

SOFTWARE

Open Access



specificity: an R package for analysis of feature specificity to environmental and higher dimensional variables, applied to microbiome species data

John L. Darcy^{1*}, Anthony S. Amend^{2,3}, Sean O. I. Swift², Pacifica S. Sommers⁴ and Catherine A. Lozupone¹

Abstract

Background: Understanding the factors that influence microbes' environmental distributions is important for determining drivers of microbial community composition. These include environmental variables like temperature and pH, and higher-dimensional variables like geographic distance and host species phylogeny. In microbial ecology, "specificity" is often described in the context of symbiotic or host parasitic interactions, but specificity can be more broadly used to describe the extent to which a species occupies a narrower range of an environmental variable than expected by chance. Using a standardization we describe here, Rao's (Theor Popul Biol, 1982. [https://doi.org/10.1016/0040-5809\(82\)90004-1](https://doi.org/10.1016/0040-5809(82)90004-1), Sankhya A, 2010. <https://doi.org/10.1007/s13171-010-0016-3>) Quadratic Entropy can be conveniently applied to calculate specificity of a feature, such as a species, to many different environmental variables.

Results: We present our R package *specificity* for performing the above analyses, and apply it to four real-life microbial data sets to demonstrate its application. We found that many fungi within the leaves of native Hawaiian plants had strong specificity to rainfall and elevation, even though these variables showed minimal importance in a previous analysis of fungal beta-diversity. In Antarctic cryoconite holes, our tool revealed that many bacteria have specificity to co-occurring algal community composition. Similarly, in the human gut microbiome, many bacteria showed specificity to the composition of bile acids. Finally, our analysis of the Earth Microbiome Project data set showed that most bacteria show strong ontological specificity to sample type. Our software performed as expected on synthetic data as well.

Conclusions: *specificity* is well-suited to analysis of microbiome data, both in synthetic test cases, and across multiple environment types and experimental designs. The analysis and software we present here can reveal patterns in microbial taxa that may not be evident from a community-level perspective. These insights can also be visualized and interactively shared among researchers using *specificity*'s companion package, *specificity.shiny*.

Keywords: Species distributions, Biogeography, Multi-omic data

Introduction

The word "specificity" has uses across multiple disciplines. In ecology, and especially for microbes, "specificity" is often used in the context of symbiotic interactions; for example the specificity of a parasitic species may be the degree to which it associates with a narrow consortia of host species [1–3]. In pharmacology and biochemistry,

*Correspondence: darcyj@colorado.edu

¹ Division of Biomedical Informatics and Personalized Medicine, University of Colorado School of Medicine, Aurora, CO, USA
Full list of author information is available at the end of the article



© The Author(s) 2022. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

specificity can describe the “narrowness of the range of substances with which an antibody or other agent acts” [4]. Synthesizing these definitions, we arrive at a general concept of specificity, where a feature (e.g. a species) is specific to some variable (e.g. elevation) when it occupies or is otherwise associated with a limited breadth of that variable.

This definition is consistent across multiple variable types. For example, a species that is found only across a narrow band of elevation, perhaps between 200 and 500 m above sea level, would have stronger specificity to elevation than a species that is found between sea level and 1000 m. This is similar to a parasite that is found only within a narrow clade of hosts; it has stronger host specificity than a parasite that is found across a much wider phylogenetic range [2]. This concept can be expanded even farther, to diversity of some co-occurring feature class. For example, to metabolites that co-occur with bacteria in the human gut microbiome (microbes within the human gut). Under our definition, a microbe may have specificity to a narrow range of metabolomic compositions. Furthermore, specificity as we describe it here, is not the same as bipartite network specialization like $H2'$, d' , and $NODF$ metrics [5, 6]. Those metrics apply to strictly categorical contingency data, for example a matrix of observation counts where columns are pollinator species and rows are plant species. Instead, our generalized specificity approach is best suited to continuous data.

Our generalized specificity analysis has several benefits over modeling a microbe’s relative abundance using a variable of interest, since specificity analysis has no underlying model. First, high-throughput sequencing (HTS) microbiome data notoriously contain many zeroes, corresponding to the lack of an observation of species in a samples. Disregarding the difficulties in modeling such data, which can certainly be overcome [7], these data are perfect for the “specificity approach”. This is because the alternative hypothesis of a specificity analysis (the focal species encounters less environmental heterogeneity than expected by chance) includes cases where the focal species only occupies a limited range of the variable of interest, being absent (zero) everywhere else. A further consideration in modeling approaches is non-monotonic relationships between species and environmental variables. For example, a species may have specificity to intermediate elevations, so its density function of elevation would be non-monotonic, or even multimodal; and that’s just one species. Within a HTS microbiome dataset, species may be expected to run the gamut of distribution shapes and modalities. Variables of interest also present their own challenges to modeling, since variables may be vectors (e.g. elevation, pH), distance or dissimilarity

matrices (e.g. geographic distance, beta-diversity), phylogenies, or even sample-type ontologies. The generalized specificity approach we present here can accommodate all of the aforementioned variable types, unlike other approaches where the statistics used to understand microbe-environment relationships are restricted by variable type. Furthermore, our approach does not produce a model, or answer the question “across what range of the variable does the species occur”. Instead, we quantify the extent to which the species occupies a limited breadth of that variable without the need for such a model.

Meaningfully applying this general idea of specificity to multiple data types is challenging because of the different specificity metrics available to different kinds of data. With host phylogenetic data, specificity may be calculated as phylodiversity [8], or host phylogenetic entropy [9], or host richness [10]. However, with other data types these metrics are not useful—one cannot calculate phylogenetic entropy of elevation, for example. Per our definition above, specificity must be a measure of the breadth (i.e. heterogeneity, diversity) of an environmental variable occupied by the focal species. With a variable like elevation, a naive specificity metric may be as simple as the variance in elevation where the focal species is present, or weighted variance for a more intuitive approach. However, such a metric would not be applicable to phylogenetic data sets because it is limited to 1-dimensional data types (i.e. column vectors). Furthermore, we wanted our general idea of specificity to be useful for dissimilarity matrices. We found that Rao’s Quadratic Entropy [11–13] is a convenient diversity metric that can be applied to all abovementioned data, with a modicum of standardization (detailed in our “Methods” section).

Here, we present a software package written in R and C++ that implements a generalized specificity analysis. Our package, *specificity*, calculates specificity values for each species in a sample-by-species matrix. In microbiology, this data structure often appears as a table of OTUs (operational taxonomic units; a substitute for species) or ASVs (Amplicon Sequence Variants; OTUs represented by unique sequences after applying a denoiser such as DADA2 [14]). We simulated species distributions with varying strengths of specificity, and used those simulated data to validate our implementation. Our simulations were also used to ensure that specificity is not sensitive to occupancy (i.e. in how many samples a species appears), which is a significant improvement compared to the standardized effect size (SES) method [2, 15], and methods that use un-weighted (presence-absence) species data [10]. Our simulations also confirmed that the specificity we calculate here is scale-invariant with regard to environmental/phylogenetic data, and also to focal species abundance data. To illustrate how specificity can be used,

we applied our software to four previously-published microbiological data sets, each from different environments: fungi living within the leaves of native Hawaiian plants, human gut microbiome bacteria, bacteria living within Antarctic glaciers, and the global Earth Microbiome Project data set.

