


METHODOLOGY

Open Access



Disentangling environmental effects in microbial association networks

Ina Maria Deutschmann^{1*} , Gipsi Lima-Mendez², Anders K. Krabberød³, Jeroen Raes^{4,5}, Sergio M. Vallina⁶, Karoline Faust^{5*†} and Ramiro Logares^{1*†}

Abstract

Background: Ecological interactions among microorganisms are fundamental for ecosystem function, yet they are mostly unknown or poorly understood. High-throughput-omics can indicate microbial interactions through associations across time and space, which can be represented as association networks. Associations could result from either ecological interactions between microorganisms, or from environmental selection, where the association is environmentally driven. Therefore, before downstream analysis and interpretation, we need to distinguish the nature of the association, particularly if it is due to environmental selection or not.

Results: We present EnDED (environmentally driven edge detection), an implementation of four approaches as well as their combination to predict which links between microorganisms in an association network are environmentally driven. The four approaches are sign pattern, overlap, interaction information, and data processing inequality. We tested EnDED on networks from simulated data of 50 microorganisms. The networks contained on average 50 nodes and 1087 edges, of which 60 were true interactions but 1026 false associations (i.e., environmentally driven or due to chance). Applying each method individually, we detected a moderate to high number of environmentally driven edges—87% sign pattern and overlap, 67% interaction information, and 44% data processing inequality. Combining these methods in an intersection approach resulted in retaining more interactions, both true and false (32% of environmentally driven associations). After validation with the simulated datasets, we applied EnDED on a marine microbial network inferred from 10 years of monthly observations of microbial-plankton abundance. The intersection combination predicted that 8.3% of the associations were environmentally driven, while individual methods predicted 24.8% (data processing inequality), 25.7% (interaction information), and up to 84.6% (sign pattern as well as overlap). The fraction of environmentally driven edges among negative microbial associations in the real network increased rapidly with the number of environmental factors.

* Correspondence: ina.m.deutschmann@gmail.com;
karoline.faust@kuleuven.be; ramiro.logares@icm.csic.es

Karoline Faust and Ramiro Logares are shared author.

¹Institute of Marine Sciences, CSIC, Passeig Marítim de la Barceloneta, 37-49,
08003 Barcelona, Spain

⁵KU Leuven Department of Microbiology, Immunology and Transplantation,
Rega Institute, Laboratory of Molecular Bacteriology, Herestraat 49, 3000
Leuven, Belgium

Full list of author information is available at the end of the article



© The Author(s). 2021 **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.

Conclusions: To reach accurate hypotheses about ecological interactions, it is important to determine, quantify, and remove environmentally driven associations in marine microbial association networks. For that, EnDED offers up to four individual methods as well as their combination. However, especially for the intersection combination, we suggest using EnDED with other strategies to reduce the number of false associations and consequently the number of potential interaction hypotheses.

Keywords: Microbial interactions, Association network, Effect of indirect dependencies, Environmentally driven edge detection

Background

Association networks to generate microbial interaction hypotheses

There is a myriad of microorganisms on Earth; current estimates indicate $\approx 10^{12}$ microbial species [1], and $\approx 10^{30}$ microbial cells [2, 3]. Microorganisms have crucial roles in the biosphere by contributing to global biogeochemical cycles [4] and underpinning diverse food webs. The importance of microbes for the functioning of ecosystems cannot be understood without considering their ecological interactions [5, 6]. These allow transferring carbon and energy to upper trophic levels, and the recycling of nutrients and energy [7]. Furthermore, ecological interactions influence microbial community turnover and composition. These interactions include win-win (e.g., mutual cross-feeding and cooperation), win-loss (e.g., predator-prey and host-parasite), and loss-loss (e.g., resource competition) relationships [8]. Although microbial communities are highly interconnected [9], our knowledge about ecological interactions in the microbial world is still limited [6, 10].

Previous studies have shown relationships between a restricted number of microorganisms. However, we need a large number of interactions to understand the functioning of complex ecosystems. This is challenging, in part, due to the vast number of possible interactions—given n microorganisms, there are $\binom{n}{2} = n(n-1)/2$ potential pairwise interactions. Thus, it is unfeasible to test them experimentally within a reasonable amount of time and cost. The problem of having a large number of potential interactions can be partially circumvented with omics technologies coupled to network analyses.

Omics can identify and quantify a large number of microorganisms from a given sample. Typically, the relative abundance for each identified organism per sample is estimated. There are multiple methods to determine associations (normally based on correlations) between microorganisms using their abundances (e.g., eLSA [11, 12], CoNet [13], SPIEC-EASI [14], or FlashWeave [15]). These abundance-based associations compose a network, where nodes represent microorganisms and edges represent either co-presence (positive association) or mutual exclusion (negative association) relationships, which constitute microbial interaction hypotheses.

Challenges in using networks as a representation of the microbial ecosystem

Although networks play an essential role in understanding complex systems, microbial ecological networks are not yet as developed in terms of inference and biological interpretation [16]. Network inference from -omics data is difficult [9, 17] because of both technical and interpretation challenges. One challenge is the compositional nature of the data produced by DNA sequencers [18]. There are several network tools [17] that consider this, e.g., SPIEC-EASI [14]. Other difficulties include data based on a small number of samples relative to the number of microorganisms they contain, i.e., a low sample-to-microorganisms ratio; plus sparse data—too many zeros in the dataset that can wrongly associate microorganisms [19]. A zero indicates either the absence of a microorganism (structural zero), or an insufficient detection level or sequencing depth (sampling zero). Thus, we should remove microorganisms appearing in just a few samples.

Interpretation of association networks is challenging because they are not equivalent to ecological networks. Edges in ecological networks represent observed ecological interactions between different microorganisms like parasitism or competition [20]. Ecological networks are directed graphs, where the directed edges (arcs) point from a start node (source) to an end node (target). In contrast, association networks are undirected. Although association networks provide ecological insight, they do not necessarily encode causal relationships or observed ecological interactions. Unless edges are verified with experiments or additional information, one should be careful when attributing biological meaning to network properties [21]. In addition, networks with too many edges (dense networks or hairballs) make interpretation more challenging. We can reduce network density when lowering the corrected p value for inferred edges [22], or increasing the cut-off for other criteria such as the association strength, prevalence, or abundance filtering [21]. Another strategy is agglomeration using taxonomic or ecological (functional) groupings [23].

The interpretation challenge addressed in this study are indirect dependencies (associations) caused by

environmental factors. For most microbial association networks, an edge indicates one of the following three alternatives:

1. Ecological interaction between two microorganisms,
2. Similar or contrary dependence (i.e., preference) to environmental factor/s or a third microorganisms,
3. Association by chance.

Indirect associations occur when two microorganisms are both dependent on an abiotic environmental factor (e.g., same nutrients and temperature requirements) or biotic factor (e.g., same prey or predator), but do not interact with one another. Here, indirect association describes the computational effect of indirect dependencies, and observing an association when in fact there is none.

Removing indirect dependencies including environmental effects

To distinguish between direct and indirect interactions, several network construction tools use a probabilistic graphical model [14, 24], e.g., SPIEC-EASI [14, 25], miic [26], or FlashWeave [15]. FlashWeave can also integrate metadata to avoid indirect associations driven by environmental factors but currently does not support missing data. The tool ARACNE [27] aims to eliminate indirect associations by using an information theoretic property (the *Data Processing Inequality*, DPI, in “Methods” section). The extension TimeDelay-ARACNE [28] tries to extract dependencies between different times. Another approach including time delay is implemented in the tool MIDER [29], which combines mutual information-based distances and entropy reduction to detect indirect interactions (*Mutual Information*, MI, in Methods). PREMER [30], an successor of MIDER, allows to include previous knowledge, e.g., associations known to be absent.

There are also several prior network construction approaches to reduce indirect associations, e.g., a high prevalence filter that preserves microorganisms present in many samples [31]. However, this will keep generalists while removing specialists. Another approach divides datasets displaying a great environmental heterogeneity into subdatasets of similar environmental conditions [21]. For example, a previous work [32] constructed two networks representing bacterial soil communities from two different sections of a pH, temperature, and humidity gradient. Another work [23] constructed ocean depth-specific networks to account for environmental differences between the surface layer and the deep chlorophyll maximum layer. In addition to dividing samples, an algorithm aiming to correct for habitat filtering effects [33], subtracts, for a given habitat, the mean

abundance from each microorganisms within each sample. However, this approach is limited to the identified habitat groups that should have a similar sample size.

In contrast, there are methods accounting for indirect dependencies after network construction. For instance, global silencing, [34] and network deconvolution [35] aim to recover true direct associations from observed correlations. Both techniques are sensitive to missing variables [36]. Another method, called *sign pattern*, SP, uses environmental triplets [23]. An environmental triplet contains two microorganisms and one environmental factor, which are associated to each other. SP combines the signs of association scores (positive or negative) to determine if a microbial association should be classified as indirect (SP in “Methods” section). Its major drawback is edge removal where microorganisms with similar environmental preference interact. Along SP and network deconvolution, the *interaction information*, II, was applied in [23]. Within an environmental triplet, the II method aims to indicate whether an edge is due entirely to shared environmental preferences ($II < 0$) or whether environmental preferences and true interactions are entangled ($II > 0$). However, II cannot determine which associations in a triplet are indirect (II in “Methods” section). Here, we study several indirect edge detection methods: SP, *overlap*, (OL, developed here), II, DPI, and their combination.

EnDED is an implementation of four methods and their combination

This article presents EnDED, which implements four approaches, and their combination, to indicate environmentally driven (indirect) associations in microbial networks. The four methods are sign pattern [23], overlap (developed here), interaction information [23, 37], and data processing inequality [27, 38]. SP requires an association score that represents co-occurrence when it is positive, and mutual exclusion when it is negative. OL requires temporal data with a known start and end of the association to determine whether the microbial association occurs in a time window when both microorganisms are associated to the same environmental factor. The II method indicates the existence of one indirect dependency between three components that are associated with each other. The DPI method states that the association with the smallest mutual information is the indirect association. Here, we evaluate each method and their combination on how well they detect environmentally driven associations on association networks from simulated data including two environmental factors. Combining methods in an intersection approach retains more true interactions than each method on its own. A union approach was discarded because it would have retained the smallest number of true interactions. We are able to

disentangle and filter environmentally-driven edges from microbial association networks (0.95–0.96 in positive predictive value and 0.35–0.83 in accuracy). We also applied EnDED to disentangle and filter environmentally driven edges from a real marine microbial association network based on 10 years of monthly sampling including ten environmental factors. EnDED contributed to both, generating more reliable hypotheses on microbial interactions, and facilitating network analysis by removing edges from dense “hairball” networks. EnDED is publicly available [39].

