

## RESOURCE ARTICLE

# SCNIC: Sparse correlation network investigation for compositional data

Michael Shaffer<sup>1</sup> | Kumar Thurimella<sup>1,3</sup>  | John D. Sterrett<sup>2</sup> | Catherine A. Lozupone<sup>1</sup>

<sup>1</sup>Department of Biomedical Informatics, University of Colorado Anschutz Medical Campus, Aurora, Colorado, USA

<sup>2</sup>Department of Integrative Physiology, University of Colorado, Boulder, Colorado, USA

<sup>3</sup>Department of Chemical Engineering and Biotechnology, University of Cambridge, Cambridge, UK

**Correspondence**

Catherine A. Lozupone, Department of Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO, USA.

Email: [catherine.lozupone@cuanschutz.edu](mailto:catherine.lozupone@cuanschutz.edu)

**Funding information**

U.S. National Library of Medicine, Grant/Award Number: 4 T15 LM009451-10; National Science Foundation

**Handling Editor:** Alex J. Dumbrell

**Abstract**

Microbiome studies are often limited by a lack of statistical power due to small sample sizes and a large number of features. This problem is exacerbated in correlative studies of multi-omic datasets. Statistical power can be increased by finding and summarizing modules of correlated observations, which is one dimensionality reduction method. Additionally, modules provide biological insight as correlated groups of microbes can have relationships among themselves. To address these challenges, we developed SCNIC: Sparse Cooccurrence Network Investigation for compositional data. SCNIC is open-source software that can generate correlation networks and detect and summarize modules of highly correlated features. Modules can be formed using either the Louvain Modularity Maximization (LMM) algorithm or a Shared Minimum Distance algorithm (SMD) that we newly describe here and relate to LMM using simulated data. We applied SCNIC to two published datasets and we achieved increased statistical power and identified microbes that not only differed across groups, but also correlated strongly with each other, suggesting shared environmental drivers or cooperative relationships among them. SCNIC provides an easy way to generate correlation networks, identify modules of correlated features and summarize them for downstream statistical analysis. Although SCNIC was designed considering properties of microbiome data, such as compositionality and sparsity, it can be applied to a variety of data types including metabolomics data and used to integrate multiple data types. SCNIC allows for the identification of functional microbial relationships at scale while increasing statistical power through feature reduction.

**KEYWORDS**

bioinformatics/phyloinformatics, microbial ecology, network analysis, species interactions

## 1 | BACKGROUND

Microbial communities play important roles in environmental and human health systems and can often reach great complexity. In

these rich ecosystems, microbes interact with each other, forming relationships based on predator–prey dynamics (Corno et al., 2013), competition for resources (Burkepille et al., 2006), cross-feeding of small compounds, (LaSarre et al., 2017) and other factors. Identifying

Michael Shaffer and Kumar Thurimella contributed equally to this work.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2022 The Authors. *Molecular Ecology Resources* published by John Wiley & Sons Ltd.

correlated pairs of microbes can suggest potential interactions or shared environmental preferences. Accordingly, studies have identified complex networks of co-occurring microbes in a variety of different environments ranging from the human mouth and gut (Faust et al., 2012) to soil (Barberán et al., 2012) and stream ecosystems (Widder et al., 2014).

To detect correlations between microbes, a variety of methods have been developed. While traditional correlation metrics are used by some (Bray & Curtis, 1957; Pearson, 1909; Spearman, 1904), newer methods have been developed that take into account the properties of 16S rRNA sequencing data (Haas et al., 2011; Huse et al., 2010; Kuczynski et al., 2010). A recent review tested these methods on a variety of models and identified some methods that performed better than others in ways that can depend on underlying data characteristics (Weiss et al., 2016). Although these tools are useful for finding pairwise relationships between organisms, less attention has been given towards developing methods for finding correlations among groups of microbes.

One way to explore complex interactions is to form networks in which correlated organisms are joined with an edge, and highly correlated sets of microbes are defined. Here, we refer to these sets as modules, which are synonymous to clusters or groups. There are two primary benefits of finding modules of correlated microbes. First, the combination of microbes in a module could be further explored to understand microbial interactions, such as cross-feeding relationships, or shared environmental niches (Ban et al., 2015; Barberán et al., 2012; Dugas et al., 2018; Lozupone et al., 2012). Second, considering correlation structure among microbes can aid in statistical analysis aimed at uncovering relationships between microbes and other environmental factors. Specifically, by eliminating or summarizing highly correlated features, dependence between features is decreased. Feature reduction will increase accuracy of methods that assume the independence of features such as false discovery rate technique (FDR) measurements like the Benjamini-Hochberg Correction (Benjamini & Hochberg, 2000), and statistical power is increased by reducing the number of feature comparisons.

One workflow for considering groups of correlated microbes in downstream statistical analyses requires three steps: first, correlations between microbes must be measured and used to form a network; second, modules must be identified; and third, abundance of the microbes in modules must be summarized for use in subsequent statistical analyses. One software tool that has implemented this workflow, developed for application to gene expression data, is weighted gene correlation network analysis (WGCNA) (Langfelder & Horvath, 2008). WGCNA builds correlation networks based on a correlation coefficient (such as Pearson, Spearman, or biweight midcorrelation (Wilcox, 2011)), and detects modules as subtrees in a hierarchical cluster of features (Barabási & Albert, 1999). Modules are summarized by setting module abundance to that of network hubs or an eigenvector of the abundance of all module members (Langfelder & Horvath, 2008).

Several groups have used WGCNA to find correlations within 16S rRNA sequencing data (Castillo et al., 2017; Tong et al., 2014;

Yin et al., 2017; Younge et al., 2017), but this approach may not be appropriate for several reasons (Jackson et al., 2018). First, the correlation metrics implemented in WGCNA do not account for sparsity and compositionality. Most sequencing-based microbiome datasets are sparse (i.e. there are many zeros) and compositional, meaning they only carry information on relative abundances of taxa instead of absolute abundances, which can lead to the detection of spurious correlations if proper statistical methods are not used (Gloor et al., 2017). Thus, the use of WGCNA for compositional data may be leading to the detection of spurious edges in microbiome networks. Second, the primary method WGCNA uses to pick modules assumes the correlation network will have a scale-free topology that may not be relevant to microbiome data (Broido & Clauset, 2019). Third, summarizing modules through identifying hub taxa works well in gene expression where a single transcription factor can control the expression of many genes, but may not be appropriate in microbial communities. Both the hub and eigenvector approaches to module summarization do not allow for output tables that maintain the total counts of microbial abundance per sample. Therefore, the hub and eigenvector approaches cannot be used with tools developed for microbiome data analysis that make assumptions based on total sample counts, such as ANCOM (Mandal et al., 2015) or metagenomeSeq (Paulson et al., 2013).

Optimal methods for identifying and summarizing modules of correlated features in 16S rRNA sequencing data have not been deeply explored. One study (Jackson et al., 2018) recommended an ensemble approach for correlation detection, and the Louvain modularity maximization (LMM) method (Blondel et al., 2008) to identify modules. LULU is a tool that follows a binning approach towards OTUs that co-occur, but only does so if they are highly phylogenetically related (Frøslev et al., 2017). Another tool, CoNet, uses an ensemble approach to build and visualize networks (Faust & Raes, 2016). However, no implementation of module summarization was made available for downstream statistical analysis.

To address these gaps, we have developed a tool for sparse, compositional correlation network investigation for compositional data (SCNIC), which uses methods optimized for microbiome data analysis. SCNIC is available as standalone Python software, via Bioconda (Grüning et al., 2018) and the package installer for Python (pip), and as a QIIME 2 plugin (Bolyen et al., 2019). The source code for SCNIC and the QIIME 2 plugin is freely available on GitHub (<https://github.com/lozuponelab/SCNIC>, <https://github.com/lozuponelab/q2-SCNIC>) under the BSD-3-Clause Licence.

## 2 | MATERIALS AND METHODS

### 2.1 | The SCNIC method

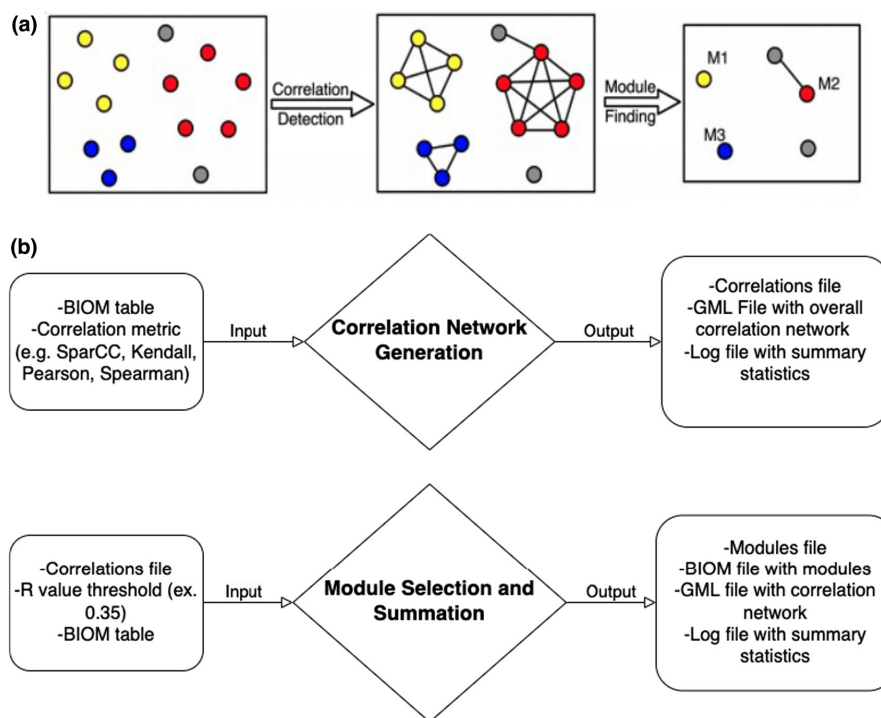
SCNIC takes a feature table containing counts of each feature in all samples as input and performs three steps: (1) a correlation network is built, (2) modules are detected in the network and (3) feature counts within a module are summed into a new single feature

