

https://doi.org/10.1093/ismeco/ycaf062 Advance access publication: 2 May 2025 Original Article

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

Aditya Mishra 🕞<sup>1,\*</sup>, Jesse McNichol 🕞<sup>2,3</sup>, Jed Fuhrman 🕞<sup>3</sup>, David Blei 🕞<sup>4,5</sup>, Christian L. Müller<sup>4,6,7</sup>

<sup>1</sup>Department of Statistics, University of Georgia, Athens, GA, 30606, United States

<sup>2</sup>Department of Biology, St. Francis Xavier University, Antigonish, NS, B2G 2W5, Canada

<sup>3</sup>Department of Biological Sciences, University of Southern California, LA, 90007, United States

<sup>4</sup>Center for Computational Mathematics, Flatiron Institute, New York, NY, 10010, United States

<sup>5</sup>Department of Statistics and Computer Science, Columbia University, New York, NY, 10027, United States

<sup>6</sup>Computational Health Center, Helmholtz Zentrum München, Munich, 85764, Germany

<sup>7</sup>Department of Statistics, LMU München, Munich, 80539, Germany

\*Corresponding author. Department of Statistics, University of Georgia, Athens, GA, 30606, USA. E-mail: 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 taxon 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' 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 alteromondales-dominated communities, each 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.

Keywords: microbiome; probabilistic model; association learning; variational inference; Tara ocean expedition

#### Introduction

Microbial species are an integral part of life on earth. Ecosystems, ranging from the human gut to the global ocean, harbor trillions of bacteria, archaea, viruses, and fungi that take on essential functional roles and have developed intricate ecological relationships within their respective habitat. Over the past decades, advances in amplicon and metagenomics sequencing techniques [1-4] and standardized experimental and bioinformatics workflows [5-7] have enabled the large-scale collection and dissemination of microbial survey data, including those from the seminal Human Microbiome Project [8], several gut-focused surveys [9-12], the Earth Microbiome Project [13], and the Tara Ocean Expedition [14]. These surveys have reached a level of maturity and complexity that ultimately allow the estimation of statistical associations between microbial abundances, typically represented as compositional counts of amplicon sequence variants (ASVs) or operational taxonomic units (OTUs), and habitat properties [14, 15], biogeochemical processes[16], and/or host health status [17, 18]. This, in turn, provides a starting point for deciphering and understanding the ecological and functional roles of different microbial clades in the ecosystem, nutrient and bio(geo)chemical dependencies, resource limitations of microbial growth, and the presence of ecological taxon–taxon interactions [19].

Conventional statistical methods for microbiome data largely focus on either statistical abundance modeling [20–25] or taxontaxon associations [19, 26–29]. The complexity and interconnected nature of microbial ecosystems would, however, benefit from integrative approaches that go beyond estimating individual statistical associations or taxon-taxon associations. To fully capture the interplay between microbial abundances, host or habitat characteristics, and taxon-taxon interactions, it is essential to estimate their contributions within a unified model.

Here, we introduce such an integrative probabilistic modeling framework that is specifically tailored to microbiome survey data and enables joint estimation of habitat-dependent drivers and broad associations patterns of microbial taxa abundances (see Fig. 1). Our approach, termed VI-MIDAS (variational inference for microbiome survey data analysis), models the observed taxon abundances by simultaneously learning taxon-specific latent representations that leverage the effects of host or environmental factors and taxon-taxon associations via an item-item

Received: 18 March 2024. Revised: 21 January 2025. Accepted: 8 April 2025

© The Author(s) 2025. Published by Oxford University Press on behalf of the International Society for Microbial Ecology.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.



**Figure 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.

interaction modeling approach, originally proposed for market basket analysis [30].

VI-MIDAS uses the parametric structure of the negative binomial distribution [25, 31] to account for the overdispersed nature of the amplicon count data and comprises two main model components: (i) a component that allows for full adjustment of taxon abundances from a user-defined subset of covariates and (ii) taxa-specific latent vectors that incorporate, e.g. spatio-temporal or environmental covariates and taxon-taxon interactions, thus providing a marginal characterization of each taxon. We resort to mean-field variational inference for parameter estimation of VI-MIDAS' intractable posterior distribution [32], thus complementing other recent variational approaches to microbiome data modeling, such as, e.g. Poisson principal component analysis [33], microbiome dynamics modeling [34], Dirichlet Multinomial modeling [35], multi-level modeling [36], and microbiome ordination [37].

To illustrate the complete workflow of the VI-MIDAS framework, we focus on integrative analysis of global marine microbiome survey data. The ocean microbiome is of fundamental importance for life on earth, being responsible for about half of all primary production (i.e. the production of chemical energy in organic compounds) and holds enormous potential for climate remediation [38]. Several initiatives such as the Tara Oceans Project [39] and the Simons CMAP [40] provide well-structured sequencing data, biogeochemical and environmental covariate data, and satellite-derived products that are amenable to statistical analysis. Here, we re-analyze Tara expedition data (http:// ocean-microbiome.embl.de/companion.html), originally considered in [14] to study the structure and function of the global ocean microbiome. The expedition collected ocean water samples from 68 distinct geographical locations at varying levels of depth. We will make extensive use of this dataset to motivate and describe

the details of the VI-MIDAS framework as well as the learned representations and associations of the global ocean microbiome.