Results

Simulated data

To evaluate EnDED’s performance in removing environmentally driven associations, we simulated 1000 abundance time-series datasets with 50 microorganisms and known true interactions between them. We obtained another 1000 datasets with noise (hereafter dwn). We constructed the networks (hereafter simulated networks) with the tool eLSA [11, 12] (see “Methods” section). The simulated networks contained on average (computed as the median) 50 nodes and 1087 edges (1063 dwn), of which 60 (59 dwn) were true interactions (edges present in the inferred and true network) and 1026 (1005 dwn) false associations (edges present in the inferred but absent in the true network). Networks inferred from simulated data without noise contained on average one more true interaction but also 21 more false interactions than the networks inferred from simulated data with noise.

A simple approach to discriminate true interactions (desired) from false associations (undesired) would be to use a threshold for the association strength, which could be suitable if the values for true interactions and false associations are (i) following different distributions, and (ii) the distributions are mainly non-overlapping. We tested the former requirement with a two-sample Kolmogorov-Smirnov test with the R [40] function `ks.test`. Using a 95% (99%, 99.9%) confidence level, the distributions were significantly different for 358 (192, 66) simulated datasets and 355 (173, 68) simulated datasets with noise, which is slightly more than one third of them. This indicates that an association strength cut-off is unsuitable to separate true interactions from false associations. More sophisticated approaches than a simple threshold include the methods implemented in EnDED: SP, OL, II, DPI, and their combination.

Combining the methods in an intersection approach (hereafter referred to as intersection combination), we classified on average 348 (228 dwn), that is 32% (22% dwn) of the associations, to be environmentally driven. The number of correctly detected false associations was on average 332 (219 dwn), i.e., 96% of the removed edges. The resulting networks contained on average 737

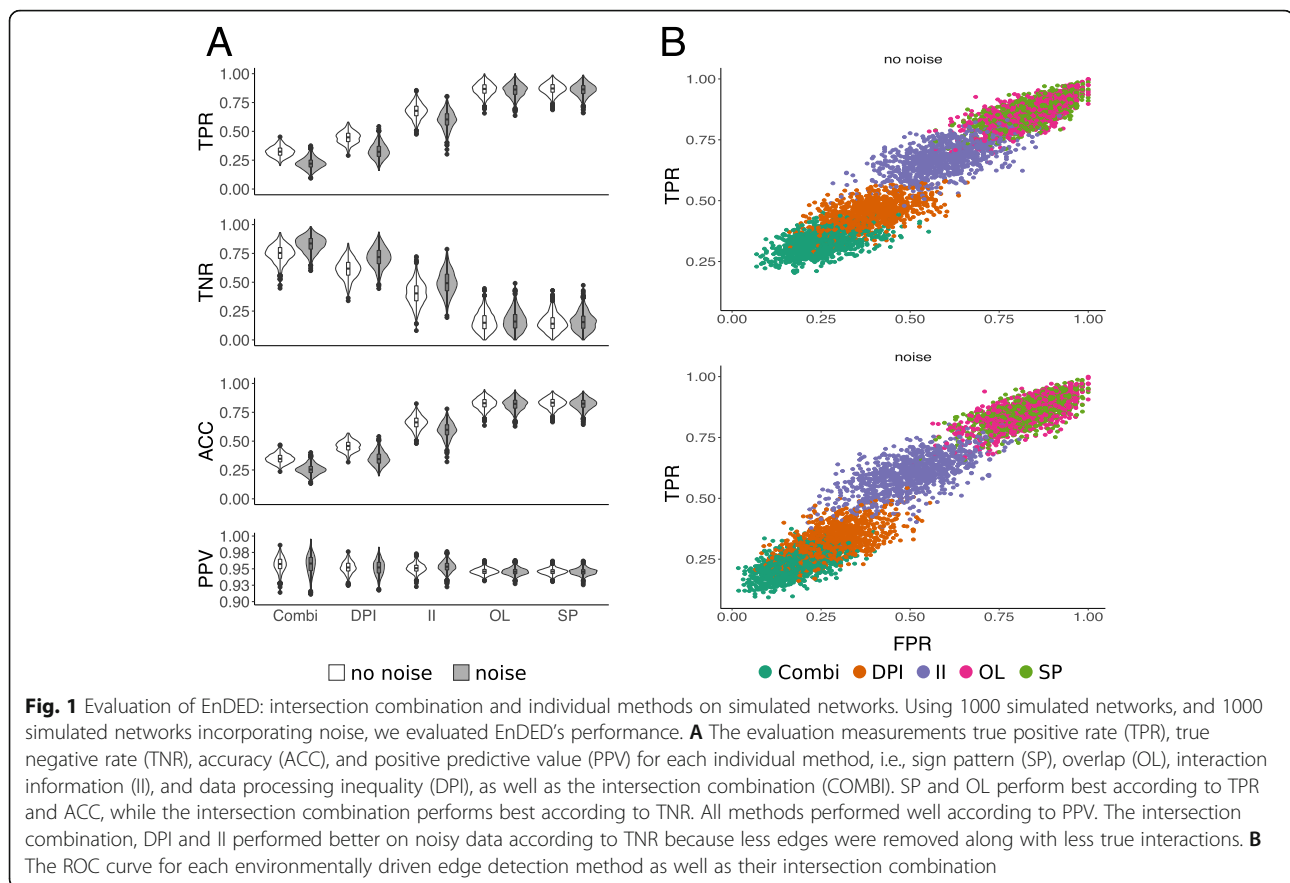
(828 dwn) edges. When each method was individually applied more edges were removed: 87% (86% dwn) for SP and OL, 67% (60% dwn) for II, and 44% (32% dwn) for DPI. The fraction of correctly removed edges for individual methods was on average 95%. Comparing the methods on correctly detected false associations, the greatest agreement was observed between SP and OL, whereas DPI appeared to be the most conservative in not agreeing with other methods and, subsequently, reducing the number of detected edges in the intersection combination approach (Supplementary Table S1). Individual methods removed more edges from the network than the intersection combination, where all methods must agree. However, a method’s performance is not solely determined by the number of removed edges.

To evaluate the removal of environmentally driven edges, we scored the different approaches based on five evaluation measurements (see “Methods” section): the true positive rate, TPR, true negative rate, TNR, false positive rate, FPR, positive predicted value, PPV, and accuracy, ACC, (Fig. 1 and Supplementary Table S2). In order to determine these measurements, we first determined true and false positives, as well as true and false negatives. A true positive is a false association in the network that is correctly removed by a method, and a false negative is a false association that is incorrectly not removed. A false positive is a true interaction in the network that is incorrectly removed by a method, and a true negative is a true interaction that correctly is not removed by a method. The ideal method maximizes true positives and true negatives and minimizes false positives and false negatives.

The intersection combination under-performed compared to each individual method, SP and OL perform best, and II performs better than DPI according to TPR and ACC (Fig. 1). However, applying each method individually has the drawback of removing more true interactions. On average, there are 60 (59 dwn) true interactions in the simulated networks. The individual methods removed 86% (85% dwn) (SP), 85% (84% dwn) (OL), 60% (51% dwn) (II), and 38% (28% dwn) (DPI). Therefore, although the intersection combination removed fewer edges, it outperformed the others according to the TNR because it eliminated fewer of the true interactions, 25% (16% dwn). All methods had high PPV values with half of all measured PPV above ≈ 0.95 . According to PPV, intersection combination performed best and SP and OL performed worst (Fig. 1).

Real data

After testing EnDED’s performance on simulated networks, we applied it to a real microbial association network, which was constructed from 10 years of monthly samples from January 2004 to December 2013 at the



Blanes Bay Microbial Observatory (BBMO) [41]. These samples included bacteria and eukaryotes of two size-fractions: picoplankton (0.2–3 μm) and nanoplankton (3–20 μm). We estimated community composition via metabarcoding of the 16S and 18S rRNA gene, and inferred an association network, hereafter referred to as BBMO network (see “Methods” section). The BBMO network contained 762 nodes including 754 ASVs (Amplicon Sequence Variants) and 8 environmental factors, and 30498 edges including 29820 microbial edges and 607 edges between a microorganism and an environmental factor. The network contained more positive (24458, 82.0%) than negative (5362, 18.0%) microbial associations (Fig. 2).

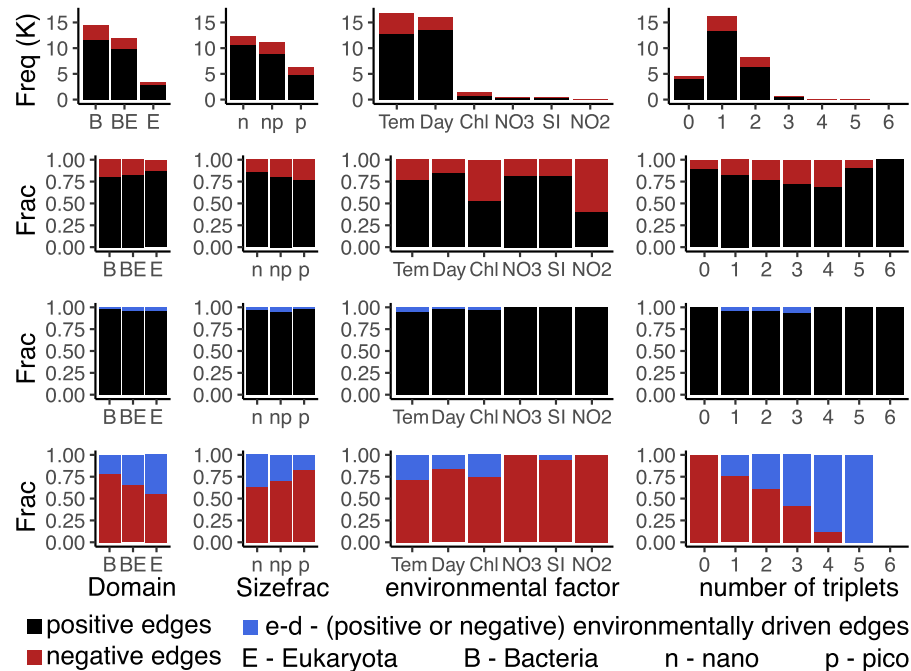
We found that 25230 (84.6%) of the network edges were in at least one and in maximum six environmental triplets (Fig. 2 and Supplementary Table S3). Overall, we detected 35166 environmental triplets within the BBMO network. Of the ten considered environmental factors, PO_4^{3-} and salinity were not associated to any microorganism in the network, and turbidity and NH_4^+ were not found within a triplet. Thus, six environmental factors remained: temperature (1831 environmentally driven edges were removed due to Temperature) and day length (652 removed edges) were the top two

environmental factors affecting microbial associations, followed by total chlorophyll (175), SiO_2 (5), and NO_3^- (1); no edge was removed due to NO_2^- .