Methods

RQE

We calculate specificity using Rao’s metric [11, 12]. It is sometimes abbreviated FDQ for quadratic functional diversity, but since we use the same mathematics in a non-functional context, here we simply refer to the metric as *RQE* (Rao’s Quadratic Entropy, similar to the use of “QE” by its inventor. *RQE* is the sum of the element-wise product of two square matrices (excluding the diagonal). In our use, the first matrix (*D*) is a dissimilarity matrix containing differences between samples (Fig. 1). For example, in the context of phylogenetic specificity these differences are phylogenetic distances (i.e. cophenetic distances) between hosts. Samples from the same host species have 0 distance. The second matrix (*W*) contains all pairwise products of weights for the focal species. Given a column vector of species weights *p* from a site-by-species matrix (“OTU table”), *W_{ij}* contains the product of the abundances (weights) of the focal species at sample *i* and sample *j*: *p_ip_j*. Via *D* ∘ *W* (or *D_{ij}p_ip_j* for a single pairwise product; Eq. 1), we weight matrix *D* to

up-weight distances between samples where the focal species occurred, and down-weight distances between samples where the focal species was absent in either. We use the term “weights” to describe *p* because the values within could be relative abundances or any other metric that describes the importance of a species within a sample. Conversely, we have chosen to focus this manuscript on “species”, but note that *p* could be a vector of weights for any feature (a type of rock, a metabolite, etc).

$$RQE = \sum_{i \neq j} D_{ij} p_i p_j \tag{1}$$

With *RQE*, a focal species with strong specificity has relatively high weights for low differences. This metric was originally developed for phylogenetic distances, but here we apply it to many different *D* matrices, including euclidean transformations of 1-dimensional data (e.g. pairwise elevational difference), or more complex 2-dimensional data like Bray–Curtis dissimilarity between host metabolomic profiles.

As such, a species with “perfect” specificity will always have *RQE* = 0. For example, consider a focal species *S* that can be found in habitats A, B, C, or D, with multiple samples collected for each category (Fig. 1). If *S* is only found in samples from habitat A, matrices *D* and *W* will contain zeroes in opposite positions, resulting in *RQE* = 0. Note that weights near zero can also act similarly to zero since this is a weighted metric. In this way, a focal species can occupy every single sample (all values of *p* are nonzero positive) and still have *RQE* near zero. This is important so that spurious species detections do not significantly contribute to specificity. For example, in a DNA sequencing experiment, small amounts of contamination may occur during DNA extraction or library preparation. The magnitude of that contamination is expected to be small compared to the signal in an actual sample, but may result in spurious species detections. However, because these contaminants would be expected to be rare in the sample, their weights would be low in samples where they are noise, and high in samples where they are signal.

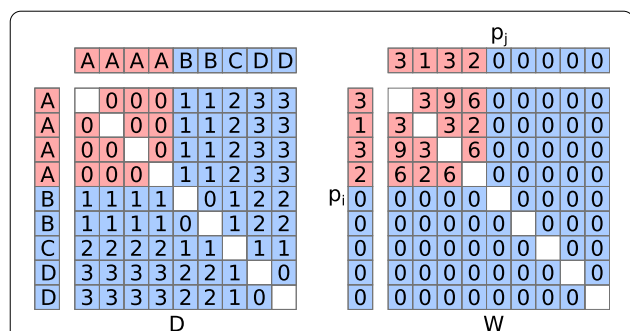


Fig. 1 *RQE* as it applies to specificity. In this example, two matrices are shown, *D* and *W*. *D* is an environmental dissimilarity matrix, describing how different are several environment types, A through D, with multiple samples represented for each environment type. Note that diagonals are empty because they are not used; see Eq. 1. Matrix *W* is the pairwise product of species weights *p* (Eq. 1). In this example, the focal species is perfectly specific to habitat A, which can be seen in *p*. Data corresponding to species detections are colored in red, and species absences in blue. The product *D* ∘ *W* (= *D_{ij}p_ip_j*) will be all zeroes for this example, because this example shows perfect specificity. Thus, the sum of that product, *RQE*, will be zero. If *p* had relatively small values instead of its zeroes, for example 0.25, those small values would still down-weight their corresponding larger differences in *D* and produce a signal of specificity, compared to random permutations of *p* which produce much higher *RQE* values

Standardization

While empirical *RQE* is calculated as described above, it must be standardized in order to compare effect size between different variables and different species, because the metric’s scale is dependent on the scales of *D* and *p*. Phylogenetic specificity approaches have previously used a standardized effect size (SES) approach [2], but we found that SES has unfortunate properties when used with our generalized approach to specificity. Critically, SES was

highly sensitive to occupancy, which is the number of samples a species occupies. One would expect the strength of specificity to be lower when a species occupies more samples, because this means the species must occur in a broader range of habitat. However, SES counter-productively yields stronger specificity for more occupant species, and also for species with more even distributions. This is because SES is standardized using a distribution of values generated with permuted species weights. If all weights are similar (high evenness), the standard deviation of that distribution will be small, leading to a strong SES. SES is also undesirable because for a given species, it is tightly associated with that species' p value (the probability of SES being as strong or stronger, despite the null hypothesis being true), enough so that a suggested remedy for yet other problems with SES is to use a probit transformed p value in its place [15].

We standardize RQE by leveraging the fact that for perfect specificity, empirical RQE (RQE_{emp} , Fig. 1) equals zero. Our statistic, which we simply call $Spec$, ranges from -1 to 1 , with 0 as the null hypothesis that species weights are randomly ordered with regard to sample identity. Similar to SES, RQE_{sim} is a vector of RQE values calculated using random permutations of p . The distribution of RQE_{sim} is used for calculating p values, and its central tendency is defined as $Spec = 0$. This central tendency can be calculated several ways using our software, but the default is to use the mean of RQE_{sim} since in our testing, mean and median showed highly concordant $Spec$ results (Additional file 1: Fig. S1). The equation for $Spec$ (Eq. 2) is a piecewise function, with the two parts corresponding to specificity and generality, respectively. In the case of specificity, $Spec$ simply scales RQE_{emp} relative to the center of RQE_{sim} such that perfect specificity returns -1 , and the null hypothesis returns 0 . In this case, "null hypothesis" refers to RQE_{emp} being the expected value of RQE_{sim} .

$$Spec = \begin{cases} RQE_{emp} \leq \overline{RQE_{sim}}, & \frac{RQE_{emp} - \overline{RQE_{sim}}}{\overline{RQE_{sim}}} \\ RQE_{emp} > \overline{RQE_{sim}}, & \frac{RQE_{emp} - \overline{RQE_{sim}}}{\overline{RQE_{max}} - \overline{RQE_{sim}}} \end{cases} \quad (2)$$

The case of generality is slightly more complicated, since there is no intuitive maximum theoretical RQE value. "Generality" in this context refers to species that encounter greater environmental heterogeneity than expected by chance. We find that maximum value computationally, and standardize $Spec$ as a proportion of that value (see Eq. 2). For each p there exists an optimized permutation that yields the highest possible RQE value, RQE_{max} . We use a genetic algorithm (GA) with Population Based Incremental Learning [16] to search permutations of p that create RQE_{max} . Our GA begins with a population of surrogate vectors initialized via random permutations

of p (default 150), and random swaps of p (a swap being the pairwise substitution of two values within the vector; default 150), and also p itself. Each generation, the GA calculates RQE for each vector in the population, then keeps some of the vectors with the highest RQE value (default 5). The next generation is composed of those kept vectors, and random swaps thereof until the total original population size is met (default 301). Our swapping algorithm can also use a stochastic number of swaps per vector per generation (including initialization), drawn at random from a user-defined set (default 1,1,2,3). In addition to swapping, mutation can be performed by crossover via the PMX algorithm [17], which is used because it incorporates both order and position of both parents, which is required for this problem. However, in our testing we found that crossover did not improve GA efficiency, so the default operation is not to perform PMX. The GA runs for a fixed number of generations (default 400), or until a number of generations have passed with no improvement (default 10). These parameters were chosen because they performed well on the data sets we analyze here, meaning that species reached the early termination condition.

Our GA is relatively computationally intensive, consuming the majority of computational time for a given specificity analysis even though it is only used for a minority of species. This is unlikely to be a concern on smaller data sets (i.e. a few hundred samples), but since many users may not be interested in "general" species, another option is to scale $Spec$ for all species using the top half of Eq. 2 instead. This is considerably faster, and the user can either discard "general" species as uninteresting, or choose to interpret $Spec > 0$ within an ordinal framework (a brief analysis showed the results of this approach and those of the GA are strongly correlated; Additional file 1: Fig. S2).