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

#### **Materials and methods**

# Tara ocean data and ecologically relevant taxa re-classification

We consider the processed Tara expedition data, as provided at http://ocean-microbiome.embl.de/companion.html. The expedition collected water samples from 68 distinct geographical locations (Fig. 2B) across different depths, resulting in n = 139 distinct samples. Across these samples, the original data comprises microbial taxa abundances profiles of more than 35 000 bacterial



**Figure 2.** Illustration of the Tara ocean data: (A) Taxon abundance profiles, agglomerated to expert-derived ERCs for two samples (marked as 1 and 2 in Fig. 2B). (B) Tara ocean sample locations. (C) Environmental features associated with the samples marked as 1 and 2 in Fig. 2B; (D) Abundance profiles log(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.

taxa in form of metagenomic OTUs (mOTUs) (derived using the miTAGS framework [41]).

Here, we focus on the most abundant taxa by taking the union of all mOTUs that, in each individual sample, contribute to 40% of the total library size. This filtering allows us to cover the abundance profiles of the q = 1378 taxa with the most significant variability and reduces the number of excess zero counts. To account for the highly variable sequencing depth across the samples, we normalize the abundance data with respect to the lowest library size via common-sum scaling [31]. Fig. 2D shows the log-transformed abundance profiles  $\mathbf{W} \in \mathbb{R}^{n \times q}$ . Since the original taxonomic affiliations of the miTAGS are difficult to interpret, we next developed a partitioning of the selected taxa into ecologically relevant classes (ERCs). The original full taxonomy strings are too long to understand at a glance, and parsing by taxonomic level is not a good option since taxa vary widely in the depth of their annotations. For example, cyanobacteria should be annotated at the genus level or higher, but many other abundant but less described taxa do not have any taxonomic information at that level. We manually curated the data to provide a short relevant taxonomic indicator that provides a rough indicator of the ecological niche of an organism while remaining short enough to be interpreted at a glance [42]. Some taxonomies have been altered to preserve the updated SILVA taxonomy (i.e. Betaproteobacteria

is now Burkholderiales). New SILVA 138 [43] taxonomies have been used wherever possible (i.e. when the original ID was still in SILVA 138), but in cases where there was only the SILVA 108 taxonomic information, we have used our best guess. For example, if an organism had the same classification as other organisms in SILVA 108, we have often given it the same name as its counterparts in SILVA 138. We present all our findings in terms of these 29 ERCs.

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

### Generative modeling in VI-MIDAS

We seek to model the abundance profiles of q microbial taxa where we denote a single sample by the random variable  $\mathbf{w} \in \mathbb{R}^{q}$ and the observed data from n samples by  $\mathbf{W} = [w_{ij}]_{n \times q} \in \mathbb{R}^{n \times q}$ . For concreteness, we illustrate model building and analysis using the Tara abundance profiles (see Fig. 2D) of q = 1378 taxa but the

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

modeling strategy is applicable to any multimodal microbiome survey.

#### Distributional model

VI-MIDAS posits that the overdispersed microbial count data **W** are reasonably well modeled with the negative binomial distribution [18, 44, 45]. While other generative statistical modeling approaches are available, including the Dirichlet Multinomial (mixture) framework [20, 46], latent Dirichlet allocation [47], and Poisson distribution models [21, 24, 48], we found the negative binomial model to be an excellent choice for the Tara ocean data (see Fig. S1 B of the Supplementary material for the overdispersion analysis). Using the negative binomial distribution with mean and dispersion parameterization [44], VI-MIDAS models the jth taxa in the ith sample as:

$$p(\omega_{ij};\tau_j\mu_{ij},\phi_j) = \text{NB}(\omega_{ij};\tau_j\mu_{ij},\phi_j)$$
$$= \binom{\omega_{ij}+\phi_j-1}{\omega_{ij}} \left(\frac{\tau_j\mu_{ij}}{\tau_j\mu_{ij}+\phi_j}\right)^{\omega_{ij}} \left(\frac{\phi_j}{\tau_j\mu_{ij}+\phi_j}\right)^{\phi_j}.$$
(1)

Here, the mean parameter  $\tau_j \mu_{ij}$  is the product of a taxon-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  $\Phi = [\phi_1, \ldots, \phi_q]$  and  $\tau = [\tau_1, \ldots, \tau_q]$ , respectively. The shape parameter  $\tau$  accounts for the disparity in abundance among microbial taxa. The generative model (1) of VI-MIDAS implies  $\mathbb{E}(\omega_{ij}) = \tau_j \mu_{ij}$  and  $\operatorname{Var}(\omega_{ij}) = \tau_j \mu_{ij} + \frac{\tau_j^2 \mu_{ij}^2}{\phi_{ij}^2}$ . This variance structure reflects the overdispersion characteristic of the negative binomial distribution , making this framework well-suited for modeling overdispersed count data.

#### Modeling strategy and model components

One novelty in VI-MIDAS is the combination of ideas from generalized linear modeling [44] and compositional data analysis [49] to associate the microbial relative count data with spatiotemporal, environmental, and taxa information. Specifically, we model the log-transformed mean parameter  $\boldsymbol{\mu} = [\mu_{ij}]_{n \times q}$  of the generative model (1) with two components, a consistent zero-aware geometric mean estimate  $t_i$  and a linear predictor  $\boldsymbol{\eta} = [\eta_{ij}]_{n \times q} \in \mathbb{R}^{n \times q}$ 

as follows:

$$\log \mu_{ij} = \log t_i + \eta_{ij} \,. \tag{2}$$

The sample-wise parameter t<sub>i</sub> is estimated by a zero-aware geometric mean estimator, introduced in [50], 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<sub>i</sub> 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  $\eta$  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  $\eta$  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  $\eta_{ij}^{[P]}$ ,  $\eta_{ij}^{[D]}$ , and  $\eta_{ij}^{[S]}$ , respectively. Lastly, statistical associations among co-occurring taxa are included via a latent interaction term  $\eta_{ij}^{[I]}$ , leading to following model:

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

 $\eta_{ii}^{[E]} = \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\gamma}_{.j} \,.$ 

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 [14], 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 hyper-parameter 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 [51]. We denote the model matrix indicating the *r* distinct regions of the *n* samples by  $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_n]^T \in \{0, 1\}^{n \times r}$ . Here,  $\mathbf{r}_i$  is a one-hot encoded vector, indicating membership of the *i*-th sample to one of the *r* provinces. The matrix  $\mathbf{R}$  connects to the joint latent space via the coefficient matrix  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}]_{r \times k} \in \mathbb{R}^{r \times k}$ , leading to

$$\boldsymbol{\eta}_{ij}^{[\mathrm{P}]} = \mathbf{r}_i \boldsymbol{\alpha} \boldsymbol{\beta}_{.j} \,. \tag{5}$$

Similarly, the Tara data include samples across d = 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  $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_n]^T \in \{0, 1\}^{n \times d}$  ( $\mathbf{d}_i$  is one-hot encoded vectors representing membership of the *i*-th sample in one of d ocean depth) and connect it to the joint latent space via the coefficient matrix  $\boldsymbol{\delta} = [\boldsymbol{\delta}]_{d \times k} \in \mathbb{R}^{d \times k}$ , leading to

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

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 s = 4 seasons: the 1<sup>st</sup> (Q1, January to March), 2<sup>nd</sup> (Q2, April to June), 3<sup>rd</sup> (Q3, July to September), and 4<sup>th</sup> (Q4, October to December) yearly quarter, and construct the season indicator matrix  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^T \in \{0, 1\}^{n \times s}$ , where each row  $\mathbf{s}_i$  represents the membership of the ith sample to one of the s seasonal categories. The coefficient matrix  $\boldsymbol{\vartheta} = [\vartheta]_{s \times k} \in \mathbb{R}^{s \times k}$  couples  $\mathbf{S}$  to the latent space  $\boldsymbol{\beta}$ , leading to

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

In summary, the coupling of the described features to a shared latent space via the coefficient matrices  $\alpha$ ,  $\delta$ ,  $\vartheta$  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 [52]. While discovering detailed ecological interactions among taxa, such as, e.g. competition, mutualism, parasitism, 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 [30]. In SHOPPER, Ruiz et al. [30] 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 itemitem interaction component. Following their setup, the "interaction," or, in the biological context, association of the jth taxa with any *m*th taxa is given by  $\boldsymbol{\rho}_i^{\mathrm{T}} \boldsymbol{\beta}_m$  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  $\eta^{[1]}$  for the jth taxon in the ith sample are thus given by

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

where  $a_i = \sum_{m=1}^{q} \mathbf{1}_{\omega_{im}\neq 0}$  is the total number of taxa present in the ith sample. Note that the interaction term  $\rho^T \boldsymbol{\beta}$  is not symmetric. Note that, while some ecological interactions, such as parasitism, are directed and asymmetric, thus making asymmetry biologically plausible, the model parameter  $\rho$  and  $\beta$  do not allow to estimate *directionality*. Consequently, following the SHOPPER model [30], we derive a symmetrized interaction matrix  $\mathbf{I} = [I_{i,j}] \in \mathbb{R}^{q \times q}$  with each entry being computed as:

$$I_{i,j} = \left(\boldsymbol{\rho}_{\cdot i}^{\mathrm{T}} \boldsymbol{\beta}_{\cdot j} + \boldsymbol{\rho}_{\cdot j}^{\mathrm{T}} \boldsymbol{\beta}_{\cdot i}\right) / 2 \tag{9}$$

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  $\ell = \{\alpha, \vartheta, \beta, \gamma, \rho, \tau, \Phi\}$  (see Table S1 of the Supplementary material). Given the microbial abundance data **W**, the (direct) covariates **X**, and the model parameters  $\ell$ , we integrate the generative model (1) into a Bayesian framework where the posterior distribution reads:

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

where  $p(\mathbf{W}; \boldsymbol{\ell}, \mathbf{X}, \mathbf{t}) = \prod_{i,j} p(w_{ij}; \tau_j \mu_{ij}, \phi_j)$  denotes the likelihood of  $\mathbf{W}$  and  $p(\boldsymbol{\ell}) = p(\boldsymbol{\alpha})p(\boldsymbol{\delta})p(\boldsymbol{\beta})p(\boldsymbol{\gamma})p(\boldsymbol{\phi})p(\boldsymbol{\Phi})p(\boldsymbol{\tau})p(\boldsymbol{\vartheta})$  the prior distribution,

respectively. To achieve good generalizabilty and interpretability of VI-MIDAS' over-parameterized model, we place sparsityinducing 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.  $p(\phi_j) = \text{inverse-Cauchy}(0, \upsilon)$  and  $p(\Phi) = \prod_j p(\phi_j)$ , and a Uniform (1,2) prior for the shape parameter  $\tau$ , i.e.  $\tau_j \sim$ 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) [32, 53, 54] and approximate the posterior with a variational posterior distribution of the latent variable  $\boldsymbol{\ell}$ . Briefly, let  $q(\boldsymbol{\ell}; \boldsymbol{\nu})$  be the variational posterior distribution with parameter  $\boldsymbol{\nu}$ . VI approximates sampling of the posterior by minimizing the Kullback-Leibler (KL) divergence,

