# 1     Using null models to infer microbial co-occurrence networks

2     **Nora Connor, Albert Barberán & Aaron Clauset**

3     Affiliations: 1. Department of Computer Science, University of Colorado, Boulder CO, USA. 2. Cooperative

4     Institute for Research in Environmental Sciences, University of Colorado, Boulder, CO, USA. 3. BioFrontiers

5     Institute, University of Colorado, Boulder, CO, USA. 4. Santa Fe Institute, Santa Fe NM, USA.

6     ***Abstract***

7

8     Although microbial communities are ubiquitous in nature, relatively little is known

9     about the structural and functional roles of their constituent organisms' underlying

10     interactions. A common approach to study such questions begins with extracting a

11     network of statistically significant pairwise co-occurrences from a matrix of observed

12     operational taxonomic unit (OTU) abundances across sites. The structure of this

13     network is assumed to encode information about ecological interactions and processes,

14     resistance to perturbation, and the identity of keystone species. However, common

15     methods for identifying these pairwise interactions can contaminate the network with

16     spurious patterns that obscure true ecological signals. Here, we describe this problem

17     in detail and develop a solution that incorporates null models to distinguish ecological

18     signals from statistical noise. We apply these methods to the initial OTU abundance

19     matrix and to the extracted network. We demonstrate this approach by applying it to a

20     large soil microbiome data set and show that many previously reported patterns for

21     these data are statistical artifacts. In contrast, we find the frequency of three-way

22     interactions among microbial OTUs to be highly statistically significant. These results

23    demonstrate the importance of using appropriate null models when studying

24    observational microbiome data, and suggest that extracting and characterizing three-

25    way interactions among OTUs is a promising direction for unraveling the structure and

26    function of microbial ecosystems.

27

28    ***Author Summary***

29    Microbes are ubiquitous in the environment. We know that microbial communities –

30    the groups of microbes that live together, interact, and depend on one another – vary

31    across environments.  Multiple processes, ranging from competition between microbes

32    to environmental stress, are believed to alter microbial community composition. Here,

33    we describe a set of statistical techniques that can more accurately identify the

34    underlying taxa relationships that structure the observed abundances of microbes

35    across habitats. Using a large data set of soil samples collected across North and South

36    America, we both illustrate the statistical artifacts that incorrect methods can introduce

37    and describe proper techniques based on appropriate null models for studying how the

38    abundances of taxa vary across soil samples. These tools improve our ability to

39    distinguish ecologically meaningful interactions from simple statistical noise in such

40    observational data. Our application of these tools suggests some previous claims about

41    the network structure of microbial communities may be statistical artifacts.

42    Furthermore, we find that three-way interactions among microbial taxa are

43    significantly more common than we would expect at random, and thus may provide a

44    novel means for identifying ecologically meaningful interactions.

### *Introduction*

Microbes play essential roles in many, if not most, ecosystems. They play particularly important roles in regulating agricultural systems (e.g. Navarrete *et al.* 2015), human health (for a review, see Cho & Blaser 2012), and may even have an effect on mental health and behavior (Yano *et al.* 2015). Yet despite the importance of microbes and the recent technological advances in the field, essential questions remain about the composition and ecological structure of these microbial communities. For instance, how do communities change in response to internal dynamics and external perturbations, and how could we design communities with novel functionality? Deeper insights into the variables that shape the structure and function of microbial communities would have wide-ranging significance, both practical and theoretical.

One difficulty in scientifically addressing questions about microbial communities comes from the inability to culture the vast majority of microbes in a laboratory environment (Rappe & Giovannoni 2003). Instead, microbial community composition must be inferred from sequence data obtained by environmental DNA sampling. This limitation restricts our ability to test for causal mechanisms that drive a microbial community's structure and composition. Instead, observational data is often drawn from multiple samples across time or habitats (Barberán *et al.* 2015, Faust & Raes 2012, Peura *et al.* 2015, Steele *et al.* 2011, Kara *et al*. 2013). Complicating these efforts is a lack of robust statistical methods for analyzing these observational data in a way that reliably controls for plausible sources of variability and the spurious co-occurrence network patterns

68    they can produce. Here, we present and test methods for extracting statistically

69    significant co-occurrence patterns among microbes and for interpreting the induced

70    network structure.