Hypothesis testing

For the $Spec$ calculation above, a p value may be calculated as the proportion of RQE_{sim} values that are lower than RQE_{emp} . The default operation of our software is to adjust p values calculated for different species from the same variable for multiple hypothesis testing by applying the Benjamini–Hochberg procedure [18].

Features of $Spec$

$Spec$ captures signal of specificity to simulated vector, matrix, and phylogenetic data (Additional file 1: Figs. S3, S4). It is insensitive to species occupancy (Additional file 1: Figs. S5, S6) and is insensitive to the number of samples within a data set (Additional file 1: Fig. S7). $Spec$ is also scale-invariant, independently in regard to p and D inputs (Additional file 1: Fig. S8). It is sensitive to

multimodality, and multimodal species distributions are still detected as exhibiting significant specificity by *Spec* (Additional file 1: Fig. S9).

Validation analyses

Species were simulated with varying levels of specificity by drawing from a normal distribution centered on an artificial “optimum” environmental location (e.g. elevation of 300 m). Varying specificities were achieved by widening the standard deviation of that distribution, or by mixing the normal distribution with varying proportions of a uniform distribution. Multimodal specificity was simulated similarly by combining multiple distributions. Specificity of simulated species was analyzed using our software. Occupancy of simulated species was increased or decreased by randomly substituting simulated weights with zero, and specificity was analyzed across an occupancy gradient using that approach. Real data (see Endophyte analysis, below) were randomly downsampled to test the sensitivity of *Spec* to sample size. Real data were also re-ordered to create simulated high-specificity species that use empirical distributions of weights, and then those simulated species were subjected to a swapping algorithm that gradually introduced entropy into the species. The swapping algorithm swaps values from two randomly selected positions in p (Eq. 1). This was done recursively for 1000 generations (2 swaps per generation), saving p each time. Our software was then run on all vectors simulated this way.

Analysis of endophyte data

Data from Hawaiian foliar endophytic fungi [19] were downloaded from FigShare. These are illumina MiSeq data of the Internal Transcribed Spacer (ITS) region of fungal ribosomal RNA, from 760 samples collected from the leaves of native Hawaiian plants across five islands in the Hawaiian archipelago. This data set is also included in our R package. The features under investigation in this analysis were fungal OTUs. Data were transformed (“closed”) using total sum scaling, and fungal OTUs present in fewer than 10 samples were excluded from specificity analysis, because low-occupancy data can be unreliable (Additional file 1: Fig. S6). The remainder (416 OTUs) were run through our software using default settings except run in parallel using 20 CPU cores. Variables used in this analysis were NDVI (an index of vegetation density), elevation, evapotranspiration, rainfall, host plant phylogeny, and geographic distance between sample sites.

Analysis of Antarctic bacteria data

Data from Antarctic cryoconite hole bacteria [20] were obtained from the authors. Cryoconite holes are

isolated melt pools on the surface of glaciers, caused by debris from nearby slopes falling onto the glacier, and then melting into its surface. These holes form discrete microbial communities that have been described as “natural microcosms” [21]. This data set comprises 90 samples across three adjacent glaciers, and features are bacterial (16S rRNA) and eukaryal (18S rRNA) Amplicon Sequence Variants (ASVs; a type of OTU). Taxonomy was assigned to 18S rRNA ASVs using dada2 [14], and Bray–Curtis beta-diversity was calculated for only those ASVs that were determined to be algae. Analyses on 16S rRNA data were run and visualized as above, with variables N (Nitrogen), P (Phosphorus), pH, geographic distance, fungal Bray–Curtis beta-diversity (calculated from 18S rDNA data), and algal beta-diversity. Bacterial associations with individual glaciers (e.g. “ASV4 is found predominantly on Canada Glacier”) were computed using Dunn’s test [22], which is a nonparametric post-hoc test of difference in means.

Analysis of Human microbiome data

Data from Franzosa et al. [23] were downloaded as supplemental data from the online version of the article. These data contain both gut bacterial and archaeal species composition data as well as corresponding metabolomic data, collected from 220 adults with Crohn’s disease, ulcerative colitis, or healthy controls. Data were downloaded in a processed state, after the following procedures had been completed: species composition data from this study were derived from metagenomic data, which were assigned taxonomy and grouped into OTUs using MetaPhlan2 [24], and excluded samples that did not meet a 0.1% relative abundance threshold in at least 5 samples. Metabolite data were measured using positive and negative ion mode LC/MS, and were reported as parts per million. Metabolite identities were assigned programatically, and were clustered into broad classes per the Human Metabolome Database [25]. We subset the matrix of metabolome data by those classes, and used Euclidean distance to calculate the extent to which any two given samples differed in metabolomic composition within a given class (e.g. “how different is the composition of bile acids between sample A and sample B?”). Metabolite classes were excluded if they were totally absent in any sample, or if they contained fewer than 10 metabolites, which left 83 classes. specificity was used to calculate *Spec* for microbial OTUs to each metabolite class distance matrix.

Analysis of Earth Microbiome Project data

Data from the Earth Microbiome Project (EMP) [26] were compiled and downloaded from Qiita [27]. These data comprise a global sampling of 16S rRNA ASVs

produced by multiple studies. All of the studies followed a uniform protocol for collection, processing, and analysis of microbial data. A major component of the EMP is a rigid sample type ontology. The EMP Ontology (EMPO) was designed to categorically represent two main drivers of bacterial community composition: host association and salinity, for each sample that was collected. At the broadest level (EMPO1), samples were categorized as being host-associated or free living. At the intermediary level (EMPO2), samples were further divided into saline, non-saline, animal, plant, and fungus. The finest level (EMPO3) separated samples into 22 discrete substrate types (e.g. saline water, plant corpus, animal distal gut).

Because this data set is so large (28,842 samples and 309,469 ASVs), ASVs were excluded that were not present within at least 30 samples, samples were discarded if they had fewer than 5000 reads, and samples without ontological data were discarded (leaving 25,188 samples and 7014 ASVs). ASVs were excluded due to low occupancy to avoid spurious ASVs and to avoid low-abundance ASVs that do not perform well with *Spec*, and more importantly to keep computation size manageable for this massive data set. Samples were also discarded mainly due to computational concerns, with low-count samples being dropped first due to lower confidence in their proportional abundance calculations. The EMP ontology was transformed into a phylogeny using *specifity*'s “*onto2nwk*” function, which makes a cladogram within which all branch-lengths were set to 1. Specificity analysis was run using the ASV table and the ontological data. Database matches for individual species of interest were manually obtained using nucleotide BLAST [28] via the NCBI web portal, using the 16S rRNA sequence database as reference.

Implementation

specifity was written in the R programming language, with some functions written in C++ using *Rcpp* [29]. The general format of the package follows standard R package structure [30]. Unit testing was done using *testthat* [31]. *specifity.shiny* was written entirely in R, and uses the *shiny* [32] interactive web application framework. Both packages are free and open source software, licensed under the Gnu Public License (GPL). Installation is easily done using the “*install_github*” function of the R package *remotes* [33]; see data “Availability of data and materials” section for details.

Results and discussion

Hawaiian endophyte specificity analysis

We found that foliar endophytic fungi (FEF) from within the leaves of native Hawaiian plants exhibited strong and statistically significant specificity to several

environmental variables (Fig. 2), including variables that were only weakly associated with FEF community composition [19]. For example, in the original paper, rainfall and elevation were relatively weak predictors of FEF phylogenetic beta-diversity, but many FEF species show strong specificity to those variables in the present analysis. This reflects a fundamental difference between community-centric approaches (e.g. FEF community composition) vs. species-centric approaches like (e.g. specificity analysis). The signal of individual species is lost when a community is aggregated into a beta-diversity matrix or similar, and consequentially individual species within the community may even respond to environmental variables orthogonally to the community as a whole. Species that were strongly specific to rainfall or elevation are examples of this.

