# Disentangling environmental effects in microbial association networks

— **Source link** ↗

Ina M. Deutschmann, Gipsi Lima-Mendez, Anders K. Krabberød, Jeroen Raes ...+3 more authors

**Institutions:** Spanish National Research Council, Université catholique de Louvain, University of Oslo, Katholieke Universiteit Leuven

# Disentangling environmental effects in microbial association networks

Ina Maria Deutschmann[1*], Gipsi Lima-Mendez[2], Anders K. Krabberød[3], Jeroen Raes[4,5], Sergio M. Vallina[6], Karoline Faust[5*†] and Ramiro Logares[1*]

[1] Institute of Marine Sciences, CSIC, Passeig Marítim de la Barceloneta, 37, 08003, Barcelona, Spain.

[2] Louvain Institute of Biomolecular Science and Technology (IBST), Université catholique de Louvain, Croix du sud 4-5/L7.07.02, 1348, Louvain-la-Neuve, Belgium.

[3] Department of Biosciences/Section for Genetics and Evolutionary Biology (EVOGENE), University of Oslo, p.b. 1066 Blindern, N-0316, Oslo, Norway.

[4] VIB Center for Microbiology, Herestraat 49-1028, 3000, Leuven, Belgium.

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

[6] Spanish Institute of Oceanography (IEO), Ave Principe de Asturias 70 Bis, 33212, Gijon, Spain.

[*] Corresponding authors:
Ina Maria Deutschmann (ina.m.deutschmann@gmail.com)
Karoline Faust (karoline.faust@kuleuven.be)
Ramiro Logares (ramiro.logares@icm.csic.es)

[†] Shared last authors

# 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 associations are 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 (**En**vironmentally-**D**riven **E**dge **D**etection), 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.

**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 (Locey & Lennon, 2016), and $\approx 10^{30}$ microbial cells (Whitman *et al.*, 1998; Kallmeyer *et al.*, 2012). Microorganisms have crucial roles in the biosphere by contributing to global biogeochemical cycles (Falkowski *et al.*, 2008) and underpinning diverse food webs. The importance of microbes for the functioning of ecosystems cannot be understood without considering their ecological interactions (DeLong, 2009; Krabberød *et al.*, 2017). These allow transferring carbon and energy to upper trophic levels, and the recycling of nutrients and energy (Worden *et al.*, 2015). 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 (Faust & Raes, 2012). Although microbial communities are highly interconnected (Layeghifard *et al.*, 2017), our knowledge about ecological interactions in the microbial world is still limited (Krabberød *et al.*, 2017; Bjorbækmo *et al.*, 2019).

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 (Xia *et al.*, 2011, 2013), CoNet (Faust & Raes, 2016), SPIEC-EASI (Kurtz *et al.*, 2015), or FlashWeave (Tackmann *et al.*, 2019)). 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

3

87    microbial interaction hypotheses.

88

**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 (Lv *et al.*, 2019). Network inference from -omics data is difficult (Li *et al.*, 2016; Layeghifard *et al.*, 2017) because of both technical and interpretation challenges. One challenge is the compositional nature of the data produced by DNA sequencers (Gloor *et al.*, 2017). There are several network tools (Li *et al.*, 2016) that consider this, e.g., SPIEC-EASI (Kurtz *et al.*, 2015). 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 (Aitchison, 1981). 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 (Xiao *et al.*, 2017). 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 (Röttjers & Faust, 2018). 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 (Weiss *et al.*, 2016), or increasing the cut-off for other criteria such as the association strength, prevalence, or abundance filtering (Röttjers & Faust, 2018). Another strategy is agglomeration using taxonomic or ecological (functional) groupings (Lima-Mendez *et al.*, 2015).

The interpretation challenge addressed in this study are indirect dependencies (associations) caused by environmental factors. For most microbial association networks, an

4

118  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 (Kurtz *et al.*, 2015; Yang *et al.*, 2017), e.g. SPIEC-EASI (Kurtz *et al.*, 2015, 2019), miic (Verny *et al.*, 2017), or FlashWeave (Tackmann *et al.*, 2019). FlashWeave can also integrate metadata to avoid indirect associations driven by environmental factors but currently does not support missing data. The tool ARACNE (Margolin *et al.*, 2006) aims to eliminate indirect associations by using an information theoretic property (the *Data Processing Inequality*, DPI, in Methods). The extension TimeDelay-ARACNE (Zoppoli *et al.*, 2010) tries to extract dependencies between different times. Another approach including time-delay is implemented in the tool MIDER (Villaverde *et al.*, 2014), which combines mutual information-based distances and entropy reduction to detect indirect interactions (*Mutual Information*, MI, in Methods). PREMER (Villaverde *et al.*, 2018), a successor of MIDER, allows to include previous knowledge, e.g., known non-existent associations.

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 (Pascual-García *et al.*, 2014). However, this will keep generalist while removing specialist. Another approach divides datasets displaying a great environmental heterogeneity into sub datasets of similar environmental conditions (Röttjers & Faust, 2018). For example, a previous work (Mandakovic *et al.*, 2018) constructed two networks representing bacterial

149  soil communities from two different sections of a pH, temperature, and humidity gradient.
150  Another work (Lima-Mendez *et al.*, 2015) constructed ocean depth-specific networks to
151  account for environmental differences between the surface layer and the deep chlorophyll
152  maximum layer. In addition to dividing samples, an algorithm aiming to correct for habitat
153  filtering effects (Brisson *et al.*, 2019), subtracts, for a given habitat, the mean abundance from
154  each microorganisms within each sample. However, this approach is limited to the identified
155  habitat groups that should have a similar sample size.

156  In contrast, there are methods accounting for indirect dependencies after network
157  construction. For instance, global silencing, (Barzel & Barabási, 2013) and network
158  deconvolution (Feizi *et al.*, 2013) aim to recover true direct associations from observed
159  correlations. Both techniques are sensitive to missing variables (Alipanahi & Frey, 2013).
160  Another method, called *Sign Pattern*, SP, uses environmental triplets (Lima-Mendez *et al.*,
161  2015). An environmental triplet contains two microorganisms and one environmental factor,
162  which are associated to each other. SP combines the signs of association scores (positive or
163  negative) to determine if a microbial association should be classified as indirect (SP in
164  Methods). Its major drawback is edge removal where microorganisms with similar
165  environmental preference interact. Along SP and network deconvolution, the *Interaction*
166  *Information*, II, was applied in (Lima-Mendez *et al.*, 2015). Within an environmental triplet,
167  the II method aims to indicate whether an edge is due entirely to shared environmental
168  preferences (II<0) or whether environmental preferences and true interactions are entangled
169  (II>0). However, II cannot determine which associations in a triplet is indirect (II in
170  Methods). Here, we study several indirect edge detection methods: SP, *Overlap*, (OL,
171  developed here), II, DPI, and their combination.

172

**EnDED is an implementation of four methods and their combination**
174  This article presents EnDED, which implements four approaches, and their combination, to
175  indicate environmentally-driven (indirect) associations in microbial networks. The four
176  methods are: Sign Pattern (Lima-Mendez *et al.*, 2015), Overlap (developed here), Interaction
177  Information (Lima-Mendez *et al.*, 2015; Ghassami & Kiyavash, 2017), and Data Processing
178  Inequality (Cover & Thomas, 2001; Margolin *et al.*, 2006). SP requires an association score
179  that represents co-occurrence when it is positive, and mutual-exclusion when it is negative.

180  OL requires temporal data with a known start and end of the association to determine whether

181  the microbial association occurs in a time window when both microorganisms are associated

182  to the same environmental factor. The II method indicates the existence of one indirect

183  dependency between three components that are associated with each other. The DPI method

184  states that the association with the smallest mutual information is the indirect association.

185  Here, we evaluate each method and their combination on how well they detect

186  environmentally-driven associations on association networks from simulated data including

187  two environmental factors. Combining methods in an intersection approach retains more true

188  interactions than each method on its own. A union approach was discarded because it would

189  have retained the smallest number of true interactions. We are able to disentangle and filter

190  environmentally-driven edges from microbial association networks (0.95-0.96 in positive

191  predictive value and 0.35-0.83 in accuracy). We also applied EnDED to disentangle and filter

192  environmentally-driven edges from a real marine microbial association network based on ten

193  years of monthly sampling including ten environmental factors. EnDED contributed to both,

194  generating more reliable hypotheses on microbial interactions, and facilitating network

195  analysis by removing edges from dense "hairball" networks. EnDED is publicly available

196  (Deutschmann, 2019).

197

198  ## Results

199  **Simulated data**

200  To evaluate EnDED's performance in removing environmentally-driven associations, we

201  simulated 1000 abundance time-series datasets with 50 microorganisms and known true

202  interactions between them. We obtained another 1000 datasets with noise (hereafter dwn).

203  We constructed the networks (hereafter simulated networks) with the tool eLSA (Xia *et al.*,

204  2011, 2013) (see methods). The simulated networks contained on average (computed as the

205  median) 50 nodes and 1087 edges (1063 dwn), of which 60 (59 dwn) were true interactions

206  (edges present in the inferred and true network) and 1026 (1005 dwn) false associations

207  (edges present in the inferred but absent in the true network). Networks inferred from

208  simulated data without noise contained on average one more true interaction but also 21 more

209  false interactions than the networks inferred from simulated data with noise.

210      A simple approach to discriminate true interactions (desired) from false associations

211     (undesired) would be to use a threshold for the association strength, which could be suitable

212     if the values for true interactions and false associations are i) following different distributions,

213     and ii) the distributions are mainly non-overlapping. We tested the former requirement with

214     a two-sample Kolmogorov-Smirnov test with the R (R Core Team, 2019) function *ks.test*.

215     Using a 95% (99%, 99.9%) confidence level, the distributions were significantly different

216     for 358 (192, 66) simulated datasets and 355 (173, 68) simulated datasets with noise, which

217     is slightly more than one third of them. This indicates that an association strength cut-off is

218     unsuitable to separate true interactions from false associations. More sophisticated

219     approaches than a simple threshold include the methods implemented in EnDED: SP, OL, II,

220     DPI, and their combination.

221          Combining the methods in an intersection approach (hereafter referred to as

222     intersection combination), we classified on average 348 (228 dwn), that is 32% (22% dwn)

223     of the associations, to be environmentally-driven. The number of correctly detected false

224     associations was on average 332 (219 dwn), i.e., 96% of the removed edges. The resulting

225     networks contained on average 737 (828 dwn) edges. When each method was individually

226     applied more edges were removed: 87% (86% dwn) for SP and OL, 67% (60% dwn) for II,

227     and 44% (32% dwn) for DPI. The fraction of correctly removed edges for individual methods

228     was on average 95%. Comparing the methods on correctly detected false associations, the

229     greatest agreement was observed between SP and OL, whereas DPI appeared to be the most

230     conservative in not agreeing with other methods and, subsequently, reducing the number of

231     detected edges in the intersection combination approach (Supplementary Table

232     S1).Individual methods removed more edges from the network than the intersection

233     combination, where all methods must agree. However, a method's performance is not solely

234     determined by the number of removed edges.

235          To evaluate the removal of environmentally-driven edges, we scored the different

236     approaches based on five evaluation measurements (see Methods): the true positive rate,

237     TPR, true negative rate, TNR, false positive rate, FPR, positive predicted value, PPV, and

238     accuracy, ACC, (Figure 1 and Supplementary Table S2). In order to determine these

239     measurements, we first determined true and false positives, as well as true and false

240     negatives. A true positive is a false association in the network that is correctly removed by a

241     method, and a false negative is a false association that is incorrectly not removed. A false

8

242 positive is a true interaction in the network that is incorrectly removed by a method, and a
243 true negative is a true interaction that correctly is not removed by a method. The ideal method
244 maximizes true positives and true negatives and minimizes false positives and false
245 negatives.

246      The intersection combination under-performed compared to each individual method,
247 SP and OL perform best, and II performs better than DPI according to TPR, FPR and ACC
248 (Figure 1). However, applying each method individually has the drawback of removing more
249 true interactions. On average there are 60 (59 dwn) true interactions in the simulated
250 networks. The individual methods removed 86% (85% dwn) (SP), 85% (84% dwn) (OL),
251 60% (51% dwn) (II), and 38% (28% dwn) (DPI). Therefore, although the intersection
252 combination removed fewer edges, it outperformed the others according to the TNR because
253 it eliminated fewer of the true interactions, 25% (16% dwn). All methods had high PPV
254 values with half of all measured PPV above ≈0.95. According to PPV, intersection
255 combination performed best and SP and OL performed worst (Figure 1).

256

257 **Real data**
258 After testing EnDED's performance on simulated networks, we applied it to a real microbial
259 association network, which was constructed from 10 years of monthly samples from January
260 2004 to December 2013 at the Blanes Bay Microbial Observatory (BBMO) (Gasol *et al.*,
261 2016). These samples included bacteria and eukaryotes of two size-fractions: picoplankton
262 (0.2-3 µm) and nanoplankton (3-20 µm). We estimated community composition via
263 metabarcoding of the 16S and 18S rRNA gene, and inferred an association network, hereafter
264 referred to as BBMO network (see Methods). The BBMO network contained 762 nodes
265 including 754 ASVs and eight of the ten available environmental factors, and 30498 edges
266 including 29820 microbial edges and 607 edges between a microorganism and an
267 environmental factor. The network contained more positive (24458, 82.0%) than negative
268 (5362, 18.0%) microbial associations (Figure 2).

269      We found that 25230 (84.6%) of the network edges were in at least one and in
270 maximum six environmental triplets (Figure 2 and Supplementary Table S3). Overall, we
271 detected 35166 environmental triplets within the BBMO network. Of the ten considered
272 environmental factors, $PO_4^{3-}$ and salinity were not associated to any microorganism in the

273    network, and turbidity and $NH_4^+$ were not found within a triplet. Thus, six environmental

274    factors remained: Temperature (1831 environmentally-driven edges were removed due to

275    Temperature) and day length (652 removed edges) were the top two environmental factors

276    affecting microbial associations, followed by total chlorophyll (175), $SiO_2$ (5) and $NO_3^-$ (1);

277    no edge was removed due to $NO_2^-$.

