon-taxon interactions. We leverage mean-field variational inference for posterior VI-MIDAS model parameter estimation and illustrate model building and analysis using Tara Ocean Expedition survey data. Using VI-MIDAS' latent embedding model and tools from network analysis, we show that marine microbial communities can be broadly categorized into five modules, including SAR11-, Nitrosopumilus-, and Alteromonadales-dominated communities, each associated with specific environmental and spatiotemporal signatures. VI-MIDAS also finds evidence for largely positive taxon-taxon associations in SAR11 or Rhodospirillales clades, and negative associations with Alteromonadales and Flavobacteriales classes. Our results indicate that VI-MIDAS provides a powerful integrative statistical analysis framework for discovering broad patterns of associations between microbial taxa and context-specific covariate data from microbiome survey data.

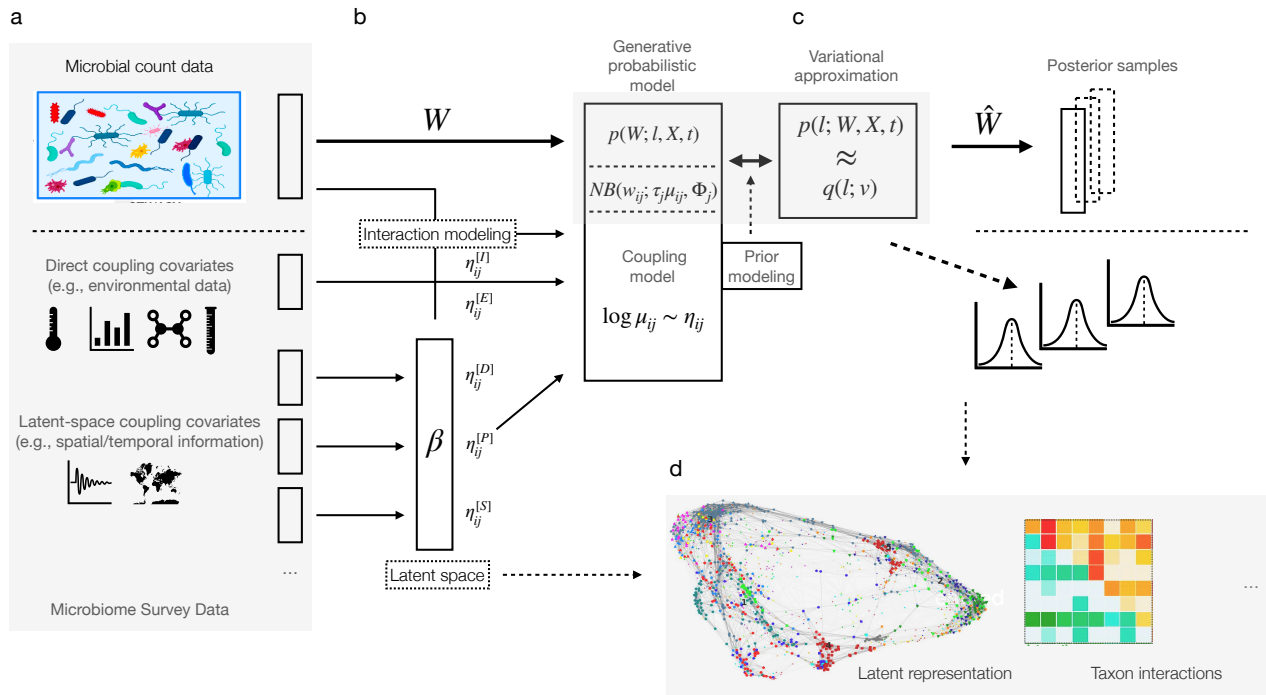
**Key words:** Microbiome; Probabilistic model; Association learning; Variational inference; Tara ocean expedition

## 7 Introduction

8 Microbial species are an integral part of life on earth.  
9 Ecosystems, ranging from the human gut to the global ocean,  
10 harbor trillions of bacteria, archaea, viruses, and fungi that  
11 take on essential functional roles and have developed intricate  
12 ecological relationships within their respective habitat. Over  
13 the past decades, advances in amplicon and metagenomics  
14 sequencing techniques [74, 54, 52, 70] and standardized  
15 experimental and bioinformatics workflows [63, 10, 9] have  
16 enabled the large-scale collection and dissemination of  
17 microbial survey data, including those from the seminal  
18 Human Microbiome Project [69], several gut-focused surveys  
19 [28, 64, 32, 45], the Earth Microbiome Project [25], and  
20 the Tara Ocean Expedition [67]. These surveys have reached  
21 a level of maturity and complexity that ultimately allow  
22 the estimation of statistical associations between microbial  
23 abundances, typically represented as compositional counts of  
24 Amplicon Sequence Variants (ASVs) or Operational Taxonomic  
25 Units (OTUs), and habitat properties [67, 7], biogeochemical  
26 processes [29], and/or host health status [23, 48]. This, in turn,

provides a starting point for deciphering and understanding 27  
the ecological and functional roles of different microbial clades 28  
in the ecosystem, nutrient and bio(geo)chemical dependencies, 29  
resource limitations of microbial growth, and the presence of 30  
ecological taxon-taxon interactions [18]. 31

Here, we introduce an integrative probabilistic modeling 32  
framework that is specifically tailored to microbiome survey 33  
data and enables joint estimation of habitat-dependent drivers 34  
and broad associations patterns of microbial taxa abundances 35  
(see Figure 1). Our approach, termed VI-MIDAS (Variational 36  
Inference for Microbiome survey Data analysis), models the 37  
observed taxon abundances by *simultaneously* learning taxon- 38  
specific latent representations that leverage the effects of host 39  
or environmental factors *and* taxon-taxon associations via an 40  
item-item interaction modeling *ansatz*, originally proposed for 41  
market basket analysis [61]. As such, VI-MIDAS seamlessly 42  
extends common statistical methods for microbiome data that 43  
only focus on either statistical abundance modeling [31, 39, 13, 44  
79, 75, 49] or microbial association estimation [37, 18, 77, 57, 45  
26]. 46



**Fig. 1.** Overview of the VI-MIDAS framework. **a.** VI-MIDAS integrates microbiome survey data in form of microbial abundance data  $W$ , host-associated, habitat or environmental data, and spatio-temporal information. **b.** Different data sources are coupled directly or indirectly through a latent space  $\beta$  to a generative model. An additional latent space taxon interaction model is included. The generative probabilistic model (e.g., Negative Binomial (NB) model) integrates covariate data via a coupling model. **c.** Variational approximation and mean-field estimation are used for Bayesian parameter estimation, resulting in posterior microbial abundance samples  $\hat{W}$  and model parameter distributions. **d.** Model components, such as estimated latent representation and taxon-taxon interactions, can be used for data understanding, visualization, and downstream analysis.

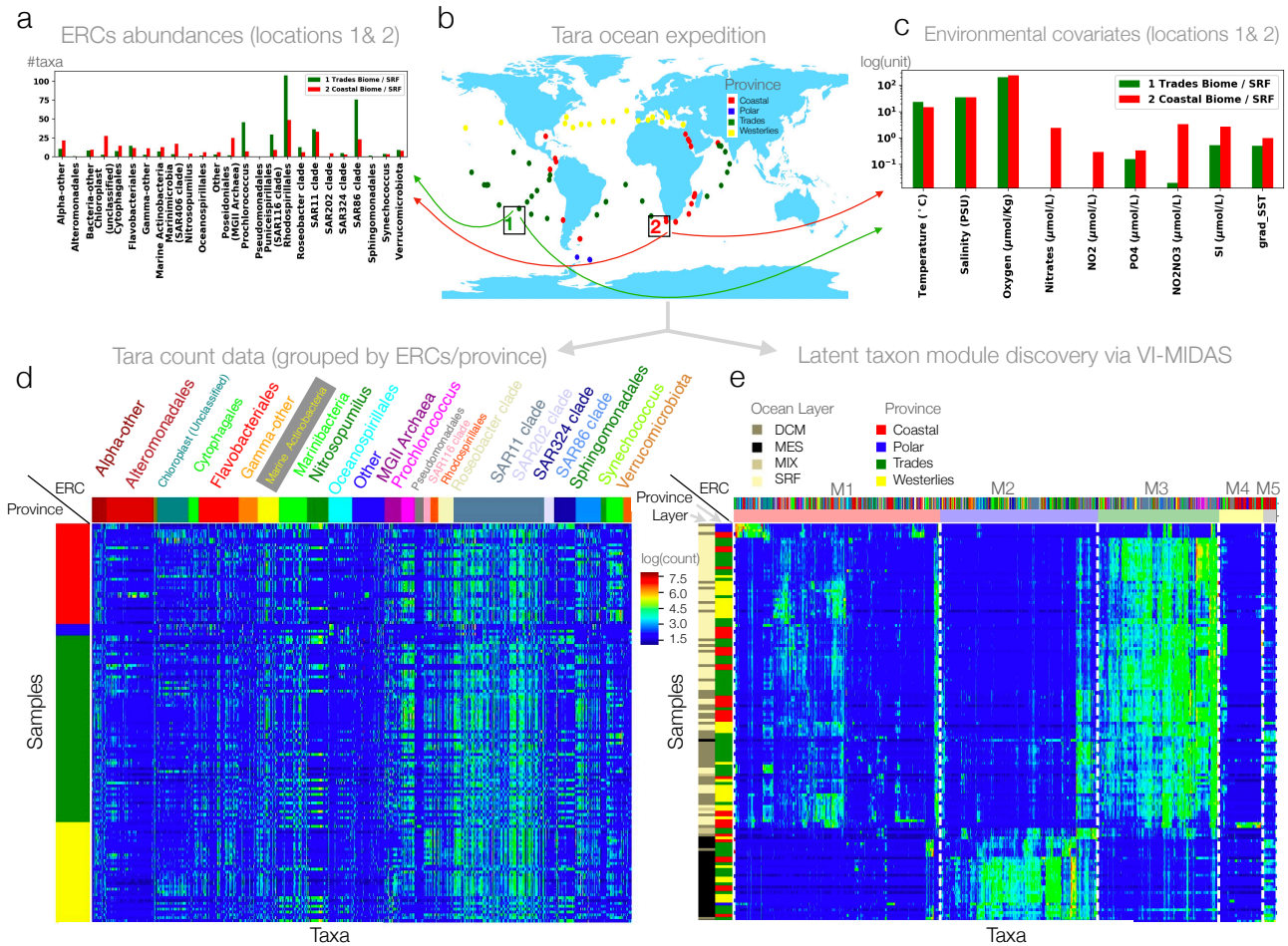
47 VI-MIDAS uses the parametric structure of the Negative  
 48 binomial distribution [46, 49] to account for the overdispersed  
 49 nature of the amplicon count data and comprises two main  
 50 model components: (i) a component that allows for full  
 51 adjustment of taxon abundances from a user-defined subset of  
 52 covariates and (ii) taxa-specific latent vectors that incorporate,  
 53 e.g., spatio-temporal or environmental covariates and taxon-  
 54 taxon interactions, thus providing a marginal characterization  
 55 of each taxon. We resort to mean-field variational inference  
 56 for parameter estimation of VI-MIDAS' intractable posterior  
 57 distribution [8], thus complementing other recent variational  
 58 approaches to microbiome data modeling, such as, e.g., Poisson  
 59 principal component analysis [14], microbiome dynamics  
 60 modeling [24], Dirichlet Multinomial modeling [30], multi-level  
 61 modeling [42], and microbiome ordination [78].