(identified as "module-x" where x is whole numbered consecutively starting at zero; Figure 1). The modules are ordered based on size, where the lower numbered modules have a larger number of members compared to higher numbered modules. To summarize modules, SCNIC uses a sum of count data from all features in a module. There is no maximum or minimum size constraint on module size when modules are created. The newly generated modules are included in a new feature table alongside all features not grouped into a module. This maintains the total counts per sample, allowing for downstream analyses with tools that have assumptions related to total sample counts. SCNIC produces a graph modelling language (GML) format (Himsolt, 1997) file compatible with Cytoscape (Shannon et al., 2003) for network visualization in which the edges in the correlation network represent the positive correlations which are stronger than a user specified *R*-value cutoff (between 0 and 1), a file describing which features compose each defined module, and a feature table in the Biological Observation Matrix (BIOM) (McDonald et al., 2012; Figure 1).

SCNIC allows users to choose between multiple methods for detecting correlations and of defining modules of co-occurring microbes. For correlations, SCNIC can implement traditional correlation metrics (including Pearson's *r*, Spearman's  $\rho$  and Kendall's  $\tau$ ) or the compositionality- and sparsity-aware correlation metric from

SparCC (Friedman & Alm, 2012; Watts et al., 2019) to correct for aspects of microbiome data. SparCC has been shown to perform well in detecting correlations compared to other correlation measures (Weiss et al., 2016). Specifically, SparCC performs well in communities with an inverse Simpson index above 13 (which would be indicative of a high number of successful species, a complex food web, and many ecological niches, as would be seen in many high biomass microbial communities such as gut or soil microbiomes) (Fernandes et al., 2014; Watts et al., 2019), and it thus was chosen as the default metric.

To define modules of co-correlated features, we implement two methods: (1) Louvain modularity maximization (LMM) and (2) a novel shared minimum distance (SMD) module detection algorithm; unlike WGCNA, neither of these algorithms make assumptions about network topology. LMM was previously proposed as a method for clustering correlation networks of microbes into modules (Blondel et al., 2008). LMM works by first assigning one module per feature. Each pair of adjacent modules are joined and the change in modularity (defined by the number of edges within the module compared to outside) is calculated for each module. The pair which increases the mean modularity of the network the most is then joined. This process is repeated until the modularity of the network is not increased. LMM uses two parameters provided by the user: The first



**FIGURE 1** SCNIC schematic and data flow. (a) The basic process of SCNIC involves first identifying pairwise correlations between species and using them to build a correlation network. Modules of correlated features are identified and then summarized for downstream statistical analysis, or multi-omic analysis between modules of microbes and other feature types. (b) The input to SCNIC comes in the form of a count table in BIOM format. The first step takes the table and generates a correlation table and network. The table is in a tab delimited format and the network is in GML format and can be used to visualize the network in Cytoscape. Modules are detected and summarized in the final step which generates a module membership file indicating which features are in each module. The collapsed BIOM table contains the same total counts per sample as the original table, but with less features. All features not included in modules are retained with their original counts and all modules have a total count per sample of the sum of the counts of all features in that module.

parameter,  $R$ -value, defines the minimum correlation coefficient for defining an edge between features. The second parameter, gamma (also referred to as resolution), controls the size of modules detected, with large gamma values yielding larger modules.

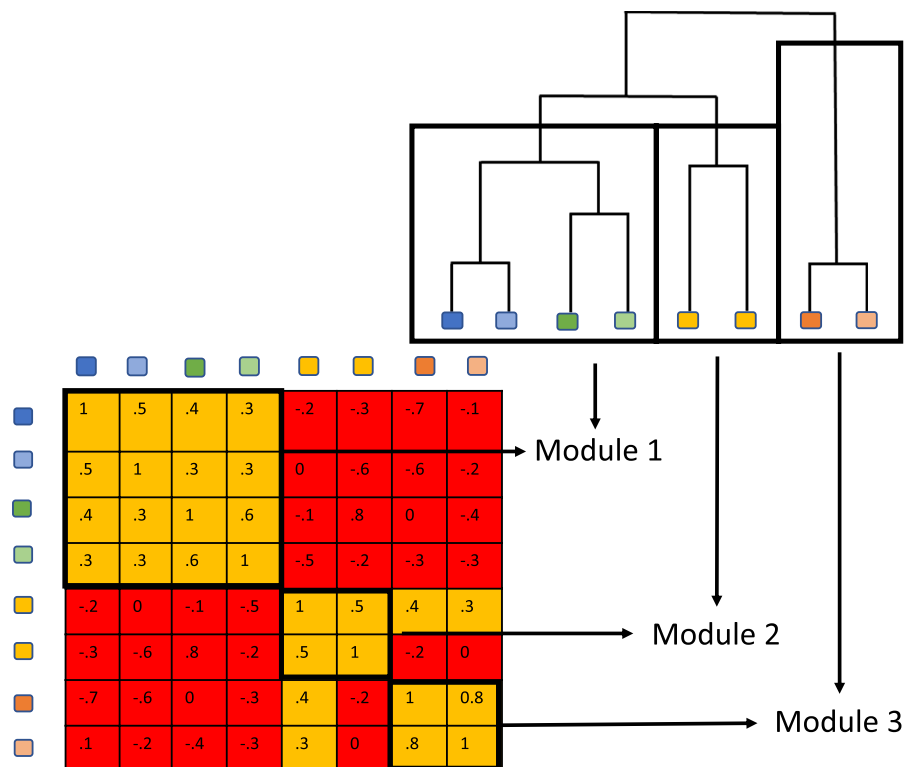
WGCNA and LMM have a potential weakness in that modules can contain pairs of taxa that are not strongly correlated (e.g., if they are several steps away from each other in the network). To address this weakness, we also implement the SMD method to ensure that correlations between all pairs of features in the module have an  $R$ -value greater than the user provided minimum (Figure 2). Specifically, the SMD method defines modules by first applying complete linkage hierarchical clustering to correlation coefficients to make a tree of features. Next, SMD defines modules as subtrees where correlations between all pairs of tips have an  $R$ -value above the specified value. SMD has been set as the default method in SCNIC because of the desirable property of only producing modules where all features are correlated over a user-specified threshold.

A large proportion of microbiome studies sample highly uneven communities which leads to strong compositionality-driven artefacts (Fernandes et al., 2014; Gloor et al., 2017; Tsilimigras & Fodor, 2016). Because of this, we use SparCC, specifically the implementation of FastSpar (Watts et al., 2019), as the default correlation measure. SparCC was used as the correlation metric based on analysis that suggested a high precision in the number of correct edges recovered when correlations were calculated in synthetic data (Weiss et al., 2016). SCNIC additionally includes the option of using Pearson's  $r$ , Spearman's  $\rho$  and Kendall's  $\tau$  to evaluate non-compositional or dense data types.

## 2.2 | Evaluating the SMD algorithm using simulated data

Since SMD has not been applied to microbiome module detection before, we compared SMD to LMM using simulated data. In order to evaluate the performance of SMD for module detection under different parameter settings and compare it to LMM, we simulated a wide range of networks. The simulations had networks with similar characteristics to those seen in networks generated from microbiome datasets. These included networks with power law degree distributions ( $N = 175$ ) with values of  $a$ , the exponent term of the power law formula ( $y = kX^{-a}$ ), varying between 1.8 and 2.6, as well as networks with regular degree distributions ( $N = 200$ ) with  $p$ , the probability of one node being connected to another, varying from 0.001 to 0.2. The power law and regular degree distribution networks were created using the NetworkX v2.6.3 implementations of configuration\_model and erdos\_renyi\_graph, respectively, and all had a size of 500. The networks with power law degree distributions had modularity values between 0.2 and 0.9, with higher  $a$  corresponding to higher modularity, and the networks with regular degree distributions had modularity values between 0.07 and 0.98, with lower  $p$  corresponding to higher modularity. Higher modularity scores indicate many connections within modules and fewer connections between modules. We then calculated SMD and LMM partitions (with LMM gamma = 1) of each network and compared the homogeneity between the two partitions. Because SMD modules are smaller than LMM modules, we used the homogeneity metric described by Rosenberg and Hirschberg (Rosenberg

**FIGURE 2** SMD algorithm for defining modules. SMD defines a module as a group where all features have a correlation above a given threshold. To do this SMD first uses complete linkage hierarchical clustering on correlation coefficients to create a tree. Each module is defined as a subtree where the correlation coefficients between all tips are greater than the threshold.



& Hirschberg, 2007) (implemented via Scikit-learn v0.24.2) to assess whether nodes partitioned together by SMD are a subset of the module partitioned by LMM. A score of 1 represents that all nodes in SMD modules represent sub-modules of LMM-partitioned modules, whereas a score of 0 represents that no two nodes that were classified by SMD into the same module were partitioned into a module together by the LMM method.