278         The intersection combination removed 2488 (≈8.3%) associations from the BBMO

279    network. We classified and quantified these indirect edges according to the domain of the

280    nodes (bacteria - eukaryotes, nanoplankton – picoplankton), environmental factor, and the

281    number of triplets a microbial edge was in (Figure 2 and Supplementary Table S4). Compared

282    to the intersection combination, each method individually removed more edges: 84.6% (SP

283    and OL removing all microbial edges present in a triplet), 25.7% (II), and 24.8% (DPI); that

284    is, removal was 3 to 10 times larger.

285         We also determined for each association the Jaccard index, which indicates how often

286    two microorganisms appear together in the dataset. We assumed that two microbes that

287    appear together < 50% of the time are less likely to have true contemporary ecological

288    interactions and the corresponding association is more likely to be false. We found that only

289    27.7% of the indirect associations had a Jaccard index above 0.5 compared to 61.1% of the

290    associations that were not indirect. This discrepancy was bigger for negative edges, with

291    1.2% above and 98.8% below 0.5 (Table 1). The fact that over 72.3% of environmentally-

292    driven associations had a Jaccard index equal or below 0.5 strengthened the decision of their

293    removal.

294    The intersection combination removed more negative than positive edges, 1554 and 934,

295    respectively (Figure 2). However, there were 20334 positive and 4896 negative microbial

296    associations that were found in at least one environmental triplet, so the method removed

297    31.7% of the negative and only 4.6% of the positive edges. If we randomly removed 2488

298    edges, we would expect 18.0 % to be negative (i.e. 448) and 82.0 % of them to be positive

299    (i.e. 2040). If we restrict these calculations to the 25230 microbial associations that were

300    found in at least one environmental triplet, with 20334 of them being positive and 4896 being

301    negative, we would expect to remove 19.4% (i.e. 483) of negative and 80.6% (i.e. 2005) of

302    positive edges. The probability of randomly removing less positive than negative associations

303    is nearly zero, since it follows a multivariate hypergeometric distribution:

$$P\left(k_{neg}, k_{pos}\right) = \frac{\binom{N_{neg}}{k_{neg}} \cdot \binom{N_{pos}}{k_{pos}}}{\binom{N}{n}},$$ Eq. (1)

304 where $N_{pos}$ and $N_{neg}$ are the number of positive and negative associations in the network,

305 respectively, $k_{pos}$ is the number of removed positive and $k_{neg}$ the removed negative

306 associations from the network, $N$ is the number of associations in the network, and $n$ is the

307 number of removed associations from the network. The removal of more negative edges

308 through intersection combination indicates that this removal was not random or, in other

309 words, that negative associations are more likely to represent environmentally-driven edges.

310     To evaluate the performance of EnDED on the BBMO network, we considered

311 interactions described in literature and collected in the Protist Interaction Database (PIDA)

312 (Bjorbækmo *et al.*, 2019). Studies typically compare the associations of a network to those

313 reported in the literature at the genus level (Lima-Mendez *et al.*, 2015). The ambiguity in

314 taxonomic classification and the large number of edges challenged this comparison. Thus,

315 we implemented a function to compare strings and match the taxonomic classification of a

316 microorganism in the BBMO network to those in the scientific literature (PIDA). We found

317 that only 29 (0.1%) associations were supported by interactions described in the literature

318 (Table 2). That is, 99.9% of associations in the BBMO network (before applying EnDED)

319 could not be used to evaluate EnDED's performance. These 29 associations describe eight

320 unique interactions between eight microorganisms, and 18 edges were in an environmental

321 triplet to which each method as well as their combination were applied (summary in Table

322 2). Ideally none of these described associations should be removed by EnDED. Yet, the

323 intersection combination removed five associations (Table 2). In contrast and even worse,

324 SP and OL removed all 18 edges, II eight and DPI nine edges. The additionally removed

325 edges by individual methods are associations between a diatom (*Thalassiosira*) and an

326 unknown *Flavobacteriia*. Considering only the genus level, there were 171 unique genera in

327 the BBMO network, and 700 in PIDA, combined there were 837 microbial genera, and 34

328 genera in both. Thus, 19.9% of the microbial genera found in the BBMO network were also

329 in PIDA, and 4.9% of the genera found in PIDA were also found in the BBMO network.

330

331

## 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 evaluated 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 (Lv *et al.*, 2019). 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, Röttjers and Faust simulated an increasing environmental influence and observed a decrease in retrieving true interactions from inferred associations (Röttjers & Faust, 2018). The observation holds for several network construction methods for cross-sectional data, including CoNet (Faust *et al.*, 2012), SparCC (Friedman & Alm, 2012), SPIEC-EASI (Kurtz *et al.*, 2015), 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

363    pico- and nano-plankton at the BBMO, using the same basal 10-year dataset used here,

364    indicated a strong seasonality at the community level (Giner *et al.*, 2019).

365

**Comparisons of indirect edge detection on other datasets**

367    In our BBMO network we found that the majority (84.6%) of the microbial edges was within

368    at least one environmental triplet. This was 2.6 times higher than what was found for an

369    association network inferred from data considering microorganisms and small metazoans

370    from two ocean depths across 68 stations around the world and various size fractions

371    (hereafter global interactome) (Lima-Mendez *et al.*, 2015). This global interactome contains

372    29912 (32.3%) edges that were within at least one environmental triplet (Lima-Mendez *et

373    al.*, 2015). In the previous study, 29900 edges in the global interactome (≈100% of triplets

374    and 32% of all edges) were attributed to environmental factors by SP, similarly to this study

375    as SP removed all edges within triplets in the BBMO network. II indicated 11043

376    environmentally-driven edges in the global interactome (≈37% of triplets and 12% of all

377    edges) with *p*-value below 0.05 in a permutation test with 500 iterations. In comparison, II

378    removed a higher fraction of edges in the BBMO network when considering all edges

379    (25.7%), but less when considering within the triplets (30.4%). Network deconvolution

380    suggested 22439 environmentally-driven edges (≈75% of triplets and 24% of all edges)

381    within the global interactome, and the three methods agreed for 8209 edges (≈27% of triplets

382    and 8.9% of all edges). In comparison, we detected slightly less environmentally-driven

383    associations for the BBMO network (8.3% of all edges). These differences suggest that a

384    higher environmental heterogeneity in the dataset may induce more indirect edges. Also, the

385    effects of indirect dependencies may depend on dataset type (e.g., temporal vs. spatial). These

386    possible differences and their effect on environmentally-driven edges should be further

387    investigated.

388         Using II for the BBMO network, we identified a moderate number of

389    environmentally-driven associations. DPI also identified a moderate number (24.8%, 29.3%

390    when considering only triplets), whereas SP or OL identified a ubiquitous number of

391    environmentally-driven edges (84.6%, 100% when considering only triplets). This indicates

392    that SP and OL are strict and should be used in combination with other methods in an

393    intersection approach.

13

394   In another study, the tool FlashWeave (Tackmann *et al.*, 2019) predicted direct
395 microbial interactions in the human microbiome using the Human Microbiome Project
396 (HMP) dataset, including heterogeneous microbial abundance data of 68818 samples (The
397 Human Microbiome Project Consortium: Huttenhower *et al.*, 2012). The inferred networks
398 (with and without metadata) were sparser than our networks. The network with metadata
399 contained 10.7% fewer associations compared to the network without metadata, slightly
400 more than in our results from BBMO.

401

**Factors causing indirect microbial associations**

403 From the simulated networks, we found that using the intersection combination instead of
404 each method individually, we maintained more true interactions at the cost of more false
405 associations in the network—more when considering simulated networks including noise.
406 Comparing our simulated network against the BBMO network, the intersection combination
407 classified a higher number of edges as environmentally-driven in the simulated networks
408 32% (22% dwn) than in the BBMO network (8.3%). For the simulated data, we previously
409 knew the environmental factor influencing pairwise microbial associations. For the BBMO
410 data, we used ten available environmental factors, but not all factors that could affect
411 microbial dynamics. Even though the most important factors influencing microbial seasonal
412 dynamics at BBMO were considered (Giner *et al.*, 2019), there are several factors that were
413 not measured and that could generate indirect edges. The indirect edges associated to these
414 factors were not detected in our analyses. Similarly, indirect edges associated to biotic
415 interactions (e.g., two bacteria sharing a positive edge as they are symbionts in the same
416 protists) were not considered. Future sampling for microbial interaction research should
417 expand metadata collection in order to detect (more) abiotic and biotic factors that could
418 generate indirect edges.

419   While temperature and day length (hours of light) were the top two environmental
420 factors affecting microbial associations in the BBMO network, the most important
421 environmental factors in the global interactome (Lima-Mendez *et al.*, 2015) were phosphate
422 concentration and temperature, followed by nitrite concentration and mixed-layer depth.
423 Although we considered $PO_4^{3-}$ and salinity, they were not associated to any microorganism
424 in the network, which may reflect the low variation of these environmental factors in the

14

425    studied marine site (BBMO). For instance, the standard deviation in the BBMO dataset was

426    < 1 for $PO_4^{3-}$ and salinity, in contrast to the global interactome dataset (Lima-Mendez *et al.*,

427    2015), where it was about 20-30 when considering all samples. During the Malaspina-2010

428    Circumnavigation Expedition, the concentrations of trace metals were determined for 110

429    surface water samples (Pinedo-González *et al.*, 2015). The previous study indicates

430    relationships between primary productivity and trace nutrients, more specifically for the

431    Indian Ocean Cd, the Atlantic Ocean Co, Fe, Cd, Cu, V and Mo, and the Pacific Ocean Fe,

432    Cd, and V. Thus, trace metals are further environmental factors that may play an important

433    role in regulating oceanic primary productivity.

434

435    **Limitations of EnDED**

436    EnDED detects and removes environmentally-driven indirect edges. However, its triplet

437    analysis could be extended to remove indirect edges driven by taxa, as done with gene triplets

438    (Margolin *et al.*, 2006). A recent update of the network construction tool eLSA (Xia *et al.*,

439    2011, 2013) permits to examine how a factor, such as a microorganism or environmental

440    variable, mediates the association of two other factors (Ai *et al.*, 2019), which allows the

441    study of interactions between three factors. Furthermore, triplets limit the study to first-order

442    indirect dependencies, neglecting higher-order indirect dependencies. Such limitation was

443    solved for the DPI method by examining associations in quadruplets, quintuplets, and

444    sextuplets (Jang *et al.*, 2013). Implementing higher-order DPI and adjusting the other three

445    methods to account for higher-order indirect dependencies may be promising but one needs

446    to be aware that incorporating higher-order dependencies will also increase the risk of over-

447    fitting. Further, all relevant (measured) environmental factors could be incorporated into the

448    calculation of II, which would combine environmental triplets. However, we reason that such

449    adjustments would require a larger sample size. Both II and DPI calculate MI that measures

450    the dependence between two random variables. EnDED is limited by including one function

451    to estimate the MI. A comparison of four different MI estimates revealed that obtaining the

452    true value of MI is not straightforward, and minor variations of assumptions yield different

453    estimates (Fernandes & Gloor, 2010). Lastly, the conditional mutual information, CMI,

454    which quantifies nonlinear direct relationships among variables, can be underestimated if

455    variables have tight associations in a network (Zhao *et al.*, 2016). The so-called part mutual

456    information, PMI, measurement can help overcome CMI's underestimations. Although using

457    PMI instead of CMI looks promising, calculating PMI is computationally more demanding

458    (Zhao *et al.*, 2016).

459

460    **Future Perspectives**

461    In this study, we have shown that EnDED with an intersection combination approach

462    provides less dense networks, but still with many potential interactions. We observed a trade-

463    off comparing single methods with the combination approach (intersection combination).

464    Although the latter kept more true interactions, it kept also more false associations. Inferring

465    emergent properties is a key task in microbial ecology to characterize microbial ecosystems

466    from a network-perspective. Thus, if the study aim is to explore patterns of network topology

467    rather than single edges, inferring a network comparable to the real interaction network may

468    be more useful than accuracy of single edges. However, investigations aiming to provide

469    potential interaction partners may use EnDED with the intersection combination approach

470    (e.g., (Latorre *et al.*, 2021)). Specific associations may be validated with experiments or

471    microscopy (Lima-Mendez *et al.*, 2015; Krabberød *et al.*, 2017). However, we suggest to

472    first further reduce the set of potential interaction hypotheses. To improve the selection of

473    interaction hypotheses, we propose to score associations based on re-occurrence: in time, as

474    done with microbial abundance seasonality (Giner *et al.*, 2019), or space, where an

475    association appears in different networks based on different datasets, or different regions of

476    the world. In a previous study using 313 samples, including seven size-fractions, four

477    domains (Bacteria, Archaea, Eukarya, and viruses), and two depths from 68 stations across

478    eight oceanic provinces, 14% of the 81590 predicted biotic interactions were identified as

479    local (Lima-Mendez *et al.*, 2015). Thus, re-occurrent associations may suggest a higher

480    likelihood that the association represents a true ecological interaction, reducing the number

481    of interaction hypotheses to the strongest ones. Another strategy to shortlist interaction

482    hypotheses is to incorporate additional data into the network and use a multi-layer network

483    approach. Such data could be environmental preferences such as temperature or salinity

484    optima, size of cells, presence of chloroplasts, or data obtained from High-Throughput

485    Cultivation (Faust, 2019), microbial community transcriptomes that reveal metabolic

486    pathways (McCarren *et al.*, 2010), or interactions inferred from Single-Cell genome data

487     (Yoon *et al.*, 2011; Krabberød *et al.*, 2017).