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

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

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

which is a lower bound on the logarithm of the joint probability of the observations  $log P(\mathbf{W}; \mathbf{X}, \mathbf{t})$  [53]. Replacing the joint distribution  $P(\mathbf{W}, \boldsymbol{\ell}; \mathbf{X}, \mathbf{t})$  with a product of likelihood and prior distribution  $P(\mathbf{W}, \boldsymbol{\ell}; \mathbf{X}, \mathbf{t}) = P(\mathbf{W}; \boldsymbol{\ell}, \mathbf{X}, \mathbf{t})P(\boldsymbol{\ell})$  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 [55]. 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 [56]. 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 [57] 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 [56, 58] on the full data. The procedure requires generating *m* posterior samples, denoted by the random variables  $\mathbf{W}^{rep} = [\mathbf{w}_{ij}^{rep}] \in \mathbb{R}^{n \times q}_+$ , 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 Fig. S7 in 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 **W**. Figure 3A and 3B shows the observed and estimated abundance profiles (averaged over m = 100 samples), respectively. Fig. 3C shows the count histograms of data and model (pooled across all samples and species), and Fig. 3D shows 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 [14], the out-of-sample LLPD is reduced the most when the ocean depth (D) component is ablated (LLPD = -3.3882). This reflects the well-known depth stratification of marine species between the sunlit ocean and aphotic deep ocean ecosystems. Figure S4 in Supplementary material illustrates the learned depth stratification across all taxa, as reflected in the component  $\delta\beta$ . The environmental component was identified as the second most important component with an LLPD reduction of -3.3554.

Figure 4 summarizes the estimated effects  $\delta\beta$  of the ocean depth features and the environmental effects  $\gamma$  on the abundance of species aggregated into ERCs, respectively. The ocean depth summary (Fig. 4A) reveals three distinct sets of occurrence



**Figure 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.

patterns for two different groups of ERCs. One group (right most in Fig. 4A) comprises ERCs such as Nitrosopumilius, Pseudomonadales, SAR 324 clade, and Sphingomonadales which thrive in the Mesopelagic (MES) zone. A second group includes species like Prochlorococcus, SAR 116 clade, and Synechococcus, which flourish within the ecosystem of the ocean's deep chlorophyll maximum (DCM) and surface mixed layer (SRF) zones. The third group comprises marine Actinobacteria, Verrucomicrobiota, and others that show no dependence on depth. A summary of geochemical features highlights temperature (the top row in Fig. 4B) as the primary positive factor influencing the abundance of Synechococcus, Prochlorococcus, and Puniceispirillales (SAR116 clade). Oxygen concentration emerges as the main positive driver of abundance for Cytophagales, Flavobacteriales, and Roseobacter clades, while Nitrates, Nitrites, and Phosphate are identified as key drivers for the SAR324 clades, Nitrosopumilus, and Oceanospirillales (four right most columns in Fig. 4B). The estimated patterns broadly recapitulate known biology about ocean microbial ecosystems.

#### VI-MIDAS reveals five latent microbial sub-communities

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

To quantify similarity between microbial taxa in the latent space, we first computed cosine distances of all pairs of the q latent vectors. This particular choice of distance allows us to bypass the non-identifiability issue of the parameter  $\beta$ . We used the resulting distance matrix to construct a k-nearest neighbors graph ( $k_{nn} = 10$ ). Figure 5 shows the latent space embedding using a force-directed layout of the k-nn graph. We next performed Clauset-Newman-Moore greedy modularity analysis of the nearest neighbor graph [59] and identified five distinct modules in the latent space (see M1-M5 in Fig. 5 with top five ERCs highlighted and color-coded). The latent space representation reveals several distinct microbial sub-communities, dominated by a few ERCs, including one sub-community dominated by Prochlorococcus and SAR11 clades and one dominated by Nitrosopumilus. Module 1 (M1) comprises Flaviobacteriales, SAR86 clades, and the Chloroplast class. SAR11 clade, SAR86 clade, and Flavobacteriales are heterotrophs with functional similarity in oxidizing carbon in the ocean [60]. Both SAR86 clade and SAR11 clade follow a similar seasonal pattern (in the Bermuda Atlantic Time Series oceanographic stations) and coexist in oligotrophic regions with less nutrient supply [61]. Module 2 (M2) includes Nitrosopumilus, Marinimicrobia, and SAR324 clades. Existing literature supports that SAR11 clade (a subgroup of a species), Marinimicrobia, and MGII Archaea are more abundant in deep sea water [62]. Module 3 (M3) comprises Prochlorococcus, SAR11, Marine Actinobacteria, and SAR86 clades, among others, all comprising dominant taxa of the sunlit ocean. The two smallest modules 4 and 5 (M4 and M5, respectively) are dominated by Alteromonadales and are separating M2 from M1 and M3. Interestingly, Module 4 also



Figure 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 ERCs.

comprises Synechococcus species. This module thus hints at the known metabolic dependency of certain Alteromonadales taxa on Synechococcus (a photoautotroph) [63]. Although the latent representation does separate the majority of ERCs into distinct subgroups, we nonetheless observe that taxa of certain ERCs are spread out over the latent space, indicating different niche specialization. For instance, the SAR11 clade, one of the most abundant marine microbial taxa, is present in three different modules. Likewise, taxa in the SAR86 clade are present in both modules M1 and M3. For ease of identification, Table S3 in Supplementary material 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 **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, PO<sub>4</sub>, and NO<sub>2</sub>NO<sub>3</sub> 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 (200–1000 m) of the ocean (see Fig. 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,  $PO_4$ , and  $NO_2NO_3$ . In the oxygen-depleted environment, Nitrosopumilus survives by oxidizing ammonia to nitrite, confirming the observed association pattern [64]. Marinomicrobia (SAR406 clade) in groups M1 and M2 allow us to distinguish subgroups of species that can survive in both deep and shallow water [62].