71

72    A common design for a microbial community observational study has the following

73    form. Using high-throughput sequencing technologies, genetic data is extracted from a

74    set of locations, such as soil, water, or host-associated habitats including fecal samples

75    or cheek swabs. The observed DNA sequences are then binned into operational

76    taxonomic units (OTUs), which are taxonomic categories for microbes and are based on

77    a DNA sequence similarity threshold (usually 97% for 16S rRNA gene). This step is

78    necessary due to the difficulty in objectively defining microbial species, since these taxa

79    reproduce asexually and many have the ability to transfer genes horizontally. The OTUs

80    are placed into an abundance matrix $A$, where each element $A_{i,j}$ gives the number of

81    sequences representing a particular OTU $i$ observed in a particular sample or location $j$.

82    This matrix is then used to identify pairwise interactions, under the assumption that

83    OTUs whose abundances correlate across samples are likely to be ecologically related,

84    either symbiotically or through similar environmental preferences. To obtain

85    correlation values, a similarity measure is computed for each pair of vectors of OTU

86    abundances across locations (Faust & Raes 2012), and statistically significant

87    similarities are interpreted as potential ecological interactions. The set of such pairwise

88    interactions among the sampled OTUs can be transformed into a network of microbial

89    interactions, where nodes are OTUs and significant pairwise correlations are

90    represented as edges in the network. This network's structure can then be used to

91    understand the community's organization and function.

92

93    Such microbial interaction networks have many uses, not the least of which is making

94    complex data visually interpretable. They also facilitate the investigation of underlying

95    ecological processes that shape microbial communities. Past work on microbial

96    networks has examined many of their structural properties, including an OTU's degree

97    (number of connections), an OTU's betweenness centrality (a geometric measure of its

98    network position), the network's frequency of three-way interactions (the clustering

99    coefficient), and the network's average path length (a measure of system compactness).

100   These properties have been measured for networks derived from a variety of habitats,

101   including soil (Barberán *et al.* 2012), marine (Steele *et al.* 2011), and freshwater

102   communities (Kara *et al.* 2013). For instance, nodes in a network that have high degree

103   or high centrality may be interpreted as keystone taxa (Steele *et al.* 2011, Berry &

104   Widder 2014, Williams *et al.* 2014). Recent work has shown that these keystone taxa

105   play important roles in structuring microbial communities in plant-microbe

106   interactions (Agler *et al.* 2016). A group of OTUs that tend to co-occur may correspond

107   to taxa that share an ecological niche due to habitat filtering, or that participate in a

108   symbiotic interaction (Faust & Raes 2012). Similarly, groups of OTUs that tend to

109   mutually exclude each other may represent competitive interactions within a given

110   niche. We may also compare the structure of these microbial communities with that of

111   other biological networks (Williams *et al.* 2014), e.g., in order to understand whether

112   principles from macroecology also hold for microbial communities.

113

114    Network structure can also shed light on how a microbial community may respond to

115    environmental perturbations.  A right-skewed degree distribution among OTUs may be

116    evidence for robustness to high levels of random removal of species, or sensitivity to

117    the targeted removal of the keystone taxa (Faust & Raes 2012, Peura *et al*. 2015). This

118    network property may be related, for instance, to predicting whether a person's gut

119    microbiome will recover after a course of antibiotics. Similarly, network structure can

120    facilitate the identification of community assembly processes, for instance, by

121    comparing the structural signatures of neutral processes where all taxa are

122    demographically equivalent, versus those produced by niche-structured processes like

123    niche partitioning and competitive exclusion (O'Dwyer *et al*. 2012, Levy & Borenstein

124    2013, Pholchan *et al*. 2013, Tucker *et al.* 2015).  Greater insight into assembly dynamics

125    may facilitate predictions of community response to natural or artificial perturbations

126    (Faust & Raes 2012).

127

128    The broad importance of microbial interaction networks makes it essential that they be

129    reliably and accurately extracted from OTU abundance matrices, and that patterns in

130    the resulting network structure be properly interpreted. However, within the standard

131    approach to extracting these networks from co-abundance matrices are underlying

132    statistical assumptions that can contaminate the network with spurious or misleading

133    patterns. Specifically, spurious patterns in microbial co-occurrence networks may arise

134    from matrix sparsity, the choice of correlation function, and the use of thresholds.

135    Separate problems may arise when abundance data is normalized, making it

136    compositional. Addressing the issues of compositional data is beyond the scope of this

137    paper; however, in our conclusions we offer a brief discussion of their relationship to

138    the methods described here. In the following sections we examine the consequences of

139    spurious patterns in the data and leverage the ensuing errors as a motivation for the

140    use of null models as the foundation for the statistical methods we introduce. Our

141    methods are statistically principled methods, being based on standard null models, and

142    allow us to more accurately distinguish ecological signals from statistical noise, both in

143    the abundance matrix itself and in the distribution of edges in the derived network.