488

489     ## Conclusion

490     In this paper, we presented EnDED, an analysis tool to reduce the number of environmentally

491     induced indirect edges in inferred microbial networks. Applying EnDED on simulated

492     networks indicated that false associations, driven by environmental variables instead of true

493     interactions, were ubiquitous. However, EnDED's intersection combination classified a

494     minority of associations as environmentally-driven in a real (BBMO) network. Depending

495     on the single method used, we classified a moderate to high number of associations as

496     environmentally-driven in the same network. Nevertheless, associations driven by

497     environmental factors must be determined and quantified to generate more accurate insights

498     regarding true microbial interactions. EnDED provides a step forward in this direction.

499

500     ## Methods

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

502     To evaluate the performance of EnDED, we simulated a time series using an adjusted version

503     of the standard *generalized Lotka-Volterra model*, gLV (Berry & Widder, 2014; Bashan *et*

504     *al.*, 2016). The gLV can describe the dynamics of microbial communities, by including a first

505     order approach of the microbial interactions. The model's simplicity arises from the

506     assumption of linear interactions, which facilitates implementation and allows fast numerical

507     simulations. The gLV has, however, several limitations (Gonze *et al.*, 2018). For example,

508     gLV neglects higher-order interactions and the additivity of interaction strengths is a

509     weakness because they may be combined in different ways. Also, interactions are often

510     assumed to be constant parameters, but a reducing level of a nutrient may weaken cross-

511     feeding relationships. Moreover, gLV omits the influence of environmental factors, which,

512     for example, can induce oscillations in natural communities (Benincà *et al.*, 2011). Using a

513     model that accounts for nutrients (Kettle *et al.*, 2018) is more realistic but also more complex.

514     More elaborate mechanistic models of microbial dynamics than gLV solve explicitly the

515     global cycling of nutrients and are coupled to the oceanic circulation (see (Vallina *et al.*,

516     2019) for a review), but the added complexity can hamper understanding about the ecological

517     interactions among microorganisms when compared to a simpler gLV approach. Thus, we

518   chose to use a simpler extension of the gLV to account for the influence of environmental

519   factors (Stein *et al.*, 2013; Dam *et al.*, 2016). In order to allow the growth rates to vary when

520   the environmental variables change, environmental variables can be incorporated directly

521   into the gLV (Dam *et al.*, 2016; Röttjers & Faust, 2018). We simulated a time series using

522   the Klemm-Eguíluz algorithm (Klemm & Eguíluz, 2002), and an adjusted gLV. We adjusted

523   the model by defining microbial growth rates as a function dependent on one seasonal abiotic

524   environmental factor, and added an abiotic environmental factor in the interaction matrix.

525   We then used the time series generated by the gLV to obtain temporal microbial abundance

526   data. With this simulated data, we inferred a network that contained environmentally-driven

527   associations, needed to evaluate the performance of EnDED. We repeated this procedure

528   1000 times to obtain a large set of simulated networks, and then used the determined

529   abundance tables and Poisson distribution to obtain another 1000 simulated networks

530   including noise. The addition of noise was done by randomly drawing an abundance from

531   the Poisson distribution with λ equaling the original abundance of a specific microorganisms

532   to a specific time.

533

534   Adjusting the gLV

535   To evaluate EnDED, we simulated a time series of microbial abundances with a gLV

536   including true pairwise interactions between 50 microorganisms and adjusted it by

537   incorporating two environmental factors:

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

Eq. (2)

538   where $t$ is time, $dy(t)/dt$ is the rate of change of microbial abundances as a column vector,

539   $y(t)$ is the vector of microbial abundance at time $t$, b is the growth rate vector determined

540   through microorganism's specific growth rate functions that depend on an environmental

541   factor (see equation (4)), and $A$ is the interaction matrix.

542

543   Interaction matrix

544   In the interaction matrix $A$, each coefficient $a_{ji}$ provides the linear effect that a change in the

545   abundance of microorganism $i$ has on the growth of microorganism $j$ (Novak *et al.*, 2016).

546   We simulated the interaction coefficients $a_{ji}$ with the Klemm-Eguíluz algorithm (Klemm &

18

547 Eguíluz, 2002), which generates a modular and scale-free matrix. We also set the interaction

548 probability to 0.01, the percentage of positive coefficients to 30%, and diagonal coefficients

549 to zero. Negative diagonal coefficients $a_{ii}$ (i.e., the interaction of a microorganism with itself)

550 can represent intra-specific competition and provides the carrying capacity for each

551 microorganism, preventing its explosive growth (Haydon, 1994). We set the diagonal

552 coefficients $a_{ii} = -0.5$ to avoid excessive microbial abundances in the simulations.

553

554 <u>Two abiotic environmental factors</u>

555 We adjusted the gLV by including two environmental factors. For simplicity, we assume no

556 feedback between the microorganisms and the environmental factors. That is, the

557 environmental factors affect the growth of the microorganisms but not vice-versa. The first

558 environmental factor affects the specific growth rate of each microorganism by interacting

559 with two of their traits: optimal environmental value for growth and tolerance range of

560 environmental values. We simulated the environmental factor using a periodic sinusoidal

561 function (see equation (3)), rounded to 3 digits:

$$\epsilon(t) \triangleq round(\sin(\omega \cdot t), \text{digits} = 3), \hspace{2cm} \text{Eq. (3)}$$

562 where $t$ is the time axis (months), $\omega = (-2\pi/T)$ is the signal frequency (radians) and $T =$

563 12 is the signal periodicity (months); resulting in a signal phase shift of $T/4$ (months). While

564 the first environmental factor is considered to be "external" to the microbial community, the

565 second environmental factor is considered to be "internal", and therefore it is included in the

566 interaction matrix. The interaction coefficients between the microorganisms and the second

567 environmental factor were generated by splitting the microorganisms into two groups: the

568 second abiotic environmental factor influenced positively one half and negatively the other

569 half of the microorganisms. We obtained the interaction coefficients from two uniform

570 distributions defined to range between [-0.8, -0.2] and [0.2, 0.8] respectively. As the

571 microorganisms did not influence the abiotic factor, the corresponding interaction

572 coefficients were set to zero.

573

574 <u>Species growth rate</u>

575 The external seasonal abiotic environmental variable affects the growth rate, $g$, of each

576 microorganism. This dependency is given by:

19

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

577    where $E(t)$ is the environmental parameter that affects the microorganisms growth rate $g(t)$

578    at time $t$, $g_{max}$ is the microorganism' specific maximum growth rate that determines the

579    amplitude of the growth-rate curve, $\epsilon_{opt}$ is the microorganism' specific optimal

580    environmental value that determines the peak of the growth-rate curve, and $\sigma$ is the

581    microorganism' specific ecological tolerance (niche width) determining the environmental

582    range in which the microorganism grows, which determines the length (niche spread) of the

583    growth-rate curve. We obtained the two constant parameters $g_{max}$, and $\sigma$ for each

584    microorganism from a uniform distribution ranging between 0.3 and 1 to assure positive

585    values. The values $\epsilon_{opt}$ were drawn from a uniform distribution ranging between the minimal

586    and maximal value of the seasonal environmental factor. We defined the internal abiotic

587    environmental factor, which is included in the interaction matrix, through the same function

588    with $g_{max} = 0.8$, $\epsilon_{opt} = 0.5$, and $\sigma = 0.5$. Since the growth rates depend on the

589    environmental factor, they vary seasonally. Different microorganisms will grow better or

590    worse at different times of the year following their environmental niches. This will lead to

591    an asynchrony of their growth rate responses to the environment that will translate into an

592    asynchrony of their abundances in time.

593

594    <u>Initial abundances</u>

595    To obtain the microbial abundances in time with the adjusted gLV, we simulated the initial

596    microbial abundances with a stick-breaking process such that abundances add up to 1, using

597    the function bstick (Jackson, 1993; Legendre & Legendre, 2012), and the package vegan

598    (Oksanen *et al.*, 2019). We generated uneven initial microbial abundances without

599    introducing zeros and set the initial value for the internal abiotic environmental factor

600    included in the interaction matrix to 0.001.

601

602    <u>Species abundances in time</u>

603    Once we have set the initial conditions, we simulated microbial abundances over time by

604    solving the equations given in the adjusted gLV (see equation (2)). Start time was 0, end time

605    49.5, and sample resolution 0.5 resulting in 100 samples. We used the solver function lsoda

606    (Soetaert *et al.*, 2010). The simulated abundances in time were used to construct an

607    association network, which is referred to as the simulated network.

608

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

610    <u>Microbial abundances</u>

611    Surface water (≈ 1m depth) was sampled monthly from January 2004 to December 2013, at

612    the BBMO in the North-Western Mediterranean Sea (41°40′N 2°48′E) (Gasol *et al.*, 2016).

613    About 6L of seawater were filtered and separated into picoplankton (0.2-3 µm) and

614    nanoplankton (3-20 µm), as described in (Giner *et al.*, 2019). The DNA was extracted using

615    a phenol-chloroform standard method (Schauer *et al.*, 2003), which has been modified by

616    using Amicon units (Millipore) for purification.

617    Next, community DNA was extracted, and the 18S ribosomal RNA-gene (V4 region)

618    was amplified in (Giner *et al.*, 2019) using the primer pair TAReukFWD1 and TAReukREV3

619    (Stoeck *et al.*, 2010). The 16S ribosomal RNA-gene (V4 region) was also amplified from the

620    same DNA extracts using the primers Bakt 341F (Herlemann *et al.*, 2011) and 806R (Apprill

621    *et al.*, 2015). Amplicons were sequenced in a MiSeq platform (2x250bp) at the sequencing

622    service RTL Genomics in Lubbock, Texas. Read quality control, trimming, and inference of

623    Operational Taxonomic Units (OTUs) as Amplicon Sequence Variants (ASV) was made

624    with DADA2 v1.10.1 (Callahan *et al.*, 2016) with the maximum number of expected errors

625    (MaxEE), set to 2 and 4 for the forward and reverse reads, respectively.

626    ASV sequence abundance tables were obtained for both microbial eukaryotes and

627    prokaryotes. We subsampled both tables to the lowest sequencing depth of 4907 reads, with

628    the rrarefy function from the Vegan package in R (Oksanen *et al.*, 2019) , v2.4-2. We

629    excluded 29 nanoplankton samples (March 2004, February 2005, and May 2010 to July 2012)

630    featuring suboptimal amplicon sequencing. In these, we estimated microbial abundances

631    using seasonally aware missing value imputation by weighted moving average for time series

632    as implemented in the R package imputeTS (Moritz & Gatscha, 2017), v2.8.

633    Dislodging cells or particles and filter clogging can bias the collection of DNA in

634    either small or large organismal size fractions. To reduce the bias, we divided the sequence

635    abundance sum of the nanoplankton by the picoplankton for each ASV appearing in both size

21

636    fractions and set the picoplankton abundances to zero if the ratio exceeded 2. Likewise, we

637    set the nanoplankton abundances to zero if the ratio was below 0.5.

638

639    Taxonomic classification

640    The taxonomic classification of each ASV was inferred with the naïve Bayesian classifier

641    method (Wang *et al.*, 2007) together with the SILVA version 132 (Quast *et al.*, 2012)

642    database as implemented in DADA2 (Callahan *et al.*, 2016). In addition, eukaryotic

643    microorganisms were BLASTed (Altschul *et al.*, 1990) against the Protist Ribosomal

644    Reference database [PR2, version 4.10.0; (Guillou *et al.*, 2012)]. If the taxonomic assignment

645    for eukaryotes disagreed between SILVA and PR2, we used the PR2 classification. We

646    removed microorganisms identified as either Metazoa, or Streptophyta, plastids and

647    mitochondria. In addition, we removed Archaeas since the 341F primer is not optimal for

648    recovering this domain (McNichol *et al.*, 2021). The resulting microbial sequence abundance

649    table contained microbial eukaryotic and bacterial ASVs. Rare ASVs were removed, i.e., we

650    kept only ASVs present in more than 15% of the samples and with a sequence abundance

651    sum above 100.

652

653    Environmental factors

654    We measured environmental factors that may affect the ecosystem's dynamics. We

655    considered a total of ten contextual abiotic and biotic variables: day length (hours of light),

656    temperature (C∘), turbidity (Secchi depth m), salinity, total cholorophyll (µg/l), and inorganic[4]

657    nutrients— $PO_4^{3-}$ (µM), $NH_4^+$ (µM), $NO_2^-$ (µM), $NO_3^-$ (µM), and $SiO_2$ (µM) (Giner *et al.*,

658    2019). Water temperature and salinity were sampled in situ with a CTD (Conductivity,

659    Temperature, and Depth) measuring device. Inorganic nutrients were measured with an

660    Alliance Evolution II autoanalyzer (Grasshoff *et al.*, 2009). See (Gasol *et al.*, 2016) for

661    specific details on how other variables were measured.

662

663    **Network construction**

664    We constructed association networks from the simulated and the real microbial abundance

665    tables and environmental parameters using eLSA (Xia *et al.*, 2011, 2013). We included

666    default normalization and a z-score transformation using median and median absolute

22

667    deviation. We estimated the $p$-value with a mixed approach that performs a random

668    permutation test if the theoretical $p$-values for the comparison are below 0.05; the number of

669    iterations was 2000. Although we are aware of time-delayed interactions and that eLSA (Xia

670    *et al.*, 2011, 2013) could account for them, we considered our sampling interval as too large

671    (1 month) for inferring time-delayed associations with a solid ecological basis. Thus, in our

672    study, we focused on contemporary interactions between co-occurring microbes. For the

673    BBMO dataset, the Bonferroni false discovery rate, $q$, was calculated for all edges from the

674    $p$-values using the R function *p.adjust* (R Core Team, 2019). Lastly, we used a significance

675    threshold for the $p$ and $q$ value of 0.001 as suggested in other works (Weiss *et al.*, 2016).