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 (and negatively associated with MES. Among the geochemical factors, temperature, salinity, and oxygen concentration are positively associated, whereas the concentration of nitrates,  $PO_4$ , and  $NO_2NO_3$  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.

#### Positive and Negative interactions among ERCs

VI-MIDAS includes a mechanism for learning microbial interactions adjusted for direct (here, environmental) covariates. Contrary to prominent (partial) correlation-based methods [65, 66], VI-MIDAS follows the SHOPPER utility model [30] and quantifies pairwise interactions  $I_{ij}$  between any two taxa *i* and *j* in terms of the latent variables  $\rho$  and  $\beta$  (see Eq. 9).

To get a high-level view of the estimated interactions, we aggregated the adjacency matrices of significant positive and negative interactions among taxa by ERCs (for a more detailed view of the most significant taxon-level interactions, we refer to Section 4 of the Supplementary Materials). Figure 7 illustrates the aggregated positive (lower triangle) and negative (upper triangle) interactions among ERCs. The diagonal entry highlights the maximum of the two types of interactions to avoid confusion (see



Figure 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 in each of the five modules (see main text for further information).

also Section 4 of the Supplementary materials for the matrix of ratios between positive and negative interactions). We observe that SAR11 clade and Rhodospirillales form positive interactions with almost all other ERCs. SAR11 clade and Rhodospirillales belong to the Alphaproteobacteria phylum that play a critical role in carbon and nitrogen fixation [67, 68], potentially explaining the large number of interactions. However, members of the SAR11 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 [14, 41, 69–72]. 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, 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 apply mean-field variational inference framework to learn an approximate posterior distribution of the parameters.

VI-MIDAS is available in Python and uses the probabilistic programming language Stan [57]. The implementation is available on GitHub (https://github.com/amishra-stats/vi-midas). The repository also includes Python scripts and Jupyter notebooks for VI-MIDAS' three-stage parameter estimation framework: hyperparameter tuning, component contribution analysis, and sensitivity analysis.

To illustrate the VI-MIDAS modeling and analysis workflow, we have used data from the global Tara expedition [14], connecting the available spatiotemporal and environmental characteristics with generative modeling of the amplicon count data. To ease interpretability, we also grouped the amplicon-derived taxa into expert-annotated ecologically relevant classes (ERCs) which may be of independent interest for the analysis of other marine sequencing data. Focusing on the q = 1378 most abundant taxa representing 23 ERCS, we integrated the geochemical data using the direct coupling mechanism, effectively removing influence of common environmental factors such as temperature, salinity, and elemental compositions on microbial abundances. The remaining spatiotemporal features, 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



**Figure 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 and ocean depth (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.

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 (> 1200) can be categorized into three global microbial communities (M1–M3 in Fig. 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. [73] and references therein). Furthermore, our analysis suggests two distinct Alteromonadales-dominated communities that show different depth and province dependencies (M4–M5) (see Fig. 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 [74].

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



Figure 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.

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 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 [75] or the Dirichlet-Multinomial model [35, 76]. Finally, in its current state, VI-MIDAS is built on Stan [57] 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 [77] or PyTorch [78], may enable efficient porting of VI-MIDAS into these general-purpose ecosystems. Paired with variational inference tools [79], would potentially allow for faster model adaptation and alternative optimization routines.

VI-MIDAS makes the explicit methodological choice to symmetrize the latent interaction structure, leading to the

interpretation of positive ("mutualistic") and negative ("competitive") associations (see Eq. 9). While this approach precludes the representation of asymmetric (directed) interactions such as parasitism, it enables meaningful aggregation of the overall interaction structure across ERCs. In our future work, we aim to explore explicit modeling of directionality in microbial interactions, potentially leveraging approaches designed for directed networks.

The VI-MIDAS framework was designed to adapt model complexity to dataset size, ensuring applicability to presently available microbiome survey data. VI-MIDAS employed a random search strategy for hyperparameter tuning, sampling 50 parameter combinations for latent dimension (k), sparsity-inducing priors ( $\lambda$ ), and dispersion parameters (v) (see Section 3.3 of the Supplementary material). Five-fold cross-validation based on out-of-sample (LLPD) identified optimal hyperparameter settings, with robustness further ensured through 50 random initializations to address non-convexity in ELBO optimization (see Section 3.5 of the Supplementary material). We also performed model ablation to evaluate the importance and contribution of specific model components by systematically removing them and assessing their impact on performance or interpretability. This three-pronged strategy—hyperparameter

tuning, sensitivity analysis, and model ablation—effectively balances model complexity with sample size. In our experience, VI-MIDAS performs robustly with datasets of several hundred samples. For smaller datasets, we recommend to simplify the model by removing the latent taxon-taxon interaction terms to mitigate overfitting.

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 largescale survey data, including gut microbiome surveys [12, 80, 81], thus representing a potentially valuable general-purpose tool for the integrated analysis of modern microbiome data collections.

### Acknowledgments

We would like to extend our sincere gratitude to the reviewers for their thoughtful and constructive comments, which significantly improved the quality and clarity of this paper.

### Supplementary material

Supplementary material is available at ISME Communications online.

# **Conflicts of interest**

No conflict of interest.