The intersection combination removed 2488 ($\approx 8.3\%$) associations from the BBMO network. We classified and quantified these indirect edges according to the domain of the nodes (bacteria–eukaryotes, nanoplankton–picoplankton), environmental factor, and the number of triplets a microbial edge was in (Fig. 2 and Supplementary Table S4). Compared to the intersection combination, each method individually removed more edges: 84.6% (SP and OL removing all microbial edges present in a triplet), 25.7% (II), and 24.8% (DPI); that is, removal was 3 to 10 times larger.

We also determined for each association the Jaccard index, which indicates how often two microorganisms appear together in the dataset. We assume that two microbes that appear together $< 50\%$ of the time are less likely to have true contemporary ecological interactions and the corresponding association is more likely to be false. We found that only 27.7% of the indirect associations had a Jaccard index above 0.5 compared to 61.1% of the associations that were not indirect. This discrepancy is bigger for negative edges, with 1.2% above and 98.8% below 0.5 (Table 1). The fact that over 72.3% of

A) Classification and quantification of edges in the BBMO network



B) Location of specific edges in the BBMO network

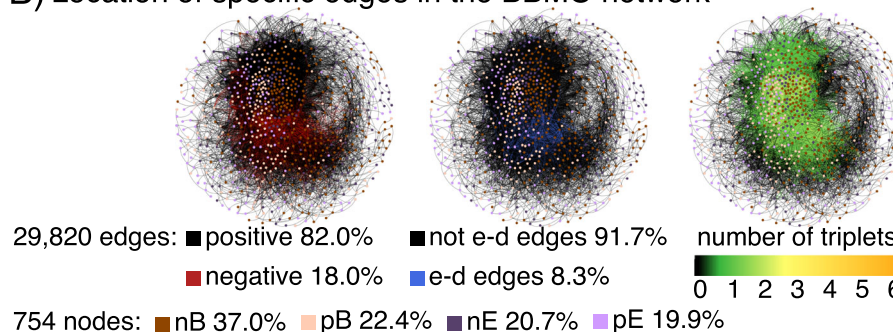


Fig. 2 Quantification of environmentally driven associations in the BBMO network **A** The first column shows the number (in thousands, K) or fraction of microbial associations divided by domain: bacteria–bacteria associations (B), bacteria–eukaryote associations (BE), and eukaryote–eukaryote associations (E). The second column shows the number (or fraction) of associations divided by size-fractions: association within the nano size fraction (n), within the pico size fraction (p), and between these two size fractions (np). The third column shows all microbial edges connected to an environmental parameter: temperature (Tem), day length (day), chlorophyll (Chl), inorganic nutrients NO_3^- (NO_3), SiO_2 (Si), and NO_2^- (NO_2). The last column shows the number of edges divided in how many triplets they have been found ranging from no triplets (0) to six triplets. The first two rows display the number of microbial associations of the BBMO network before applying EnDED. Positive associations are indicated with black, negative associations with red. The last two rows indicate in blue the fraction of environmentally driven edges among the positive (3rd row) and negative (4th row) microbial associations. **B** The left network shows in black the positive and in red the negative associations. The right network shows the number of triplets a microbial edge is in ranging from one (green) to six (orange), and no triplet (black). The middle network shows in blue the environmentally driven associations that were detected by the intersection combination of the four methods sign pattern, overlap, interaction information, and data processing inequality

environmentally driven associations have a Jaccard index equal or below 0.5 strengthens the decision of their removal.

The intersection combination removed more negative than positive edges, 1554 and 934, respectively (Fig. 2). However, there were 20334 positive and 4896 negative microbial associations that were found in at least one

environmental triplet, so the method removed 31.7% of the negative and only 4.6% of the positive edges. If we randomly removed 2488 edges, we would expect 18.0 % to be negative (i.e., 448) and 82.0% of them to be positive (i.e., 2040). If we restrict these calculations to the 25230 microbial associations that were found in at least one environmental triplet, with 20334 of them being

Table 1 Jaccard index of edges. The BBMO network before applying EnDED contained 29820 edges of which 2488 (8.3%) were environmentally driven (indirect). Considering the Jaccard index for these indirect edges, 688 (27.7% of indirect edges) score above 50%, and 1800 (72.3%) score below or equal to 50%. In contrast, 61.1% of edges not considered as indirect have a Jaccard index above 50%, and 38.9% of all not indirect edges have a Jaccard index equal or below 50%

| | All | Jaccard index > 50 | Jaccard index ≤ 50 |
|------------------------------------|----------------|--------------------|--------------------|
| BBMO network | 29 820 (100%) | 17 383 (58.3%) | 12 437 (41.7%) |
| Positive edges | 24 458 (82.0%) | 17 212 (70.4%) | 7 246 (29.6%) |
| Negative edges | 5 362 (18.0%) | 171 (3.2%) | 5 191 (96.8%) |
| Indirect (intersection) | 2 488 (8.3%) | 688 (27.7%) | 1 800 (72.3%) |
| Positive + indirect (intersection) | 934 (3.1%) | 670 (71.7%) | 264 (28.3%) |
| Negative + indirect (intersection) | 1 554 (5.2%) | 18 (1.2%) | 1 536 (98.8%) |
| Not indirect (all) | 27 332 (91.7%) | 16 695 (61.1%) | 10 637 (38.9%) |
| Not indirect (min 1 triplet) | 22 742 (76.3%) | 14 242 (62.6%) | 8 500 (37.4%) |
| Not indirect (no triplet) | 4 590 (15.4%) | 2 453 (53.4%) | 2 137 (46.6%) |
| Sign pattern | 25 230 (84.6%) | 14 930 (59.2%) | 10 300 (40.8%) |
| Overlap | 25 230 (84.6%) | 14 930 (59.2%) | 10 300 (40.8%) |
| Interaction information | 7 672 (25.7%) | 4 962 (64.7%) | 2 710 (35.3%) |
| Data processing inequality | 7 394 (24.8%) | 1 862 (25.2%) | 5 532 (74.8%) |

positive and 4896 being negative, we would expect to remove 19.4% (i.e., 483) of negative and 80.6% (i.e., 2005) of positive edges. The probability of randomly removing less positive than negative associations is nearly zero, since it follows a multivariate hypergeometric distribution:

$$P(k_{neg}, k_{pos}) = \frac{\binom{N_{neg}}{k_{neg}} \cdot \binom{N_{pos}}{k_{pos}}}{\binom{N}{n}}, \quad (1)$$

where N_{pos} and N_{neg} are the number of positive and negative associations in the network, respectively, k_{pos} is

the number of removed positive and k_{neg} the removed negative associations from the network, N is the number of associations in the network, and n is the number of removed associations from the network. The removal of more negative edges through intersection combination indicates that this removal was not random or, in other words, that negative associations are more likely to represent environmentally driven edges.

To evaluate the performance of EnDED on the BBMO network, we considered interactions described in literature and collected in the Protist Interaction Database (PIDA) [10]. Studies typically compare the associations of a network to those reported in the literature at the genus level [23]. The ambiguity in taxonomic

Table 2 Associations found in the BBMO network that have been reported in the literature. The table mentions whether or not the associations were removed or kept by EnDED. For example, the association between the ASVs classified as *Dia. Thalassiosira* and ASVs classified as *F. unknown Flavobacteriia* has been found 17 times in the network: 4 were removed and 13 were kept

| Microorganisms | EnDED | ID in PIDA |
|--|-----------|------------------------|
| Included in 1, 2, 3, or 4 triplets | | |
| <i>Dia. Thalassiosira</i> -Dino. <i>Heterocapsa</i> | 1 removed | 1665 |
| <i>Dia. Thalassiosira</i> - <i>F. unknown Flavobacteriia</i> | 4 removed | 2199 |
| | 13 kept | |
| Not included in a triplet | | |
| Dino. <i>Heterocapsa</i> -Dino. <i>Prorocentrum</i> | 1 kept | 1501, 1511 |
| Dino. <i>Gyrodinium</i> -Dino. <i>Heterocapsa</i> <i>Heterocapsa</i> | 1 kept | 1313, 1314, 1780, 1783 |
| Dino. <i>Prorocentrum</i> -Dino. <i>Gymnodinium</i> | 2 kept | 1499 |
| Dino. <i>Prorocentrum</i> -Dino. <i>Prorocentrum</i> | 4 kept | 1509, 1510 |
| Dino. <i>Prorocentrum</i> -Dino. <i>Scrippsiella</i> | 2 kept | 1513 |
| <i>F. unknown Flavobacteriia</i> - <i>Dia. Pseudo-nitzschia</i> | 1 kept | 2196 |

Abbreviations: *Dia* Diatomea, *Dino* Dinoflagellata, *F* Flavobacteriia, ID in PIDA refers to the number PIDA gave to an interaction described in the literature

classification and the large number of edges challenged this comparison. Thus, we implemented a function to compare strings and match the taxonomic classification of a microorganism in the BBMO network to those in the scientific literature (PIDA). We found that only 29 (0.1%) associations were supported by interactions described in the literature (Table 2). That is, 99.9% of associations in the BBMO network (before applying EnDED) could not be used to evaluate EnDED's performance. These 29 associations describe 8 unique interactions between 8 microorganisms, and 18 edges were in an environmental triplet to which each method as well as their combination were applied (summary in Table 2). Ideally none of these described associations should be removed by EnDED. Yet, the intersection combination removed five associations (Table 2). In contrast and even worse, SP and OL removed all 18 edges, II 8, and DPI 9 edges. The additionally removed edges by individual methods are associations between a diatom (*Thalassiosira*) and an unknown Flavobacteriia. Considering only the genus level, there were 171 unique genera in the BBMO network, and 700 in PIDA, combined there were 837 microbial genera, and 34 genera in both. Thus, 19.9% of the microbial genera found in the BBMO network were also in PIDA, and 4.9% of the genera found in PIDA were also found in the BBMO network.

Discussion

Using EnDED to disentangle environmental effects in microbial association networks

EnDED makes several indirect-edge removal techniques accessible to microbial ecologists without requiring previous programming experience. These techniques can be used individually or combined. In addition, this work systematically evaluates the different techniques and their combination to remove indirect edges from microbial association networks. Here, we tested only the union and intersection combination of all four methods, but other combination strategies are possible with EnDED. EnDED requires data of the environmental factors in order to predict if an association is environmentally driven. This is a limitation, since it may be impossible to consider all environmental factors [16]. However, EnDED can perform well if the major environmental factors, such as, e.g., temperature and nutrient concentrations for marine microorganisms, are provided. Moreover, knowledge of microbial interactions in nature is rather limited and therefore, determining the performance of EnDED for real networks is challenging and carries some degree of uncertainty. Thus, EnDED's results should be interpreted with care.

For the simulated networks, we found that each method individually removed on average a moderate to high number of edges. The intersection combination

removed fewer edges but kept more true interactions. To understand the impact of the environment, an increasing environmental influence was simulated, which was observed to be linked to a decrease in retrieving true interactions from inferred associations [21]. The observation holds for several network construction methods for cross-sectional data, including CoNet [42], SparCC [43], SPIEC-EASI [14], and Spearman correlations. In agreement with these findings, we observed a slight increase in retrieving true interactions when removing environmentally driven associations in our simulation networks.