676

677    **Intersection combination of EnDED—Environmentally-Driven Edge Detection**

678    **methods**

679    EnDED includes four methods: SP, OL, II, DPI (described below) and their intersection

680    combination (an ensemble approach of the four methods). We applied these methods to find

681    environmentally-driven associations of microorganisms that were within an environmental

682    triplet, as in (Lima-Mendez *et al.*, 2015). An environmental triplet is a special case of a closed

683    triplet where one of the nodes corresponds to an environmental factor and the other two nodes

684    correspond to microorganisms. We define the closed triplet, where there is an edge between

685    each pair of three nodes, as $T = \{v, w, f\}$ where $v$ and $w$ are two microorganisms, and $f$ is

686    an environmental component (see Figure 3).

687        For the intersection combination, all four individual methods must converge to the

688    same solution, i.e., if all methods classify the microbial edge as environmentally-driven, the

689    edge is removed from the network. If a microbial association is within several environmental

690    triplets, at least one of them must indicate the association as environmentally-driven. In sum,

691    the intersection combination retains an association in the network if no triplet classifies the

692    association as environmentally-driven.

693

694    <u>Sign Pattern</u>

695    The SP method (Lima-Mendez *et al.*, 2015) filters environmentally-driven edges from a

696    network in which a positive association score indicates co-occurrence, and a negative

697    association score indicates mutual exclusion. Let $s_{vw}$ be the sign of the association score of

23

698    the association between $v$ and $w$ (i.e., $s_{vw} = +$ or $s_{vw} = -$). A closed triplet $T$ has eight SP

699    combinations that group into two sets (see Figure 3). If the product of the three association

700    scores is positive, then the SP suggests that the edge between the two microorganisms is

701    environmentally-driven. Otherwise, if the product of the three association scores is negative,

702    SP does not suggest that the association is environmentally-driven.

703

704    <u>Overlap</u>

705    We have developed the OL method to support the SP for temporal data: a microbial edge

706    should be disregarded as environmentally-driven when the associations are misaligned in

707    time. Thus, OL requires the time when the association begins as well as how long the

708    associations lasts, i.e., duration or length of association in time, both determined by the

709    network construction tool eLSA (Xia *et al.*, 2011, 2013). Given an association between $v$ and

710    $w$, let $b_{vw}^{v}$ be the beginning of the association for $v$, $b_{vw}^{w}$ the beginning of the association for

711    $w$, and $d_{vw}$ be the duration of the association between $v$ and $w$. Although not used in the

712    BBMO network, OL can consider time-delays by assuming that the beginning of the

713    association is the minimum of the two beginnings, $b_{vw} = \min(b_{vw}^{v}, b_{vw}^{w})$, and the end of the

714    association is the maximum, $e_{vw} = \max(b_{vw}^{v} + d_{vw}, b_{vw}^{w} + d_{vw})$. We indicate two

715    microorganisms with $v$ and $w$, and the factor by $f$. The OL method calculates the overlap O

716    of the microbial association with the two microorganism-environment associations through

717    equation (5). As depicted in Figure 3, if $O > 60\%$, the microbial association is considered

718    environmentally-driven.

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

719    <u>Mutual Information and Conditional Mutual Information</u>

720    The method II employs two measurements: MI and CMI. The former is also used by DPI.

721    Thus, before describing the methods, we first describe the two measurements. MI is a

722    measure of the degree of statistical dependency between two variables (Margolin *et al.*,

723    2006). We first consider $\boldsymbol{v} = v_1, \dots, v_n$, $\boldsymbol{w} = w_1, \dots, w_n$, and $\boldsymbol{f} = f_1, \dots, f_n$ as discrete

724    random variables. The marginal probability of each discrete state (value) of the variable is

725    denoted by $p(v_i) = P(\boldsymbol{v} = v_i)$, the joint probability by $p(v_i, w_j)$, and $p(v_i, w_j, f_k)$, and

726    the conditional probability by $p(v_i|f_k)$, and $p(v_i, w_j|f_k)$. To obtain MI, we calculate the

24

727    entropy of $\boldsymbol{v}$ as

$$S(\boldsymbol{v}) = -\sum_{i=1}^{n} p(v_i) \log\big(p(v_i)\big), \qquad \text{Eq. (6)}$$

728    and the joint entropy of $\boldsymbol{v}$ and $\boldsymbol{w}$ as

$$S(\boldsymbol{v}, \boldsymbol{w}) = -\sum_{i=1, j=1}^{n} p(v_i, w_j) \log\Big(p(v_i, w_j)\Big), \qquad \text{Eq. (7)}$$

729    using the natural logarithm. The MI of $\boldsymbol{v}$ and $\boldsymbol{w}$ is defined through the sum of their entropies

730    subtracted by their joint entropy:

$$\text{MI}(\boldsymbol{v}; \boldsymbol{w}) = S(\boldsymbol{v}) + S(\boldsymbol{w}) - S(\boldsymbol{v}, \boldsymbol{w}) \qquad \text{Eq. (8)}$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} p(v_i, w_j) \log\left(\frac{p(v_i, w_i)}{p(v_i)p(w_j)}\right), \qquad \text{Eq. (9)}$$

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

732        The measurement CMI is the expected value of the MI of two random variables given

733    a third random variable. It is defined as

$$\text{CMI}(\boldsymbol{v}; \boldsymbol{w}|\boldsymbol{f}) = S(\boldsymbol{v}, \boldsymbol{f}) + S(\boldsymbol{w}, \boldsymbol{f}) - S(\boldsymbol{v}, \boldsymbol{w}, \boldsymbol{f}) - S(\boldsymbol{f}) \qquad \text{Eq. (10)}$$

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

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

734

735    Interaction Information

736    The II is calculated with microbial abundance and environmental data. In this study, as in

737    (Lima-Mendez *et al.*, 2015), II is computed as the difference of the CMI and MI:

$$\text{II} = \text{CMI} - \text{MI}. \qquad \text{Eq. (12)}$$

738    In other works (Ghassami & Kiyavash, 2017), the II is defined with a different sign

739    convention: $\text{II} = \text{MI} - \text{CMI}$. In our study, if II is positive, the method suggests that the

740    microbial association is not environmentally-driven. If II is negative, there is an

741 environmentally-driven association within the closed triplet. However, this method cannot

742 detect which of the three associations is indirect. In other works (Lima-Mendez *et al.*, 2015),

743 the microbial association is assumed to be environmentally-driven if II is negative, but here

744 we suggest to combine it with DPI (see below).

745

746 <u>Significance of Interaction Information</u>

747 We determined the significance of II following a strategy from (North *et al.*, 2002; Veech,

748 2012). We used a parameter-free permutation test and computed the $p$-value by randomizing

749 the environmental vector $\boldsymbol{f}$. Since the MI is independent of the environmental factor and

750 therefore remains constant, the significance of the II is the same as the CMI. Thus, we

751 determined the significance of CMI with 1000 permutations: we randomized the

752 environmental vector $\boldsymbol{f}$ and recalculated the CMI 1000 times, obtaining a $CMI_i$ with $i \in$

753 $\{1, \ldots, 1000\}$. Afterwards, we quantified with $c$ how many random $CMI_i$ were at least as

754 small as the original $CMI_i$: $c = |i: CMI_i \leq CMI_{original}, i \in \{1, \ldots, 1000\}|$. We calculated the

755 $p$-value as

$$p = \frac{c+1}{1000+1} \ . \qquad \text{Eq. (13)}$$

756

757 <u>Data Processing Inequality</u>

758 As mentioned above, the II method can detect if an indirect association exists within a triplet

759 but cannot determine which of the three associations is indirect. Thus, we added DPI to

760 EnDED. DPI states that if two components $v$ and $w$ interact only through a third component

761 $f$ (i.e., in a network $v$ and $w$ are connected through a path containing $f$ and there is no

762 alternative path between $v$ and $w$), then the MI of $v$ and $w$, $MI(\boldsymbol{v}; \boldsymbol{w})$ is smaller than

763 $MI(\boldsymbol{v}; \boldsymbol{f})$ and $MI(\boldsymbol{w}; \boldsymbol{f})$ (Cover & Thomas, 2001):

$$MI(\boldsymbol{v}; \boldsymbol{w}) \leq \min \{MI(\boldsymbol{v}; \boldsymbol{f}), MI(\boldsymbol{w}; \boldsymbol{f})\} \ . \qquad \text{Eq. (14)}$$

764 While DPI has been used in previous works on gene triplets (Margolin *et al.*, 2006), we used

765 the DPI method for environmental triplets. We compared the MI between the two

766 microorganisms with the MI between a microorganism and the environmental factor. If the

767 MI between the microorganisms is the smallest, then the method suggests that the edge is

768 environmentally-driven. This method complements the II method.

26

769

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 $v = (v_1, \ldots, v_n)$ between the lowest value ($v_{min}$) and highest value ($v_{max}$). It is a procedure implemented in other works (Meyer *et al.*, 2008). Given vector $v$ of length $n$ (that is the sample size) and number of bins $|B| = \lfloor \sqrt{n} \rfloor$, the discretized value $v_d$ of variable $v$ in vector $v$ is:

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

Eq. (15)

This equation assumes positive values. However, if $v$ contains negative values, $v_{min} < 0$, we adjust equation (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.

784

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.

792

**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

798    that associations not present in the interaction matrices, are environmentally-driven. We

799    consider P as the number of all false associations, both true positive and false negative

800    detected environmentally-driven edges: $P = TP + FN$, and N as the number of all true

801    interactions, i.e., all true negative and false positive detected environmentally-driven edges:

802    $N = TN + FP$. Then, we calculated the true positive rate (sensitivity), by dividing the number

803    of true positives by the number of all real positives: $TPR = (TP)/(P)$. Equivalently, we can

804    also calculate the true negative rate (specificity) by dividing the number of true negatives by

805    the number of all real negatives, $TNR = (TN)/(N)$. The false positive rate (fall out) is the

806    complementary to TNR, i.e. $FPR = 1 - TNR$. The positive predictive value (precision) can

807    be calculated by dividing the number of true positives by the sum of all predicted positives,

808    $PPV = (TP)/(TP + FP)$. The accuracy is calculated by dividing the sum of true positives

809    and true negatives by the sum of all real positives and real negatives, $ACC = (TP +$

810    $TN)/(P + N)$.

811

812    <u>Real Dataset</u>

813    *Literature based database*

814    The real network evaluation is limited since the true interactions and the microorganisms that

815    do not interact with each other are poorly known. We assessed true interactions known in the

816    literature based on the genus, which are compiled within the Protist Interaction Database,

817    PIDA (Bjorbækmo *et al.*, 2019). On October 15th 2019, PIDA contained 2448 interactions.

818    Although our dataset contains protists as well as bacteria, we were unable to evaluate

819    interactions between bacteria through PIDA.

820

821    *Jaccard index*

822    In ecology, the Jaccard index (Jaccard similarity coefficient) is normally used for

823    communities. Here, for each pair of microorganisms in the BBMO network, we computed

824    the Jaccard index as the number of samples in which both microorganisms occur, divided by

825    the number of samples in which at least one of the two microorganisms are present.

**Ethics approval and consent to participate**

Not applicable.

**Consent for publication**

Not applicable.

**Availability of data and material**

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.

**Competing interests**

The authors declare that they have no competing interests.

**Author's 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 GLM and KF (measurements for simulation dataset) and AKK (literature based database for real dataset). RL processed the amplicon data from BBMO generating ASV 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.

# References

AI, D., LI, X., PAN, H., CHEN, J., CRAM, J.A., & XIA, L.C. (2019) Explore mediated co-varying dynamics in microbial community using integrated local similarity and liquid association analysis. *BMC Genomics*, **20**, 185.

AITCHISON, J. (1981) A new approach to null correlations of proportions. *Journal of the International Association for Mathematical Geology*.

ALIPANAHI, B. & FREY, B.J. (2013) Network cleanup. *Nature Biotechnology*, **31**, 714–715.

ALTSCHUL, S.F., GISH, W., MILLER, W., MYERS, E.W., & LIPMAN, D.J. (1990) Basic local alignment search tool. *Journal of Molecular Biology*, **215**, 403–410.

APPRILL, A., MCNALLY, S., PARSONS, R., & WEBER, L. (2015) Minor revision to V4 region SSU rRNA 806R gene primer greatly increases detection of SAR11 bacterioplankton. *Aquatic Microbial Ecology*, **75**, 129–137.

BARZEL, B. & BARABÁSI, A.-L. (2013) Network link prediction by global silencing of indirect correlations. *Nature Biotechnology*, **31**, 720–725.

BASHAN, A., GIBSON, T.E., FRIEDMAN, J., CAREY, V.J., WEISS, S.T., HOHMANN, E.L., & LIU, Y.-Y. (2016) Universality of human microbial dynamics. *Nature*, **534**, 259–262.

BENINCÀ, E., DAKOS, V., VAN NES, E.H., HUISMAN, J., & SCHEFFER, M. (2011) Resonance of Plankton Communities with Temperature Fluctuations. *The American Naturalist*, **178**, E85–E95.

BERRY, D. & WIDDER, S. (2014) Deciphering microbial interactions and detecting keystone species with co-occurrence networks. *Frontiers in Microbiology*, **5**, 219.

BJORBÆKMO, M.F.M., EVENSTAD, A., RØSÆG, L.L., KRABBERØD, A.K., & LOGARES, R. (2019) The planktonic protist interactome: where do we stand after a century of research? *The ISME Journal*, DOI: 10.1038/s41396-019-0542-5.

BRISSON, V., SCHMIDT, J., NORTHEN, T.R., VOGEL, J.P., & GAUDIN, A. (2019) A New Method to Correct for Habitat Filtering in Microbial Correlation Networks. *Frontiers in Microbiology*, **10**, 585.

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**, 581–583.

COVER, T.M. & THOMAS, J.A. (2001) Inequalities in Information Theory. *Elements of Information Theory*.

DAM, P., FONSECA, L.L., KONSTANTINIDIS, K.T., & VOIT, E.O. (2016) Dynamic models of the complex microbial metapopulation of lake mendota. *npj Systems Biology and Applications*, **2**, 16007.

DELONG, E.F. (2009) The microbial ocean from genomes to biomes. *Nature*.