## 2.3 | Demonstrating the use of SCNIC

We demonstrate the use of SCNIC with two example datasets. These are (1) a study that used 16S rRNA sequencing of faecal material to compare microbiome composition in individuals with and without HIV and in men who have sex with men (MSM) who were at a high risk of contracting HIV (Noguera-Julian et al., 2016), and (2) a dataset analysing the microbiome of water samples at various depths in two of the Great Lakes. We chose these two datasets so that we could evaluate performance using datasets from both host-associated and free-living microbiomes. We also used the Great Lakes dataset to compare module size and modularity between SMD and LMM selected modules.

## 2.4 | HIV dataset

The HIV data set was retrieved from NCBI SRA accession number SRP068240, and samples from the BCNO cohort were used for these analyses. Reads were error corrected, quality trimmed, and primers were removed using default parameters in BBTtools (Bushnell et al., 2017). DADA2 (Callahan et al., 2016) was used to define amplicon sequence variants (ASVs) with reads trimmed from the left by 30 base pairs and truncated at 269. ASVs were binned into operational taxonomic units (OTUs) using USEARCH (Edgar, 2010) at 99% identity using QIIME 1 (Caporaso et al., 2010). A phylogenetic tree was made using a single representative sequence from each OTU and the SEPP protocol (Janssen et al., 2018; Mirarab et al., 2012) using QIIME 2 (Bolyen et al., 2019). We evaluated the average phylogenetic distance between OTUs in the same module using the *distance* method of Biopython (Cock et al., 2009; Talevich et al., 2012). Taxonomy was assigned using the Naive Bayes QIIME 2 feature classifier, version gg-13-8-99-515-806-nb-classifier.qza.

The original study describing these data showed a strong divergence in gut microbiome composition in MSM compared to non-MSM independent of HIV infection status and more subtle differences associated with HIV infection when controlling for MSM behaviour. The goal of our analysis was to evaluate whether comparing gut microbiome composition between HIV negative MSM and non-MSM with SCNIC modules provide additional significant taxa compared to without, and additional insights as to which taxa that differ with MSM also are in turn demonstrating co-correlated structure with each other. Co-correlation of microbes may indicate that they are a part of a broader community type, interact with each

other, or have shared environmental drivers of their prevalence. A further goal of this analysis is to examine the effects of using different *R*-value thresholds on the results. The SMD method was specifically used with SparCC *R*-value thresholds between 0.20 and 1.0, with 0.05 increments.

### 2.4.1 | Great Lakes dataset

The Great Lakes dataset was previously published as part of the Earth Microbiome Project (Thompson et al., 2017). This study evaluated patterns of microbial relative abundance across depths in Lake Michigan ( $N = 16$ ) and Lake Superior ( $N = 33$ ), with depth of samples collected ranging from 5 to 3654 meters. The study additionally recorded data on pH and salinity. The Great Lakes data set was retrieved from QIITA accession number 1041 (Gonzalez et al., 2018). ASVs were found using DADA2 with a left trim of 30 and a truncation length of 135. OTUs were subsequently picked on the ASVs using VSEARCH (Rognes et al., 2016) with a 99% identity threshold, resulting in 3871 OTUs. These steps were done with QIIME 2 (Bolyen et al., 2019). SCNIC was applied with the SMD method and 0.2, 0.4 and 0.65 *R*-value thresholds.

### 2.4.2 | Comparison of SMD to LMM using the Great Lakes dataset

To identify differences in module structure from SMD versus LMM partitions, we assessed the module size and modularity of 221 separately partitioned networks from the Great Lakes dataset using varying parameters for SCNIC. The parameters included SCNIC *R* thresholds ranging from 0.1 to 0.7 and gamma ranging from 0.15 to 0.9 for LMM.

## 2.5 | Evaluating effects of applying SCNIC to discern microbes that differ between groups in the HIV and Great Lakes datasets

OTUs/modules that differed with MSM status (HIV study) and between lakes (Great Lakes study) were identified using ANCOM (Mandal et al., 2015) for each feature. For the first study, we focused on evaluating differences in the microbiome between MSM and non-MSM without confounding by HIV infection status, by only using samples from HIV negative individuals. We chose ANCOM because it is also a tool designed specifically for working with compositional microbiome data. ANCOM was applied to the original feature table where SCNIC was not applied, as well as to feature tables output from SCNIC using SparCC at different *R*-value thresholds with the SMD algorithm.

While applying SparCC, SCNIC uses the recommended practice of the SparCC manuscript of filtering based on average relative abundance across samples (Friedman & Alm, 2012). The SparCC

manuscript suggests this filter because removing features with high abundances, even in a few samples, will upset the ability of the method to control for the number of reads per sample in its compositionality adjustment. Because this method can retain OTUs that are highly abundant in only a single sample, we removed features that had 0 values in more than 5% (~29/146) of samples before applying ANCOM but after applying SparCC. Significant differences between groups were determined as those above the *W*-value threshold determined by ANCOM.

## 2.6 | Time and memory usage by SCNIC

To evaluate the memory resources needed by SCNIC, we ran the SCNIC modules step locally on a 2015 MacBook Pro with 16GB RAM with a 2.5 GHz Quad-Core Intel Core i7 processor for both the Great Lakes dataset and an integrated microbiome-metabolome dataset with 1301 features, which will be referred to as the GT dataset. The runtime was recorded across 3 runs per method (SMD vs LMM) for each dataset using GNU Time, and memory was profiled using memory-profiler 0.60.0. The “within” step, which calculates correlations between features and creates the correlation network was not tested because it depends greatly on the correlation metric used, and the runtime and memory usage of FastSpar (likely the most computationally intensive correlation metric to be used in this step) have already been profiled (Watts et al., 2019). The modules step only utilizes a correlation matrix and as such does not scale with the number of samples, only the number of features, except when the values of the count table are being summed, which is a generally trivial calculation compared to the module generation step.

## 3 | RESULTS

### 3.1 | Comparison of LMM to SMD on real and simulated data

In order to evaluate the relationship between modules detected with SMD versus LMM, we chose modules on the Great Lakes dataset using SMD at *R*-value thresholds ranging from 0.05 to 0.7 and with LMM at the same *R*-value thresholds and gamma values ranging from 0.15 to 0.9. We found that with both LMM and SMD, modularity increased with increasing *R*-value thresholds. However, SMD produced less modular partitions and smaller modules than LMM, even when LMM was applied with very low values for the gamma parameter that controls module size (Figure 3a,b).

In order to determine whether SMD produced related modules to LMM (e.g., since SMD modules are smaller, whether they represent sub-graphs of the larger LMM modules), we calculated a homogeneity score (described in methods section) between SMD and LMM modules in simulated networks. All networks contained 500 nodes. Modularity was calculated for the LMM partitions, and the homogeneity of SMD and LMM partitions was calculated. From

our simulated networks, the homogeneity between SMD and LMM module partitions was between 0.55 and 0.87 (Figure 3c,d). Notably, we found that for networks simulated with both power law and regular node degree distributions, as modularity of LMM partitions increased, the homogeneity between SMD- and LMM-partitioned modules increased (power law network Pearson  $R = 0.87$ ,  $p < .001$ ; regular network Pearson  $R = 0.93$ ,  $p < .001$ ). Thus, when network modularity is high (i.e., there is a high number of edges within the module compared to between modules), SMD partitions tend to be sub-partitions of LMM partitions.

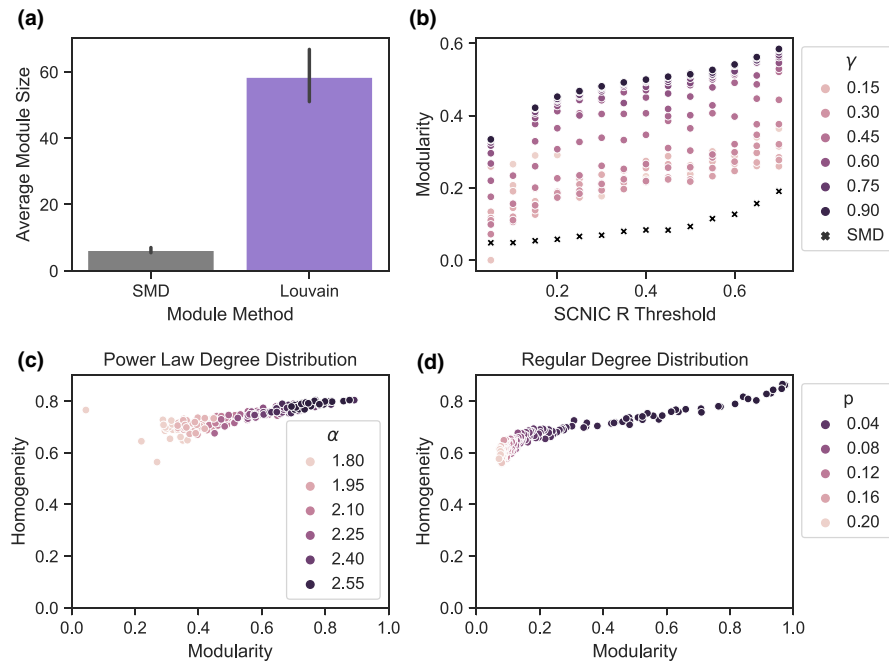
### 3.2 | *R*-value thresholds influence module size and phylogenetic relatedness of OTUs binned into a module

A key parameter in SCNIC is the *R*-value threshold used to pick modules. Use of a high *R*-value threshold would be expected to bin only very tightly correlated microbes with strong relationships, while less stringent thresholds may identify community-level patterns representing more loosely connected microbial pairs. To illustrate this concept, we binned OTUs into modules using the SMD method at *R*-value thresholds between 0.2 and 1.0 using the HIV dataset. As expected, at lower *R*-value thresholds, more OTUs were binned into modules and lower numbers of modules of smaller average size were formed as the threshold increased (Figure 4a). To illustrate the effects of *R*-values thresholds on the nature of the identified modules, we compare SCNIC outputs using *R*-value thresholds of 0.2, 0.4, and 0.65. As shown in Figure 4, which visualizes modules in Cytoscape using SCNIC output files, the *R*-value threshold influences the size and connectivity of the network. We also illustrate the effect of using different thresholds by examining the correlations between OTUs that are included in the first module output by SCNIC, which is the largest module (module-0) (Figure 5). All OTUs in module-0 are positively correlated with each other, since SCNIC only captures positive correlations.