We found that many FEF species have strong and statistically significant specificity to geographic location, which makes sense given the discrete spatial structure of the Hawaiian islands [34], and that these FEF communities only are spatially structured up to distances of 36 km [19]. But geographic specificity may be an artifact of specificity to other variables with strong geographic autocorrelation. For example, in Hawaii (and elsewhere), rainfall is a spatially structured phenomenon [35], with nearby areas experiencing dramatically different rainfall averages as a consequence of aspect and elevation. Thus,

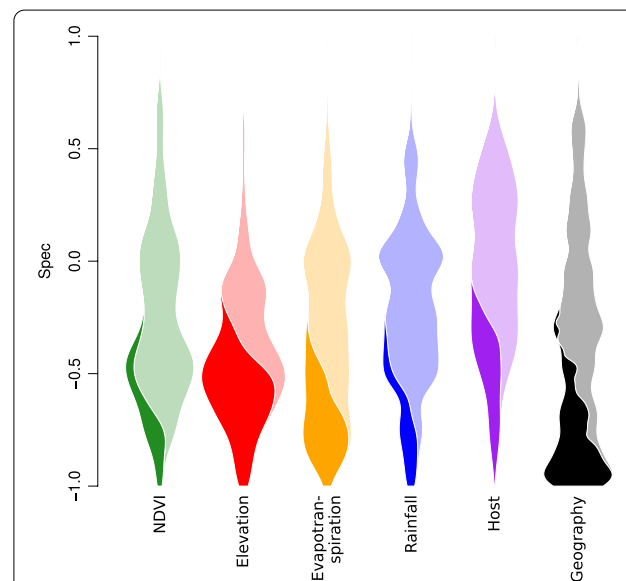


Fig. 2 Specificity of Hawaiian Foliar Endophytic Fungi. In this plot, the *Spec* values for 416 fungal species are plotted as violins for different variables. Since the number of species is the same for each variable, each violin has the same total area. Violin area is divided between species with statistically significant specificity (dark) versus species without (light)

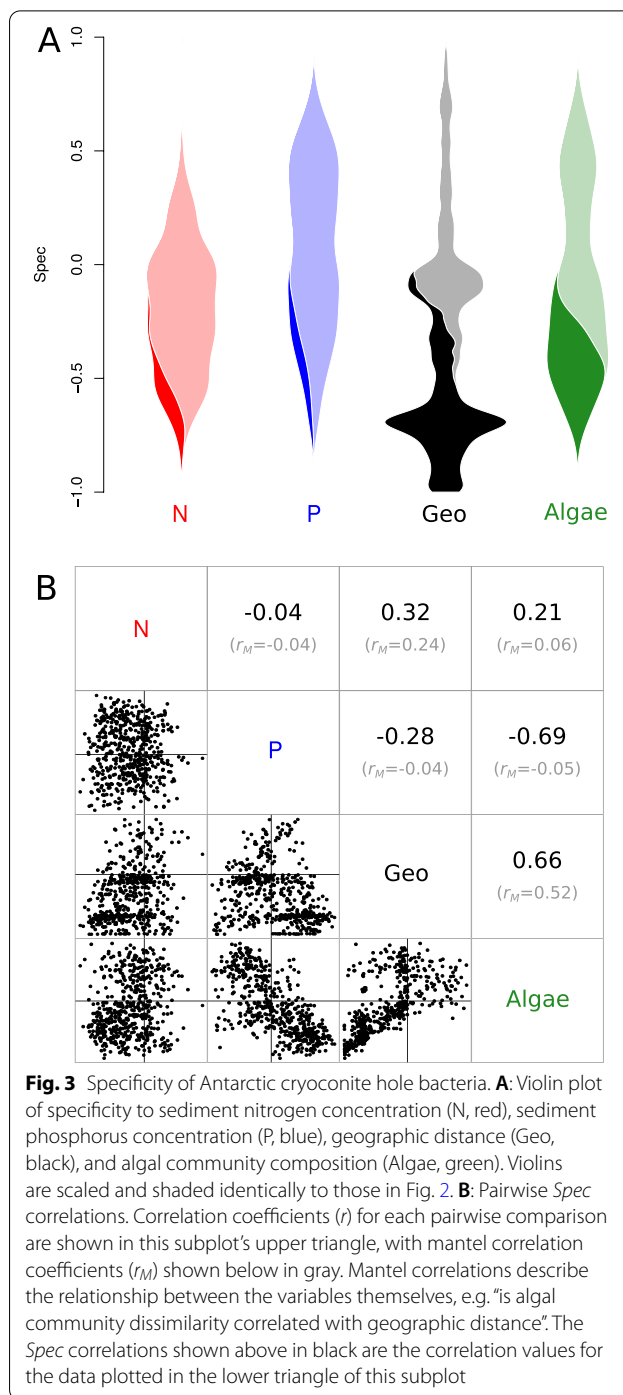
species that have strong specificity to rainfall likely also have strong specificity to geography, which is true in our analysis. Indeed, the same is true for elevation, albeit to a lesser extent Fig. 2.

One of the fungi with strongest specificity in our analysis was of the genus *Harknesia*, with closest BLAST match in the NCBI nucleotide database to *H. platyphylae*, a eucalyptus pathogen. In our data set, this fungus was found on multiple hosts, including *Metrosideros polymorpha*, which is in the same family as eucalyptus (*Myrtaceae*). This *Harknesia* was found exclusively within the interior of Hawaii island, at elevations between 1700 and 2000 m above sea level. Likely because of its strong geographic and elevational specificity, it also exhibited strong specificity to evapotranspiration, only being found in areas with 40–50 mm per year. Other fungi, such as an ASV most closely related to *Phaesosphaeria papayae*, exhibited strong specificity to elevation, evapotranspiration, and vegetation (NDVI), but were found on multiple islands across the Hawaiian archipelago. Thus, while geographic specificity can appear as specificity to geographically autocorrelated variables, this is not always the case.

Notable generalists ($Spec > 0$) in the endophyte dataset include the genus *Colletotrichum*, a globally distributed genus of plant pathogenic and endophytic fungi. Almost all agricultural crops are impacted by members of *Colletotrichum* and it is considered a ‘top ten’ fungal pathogen for molecular plant pathological research [36]. Of the 9 ASVs identified as *Colletotrichum*, none showed specialization to plant host or geography. Recently, genomic studies of this genus have provided insight into the genetic mechanisms behind host generalism and the activation of latent pathogenicity [37, 38]. Low specificity to geography and host within this dataset indicates that asymptomatic *Colletotrichum* species are widespread within the native Hawaiian flora.

Antarctic glacier bacteria specificity analysis

Similar to the FEF analysis above, bacteria living in cryoconite holes (isolated melt pools) on glaciers in Antarctica’s Taylor Valley [20, 21] exhibited strongest specificity to geographic distance (Fig. 3). This data set spanned three glaciers: Canada, Taylor, and Commonwealth, with equal sampling on each, but geographic distance accounts for distances within glaciers as well. The strong geographic specificity observed here reflects bacteria that are differentially abundant among glaciers, for example occupying only one or two of the three. The three glaciers, each flowing into a separate lake basin, are spaced along the 40-km length of the Taylor Valley, which stretches from the polar ice sheet to the coast. The cryoconite holes from the most inland glacier, Taylor Glacier, have fewer nutrients than those nearer the



coast, on Commonwealth Glacier [39]. The more inland cryoconite holes also have the lowest diversity of bacteria, while the holes nearest the coast support the most diverse bacterial communities [21]. Many bacterial species may therefore be specific to Commonwealth Glacier because it supports more species within its cryoconite holes than the other two glaciers. Besides the differences

in bacterial richness among the glaciers, the composition of the bacterial community turns over among glaciers, with different sequence variants dominating each glacier [20]. Biogeochemical differences within cryoconite holes among glaciers furthermore correspond with biogeochemical gradients along the valley in the surrounding soils [39]. The difference in dominant bacterial taxa on each glacier may primarily reflect (1) differences in which bacteria dominate the soils surrounding each glacier and therefore disperse onto each glacier, (2) a response to biogeochemical conditions within cryoconite holes on the glacier, or (3) an interaction of the two. Experimental microcosms manipulating dispersal and nutrient availability could help to parse out dominant controls on geographic specificity of bacteria in the future.