DEUTSCHMANN, I.M. (2019) *EnDED - - Environmentally-Driven Edge Detection Program*. Zenodo.

FALKOWSKI, P.G., FENCHEL, T., & DELONG, E.F. (2008) The Microbial Engines That Drive Earth's Biogeochemical Cycles. *Science*.

FAUST, K. (2019) Towards a Better Understanding of Microbial Community Dynamics through High-Throughput Cultivation and Data Integration. *mSystems*, **4**.

FAUST, K. & RAES, J. (2012) Microbial interactions: from networks to models. *Nature*

*Reviews Microbiology*, **10**, 538–550.

FAUST, K. & RAES, J. (2016) CoNet app: inference of biological association networks using Cytoscape [version 2; peer review: 2 approved]. *F1000Research*, **5**.

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*.

FEIZI, S., MARBACH, D., MÉDARD, M., & KELLIS, M. (2013) Network deconvolution as a general method to distinguish direct dependencies in networks. *Nature Biotechnology*, **31**, 726–733.

FERNANDES, A.D. & GLOOR, G.B. (2010) Mutual information is critically dependent on prior assumptions: would the correct estimate of mutual information please identify itself? *Bioinformatics*, **26**, 1135–1139.

FRIEDMAN, J. & ALM, E.J. (2012) Inferring Correlation Networks from Genomic Survey Data. *PLOS Computational Biology*, **8**, 1–11.

GASOL, J.M., CARDELÚS, C., G MORÁN, X.A., BALAGUÉ, V., FORN, I., MARRASÉ, C., MASSANA, R., PEDRÓS-ALIÓ, C., MONTSERRAT SALA, M., SIMÓ, R., VAQUÉ, D., & ESTRADA, M. (2016) Seasonal patterns in phytoplankton photosynthetic parameters and primary production at a coastal NW Mediterranean site. *Scientia Marina*.

GHASSAMI, A. & KIYAVASH, N. (2017) Interaction information for causal inference: The case of directed triangle. *2017 IEEE International Symposium on Information Theory (ISIT)*. pp. 1326–1330.

GINER, C.R., BALAGUÉ, V., KRABBERØD, A.K., FERRERA, I., REÑÉ, A., GARCÉS, E., GASOL, J.M., LOGARES, R., & MASSANA, R. (2019) Quantifying long-term recurrence in planktonic microbial eukaryotes. *Molecular Ecology*, **28**, 923–935.

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.

GONZE, D., COYTE, K.Z., LAHTI, L., & FAUST, K. (2018) Microbial communities as dynamical systems. *Current Opinion in Microbiology*, **44**, 41–49.

GRASSHOFF, K., KREMLING, K., & EHRHARDT, M. (2009) *Methods of seawater analysis*. John Wiley & Sons.

GUILLOU, L., BACHAR, D., AUDIC, S., BASS, D., BERNEY, C., BITTNER, L., BOUTTE, C., BURGAUD, G., DE VARGAS, C., DECELLE, J., DEL CAMPO, J., DOLAN, J.R., DUNTHORN, M., EDVARDSEN, B., HOLZMANN, M., KOOISTRA, W.H.C.F., LARA, E., LE BESCOT, N., LOGARES, R., MAHÉ, F., MASSANA, R., MONTRESOR, M., MORARD, R., NOT, F., PAWLOWSKI, J., PROBERT, I., SAUVADET, A.-L., SIANO, R., STOECK, T., VAULOT, D., ZIMMERMANN, P., & CHRISTEN, R. (2012) The Protist Ribosomal Reference database (PR$^2$): a catalog of unicellular eukaryote Small Sub-Unit rRNA sequences with curated taxonomy. *Nucleic Acids Research*, **41**, D597–D604.

HAYDON, D. (1994) Pivotal Assumptions Determining the Relationship between Stability and Complexity: An Analytical Synthesis of the Stability-Complexity Debate. *The American Naturalist*, **144**, 14–29.

HERLEMANN, D.P., LABRENZ, M., JÜRGENS, K., BERTILSSON, S., WANIEK, J.J., & ANDERSSON, A.F. (2011) Transitions in bacterial communities along the 2000 km salinity gradient of the Baltic Sea. *The ISME Journal*, **5**, 1571–1579.

964 JACKSON, D.A. (1993) Stopping Rules in Principal Components Analysis: A
965     Comparison of Heuristical and Statistical Approaches. *Ecology*, **74**, 2204–
966     2214.
967 JANG, I.S., MARGOLIN, A., & CALIFANO, A. (2013) hARACNe: improving the accuracy
968     of regulatory model reverse engineering via higher-order data processing
969     inequality tests. *Interface Focus*, **3**, 20130011.
970 KALLMEYER, J., POCKALNY, R., ADHIKARI, R.R., SMITH, D.C., & D'HONDT, S. (2012)
971     Global distribution of microbial abundance and biomass in subseafloor
972     sediment. *Proceedings of the National Academy of Sciences*, **109**, 16213–
973     16216.
974 KETTLE, H., HOLTROP, G., LOUIS, P., & FLINT, H.J. (2018) microPop: Modelling
975     microbial populations and communities in R. *Methods in Ecology and
976     Evolution*, **9**, 399–409.
977 KLEMM, K. & EGUÍLUZ, V.M. (2002) Growing scale-free networks with small-world
978     behavior. *Physical Review E*, **65**, 057102.
979 KRABBERØD, A.K., BJORBÆKMO, M.F.M., SHALCHIAN-TABRIZI, K., & LOGARES, R.
980     (2017) Exploring the oceanic microeukaryotic interactome with metaomics
981     approaches. *Aquatic Microbial Ecology*, **79**, 1–12.
982 KURTZ, Z.D., BONNEAU, R., & MÜLLER, C.L. (2019) Disentangling microbial
983     associations from hidden environmental and technical factors via latent
984     graphical models. *bioRxiv*, DOI: 10.1101/2019.12.21.885889.
985 KURTZ, Z.D., MÜLLER, C.L., MIRALDI, E.R., LITTMAN, D.R., BLASER, M.J., & BONNEAU,
986     R.A. (2015) Sparse and Compositionally Robust Inference of Microbial
987     Ecological Networks. *PLOS Computational Biology*.
988 LATORRE, F., DEUTSCHMANN, I.M., LABARRE, A., OBIOL, A., KRABBERØD, A.K.,
989     PELLETIER, E., SIERACKI, M.E., CRUAUD, C., JAILLON, O., MASSANA, R., &
990     LOGARES, R. (2021) Niche adaptation promoted the evolutionary
991     diversification of tiny ocean predators. *Proc Natl Acad Sci USA*, **118**,
992     e2020955118.
993 LAYEGHIFARD, M., HWANG, D.M., & GUTTMAN, D.S. (2017) Disentangling Interactions
994     in the Microbiome: A Network Perspective. *Trends in Microbiology*.
995 LEGENDRE, P. & LEGENDRE, L.F. (2012) *Numerical ecology*, vol. 24. Elsevier.
996 LI, C., LIM, K.M.K., CHNG, K.R., & NAGARAJAN, N. (2016) Predicting microbial
997     interactions through computational approaches. *Methods*.
998 LIMA-MENDEZ, G., FAUST, K., HENRY, N., DECELLE, J., COLIN, S., CARCILLO, F.,
999     CHAFFRON, S., IGNACIO-ESPINOSA, J.C., ROUX, S., VINCENT, F., BITTNER, L.,
1000    DARZI, Y., WANG, J., AUDIC, S., BERLINE, L., BONTEMPI, G., CABELLO, A.M.,
1001    COPPOLA, L., CORNEJO-CASTILLO, F.M., d'OVIDIO, F., DE MEESTER, L., FERRERA,
1002    I., GARET-DELMAS, M.-J., GUIDI, L., LARA, E., PESANT, S., ROYO-LLONCH, M.,
1003    SALAZAR, G., SÁNCHEZ, P., SEBASTIAN, M., SOUFFREAU, C., DIMIER, C.,
1004    PICHERAL, M., SEARSON, S., KANDELS-LEWIS, S., GORSKY, G., NOT, F., OGATA,
1005    H., SPEICH, S., STEMMANN, L., WEISSENBACH, J., WINCKER, P., ACINAS, S.G.,
1006    SUNAGAWA, S., BORK, P., SULLIVAN, M.B., KARSENTI, E., BOWLER, C., DE
1007    VARGAS, C., & RAES, J. (2015) Determinants of community structure in the
1008    global plankton interactome. *Science*, **348**, 1262073.
1009 LOCEY, K.J. & LENNON, J.T. (2016) Scaling laws predict global microbial diversity.
1010    *Proceedings of the National Academy of Sciences*, **113**, 5970–5975.

1011 Lv, X., Zhao, K., Xue, R., Liu, Y., Xu, J., & Ma, B. (2019) Strengthening Insights in
1012       Microbial Ecological Networks from Theory to Applications. *mSystems*, **4**,
1013       e00124-19.
1014 Mandakovic, D., Rojas, C., Maldonado, J., Latorre, M., Travisany, D., Delage,
1015       E., Bihouée, A., Jean, G., Díaz, F.P., Fernández-Gómez, B., Cabrera, P.,
1016       Gaete, A., Latorre, C., Gutiérrez, R.A., Maass, A., Cambiazo, V.,
1017       Navarrete, S.A., Eveillard, D., & González, M. (2018) Structure and co-
1018       occurrence patterns in microbial communities under acute environmental
1019       stress reveal ecological factors fostering resilience. *Scientific Reports*, **8**,
1020       5875.
1021 Margolin, A.A., Nemenman, I., Basso, K., Wiggins, C., Stolovitzky, G., Favera,
1022       R.D., & Califano, A. (2006) ARACNE: An Algorithm for the Reconstruction
1023       of Gene Regulatory Networks in a Mammalian Cellular Context. *BMC
1024       Bioinformatics*, **7**, S7.
1025 McCarren, J., Becker, J.W., Repeta, D.J., Shi, Y., Young, C.R., Malmstrom, R.R.,
1026       Chisholm, S.W., & DeLong, E.F. (2010) Microbial community
1027       transcriptomes reveal microbes and metabolic pathways associated with
1028       dissolved organic matter turnover in the sea. *Proceedings of the National
1029       Academy of Sciences*, **107**, 16420–16427.
1030 McNichol, J., Berube, P.M., Biller, S.J., Fuhrman, J.A., & Gilbert, J.A. (2021)
1031       Evaluating and Improving Small Subunit rRNA PCR Primer Coverage for
1032       Bacteria, Archaea, and Eukaryotes Using Metagenomes from Global Ocean
1033       Surveys. *mSystems*, **6**, e00565-21.
1034 Meyer, P.E., Lafitte, F., & Bontempi, G. (2008) minet: A R/Bioconductor Package
1035       for Inferring Large Transcriptional Networks Using Mutual Information.
1036       *BMC Bioinformatics*, **9**, 461.
1037 Moritz, S. & Gatscha, S. (2017) *imputeTS: Time Series Missing Value Imputation*.
1038 North, B.V., Curtis, D., & Sham, P.C. (2002) A Note on the Calculation of Empirical
1039       P Values from Monte Carlo Procedures. *The American Journal of Human
1040       Genetics*, **71**, 439–441.
1041 Novak, M., Yeakel, J.D., Noble, A.E., Doak, D.F., Emmerson, M., Estes, J.A.,
1042       Jacob, U., Tinker, M.T., & Wootton, J.T. (2016) Characterizing Species
1043       Interactions to Understand Press Perturbations: What Is the Community
1044       Matrix? *Annual Review of Ecology, Evolution, and Systematics*, **47**, 409–
1045       432.
1046 Oksanen, J., Blanchet, F.G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D.,
1047       Minchin, P.R., O'Hara, R.B., Simpson, G.L., Solymos, P., Stevens, M.H.H.,
1048       Szoecs, E., & Wagner, H. (2019) *vegan: Community Ecology Package*.
1049 Pascual-García, A., Tamames, J., & Bastolla, U. (2014) Bacteria dialog with Santa
1050       Rosalia: Are aggregations of cosmopolitan bacteria mainly explained by
1051       habitat filtering or by ecological interactions? *BMC Microbiology*, **14**, 284.
1052 Pinedo-González, P., West, A.J., Tovar-Sánchez, A., Duarte, C.M., Marañón, E.,
1053       Cermeño, P., González, N., Sobrino, C., Huete-Ortega, M., Fernández, A.,
1054       López-Sandoval, D.C., Vidal, M., Blasco, D., Estrada, M., & Sañudo-
1055       Wilhelmy, S.A. (2015) Surface distribution of dissolved trace metals in the
1056       oligotrophic ocean and their influence on phytoplankton biomass and
1057       productivity. *Global Biogeochemical Cycles*, **29**, 1763–1781.