In our BBMO dataset, the intersection combination removed a modest number of the edges—a much higher fraction of negative than positive edges. We argue that several negative associations are probably due to different environmental preference (different niches) of microorganisms. The Jaccard index representing a level of microbial co-occurrence scored equal or below 50% for most negative associations. These may partially represent microorganisms adapted to different seasons. Previous work on the eukaryotic pico- and nano-plankton at the BBMO, using the same basal 10-year dataset used here, indicated a strong seasonality at the community level [44].

Comparisons of indirect edge detection on other datasets

In our BBMO network, we found that the majority (84.6%) of the microbial edges was within at least one environmental triplet. This was 2.6 times higher than what was found for an association network inferred from data considering microorganisms and small metazoans from two ocean depths across 68 stations around the world and various size fractions (hereafter global interactome) [23]. This global interactome contains 29,912 (32.3%) edges that were within at least one environmental triplet [23]. In the previous study, 29,900 edges in the global interactome ($\approx 100\%$ of triplets and 32% of all edges) were attributed to environmental factors by SP, similarly to this study as SP removed all edges within triplets in the BBMO network. II indicated 11,043 environmentally driven edges in the global interactome ($\approx 37\%$ of triplets and 12% of all edges) with p value below 0.05 in a permutation test with 500 iterations. In comparison, II removed a higher fraction of edges in the BBMO network when considering all edges (25.7%), but less when considering within the triplets (30.4%). Network deconvolution suggested 22,439 environmentally driven edges ($\approx 75\%$ of triplets and 24% of all edges) within the global interactome, and the three methods agreed for 8209 edges ($\approx 27\%$ of triplets and 8.9% of all edges). In comparison, we detected slightly less environmentally driven associations for the BBMO network (8.3% of all edges). These differences suggest that a

higher environmental heterogeneity in the dataset may induce more indirect edges. Also, the effects of indirect dependencies may depend on dataset type (e.g., temporal vs. spatial). These possible differences and their effect on environmentally driven edges should be further investigated.

Using II for the BBMO network, we identified a moderate number of environmentally driven associations. DPI also identified a moderate number (24.8%, 29.3% when considering only triplets), whereas SP or OL identified a high number of environmentally driven edges (84.6%, 100% when considering only triplets). This indicates that SP and OL are strict and should be used in combination with other methods in an intersection approach.

In another study, the tool FlashWeave [15] predicted direct microbial interactions in the human microbiome using the Human Microbiome Project (HMP) dataset, including heterogeneous microbial abundance data of 68,818 samples [45]. The inferred networks (with and without metadata) were sparser than our networks. The network with metadata contained 10.7% fewer associations compared to the network without metadata, slightly more than in our results from BBMO.

Factors causing indirect microbial associations

From the simulated networks, we found that using the intersection combination instead of each method individually, we maintained more true interactions at the cost of more false associations in the network—more when considering simulated networks including noise. Comparing our simulated network with the BBMO network, the intersection combination classified a higher number of edges as environmentally driven in the simulated networks 32% (22% down) than in the BBMO network (8.3%). For the simulated data, we previously knew the environmental factor influencing pairwise microbial associations. For the BBMO data, we used ten available environmental factors, but not all factors that could affect microbial dynamics. Even though the most important factors influencing microbial seasonal dynamics at BBMO were considered [44], there are several factors that were not measured and that could generate indirect edges. The indirect edges associated to these factors were not detected in our analyses. Similarly, indirect edges associated to biotic interactions (e.g., two bacteria sharing a positive edge as they are symbionts in the same protists) were not considered. Future sampling for microbial interaction research should expand metadata collection in order to detect (more) abiotic and biotic factors that could generate indirect edges.

While temperature and day length (hours of light) were the top two environmental factors affecting microbial associations in the BBMO network, the most

important environmental factors in the global interactome [23] were phosphate concentration and temperature, followed by nitrite concentration and mixed-layer depth. Although we considered PO_4^{3-} and salinity, they were not associated to any microorganism in the network, which may reflect the low variation of these environmental factors in the studied marine site (BBMO). For instance, the standard deviation in the BBMO dataset was < 1 for PO_4^{3-} and salinity, in contrast to the global interactome dataset [23], where it was about 20–30 when considering all samples. During the Malaspina 2010 Circumnavigation Expedition, the concentrations of trace metals were determined for 110 surface water samples [46]. The previous study indicates relationships between primary productivity and trace nutrients, more specifically for the Indian Ocean Cd; the Atlantic Ocean Co, Fe, Cd, Cu, V, and Mo; and the Pacific Ocean Fe, Cd, and V. Thus, trace metals are further environmental factors that may play an important role in regulating oceanic primary productivity.

Limitations of EnDED

EnDED detects and removes environmentally driven indirect edges. However, its triplet analysis could be extended to remove indirect edges driven by taxa, as done with gene triplets [27]. A recent update of the network construction tool eLSA [11, 12] permits to examine how a factor, such as a microorganism or environmental variable, mediates the association of two other factors [47], which allows the study of interactions between three factors. Furthermore, triplets limit the study to first-order indirect dependencies, neglecting higher-order indirect dependencies. Such limitation was solved for the DPI method by examining associations in quadruplets, quintuplets, and sextuplets [48]. Implementing higher-order DPI and adjusting the other three methods to account for higher-order indirect dependencies may be promising but one needs to be aware that incorporating higher-order dependencies will also increase the risk of overfitting. Further, all relevant (measured) environmental factors could be incorporated into the calculation of II, which would combine environmental triplets. However, we reason that such adjustments would require a larger sample size. Both II and DPI calculate MI that measures the dependence between two random variables. EnDED is limited by including one function to estimate the MI. A comparison of four different MI estimates revealed that obtaining the true value of MI is not straightforward, and minor variations of assumptions yield different estimates [49]. Lastly, the conditional mutual information, CMI, which quantifies non-linear direct relationships among variables, can be underestimated if variables have tight associations in a network [50]. The so-called part mutual information, PMI, measurement

can help overcome CMI's underestimations. Although using PMI instead of CMI looks promising, calculating PMI is computationally more demanding [50].

Future perspectives

In this study, we have shown that EnDED with an intersection combination approach provides less dense networks, but still with many potential interactions. We observed a trade-off comparing single methods with the combination approach (intersection combination). Although the latter kept more true interactions, it kept also more false associations. Inferring emergent properties is a key task in microbial ecology to characterize microbial ecosystems from a network perspective. Thus, if the study aim is to explore patterns of network topology rather than single edges, inferring a network comparable to the real interaction network may be more useful than accuracy of single edges. However, investigations aiming to provide potential interaction partners may use EnDED with the intersection combination approach (e.g., [51]). Specific associations may be validated with experiments or microscopy [6, 23]. However, we suggest to first further reduce the set of potential interaction hypotheses. To improve the selection of interaction hypotheses, we propose to score associations based on re-occurrence: in time, as done with microbial abundance seasonality [44], or space, where an association appears in different networks based on different datasets, or different regions of the world. In a previous study using 313 samples, including 7 size-fractions, 4 domains (Archaea, Bacteria, Eukarya, and viruses), and 2 depths from 68 stations across 8 oceanic provinces, 14% of the 81,590 predicted biotic interactions were identified as local [23]. Thus, re-occurent associations may suggest a higher likelihood that the association represents a true ecological interaction, reducing the number of interaction hypotheses to the strongest ones. Another strategy to shortlist interaction hypotheses is to incorporate additional data into the network and use a multi-layer network approach. Such data could be environmental preferences such as temperature or salinity optima, size of cells, presence of chloroplasts, or data obtained from high-throughput cultivation [52], microbial community transcriptomes that reveal metabolic pathways [53], or interactions inferred from single-cell genome data [6, 54].

Conclusion

In this paper, we present EnDED, an analysis tool to reduce the number of environmentally induced indirect edges in inferred microbial networks. Applying EnDED on simulated networks indicated that false associations, driven by environmental variables instead of true interactions, were ubiquitous. However, EnDED's intersection

combination classified a minority of associations as environmentally driven in a real (BBMO) network. Depending on the single method used, we classified a moderate to high number of associations as environmentally driven in the same network. Nevertheless, associations driven by environmental factors must be determined and quantified to generate more accurate insights regarding true microbial interactions. EnDED provides a step forward in this direction.

Methods

Simulated dataset: time series based on an adjusted generalized Lotka-Volterra model

To evaluate the performance of EnDED, we simulated a time series using an adjusted version of the standard *generalized Lotka-Volterra model*, gLV [55, 56]. The gLV models the dynamics of microbial communities assuming that it is well described by pair-wise interactions. The model's simplicity arises from the assumption of linear interactions, which facilitates implementation and allows fast numerical simulations. The gLV has, however, several limitations [57]. For example, gLV neglects higher-order interactions and the additivity of interaction strengths is a weakness because they may be combined in different ways. Also, interaction strengths are often assumed to be constant parameters, but a reducing level of a nutrient may weaken cross-feeding relationships. Moreover, gLV omits the influence of environmental factors, which, for example, can induce oscillations in natural communities [58]. Using a model that accounts for nutrients [59] is more realistic but also more complex. More elaborate mechanistic models of microbial dynamics than gLV solve explicitly the global cycling of nutrients and are coupled to the oceanic circulation (see [60] for a review), but the added complexity can hamper understanding about the ecological interactions among microorganisms when compared to a simpler gLV approach. Thus, we chose to use a simpler extension of the gLV to account for the influence of environmental factors [61, 62]. In order to allow the growth rates to vary when the environmental variables change, environmental variables can be incorporated directly into the gLV [21, 61]. We simulated a time series using the Klemm-Eguíluz algorithm [63], and an adjusted gLV. We adjusted the model by defining microbial growth rates as a function dependent on one seasonal abiotic environmental factor, and added an abiotic environmental factor in the interaction matrix. We then used the time series generated by the gLV to obtain temporal microbial abundance data. With this simulated data, we inferred a network that contained environmentally driven associations, needed to evaluate the performance of EnDED. We repeated this procedure 1000 times to obtain a large set of simulated networks, and then used the determined

abundance tables and Poisson distribution to obtain another 1000 simulated networks including noise. The addition of noise was done by randomly drawing an abundance from the Poisson distribution with λ equaling the original abundance of a specific microorganisms to a specific time.

Adjusting the gLV

To evaluate EnDED, we simulated a time series of microbial abundances with a gLV including true pairwise interactions between 50 microorganisms and adjusted it by incorporating two environmental factors:

$$\frac{dy(t)}{dt} = y(t)[b + Ay(t)] \quad (2)$$

where t is time, $dy(t)/dt$ is the rate of change of microbial abundances as a column vector, $y(t)$ is the vector of microbial abundance at time t , b is the growth rate vector determined through microorganisms' specific growth rate functions that depend on an environmental factor (see Eq. (4)), and A is the interaction matrix.