Strong bacterial specificity to co-occurring algal communities was expected, given the strong correlation between bacterial and eukaryal diversity previously observed in this supraglacial system [20] and elsewhere [40]. In our analysis, we found that specificity to algae was strongly negatively correlated with specificity to phosphorus ($r = -0.69$); even though those two variables were not strongly correlated with each other ($r_M = -0.05$). In other words, bacteria that are specific to algal community composition are not specific to sediment phosphorus concentration, and vice-versa. Using post-hoc tests, we found that bacteria with strong specificity to phosphorus concentration were predominantly associated with Taylor Glacier (but not exclusively), and that bacteria with strong specificity to algal community composition were predominantly found on Commonwealth Glacier.

Similarly, the correlation in *Spec* between geographic distance and algae (Fig. 3B) highlights a feature of specificity analysis using *Spec*: when comparing specificity of the same features to two different variables, *Spec* is likely to be strongly correlated when the variables have a linear relationship. That linear relationship can be seen in the Mantel correlation between pairwise geographic distance and algal beta-diversity ($r_M = 0.52$; Fig. 3B). But with variables that are weakly correlated, *Spec* may or may not be correlated between the variables. For example, difference in phosphorus concentration is not correlated with algal beta-diversity, but bacterial specificity to those variables is correlated (Fig. 3B).

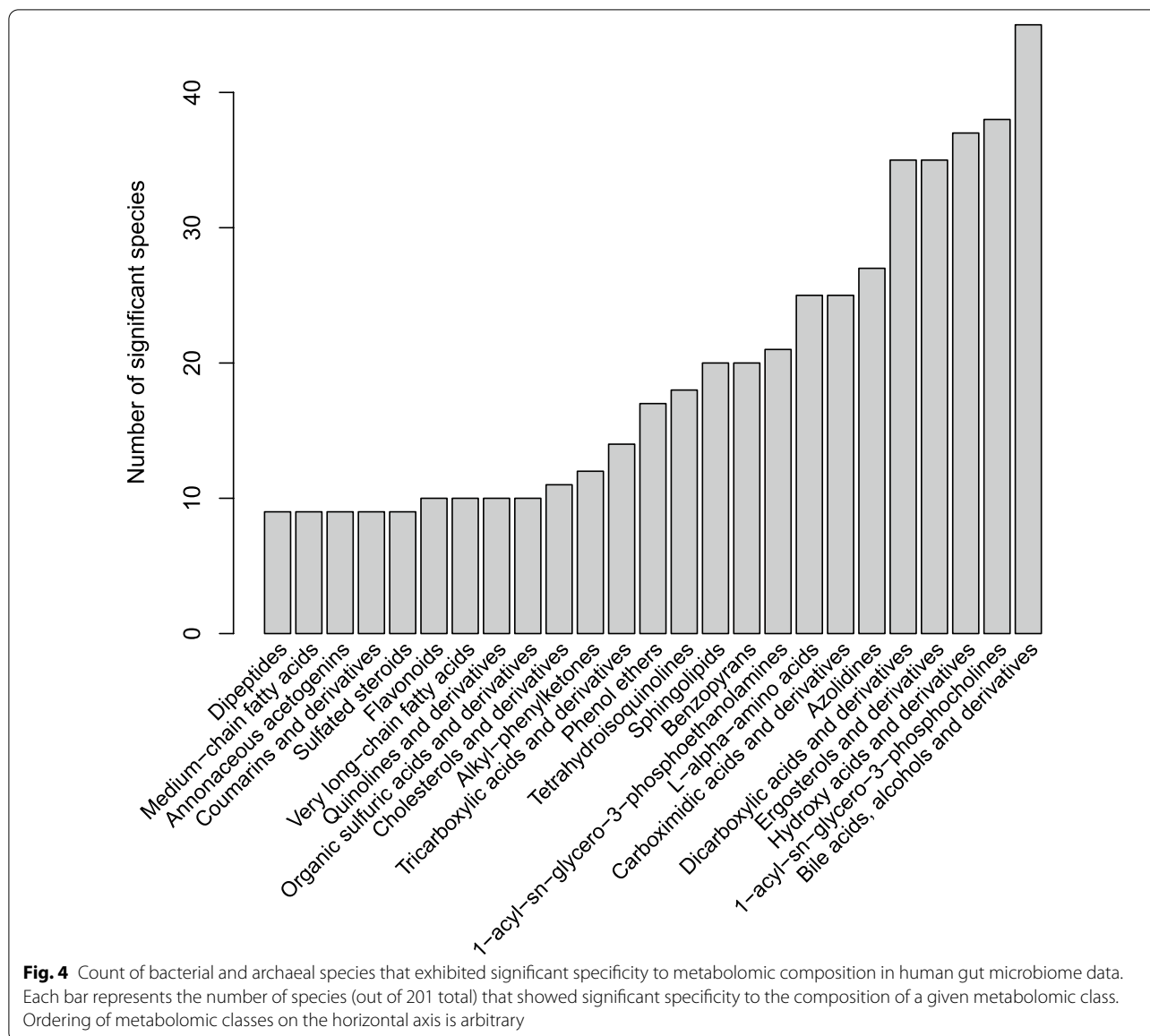
Human gut microbiome metabolomic composition specificity

In this analysis, we asked whether bacteria and archaea in the human gut microbiome had specificity to paired metabolomic data. We computed bacterial specificity to compositional dissimilarity of 83 different metabolomic classes, and of those 83, microbes showed statistically

significant specificity to only 25 (Fig. 4). The interpretation of this analysis is similar to that of the specificity to algal community composition above. Specificity to a certain composition of paired data is a more abstract concept than specificity to elevation or even to the EMP ontology, but this type of specificity makes intuitive ecological sense in the context of species-habitat association. Different microbes have different environmental needs, both within the human gut microbiome [41] and elsewhere [42]. As such, those microbes can be expected to be found in environments that meet those needs. Similarly, microbes in the human gut influence their environment [43], and as such can be expected to be found in environments that are changed by their presence. Since different microbes interact with different sorts of metabolites, differential specificity to metabolite classes is an expected outcome.

We found that more microbial species had significant specificity to the composition of bile acids, alcohols, and derivatives, compared to other metabolomic classes (Fig. 4). This result is not surprising, since bile acids strongly interact with the gut microbiome, and are also created and manipulated by it [44]. Furthermore, the experimental design for these data contained subjects with Crohn's disease, with ulcerative colitis, and healthy controls; and bile acids play a significant role in both Crohn's disease and ulcerative colitis [45]. Microbes in this analysis could be specific to either of those two conditions and their plausibly co-occurring bile acids, alcohols, and derivatives, or to subclasses thereof. Composition of this metabolomic class did not strongly correlate with composition of other metabolomic classes; its highest Mantel correlation was with Benzopyrans ($r_M = 0.46$).

Species with the strongest specificity to bile acids *etc.* were *Bacteroides plebeius*, an unclassified *Methanobrevibacter* species, *Odoribacter laneus*, *Methanosphaera stadtmanae*, and *Ruminococcus callidus* (all $Spec < -0.60$), although many other species showed significant specificity to this metabolomic class as well. *B. plebeius* was initially isolated on bile media [46], and was found to be associated with primary bile acids (as opposed to secondary bile acids) in patients with pediatric Crohn's disease [47]. *Methanobrevibacter sp.* (likely *M. smithii*) and *Methanosphaera stadtmanae* (both Archaea) are the predominant methanogens found in the human gut [48]. *M. smithii* is known to grow in the presence of bile salts [49], and may be a biomarker against inflammatory bowel disease (IBD) or for its remission [50]. The metabolomic class with the second most specific microbes was 1-acyl-sn-glycerol-3-phosphocholines (also called 2-lysolecithins). These compounds are derivatives of phosphatidylcholine,



which is used as a treatment for ulcerative colitis, but is also found naturally in some foods [51]. This finding is different than showing some bacterial species' relative abundances correspond to the amount of phosphatidylcholine derivatives; instead this analysis focuses on the composition of phosphatidylcholine derivatives; albeit with the amount of those compounds as a component since Euclidean distance was used.