1058 QUAST, C., PRUESSE, E., YILMAZ, P., GERKEN, J., SCHWEER, T., YARZA, P., PEPLIES, J.,
1059 & GLÖCKNER, F.O. (2012) The SILVA ribosomal RNA gene database project:
1060 improved data processing and web-based tools. *Nucleic Acids Research*, **41**,
1061 D590–D596.
1062 R CORE TEAM (2019) *R: A Language and Environment for Statistical Computing*.
1063 Vienna, Austria: R Foundation for Statistical Computing.
1064 RÖTTJERS, L. & FAUST, K. (2018) From hairballs to hypotheses–biological insights
1065 from microbial networks. *FEMS Microbiology Reviews*, **42**, 761–780.
1066 SCHAUER, M., BALAGUÉ, V., PEDRÓS-ALIÓ, C., & MASSANA, R. (2003) Seasonal changes
1067 in the taxonomic composition of bacterioplankton in a coastal oligotrophic
1068 system. *Aquatic Microbial Ecology*, **31**, 163–174.
1069 SOETAERT, K., PETZOLDT, T., & SETZER, R.W. (2010) Solving Differential Equations in
1070 R: Package deSolve. *Journal of Statistical Software*, **33**, 1–25.
1071 STEIN, R.R., BUCCI, V., TOUSSAINT, N.C., BUFFIE, C.G., RÄTSCH, G., PAMER, E.G.,
1072 SANDER, C., & XAVIER, J.B. (2013) Ecological Modeling from Time-Series
1073 Inference: Insight into Dynamics and Stability of Intestinal Microbiota. *PLOS*
1074 *Computational Biology*, **9**, 1–11.
1075 STOECK, T., BASS, D., NEBEL, M., CHRISTEN, R., JONES, M.D.M., BREINER, H.-W., &
1076 RICHARDS, T.A. (2010) Multiple marker parallel tag environmental DNA
1077 sequencing reveals a highly complex eukaryotic community in marine anoxic
1078 water. *Molecular Ecology*, **19**, 21–31.
1079 TACKMANN, J., RODRIGUES, J.F.M., & VON MERING, C. (2019) Rapid Inference of
1080 Direct Interactions in Large-Scale Ecological Networks from Heterogeneous
1081 Microbial Sequencing Data. *Cell Systems*, **9**, 286-296.e8.
1082 THE HUMAN MICROBIOME PROJECT CONSORTIUM: HUTTENHOWER, C., GEVERS, D.,
1083 KNIGHT, R., ABUBUCKER, S., BADGER, J.H., CHINWALLA, A.T., CREASY, H.H.,
1084 EARL, A.M., FITZGERALD, M.G., FULTON, R.S., GIGLIO, M.G., HALLSWORTH-
1085 PEPIN, K., LOBOS, E.A., MADUPU, R., MAGRINI, V., MARTIN, J.C., MITREVA, M.,
1086 MUZNY, D.M., SODERGREN, E.J., VERSALOVIC, J., WOLLAM, A.M., WORLEY, K.C.,
1087 WORTMAN, J.R., YOUNG, S.K., ZENG, Q., AAGAARD, K.M., ABOLUDE, O.O.,
1088 ALLEN-VERCOE, E., ALM, E.J., ALVARADO, L., ANDERSEN, G.L., ANDERSON, S.,
1089 APPELBAUM, E., ARACHCHI, H.M., ARMITAGE, G., ARZE, C.A., AYVAZ, T., BAKER,
1090 C.C., BEGG, L., BELACHEW, T., BHONAGIRI, V., BIHAN, M., BLASER, M.J., BLOOM,
1091 T., BONAZZI, V., PAUL BROOKS, J., BUCK, G.A., BUHAY, C.J., BUSAM, D.A.,
1092 CAMPBELL, J.L., CANON, S.R., CANTAREL, B.L., CHAIN, P.S.G., CHEN, I.-M.A.,
1093 CHEN, L., CHHIBBA, S., CHU, K., CIULLA, D.M., CLEMENTE, J.C., CLIFTON, S.W.,
1094 CONLAN, S., CRABTREE, J., CUTTING, M.A., DAVIDOVICS, N.J., DAVIS, C.C.,
1095 DESANTIS, T.Z., DEAL, C., DELEHAUNTY, K.D., DEWHIRST, F.E., DEYCH, E.,
1096 DING, Y., DOOLING, D.J., DUGAN, S.P., MICHAEL DUNNE, W., SCOTT DURKIN, A.,
1097 EDGAR, R.C., ERLICH, R.L., FARMER, C.N., FARRELL, R.M., FAUST, K.,
1098 FELDGARDEN, M., FELIX, V.M., FISHER, S., FODOR, A.A., FORNEY, L.J., FOSTER,
1099 L., DI FRANCESCO, V., FRIEDMAN, J., FRIEDRICH, D.C., FRONICK, C.C., FULTON,
1100 L.L., GAO, H., GARCIA, N., GIANNOUKOS, G., GIBLIN, C., GIOVANNI, M.Y.,
1101 GOLDBERG, J.M., GOLL, J., GONZALEZ, A., GRIGGS, A., GUJJA, S., KINDER HAAKE,
1102 S., HAAS, B.J., HAMILTON, H.A., HARRIS, E.L., HEPBURN, T.A., HERTER, B.,
1103 HOFFMANN, D.E., HOLDER, M.E., HOWARTH, C., HUANG, K.H., HUSE, S.M.,
1104 IZARD, J., JANSSON, J.K., JIANG, H., JORDAN, C., JOSHI, V., KATANCIK, J.A.,

KEITEL, W.A., KELLEY, S.T., KELLS, C., KING, N.B., KNIGHTS, D., KONG, H.H., KOREN, O., KOREN, S., KOTA, K.C., KOVAR, C.L., KYRPIDES, N.C., LA ROSA, P.S., LEE, S.L., LEMON, K.P., LENNON, N., LEWIS, C.M., LEWIS, L., LEY, R.E., LI, K., LIOLIOS, K., LIU, B., LIU, Y., LO, C.-C., LOZUPONE, C.A., DWAYNE LUNSFORD, R., MADDEN, T., MAHURKAR, A.A., MANNON, P.J., MARDIS, E.R., MARKOWITZ, V.M., MAVROMATIS, K., McCORRISON, J.M., McDONALD, D., McEWEN, J., McGUIRE, A.L., McINNES, P., MEHTA, T., MIHINDUKULASURIYA, K.A., MILLER, J.R., MINX, P.J., NEWSHAM, I., NUSBAUM, C., O'LAUGHLIN, M., ORVIS, J., PAGANI, I., PALANIAPPAN, K., PATEL, S.M., PEARSON, M., PETERSON, J., PODAR, M., POHL, C., POLLARD, K.S., POP, M., PRIEST, M.E., PROCTOR, L.M., QIN, X., RAES, J., RAVEL, J., REID, J.G., RHO, M., RHODES, R., RIEHLE, K.P., RIVERA, M.C., RODRIGUEZ-MUELLER, B., ROGERS, Y.-H., ROSS, M.C., RUSS, C., SANKA, R.K., SANKAR, P., FAH SATHIRAPONGSASUTI, J., SCHLOSS, J.A., SCHLOSS, P.D., SCHMIDT, T.M., SCHOLZ, M., SCHRIML, L., SCHUBERT, A.M., SEGATA, N., SEGRE, J.A., SHANNON, W.D., SHARP, R.R., SHARPTON, T.J., SHENOY, N., SHETH, N.U., SIMONE, G.A., SINGH, I., SMILLIE, C.S., SOBEL, J.D., SOMMER, D.D., SPICER, P., SUTTON, G.G., SYKES, S.M., TABBAA, D.G., THIAGARAJAN, M., TOMLINSON, C.M., TORRALBA, M., TREANGEN, T.J., TRUTY, R.M., VISHNIVETSKAYA, T.A., WALKER, J., WANG, L., WANG, Z., WARD, D.V., WARREN, W., WATSON, M.A., WELLINGTON, C., WETTERSTRAND, K.A., WHITE, J.R., WILCZEK-BONEY, K., WU, Y., WYLIE, K.M., WYLIE, T., YANDAVA, C., YE, L., YE, Y., YOOSEPH, S., YOUMANS, B.P., ZHANG, L., ZHOU, Y., ZHU, Y., ZOLOTH, L., ZUCKER, J.D., BIRREN, B.W., GIBBS, R.A., HIGHLANDER, S.K., METHÉ, B.A., NELSON, K.E., PETROSINO, J.F., WEINSTOCK, G.M., WILSON, R.K., & WHITE, O. (2012) Structure, function and diversity of the healthy human microbiome. *Nature*, **486**, 207–214.

VALLINA, S.M., MARTINEZ-GARCIA, R., SMITH, S.L., & BONACHELA, J.A. (2019) Models in Microbial Ecology. *Encyclopedia of Microbiology (Fourth Edition)*, Fourth Edition ed. (Schmidt, T.M. ed). Oxford: Academic Press, pp. 211–246.

VEECH, J.A. (2012) Significance testing in ecological null models. *Theoretical Ecology*, **5**, 611–616.

VERNY, L., SELLA, N., AFFELDT, S., SINGH, P.P., & ISAMBERT, H. (2017) Learning causal networks with latent variables from multivariate information in genomic data. *PLOS Computational Biology*, **13**, 1–25.

VILLAVERDE, A.F., BECKER, K., & BANGA, J.R. (2018) PREMER: A Tool to Infer Biological Networks. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, **15**, 1193–1202.

VILLAVERDE, A.F., ROSS, J., MORÁN, F., & BANGA, J.R. (2014) MIDER: Network Inference with Mutual Information Distance and Entropy Reduction. *PLOS ONE*, **9**, 1–15.

WANG, Q., GARRITY, G.M., TIEDJE, J.M., & COLE, J.R. (2007) Naïve Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Applied and Environmental Microbiology*, **73**, 5261–5267.

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**, 1669–1681.

1152 WHITMAN, W.B., COLEMAN, D.C., & WIEBE, W.J. (1998) Prokaryotes: The unseen
1153       majority. *Proceedings of the National Academy of Sciences*, **95**, 6578–6583.
1154 WORDEN, A.Z., FOLLOWS, M.J., GIOVANNONI, S.J., WILKEN, S., ZIMMERMAN, A.E., &
1155       KEELING, P.J. (2015) Rethinking the marine carbon cycle: Factoring in the
1156       multifarious lifestyles of microbes. *Science*, **347**.
1157 XIA, L.C., AI, D., CRAM, J., FUHRMAN, J.A., & SUN, F. (2013) Efficient statistical
1158       significance approximation for local similarity analysis of high-throughput
1159       time series data. *Bioinformatics*, **29**, 230–237.
1160 XIA, L.C., STEELE, J.A., CRAM, J.A., CARDON, Z.G., SIMMONS, S.L., VALLINO, J.J.,
1161       FUHRMAN, J.A., & SUN, F. (2011) Extended local similarity analysis (eLSA) of
1162       microbial community and other time series data with replicates. *BMC
1163       Systems Biology*, **5**, S15.
1164 XIAO, Y., ANGULO, M.T., FRIEDMAN, J., WALDOR, M.K., WEISS, S.T., & LIU, Y.-Y. (2017)
1165       Mapping the ecological networks of microbial communities. *Nature
1166       Communications*, **8**, 2042.
1167 YANG, Y., CHEN, N., & CHEN, T. (2017) Inference of Environmental Factor-Microbe
1168       and Microbe-Microbe Associations from Metagenomic Data Using a
1169       Hierarchical Bayesian Statistical Model. *Cell Systems*, **4**, 129-137.e5.
1170 YOON, H.S., PRICE, D.C., STEPANAUSKAS, R., RAJAH, V.D., SIERACKI, M.E., WILSON,
1171       W.H., YANG, E.C., DUFFY, S., & BHATTACHARYA, D. (2011) Single-Cell
1172       Genomics Reveals Organismal Interactions in Uncultivated Marine Protists.
1173       *Science*, **332**, 714–717.
1174 ZHAO, J., ZHOU, Y., ZHANG, X., & CHEN, L. (2016) Part mutual information for
1175       quantifying direct associations in networks. *Proceedings of the National
1176       Academy of Sciences*, **113**, 5130–5135.
1177 ZOPPOLI, P., MORGANELLA, S., & CECCARELLI, M. (2010) TimeDelay-ARACNE:
1178       Reverse engineering of gene networks from time-course data by an
1179       information theoretic approach. *BMC Bioinformatics*, **11**, 154.
1180
1181

# Figures

Figure 1: **Evaluation of EnDED: intersection combination and individual methods on simulated networks**. Using 1000 simulated networks, and 1000 simulated networks incorporating noise, we evaluated EnDED's performance. Plot A) displays the evaluation measurements true positive rate (TRP), true negative rate (TNR), accuracy (ACC), and positive predictive value (PPV) for each individual method, i.e., Sign Pattern (SP), Overlap (OL), Interaction Information (II), and Data Processing Inequality (DPI), as well as the intersection combination (Combi). SP and OL perform best according to TRP and ACC, while the intersection combination performs best according to TNR. All methods performed well according to PPV. The intersection combination, DPI and II performed better on noisy data according to TNR because less edges were removed along with less true interactions. Plot B) displays the ROC curve for each environmentally-driven edge detection method as well as their intersection combination.

Figure 2: **Quantification of environmentally-driven associations in the BBMO network.** For A) the first column shows the number and 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 and 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 $NO_3^-$ (NO3), $SiO_2$ (Si), and $NO_2^-$ (NO2). The last column shows the number and fraction 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 and fraction 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 (third row) and negative (fourth 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.

Figure 3: **EnDED Methods Overview**. EnDED is an implementation of four methods aiming to determine whether an edge between two microorganisms is indirect through the action of an environmental factor. The four methods are: Sign Pattern, Overlap, Interaction Information, and Data Processing Inequality (see Methods). Each method can be used individually or in combination. Here, we show the intersection combination approach, i.e., only if all methods classify an edge as indirect, it is removed from the network. Otherwise, the edge is classified as not indirect and kept in the network.