Interaction matrix

In the interaction matrix A , each coefficient a_{ji} provides the linear effect that a change in the abundance of microorganism i has on the growth of microorganism j [64]. We simulated the interaction coefficients a_{ji} with the Klemm-Eguíluz algorithm [63], which generates a modular and scale-free matrix. To simulate interaction coefficients, we set the interaction probability to 0.01, the percentage of positive coefficients to 30%, and diagonal coefficients to 0. Negative diagonal coefficients a_{ji} (i.e., the interaction of a microorganism with itself) can represent intra-specific competition and provides the carrying capacity for each microorganism, preventing its explosive growth [65]. Then, after simulating interaction coefficients, we set the diagonal coefficients $a_{ii} = -0.5$ to avoid excessive microbial abundances in the simulations.

Two abiotic environmental factors

We adjusted the gLV by including two environmental factors. For simplicity, we assume no feedback between the microorganisms and the environmental factors. That is, the environmental factors affect the growth of the microorganisms but not vice-versa. The first environmental factor affects the specific growth rate of each microorganism by interacting with two of their traits: optimal environmental value for growth and tolerance range of environmental values. We simulated the environmental factor using a periodic sinusoidal function (see Eq. (3)), rounded to 3 digits:

$$\epsilon(t) \triangleq \text{round}(\sin(\omega \bullet t), \text{digits} = 3) \quad (3)$$

where t is the time axis (months), $\omega = (-2\pi/T)$ is the signal frequency (radians) and $T = 12$ is the signal periodicity (months); resulting in a signal phase shift of $T/4$ (months). While the first environmental factor is considered to be “external” to the microbial community, the second environmental factor is considered to be “internal”, and therefore, it is included in the interaction matrix. The interaction coefficients between the microorganisms and the second environmental factor were generated by splitting the microorganisms into two groups: the second abiotic environmental factor influenced positively one half and negatively the other half of the microorganisms. We obtained the interaction coefficients from two uniform distributions defined to range between $[-0.8, -0.2]$ and $[0.2, 0.8]$ respectively. As the microorganisms did not influence the abiotic factor, the corresponding interaction coefficients were set to 0.

Species growth rate

The external seasonal abiotic environmental variable affects the growth rate, g , of each microorganism. This dependency is given :

$$g(t) \triangleq g_{\max}^2 \exp\left(-\frac{1}{2} \frac{(\epsilon_{\text{opt}} - \epsilon(t))^2}{\sigma^2}\right), \quad (4)$$

where $E(t)$ is the environmental parameter that affects the microorganism's growth rate $g(t)$ at time t , g_{\max} is the microorganism's specific maximum growth rate that determines the amplitude of the growth-rate curve, ϵ_{opt} is the microorganism's specific optimal environmental value that determines the peak of the growth-rate curve, and σ is the microorganism's specific ecological tolerance (niche width) determining the environmental range in which the microorganism grows, which determines the length (niche spread) of the growth-rate curve. We obtained the two constant parameters g_{\max} and σ for each microorganism from a uniform distribution ranging between 0.3 and 1 to assure positive values. The values ϵ_{opt} were drawn from a uniform distribution ranging between the minimal and maximal value of the seasonal environmental factor. We defined the internal abiotic environmental factor, which is included in the interaction matrix, through the same function with $g_{\max} = 0.8$, $\epsilon_{\text{opt}} = 0.5$, and $\sigma = 0.5$. Since the growth rates depend on the environmental factor, they vary seasonally. Different microorganisms will grow better or worse at different times of the year following their environmental niches. This will lead to an asynchrony of their growth rate responses to the environment that will translate into an asynchrony of their abundances in time.

Initial abundances

To obtain the microbial abundances in time with the adjusted gLV, we simulated the initial microbial abundances with a stick-breaking process such that abundances add up to 1, using the function `bstick` [66, 67], and the package `vegan` [68]. We generated uneven initial microbial abundances without introducing zeros and set the initial value for the internal abiotic environmental factor included in the interaction matrix to 0.001.

Species abundances in time

Once we have set the initial conditions, we simulated microbial abundances over time by solving the equations given in the adjusted gLV (see Eq. (2)). Start time was 0, end time 49.5, and sample resolution 0.5 resulting in 100 samples. We used the solver function `lsoda` [69]. The simulated abundances in time were used to construct an association network, which is referred to as the simulated network.

Real dataset: Blanes Bay Microbial Observatory (BBMO) time series

Microbial abundances

Surface water (≈ 1 m depth) was sampled monthly from January 2004 to December 2013, at the Blanes Bay Microbial Observatory (BBMO; <http://bbmo.icm.csic.es>) in the North-Western Mediterranean Sea $41^{\circ}40' \text{ N } 2^{\circ}48' \text{ E}$ [41]. About 6 L of seawater were filtered and separated into picoplankton (0.2–3 μm) and nanoplankton (3–20 μm), as described in [44]. The DNA was extracted using a phenol-chloroform standard method [70], which has been modified by using Amicon units (Millipore) for purification.

Next, community DNA was extracted, and the 18S ribosomal RNA-gene (V4 region) was amplified in [44] using the primer pair TAREukFWD1 and TAREukREV3 [71]. The 16S ribosomal RNA-gene (V4 region) was also amplified from the same DNA extracts using the primers Bakt 341F [72] and 806R [73]. Amplicons were sequenced in a MiSeq platform (2×250 bp) at the sequencing service RTL Genomics in Lubbock, TX, USA. Read quality control, trimming, and inference of operational taxonomic units (OTUs) as amplicon sequence variants (ASV) was made with DADA2 v1.10.1 [74] with the maximum number of expected errors (MaxEE), set to 2 and 4 for the forward and reverse reads, respectively.

ASV sequence abundance tables were obtained for both microbial eukaryotes and prokaryotes. We subsampled both tables to the lowest sequencing depth of 4907 reads, with the `rrarefy` function from the `Vegan` package in R [68], v2.4–2. We excluded 29 nanoplankton samples (March 2004, February 2005, and May 2010 to July

2012) featuring suboptimal amplicon sequencing. In these, we estimated microbial abundances using seasonally aware missing value imputation by weighted moving average for time series as implemented in the R package `imputeTS` [75], v2.8.

Dislodging cells or particles and filter clogging can bias the collection of DNA in either small or large organismal size fractions. To reduce the bias, we divided the sequence abundance sum of the nanoplankton by the picoplankton for each ASV appearing in both size fractions and set the picoplankton abundances to 0 if the ratio exceeded 2. Likewise, we set the nanoplankton abundances to 0 if the ratio was below 0.5.

Taxonomic classification

The taxonomic classification of each ASV was inferred with the naïve Bayesian classifier method [76] together with the SILVA version 132 [77] database as implemented in DADA2 [74]. In addition, eukaryotic microorganisms were BLASTed [78] against the Protist Ribosomal Reference database [PR2, version 4.10.0 [79];]. If the taxonomic assignment for eukaryotes disagreed between SILVA and PR2, we used the PR2 classification. We removed microorganisms identified as either Metazoa or Streptophyta, plastids, and mitochondria. In addition, we removed Archaea since the 341F primer is not optimal for recovering this domain [80]. The resulting microbial relative abundance table contained microbial eukaryotic and bacterial ASVs. Rare ASVs were removed, i.e., we kept only ASVs present in more than 15% of the samples and with a sequence abundance sum above 100.

Environmental factors

We measured environmental factors that may affect the ecosystem's dynamics. We considered a total of ten contextual abiotic and biotic variables: day length (hours of light), temperature ($^{\circ}\text{C}$), turbidity (Secchi depth m), salinity, total chlorophyll ($\mu\text{g/l}$), and inorganic nutrients— PO_4^{3-} (μM), NH_4^+ (μM), NO_2^- (μM), NO_3^- (μM), and SiO_2 (μM) [44]. Water temperature and salinity were sampled in situ with a CTD (conductivity, temperature, and depth) measuring device. Inorganic nutrients were measured with an Alliance Evolution II autoanalyzer [81]. See [41] for specific details on how other variables were measured.

Network construction

We constructed association networks from the simulated and the real microbial abundance tables and environmental parameters using eLSA [11, 12]. We included default normalization and a z -score transformation using median and median absolute deviation. We estimated the p value with a mixed approach that performs a

random permutation test if the theoretical p values for the comparison are below 0.05; the number of iterations was 2000. Although we are aware of time-delayed interactions and that eLSA [11, 12] could account for them, we considered our sampling interval as too large (1 month) for inferring time-delayed associations with a solid ecological basis. Thus, in our study, we focused on contemporary interactions between co-occurring microbes. For the BBMO dataset, the Bonferroni false discovery rate, q , was calculated for all edges from the p values using the R function `p.adjust` [40]. Lastly, we used a significance threshold for the p and q value of 0.001 as suggested in other works [22].