62 To illustrate the complete workflow of the VI-MIDAS  
 63 framework, we focus on integrative analysis of global marine  
 64 microbiome survey data. The ocean microbiome is of  
 65 fundamental importance for life on earth, being responsible for  
 66 about half of all primary production (i.e., the production of  
 67 chemical energy in organic compounds) and holds enormous  
 68 potential for climate remediation [50]. Several initiatives such  
 69 as the Tara Oceans Project [56] and the Simons CMAP [4]  
 70 provide well-structured sequencing data, biogeochemical and  
 71 environmental covariate data, and satellite-derived products  
 72 that are amenable to statistical analysis. Here, we re-analyze  
 73 Tara expedition data <sup>1</sup>, originally considered in [67] to study

74 the structure and function of the global ocean microbiome. 74  
 75 The expedition collected ocean water samples from 68 distinct 75  
 76 geographical locations at varying levels of depth. We will 76  
 77 make extensive use of this dataset to motivate and describe 77  
 78 the details of the VI-MIDAS framework as well as the 78  
 79 learned representations and associations of the global ocean 79  
 80 microbiome. 80

81 We start with an overview of the Tara Oceans data 81  
 82 under study, introduce the generative model components of 82  
 83 VI-MIDAS, and show how different data types enter the 83  
 84 modeling framework. We then give a high-level overview 84  
 85 of the variational parameter estimation procedure, including 85  
 86 the selection of VI-MIDAS' hyperparameters, such as the 86  
 87 choice of the priors and the dimensionality of the latent 87  
 88 representation. Following model parameter inference, we 88  
 89 illustrate how standard modularity analysis of VI-MIDAS' 89  
 90 learned latent representation of the Tara data identifies five 90  
 91 distinct groups of microbial consortia. We analyze the inferred 91  
 92 modules in terms of their composition of ecologically relevant 92  
 93 clades and discuss the derived module-specific environmental 93  
 94 and spatiotemporal signatures. Finally, we highlight the 94  
 95 emerging interaction pattern among ecologically relevant clades 95  
 96 and discuss the framework in the larger context of other 96  
 97 microbiome survey data. Further methodological details are 97  
 98 summarized in the Supplemental Material. Code for the 98  
 99 presented VI-MIDAS workflow is available at <http://github.com/amishra-stats/vi-midas> and requires minimal adjustment 99  
 100 to analyze other microbiome survey data. 100  
 101

<sup>1</sup> <http://ocean-microbiome.embl.de/companion.html>



**Fig. 2.** Illustration of the Tara ocean data: **a.** Taxon abundance profiles, agglomerated to expert-derived ecologically relevant classes (ERCs) for two samples (red and green, marked as 1 and 2 in Figure 2b). **b.** Tara ocean sample locations. **c.** Environmental features associated with the samples marked as 1 and 2 in Figure 2 (b); **d.** Abundance profiles  $\log(\mathbf{W} + 1)$  of  $q = 1379$  taxa at  $n = 139$  distinct locations with rows highlighting province of the sample and columns grouped by ERC. **e.** Abundance profiles clustered into five modules (M1-M5) as identified by modularity analysis of the latent space  $\beta$  (see Section Modularity analysis for more details). The dashed vertical lines separate the latent modules. The five microbial modules (M1-M5) comprise 524, 400, 307, 112 and 35 taxa/OTUs, respectively. The first column shows ocean depth layer, the second column the province indicator.

## 102 Materials and Methods

### 103 Tara ocean data and ecologically relevant taxa 104 re-classification

105 We consider the processed Tara expedition data, as  
106 provided at <http://ocean-microbiome.embl.de/companion.html>.  
107 The expedition collected water samples from 68 distinct  
108 geographical locations (Figure 2b) across different depths,  
109 resulting in  $n = 139$  distinct samples. Across these samples, the  
110 original data comprises microbial taxa abundances profiles of  
111 more than 35,000 bacterial taxa in form of metagenomic OTUs  
112 (mOTUs) (derived using the miTAGS framework [68]).

113 Here, we focus on the most abundant taxa by taking the  
114 union of all mOTUs that, in each individual sample, contribute  
115 to 40% of the total library size. This filtering allows us to  
116 cover the abundance profiles of the  $q = 1378$  taxa with the  
117 most significant variability and reduces the number of excess  
118 zero counts. To account for the highly variable sequencing  
119 depth across the samples, we normalize the abundance data  
120 with respect to the lowest library size via common-sum scaling  
121 [46]. Figure 2d shows the log-transformed abundance profiles

122  $\mathbf{W} \in \mathbb{R}^{n \times q}$ . Since the original taxonomic affiliations of  
123 the miTAGS are difficult to interpret, we next developed  
124 a partitioning of the selected taxa into ecologically relevant  
125 classes (ERCs). The original full taxonomy strings are too long  
126 to understand at a glance, and parsing by taxonomic level is  
127 not a good option since taxa vary widely in the depth of their  
128 annotations. For example, cyanobacteria should be annotated  
129 at the genus level or higher, but many other abundant but  
130 less described taxa do not have any taxonomic information at  
131 that level. We manually curated the data to provide a short  
132 relevant taxonomic indicator that provides a rough indicator  
133 of the ecological niche of an organism while remaining short  
134 enough to be interpreted at a glance. Some taxonomies have  
135 been altered to preserve the updated SILVA taxonomy (i.e.,  
136 Betaproteobacteria is now Burkholderiales). New SILVA 138  
137 [58] taxonomies have been used wherever possible (i.e., when  
138 the original ID was still in SILVA 138), but in cases where  
139 there was only the SILVA 108 taxonomic information, we have  
140 used our best guess. For example, if an organism had the same  
141 classification as other organisms in SILVA 108, we have often

**Table 1.** Environmental and spatiotemporal variables included in the VI-MIDAS model

Model component	Variables	Description
Environmental $\eta_{ij}^{[E]}$	<b>Environmental covariates</b>	Sea surface temperature (and its gradient), salinity, chlorophyll, nitrate, Nitrogen Dioxide, Phosphate, Silicon, and oxygen concentration.
Spatial (Depth)	<b>SRF</b>	Surface water layer; up to 5 m below the surface
	<b>DCM</b>	Deep chlorophyll maximum; approximately 17 m to 188 m below the surface; region below the surface with maximum chlorophyll concentration
$\eta_{ij}^{[D]}$	<b>MIX</b>	Subsurface epipelagic mixed layer; approximately 25 m to 150 m below the surface
	<b>MES</b>	Mesopelagic zone; approximately 250 m to 1000 m below the surface
Spatial (Longhurst Province)	<b>Polar biome</b>	Polar region in the northern and southern hemisphere characterized by low taxonomic diversity at all trophic levels.
	<b>Westerlies biome</b>	High-latitude region below the westerly winds
	<b>Trades biome</b>	Low-latitude region below the easterly trades characterized by high taxonomic diversity
$\eta_{ij}^{[P]}$	<b>Coastal biome</b>	Region in the upper part of the continental slope
Seasonal $\eta_{ij}^{[S]}$	<b>Q1, Q2, Q3, Q4</b>	Derived indicator of seasonal quarter when sample was taken (January to March; April-June; July-September; October-December)

142 given it the same name as its counterparts in SILVA 138. We  
143 present all our findings in terms of these 29 ERCs.

144 Each Tara sample also contains environmental and  
145 spatiotemporal information, including geolocation, the derived  
146 Longhurst province (biome) indicator, sampling date, ocean  
147 depth information (depth from sea surface), environmental  
148 covariates, such as, e.g., sea surface temperature (SST), and  
149 biogeochemical features such as salinity, chlorophyll, nitrate,  
150 and oxygen concentration (see Figure 2c for illustration).  
151 Table 1 summarizes the measured covariates and derived  
152 spatiotemporal indicator variables that are included in the  
153 VI-MIDAS framework and their corresponding mathematical  
154 representation.

## 155 Generative Modeling in VI-MIDAS

156 We seek to model the abundance profiles of  $q$  microbial taxa  
157 where we denote a single sample by the random variable  $\mathbf{w} \in \mathbb{R}^q$   
158 and the observed data from  $n$  samples by  $\mathbf{W} = [w_{ij}]_{n \times q} \in$   
159  $\mathbb{R}^{n \times q}$ . For concreteness, we illustrate model building and  
160 analysis using the Tara abundance profiles (see Figure 2(d))  
161 of  $q = 1378$  taxa but the modeling strategy is applicable to any  
162 multimodal microbiome survey.

### 163 Distributional model

VI-MIDAS posits that the overdispersed microbial count data  
 $\mathbf{W}$  are reasonably well modeled with the Negative Binomial  
distribution [11, 44, 48]. While other generative statistical  
modeling approaches are available, including the Dirichlet  
Multinomial (mixture) framework [31, 71], latent Dirichlet  
allocation [62], and Poisson distribution models [39, 5, 75],  
we found the Negative Binomial model to be an excellent  
choice for the Tara ocean data (see Figure S1 (b) of the  
Supplementary Material for the over-dispersion analysis). Using  
the Negative Binomial distribution with mean and dispersion  
parameterization [11], VI-MIDAS models the  $j$ th taxa in the  
 $i$ th sample as:

$$\begin{aligned}
 p(w_{ij}; \tau_j \mu_{ij}, \phi_j) &= \text{NB}(w_{ij}; \tau_j \mu_{ij}, \phi_j) \\
 &= \binom{w_{ij} + \phi_j - 1}{w_{ij}} \left( \frac{\tau_j \mu_{ij}}{\tau_j \mu_{ij} + \phi_j} \right)^{w_{ij}} \left( \frac{\phi_j}{\tau_j \mu_{ij} + \phi_j} \right)^{\phi_j}. \quad (1)
 \end{aligned}$$

164 Here, the mean parameter  $\tau_j \mu_{ij}$  is the product of a taxon-  
165 specific shape parameter  $\tau_j \in (0, 1)$  and the entry-specific

parameter  $\mu_{ij} \in \mathbb{R}^+$ . The parameter  $\phi_j \in \mathbb{R}^+$  is the taxon-  
specific dispersion parameter. Let us denote the dispersion and  
shape parameters for  $q$  outcomes by  $\Phi = [\phi_1, \dots, \phi_q]$  and  
 $\tau = [\tau_1, \dots, \tau_q]$ , respectively. The shape parameter  $\tau$  accounts  
for the disparity in abundance among microbial taxa. The  
generative model (1) of VI-MIDAS thus implies  $\mathbb{E}(w_{ij}) = \tau_j \mu_{ij}$   
and  $\text{Var}(w_{ij}) = \tau_j \mu_{ij} + \frac{\tau_j^2 \mu_{ij}^2}{\phi_j}$ . Consequently,  $\text{Var}(w_{ij}) >$   
 $\mathbb{E}(w_{ij})$ , thus making the parametric framework (1) suitable for  
modeling the overdispersed count data.

### 175 Modeling strategy and model components

176 One novelty in VI-MIDAS is the combination of ideas  
177 from generalized linear modeling [11] and compositional data  
178 analysis [2] to associate the microbial relative count data  
179 with spatiotemporal, environmental, and taxa information.  
180 Specifically, we model the log-transformed mean parameter  
181  $\boldsymbol{\mu} = [\mu_{ij}]_{n \times q}$  of the generative model (1) with two components,  
182 a consistent zero-aware geometric mean estimate  $t_i$  and a linear  
183 predictor  $\boldsymbol{\eta} = [\eta_{ij}]_{n \times q} \in \mathbb{R}^{n \times q}$  as follows:

$$\log \mu_{ij} = \log t_i + \eta_{ij}, \quad (2)$$

The sample-wise parameter  $t_i$  is estimated by a zero-aware  
geometric mean estimator, introduced in [16], which provides  
a principled approximation to the geometric means across all  
 $n$  samples in the presence of excess zeros. We detail the  
exact formulation of  $t_i$  and its approximation guarantees in  
Section 3.1 of the Supplementary Material. Including  $\mathbf{O} =$   
 $[\log t_1, \dots, \log t_n]$  as an offset term in the model is necessary  
since we do not have access to absolute microbial abundance  
data, thus requiring transforming the compositional data  
appropriately. The second term  $\boldsymbol{\eta}$  effectively models centered  
log-ratio (clr) transformed (rather than the original count)  
data and is the key component to couple habitat (or host)  
information to the microbial abundance profiles. VI-MIDAS  
introduces a novel decomposition of the component  $\boldsymbol{\eta}$  that  
allows the incorporation of three distinct coupling mechanisms:  
(i) a direct coupling term for covariates, (ii) an indirect coupling  
term for covariates via a latent space representation, and (iii)  
a latent taxon-taxon interaction term.