Microbes co-occurring in the same environmental niche have previously been observed to be phylogenetically closer on average (Faust et al., 2012). This is likely because phylogenetic relatedness has been correlated with functional relatedness, such as through having more shared genome content, leading towards success in similar environments (Zaneveld et al., 2010). We show that increasing the *R*-value threshold results in modules that contain OTUs that are more phylogenetically similar on average (Figure 5).

### 3.3 | Use of SCNIC influences the detection of OTUs that differ between MSM and non-MSM

We next evaluated the effects of applying SCNIC with default SparCC and SMD parameters and varying *R*-value thresholds on downstream statistical analysis results. To investigate differential abundance based on MSM status in the HIV dataset we used ANCOM (Mandal



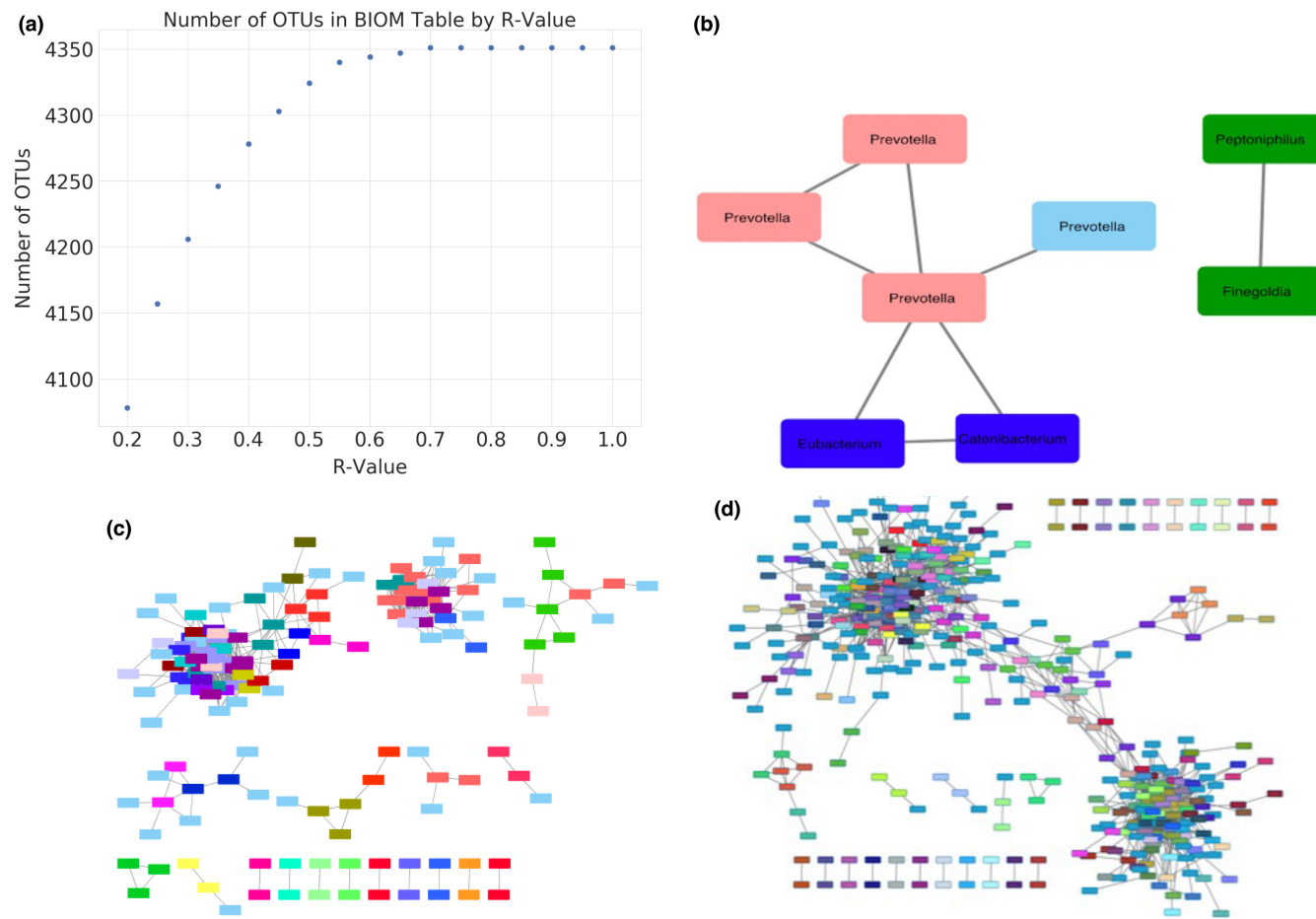
**FIGURE 3** Comparison of the SMD and LMM algorithms for module selection. Panels A and B show a comparison of module size and modularity between SMD and LMM selected modules on the Great Lakes data at  $R$  thresholds ranging from 0.1 to 0.7 and  $\gamma$  ranging from 0.15 and 0.9 for LMM. Panel a compares average module size and Panel b shows modularity across these parameters. In Panels C and D the homogeneity of LMM versus SMD selected module was calculated using simulated networks formed using both Power-Law (Panel c) and Regular (Panel d) degree distributions. Nodes are coloured by the  $\alpha$  (exponent) parameter for power law degree distributions in Panel c and by  $p$  (probability of one node being connected to any other node) in regular degree distributions in Panel d. A homogeneity of 1 denotes that all SMD modules are sub-modules of the LMM modules, whereas a homogeneity of 0 denotes that no two nodes grouped into a SMD module were partitioned into the same module by LMM.

et al., 2015). After removing taxa that were not present in at least 5% of the samples the OTU table had 317 samples and 639 OTUs. We found that 12 OTUs were significantly different between MSM and non-MSM without using SCNIC. Using SCNIC at  $R$ -values of 0.2, 0.4, and 0.65 and running ANCOM on the filtered output feature table, we found that most of the significant features were modules at an  $R$ -value of 0.2 and 0.4 but not 0.65 (e.g., 14 of the 15 significant features were modules at  $R = 0.2$ ) (Table 1). This was the case even though the vast majority of OTUs were not a part of modules at the 0.4  $R$ -value threshold (Figure 4a). The majority of 12 of the OTUs that were significant without running SCNIC, were grouped into modules with each other and with OTUs that were not individually significant without running SCNIC. These significant modules contained 74, 26, and 1 new OTU at  $R$ -values of 0.2, 0.4 and 0.65 respectively. Using SCNIC also resulted in the identification of 1, 5 and 25 (at  $R$ -values of 0.2, 0.4 and 0.65) OTUs that were individually significant that were not significant without running SCNIC, with no OTUs that were individually significant losing significance because they were binned in a module, indicating an increase in statistical power resulting from running a test like ANCOM that controls the FDR.

Considering correlation structure of significant features can help in understanding the broader community context of bacteria that differ with MSM status. In module-0 for each of the  $R$ -values, which significantly differed by MSM status in all cases, *Prevotella* was the dominant genus (Figure 5). At an  $R$ -value of 0.65, all OTUs in module-0

were assigned to the genus *Prevotella* (Figure 5c). However, at an  $R$ -value of 0.4 module-0 included seven *Prevotella* OTUs, one *Dialister*, and an unidentified member of the *Paraprevotellaceae* family. At the  $R$ -value of 0.2, *Prevotella* accounted for 13 of the 25 OTUs and 11 of the 12 pre-SCNIC significant OTUs were all found in this module. This suggests that individual OTUs that differ with MSM status may in some cases be a part of a consortium of diverse members that collectively display features that may contribute to differences in microbiome function.

To further explore this concept, we investigated the results generated with an  $R$ -value of 0.4, as the significant features maintain a strong level of correlation while being phylogenetically diverse. When running ANCOM on this feature table, we found that these individually significant OTUs tended to be joined into modules with other highly co-correlated microbes and that these modules significantly differed with MSM (Figure 6). Of particular note, we observe that the modules and taxa that are significantly related to MSM do not all correlate with each other. At the  $R$ -value of 0.4, module-36 contains two taxa, *Erysipelotrichaceae* and *Clostridium* that are negatively correlated with the other significant taxa and modules (Figure 6). Module-2 contains *Eubacterium*, *Catenibacterium* and *Prevotella* which are phylogenetically heterogeneous but mutually co-occurring. A follow up experiment, which leverages insights that SCNIC generates, may combine different strains of microbes to assemble a community type to test for functional correlates of disease.