In our analysis, we asked which metabolomic classes had the most microbial species specific to them (Fig. 4). The results of this type of analysis are intended to mirror common beta-diversity analyses used in microbial ecology, which ask to what extent variables explain

differences in microbiome community structure [19, 20, 52]. However, more complex questions can be asked of these data, using the results of specificity as a starting point for feature set reduction or variable selection. For example, given an individual bacterial species of interest, the variables to which it is specific may be used in a random forests model to predict its presence. For the purpose of variable selection, specificity has very low computational resource requirements when used with only the top half of Eq. 2 (using option `denom_type="sim_center"`), and can be run on personal computer hardware. This mode produces the same p values as the more comprehensive mode, and produces the same $Spec$ values for any species with $Spec < 0$. In

addition to variable selection, specificity has application in detecting species that may be common lab contaminants, as shown in our EMP analysis below.

Earth Microbiome Project (EMP) ontological specificity

As expected, the vast majority (6909/7014) of bacterial ASVs we analyzed within the EMP data set exhibited significant and strong specificity to the EMP ontology (Fig. 5). Given the distinct community-level differences in microbiomes across this ontology [26], it is not a surprise that most microbial species exhibit the same pattern. Instead, the species that buck this trend are of interest as potential cosmopolitan taxa, or possible bioinformatic failures. The sequence data obtained from the EMP data set are only 91 base-pairs in length, and even though they were clustered as exact sequence variants [14], it is possible that environmentally divergent ecotypes [53] were clustered together into the same ASV in this way.

One such highly distributed ASV was found in almost all EMP ontology categories except for saline, hypersaline, and non-saline water ($Spec = -0.13$, $P = 0.35$). It was 100% identical to multiple *Actinomadura* species, including *A. algeriensis* from Saharan soil [54], the human pathogen *A. madurae* [55], the root endophyte *A. syzygii* [56], *A. maheshkhaliensis* from mangrove rhizosphere soil [57], *A. apis* from honeybees [58], and others. These environment types span the EMP ontology, suggesting that the highly distributed ASV in question is a spurious combination of multiple *Actinomadura* ecotypes, each of

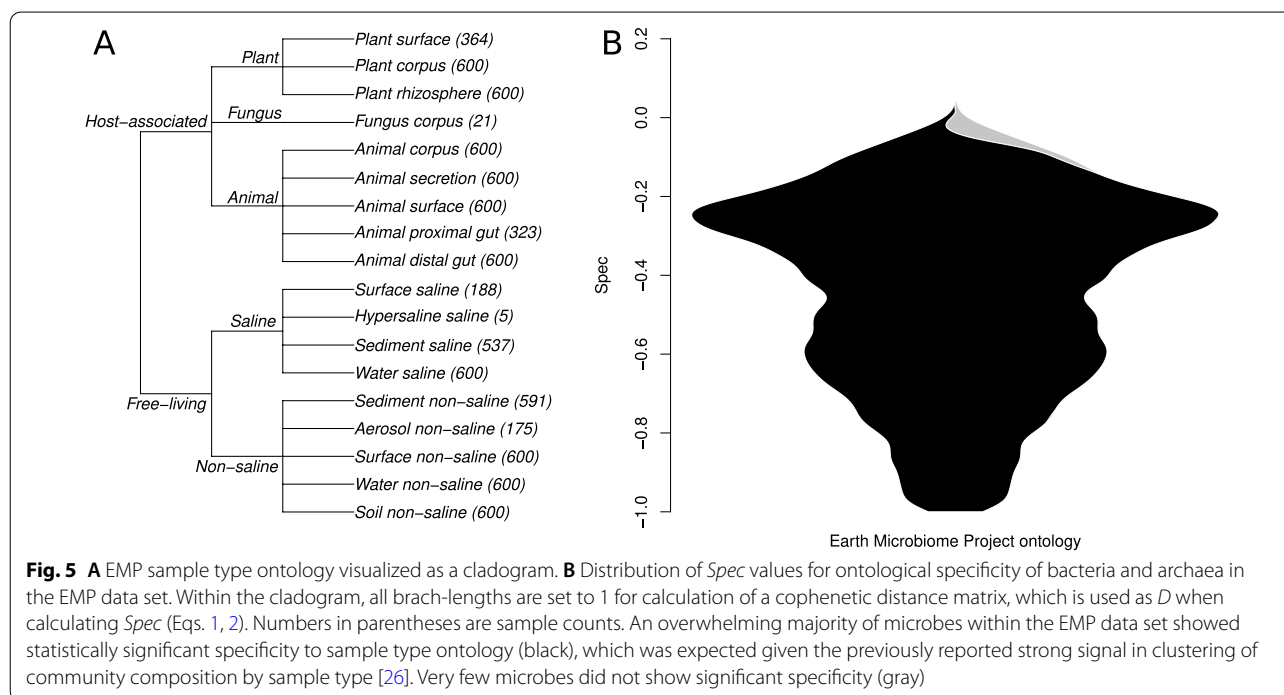
which would likely have a strong specificity signal if analyzed independently. A counterexample is a strongly specific ASV, found exclusively in animal distal guts. It was 100% identical to multiple species as well, although each was originally isolated from similar environments: *Oxobacter pfennigii* (cattle rumen [59]), *Proteiniclasticum ruminis* (yak rumen [60]), and *Lutispora thermophila* (solid waste bioreactor [61]). These examples illustrate that like with every analysis, results can only be as good as input data. Users with very short-read marker gene data are likely already aware of this limitation, so we will not belabor the point.

Interactive data visualization

In addition to using our *specificity* R package to calculate *Spec* and produce the figures shown above, we also used its companion package, *specificity.shiny*, to explore data and identify interesting features (Additional file 1: Fig. S10). With this tool, users can easily create interactive visualizations from specificity analyses, and share them over the internet. *specificity.shiny* was used in the preparation of this manuscript, to share results between authors.

Conclusion

Our R package, *specificity*, enables specificity analysis of microbiome data in the context of multiple variable types. Here, we've shown examples of specificity to geographic variables like elevation and rainfall (Fig. 2),



host phylogenetic specificity (Fig. 2), specificity to co-occurring microbial community structure (Fig. 3) and metabolomic structure (Fig. 4), and specificity to sample type ontology (Fig. 5). Our validation analyses show that our statistic, *Spec*, performs intuitively and is sensitive to specificity in both empirical and simulated data (Additional file 1: Figs. S3–S9). Our companion package, *specificity.shiny*, can be used to explore results and collaborate on specificity analyses (Additional file 1: Fig. S10), which was done by the authors on the four example analyses we presented here. Both *specificity* and *specificity.shiny* are available from the authors' GitHub repository, along with installation instructions and a tutorial vignette.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40793-022-00426-0>.

Additional file 1. All supplementary figures for this manuscript.

Acknowledgements

The authors thank J. Siebert and C. Martin for many helpful discussions and data wrangling.

Author contributions

All authors read and approved the final manuscript.

Funding

Funding was provided by an NIH NLM Computational Biology training Grant (5 T15 LM009451-12) an NSF award (1255972). Funding bodies had no role in study design, analysis, interpretation, or in the preparation of this manuscript.

Availability of data and materials

R packages *specificity* and *specificity.shiny* can be downloaded from GitHub: <https://github.com/darcyj/specificity>, <https://github.com/darcyj/specificity.shiny>. In addition to software, *specificity*'s GitHub repository contains thorough documentation, including guidance on installation, and a full tutorial vignette for using *specificity* with examples from an included data set. Code and data to replicate the analyses shown here can be found on GitHub as well: https://github.com/darcyj/specificity_analyses.

Declarations

Consent for publication

All authors have given consent for this manuscript to be published.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Division of Biomedical Informatics and Personalized Medicine, University of Colorado School of Medicine, Aurora, CO, USA. ²School of Life Sciences, University of Hawai'i at Mānoa, Honolulu, HI, USA. ³Pacific Biosciences Research Center, University of Hawai'i at Mānoa, Honolulu, HI, USA. ⁴Department of Ecology and Evolutionary Biology, University of Colorado, Boulder, Boulder, CO, USA.