In our ocean application, the first component, denoted by  
 $\eta_{ij}^{[E]}$ , includes all relevant environmental attributes (see first  
row in Table 1). All spatiotemporal features, i.e., the Longhurst  
Province indicator, the Depth information, and the Seasonal

indicator (see second to last row in Table 1) are handled by the latent coupling term and are denoted by  $\eta_{ij}^{[D]}$ ,  $\eta_{ij}^{[P]}$ , and  $\eta_{ij}^{[S]}$ , respectively. Lastly, statistical associations among co-occurring taxa are included via a latent interaction term  $\eta_{ij}^{[I]}$ , leading the following model:

$$\eta_{ij} = \eta_{ij}^{[E]} + (\eta_{ij}^{[P]} + \eta_{ij}^{[D]} + \eta_{ij}^{[S]}) + \eta_{ij}^{[I]}. \quad (3)$$

The following paragraphs detail the parametric form of each of the components, the nature of the underlying covariate data, and their biological relevance.

Direct coupling of environmental features

Let us denote the  $p$  covariates in the direct coupling term by  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T = [x_{ij}]_{n \times p}$ . VI-MIDAS models the direct component for the  $j$ th taxa in the  $i$ th sample via

$$\eta_{ij}^{[E]} = \mathbf{x}_i^T \boldsymbol{\gamma}_{.j}. \quad (4)$$

with  $\boldsymbol{\gamma} = [\gamma_{ij}]_{p \times q} \in \mathbb{R}^{p \times q}$  denoting the matrix of all coefficients. For the Tara data, we opted to model  $\eta_{ij}^{[E]}$  using following  $p = 9$  covariates: sea surface temperature (SST) (and its gradient grad SST), salinity, chlorophyll, nitrate, nitrogen dioxide, phosphate, silicon, and oxygen concentration. All variables are mean-centered prior to incorporation into the model. In the original Tara analysis [67], temperature and oxygen have been identified as key drivers of taxonomic compositions. The VI-MIDAS analysis will allow a refined picture of the these general tendencies.

Latent space coupling of spatiotemporal features

VI-MIDAS offers a second mechanism for including variables of interest through latent space modeling. We denote  $q$  taxa-specific shared latent variables of size  $k$  by  $\boldsymbol{\beta} = [\beta_{ij}]_{k \times q} \in \mathbb{R}^{k \times q}$ . The size factor  $k$  is an application-specific hyperparameter that controls the expressiveness of the latent space. Features are then coupled to the latent space in a multiplicative fashion.

For the Tara data, we illustrate this mechanism by coupling all available spatial and temporal indicators to the latent space component. We first consider the  $r = 4$  primary provinces (or biomes): polar, Westerlies, coastal, and Trades [43]. We denote the model matrix indicating the  $r$  distinct regions of the  $n$  samples by  $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]^T \in \mathbb{R}^{n \times r}$  and connect it to the joint latent space via the coefficient matrix  $\boldsymbol{\alpha} = [\alpha]_{r \times k} \in \mathbb{R}^{r \times k}$ , leading to

$$\eta_{ij}^{[P]} = \mathbf{r}_i \boldsymbol{\alpha} \boldsymbol{\beta}_{.j}. \quad (5)$$

Similarly, the Tara data includes samples across  $b = 4$  ocean depths: surface water (SRF), deep chlorophyll maximum (DCM), the subsurface epipelagic mixed layer (MIX), and the mesopelagic zone (MES). We denote the depth indicator matrix of the  $n$  samples by  $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_n]^T \in \mathbb{R}^{n \times d}$  and connect it to the joint latent space via the coefficient matrix  $\boldsymbol{\delta} = [\delta]_{b \times k} \in \mathbb{R}^{b \times k}$ , leading to

$$\eta_{ij}^{[D]} = \mathbf{d}_i \boldsymbol{\delta} \boldsymbol{\beta}_{.j}. \quad (6)$$

Finally, by parsing the sampling dates at the different Tara locations, we can associate a temporal indicator with each sample. Here, we group the samples into  $m = 4$  seasons: the 1<sup>st</sup> (Q1, January-March), 2<sup>nd</sup> (Q2, April-June), 3<sup>rd</sup> (Q3, July-September), and 4<sup>th</sup> (Q4, October-December) yearly quarter, and construct the season indicator matrix  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^T \in$

$\mathbb{R}^{n \times s}$ . The coefficient matrix  $\boldsymbol{\vartheta} = [\vartheta]_{m \times k} \in \mathbb{R}^{m \times k}$  couples  $\mathbf{S}$  to the latent space  $\boldsymbol{\beta}$ , leading to

$$\eta_{ij}^{[S]} = \mathbf{s}_i \boldsymbol{\vartheta} \boldsymbol{\beta}_{.j}. \quad (7)$$

In summary, the coupling of the described features to a shared latent space via the coefficient matrices  $\boldsymbol{\alpha}, \boldsymbol{\delta}, \boldsymbol{\vartheta}$  allows to quantify to what extent spatiotemporal information influences each taxon's (latent) abundance after discounting the contribution of the environmental component.

Latent modeling of taxon-taxon associations

It is well-established that the abundances of species in an ecosystem are not only driven by environmental or spatiotemporal factors but also by interactions among the species themselves [41]. While discovering detailed ecological interactions among taxa, such as, e.g., competition, mutualism, or commensalism, is beyond the reach of coarse-grained statistical models, VI-MIDAS' latent space modeling offers a principled mechanism to assess the influence of taxa *co-occurrences* on their respective abundances. We achieve this by borrowing recent ideas from market basket analysis and adopt the so-called SHOPPER utility model for interaction analysis [61]. In SHOPPER, Ruiz et al. [61] proposed a probabilistic model based on the basket data from a supermarket to learn about the latent characteristic of each item and exchangeable/complementary interactions among items. The approach uses item-specific latent variables to define an item-item interaction component. Following their setup, the "interaction", or, in the biological context, association of the  $j$ th taxa with any  $m$ th taxa is given by  $\boldsymbol{\rho}_{.j}^T \boldsymbol{\beta}_{.m}$  where  $\boldsymbol{\rho} = [\rho]_{k \times q} \in \mathbb{R}^{k \times q}$  comprises length- $k$  latent variables for each of the  $q$  taxa. The entries of VI-MIDAS' interaction component  $\eta_{ij}^{[I]}$  for the  $j$ th taxon in the  $i$ th sample are thus given by

$$\eta_{ij}^{[I]} = \begin{cases} 0, & w_{ij} = 0 \\ \frac{1}{a_i - 1} \boldsymbol{\rho}_{.j}^T \sum_{m \neq j} \mathbf{1}_{w_{im} \neq 0} \boldsymbol{\beta}_{.m}, & w_{ij} \neq 0, \end{cases} \quad (8)$$

where  $a_i = \sum_{m=1}^q \mathbf{1}_{w_{im} \neq 0}$  is the total number of taxa present in the  $i$ th sample. Note that the interaction term  $\boldsymbol{\rho}^T \boldsymbol{\beta}$  is not symmetric. However, we can derive a symmetrized  $\mathbf{I} = [I_{i,j}] \in \mathbb{R}^{q \times q}$  with each entry being computed as:

$$I_{i,j} = (\boldsymbol{\rho}_{.i}^T \boldsymbol{\beta}_{.j} + \boldsymbol{\rho}_{.j}^T \boldsymbol{\beta}_{.i})/2 \quad (9)$$

This allows easier downstream network analysis of potentially *positive* (mutualistic) and *negative* (competitive) associations among the taxa, or in our case, among the ecologically relevant clades.

## Variational inference in VI-MIDAS

The generality and flexibility of VI-MIDAS poses a considerable challenge for fast and accurate model parameter estimation. We introduce a variational inference framework that makes estimation in VI-MIDAS feasible and illustrate its performance and parameter sensitivities using the Tara data. For ease of presentation, we summarize the key ingredients below and refer to the extensive Supplementary Information and the documented code base available at <https://github.com/amishra-stats/vi-midas> for details.

Bayesian model and variational approximation

We begin by denoting all (latent) parameters in the VI-MIDAS framework by  $\boldsymbol{\ell} = \{\boldsymbol{\alpha}, \boldsymbol{\vartheta}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\rho}, \boldsymbol{\tau}, \boldsymbol{\Phi}\}$  (see Table S1 of the

Supplementary Material). Given the microbial abundance data  $\mathbf{W}$ , the (direct) covariates  $\mathbf{X}$ , and the model parameters  $\ell$ , we integrate the generative model (1) into a Bayesian framework where the posterior distribution reads:

$$p(\ell; \mathbf{W}, \mathbf{X}, \mathbf{t}) = \frac{p(\mathbf{W}; \ell, \mathbf{X}, \mathbf{t})p(\ell)}{p(\mathbf{W}; \mathbf{X}, \mathbf{t})}, \quad (10)$$

where  $p(\mathbf{W}; \ell, \mathbf{X}, \mathbf{t}) = \prod_{i,j} p(w_{ij}; \tau_j \mu_{ij}, \phi_j)$  denotes the likelihood of  $\mathbf{W}$  and  $p(\ell) = p(\alpha)p(\delta)p(\beta)p(\gamma)p(\rho)p(\Phi)p(\tau)p(\vartheta)$  the prior distribution, respectively. To achieve good generalizability and interpretability of VI-MIDAS' overparameterized model, we place sparsity-inducing Laplace priors with scale parameter  $\lambda$  on each of the unconstrained latent variables in the set  $\{\alpha, \delta, \beta, \gamma, \rho, \vartheta\}$ . For example, the prior on  $\alpha$  reads  $p(\alpha) = \prod_{i,j} p(\alpha_{ij})$  with  $p(\alpha_{ij}) = \text{Laplace}(0, \lambda)$ . Furthermore, we place an inverse-Cauchy prior on the dispersion parameter  $\Phi$ , i.e.,  $p(\phi_j) = \text{inverse-Cauchy}(0, \nu)$  and  $p(\Phi) = \prod_j p(\phi_j)$ , and a Uniform(1,2) prior for the shape parameter  $\tau$ , i.e.,  $\tau_j \sim \text{Beta}(1,1)$  and  $p(\tau) = \prod_j p(\tau_j)$ . Choosing suitable hyperparameters for the priors will be discussed below.

In the high-dimensional setting, computing the posterior distribution is challenging because of the intractable form of the marginal distribution  $p(\mathbf{W}; \mathbf{X}, \mathbf{t})$  and the non-conjugate priors on the model parameters. Markov Chain Monte Carlo (MCMC) sampling provides a helpful paradigm for obtaining samples from the posterior distribution in the Bayesian framework. However, since MCMC lacks computational efficiency in large/high-dimensional problems, we use mean-field Variational Inference (VI) [34, 72, 8] and approximate the posterior with a variational posterior distribution of the latent variable  $\ell$ . Briefly, let  $q(\ell; \nu)$  be the variational posterior distribution with parameter  $\nu$ . VI approximates sampling of the posterior by minimizing the Kullback-Leibler (KL) divergence,

$$\min_{\nu} \text{KL}(q(\ell; \nu) || p(\ell; \mathbf{W}, \mathbf{X}, \mathbf{t}))$$

such that  $\text{supp}(q(\ell; \nu)) \subseteq \text{supp}(p(\ell; \mathbf{W}, \mathbf{X}, \mathbf{t}))$ . It can be shown that the above optimization problem simplifies to maximizing the evidence lower bound (ELBO) given by

$$\mathcal{L}(\nu) = \mathbb{E}_{q(\ell; \nu)}[\log P(\mathbf{W}, \ell; \mathbf{X}, \mathbf{t})] - \mathbb{E}_{q(\ell; \nu)}[\log q(\ell; \nu)], \quad (11)$$

which is a lower bound on the logarithm of the joint probability of the observations  $\log P(\mathbf{W}; \mathbf{X}, \mathbf{t})$  [34]. Replacing the joint distribution  $P(\mathbf{W}, \ell; \mathbf{X}, \mathbf{t})$  with a product of likelihood and prior distribution  $P(\mathbf{W}, \ell; \mathbf{X}, \mathbf{t}) = P(\mathbf{W}; \ell, \mathbf{X}, \mathbf{t})P(\ell)$  further simplifies the objective.