**FIGURE 4** SCNIC feature reduction and visualization of SCNIC networks. (a) We used SCNIC to select modules using the OTU table from the HIV dataset, the SMD module selection algorithm, and SparCC  $R$ -values ranging from 0.2 to 1.0, in increments of 0.05. The  $R$ -value is plotted against the number of features in the resulting BIOM table produced by SCNIC. As the  $R$ -value increases the number of modules decreases and the number of single features (modules + OTUs not included in modules) increases. After the  $R$ -value of 0.65, the number of features in the resulting file remained the same at 4351 features, because there were no modules that were created past a SparCC  $R$  of 0.65. The Cytoscape output allows for easy exploration and visualization of relationships between different OTUs/taxa in an interactive interface. (b)  $R = 0.65$  (c)  $R = 0.4$  (d)  $R = 0.2$ . As the  $R$ -value increases, the size of the network decreases as SCNIC does not include singletons (features with no significant positive correlations) in the resulting network file. Correlation network where edges are correlations with a  $R$ -value greater than the threshold set. Nodes are OTUs and node colour represents module membership (i.e., module-0 is pink in Panel b).

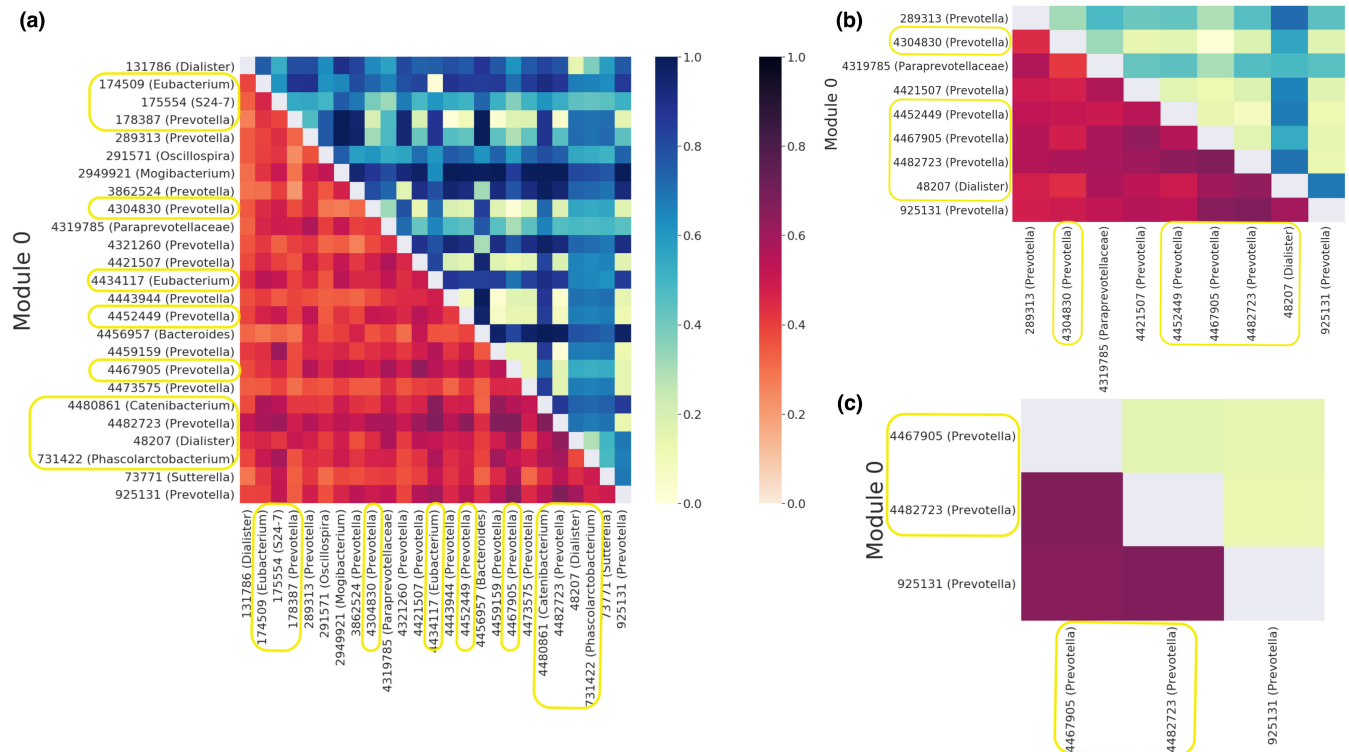
### 3.4 | Use of SCNIC results in the detection of lake associated taxa

To test consistency in patterns across different datasets, we also tested the effects of applying SCNIC with default SparCC and SMD parameters and varying  $R$ -value thresholds on results from the Great Lakes dataset. Specifically, we identified features that significantly differed between Lake Michigan ( $N = 16$ ) and Lake Superior ( $N = 33$ ) using ANCOM (Mandal et al., 2015).

We began with a table of 3871 OTUs, and 725 of these remained after removing OTUs not present in at least 5% of the samples. We found that 168 OTUs were significantly different between lakes without using SCNIC using ANCOM. Using SCNIC at  $R$ -values of 0.2, 0.4, and 0.65 and running ANCOM on the filtered output OTU table, we found that most significant features were modules at an  $R$  threshold of 0.4 but not 0.2 or 0.65 (Table 1). Use of SCNIC resulted in the

detection of individual OTUs that were now significant that were not before (3 and 13 for  $R$ -value thresholds of 0.4 and 0.65 respectively). Application of SCNIC also identified many additional OTUs that became of interest because they were now part of significant modules (139, 64, and 12 OTUs at 0.2, 0.4, and 0.65 respectively; Table 1). However, unlike for the HIV dataset, several OTUs that were individually significant were no longer significant with ANCOM after applying SCNIC and this effect was the most pronounced with lower  $R$ -value thresholds (64, 14, and 6 OTUs that were significant with SCNIC were no longer significant after applying SCNIC at 0.2, 0.4, and 0.65  $R$ -value thresholds respectively; Table 1). This is likely because microbes that differed between lakes were binned with loosely correlated microbes that did not, leading to a loss of signal. Thus, in this case, only SCNIC with a moderate to high  $R$ -value threshold appeared to balance the benefit of the increased power and disadvantages of loss of signal from binning loosely correlated features.





**FIGURE 5** Module-0 across different  $R$ -values. Module-0 expanded to view taxonomy and correlations among them at  $R$ -values of 0.2 (a), 0.4 (b), and 0.65 (c). The heatmap in the lower triangle corresponds to the correlation found by SparCC coloured on a light red (low correlation) to dark red (high correlation) spectrum as defined in the colour bar on the right. The heatmap in the upper triangle represents the phylogenetic distance between organism pairs coloured on a yellow (small phylogenetic distance) to dark blue (high phylogenetic distance) spectrum as defined in the colour bar on the right. As the  $R$ -value increases, the species in module-0 become more phylogenetically similar. Module-0 has 11, 5 and 2 of the significant Pre-SCNIC OTUs at  $R$ -values of 0.2, 0.4 and 0.65, and are highlighted in a yellow border.

**TABLE 1** Significant SCNIC modules and features across  $R$ -values in the HIV and Great Lakes datasets ANCOM analysis of the HIV and Great Lakes datasets after using SCNIC at  $R$ -value thresholds at 0.2, 0.4, and 0.65

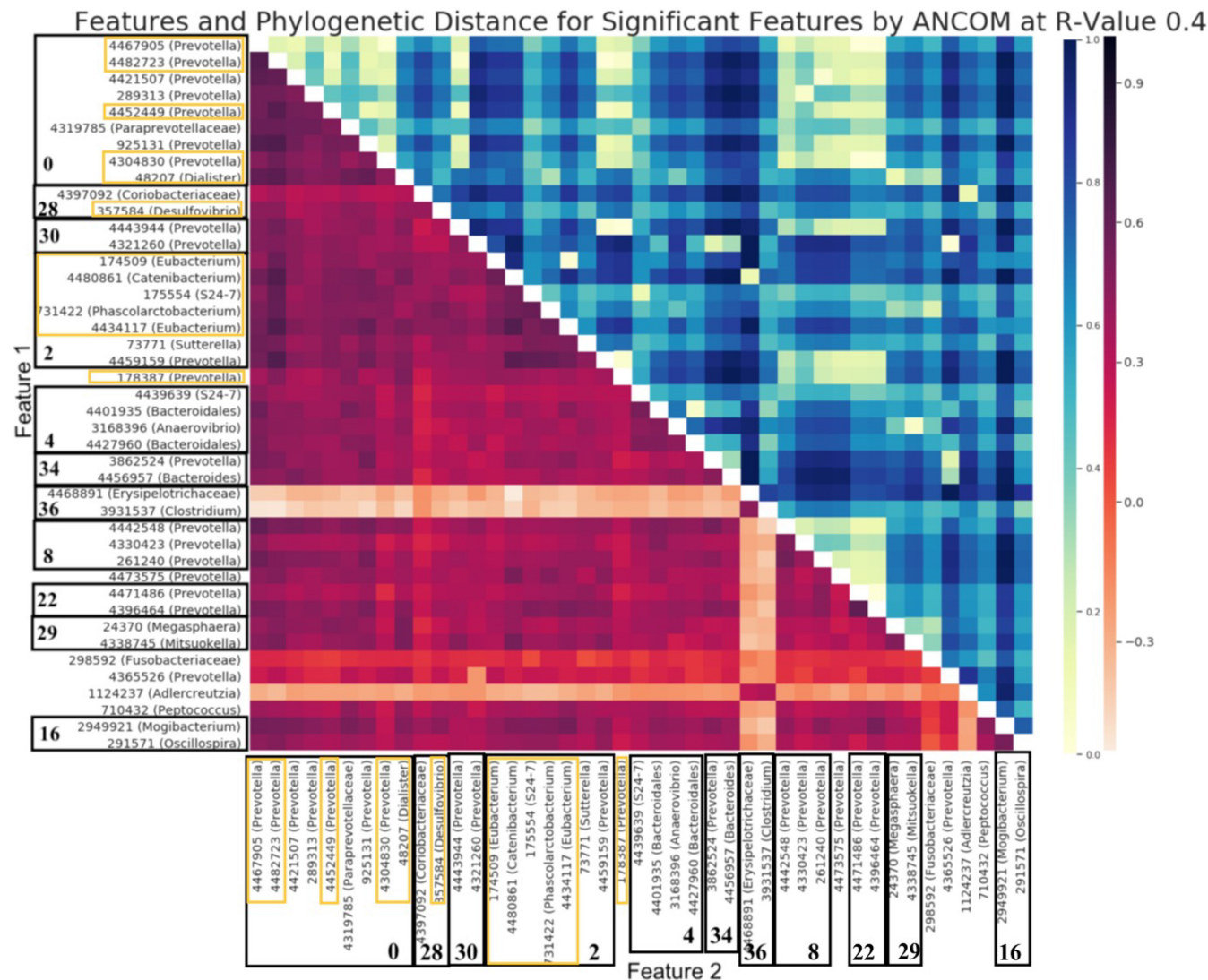
HIV dataset					
$R$ -value	New OTUs in sig. modules	New significant OTUs	Lost significant OTUs	# of significant modules	Total significant features
0.2	74	1	0	14	15
0.4	26	5	0	11	17
0.65	1	25	0	2	35
Great Lakes dataset					
$R$ -value	New OTUs in sig. modules	New significant OTUs	Lost significant OTUs	# of significant modules	Total significant features
0.2	139	0	64	1	25
0.4	64	3	14	29	33
0.65	12	13	6	24	74