# Funding

The research of Aditya Mishra is supported by the generous support of a startup grant from the University of Georgia. Initial support was provided by the Flatiron Institute, Simons Foundation. The research of David Blei is supported by funding from the National Science Foundation (NSF) under grants IIS-2127869 and DMS-2311108, the Office of Naval Research (ONR) under grant N000142412243, and the Simons Foundation. The research of Jed Fuhrman is supported by the Simons Collaboration on Computational Biogeochemical Modeling of Marine Ecosystems (CBIOMES) through grant 549943 to J.A.F., and by the NSF under grant EF-2125142 to J.A.F. The research of Christian L. Müller is supported by the Simons Collaboration on CBIOMES through grant 986803.

# Data availability

We have used microbial species abundance data from the Tara Ocean Expedition, available at (http://ocean-microbiome.embl. de/companion.html). The source code required to reproduce the results in this article is freely available at (https://github.com/ amishra-stats/vi-midas).

### References

- Fox G, Woese C. Phylogenetic structure of the prokaryotic domain. PNAS 1977;74:5088–90. https://doi.org/10.1073/ pnas.74.11.5088
- Pace NR, Stahl DA, Lane DJ et al. The analysis of natural microbial populations by ribosomal RNA sequences. In: Marshall K.C. (ed.), Advances in Microbial Ecology, Vol. 9. Boston, MA: Springer, 1986, 1–55.

- Olsen GJ, Lane DJ, Giovannoni SJ et al. Microbial ecology and evolution: a ribosomal RNA approach. Ann Rev Microbiol 1986;40: 337–65. https://doi.org/10.1146/annurev.mi.40.100186.002005
- Tyson GW, Chapman J, Hugenholtz P et al. Community structure and metabolism through reconstruction of microbial genomes from the environment. Nature 2004;428:37–43. https:// doi.org/10.1038/nature02340
- Schloss PD, Westcott SL, Ryabin T et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol* 2009;**75**:7537–41. https://doi.org/10.1128/ AEM.01541-09
- Callahan BJ, Sankaran K, Fukuyama JA et al. Bioconductor workflow for microbiome data analysis: from raw reads to community analyses. F1000Research 2016;5. https://doi.org/10.12688/ f1000research.8986.1
- Bolyen E, Rideout JR, Dillon MR et al. Reproducible, interactive, scalable and extensible microbiome data science using qiime 2. Nat Biotechnol 2019;37:852–7. https://doi.org/10.1038/ s41587-019-0209-9
- Turnbaugh PJ, Ley RE, Hamady M et al. The human microbiome project. Nature 2007;449:804–10. https://doi.org/10.1038/ nature06244
- Goodrich JK, Waters JL, Poole AC et al. Human genetics shape the gut microbiome. Cell 2014;159:789–99. https://doi.org/10.1016/j. cell.2014.09.053
- Scholtens S, Smidt N, Swertz MA et al. Cohort profile: lifeLines, a three-generation cohort study and biobank. Int J Epidemiol 2015;44:1172–80. https://doi.org/10.1093/ije/dyu229
- Ikram MA, Brusselle GGO, Murad SD et al. The Rotterdam study: 2018 update on objectives, design and main results. Eur J Epidemiol 2017;32:807–50. https://doi.org/10.1007/ s10654-017-0321-4
- McDonald D, Hyde E, Debelius JW et al. American gut: an open platform for citizen science microbiome research. mSystems 2018;3:10–1128. https://doi.org/10.1128/mSystems.00031-18
- Gilbert JA, Jansson JK, Knight R. The earth microbiome project: successes and aspirations. BMC Biol 2014;12:69. https://doi. org/10.1186/s12915-014-0069-1
- 14. Sunagawa S, Coelho LP, Chaffron S et al. Structure and function of the global ocean microbiome. *Science* 2015;**348**:1261359.
- Bien J, Yan X, Simpson L et al. Tree-aggregated predictive modeling of microbiome data. Sci Rep 2021;11:1–13.
- Guidi L, Chaffron S, Bittner L et al. Plankton networks driving carbon export in the oligotrophic ocean. Nature 2016;532: 465–70.
- Gevers D, Kugathasan S, Denson LA et al. The treatment-naive microbiome in new-onset Crohn's disease. Cell Host Microbe 2014;15:382–92. https://doi.org/10.1016/j.chom.2014.02.005
- Mishra A, Müller CL. Robust regression with compositional covariates. Comput Stat Data Anal 2022;165:107315. https://doi. org/10.1016/j.csda.2021.107315
- Faust K, Raes J. CoNet app: inference of biological association networks using Cytoscape. F1000Research 2016;5:1519. https:// doi.org/10.12688/f1000research.9050.1
- Holmes I, Harris K, Quince C. Dirichlet multinomial mixtures: generative models for microbial metagenomics. PLoS One 2012;7:e30126. https://doi.org/10.1371/journal.pone.0030126
- Lee S, Chugh PE, Shen H et al. Poisson factor models with applications to non-normalized microrna profiling. Bioinformatics 2013;29:1105–11. https://doi.org/10.1093/bioinformatics/btt091
- Chen J, Li H. Variable selection for SPARSE DIRICHLETmultinomial regression with an application to microbiome data analysis 1. Ann Appl Stat 2013;7:418–42.