Intersection combination of EnDED—environmentally driven edge detection methods

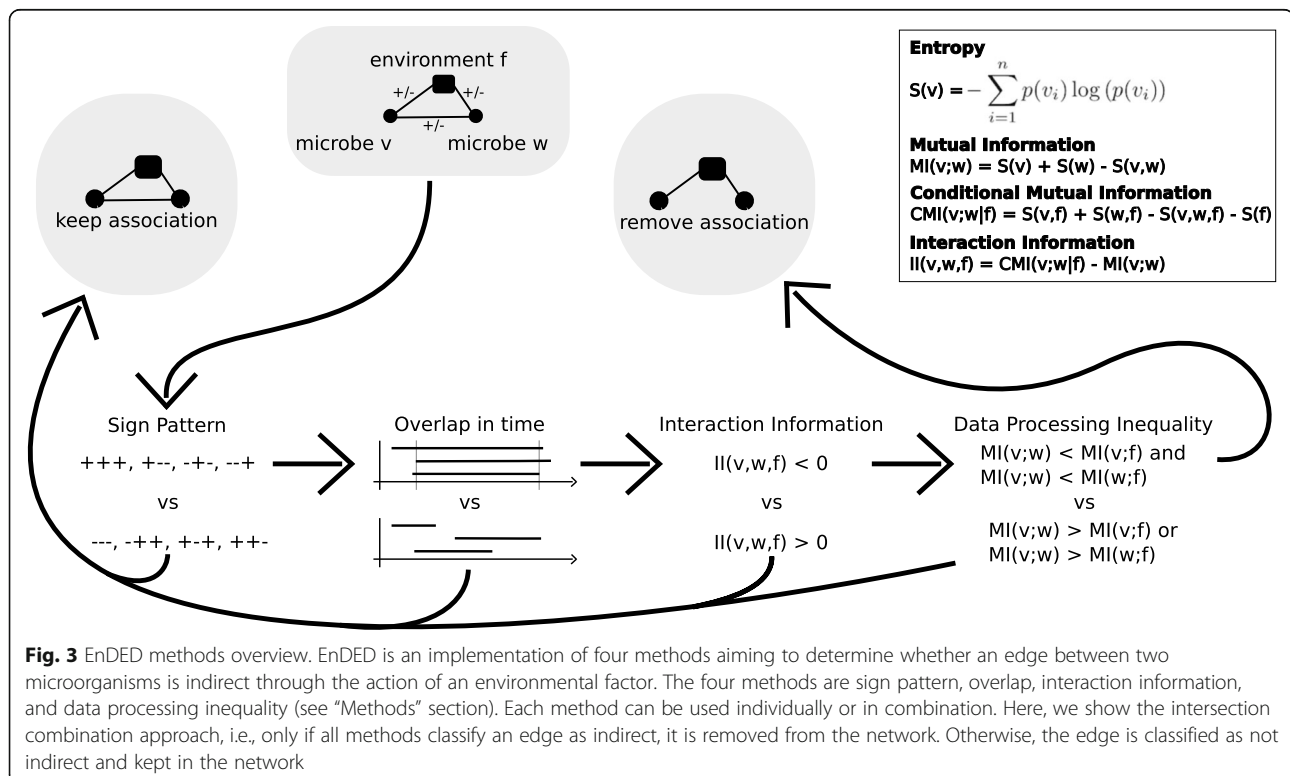
EnDED includes four methods: SP, OL, II, DPI (described below), and their intersection combination (an ensemble approach of the four methods). We applied these methods to find environmentally driven associations of microorganisms that were within an environmental triplet, as in [23]. An environmental triplet is a special case of a closed triplet where one of the nodes corresponds to an environmental factor and the other two nodes correspond to microorganisms. We define the closed triplets, where there is an edge between each pair of three nodes, as $T = \{v, w, f\}$ where v and w are two

microorganisms, and f is an environmental component (see Fig. 3).

For the intersection combination, all four individual methods must converge to the same solution, i.e., if all methods classify the microbial edge as environmentally driven, the edge is removed from the network. If a microbial association is within several environmental triplets, at least one of them must indicate the association as environmentally driven. In sum, the intersection combination retains an association in the network if no triplet classifies the association as environmentally driven.

Sign pattern

The SP method [23] filters environmentally driven edges from a network in which a positive association score indicates co-occurrence, and a negative association score indicates mutual exclusion. Let s_{vw} be the sign of the association score of the association between v and w (i.e., $s_{vw} = +$ or $s_{vw} = -$). A closed triplet T has eight SP combinations that group into two sets (see Fig. 3). If the product of the three association scores is positive, then the SP suggests that the edge between the two microorganisms is environmentally driven. Otherwise, if the product of the three association scores is negative, SP does not suggest that the association is environmentally driven.



Overlap

We have developed the OL method to support the SP for temporal data: a microbial edge should be disregarded as environmentally driven when the associations are misaligned in time. Thus, OL requires the time when the association begins as well as how long the associations lasts, i.e., duration or length of association in time, both determined by the network construction tool eLSA [11, 12]. Given an association between v and w , let b_{vw}^v be the beginning of the association for v , b_{vw}^w the beginning of the association for w , and d_{vw} be the duration of the association between v and w . Although not used in the BBMO network, OL can consider time-delays by assuming that the beginning of the association is the minimum of the two beginnings, $b_{vw} = \min(b_{vw}^v, b_{vw}^w)$, and the end of the association is the maximum, $e_{vw} = \max(b_{vw}^v + d_{vw}, b_{vw}^w + d_{vw})$. We indicate two microorganisms with v and w , and the factor by f . The OL method calculates the overlap O of the microbial association with the two microorganism-environment associations through Eq. (5). As depicted in Fig. 3, if $O > 60\%$, the microbial association is considered environmentally driven.

$$O = 100 \frac{\min(e_{vw}, e_{vf}, e_{wf}) - \max(b_{vw}, b_{vf}, b_{wf})}{e_{vw} - b_{vw}} \quad (5)$$

Mutual information and conditional mutual information

The method II employs two measurements: MI and CMI. The former is also used by DPI. Thus, before describing the methods, we first describe the two measurements. MI is a measure of the degree of statistical dependency between two variables [27]. We first consider $\mathbf{v} = v_1, \dots, v_n$, $\mathbf{w} = w_1, \dots, w_n$, and $\mathbf{f} = f_1, \dots, f_n$ as discrete random variables. The marginal probability of each discrete state (value) of the variable is denoted by $p(v_i) = P(\mathbf{v} = v_i)$, the joint probability by $p(v_i, w_j)$, and $p(v_i, w_j, f_k)$, and the conditional probability by $p(v_i|f_k)$, and $p(v_i, w_j|f_k)$. To obtain MI, we calculate the entropy of \mathbf{v} as

$$S(\mathbf{v}) = - \sum_{i=1}^n p(v_i) \log(p(v_i)), \quad (6)$$

and the joint entropy of \mathbf{v} and \mathbf{w} as

$$S(\mathbf{v}, \mathbf{w}) = - \sum_{i=1, j=1}^n p(v_i, w_j) \log(p(v_i, w_j)), \quad (7)$$

using the natural logarithm. The MI of \mathbf{v} and \mathbf{w} is defined through the sum of their entropies subtracted by their joint entropy:

$$\text{MI}(\mathbf{v}; \mathbf{w}) = S(\mathbf{v}) + S(\mathbf{w}) - S(\mathbf{v}, \mathbf{w}) \quad (8)$$

$$= \sum_{i=1}^n \sum_{j=1}^n p(v_i, w_j) \log\left(\frac{p(v_i, w_j)}{p(v_i)p(w_j)}\right), \quad (9)$$

with marginal probabilities $p(v_i) = \sum_{j=1}^n p(v_i, w_j)$, and $p(w_j) = \sum_{i=1}^n p(v_i, w_j)$.

The measurement CMI is the expected value of the MI of two random variables given a third random variable. It is defined as

$$\text{CMI}(\mathbf{v}; \mathbf{w}|\mathbf{f}) = S(\mathbf{v}, \mathbf{f}) + S(\mathbf{w}, \mathbf{f}) - S(\mathbf{v}, \mathbf{w}, \mathbf{f}) - S(\mathbf{f}) \quad (10)$$

$$= \sum_{k=1}^n p(f_k) \sum_{i=1}^n \sum_{j=1}^n p(v_i, w_j|f_k) \log\left(\frac{p(v_i, w_j|f_k)}{p(v_i|f_k)p(w_j|f_k)}\right) \quad (11)$$

$$= \sum_{k=1}^n \sum_{i=1}^n \sum_{j=1}^n p(v_i, w_j, f_k) \log\left(\frac{p(f_k)p(v_i, w_j, f_k)}{p(v_i, f_k)p(w_j, f_k)}\right)$$

Interaction information

The II is calculated with microbial abundance and environmental data. In this study, as in [23], II is computed as the difference of the CMI and MI:

$$\text{II} = \text{CMI} - \text{MI} \quad (12)$$

In other works [37], the II is defined with a different sign convention: $\text{II} = \text{MI} - \text{CMI}$. In our study, if II is positive, the method suggests that the microbial association is not environmentally driven. If II is negative, there is an environmentally driven association within the closed triplet. However, this method cannot detect which of the three associations is indirect. In other works [23], the microbial association is assumed to be environmentally driven if II is negative, but here we suggest to combine it with DPI (see below).

Significance of interaction information

We determined the significance of II following a strategy from [82, 83]. We used a parameter-free permutation test and computed the p value by randomizing the environmental vector \mathbf{f} . Since the MI is independent of the environmental factor and therefore remains constant, the significance of the II is the same as the CMI. Thus, we determined the significance of CMI with 1000 permutations: we randomized the environmental vector \mathbf{f} and recalculated the CMI 1000×, obtaining a CMI_i with $i \in \{1, \dots, 1000\}$. Afterwards, we quantified with c how many random CMI_i were at least as small as the original CMI_i : $c = |\{i : \text{CMI}_i \leq \text{CMI}_{\text{original}}\}|$. We calculated the p value as

$$p = \frac{c + 1}{1000 + 1} \quad (13)$$

Data processing inequality

As mentioned above, the II method can detect if an indirect association exists within a triplet but cannot determine which of the three associations is indirect. Thus, we added DPI to EnDED. DPI states that if two components v and w interact only through a third component f (i.e., in a network v and w are connected through a path containing f and there is no alternative path between v and w), then the MI of v and w , $MI(v; w)$ is smaller than $MI(v; f)$ and $MI(w; f)$ [38]:

$$MI(v; w) \leq \min\{MI(v; f), MI(w; f)\} \quad (14)$$

While DPI has been used in previous works on gene triplets [27], we used the DPI method for environmental triplets. We compared the MI between the two microorganisms with the MI between a microorganism and the environmental factor. If the MI between the microorganisms is the smallest, then the method suggests that the edge is environmentally driven. This method complements the II method.

Equal width discretization

To compute the MI, CMI, and subsequently II, we discretized the abundance data and environmental parameters. EnDED uses the equal width discretization algorithm, which creates equal sized ranges (also called bins or buckets) for an abundance vector $\mathbf{v} = (v_1, \dots, v_n)$ between the lowest value (v_{\min}) and highest value (v_{\max}). It is a procedure implemented in other works [84]. Given vector \mathbf{v} of length n (that is sample size) and number of bins $|B| = \lfloor \sqrt{n} \rfloor$, the discretized value v_d of variable v in vector \mathbf{v} :

$$v_d = \left\lceil \frac{(v - v_{\min}) \cdot |B|}{v_{\max} - v_{\min}} \right\rceil. \quad (15)$$

This equation assumes positive values. However, if \mathbf{v} contains negative values, $v_{\min} < 0$, we adjust Eq. (15) by substituting v_{\max} for $v'_{\max} = v_{\max} - v_{\min}$. This method does not fill in missing values, and it is limited by the presence of outliers as most values would go within the same bin. We can solve this problem with a different discretization method (where bins have the same number of elements) but we have not implemented it in the current version of EnDED.

Applying EnDED to networks constructed from simulated and real data

We applied EnDED to association networks constructed from time series of simulated abundances and estimated

microbial abundances from sequence data. The simulated networks were based on a gLV, while the real network was based on data from the BBMO. For the methods II and DPI, we also included the corresponding abundance tables, and environmental factors. EnDED was run with the OL threshold of 60%. We set the significance threshold for the II score to 0.05 and used 1000 iterations.