Note: MSM was used as the categorical variable for differential abundance in the HIV analysis and the Great Lakes analysis tested for taxa that differed between lakes.

### 3.5 | Time and memory resources used by SCNIC

SCNIC's module generation step can be run locally on a laptop computer. For both the Great Lakes and GT datasets, which had 764 and 1301 features respectively, SCNIC's module generation

step ran in <2 min when using the SMD method and <30s when using the Louvain method. The Great Lakes dataset used <200 mebibytes (MiB) memory, and the GT dataset required a maximum of 300 MiB memory. Table S1 shows time and memory usage during SCNIC.



**FIGURE 6** All significant features at R-value 0.4 found by ANCOM. Each of the borders in the y-axis represents the different modules, with the module number bolded. The Pre-SCNIC OTUs that were significant are highlighted in a yellow border. The heatmap in the lower triangle corresponds to the correlation found by SparCC coloured on a light red (low correlation) to dark red (high correlation) spectrum as defined in the colour bar on the right. The heatmap in the upper triangle represents the phylogenetic distance between organism pairs coloured on a yellow (small phylogenetic distance) to dark blue (high phylogenetic distance) spectrum as defined in the colour bar on the right.

## 4 | DISCUSSION

SCNIC provides a method to measure correlations, find and visualize modules of correlated features, and summarize modules by summing their counts for use in downstream statistical analysis as one method for dimensionality reduction. Using SCNIC with the SMD algorithm for module detection aids in feature reduction in 16S rRNA sequencing data while ensuring a minimum strength of association within modules. As expected, our workflow identified modules in which OTUs tended to be phylogenetically related, especially at relatively high values of  $R$ . Using SCNIC, we overall achieved increased statistical power from performing less comparisons, but use of low  $R$ -value thresholds had the potential to lead to loss of significance by binning loosely correlated features. In these analyses we used

SparCC to calculate pairwise correlations in compositional data but Spearman and Pearson are also implemented for cases when the underlying data do not match those well suited for SparCC (e.g., if they are not sparse or with an inverse Simpson index above 13). In these analyses, we also used OTUs as features; however, other microbiome features can be used with SCNIC, such as ASVs, genera, or species defined with a taxonomic classifier, as well as other data types such as metabolome data. SCNIC has also been used in previously published work to perform feature reduction prior to random forest analysis with the microbiome and diverse other data types (Armstrong et al., 2021).

SCNIC complements existing methods because these either: (1) form correlation networks of microbes for visualization but do not have functionality for selecting and summarizing modules for

downstream statistical analysis (Faust & Raes, 2016), (2) can select and summarize modules for downstream statistical analysis but are designed for gene expression and not microbiome data (Langfelder & Horvath, 2008), only summarize features if they are phylogenetically related (Frøslev et al., 2017), or suggest methods for finding modules of correlated microbes but do not provide a convenient implementation (Blondel et al., 2008). SCNIC is available both as a stand-alone application and as a QIIME 2 plugin for easy integration with existing microbiome workflows.

SCNIC implements both the LMM algorithm, which had been previously recommended for selecting modules of correlated microbes (Baldassano & Bassett, 2016; Jackson et al., 2018), and a novel SMD algorithm. The advantage of the SMD algorithm is that all pairs of features in the module have an *R*-value greater than the user-provided minimum threshold. Using real and simulated data, we showed that SMD produced smaller modules that generally represent sub-graphs of the larger LMM modules. Since the use of lower *R*-value thresholds similarly produced larger modules including more weakly correlated modules, we speculate that use of LMM might result in a similar trend of identifying more OTUs within significant modules, but with the disadvantage of individually significant OTUs being lost because they are combined with loosely correlated microbes that are not related to the outcome being tested.

We illustrate here that varying the *R*-value threshold input by the user has a great impact on the results. However, we have avoided giving specific *R*-value threshold recommendations here, because optimal *R*-values may vary across datasets and data types. Using higher *R*-values thresholds was more likely to identify highly phylogenetically related microbes that likely share overlapping functionality, and in principle could also identify diverse organisms with overlapping niches or highly complementary metabolic functions. Using a lower *R*-value threshold bins a broader community of more loosely correlated features with the risk of bringing together features which should not be grouped and losing significance of OTUs – as was illustrated in the Great Lakes dataset analysis conducted here. By summarizing correlated features, SCNIC can mitigate over-correction in multiple test adjustments by reducing the number of taxa and false discovery rate for downstream analysis. However, further work with both real and simulated datasets is required to determine the degree to which network characteristics that are inherent in different microbiome datasets may influence the optimal methods for both selecting and summarizing modules in order to enhance statistical power.

The results of our HIV dataset analysis confirm original findings, as well as those of another study (Armstrong et al., 2018), but included many new significantly associated taxa. SCNIC also assists in interpretation of microbiome data by identifying correlations among these taxa. Our results recapitulated those of the original publication of these data and previous HIV microbiome studies that all found enrichment of *Prevotella* with MSM status (Armstrong et al., 2018; Dillon et al., 2014; Lozupone et al., 2013; Noguera-Julian et al., 2016). However, our analyses provide additional insight by identifying correlations between differentiating

taxa. For instance, in module-0, which was more abundant in MSM samples, OTUs assigned taxonomically to the *Prevotella* genus are correlated with two OTUs identified as *Eubacterium bifforme* (which has recently been renamed *Holdemanella biformis* (De Maesschalck et al., 2014)). *Prevotella copri* has previously been associated with increased inflammation (Dillon et al., 2014) while in vitro stimulations of human immune cells have found that *P. copri* did not induce particularly high levels of inflammation but *E. biforme* did (Lozupone et al., 2013). This strong correlation between *P. copri* and *E. biforme* in MSM could explain the increased inflammation seen in individuals with higher levels of *P. copri*, with *E. biforme* being the true driver. Indeed, MSM status has previously been associated with increased inflammation (Gianella et al., 2012; Palmer et al., 2014). With the use of SCNIC, this correlation highlighted a route of mechanistic understanding which could be functionally followed up on in further experimental studies.

SCNIC detected multiple significant modules, of which none of the OTUs within were significant when analysed separately. Module-20, which was associated with MSM status, is the fourth most significant feature at *R*-value of 0.2, and is made up of *Acidaminococcus*, *Megasphaera*, and *Mitsuokella* species. These are all from the *Veillonellaceae* family which is likely the explanation for their correlation. Members of the *Veillonellaceae* family have been linked with inflammation (Bajaj et al., 2012).

By increasing statistical power and providing context for the relationships between significant taxa, SCNIC modules open new opportunities for analysis. When a module is associated with a variable of interest, the correlations within the module may imply functional relationships. These can be further investigated with in vitro and in vivo experiments. Studies which aim to test hypotheses generated from correlative analysis will commonly use a single significantly associated microbes. This often does not adequately represent in vivo systems because microbes in isolation often do not affect a disease state or their environment. SCNIC can enhance these confirmatory studies by identifying groups of microbes that may grow better than individual microbes and may better elicit relevant phenotypes than when grown separately.

#### AUTHOR CONTRIBUTIONS

M.S. coded the initial implementation of SCNIC and made major contributions to its conceptualization and design. K.T. improved the SCNIC implementation and performed the HIV case study. J.S. performed the comparisons of SMD to LMM with real and simulated data. C.L. conceptualized SCNIC and guided its implementation and design. K.T., M.S., J.S. and C.L. all wrote the manuscript together.

#### ACKNOWLEDGEMENTS

We would like to thank Elmar Pruesse for input on the design of SCNIC. We thank Jennifer Fouquier, Abigail Armstrong and Casey Martin for beta testing SCNIC. We thank Jennifer Fouquier for reading and commenting on the manuscript. Funding for KT came from the University of Colorado School of Medicine Research Track. Funding for MS came from NIH NLM 4T15 LM009451-10. Funding

for JS came from the National Science Foundation funding of the Interdisciplinary Quantitative Biology program and the William J. Freytag Fellowship.