- Zhang X, Mallick H, Tang Z et al. Negative binomial mixed models for analyzing microbiome count data. BMC bioinformatics 2017;18:4. https://doi.org/10.1186/s12859-016-1441-7
- Tianchen X, Demmer RT, Li G. Zero-inflated poisson factor model with application to microbiome read counts. *Biometrics* 2021;**77**:91–101.
- Mishra AK, Müller CL. Negative binomial factor regression with application to microbiome data analysis. Stat Med 2022;41: 2786–803. https://doi.org/10.1002/sim.9384
- Kurtz ZD, Müller CL, Miraldi ER. et al. Sparse and compositionally robust inference of microbial ecological networks. PLoS Comput Biol 2015;11:e1004226. https://doi.org/10.1371/journal.pcbi.1004226
- Yoon G, Gaynanova I, Müller CL. Microbial networks in springsemi-parametric rank-based correlation and partial correlation estimation for quantitative microbiome data. Front Genet 2019;10:516. https://doi.org/10.3389/fgene.2019.00516
- Peschel S, Müller CL, von Mutius E et al. NetCoMi: network construction and comparison for microbiome data in R. Brief Bioinform 2021;22:bbaa290.
- 29. Gleich SJ, Cram JA, Weissman JL. et al. Netgam: using generalized additive models to improve the predictive power of ecological network analyses constructed using time-series data. *ISME Commun* 2022;**2**:1–9.
- Ruiz FJR, Athey S, Blei DM. Shopper: a probabilistic model of consumer choice with substitutes and complements. The Annals of Applied Statistics 2020;14:1–27.
- McMurdie PJ, Holmes S. Phyloseq: an r package for reproducible interactive analysis and graphics of microbiome census data. PLoS One 2013;8:1–11.
- Blei DM, Kucukelbir A, McAuliffe JD. Variational inference: a review for statisticians. J Am Stat Assoc 2017;112:859–77. https:// doi.org/10.1080/01621459.2017.1285773
- Chiquet J, Mariadassou M, Robin S. Variational inference for probabilistic poisson pca. Ann Appl Stat 2018;12:2674–98. https:// doi.org/10.1214/18-AOAS1177
- Gibson T, Gerber G. Robust and scalable models of microbiome dynamics. In: International Conference on Machine Learning. PMLR, 2018, 1763–72.
- Harrison JG, John Calder W, Shastry V et al. Dirichletmultinomial modelling outperforms alternatives for analysis of microbiome and other ecological count data. Mol Ecol Resour 2020;20:481–97. https://doi.org/10.1111/1755-0998.13128
- Liu T, Peirong X, Yueyao D et al. Mzinbva: variational approximation for multilevel zero-inflated negative-binomial models for association analysis in microbiome surveys. Brief Bioinform 2022;23:bbab443. https://doi.org/10.1093/bib/bbab443
- Zeng Y, Zhao H, Wang T. Model-based microbiome data ordination: a variational approximation approach. J Comput Graph Stat 2021;30:1036–48. https://doi.org/10.1080/10618600.2021. 1882467
- Moran MA. The global ocean microbiome. Science 2015; 350: aac8455.
- Pesant S, Not F, Picheral M et al. Open science resources for the discovery and analysis of tara oceans data. Scientific data 2015;2: 1–16.
- Ashkezari MD, Hagen NR, Denholtz M et al. Simons collaborative marine atlas project (simons cmap): an open-source portal to share, visualize, and analyze ocean data. *Limnol Oceanogr Methods* 2021;19:488–96.
- Sunagawa S, Mende DR, Zeller G et al. Metagenomic species profiling using universal phylogenetic marker genes. Nat Methods 2013;10:1196–9. https://doi.org/10.1038/nmeth.2693

- 42. McNichol J. Sunagawa miTAG annotations. 2020. Retrieved from: https://github.com/jcmcnch/Sunagawa-miTAG-annotations.
- Quast C, Pruesse E, Yilmaz P et al. The silva ribosomal RNA gene database project: improved data processing and web-based tools. Nucleic Acids Res 2012;41:D590–6. https://doi.org/10.1093/ nar/gks1219
- Cameron AC, Trivedi PK. Regression analysis of count data. Vol. 53. Cambridge University Press, 2013. https://doi.org/10. 1017/CB09781139013567
- Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with deseq2. *Genome* Biol 2014;15:1–21.
- 46. Duncan Wadsworth W, Argiento R, Guindani M et al. An integrative Bayesian dirichlet-multinomial regression model for the analysis of taxonomic abundances in microbiome data. BMC Bioinformatics 2017;18:1–12.
- Sankaran K, Holmes SP. Latent variable modeling for the microbiome. Biostatistics 2019;20:599–614. https://doi.org/10.1093/ biostatistics/kxy018
- Sohn MB, Li H. A glm-based latent variable ordination method for microbiome samples. Biometrics 2018;74:448–57. https://doi. org/10.1111/biom.12775
- Aitchison J. The Statistical Analysis of Compositional Data. Caldwell, NJ: Blackburn Press, 2003.
- 50. De La Cruz R, Kreft JU. Geometric mean extension for data sets with zeros. arXiv preprint arXiv:1806.06403. 2018.
- 51. Longhurst AR. Ecological geography of the sea. Elsevier, 2010.
- Lima-Mendez G, Faust K, Henry N et al. Determinants of community structure in the global plankton interactome. Science 2015;348:1262073. https://doi.org/10.1126/science.1262073
- Jordan MI, Ghahramani Z, Jaakkola TS et al. An introduction to variational methods for graphical models. Mach Learn 1999;37: 183–233. https://doi.org/10.1023/A:1007665907178
- Wainwright MJ, Jordan MI et al. Graphical models, exponential families, and variational inference. Found Trends Mach Learn 2008;1:1–305.
- 55. Kucukelbir A, Tran D, Ranganath R et al. Automatic differentiation variational inference. J Mach Learn Res 2017;**18**:430–74.
- Gelman AB, Carlin JB, Stern HS et al. Bayesian Data Analysis, 3rd edn. Boca Raton FL: CRC Press. [Google Scholar], 2013, https:// doi.org/10.1201/b16018
- Carpenter B, Gelman A, Hoffman MD et al. Stan: a probabilistic programming language. J Stat Softw 2017;76:1–32. https://doi. org/10.18637/jss.v076.i01
- Rubin DB. Bayesianly justifiable and relevant frequency calculations for the applied statistician. Ann Stat 1984;12:1151-72. https://doi.org/10.1214/aos/1176346785
- Clauset A, Newman MEJ, Moore C. Finding community structure in very large networks. Phys Rev E 2004;70:066111. https://doi. org/10.1103/PhysRevE.70.066111
- Aldunate M, De la Iglesia R, Bertagnolli AD et al. Oxygen modulates bacterial community composition in the coastal upwelling waters off Central Chile. Deep-Sea Res II Top Stud Oceanogr 2018;156:68–79. https://doi.org/10.1016/j.dsr2.2018.02.001
- West NJ, Lepère C, de O Manes C-L et al. Distinct spatial patterns of sar11, sar86, and Actinobacteria diversity along a transect in the ultra-oligotrophic south pacific ocean. Front Microbiol 2016;7:234.
- 62. Yilmaz P, Yarza P, Rapp JZ et al. Expanding the world of marine bacterial and archaeal clades. Front Microbiol 2016;**6**:1524.
- 63. Qiang Zheng Y, Wang RX, Lang AS *et al.* Dynamics of heterotrophic bacterial assemblages within Synechococcus cultures. *Appl Environ Microbiol* 2018;**84**:e01517–7.