Model estimation, hyperparameter tuning, and posterior estimates

The non-convexity of the variational objective and the large number of model parameters require careful assessment of all aspects of model parameter estimation, hyperparameter tuning, and generalization capability. To estimate the parameters of the variational posterior distribution, we employ stochastic gradient descent within the automatic differentiation variational inference (ADVI) framework [36]. The key steps of ADVI are outlined in Algorithm 1 of the Supplementary Material. A prerequisite for model parameter estimation is the identification of suitable model *hyperparameters*. In VI-MIDAS, the key hyperparameters are the scale of the sparsity-inducing Laplace prior, the scale of the inverse-Cauchy prior, and the intrinsic dimensionality  $k$  of the latent

space  $\beta$ , respectively. VI-MIDAS tunes these parameters via random search (see Section 3.3 of the Supplementary Material for details) where the out-of-sample log-likelihood posterior predictive density (LLPD) is used for assessing optimality of the hyperparameters [22]. Due to the non-convexity of the objective and the use of stochastic optimization in VI initialization, we further evaluate the suitability of hyperparameter setting across fifty random initializations and select the hyperparameter set leading to the best averaged LLPD (see Section 3.5 of the Supplementary Material). The computational workflow is implemented in Python using the probabilistic programming language Stan [12] and is available in the GitHub repository (<https://github.com/amishra-stats/vi-midas>).

After hyperparameter tuning, we re-estimate the final model parameters on complete data. VI-MIDAS generates  $m = 100$  posterior samples of each of the latent variables in the set  $\ell$  and estimates the model parameters  $\ell$  using the mean of the samples from the variational posterior distribution. The model fit is numerically evaluated using the posterior predictive check [60, 22] on the full data. The procedure requires generating  $m$  posterior samples, denoted by the random variables  $\mathbf{W}^{rep} = [w_{ij}^{rep}] \in \mathbb{R}_+^{n \times q}$ , and then computing the p-value of the model fit as p-value :=  $p(t(\mathbf{W}^{rep}) < t(\mathbf{W}))$ , where  $t$  is the test statistic. In practice, we use the test statistics  $t(\mathbf{W}^{rep}) = \mathbf{E}(\log p(\mathbf{W}^{rep} | \ell))$  and  $t(\mathbf{W}) = \mathbf{E}(\log p(\mathbf{W} | \ell))$ .

## Results

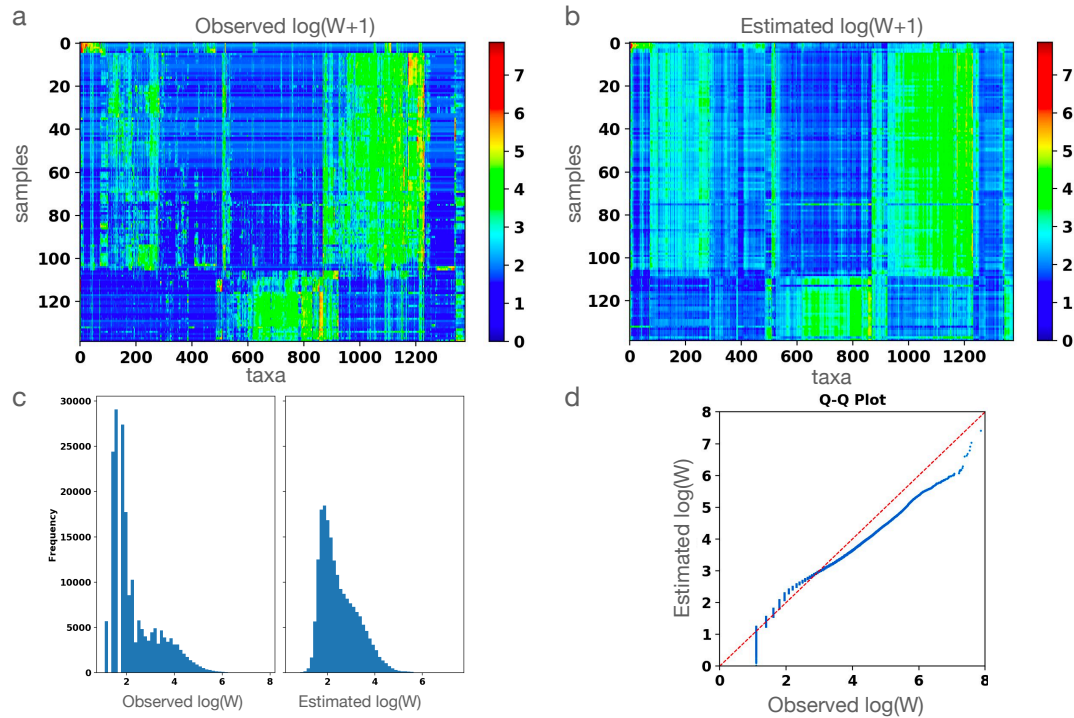
VI-MIDAS recapitulates broad statistical patterns of the observed species abundances

VI-MIDAS' hyperparameter tuning revealed that the setting  $k = 200$ ,  $\lambda = 0.246$ , and  $\nu = 0.10063$  achieved the highest average LLPD of 3.332 on the Tara data (see Figure S7 in the Supplementary Material). For this setting, a posterior predictive check on the generated samples achieved a p-value = 0.53. We thus fail to reject the null hypothesis that the posterior samples are different from the observed  $\mathbf{W}$ . Figure 3a and 3b the observed and estimated abundance profiles (averaged over  $m = 100$  samples), respectively. Figure 3c shows the count histograms of data and model (pooled across all samples and species), and Figure 3d the Q-Q plot. We observe that, apart from the low-abundance tail of the distribution, VI-MIDAS broadly recapitulates the statistical abundances patterns across all samples and species.

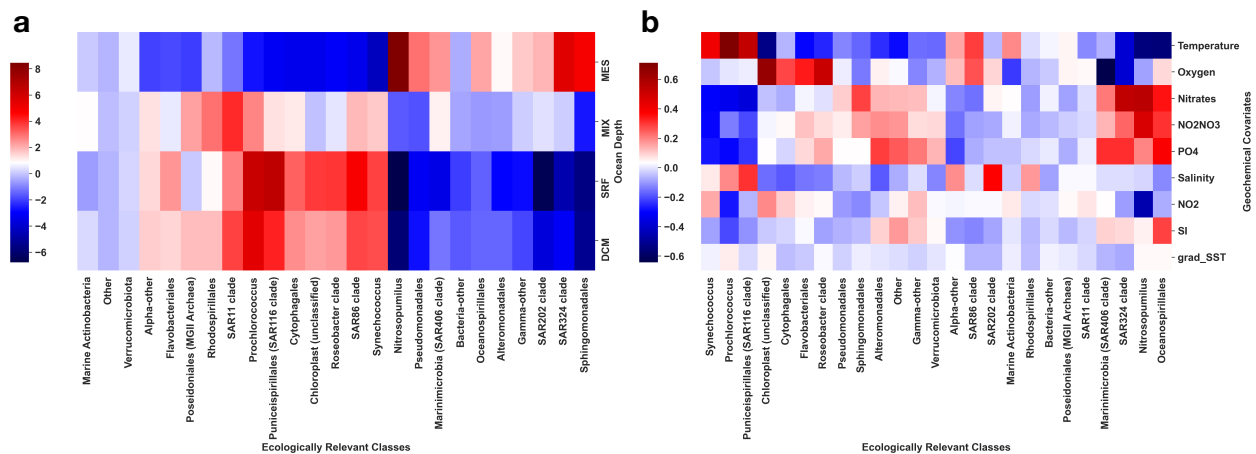
VI-MIDAS identifies depth and environmental features as main drivers

We next assessed the contribution of each model component toward explaining the species abundance patterns in the Tara data. The modularity of the VI-MIDAS framework facilitates an "ablation" study (see Section 3.4 of the Supplementary Material) where each model component is excluded, followed by a re-evaluation of the out-of-sample LLPD. Table S4 (see Supplementary Materials) shows the LLPDs of the full model and the model after ablation of the environmental(E), province(P), ocean depth(D), seasonality(S), and latent interaction (I) component, respectively.

Firstly, the ablation study confirmed that all components helped improve model generalization since every ablated model has reduced out-of-sample LLPDs. While the seasonality component(S) shows comparatively little influence on explaining the abundance pattern in the current model, as previously observed for this dataset [67], the out-of-sample LLPD is reduced the most when the ocean depth(D) component is



**Fig. 3.** Comparison of observed abundances and VI-MIDAS posterior samples: **a.** Heatmap showing the abundance profile  $\log(W + 1)$  of 1378 species for  $n = 139$  samples. **b.** Expected value of the abundance using the hyperparameter corresponding to best model fit. **c.** Histograms of observed and estimated species abundances. **d.** Q-Q plot comparing the observed and estimated abundance profile of the species.



**Fig. 4.** Summary of the estimated average effect sizes of the influence of **a.** ocean depth (VI-MIDAS model component  $\delta\beta$ ) and **b.** environmental covariates (VI-MIDAS model component  $\gamma$ ) on all ecologically relevant classes (ERCs).

441 ablated (LLPD=-3.3882). This reflects the well-known depth  
 442 stratification of marine species between the sunlit ocean and  
 443 aphotic deep ocean ecosystems. Figure S3 in the Supplementary  
 444 material illustrates the learned depth stratification across all  
 445 taxa, as reflected in the component  $\delta\beta$ . The environmental  
 446 component was identified as the second most important  
 447 component with an LLPD reduction of -3.3554.

448 Figure 4 summarizes the estimated effects  $\delta\beta$  of the  
 449 ocean depth features and the environmental effects  $\gamma$  on  
 450 the abundance of species aggregated into ERCs, respectively.  
 451 The ocean depth summary (Fig. 4a) reveals three distinct  
 452 sets of occurrence patterns for two different groups of ERCs.

453 One group (right most in Fig. 4a) comprises ERCs such as  
 454 Nitrosopumilus, Pseudomonadales, SAR 324 clade, and  
 455 Sphingomonadales which thrive in the Mesopelagic (MES)  
 456 zone. A second group includes species like Prochlorococcus,  
 457 SAR 116 clade, and Synechococcus, which flourish within  
 458 the ecosystem of the ocean's Deep Chlorophyll Maximum  
 459 (DCM) and Surface Mixed Layer (SRF) zones. The third  
 460 group comprises marine Actinobacteria, Verrucomicrobiota,  
 461 and others that show no dependence on depth. A summary  
 462 of geochemical features highlights temperature (the top  
 463 row in Fig. 4b) as the primary positive factor influencing  
 464 the abundance of Synechococcus, Prochlorococcus, and

465 Puniceispirillales (SAR116 clade). Oxygen concentration  
466 emerges as the main positive driver of abundance for  
467 Cytophagales, Flavobacteriales, and Roseobacter clades, while  
468 Nitrates, Nitrites, and Phosphate are identified as key drivers  
469 for the SAR324 clades, Nitrosopumilus, and Oceanospirillales  
470 (four right most columns in Fig. 4b). The estimated patterns  
471 broadly recapitulate known biology about ocean microbial  
472 ecosystems.