1213 # Tables

1214 Table 1: **Jaccard index of edges.** The BBMO network before applying EnDED contained 29820 edges of which 2488
1215 (8.3%) were environmentally-driven (indirect). Considering the Jaccard index for these indirect edges, 688 (27.7% of indirect
1216 edges) score above 50%, and 1800 (72.3%) score below or equal to 50%. In contrast, 61.1% of edges not considered as
1217 indirect have a Jaccard index above 50%, and 38.9% of all not indirect edges have a Jaccard index equal or below 50%.
1218

|  | All edges | 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%) |

1219

1220 Table 2: **Interactions found in the BBMO network that have been reported in the literature.** The table mentions whether
1221 or not the associations were removed or kept by EnDED via the combination interaction approach. For example, the
1222 association between the ASVs classified as Dia. Thalassiosira and ASVs classified as F. unknown Flavobacteriia has been
1223 found 17 times in the network: 4 were removed and 13 were kept.
1224

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

*Abbreviations indicate Dia - Diatomea; Dino - Dinoflagellata; C - Ciliophora; F - Flavobacteriia; ID in PIDA refers to the number PIDA gave to an interaction described in the literature.*

1225

1226 # Supplementary Material

1227 Supplementary Table S1: **Comparison between methods on correctly detecting false associations.** We
1228 computed the fraction (in percentage) of correctly detected false associations for each of the 1000 simulated
1229 datasets. There are only few edges that are detected by only one approach (first four rows). The most prominent
1230 groupings are highlighted in gray, e.g., SP, OL, and II agree on average on a third of edges. Combi refers to
1231 intersection combination of all four methods, SP to Sign Pattern, OL to Overlap, II to Interaction Information, and
1232 DPI to Data Processing Inequality. Less prominent groupings are aggregated with others.

| Statistic | Minimum | 1st Quartile | Median | Mean | 2nd Quartile | Maximum |
|---|---|---|---|---|---|---|
| SP | 0 | 0 | 0.2 | 0.3 | 0.5 | 3.7 |
| OL | 0 | 0 | 0.1 | 0.2 | 0.3 | 2.0 |
| II | 0 | 0.7 | 1.3 | 1.4 | 2.0 | 6.0 |
| DPI | 0 | 0.1 | 0.3 | 0.4 | 0.6 | 2.6 |
| SP and OL | 4.9 | 12.2 | 14.9 | 15.0 | 17.5 | 30.0 |
| SP, OL, and II | 19.1 | 29.5 | 32.6 | 32.8 | 36.2 | 49.6 |
| SP, OL, and DPI | 2.6 | 7.1 | 8.9 | 9.1 | 10.8 | 22.1 |
| SP, OL, II, DPI, and Combi | 22.4 | 32.1 | 35.6 | 35.5 | 38.6 | 48.6 |
| other | 0.4 | 3.3 | 4.9 | 5.1 | 6.6 | 15.4 |

1233

1234 Supplementary Table S2: **Performance of environmentally-driven edge detection methods on simulated networks.**
1235 These include 50 microorganisms and 1225 possible associations. Values display median (standard deviation) for simulated
1236 networks and simulated networks incorporating noise. Combi refers to intersection combination of all four methods, SP to
1237 Sign Pattern, OL to Overlap, II to Interaction Information, and DPI to Data Processing Inequality. The methods with highest
1238 (TP, TN, TPR, TNR, PPV, ACC) or lowest (FP, FN, FPR) median, respectively, are highlighted in gray.

| Method | Combi | SP | OL | II | DPI |
|---|---|---|---|---|---|
| without noise | | | | | |
| number of nodes | 50 (0.045) | 47 (6.6) | 48 (5.6) | 50 (0.94) | 50 (0.1) |
| number of edges | 737 (50) | 140 (52) | 144 (58) | 354 (67) | 601 (60) |
| TP | 332 (47) | 893 (64) | 888 (69) | 696 (72) | 459 (53) |
| TN | 45 (5.1) | 8 (4.3) | 9 (4.7) | 24 (5.8) | 37 (5.5) |
| FP | 15 (4.6) | 51 (5.8) | 51 (6.2) | 36 (6.4) | 23 (5.2) |
| FN | 692 (48) | 131 (49) | 136 (54) | 330 (63) | 564 (56) |
| TPR | 0.32 (0.04) | 0.87 (0.05) | 0.87 (0.05) | 0.68 (0.06) | 0.45 (0.05) |
| TNR | 0.75 (0.07) | 0.14 (0.07) | 0.15 (0.08) | 0.4 (0.10) | 0.62 (0.08) |
| FPR | 0.25 (0.07) | 0.86 (0.07) | 0.85 (0.08) | 0.6 (0.10) | 0.38 (0.08) |
| PPV | 0.96 (0.011) | 0.95 (0.005) | 0.95 (0.005) | 0.95 (0.007) | 0.95 (0.009) |
| ACC | 0.35 (0.04) | 0.83 (0.04) | 0.83 (0.048) | 0.66 (0.057) | 0.46 (0.046) |
| with noise | | | | | |
| number of nodes | 50 (0.08) | 47 (5.6) | 48 (4.9) | 50 (0.47) | 50 (0.12) |
| number of edges | 828 (56) | 144 (53) | 149 (59) | 428 (79) | 717 (73) |
| TP | 219 (48) | 864 (69) | 860 (72) | 605 (81) | 324 (64) |
| TN | 49 (5) | 9 (4.6) | 9 (4.9) | 29 (6.3) | 42 (5.8) |
| FP | 10 (3.9) | 50 (6.1) | 50 (6.4) | 30 (6.6) | 17 (5.1) |
| FN | 779 (53) | 137 (50) | 139 (55) | 398 (75) | 674 (69) |
| TPR | 0.22 (0.05) | 0.86 (0.05) | 0.86 (0.06) | 0.6 (0.08) | 0.32 (0.06) |
| TNR | 0.84 (0.07) | 0.15 (0.08) | 0.16 (0.08) | 0.49 (0.1) | 0.72 (0.09) |
| FPR | 0.16 (0.07) | 0.85 (0.08) | 0.84 (0.08) | 0.51 (0.1) | 0.28 (0.09) |
| PPV | 0.96 (0.014) | 0.95 (0.005) | 0.95 (0.005) | 0.95 (0.007) | 0.95 (0.012) |
| ACC | 0.25 (0.04) | 0.82 (0.05) | 0.82 (0.05) | 0.6 (0.07) | 0.34 (0.06) |

*SP - Sign Pattern; OL - Overlap; II - Interaction Information; DPI - Data Processing Inequality; Combi-intersection combination*

1239

1240 Supplementary Table S3: **Number of triplets a microbial edge is part of in the BBMO network.** SP and OL not listed
1241 below because they remove 100% of microbial associations that are within at least one triplet. The total number of edges
1242 (all) is given along the number of positive (pos) and negative (neg) edges. Combi refers to intersection combination of all
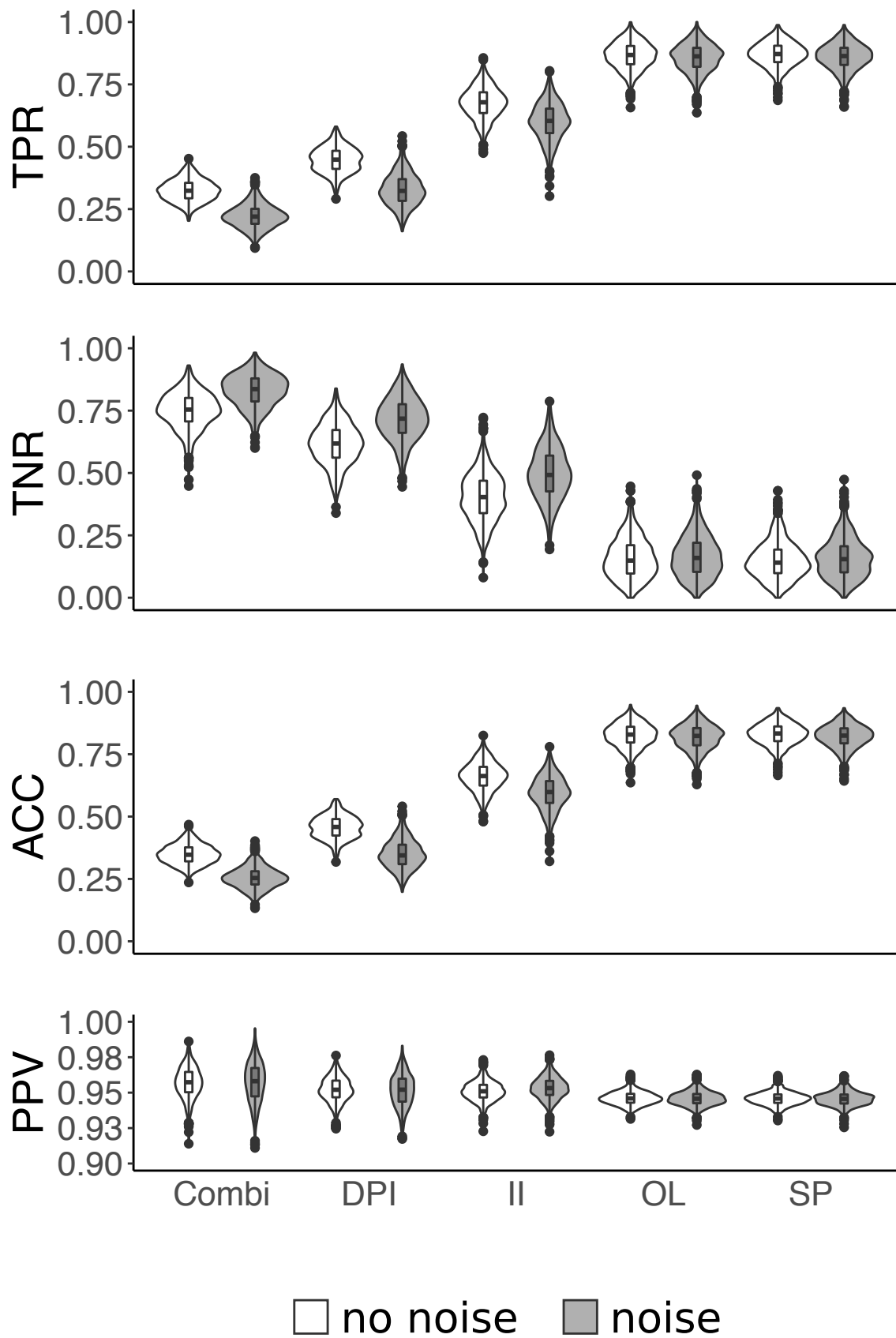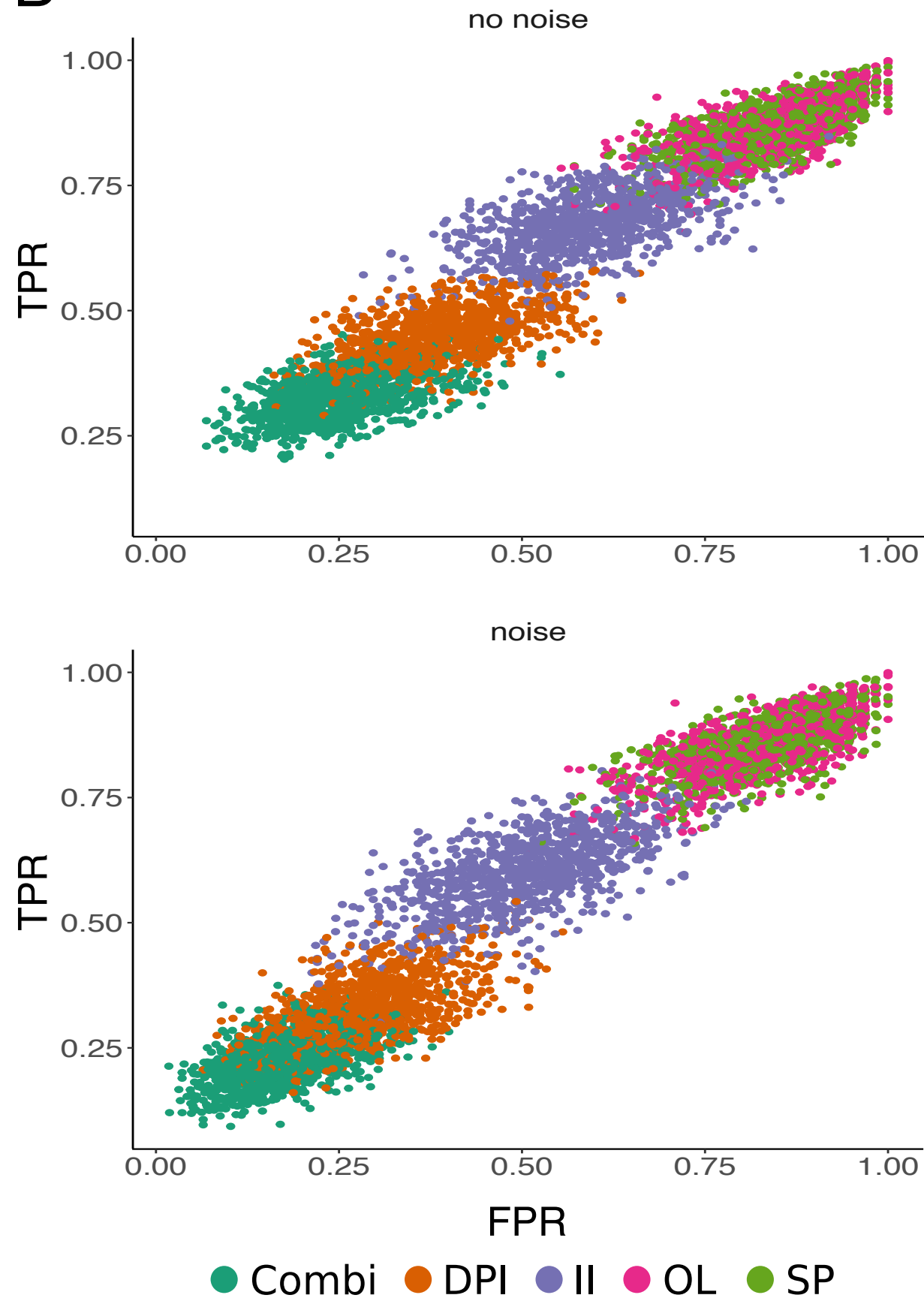1243 four methods, II to Interaction Information, and DPI to Data Processing Inequality.
1244

| Triplets | all | pos (%) | neg (%) | Combi (%) | II (%) | DPI (%) |
|---|---|---|---|---|---|---|
| 0 | 4 590 | 4 124 (89.8) | 466 (10.2) | NA | NA | NA |
| 1 | 16 193 | 13 369 (82.6) | 2 824 (17.4) | 1 276 (7.9) | 3 851 (23.8) | 4 560 (28.2) |
| 2 | 8 266 | 6 404 (77.5) | 1 862 (22.5) | 1 048 (12.7) | 3 335 (40.3) | 2 585 (31.3) |
| 3 | 667 | 484 (72.6) | 183 (27.4) | 140 (21.0) | 388 (58.2) | 222 (33.3) |
| 4 | 81 | 56 (69.1) | 25 (30.9) | 22 (27.2) | 75 (92.6) | 25 (30.9) |
| 5 | 22 | 20 (90.9) | 2 (9.1) | 2 (9.1) | 22 (100) | 2 (9.1) |
| 6 | 1 | 1 (100) | NA | NA | 1 (100) | NA |