## CONFLICT OF INTEREST

The authors declare no conflicts of interest.

## DATA AVAILABILITY STATEMENT

The Noguera-Julian et al. data set is available from NCBI SRA accession number SRP068240. The lakes dataset is available from QIITA accession number 1041; [dataset] IrsiCaixa Foundation. human gut metagenome, Human feces metagenome 16s rDNA sequencing. 2015/12. In: BioProject [Internet]. Bethesda, MD: National Library of Medicine (US), National Center for Biotechnology Information; 2011-. Available from: <http://www.ncbi.nlm.nih.gov/bioproject/PRJNA307231>. NCBI:BioProject: PRJNA307231; [dataset] Karl J Rockne; 2016; Great lakes Microbiome; QIITA; <https://qiita.ucsd.edu/study/description/1041>

## ORCID

Kumar Thurimella  <https://orcid.org/0000-0002-0819-4378>

## REFERENCES

- Armstrong, A. J. S., Quinn, K., Fouquier, J., Li, S. X., Schneider, J. M., Nusbacher, N. M., Doenges, K. A., Fiorillo, S., Marden, T. J., Higgins, J., Reisdorph, N., Campbell, T. B., Palmer, B. E., & Lozupone, C. A. (2021). Systems analysis of gut microbiome influence on metabolic disease in HIV-positive and high-risk populations. *mSystems*, *6*(3), e01178–20.
- Armstrong, A. J. S., Shaffer, M., Nusbacher, N. M., Griesmer, C., Fiorillo, S., Schneider, J. M., Preston Neff, C., Li, S. X., Fontenot, A. P., Campbell, T., Palmer, B. E., & Lozupone, C. A. (2018). An exploration of Prevotella-rich microbiomes in HIV and men who have sex with men. *Microbiome*, *6*(1), 198.
- Bajaj, J. S., Ridlon, J. M., Hylemon, P. B., Thacker, L. R., Heuman, D. M., Smith, S., Sikaroodi, M., & Gillevet, P. M. (2012). Linkage of gut microbiome with cognition in hepatic encephalopathy. *American journal of physiology-gastrointestinal and liver. Physiology*, *302*(1), G168–G175.
- Baldassano, S. N., & Bassett, D. S. (2016). Topological distortion and reorganized modular structure of gut microbial co-occurrence networks in inflammatory bowel disease. *Scientific Reports*, *6*, 26087.
- Ban, Y., An, L., & Jiang, H. (2015). Investigating microbial co-occurrence patterns based on metagenomic compositional data. *Bioinformatics*, *31*(20), 3322–3329.
- Barabási, A.-L., & Albert, R. (1999). Emergence of scaling in random networks. *Science*, *286*(5439), 509–512.
- Barberán, A., Bates, S. T., Casamayor, E. O., & Fierer, N. (2012). Using network analysis to explore co-occurrence patterns in soil microbial communities. *The ISME Journal*, *6*(2), 343–351.
- Benjamini, Y., & Hochberg, Y. (2000). On the adaptive control of the false discovery rate in multiple testing with independent statistics. *Journal of Educational and Behavioral Statistics*, *25*(1), 60–83.
- Blondel, V. D., Guillaume, J. L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, *2008*(10), P10008.
- Bolyen, E., Rideout, J. R., Dillon, M. R., Bokulich, N. A., Abnet, C. C., al-Ghalith, G. A., Alexander, H., Alm, E. J., Arumugam, M., Asnicar, F., Bai, Y., Bisanz, J. E., Bittinger, K., Brejnrod, A., Brislawn, C. J., Brown, C. T., Callahan, B. J., Caraballo-Rodríguez, A. M., Chase, J., ... Caporaso, J. G. (2019). Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nature Biotechnology*, *37*(8), 852–857.
- Bray, J. R., & Curtis, J. T. (1957). An ordination of the upland Forest communities of southern Wisconsin. *Ecological Monographs*, *27*(4), 325–349.
- Broido, A. D., & Clauset, A. (2019). Scale-free networks are rare. *Nature Communications*, *10*(1), 1017.
- Burkepile, D. E., Parker, J. D., Woodson, C. B., Mills, H. J., Kubanek, J., Sobczyk, P. A., & Hay, M. E. (2006). Chemically mediated competition between microbes and animals: Microbes as consumers in food webs. *Ecology*, *87*(11), 2821–2831.
- Bushnell, B., Rood, J., & Singer, E. (2017). BBMerge – Accurate paired shotgun read merging via overlap. *PLoS One*, *12*(10), e0185056.
- Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J. A., & Holmes, S. P. (2016). DADA2: High-resolution sample inference from Illumina amplicon data. *Nature Methods*, *13*(7), 581–583.
- Caporaso, J. G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F. D., Costello, E. K., Fierer, N., Peña, A. G., Goodrich, J. K., Gordon, J. I., Huttley, G. A., Kelley, S. T., Knights, D., Koenig, J. E., Ley, R. E., Lozupone, C. A., McDonald, D., Muegge, B. D., Pirrung, M., ... Knight, R. (2010). QIIME allows analysis of high-throughput community sequencing data. *Nature Methods*, *7*(5), 335–336.
- Castillo, J. D., Vivanco, J. M., & Manter, D. K. (2017). Bacterial microbiome and nematode occurrence in different potato agricultural soils. *Microbial Ecology*, *74*(4), 888–900.
- Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B., & de Hoon, M. J. L. (2009). Biopython: Freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics*, *25*(11), 1422–1423.
- Corno, G., Villiger, J., & Pernthaler, J. (2013). Coaggregation in a microbial predator–prey system affects competition and trophic transfer efficiency. *Ecology*, *94*(4), 870–881.
- De Maesschalck, C., et al. (2014). Faecalibacterium acidiformans gen. Nov., sp. nov., isolated from the chicken caecum, and reclassification of streptococcus pleomorphus (Barnes et al. 1977), Eubacterium bifforme (Eggerth 1935) and Eubacterium cylindroides (Cato et al. 1974) as Faecalibacterium pleomorphus comb. nov., Holdemanella bifformis gen. Nov., comb. nov. and Faecalitalea cylindroides gen. Nov., comb. nov., respectively, within the family Erysipelotrichaceae. *International Journal of Systematic and Evolutionary Microbiology*, *64*(Pt\_11), 3877–3884.
- Dillon, S. M., Lee, E. J., Kotter, C. V., Austin, G. L., Dong, Z., Hecht, D. K., Gianella, S., Siewe, B., Smith, D. M., Landay, A. L., Robertson, C. E., Frank, D. N., & Wilson, C. C. (2014). An altered intestinal mucosal microbiome in HIV-1 infection is associated with mucosal and systemic immune activation and endotoxemia. *Mucosal Immunology*, *7*(4), 983–994.
- Dugas, L. R., Bernabé, B. P., Priyadarshini, M., Fei, N., Park, S. J., Brown, L., Plange-Rhule, J., Nelson, D., Toh, E. C., Gao, X., Dong, Q., Sun, J., Kliethermes, S., Gittel, N., Luke, A., Gilbert, J. A., & Layden, B. T. (2018). Decreased microbial co-occurrence network stability and SCFA receptor level correlates with obesity in African-origin women. *Scientific Reports*, *8*(1), 17135.
- Edgar, R. C. (2010). Search and clustering orders of magnitude faster than BLAST. *Bioinformatics*, *26*(19), 2460–2461.
- Faust, K., & Raes, J. (2016). CoNet app: Inference of biological association networks using Cytoscape. *F1000Research*, *5*, 1519.
- Faust, K., Sathirapongsasuti, J. F., Izard, J., Segata, N., Gevers, D., Raes, J., & Huttenhower, C. (2012). Microbial co-occurrence relationships in the human microbiome. *PLoS Computational Biology*, *8*(7), e1002606.