Received: 7 February 2022 Accepted: 30 May 2022

Published online: 25 June 2022

References

- Combes C. Parasitism : the Ecology and Evolution of Intimate Interactions. Translated by Isauere de Buron and Vincent A. Connors, with a New Foreword by Daniel Simberloff; 2001. Interspecific interactions.
- Poulin R, Krasnov BR, Mouillot D. Host specificity in phylogenetic and geographic space; 2011. <https://doi.org/10.1016/j.pt.2011.05.003>. ISSN: 14714922. Trends in Parasitology.
- Shefferson RP, Bunch W, Cowden CC, Lee YI, Kartzinel TR, Yukawa T, Downing J, Jiang H. Does evolutionary history determine specificity in broad ecological interactions? J Ecol. 2019. <https://doi.org/10.1111/1365-2745.13170>.
- Oxford English Dictionary. Oxford English Dictionary Online; 2017. ISBN: 15424715. Oxford English Dictionary.
- Bascompte J. Mutualistic networks, vol. 7. New York: Wiley Online Library; 2009. p. 429–36.
- Dormann CF, Fruend J, Gruber B, Dormann MC, LazyData TR. Package 'barpartite'; online PDF; 2017. <https://doi.org/10.1002/sim.4177>.
- Zhang X, Yi N. NBZIMM: negative binomial and zero-inflated mixed models, with application to microbiome/metagenomics data analysis. BMC Bioinform. 2020. <https://doi.org/10.1186/s12859-020-03803-z>.
- Faith DP. Conservation evaluation and phylogenetic diversity. Biol Conserv. 1992;61(1):1–10. [https://doi.org/10.1016/0006-3207\(92\)91201-3](https://doi.org/10.1016/0006-3207(92)91201-3).
- Allen B, Kon M, Bar-Yam Y. A new phylogenetic diversity measure generalizing the Shannon index and its application to phyllostomid bats. Am Nat. 2009. <https://doi.org/10.1086/600101>.
- Costello MJ. Parasite rates of discovery, global species richness and host specificity. In: Integrative and comparative biology; 2016. <https://doi.org/10.1093/icb/icw084>. ISSN: 15577023.
- Rao CR. Diversity and dissimilarity coefficients: a unified approach. Theor Popul Biol. 1982. [https://doi.org/10.1016/0040-5809\(82\)90004-1](https://doi.org/10.1016/0040-5809(82)90004-1).
- Rao CR. Quadratic entropy and analysis of diversity. Sankhya A. 2010. <https://doi.org/10.1007/s13171-010-0016-3>.
- Botta-Dukát Z. Rao's quadratic entropy as a measure of functional diversity based on multiple traits. J Veg Sci. 2005. <https://doi.org/10.1111/j.1654-1103.2005.tb02393.x>.
- Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: high-resolution sample inference from Illumina amplicon data. Nat Methods. 2016;13:581–3. <https://doi.org/10.1038/nmeth.3869>.
- Botta-Dukát Z. Cautionary note on calculating standardized effect size (SES) in randomization test. Community Ecol. 2018. <https://doi.org/10.1556/168.2018.19.1.8>.
- Baluja S, Caruana R. Removing the genetics from the standard genetic algorithm. In: Machine learning proceedings 1995. Amsterdam: Elsevier; 1995. p. 38–46.
- Goldberg DE, Lingle R. Alleles, loci, and the traveling salesman problem. In: Proceedings of an international conference on genetic algorithms and their applications, vol. 154. Carnegie-Mellon University Pittsburgh; 1985. p. 154–159.
- Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. J R Stat Soc Ser B (Methodol). 1995;57(1):289–300.
- Darcy JL, Swift SOI, Cobian GM, Zahn GL, Perry BA, Amend AS. Fungal communities living within leaves of native Hawaiian dicots are structured by landscape-scale variables as well as by host plants. Mol Ecol. 2020. <https://doi.org/10.1111/mec.15544>.
- Sommers P, Darcy JL, Porazinska DL, Gendron EMS, Fountain AG, Zamora F, Vincent K, Cawley KM, Solon AJ, Vimercati L, Ryder J, Schmidt SK. Comparison of microbial communities in the sediments and water columns of frozen cryoconite holes in the mcmurdo dry valleys. Antarct Front Microbiol. 2019;10:65. <https://doi.org/10.3389/fmicb.2019.00065>.
- Sommers P, Porazinska DL, Darcy JL, Zamora F, Fountain AG, Schmidt SK. Experimental cryoconite holes as mesocosms for studying community ecology. Polar Biol. 2019;42(11):1973–84.
- Dunn OJ. Multiple comparisons using rank sums. Technometrics. 1964;6(3):241–52.
- Franzosa EA, Sirota-Madi A, Avila-Pacheco J, Fornoel N, Haiser HJ, Reinker S, Vatanen T, Hall AB, Mallick H, McIver LJ. Gut microbiome structure and metabolic activity in inflammatory bowel disease. Nat Microbiol. 2019;4(2):293–305.

24. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C, Segata N. MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nat Methods*. 2015;12(10):902–3.
25. Wishart DS, Feunang YD, Marcu A, Guo AC, Liang K, Vázquez-Fresno R, Sajed T, Johnson D, Li C, Karu N. HMDB 4.0: the human metabolome database for 2018. *Nucleic Acids Res*. 2018;46(D1):608–17.
26. Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A, Gibbons SM, Ackermann G. A communal catalogue reveals Earth's multiscale microbial diversity. *Nature*. 2017;551(7681):457–63.
27. Gonzalez A, Navas-Molina JA, Kosciulek T, McDonald D, Vázquez-Baeza Y, Ackermann G, DeReus J, Janssen S, Swafford AD, Orchanian SB. Qiita: rapid, web-enabled microbiome meta-analysis. *Nat Methods*. 2018;15(10):796–8.
28. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*. 1990;215(3):403–10. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
29. Eddelbuettel D, François R. Rcpp: seamless R and C++ integration. *J Stat Softw*. 2011;40(8):1–18. <https://doi.org/10.18637/jss.v040.i08>.
30. Wickham H. *R Packages*. 1st ed. Newton: O'Reilly Media Inc.; 2015.
31. Wickham H. testthat: get started with testing. *The R J*. 2011;3:5–10.
32. Chang W, Cheng J, Allaire J, Sievert R, Schloerke B, Xie Y, Allen J, McPherson J, Dipert A, Borges B. Shiny: web application framework for R; 2021. R package version 1.7.1. <https://CRAN.R-project.org/package=shiny>.
33. Hester J, Csardi G, Wickham H, Chang W, Morgan M, Tenenbaum D. remotes: R Package Installation from Remote Repositories, Including 'GitHub'; 2020.
34. Tipton L, Zahn GL, Darcy JL, Amend AS, Hynson NA. Hawaiian fungal amplicon sequence variants reveal otherwise hidden biogeography. *Microb Ecol*. 2021;83:48–57.
35. Giambelluca TW, Chen Q, Frazier AG, Price JP, Chen Y-L, Chu P-S, Eischeid JK, Delporte DM. Online rainfall atlas of Hawaii. *Bull Am Meteorol Soc*. 2013;94(3):313–6.
36. Dean R, Van Kan JA, Pretorius ZA, Hammond-Kosack KE, Di Pietro A, Spanu PD, Rudd JJ, Dickman M, Kahmann R, Ellis J. The Top 10 fungal pathogens in molecular plant pathology. *Mol Plant Pathol*. 2012;13(4):414–30.
37. Gan P, Ikeda K, Irieda H, Narusaka M, O'Connell RJ, Narusaka Y, Takano Y, Kubo Y, Shirasu K. Comparative genomic and transcriptomic analyses reveal the hemibiotrophic stage shift of Colletotrichum fungi. *New Phytol*. 2013;197(4):1236–49. <https://doi.org/10.1111/nph.12085>.
38. Gan P, Narusaka M, Kumakura N, Tsushima A, Takano Y, Narusaka Y, Shirasu K. Genus-wide comparative genome analyses of Colletotrichum species reveal specific gene family losses and gains during adaptation to specific infection lifestyles. *Genome Biol Evol*. 2016;8(5):1467–81. <https://doi.org/10.1093/gbe/evw089>.
39. Porazinska DL, Fountain AG, Nylen TH, Tranter M, Virginia RA, Wall DH. The biodiversity and biogeochemistry of cryoconite holes from McMurdo Dry Valley Glaciers, Antarctica. *Arct Antarct Alp Res*. 2004;36(1):84–91. [https://doi.org/10.1657/1523-0430\(2004\)036\[0084:TBABOC\]2.0.CO;2](https://doi.org/10.1657/1523-0430(2004)036[0084:TBABOC]2.0.CO;2).
40. Darcy JL, King AJ, Gendron EMS, Schmidt SK. Spatial autocorrelation of microbial communities atop a debris-covered glacier is evidence of a supraglacial chronosequence. *FEMS Microbiol Ecol*. 2017. <https://doi.org/10.1093/femsec/fix095>.
41. Costello EK, Stagaman K, Dethlefsen L, Bohannan BJ, Relman DA. The application of ecological theory toward an understanding of the human microbiome. *Science*. 2012;336(6086):1255–62.
42. Nemergut DR, Schmidt SK, Fukami T, O'Neill SP, Bilinski TM, Stanish LF, Knelman JE, Darcy JL, Lynch RC, Wickey P. Patterns and processes of microbial community assembly. *Microbiol Mol Biol Rev*. 2013;77(3):342–56.
43. Kelly CJ, Zheng L, Campbell EL, Saeedi B, Scholz CC, Bayless AJ, Wilson KE, Glover LE, Kominsky DJ, Magnuson A, Weir TL, Ehrentraut SF, Pickel C, Kuhn KA, Lanis JM, Nguyen V, Taylor CT, Colgan SP. Crosstalk between microbiota-derived short-chain fatty acids and intestinal epithelial HIF augments tissue barrier function. *Cell Host Microbe*. 2015;17(5):662–71. <https://doi.org/10.1016/j.chom.2015.03.005>.
44. Bile Acids and the Gut Microbiome **30**. <https://doi.org/10.1097/MOG.000000000000057>.
45. Tiraterra E, Franco P, Porru E, Katsanos KH, Christodoulou DK, Roda G. Role of bile acids in inflammatory bowel disease. *Ann Gastroenterol*. 2018;31(3):266–72. <https://doi.org/10.20524/aog.2018.0239>.
46. Kitahara M, Sakamoto M, Ike M, Sakata S, Benno Y. *Bacteroides plebeius* sp. nov. and *Bacteroides coprocola* sp. nov., isolated from human faeces. *Int J Syst Evol Microbiol*. 2005;55(5):2143–7. <https://doi.org/10.1099/ijso.63788-0>.
47. Connors J, Dunn KA, Allott J, Bandsma R, Rashid M, Otley AR, Bielawski JP, Van Limbergen J. The relationship between fecal bile acids and microbiome community structure in pediatric Crohn's disease. *ISME J*. 2020;14(3):702–13. <https://doi.org/10.1038/s41396-019-0560-3>.
48. Brown JR-M, Flemer B, Joyce SA, Zulquernain A, Sheehan D, Shanahan F, O'Toole PW. Changes in microbiota composition, bile and fatty acid metabolism, in successful faecal microbiota transplantation for Clostridioides difficile infection. *BMC Gastroenterol* 2018;18(1):131. <https://doi.org/10.1186/s12876-018-0860-5>
49. Miller TL, Wolin MJ, de Macario EC, Macario AJ. Isolation of Methanobrevibacter smithii from human feces. *Appl Environ Microbiol*. 1982;43(1):227–32. <https://doi.org/10.1128/aem.43.1.227-232.1982>.
50. Ghavami SB, Rostami E, Sephay AA, Shahrokh S, Balaii H, Aghdaii HA, Zali MR. Alterations of the human gut Methanobrevibacter smithii as a biomarker for inflammatory bowel diseases. *Microbial Pathog*. 2018;117:285–9. <https://doi.org/10.1016/j.micpath.2018.01.029>.
51. Kokkinidis DG, Bosdelekidou EE, Iliopoulou SM, Tassos AG, Texakalidis PT, Economopoulos KP, Kousoulis AA. Emerging treatments for ulcerative colitis: a systematic review. *Scand J Gastroenterol*. 2017;52(9):923–31.
52. Shaw L, Ribeiro ALR, Levine AP, Pontikos N, Balloux F, Segal AW, Roberts AP, Smith AM. The human salivary microbiome is shaped by shared environment rather than genetics: evidence from a large family of closely related individuals. *mBio*. 2017;8(5):01237–17. <https://doi.org/10.1128/mBio.01237-17>.
53. Cohan FM. What are bacterial species? *Annu Rev Microbiol*. 2002;56(1):457–87.
54. Lahoum A, Bouras N, Mathieu F, Schumann P, Spröer C, Klenk H-P, Sabaou N. *Actinomadura algeriensis* sp. nov., an actinobacterium isolated from Saharan soil. *Antonie Van Leeuwenhoek*. 2016;109(1):159–65. <https://doi.org/10.1007/s10482-015-0617-x>.
55. Izri A, Aljundi M, Billard-Pomares T, Fofana Y, Marteau A, Ferreira TG, Brun S, Caux F, Akhondji M. Molecular identification of *Actinomadura madurae* isolated from a patient originally from Algeria; observations from a case report. *BMC Infect Dis*. 2020;20(1):829. <https://doi.org/10.1186/s12879-020-05552-z>.
56. Rachniyom H, Matsumoto A, Indananda C, Duangmal K, Takahashi Y, Thamchaipenat A. *Actinomadura syzygii* sp. nov., an endophytic actinomycete isolated from the roots of a jambolan plum tree (*Syzygium cumini* L. Skeels). *Int J Syst Evol Microbiol*. 2015;65(Pt 6):1946–9. <https://doi.org/10.1099/ijso.0.000203>.
57. Ara I, Matsumoto A, Abdal Bakir M, Kudo T, Omura S, Takahashi Y. *Actinomadura maheshkhaliensis* sp. nov., a novel actinomycete isolated from mangrove rhizosphere soil of Maheshkhali, Bangladesh. *J Gen Appl Microbiol*. 2008;54(6):335–42. <https://doi.org/10.2323/jgam.54.335>.
58. Promnuan Y, Kudo T, Ohkuma M, Chantawannakul P. *Actinomadura apis* sp. nov., isolated from a honey bee (*Apis mellifera*) hive, and the reclassification of *Actinomadura cremeta* subsp. *rifamycinii* Gauze et al. 1987 as *Actinomadura rifamycinii* (Gauze et al. 1987) sp. nov., comb. nov. *Int J Syst Evol Microbiol*. 2011;61(Pt 9):2271–7. <https://doi.org/10.1099/ijso.0.026633-0>.
59. Krumholz LR, Bryant MP. *Clostridium pfnennigii* sp. nov. uses methoxyl groups of monobenzenoids and produces butyrate. *Int J Syst Evol Microbiol*. 1985;35(4):454–6. <https://doi.org/10.1099/00207173-35-4-454>.
60. Zhang K, Song L, Dong X. *Proteiniclasticum ruminis* gen. nov., sp. nov., a strictly anaerobic proteolytic bacterium isolated from yak rumen. *Int J Syst Evol Microbiol*. 2010;60(9):2221–5. <https://doi.org/10.1099/ijso.0.011759-0>.
61. Shiratori H, Ohiwa H, Ikeno H, Ayame S, Kataoka N, Miya A, Beppu T, Ueda K. *Lutispora thermophila* gen. nov., sp. nov., a thermophilic, spore-forming bacterium isolated from a thermophilic methanogenic bioreactor digesting municipal solid wastes. *Int J Syst Evol Microbiol*. 2008;58(4):964–9. <https://doi.org/10.1099/ijso.0.65490-0>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.