Evaluation of EnDED's performance

Simulated network

We evaluated EnDED with the simulated interaction matrices, which revealed the number of true positives (TP), true negatives (TN), false negatives (FN), and false positives (FP) before and after removing associations that were classified as environmentally driven. We assumed that associations not present in the interaction matrices, are environmentally driven. We consider P as the number of all false associations, both true positive and false negative detected environmentally driven edges: $P = TP + FN$, and N as the number of all true interactions, i.e., all true negative and false positive detected environmentally driven edges: $N = TN + FP$. Then, we calculated the true positive rate (sensitivity), by dividing the number of true positives by the number of all real positives: $TPR = (TP)/(P)$. Equivalently, we can also calculate the true negative rate (specificity) by dividing the number of true negatives by the number of all real negatives, $TNR = (TN)/(N)$. The false positive rate (fall out) is the complementary to TNR, i.e., $FPR = 1 - TNR$. The positive predictive value (precision) can be calculated by dividing the number of true positives by the sum of all predicted positives, $PPV = (TP)/(TP + FP)$. The accuracy is calculated by dividing the sum of true positives and true negatives by the sum of all real positives and real negatives, $ACC = (TP + TN)/(P + N)$.

Real dataset

Literature based database The real network evaluation is limited since the true interactions and the microorganisms that do not interact with each other are poorly known. We assessed true interactions known in the literature based on the genus, which are compiled within the Protist Interaction Database, PIDA [10]. On 15 October 2019, PIDA contained 2448 interactions. Although our dataset contains protists as well as bacteria, we were unable to evaluate interactions between bacteria through PIDA.

Jaccard index In ecology, the Jaccard index (Jaccard similarity coefficient) is often used for communities. Here, for each pair of microorganisms in the BBMO network, we computed the Jaccard index as the number of

samples in which both microorganisms occur, divided by the number of samples in which at least one of the two microorganisms is present.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40168-021-01141-7>.

Additional file 1: Table S1: Comparison between methods on correctly detecting false associations. We computed the fraction (in percentage) of correctly detected false associations for each of the 1000 simulated datasets. There are only few edges that are detected by only one approach (first four rows). The most prominent groupings are highlighted in grey, e.g., SP, OL, and II agree on average on a third of edges. Less prominent groupings are aggregated with others. **Table S2:** Performance of environmentally-driven edge detection methods on simulated networks. These include 50 microorganisms and 1225 possible associations. Values display median (standard deviation) for simulated networks and simulated networks incorporating noise. Combi refers to intersection combination of all four methods. The methods with highest or lowest median, respectively, are indicated with an asterisk. **Table S3:** Number of triplets an microbial edge is part of in the BBMO network. SP and OL not listed below because they remove 100% of microbial associations that are within at least one triplet. **Table S4:** The BBMO network based on real data. It contained bacteria and eukaryotes from the picoplankton and nanoplankton. This table summarizes the number and fractions of microbial associations classified by EnDED as environmentally-driven. Combi refers to the intersection combination of all four methods, II to Interaction Information, and DPI to Data Processing Inequality. Both methods, sign Pattern and Overlap, are not shown because both remove all microbial edges found in at least one triplet. For example, 349 (14.9%) associations between bacteria from the picoplankton with eukaryotes from the nanoplankton were classified by intersection combination as environmentally-driven (indirect), II classified 30.6% and DPI 37.2% as environmentally-driven.

Acknowledgements

We thank all members of the Blanes Bay Microbial Observatory sampling team and the multiple projects funding this collaborative effort over the years. We also thank collaborators at www.thepapermill.eu for help with critical reading in the early stages of the manuscript. Part of the analyses have been performed at the Marbits bioinformatics core at ICM-CSIC (<https://marbits.icm.csic.es>).

Authors' contributions

IMD, GLM, JR, KF, and RL designed and conceived the project. IMD performed data analysis, data simulation, and implementation of EnDED. IMD received substantial feedback on established indirect detection methods from GLM and KF, on data simulation from SMV and KF, on network construction from AKK, and on evaluation of EnDED from AKK (literature based database for real dataset), GLM and KF (measurements for simulation dataset). RL processed the amplicon data from BBMO generating OTU tables. AKK ran the eLSA network construction tool for the BBMO dataset and IMD ran the tool for the simulation datasets. RL provided funding for the project. The original draft was written by IMD. IMD, GLM, AKK, SMV, KF, and RL contributed substantially to manuscript revisions. All authors approved the final version of the manuscript.

Funding

This project and IMD received funding from the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement no. 675752 (SingeK: <http://www.singeke.eu>). RL was supported by a Ramón y Cajal fellowship (RYC-2013-12554, MINECO, Spain). This work was also supported by the projects INTERACTOMICS (CTM2015-69936-P, MINECO, Spain), MINIME (PID2019-105775RB-I00, AEI, Spain) and MicroEcoSystems (240904, RCN, Norway) to RL. We thank the CSIC Open Access Publication Support Initiative through the Unit of Information Resources for Research (URIC) for helping to cover publication fees.

Availability of data and materials

EnDED is publicly available: <https://github.com/InaMariaDeutschmann/EnDED>. This repository contains the file "FromDataSimulationToEvaluatingEnDED.RMD", which contains R code to generate simulated abundance tables, commands to run eLSA network construction and EnDED, as well as the command to run a C++ program (included as well) and R code used for evaluation. The repository folder BBMO data contains the BBMO abundance table, the taxonomic classification table, and the BBMO network including results of EnDED.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Institute of Marine Sciences, CSIC, Passeig Marítim de la Barceloneta, 37-49, 08003 Barcelona, Spain. ²Research Unit in Biology of Microorganisms (URBM), University of Namur, 61 Rue de Bruxelles, 5000 Namur, Belgium. ³Department of Biosciences/Section for Genetics and Evolutionary Biology (EVOGENE), University of Oslo, p.b. 1066 Blindern, N-0316 Oslo, Norway. ⁴VIB Center for Microbiology, Herestraat 49-1028, 3000 Leuven, Belgium. ⁵KU Leuven Department of Microbiology, Immunology and Transplantation, Rega Institute, Laboratory of Molecular Bacteriology, Herestraat 49, 3000 Leuven, Belgium. ⁶Spanish Institute of Oceanography (IEO - CSIC), Ave Príncipe de Asturias 70 Bis, 33212 Gijón, Spain.

Received: 11 August 2020 Accepted: 20 July 2021

Published online: 26 November 2021

References

- Locey KJ, Lennon JT. Scaling laws predict global microbial diversity. *Proc Natl Acad Sci*. 2016;113:5970–5.
- Kallmeyer J, Pockalny R, Adhikari RR, Smith DC, D'Hondt S. Global distribution of microbial abundance and biomass in subseafloor sediment. *Proc Natl Acad Sci*. 2012;109:16213–6.
- Whitman WB, Coleman DC, Wiebe WJ. Prokaryotes: the unseen majority. *Proc Natl Acad Sci*. 1998;95:6578–83.
- Falkowski PG, Fenchel T, DeLong EF. The microbial engines that drive Earth's biogeochemical cycles. *Science*. 2008;320:1034–9.
- DeLong EF. The microbial ocean from genomes to biomes. *Nature*. 2009;459:200–6.
- Krabberød AK, Bjørnbækmo MFM, Shalchian-Tabrizi K, Logares R. Exploring the oceanic microeukaryotic interactome with metaomics approaches. *Aquat Microb Ecol*. 2017;79:1–12.
- Worden AZ, Follows MJ, Giovannoni SJ, Wilken S, Zimmerman AE, Keeling PJ. Rethinking the marine carbon cycle: Factoring in the multifarious lifestyles of microbes. *Science*. 2015;347. <https://doi.org/10.1126/science.1257594>.
- Faust K, Raes J. Microbial interactions: from networks to models. *Nat Rev Microbiol*. 2012;10:538–50.
- Layeghifard M, Hwang DM, Guttman DS. Disentangling interactions in the microbiome: a network perspective. *Trends Microbiol*. 2017;25:217–28.
- Bjørnbækmo MFM, Evenstad A, Røsaeg LL, Krabberød AK, Logares R. The planktonic protist interactome: where do we stand after a century of research? *ISME J*. 2019. <https://doi.org/10.1038/s41396-019-0542-5>.
- Xia LC, Steele JA, Cram JA, Cardon ZG, Simmons SL, Vallino JJ, et al. Extended local similarity analysis (eLSA) of microbial community and other time series data with replicates. *BMC Syst Biol*. 2011;5:S15.
- Xia LC, Ai D, Cram J, Fuhrman JA, Sun F. Efficient statistical significance approximation for local similarity analysis of high-throughput time series data. *Bioinformatics*. 2013;29:230–7.
- Faust K, Raes J. CoNet app: inference of biological association networks using Cytoscape [version 2; peer review: 2 approved]. *F1000Res*. 2016;5. <https://doi.org/10.12688/f1000research.9050.2>.