- Fernandes, A. D., Reid, J. N. S., Macklaim, J. M., McMurrough, T. A., Edgell, D. R., & Gloor, G. B. (2014). Unifying the analysis of high-throughput sequencing datasets: Characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome*, 2(1), 15.
- Friedman, J., & Alm, E. J. (2012). Inferring correlation networks from genomic survey data. *PLoS Computational Biology*, 8(9), e1002687.
- Frøslev, T. G., Kjølner, R., Bruun, H. H., Ejrnæs, R., Brunbjerg, A. K., Pietroni, C., & Hansen, A. J. (2017). Algorithm for post-clustering curation of DNA amplicon data yields reliable biodiversity estimates. *Nature Communications*, 8(1), 1188.
- Gianella, S., Strain, M. C., Rought, S. E., Vargas, M. V., Little, S. J., Richman, D. D., Spina, C. A., & Smith, D. M. (2012). Associations between virologic and immunologic dynamics in blood and in the male genital tract. *Journal of Virology*, 86(3), 1307–1315.
- Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V., & Egozcue, J. J. (2017). Microbiome datasets are compositional: And this is not optional. *Frontiers in Microbiology*, 8, 2224.
- Gonzalez, A., Navas-Molina, J. A., Kosciulek, T., McDonald, D., Vázquez-Baeza, Y., Ackermann, G., DeReus, J., Janssen, S., Swafford, A. D., Orchanian, S. B., Sanders, J. G., Shorenstein, J., Holste, H., Petrus, S., Robbins-Pianka, A., Brislawn, C. J., Wang, M., Rideout, J. R., Bolyen, E., ... Knight, R. (2018). Qiita: Rapid, web-enabled microbiome meta-analysis. *Nature Methods*, 15(10), 796–798.
- Grüning, B., et al. (2018). Bioconda: Sustainable and comprehensive software distribution for the life sciences. *Nature Methods*, 15(7), 475–476.
- Haas, B. J., Gevers, D., Earl, A. M., Feldgarden, M., Ward, D. V., Giannoukos, G., Ciulla, D., Tabbaa, D., Highlander, S. K., Sodergren, E., Methé, B., DeSantis, T. Z., The Human Microbiome Consortium, Petrosino, J. F., Knight, R., & Birren, B. W. (2011). Chimeric 16S rRNA sequence formation and detection in sanger and 454-pyrosequenced PCR amplicons. *Genome Research*, 21(3), 494–504.
- Himsolt, M. 1997 GML: A portable graph file format. Technical report, Universität Passau.
- Huse, S. M., Welch, D. M., Morrison, H. G., & Sogin, M. L. (2010). Ironing out the wrinkles in the rare biosphere through improved OTU clustering. *Environmental Microbiology*, 12(7), 1889–1898.
- Jackson, M. A., Bonder, M. J., Kuncheva, Z., Zierer, J., Fu, J., Kurilshikov, A., Wijmenga, C., Zhernakova, A., Bell, J. T., Spector, T. D., & Steves, C. J. (2018). Detection of stable community structures within gut microbiota co-occurrence networks from different human populations. *PeerJ*, 6, e4303.
- Janssen, S., McDonald, D., Gonzalez, A., Navas-Molina, J. A., Jiang, L., Xu, Z. Z., Winker, K., Kado, D. M., Orwoll, E., Manary, M., Mirarab, S., & Knight, R. (2018). Phylogenetic placement of exact amplicon sequences improves associations with clinical information. *mSystems*, 3(3), e00021–18.
- Kuczynski, J., Liu, Z., Lozupone, C., McDonald, D., Fierer, N., & Knight, R. (2010). Microbial community resemblance methods differ in their ability to detect biologically relevant patterns. *Nature Methods*, 7(10), 813–819.
- Langfelder, P., & Horvath, S. (2008). WGCNA: An R package for weighted correlation network analysis. *BMC Bioinformatics*, 9, 559.
- LaSarre, B., McCully, A. L., Lennon, J. T., & McKinlay, J. B. (2017). Microbial mutualism dynamics governed by dose-dependent toxicity of cross-fed nutrients. *The ISME Journal*, 11(2), 337–348.
- Lozupone, C., Faust, K., Raes, J., Faith, J. J., Frank, D. N., Zaneveld, J., Gordon, J. I., & Knight, R. (2012). Identifying genomic and metabolic features that can underlie early successional and opportunistic lifestyles of human gut symbionts. *Genome Research*, 22(10), 1974–1984.
- Lozupone, C. A., Li, M., Campbell, T. B., Flores, S. C., Linderman, D., Gebert, M. J., Knight, R., Fontenot, A. P., & Palmer, B. E. (2013). Alterations in the gut microbiota associated with HIV-1 infection. *Cell Host & Microbe*, 14(3), 329–339.
- Mandal, S., et al. (2015). Analysis of composition of microbiomes: A novel method for studying microbial composition. *Microbial Ecology in Health and Disease*, 26, 27663.
- McDonald, D., Clemente, J. C., Kuczynski, J., Rideout, J. R., Stombaugh, J., Wendel, D., Wilke, A., Huse, S., Hufnagle, J., Meyer, F., Knight, R., & Caporaso, J. G. (2012). The biological observation matrix (BIOM) format or: How I learned to stop worrying and love the ome-ome. *Gigascience*, 1(1), 7.
- Mirarab, S., Nguyen, N., & Warnow, T. (2012). SEPP: SATé-enabled phylogenetic placement. *Pacific Symposium on Biocomputing*, 2012, 247–258.
- Noguera-Julian, M., Rocafort, M., Guillén, Y., Rivera, J., Casadellà, M., Nowak, P., Hildebrand, F., Zeller, G., Parera, M., Bellido, R., Rodríguez, C., Carrillo, J., Mothe, B., Coll, J., Bravo, I., Estany, C., Herrero, C., Saz, J., Sirera, G., ... Paredes, R. (2016). Gut microbiota linked to sexual preference and HIV infection. *eBioMedicine*, 5, 135–146.
- Palmer, C. D., Tomassilli, J., Sirignano, M., Romero-Tejeda, M., Arnold, K. B., Che, D., Lauffenburger, D. A., Jost, S., Allen, T., Mayer, K. H., & Altfeld, M. (2014). Enhanced immune activation linked to endotoxemia in HIV-1 seronegative MSM. *AIDS*, 28(14), 2162–2166.
- Paulson, J. N., Stine, O. C., Bravo, H. C., & Pop, M. (2013). Differential abundance analysis for microbial marker-gene surveys. *Nature Methods*, 10(12), 1200–1202.
- Pearson, K. (1909). Determination of the coefficient of correlation. *Science*, 30(757), 23–25.
- Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: A versatile open source tool for metagenomics. *PeerJ*, 4, e2584.
- Rosenberg, A., & Hirschberg, J. (2007). *V-measure: A conditional entropy-based external cluster evaluation measure*. Association for Computational Linguistics.
- Shannon, P., Markiel, A., Ozier, O., Baliga, N. S., Wang, J. T., Ramage, D., Amin, N., Schwikowski, B., & Ideker, T. (2003). Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Research*, 13(11), 2498–2504.
- Spearman, C. (1904). Measurement of association, part II. Correction of 'systematic deviations'. *The American Journal of Psychology*, 15, 88–101.
- Talevich, E., Invergo, B. M., Cock, P. J., & Chapman, B. A. (2012). BioPhylo: A unified toolkit for processing, analyzing and visualizing phylogenetic trees in biopython. *BMC Bioinformatics*, 13(1), 209.
- Thompson, L. R., Sanders, J. G., McDonald, D., Amir, A., Ladau, J., Locey, K. J., Prill, R. J., Tripathi, A., Gibbons, S. M., Ackermann, G., Navas-Molina, J. A., Janssen, S., Kopylova, E., Vázquez-Baeza, Y., González, A., Morton, J. T., Mirarab, S., Zech Xu, Z., Jiang, L., ... Zhao, H. (2017). A communal catalogue reveals Earth's multiscale microbial diversity. *Nature*, 551(7681), 457–463.
- Tong, M., McHardy, I., Ruegger, P., Goudarzi, M., Kashyap, P. C., Haritunians, T., Li, X., Graeber, T. G., Schwager, E., Huttenhower, C., Fornace, A. J., Jr., Sonnenburg, J. L., McGovern, D. P. B., Borneman, J., & Braun, J. (2014). Reprogramming of gut microbiome energy metabolism by the FUT2 Crohn's disease risk polymorphism. *The ISME Journal*, 8(11), 2193–2206.
- Tsilimigras, M. C. B., & Fodor, A. A. (2016). Compositional data analysis of the microbiome: Fundamentals, tools, and challenges. *Annals of Epidemiology*, 26(5), 330–335.
- Watts, S. C., Ritchie, S. C., Inouye, M., & Holt, K. E. (2019). FastSpar: Rapid and scalable correlation estimation for compositional data. *Bioinformatics*, 35(6), 1064–1066.
- Weiss, S., van Treuren, W., Lozupone, C., Faust, K., Friedman, J., Deng, Y., Xia, L. C., Xu, Z. Z., Ursell, L., Alm, E. J., Birmingham, A., Cram, J. A., Fuhrman, J. A., Raes, J., Sun, F., Zhou, J., & Knight, R. (2016). Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *The ISME Journal*, 10(7), 1669–1681.
- Widder, S., Besemer, K., Singer, G. A., Ceola, S., Bertuzzo, E., Quince, C., Sloan, W. T., Rinaldo, A., & Battin, T. J. (2014). Fluvial

network organization imprints on microbial co-occurrence networks. *Proceedings of the National Academy of Sciences of the United States of America*, 111(35), 12799–12804.

- Wilcox, R. R. (2011). *Introduction to robust estimation and hypothesis testing*. Academic Press.
- Yin, J., Han, H., Li, Y., Liu, Z., Zhao, Y., Fang, R., Huang, X., Zheng, J., Ren, W., Wu, F., Liu, G., Wu, X., Wang, K., Sun, L., Li, C., Li, T., & Yin, Y. (2017). Lysine restriction affects feed intake and amino acid metabolism via gut microbiome in piglets. *Cellular Physiology and Biochemistry*, 44(5), 1749–1761.
- Younge, N., Yang, Q., & Seed, P. C. (2017). Enteral high fat-polyunsaturated fatty acid blend alters the pathogen composition of the intestinal microbiome in premature infants with an Enterostomy. *The Journal of Pediatrics*, 181, 93–101.e6.
- Zaneveld, J. R., Lozupone, C., Gordon, J. I., & Knight, R. (2010). Ribosomal RNA diversity predicts genome diversity in gut bacteria and their relatives. *Nucleic Acids Research*, 38(12), 3869–3879.

## SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

**How to cite this article:** Shaffer, M., Thurimella, K., Sterrett, J. D., & Lozupone, C. A. (2023). SCNIC: Sparse correlation network investigation for compositional data. *Molecular Ecology Resources*, 23, 312–325. <https://doi.org/10.1111/1755-0998.13704>