473 VI-MIDAS reveals five latent microbial sub-communities

474 The generative model (1) of VI-MIDAS includes the taxon-  
475 specific latent variables  $\beta \in \mathbb{R}^{k \times q}$  to integrate spatiotemporal  
476 features and taxon-taxon associations. For the Tara data,  
477 VI-MIDAS' hyperparameter tuning scheme identified  $k =$   
478 200 as best latent dimension. After model estimation, the  
479 resulting  $k$ -dimensional latent vectors can be thought of as  
480 representing the hidden *marginal* characteristics of each of  
481 the  $q$  taxa after discounting spatiotemporal and species-species  
482 association effects, and adjusted for environmental covariates.  
483 The latent space representation thus provides an excellent  
484 opportunity to partition the different taxa into coherent sub-  
485 groups (or modules) that likely reflect functionality or niche  
486 occupation in the global ocean, independent of environmental,  
487 taxonomic or phylogenetic relatedness.

488 To quantify similarity between microbial taxa in the latent  
489 space, we first computed cosine distances of all pairs of the  
490  $q$  latent vectors. This particular choice of distance allows  
491 us to bypass the non-identifiability issue of the parameter  
492  $\beta$ . We used the resulting distance matrix to construct a  $k$ -  
493 nearest neighbors graph ( $k_{nn} = 10$ ). Figure 5 shows the  
494 latent space embedding using a force-directed layout of the  
495  $k$ -nn graph. The latent space representation reveals several  
496 distinct microbial sub-communities, dominated by a few ERCs,  
497 including one sub-community dominated by Prochlorococcus  
498 and SAR11 clades and one dominated by Nitrosopumilus.  
499 We next performed Clauset-Newman-Moore greedy modularity  
500 analysis of the nearest neighbor graph [15] and identified five  
501 distinct modules in the latent space (see M1-M5 in Fig. 5 with  
502 top five ERCs highlighted and color-coded). Module 1 (M1)  
503 comprises Flaviobacteriales, SAR86 clades, and the Chloroplast  
504 class. SAR11 clade, SAR86 clade, and Flavobacteriales are  
505 heterotrophs with functional similarity in oxidizing carbon in  
506 the ocean [3]. Both SAR86 clade and SAR11 clade follow  
507 a similar seasonal pattern (in the Bermuda Atlantic Time  
508 Series oceanographic stations) and coexist in oligotrophic  
509 regions with less nutrient supply [73]. Module 2 (M2) includes  
510 Nitrosopumilus, Marinimicrobia, and SAR324 clades. Existing  
511 literature supports that SAR11 clade (a subgroup of a species),  
512 Marinimicrobia, and MGII Archaea are more abundant in deep  
513 sea water [76]. Module 3 (M3) comprises Prochlorococcus,  
514 SAR11, Marine Actinobacteria, and SAR86 clades, among  
515 others, all comprising dominant taxa of the sunlit ocean. The  
516 two smallest modules 4 and 5 (M4 and M5, respectively) are  
517 dominated by Alteromonadales and are separating M2 from M1  
518 and M3. Interestingly, Module 4 also comprises Synechococcus  
519 species. This module thus hints at the known metabolic  
520 dependency of certain Alteromonadales taxa on Synechococcus  
521 (a photoautotroph) [80]. Although the latent representation  
522 does separate the majority of ERCs into distinct subgroups, we  
523 nonetheless observe that taxa of certain ERCs are spread out  
524 over the latent space, indicating different niche specialization.  
525 For instance, the SAR11 clade, one of the most abundant  
526 marine microbial taxa, is present in three different modules.  
527 Likewise, taxa in the SAR86 clade are present in both modules

M1 and M3. For ease of identification, Table S3 summarizes  
each module in terms of the composition of the ERCs and their  
abundance.

Global associations between biogeography and latent microbial  
sub-communities

VI-MIDAS' integrative model also enables a quantitative  
description of the identified microbial sub-communities in  
terms of the direct and indirect coupling covariates. Figure 6  
illustrates how the compositions of ERCs in each of the five  
modules are related to the most important environmental and  
spatial covariates.

Using the mean of the posterior sample from the VI-MIDAS  
model, we used the estimated  $\gamma$  as the effect sizes of the  
environmental features  $\mathbf{X}$ ,  $\delta\beta$  as effect sizes of depth, and  $\alpha\beta$  as  
the effect sizes of the  $r$  provinces, respectively Figure 6 reports  
the average effect sizes of association to the four modules.

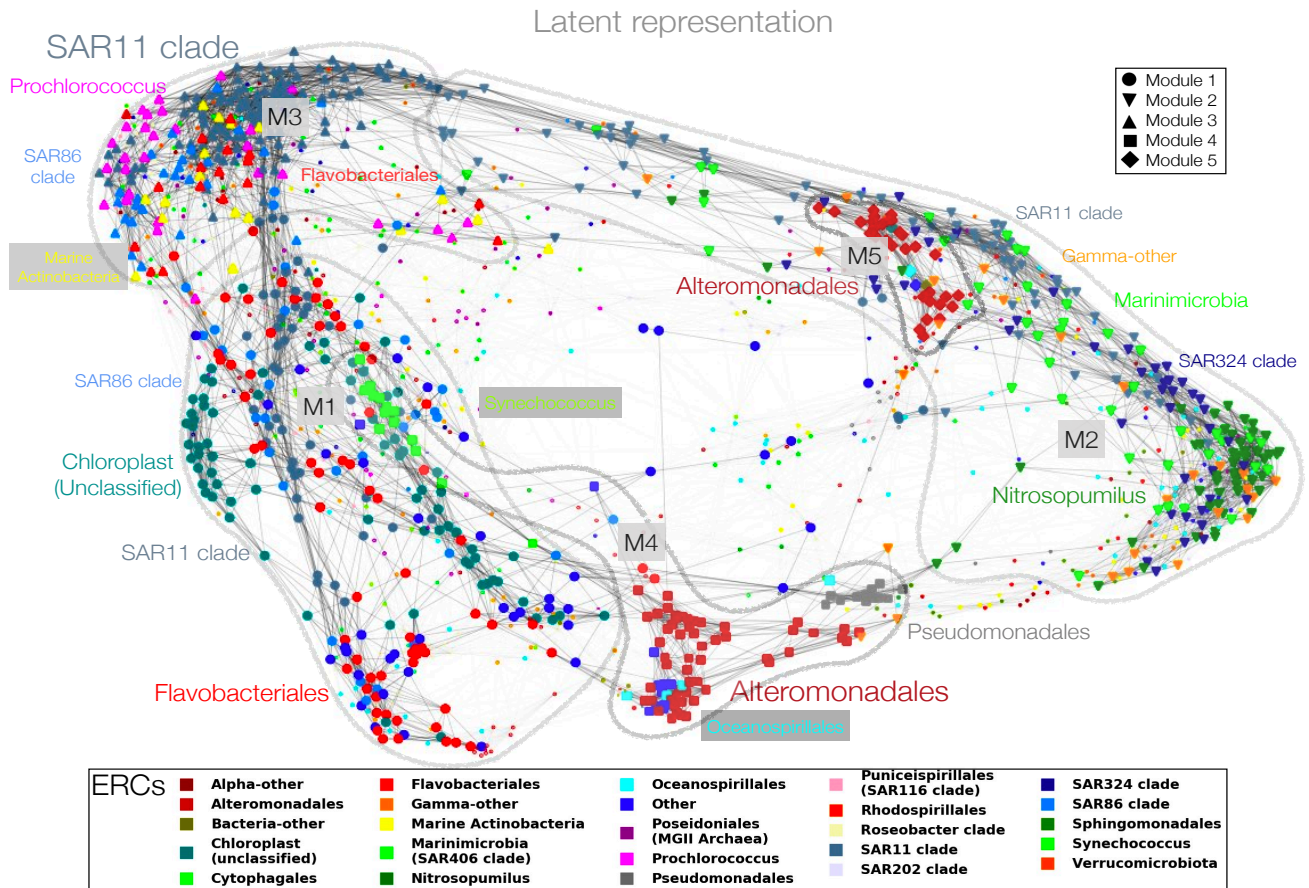
The module M1 represents taxa coexisting in the SRF  
and DCM zone of the ocean. The abundance of taxa  
in the module is associated with a higher concentration  
of oxygen,  $\text{PO}_4$ , and  $\text{NO}_2\text{NO}_3$  and lower temperature  
and salinity. In addition to representing the taxa SAR11  
clade, SAR86 clade, Chloroplast, and Flavobacteriales, the  
module also includes Synechococcus, Oceanospirillales, and  
Poseidoniales. Synechococcus is a unicellular prokaryotic  
autotrophic picoplankton that participates in the marine  
ecosystem as a primary producer via photosynthesis. Similarly,  
Chloroplast sequences are a signature of eukaryotic phytoplankton,  
though their host eukaryote is not identified in the TARA  
Oceans dataset. The presence of both taxa in M1 thus  
is consistent with environments that have higher oxygen  
concentrations due to photosynthesis and gas exchange with  
the atmosphere.

Module M2 mainly represents the species coexisting in the  
MES zone (200 m to 1000 m) of the ocean (see Figure 2 (e)).  
M2 almost exclusively represents the ERCs Nitrosopumilus  
and SAR324 clade. The abundance of the species in the  
group is associated with a lower concentration of oxygen and  
temperature, and higher concentrations of nitrates,  $\text{PO}_4$ , and  
 $\text{NO}_2\text{NO}_3$ . In the oxygen-depleted environment, Nitrosopumilus  
survives by oxidizing ammonia to nitrite, confirming the  
observed association pattern [6]. Marinimicrobia (SAR406  
clade) in groups M1 and M2 allow us to distinguish subgroups  
of species that can survive in both deep and shallow water [76].

Module M3 comprises the highest mean abundance of  
all taxa is highest, primarily representing the taxa SAR11  
clade, SAR86 clade, and Prochlorococcus (cyanobacteria).  
The abundance of the species in the group is positively  
associated with depth indicators {DCM, MIX, SRF} and  
negatively associated with MES. Among the geochemical  
factors, temperature, salinity, and oxygen concentration are  
positively associated, whereas the concentration of nitrates,  
 $\text{PO}_4$ , and  $\text{NO}_2\text{NO}_3$  is negatively associated with the taxa.

Module M4 primarily represents Alteromonadales (Proteobacteria) and some Pseudomonadales (Proteobacteria) and Synechococcus. Their abundance is associated with factors such as lower salinity and higher oxygen concentration. Module M5 also primarily represents Alteromonadales. Based on its association with the ocean depth indicators and geochemical features, we conclude that these taxa can survive in a deep-sea environment characterized by lower temperatures and oxygen concentrations. Associative patterns of Alteromonadales in M4 and M5 differ significantly, suggesting distinct ERC sub-groups that populate different niches.





**Fig. 5.** Low-dimensional embedding of the latent representation  $\beta$  using a  $k$ -nearest-neighbor ( $k_{nn} = 10$ ) graph of cosine distances. Modularity analysis reveals five distinct graph modules. We highlight 825 out of a total of 1378 taxa, comprising the top five ERCs (color-coded) in each of the five modules (see main text for further information).

591 Positive and Negative interactions among ERCs  
 592 VI-MIDAS includes a mechanism for learning microbial  
 593 interactions adjusted for direct (here, environmental) covariates.  
 594 Contrary to prominent (partial) correlation-based methods  
 595 [21, 38], VI-MIDAS follows the SHOPPER utility model [61]  
 596 and quantifies pairwise interactions  $I_{ij}$  between any two taxa  $i$   
 597 and  $j$  in terms of the latent variables  $\rho$  and  $\beta$  (see Eq. 9).