14. Kurtz ZD, Müller CL, Miraldi ER, Littman DR, Blaser MJ, Bonneau RA. Sparse and compositionally robust inference of microbial ecological networks. *PLoS Comput Biol*. 2015;11:1–25.
15. Tackmann J, Rodrigues JFM, von Mering C. Rapid inference of direct interactions in large-scale ecological networks from heterogeneous microbial sequencing data. *Cell Syst*. 2019;9:286–296.e8.
16. Lv X, Zhao K, Xue R, Liu Y, Xu J, Ma B. Strengthening insights in microbial ecological networks from theory to applications. *mSystems*. 2019;4:e00124–19.
17. Li C, Lim KMK, Chng KR, Nagarajan N. Predicting microbial interactions through computational approaches. *Methods*. 2016;102:12–9.
18. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are compositional: and this is not optional. *Front Microbiol*. 2017;8:2224.
19. Aitchison J. A new approach to null correlations of proportions. *J Int Assoc Math Geol*. 1981;13:175–89.
20. Xiao Y, Angulo MT, Friedman J, Waldor MK, Weiss ST, Liu Y-Y. Mapping the ecological networks of microbial communities. *Nat Commun*. 2017;8:2042.
21. Röttgers L, Faust K. From hairballs to hypotheses—biological insights from microbial networks. *FEMS Microbiol Rev*. 2018;42:761–80.
22. Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, et al. Correlation detection strategies in microbial data sets vary widely in sensitivity and precision. *ISME J*. 2016;10:1669–81.
23. Lima-Mendez G, Faust K, Henry N, Decelle J, Colin S, Carcillo F, et al. Determinants of community structure in the global plankton interactome. *Science*. 2015;348:1262073.
24. Yang Y, Chen N, Chen T. Inference of environmental factor-microbe and microbe-microbe associations from metagenomic data using a hierarchical Bayesian statistical model. *Cell Systems*. 2017;4:129–137.e5.
25. Kurtz ZD, Bonneau R, Müller CL. Disentangling microbial associations from hidden environmental and technical factors via latent graphical models. *bioRxiv*. 2019. <https://doi.org/10.1101/2019.12.21.885889>.
26. Verry L, Sella N, Affeldt S, Singh PP, Isambert H. Learning causal networks with latent variables from multivariate information in genomic data. *PLoS Comput Biol*. 2017;13:1–25.
27. Margolin AA, Nemenman I, Basso K, Wiggins C, Stolovitzky G, Favera RD, et al. ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC Bioinformatics*. 2006;7:57.
28. Zoppoli P, Morganello S, Ceccarelli M. TimeDelay-ARACNE: reverse engineering of gene networks from time-course data by an information theoretic approach. *BMC Bioinformatics*. 2010;11:154.
29. Villaverde AF, Ross J, Morán F, Banga JR. MIDER: network inference with mutual information distance and entropy reduction. *PLoS One*. 2014;9:1–15.
30. Villaverde AF, Becker K, Banga JR. PREMER: a tool to infer biological networks. *IEEE/ACM Trans Comput Biol Bioinformatics*. 2018;15:1193–202.
31. Pascual-García A, Tamames J, Bastolla U. Bacteria dialog with Santa Rosalia: are aggregations of cosmopolitan bacteria mainly explained by habitat filtering or by ecological interactions? *BMC Microbiol*. 2014;14:284.
32. Mandakovic D, Rojas C, Maldonado J, Latorre M, Travisany D, Delage E, et al. Structure and co-occurrence patterns in microbial communities under acute environmental stress reveal ecological factors fostering resilience. *Sci Rep*. 2018;8:5875.
33. Brisson V, Schmidt J, Northen TR, Vogel JP, Gaudin A. A new method to correct for habitat filtering in microbial correlation networks. *Front Microbiol*. 2019;10:585.
34. Barzel B, Barabási A-L. Network link prediction by global silencing of indirect correlations. *Nat Biotechnol*. 2013;31:720–5.
35. Feizi S, Marbach D, Médard M, Kellis M. Network deconvolution as a general method to distinguish direct dependencies in networks. *Nat Biotechnol*. 2013;31:726–33.
36. Alipanahi B, Frey BJ. Network cleanup. *Nat Biotechnol*. 2013;31:714–5.
37. Ghassami A, Kiyavash N. Interaction information for causal inference: The case of directed triangle. In: 2017 IEEE International Symposium on Information Theory (ISIT); 2017. p. 1326–30.
38. Cover TM, Thomas JA. Inequalities in information theory. *Elements of Information Theory*. 2001:482–509. <https://doi.org/10.1002/0471200611.ch16>.
39. Deutschmann IM. EnDED - - Environmentally-driven edge detection program. Zenodo. 2019. <https://doi.org/10.5281/zenodo.3271730>.
40. Core R. Team. R: A Language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2019; <https://www.R-project.org/>.
41. Gasol JM, Cardelús C, G Morán XA, Balagué V, Forn I, Marrasé C, et al. Seasonal patterns in phytoplankton photosynthetic parameters and primary production at a coastal NW Mediterranean site. *Sci Mar*. 2016;80:63–77.
42. Faust K, Sathirapongsasuti JF, Izard J, Segata N, Gevers D, Raes J, et al. Microbial co-occurrence relationships in the human microbiome. *PLoS Comput Biol*. 2012;8:1–17.
43. Friedman J, Alm EJ. Inferring correlation networks from genomic survey data. *PLoS Comput Biol*. 2012;8:1–11.
44. Giner CR, Balagué V, Krabberød AK, Ferrera I, Reñé A, Garcés E, et al. Quantifying long-term recurrence in planktonic microbial eukaryotes. *Mol Ecol*. 2019;28:923–35.
45. The Human Microbiome Project Consortium, Huttenhower C, Gevers D, Knight R, Abubucker S, Badger JH, et al. Structure, function and diversity of the healthy human microbiome. *Nature*. 2012;486:207–14.
46. Pinedo-González P, West AJ, Tovar-Sánchez A, Duarte CM, Marañón E, Cermeño P, et al. Surface distribution of dissolved trace metals in the oligotrophic ocean and their influence on phytoplankton biomass and productivity. *Glob Biogeochem Cycles*. 2015;29:1763–81.
47. Ai D, Li X, Pan H, Chen J, Cram JA, Xia LC. Explore mediated co-varying dynamics in microbial community using integrated local similarity and liquid association analysis. *BMC Genomics*. 2019;20:185.
48. Jang IS, Margolin A, Califano A. hARACNE: improving the accuracy of regulatory model reverse engineering via higher-order data processing inequality tests. *Interface Focus*. 2013;3:20130011.
49. Fernandes AD, Gloor GB. Mutual information is critically dependent on prior assumptions: would the correct estimate of mutual information please identify itself? *Bioinformatics*. 2010;26:1135–9.
50. Zhao J, Zhou Y, Zhang X, Chen L. Part mutual information for quantifying direct associations in networks. *Proc Natl Acad Sci*. 2016;113:5130–5.
51. Latorre F, Deutschmann IM, Labarre A, Obiol A, Krabberød AK, Pelletier E, et al. Niche adaptation promoted the evolutionary diversification of tiny ocean predators. *Proc Natl Acad Sci U S A*. 2021;118:e2020955118.
52. Faust K. Towards a better understanding of microbial community dynamics through high-throughput cultivation and data integration. *mSystems*. 2019;4. <https://doi.org/10.1128/mSystems.00101-19>.
53. McCarren J, Becker JW, Repeta DJ, Shi Y, Young CR, Malmstrom RR, et al. Microbial community transcriptomes reveal microbes and metabolic pathways associated with dissolved organic matter turnover in the sea. *Proc Natl Acad Sci*. 2010;107:16420–7.
54. Yoon HS, Price DC, Stepanauskas R, Rajah VD, Sieracki ME, Wilson WH, et al. Single-cell genomics reveals organismal interactions in uncultivated marine protists. *Science*. 2011;332:714–7.
55. Bashan A, Gibson TE, Friedman J, Carey VJ, Weiss ST, Hohmann EL, et al. Universality of human microbial dynamics. *Nature*. 2016;534:259–62.
56. Berry D, Widder S. Deciphering microbial interactions and detecting keystone species with co-occurrence networks. *Front Microbiol*. 2014;5:219.
57. Gonze D, Coyte KZ, Lahti L, Faust K. Microbial communities as dynamical systems. *Curr Opin Microbiol*. 2018;44:41–9.
58. Benincà E, Dakos V, Van Nes EH, Huisman J, Scheffer M. Resonance of plankton communities with temperature fluctuations. *Am Nat*. 2011;178:E85–95.
59. Kettle H, Holtrop G, Louis P, Flint HJ. microPop: Modelling microbial populations and communities in R. *Methods Ecol Evol*. 2018;9:399–409.
60. Vallina SM, Martinez-Garcia R, Smith SL, Bonachela JA. Models in Microbial Ecology. In: Schmidt TM, editor. *Encyclopedia of Microbiology* (Fourth Edition). 4th ed. Oxford: Academic Press; 2019. p. 211–46. <https://doi.org/10.1016/B978-0-12-809633-8.20789-9>.
61. Dam P, Fonseca LL, Konstantinidis KT, Voit EO. Dynamic models of the complex microbial metapopulation of lake mendota. *NPJ Syst Biol Appl*. 2016;2:16007.
62. Stein RR, Bucci V, Toussaint NC, Buffie CG, Ratsch G, Pamer EG, et al. Ecological modeling from time-series inference: insight into dynamics and stability of intestinal microbiota. *PLoS Comput Biol*. 2013;9:1–11.
63. Klemm K, Eguíluz VM. Growing scale-free networks with small-world behavior. *Phys Rev E*. 2002;65:057102.
64. Novak M, Yeakel JD, Noble AE, Doak DF, Emerson M, Estes JA, et al. Characterizing species interactions to understand press perturbations: what is the community matrix? *Annu Rev Ecol Syst*. 2016;47:409–32.
65. Haydon D. Pivotal assumptions determining the relationship between stability and complexity: an analytical synthesis of the stability-complexity debate. *Am Nat*. 1994;144:14–29.

66. Jackson DA. Stopping rules in principal components analysis: a comparison of heuristic and statistical approaches. *Ecology*. 1993;74:2204–14.
67. Legendre P, Legendre LF. *Numerical ecology*. Elsevier; 2012.
68. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGinn D, et al. *vegan: Community Ecology Package*. 2019. <https://CRAN.R-project.org/package=vegan>.
69. Soetaert K, Petzoldt T, Setzer RW. Solving differential equations in R: Package deSolve. *J Stat Softw*. 2010;33:1–25.
70. Schauer M, Balagué V, Pedrós-Alió C, Massana R. Seasonal changes in the taxonomic composition of bacterioplankton in a coastal oligotrophic system. *Aquat Microb Ecol*. 2003;31:163–74.
71. Stoeck T, Bass D, Nebel M, Christen R, Jones MDM, Breiner H-W, et al. Multiple marker parallel tag environmental DNA sequencing reveals a highly complex eukaryotic community in marine anoxic water. *Mol Ecol*. 2010;19:21–31.
72. Herlemann DP, Labrenz M, Jürgens K, Bertilsson S, Waniek JJ, Andersson AF. Transitions in bacterial communities along the 2000 km salinity gradient of the Baltic Sea. *ISME J*. 2011;5:1571–9.
73. Apprill A, McNally S, Parsons R, Weber L. Minor revision to V4 region SSU rRNA 806R gene primer greatly increases detection of SAR11 bacterioplankton. *Aquat Microb Ecol*. 2015;75:129–37.
74. 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.
75. Moritz S, Gatscha S, imputeTS: Time series missing value imputation. 2017. <https://github.com/SteffenMoritz/imputeTS>.
76. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naïve Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol*. 2007;73:5261–7.
77. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res*. 2012;41:D590–6.
78. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol*. 1990;215:403–10.
79. Guillou L, Bachar D, Audic S, Bass D, Berney C, Bittner L, et al. The Protist Ribosomal Reference database (PR²): a catalog of unicellular eukaryote Small Sub-Unit rRNA sequences with curated taxonomy. *Nucleic Acids Res*. 2012;41:D597–604.
80. McNichol J, Berube PM, Biller SJ, Fuhrman JA, Gilbert JA. Evaluating and improving small subunit rRNA PCR primer coverage for bacteria, archaea, and eukaryotes using metagenomes from global ocean surveys. *mSystems*. 2021;6:e00565–21.
81. Grasshoff K, Kremling K, Ehrhardt M. *Methods of seawater analysis*. John Wiley & Sons; 2009.
82. North BV, Curtis D, Sham PC. A note on the calculation of empirical P values from Monte Carlo procedures. *Am J Hum Genet*. 2002;71:439–41.
83. Veech JA. Significance testing in ecological null models. *Theoretical Ecol*. 2012;5:611–6.
84. Meyer PE, Lafitte F, Bontempi G. minet: A R/Bioconductor package for inferring large transcriptional networks using mutual information. *BMC Bioinformatics*. 2008;9:461.

Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