1245

1246 Supplementary Table S4: **The BBMO network based on real data.** The BBMO network contained bacteria (B) and
1247 eukaryotes (E) from the picoplankton (p) and nanoplankton (n). This table summarizes the number and fraction of microbial
1248 associations classified by EnDED as environmentally-driven. Combi refers to the intersection combination of all four
1249 methods, II to Interaction Information, and DPI to Data Processing Inequality. Both methods, Sign Pattern and Overlap, are
1250 not shown because both remove all microbial edges found in at least one triplet. For example (last row), 349 (14.9%)
1251 associations between bacteria from the picoplankton with eukaryotes from the nanoplankton were classified by intersection
1252 combination as environmentally-driven (indirect), II classified 30.6% and DPI 37.2% as environmentally-driven.
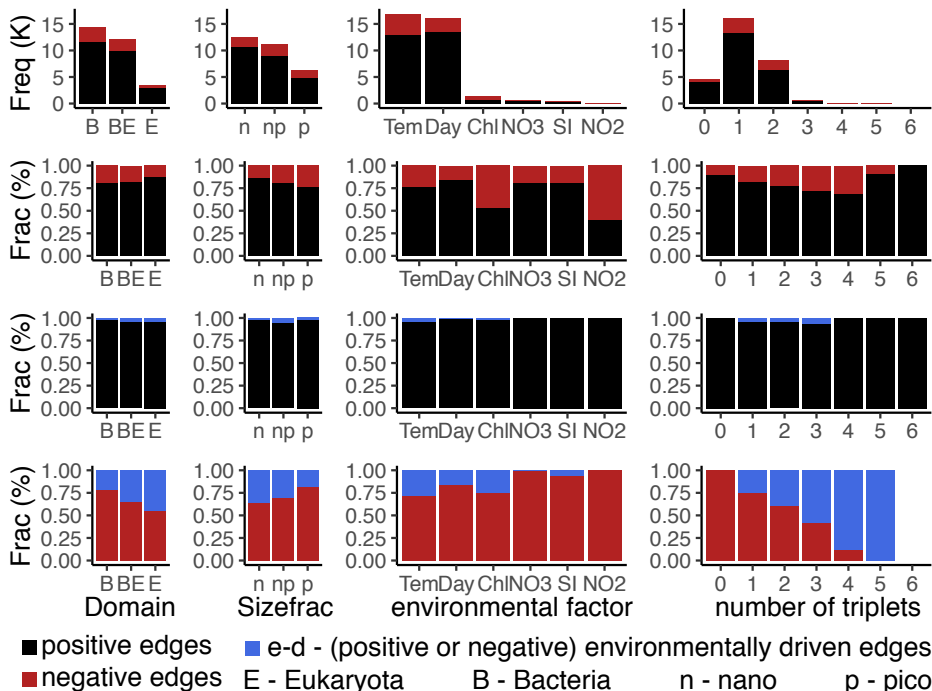1253

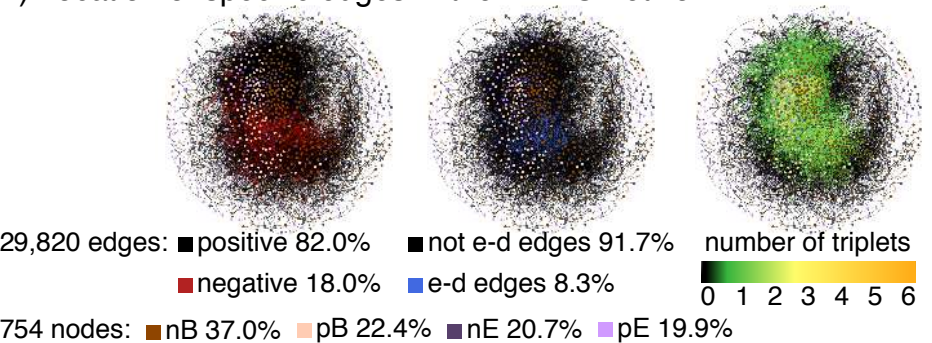| Type | edges | positive | negative | triplets | Combi | II | DPI |
|---|---|---|---|---|---|---|---|
| nB | 6 377 | 5 453 (85.5) | 924 (14.5) | 5 150 (80.8) | 376 (5.9) | 1 512 (23.7) | 1 080 (16.9) |
| n+pB | 5 191 | 4 069 (78.4) | 1 122 (21.6) | 4 824 (92.9) | 440 (8.5) | 1 381 (26.6) | 1 678 (32.3) |
| pB | 2 832 | 2 053 (72.5) | 779 (27.5) | 2 160 (76.3) | 125 (4.4) | 569 (20.1) | 631 (22.3) |
| nE | 1 319 | 1 163 (88.2) | 156 (11.8) | 1 016 (77.0) | 113 (8.6) | 350 (26.5) | 254 (19.3) |
| n+pE | 1 165 | 976 (83.8) | 189 (16.2) | 1 006 (86.4) | 158 (13.6) | 353 (30.3) | 370 (31.8) |
| pE | 895 | 820 (91.6) | 75 (8.4) | 543 (60.7) | 44 (4.9) | 153 (17.1) | 113 (12.6) |
| nB+E | 4 703 | 4 080 (86.8) | 623 (13.2) | 4 120 (87.6) | 438 (9.3) | 1 345 (28.6) | 1 043 (22.2) |
| pB+E | 2 520 | 1 908 (75.7) | 612 (24.3) | 1 980 (78.6) | 204 (8.1) | 626 (24.8) | 647 (25.7) |
| nB+pE | 2 483 | 2 100 (84.6) | 383 (15.4) | 2 222 (89.5) | 241 (9.7) | 668 (26.9) | 709 (28.6) |
| pB+nE | 2 335 | 1 836 (78.6) | 499 (21.4) | 2 209 (94.6) | 349 (14.9) | 715 (30.6) | 869 (37.2) |

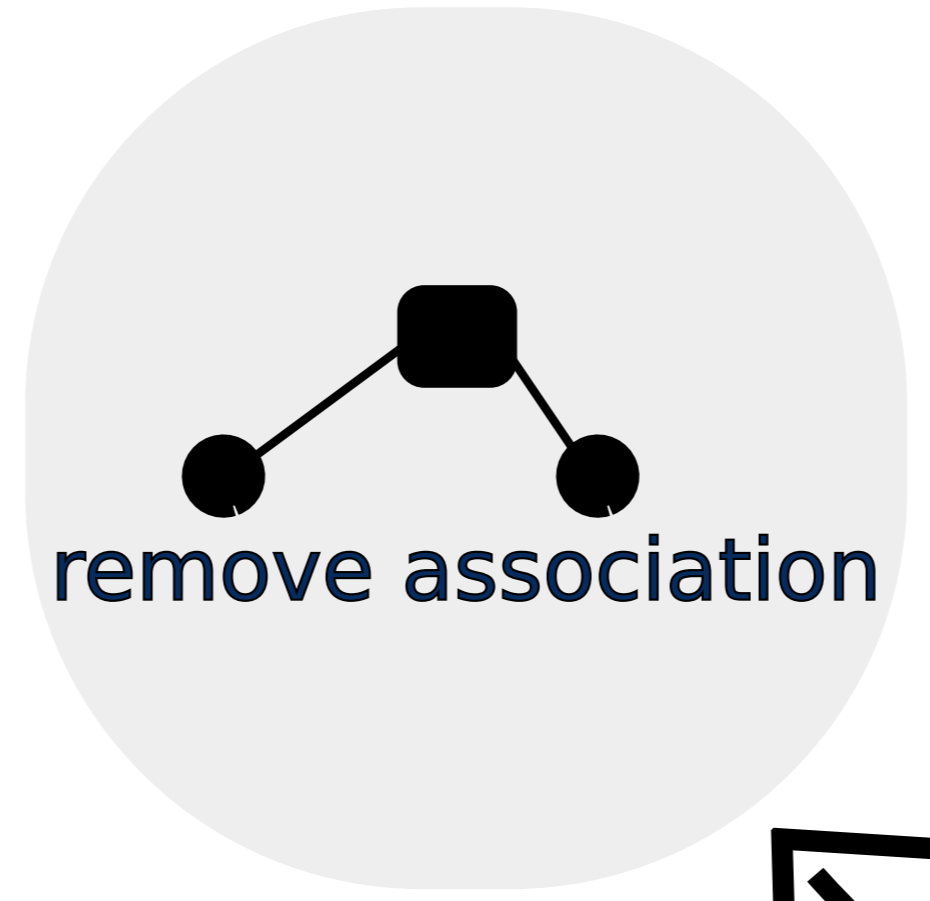*B - Bacteria; E - Eukaryotes; n - nano fraction; p - pico fraction*
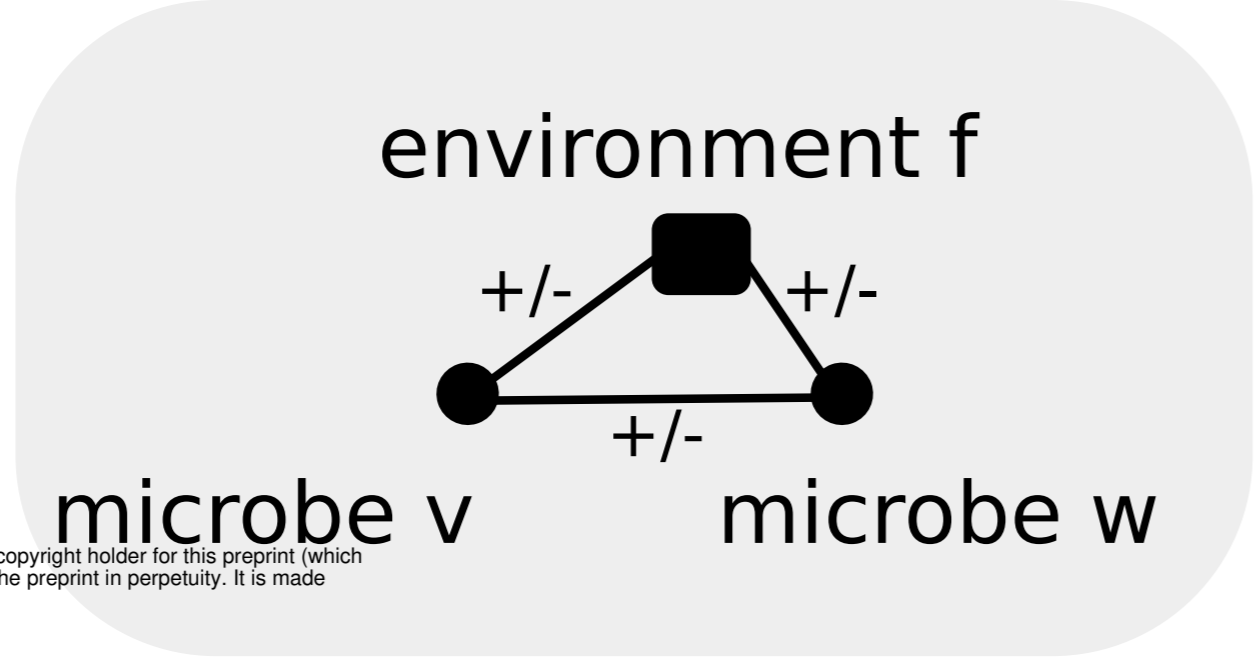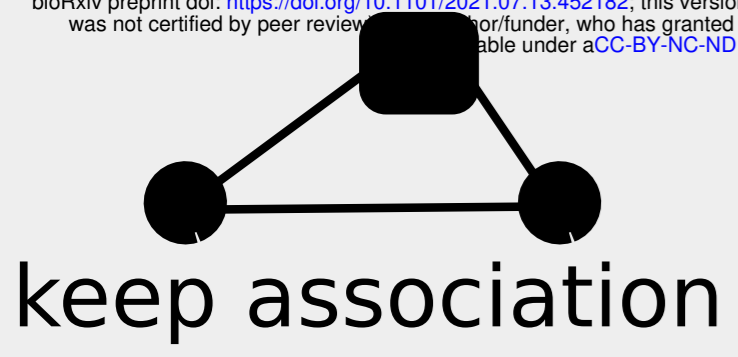
1254

**A) Classification and quantification of edges in the BBMO network**

- positive edges
- negative edges
- e-d - (positive or negative) environmentally driven edges
- E - Eukaryota
- B - Bacteria
- n - nano
- p - pico

**B) Location of specific edges in the BBMO network**

29,820 edges: positive 82.0% / not e-d edges 91.7% / number of triplets

negative 18.0% / e-d edges 8.3%

number of triplets: 0 1 2 3 4 5 6

754 nodes: nB 37.0% / pB 22.4% / nE 20.7% / pE 19.9%

keep association

environment f

+/-  +/-

+/-

microbe v     microbe w

remove association

**Entropy**

$$S(v) = -\sum_{i=1}^{n} p(v_i) \log\left(p(v_i)\right)$$

**Mutual Information**

MI(v;w) = S(v) + S(w) - S(v,w)

**Conditional Mutual Information**

CMI(v;w|f) = S(v,f) + S(w,f) - S(v,w,f) - S(f)

**Interaction Information**

II(v,w,f) = CMI(v;w|f) - MI(v;w)

Sign Pattern

+++, +--, -+-, --+

vs

---, -++, +-+, ++-

Overlap in time

vs

Interaction Information

II(v,w,f) < 0

vs

II(v,w,f) > 0

Data Processing Inequality

MI(v;w) < MI(v;f) and
MI(v;w) < MI(w;f)

vs

MI(v;w) > MI(v;f) or
MI(v;w) > MI(w;f)