598 To get a high-level view of the estimated interactions,  
 599 we aggregated the adjacency matrices of significant positive  
 600 and negative interactions among taxa by ERCs (for a more  
 601 detailed view of the most significant taxon-level interactions,  
 602 we refer to Section 4 of the Supplementary Materials). Figure 7  
 603 illustrates the aggregated positive (lower triangle) and negative  
 604 (upper triangle) interactions among ERCs. The diagonal entry  
 605 highlights the maximum of the two types of interactions to  
 606 avoid confusion (see also Section 4 of the Supplementary  
 607 Materials for the matrix of ratios between positive and  
 608 negative interactions). We observe that SAR11 clade and  
 609 Rhodospirillales form positive interactions with almost all  
 610 other ERCs. SAR11 clade and Rhodospirillales belong to the  
 611 Alphaproteobacteria phylum that play a critical role in carbon  
 612 and nitrogen fixation [40, 51], potentially explaining the large  
 613 number of interactions. However, members of the SAR11  
 614 clade also form many negative interactions with other ERCs.

Alteromonadales exhibits primarily negative interactions with  
 other ERCs (the strongest one with SAR11).

## Discussion

In recent years, multimodal and multi-omics microbiome survey  
 data have emerged for a wide range of microbial habitats [68,  
 67, 47, 27, 65, 1]. These data collections hold the promise to  
 describe and understand the functional interplay between the  
 underlying microbial ecology and the host or the environment  
 the microbiota resides in. Learning interactions among species  
 and habitat characteristics from observational data remains,  
 however, a challenging problem. To this end, we have proposed  
 VI-MIDAS (Variational Inference for Microbiome survey Data  
 analysis), a flexible and efficient probabilistic framework for  
 microbiome survey data analysis.

VI-MIDAS uses the negative binomial distributional  
 framework in combination with a principled centering  
 transformation to model overdispersed amplicon abundance  
 data and comprises three mechanisms to integrate concomitant  
 covariate data into the generative model: (i) a direct coupling  
 mechanism, (ii) an indirect latent coupling mechanism, and  
 (iii) a latent interaction term. These terms are linearly linked  
 to the probability distribution's mean parameter. Because of  
 the intractable form of the marginal distribution of data, we

638 apply mean-field variational inference framework to learn an  
639 approximate posterior distribution of the parameters.

640 VI-MIDAS is available in Python and uses the probabilistic  
641 programming language Stan [12]. The implementation  
642 is available on GitHub (<https://github.com/amishra-stats/vi-midas>). The repository also includes Python scripts and  
643 Jupyter notebooks for VI-MIDAS' three-stage parameter  
644 estimation framework: hyperparameter tuning, component  
645 contribution analysis, and sensitivity analysis.

647 To illustrate the VI-MIDAS modeling and analysis workflow,  
648 we have used data from the global Tara expedition [67],

649 connecting the available spatiotemporal and environmental  
650 characteristics with generative modeling of the amplicon count  
651 data. To ease interpretability, we also grouped the amplicon-  
652 derived taxa into expert-annotated ecologically relevant classes  
653 (ERCs) which may be of independent interest for the analysis  
654 of other marine sequencing data. Focusing on the  $q = 1378$   
655 most abundant taxa representing 23 ERCS, we integrated  
656 the geochemical data using the direct coupling mechanism,  
657 effectively removing influence of common environmental factors  
658 such as temperature, salinity, and elemental compositions on  
659 microbial abundances. The remaining spatiotemporal features,

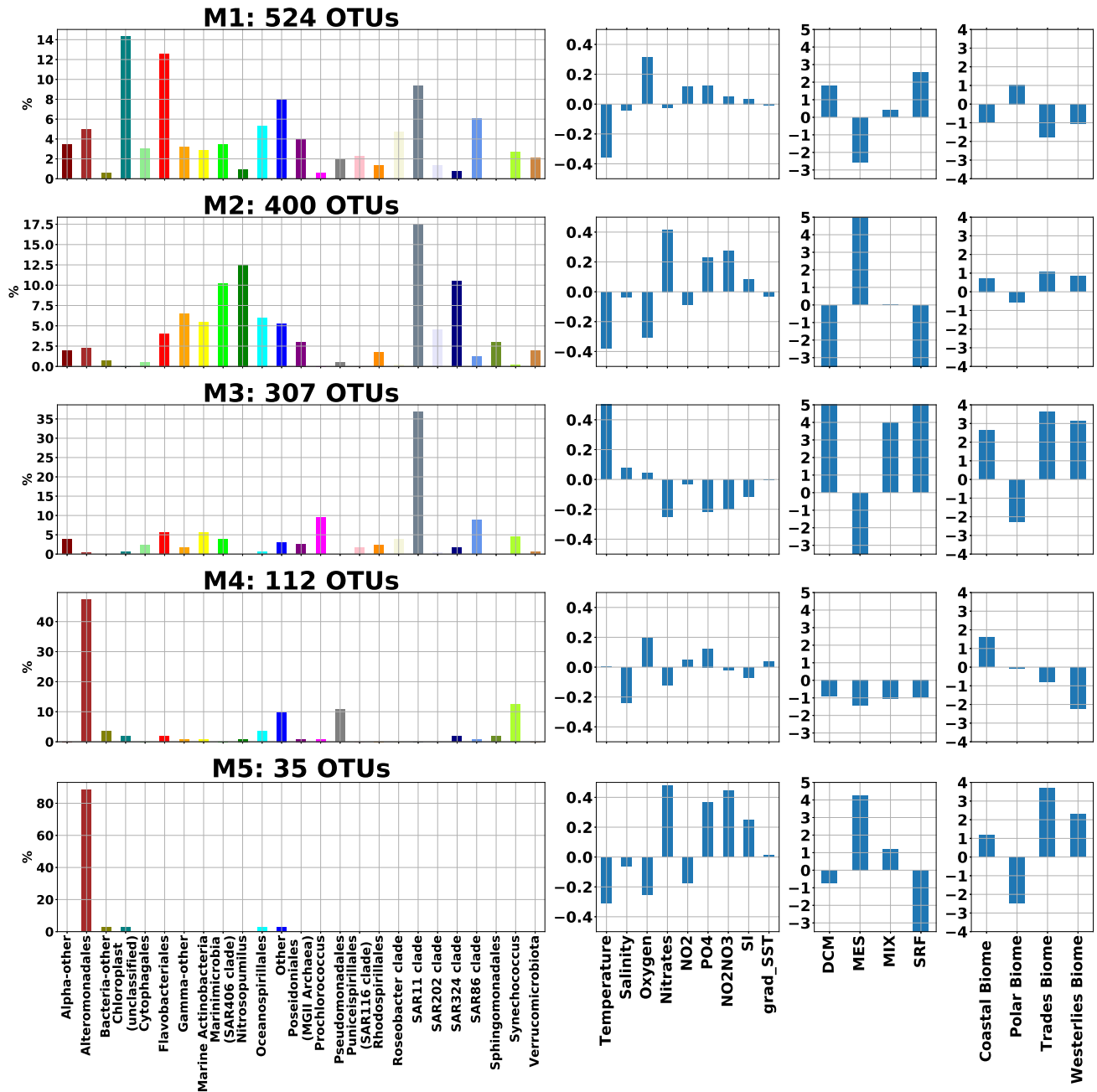
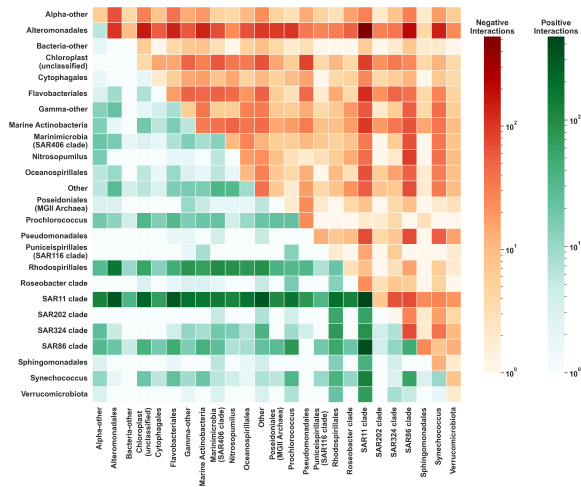


Fig. 6. Global associations between biogeography and covariates: Each row presents the average effect size of the association between the microbial abundances of taxa in a module (M1-M5) to the geochemical features, ocean depth and province/location (from left to right). A module (leftmost) is shown as the composition (in %) of the ERCS. Each module comprises different number of taxa {524, 400, 307, 112, 35}, respectively. Modules M1-M3 cover the majority of taxa, and M4-M5 two smaller Alteromonadales-dominated sub-communities.



**Fig. 7.** Summary of taxonomic interactions: The adjacency matrices of significant positive and negative interactions among taxa are grouped and aggregated by their ERCs type. Interactions summary by the ERCs types. Lower triangle reports positive interactions, the upper triangle reports negative interactions. Diagonal entries show the maximum of either (positive or negative) self-interaction.

component and link all concomitant features to the latent space representation, or alternatively, remove the latent representation altogether and directly adjust for all covariates. We will explore such modifications in future studies. Moreover, while we chose the Negative Binomial model as base distribution for the most abundant taxa, the variational formulation lends itself to other statistical models for microbial count data, including zero-inflated or hurdle- type extensions of the Negative Binomial model [19] or the Dirichlet-Multinomial model [30, 53]. Finally, in its current state, VI-MIDAS is built on Stan [12] with tailored Python code for optimization, model selection, and analysis. The advent of extensive statistical packages in modern deep learning tools, such as Tensorflow distributions [17] or PyTorch [55], may enable efficient porting of VI-MIDAS into these general-purpose ecosystems. Paired with variational inference tools [35], would potentially allow for faster model adaptation and alternative optimization routines.

In summary, VI-MIDAS provides a novel probabilistic framework for learning environment- or host-specific feature associations, latent species characterization, and species-species interactions from microbiome survey data. With minimal adjustment, the framework is readily available for the analysis of other large-scale survey data, including gut microbiome surveys [33, 45, 20], thus representing a potentially valuable general-purpose tool for the integrated analysis of modern microbiome data collections.

## Data availability

We have used microbial species abundance data from the Tara Ocean Expedition, available at (<http://ocean-microbiome.embl.de/companion.html>).

## Code availability

The source code required to reproduce the results in this article is freely available at (<https://github.com/amishra-stats/vi-midas>).

## Competing interests

No competing interest is declared.

## References

1. The integrative human microbiome project. *Nature*, 569(7758):641–648, 2019.
2. J. (John) Aitchison. *The statistical analysis of compositional data*. Blackburn Press, Caldwell, N.J., 2003.
3. Montserrat Aldunate, Rodrigo De la Iglesia, Anthony D Bertagnolli, and Osvaldo Ulloa. Oxygen modulates bacterial community composition in the coastal upwelling waters off central Chile. *Deep Sea Research Part II: Topical Studies in Oceanography*, 156:68–79, 2018.
4. Mohammad D Ashkezari, Norland R Hagen, Michael Denholtz, Andrew Neang, Tansy C Burns, Rhonda L Morales, Charlotte P Lee, Christopher N Hill, and E Virginia Armbrust. Simons collaborative marine atlas project (simons cmap): An open-source portal to share, visualize, and analyze ocean data. *Limnology and Oceanography: Methods*, 19(7):488–496, 2021.

including season, ocean province, and depth, as well as species-species associations are integrated through the latent coupling and interaction mechanism, thus delivering a latent species representation, adjusted for the influence of all available covariates. The learned VI-MIDAS’ model thus not only provides a convincing generative count model for the Tara data but also allows integrated statistical analysis of covariate feature effects and taxa abundances.