- 64. Baskaran V, Patil PK, Leo Antony M *et al.* Microbial community profiling of ammonia and nitrite oxidizing bacterial enrichments from brackishwater ecosystems for mitigating nitrogen species. Sci *Rep* 2020;**10**:1–11.
- Friedman J, Alm EJ. Inferring correlation networks from genomic survey data. PLoS Comput Biol 2012;8:e1002687. https://doi. org/10.1371/journal.pcbi.1002687
- Kurtz ZD, Müller CL, Miraldi ER *et al.* Sparse and compositionally robust inference of microbial ecological networks PLoS Computational Biology 2015;**11**:e1004226. https://doi.org/10.1371/ journal.pcbi.1004226
- Li Y, Tang K, Zhang L et al. Coupled carbon, sulfur, and nitrogen cycles mediated by microorganisms in the water column of a shallow-water hydrothermal ecosystem. Front Microbiol 2018;9:2718. https://doi.org/10.3389/fmicb.2018.02718
- Moynihan MA, Goodkin NF, Morgan KM et al. Coral-associated nitrogen fixation rates and diazotrophic diversity on a nutrientreplete equatorial reef. ISME J 2022;16:233–46. https://doi. org/10.1038/s41396-021-01054-1
- Meisel JS, Hannigan GD, Tyldsley AS et al. Skin microbiome surveys are strongly influenced by experimental design. J Invest Dermatol 2016;136:947–56. https://doi.org/10.1016/j. jid.2016.01.016
- Gobbi A, Acedo A, Imam N et al. A global microbiome survey of vineyard soils highlights the microbial dimension of viticultural terroirs. Commun Biol 2022;5:241. https://doi.org/10.1038/ s42003-022-03202-5
- Shaffer JP, Nothias L-F, Thompson LR et al. Standardized multiomics of earth's microbiomes reveals microbial and metabolite diversity. Nat Microbiol 2022;7:2128–50. https://doi.org/10.1038/ s41564-022-01266-x

- 72. The integrative human microbiome project. Nature 2019;**569**: 641–8.
- 73. Smith AN, Hennon GMM, Zinser ER et al. Comparing prochlorococcus temperature niches in the lab and across ocean basins. Limnol Oceanogr 2021;66:2632–47. https://doi. org/10.1002/lno.11777
- Rigonato J, Budinich M, Murillo AA et al. Ocean-wide comparisons of mesopelagic planktonic community structures. ISME Commun 2023;3:83. https://doi.org/10.1038/s43705-023-00279-9
- Feng CX. A comparison of zero-inflated and hurdle models for modeling zero-inflated count data. J Stat Distrib Appl 2021;8:8. https://doi.org/10.1186/s40488-021-00121-4
- 76. Ostner J, Carcy S, Müller CL. Tasccoda: Bayesian treeaggregated analysis of compositional amplicon and singlecell data. Front Genet 2021;12:766405. https://doi.org/10.3389/ fgene.2021.766405
- 77. Dillon JV, Langmore I, Tran D *et al.* Tensorflow distributions. arXiv preprint arXiv:1711.10604. 2017.
- 78. Paszke A, Gross S, Massa F et al. PyTorch: An imperative style, high-performance deep learning library. arXiv preprint arXiv:1912.01703. 2019.
- 79. Kingma DP. Variational inference & deep learning: a new synthesis. 2017.
- Integrative HMP. The integrative human microbiome project: dynamic analysis of microbiome-host omics profiles during periods of human health and disease. *Cell Host Microbe* 2014;16: 276–89.
- Forslund SK, Chakaroun R, Zimmermann-Kogadeeva M et al. Combinatorial, additive and dose-dependent drug-microbiome associations. Nature 2021;600:500–5. https://doi.org/10.1038/ s41586-021-04177-9