Modularity analysis of the similarity network of VI-MIDAS’ latent species representation revealed that the majority of taxa ( $\approx 1200$ ) can be categorized into three global microbial communities (M1-M3 in Figure 5), including a low-temperature/high-oxygen community (M1), dominated by Flavobacteriales and the Chloroplast ERC, a mesopelagic community (M2) dominated by SAR11, SAR324, and Nitrosopumilus, and a high-temperature community (M3) dominated by SAR11 and Prochlorococcus, the later of which is the most abundant clade in the oligotrophic subtropical and tropical oceans (see e.g., [66] and references therein). Furthermore, our analysis suggests two distinct Alteromonadales-dominated communities that show different depth and province dependencies (M4-M5) (see Figure 6 for further global associations overview). It is noteworthy that Alteromonadales also play a pivotal role in the latent interaction analysis, showing widespread negative associations with other ERCs. We posit that the potentially distinct role of Alteromonadales in the global ocean might be of interest for follow-up analysis on other data sets, including recent data on the global mesopelagic zone [59].

While our ablation study showed evidence that all VI-MIDAS components for the Tara data contribute to the quality of the generative model, the model is just one of several available alternatives. For covariate inclusion, we deliberately chose to directly adjust the microbial abundances for geochemical covariates to better carve out “hidden” relationships among the species. Nonetheless, the VI-MIDAS framework naturally enables other model constructions. For instance, one could have removed the direct coupling

5. Michael B. Sohn and Hongzhe Li. A glm-based latent variable ordination method for microbiome samples. *Biometrics*, 74(2):448–457, 2018.
6. Viswanathan Baskaran, Prasanna K Patil, M Leo Antony, Satheesha Avunje, Vinay T Nagaraju, Sudeep D Ghate, Suganya Nathamuni, N Dineshkumar, Shankar V Alavandi, and Kizhakedath K Vijayan. Microbial community profiling of ammonia and nitrite oxidizing bacterial enrichments from brackishwater ecosystems for mitigating nitrogen species. *Scientific reports*, 10(1):1–11, 2020.
7. Jacob Bien, Xiaohan Yan, Léo Simpson, and Christian L Müller. Tree-aggregated predictive modeling of microbiome data. *Scientific Reports*, 11(1):1–13, 2021.
8. David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
9. Evan Bolyen, Jai Ram Rideout, Matthew R Dillon, Nicholas A Bokulich, Christian C Abnet, Gabriel A Al-Ghalith, Harriet Alexander, Eric J Alm, Manimozhayan Arumugam, Francesco Asnicar, et al. Reproducible, interactive, scalable and extensible microbiome data science using qiime 2. *Nature biotechnology*, 37(8):852–857, 2019.
10. Ben J Callahan, Kris Sankaran, Julia A Fukuyama, Paul J McMurdie, and Susan P Holmes. Bioconductor workflow for microbiome data analysis: from raw reads to community analyses. *F1000Research*, 5, 2016.
11. A Colin Cameron and Pravin K Trivedi. *Regression analysis of count data*, volume 53. Cambridge university press, 2013.
12. Bob Carpenter, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of statistical software*, 76(1), 2017.
13. Jun Chen and Hongzhe Li. VARIABLE SELECTION FOR SPARSE DIRICHLET-MULTINOMIAL REGRESSION WITH AN APPLICATION TO MICROBIOME DATA ANALYSIS 1. *Ann. Appl. Stat.*, 7(1):418–442, 2013.
14. Julien Chiquet, Mahendra Mariadassou, and Stéphane Robin. Variational inference for probabilistic poisson pca. *The Annals of Applied Statistics*, 12(4):2674–2698, 2018.
15. Aaron Clauset, Mark EJ Newman, and Christopher Moore. Finding community structure in very large networks. *Physical review E*, 70(6):066111, 2004.
16. Roberto de la Cruz and Jan-Ulrich Kreft. Geometric mean extension for data sets with zeros. *arXiv preprint arXiv:1806.06403*, 2018.
17. Joshua V Dillon, Ian Langmore, Dustin Tran, Eugene Brevdo, Srinivas Vasudevan, Dave Moore, Brian Patton, Alex Alemi, Matt Hoffman, and Rif A Saurous. Tensorflow distributions. *arXiv preprint arXiv:1711.10604*, 2017.
18. Karoline Faust and Jeroen Raes. Conet app: inference of biological association networks using cytoscape. *F1000Research*, 5, 2016.
19. Cindy Xin Feng. A comparison of zero-inflated and hurdle models for modeling zero-inflated count data. *Journal of statistical distributions and applications*, 8(1):8, 2021.
20. Sofia K Forslund, Rima Chakaroun, Maria Zimmermann-Kogadeeva, Lajos Markó, Judith Aron-Wisnewsky, Trine Nielsen, Lucas Moitinho-Silva, Thomas SB Schmidt, Gwen Falony, Sara Vieira-Silva, et al. Combinatorial, additive and dose-dependent drug–microbiome associations. *Nature*, 600(7889):500–505, 2021.
21. Jonathan Friedman and Eric J Alm. Inferring correlation networks from genomic survey data. 2012.
22. AB Gelman, JB Carlin, HS Stern, DB Dunson, A Vehtari, and D Rubin. Bayesian data analysis third edition. Boca Raton. *FL: CRC Press.[Google Scholar]*, 2013.
23. Dirk Gevers, Subra Kugathasan, Lee a Denson, Yoshiaki Vázquez-Baeza, Will Van Treuren, Boyu Ren, Emma Schwager, Dan Knights, Se Jin Song, Moran Yassour, Xochitl C Morgan, Aleksandar D Kostic, Chengwei Luo, Antonio González, Daniel McDonald, Yael Haberman, Thomas Walters, Susan Baker, Joel Rosh, Michael Stephens, Melvin Heyman, James Markowitz, Robert Baldassano, Anne Griffiths, Francisco Sylvester, David Mack, Sandra Kim, Wallace Crandall, Jeffrey Hyams, Curtis Huttenhower, Rob Knight, and Ramnik J Xavier. The treatment-naive microbiome in new-onset Crohn’s disease. *Cell host & microbe*, 15(3):382–92, mar 2014.
24. Travis Gibson and Georg Gerber. Robust and scalable models of microbiome dynamics. In *International Conference on Machine Learning*, pages 1763–1772. PMLR, 2018.
25. Jack A. Gilbert, Janet K. Jansson, and Rob Knight. The earth microbiome project: successes and aspirations. *BMC Biology*, 12(1):69, 2014.
26. Samantha J Gleich, Jacob A Cram, Jake L Weissman, and David A Caron. Netgam: Using generalized additive models to improve the predictive power of ecological network analyses constructed using time-series data. *ISME Communications*, 2(1):1–9, 2022.
27. Alex Gobbi, Alberto Acedo, Nabeel Imam, Rui G Santini, Rüdiger Ortiz-Álvarez, Lea Ellegaard-Jensen, Ignacio Belda, and Lars H Hansen. A global microbiome survey of vineyard soils highlights the microbial dimension of viticultural terroirs. *Communications Biology*, 5(1):241, 2022.
28. Julia K. Goodrich, Jillian L. Waters, Angela C. Poole, Jessica L. Sutter, Omry Koren, Ran Blekhman, Michelle Beaumont, William Van Treuren, Rob Knight, Jordana T. Bell, Timothy D. Spector, Andrew G. Clark, and Ruth E. Ley. Human genetics shape the gut microbiome. *Cell*, 159(4):789–799, 2014.
29. Lionel Guidi, Samuel Chaffron, Lucie Bittner, Damien Eveillard, Abdelhalim Larhlimi, Simon Roux, Youssef Darzi, Stéphane Audic, Léo Berline, Jennifer Brum, Luis Pedro Coelho, Julio Cesar Ignacio Espinoza, Shruti Malviya, Shinichi Sunagawa, Céline Dimier, Stefanie Kandels-Lewis, Marc Picheral, Julie Poulain, Sarah Searson, Tara Oceans Coordinators, Lars Stemmann, Fabrice Not, Pascal Hingamp, Sabrina Speich, Mick Follows, Lee Karp-Boss, Emmanuel Boss, Hiroyuki Ogata, Stéphane Pesant, Jean Weissenbach, Patrick Wincker, Silvia G. Acinas, Peer Bork, Colomban de Vargas, Daniele Iudicone, Matthew B. Sullivan, Jeroen Raes, Eric Karsenti, Chris Bowler, and Gabriel Gorsky. Plankton networks driving carbon export in the oligotrophic ocean. *Nature*, 532(7600):in review, 2015.
30. Joshua G Harrison, W John Calder, Vivaswat Shastri, and C Alex Buerkle. Dirichlet-multinomial modelling outperforms alternatives for analysis of microbiome and other ecological count data. *Molecular ecology resources*, 20(2):481–497, 2020.
31. Ian Holmes, Keith Harris, and Christopher Quince. Dirichlet multinomial mixtures: generative models for microbial metagenomics. *PLoS one*, 7(2):e30126, 2012.

32. M. A. Ikram, G. G. O. Brusselle, S. D. Murad, C. M. van Duijn, O. H. Franco, A. Goedegebure, C. C. W. Klaver, T. E. C. Nijsten, R. P. Peeters, B. H. Stricker, H. Tiemeier, A. G. Uitterlinden, M. W. Vernooij, and A. Hofman. The Rotterdam Study: 2018 update on objectives, design and main results. *Eur J Epidemiol*, 32(9):807–850, 2017.
33. HMP Integrative. The integrative human microbiome project: dynamic analysis of microbiome-host omics profiles during periods of human health and disease. *Cell host & microbe*, 16(3):276–289, 2014.
34. Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
35. Diederik Pieter Kingma. Variational inference & deep learning: A new synthesis. 2017.
36. Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
37. Zachary D Kurtz, Christian L Müller, Emily R Miraldi, Dan R Littman, Martin J Blaser, and Richard A Bonneau. Sparse and compositionally robust inference of microbial ecological networks. *PLoS computational biology*, 11(5):e1004226, 2015.
38. Zachary D Kurtz, Christian L Müller, Emily R Miraldi, Dan R Littman, Martin J Blaser, and Richard A Bonneau. Sparse and compositionally robust inference of microbial ecological networks. 11(5):e1004226, 2015.
39. Seonjoo Lee, Pauline E Chugh, Haipeng Shen, R Eberle, and Dirk P Dittmer. Poisson factor models with applications to non-normalized microrna profiling. *Bioinformatics*, 29(9):1105–1111, 2013.
40. Yufang Li, Kai Tang, Lianbao Zhang, Zihao Zhao, Xiabing Xie, Chen-Tung Arthur Chen, Deli Wang, Nianzhi Jiao, and Yao Zhang. Coupled carbon, sulfur, and nitrogen cycles mediated by microorganisms in the water column of a shallow-water hydrothermal ecosystem. *Frontiers in microbiology*, 9:2718, 2018.
41. Gipsi Lima-Mendez, Karoline Faust, Nicolas Henry, Johan Decelle, Sébastien Colin, Fabrizio Carcillo, Samuel Chaffron, J Cesar Ignacio-Espinosa, Simon Roux, Flora Vincent, et al. Determinants of community structure in the global plankton interactome. *Science*, 348(6237):1262073, 2015.
42. Tiantian Liu, Peirong Xu, Yueyao Du, Hui Lu, Hongyu Zhao, and Tao Wang. Mzinbva: variational approximation for multilevel zero-inflated negative-binomial models for association analysis in microbiome surveys. *Briefings in Bioinformatics*, 23(1):bbab443, 2022.
43. Alan R Longhurst. *Ecological geography of the sea*. Elsevier, 2010.
44. Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology*, 15(12):1–21, 2014.
45. Daniel McDonald, Embriette Hyde, Justine W Debelius, James T Morton, Antonio Gonzalez, Gail Ackermann, Alexander A Aksenov, Bahar Behsaz, Caitriona Brennan, Yingfeng Chen, Lindsay DeRight Goldasich, Pieter C Dorrestein, Robert R Dunn, Ashkaan K Fahimipour, James Gaffney, Jack A Gilbert, Grant Gogul, Jessica L Green, Philip Hugenholtz, Greg Humphrey, Curtis Huttenhower, Matthew A Jackson, Stefan Janssen, Dilip V Jeste, Lingjing Jiang, Scott T Kelley, Dan Knights, Tomasz Kosciolok, Joshua Ladau, Jeff Leach, Clarisse Marotz, Dmitry Meleshko, Alexey V Melnik, Jessica L Metcalf, Hosein Mohimani, Emmanuel Montassier, Jose Navas-Molina, Tanya T Nguyen, Shyamal Peddada, Pavel Pevzner, Katherine S Pollard, Gholamali Rahnavard, Adam Robbins-Pianka, Naseer Sangwan, Joshua Shoreinstein, Larry Smarr, Se Jin Song, Timothy Spector, Austin D Swafford, Varykina G Thackray, Luke R Thompson, Anupriya Tripathi, Yoshiki Vázquez-Baeza, Alison Vrbnac, Paul Wischmeyer, Elaine Wolfe, Qiyun Zhu, and Rob Knight. American Gut: an Open Platform for Citizen Science Microbiome Research. *mSystems*, 3(3), 2018.
46. Paul J. McMurdie and Susan Holmes. phyloseq: An r package for reproducible interactive analysis and graphics of microbiome census data. *PLOS ONE*, 8(4):1–11, 2013.
47. Jacquelyn S Meisel, Geoffrey D Hannigan, Amanda S Tyldsley, Adam J SanMiguel, Brendan P Hodkinson, Qi Zheng, and Elizabeth A Grice. Skin microbiome surveys are strongly influenced by experimental design. *Journal of Investigative Dermatology*, 136(5):947–956, 2016.
48. Aditya Mishra and Christian L. Müller. Robust regression with compositional covariates. *Computational Statistics and Data Analysis*, 165:107315, 2022.
49. Aditya K Mishra and Christian L Müller. Negative binomial factor regression with application to microbiome data analysis. *Statistics in Medicine*, 41(15):2786–2803, 2022.
50. Mary Ann Moran. The global ocean microbiome. *Science*, 350(6266):aac8455, 2015.
51. Molly A Moynihan, Nathalie F Goodkin, Kyle M Morgan, Phyllis YY Kho, Adriana Lopes dos Santos, Federico M Lauro, David M Baker, and Patrick Martin. Coral-associated nitrogen fixation rates and diazotrophic diversity on a nutrient-replete equatorial reef. *The ISME journal*, 16(1):233–246, 2022.
52. GJ Olsen, DJ Lane, SJ Giovannoni, NR Pace, and DA Stahl. Microbial ecology and evolution: a ribosomal rna approach. *Annual review of microbiology*, 40:337–365, 1986.
53. Johannes Ostner, Salomé Carcy, and Christian L Müller. tascoda: Bayesian tree-aggregated analysis of compositional amplicon and single-cell data. *Frontiers in genetics*, 12:766405, 2021.
54. N.R Pace, D.A. Stahl, D.J. Lane, and G.J. Olsen. The analysis of natural microbial populations by ribosomal rna sequences. In Marshall K.C., editor, *Advances in Microbial Ecology*, volume 9, pages 1–55. Springer, Boston, MA, 1986.
55. Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
56. Stéphane Pesant, Fabrice Not, Marc Picheral, Stefanie Kandels-Lewis, Noan Le Bescot, Gabriel Gorsky, Daniele Iudicone, Eric Karsenti, Sabrina Speich, Romain Troublé, et al. Open science resources for the discovery and analysis of tara oceans data. *Scientific data*, 2(1):1–16, 2015.
57. Stefanie Peschel, Christian L Müller, Erika von Mutius, Anne-Laure Boulesteix, and Martin Depner. NetCoMi: network construction and comparison for microbiome data in R. *Briefings in Bioinformatics*, 2020.

58. Christian Quast, Elmar Pruesse, Pelin Yilmaz, Jan 069 100670. Gene W Tyson, Jarrod Chapman, Philip Hugenholtz, Gerken, Timmy Schweer, Pablo Yarza, Jörg Peplies, and 070 1007 Eric E Allen, Rachna J Ram, Paul M Richardson, Victor V Frank Oliver Glöckner. The silva ribosomal rna gene 071 1008 Solovyev, Edward M Rubin, Daniel S Rokhsar, and Jillian F database project: improved data processing and web-based 072 1009 tools. *Nucleic acids research*, 41(D1):D590–D596, 2012. 1073
59. Janaina Rigonato, Marko Budinich, Alejandro A 1010 Murillo, Manoela C Brandão, Juan J Pierella Karlusich, 1011 W Duncan Wadsworth, Raffaele Argiento, Michele 1012 Guindani, Jessica Galloway-Pena, Samuel A Shelburne, 1013 Yawouvi Dodji Soviadan, Ann C Gregory, Hisashi Endo, 1014 Florian Kokoszka, Dean Vik, et al. Ocean-wide comparisons 1077 of mesopelagic planktonic community structures. *ISME 1078 communications*, 3(1):83, 2023. 1079
60. Donald B Rubin. Bayesianly justifiable and relevant 1015 frequency calculations for the applied statistician. *The 1016 Annals of Statistics*, pages 1151–1172, 1984. 1080
61. Francisco JR Ruiz, Susan Athey, and David M Blei. 1017 Shopper: A probabilistic model of consumer choice 1076 with substitutes and complements. *arXiv preprint 1077 arXiv:1711.03560*, 2017. 1078
62. Kris Sankaran and Susan P Holmes. Latent variable 1018 modeling for the microbiome. *Biostatistics*, 20(4):599–614, 1082 2019. 1083
63. Patrick D. Schloss, Sarah L. Westcott, Thomas Ryabin, 1019 Justine R. Hall, Martin Hartmann, Emily B. Hollister, 1020 Ryan A. Lesniewski, Brian B. Oakley, Donovan H. 1021 Parks, Courtney J. Robinson, Jason W. Sahl, Blaz 1022 Stres, Gerhard G. Thallinger, David J. Van Horn, and 1023 Carolyn F. Weber. Introducing mothur: Open-source, 1085 platform-independent, community-supported software for 1086 describing and comparing microbial communities. *Applied 1087 and Environmental Microbiology*, 75(23):7537–7541, 2009. 1088
64. Salome Scholtens, Nynke Smidt, Morris A Swertz, 1024 Stephan JL Bakker, Aafje Dotinga, Judith M Vonk, Freerk 1025 van Dijk, Sander KR van Zon, Cisca Wijmenga, Bruce HR 1026 Wolffenbuttel, and Ronald P Stolk. Cohort Profile: 1027 LifeLines, a three-generation cohort study and biobank. 1028 *International Journal of Epidemiology*, 44(4):1172–1180, 1029 2015. 1090
65. Justin P Shaffer, Louis-Félix Nothias, Luke R Thompson, 1030 Jon G Sanders, Rodolfo A Salido, Sneha P Couvillion, 1031 Asker D Brejnrod, Franck Lejzerowicz, Niina Haiminen, 1032 Shi Huang, et al. Standardized multi-omics of earth's 1033 microbiomes reveals microbial and metabolite diversity. 1034 *Nature microbiology*, 7(12):2128–2150, 2022. 1091
66. Alaina N Smith, Gwenn MM Hennon, Erik R Zinser, 1035 Benjamin C Calfee, Jeremy W Chandler, and Andrew D 1036 Barton. Comparing prochlorococcus temperature niches 1037 in the lab and across ocean basins. *Limnology and 1038 Oceanography*, 66(7):2632–2647, 2021. 1092
67. Shinichi Sunagawa, Luis Pedro Coelho, Samuel Chaffron, 1039 Jens Roat Kultima, Karine Labadie, Guillem Salazar, 1040 Bardya Djahanschiri, Georg Zeller, Daniel R Mende, 1041 Adriana Alberti, et al. Structure and function of the global 1042 ocean microbiome. *Science*, 348(6237):1261359, 2015. 1093
68. Shinichi Sunagawa, Daniel R Mende, Georg Zeller, 1043 Fernando Izquierdo-Carrasco, Simon A Berger, Jens Roat 1044 Kultima, Luis Pedro Coelho, Manimozhayan Arumugam, 1045 Julien Tap, Henrik Bjørn Nielsen, et al. Metagenomic 1046 species profiling using universal phylogenetic marker genes. 1047 *Nature methods*, 10(12):1196–1199, 2013. 1094
69. Peter J. Turnbaugh, Ruth E. Ley, Micah Hamady, Claire M. 1048 Fraser-Liggett, Rob Knight, and Jeffrey I. Gordon. The 1049 human microbiome project. *Nature*, 449(7164):804–810, 1050 2007. 1095
70. Gene W Tyson, Jarrod Chapman, Philip Hugenholtz, 1051 Eric E Allen, Rachna J Ram, Paul M Richardson, Victor V 1052 Solovyev, Edward M Rubin, Daniel S Rokhsar, and Jillian F 1053 Banfield. Community structure and metabolism through 1054 reconstruction of microbial genomes from the environment. 1055 *Nature*, 428(6978):37–43, 2004. 1074
71. W Duncan Wadsworth, Raffaele Argiento, Michele 1056 Guindani, Jessica Galloway-Pena, Samuel A Shelburne, 1057 and Marina Vannucci. An integrative bayesian dirichlet- 1077 multinomial regression model for the analysis of taxonomic 1078 abundances in microbiome data. *BMC bioinformatics*, 18(1):1–12, 2017. 1080
72. Martin J Wainwright, Michael I Jordan, et al. Graphical 1059 models, exponential families, and variational inference. 1060 *Foundations and Trends® in Machine Learning*, 1(1– 1061 2):1–305, 2008. 1082
73. Nyree J West, Cécile Lepère, Carmem-Lara de O Manes, 1062 Philippe Catala, David J Scanlan, and Philippe Lebaron. 1063 Distinct spatial patterns of sar11, sar86, and actinobacteria 1064 diversity along a transect in the ultra-oligotrophic south 1065 pacific ocean. *Frontiers in microbiology*, 7:234, 2016. 1085
74. Fox G Woese C. Phylogenetic structure of the prokaryotic 1066 domain. *Pnas*, 74(11):5088–5090, 1977. 1086
75. Tianchen Xu, Ryan T Demmer, and Gen Li. Zero-inflated 1067 poisson factor model with application to microbiome read 1068 counts. *Biometrics*, 77(1):91–101, 2021. 1087
76. Pelin Yilmaz, Pablo Yarza, Josephine Z Rapp, and Frank O 1069 Glöckner. Expanding the world of marine bacterial and 1070 archaeal clades. *Frontiers in microbiology*, 6:1524, 2016. 1091
77. Grace Yoon, Irina Gaynanova, and Christian L Müller. 1071 Microbial networks in spring-semi-parametric rank- 1072 based correlation and partial correlation estimation for 1073 quantitative microbiome data. *Frontiers in genetics*, 1074 10:516, 2019. 1088
78. Yanyan Zeng, Hongyu Zhao, and Tao Wang. Model-based 1075 microbiome data ordination: A variational approximation 1076 approach. *Journal of Computational and Graphical 1077 Statistics*, pages 1–13, 2021. 1092
79. Xinyan Zhang, Himel Mallick, Zaixiang Tang, Lei Zhang, 1078 Xiangqin Cui, Andrew K Benson, and Nengjun Yi. Negative 1079 binomial mixed models for analyzing microbiome count 1080 data. *BMC bioinformatics*, 18(1):4, 2017. 1093
80. Qiang Zheng, Yu Wang, Rui Xie, Andrew S Lang, Yanting 1081 Liu, Jiayao Lu, Xiaodong Zhang, Jun Sun, Curtis A Suttle, 1082 and Nianzhi Jiao. Dynamics of heterotrophic bacterial 1083 assemblages within synechococcus cultures. *Applied and 1084 environmental microbiology*, 84(3):e01517–17, 2018. 1094