144

145    We demonstrate these techniques using a previously studied soil microbiome data set

146    from North and South America (Barberán *et al*. 2012). We find that some measures of

147    network structure are barely distinguishable from random noise, while others are more

148    plausibly the result of ecological interactions. A notable example of the latter category is

149    the network's clustering coefficient, the density of three-way OTU interactions, which

150    remains statistically significant when compared to each of our null models. We close

151    with a brief discussion of the utility of null models in studying observational data and

152    the ecological significance of triangles and modularity in microbial co-occurrence

153    networks.

154

155    ***Results***

156

157    **Two classes of null models**

158    Null models are a standard statistical approach for reliably identifying data patterns

159    that cannot be attributed to simple sources of random variation. Data distributions that

160    differ from a null model are thus potentially derived from complex processes. In our

161    case, large deviations may be interpreted as potentially caused by ecological processes.

162    One example of a null model is the common test of statistical significance, wherein we

163    measure the likelihood of observing, under the null model, a particular statistical value

164    or one more extreme. This probability is quantified by a standard $p$-value which has a

165    uniform distribution when the true data generating process is the null model. Common

166    choices for null models focus on a set of independent draws from a simple parametric

167    distribution, e.g., flipping coins or rolling dice. Null models can be substantially more

168    complicated, and in this case, numerical methods are typically required to calculate the

169    null distribution of the test statistic. If a null model is chosen well, meaning that it

170    incorporates plausible sources of random variation in the data, and the computed $p$-

171    value still low (typically below the conventional but nevertheless arbitrary threshold of

172    0.05), then a deviation between the model and the data can indicate the presence of

173    scientifically meaningful processes.

174

175    Here, we describe and study two classes of null models for inferring ecological

176    interactions from a matrix of OTU abundances. The first class facilitates the extraction

177    of significant pairwise interactions from the matrix in order to obtain a network. The

178    second class facilitates the detection of significant patterns in the distribution of edges

179    within the derived network.

180

181    In the rest of this section, we will introduce the first class of null models, in which we

182    will incorporate existing variability in the observed data to identify pairwise

183    interactions among OTUs. First, we correct the behavior of the Spearman rank

184    correlation coefficient when the OTU matrix is sparse by breaking ties randomly.

185    Second, in order to choose a threshold for significant interactions, we use matrix

186    permutations to generate artificial matrices with the same naturally high variance as

187    the data but which lack the correlations that are generated by ecological processes.

188    Applying the tie-breaking step to these artificial matrices yields a null distribution of

189    correlation scores, which provides a simple means for selecting a threshold for

190    statistically significant interactions. If any pair of OTUs in the tie-breaking model has a

191    correlation score above this threshold, we call this interaction statistically significant

192    and include it in the interaction network; any correlation below the threshold is

193    discarded.

194

195    In the second class of null models, we ask whether particular statistical patterns in the

196    distribution of these interactions across the network are likely the result of random

197    connectivity, and thus unlikely to be caused by ecological processes. Our approach here

198    builds on standard random graph models from network science, which control for the

199    average degree or the distribution of these degrees in order to construct an appropriate

200    null distribution for other network properties. Characteristics that are independent of

201    size and connectivity indicate co-existence of taxa, which may plausibly be attributed to

202    ecological interactions or functions.

203

204     The fact that some properties can be explained by the size, degree, or connectivity of

205     the network does not make them ecologically unimportant. In fact, the ecological impact

206     of overall biodiversity as well as co-occurrence patterns (i.e., functional redundancy) is

207     well established (Van Der Heijden *et al.* 2008, Philippot *et al.* 2013). In practice, these

208     null models can be used to identify more complicated statistically interesting patterns,

209     such as heterogeneous interactions among groups of microbes, that may relate to other

210     ecological processes, either known or unknown.

211

212     **The abundance matrix of microbial soil communities**

213     To illustrate the importance of examining microbial abundance data with respect to the

214     two null model classes, we apply these methods to previously collected data on soil

215     microbes sampled from 151 sites in North and South America (Lauber *et al.* 2009).

216     From soil samples, Barberán *et al.* extracted 16S rRNA sequences and binned them into

217     OTUs at a 90% rRNA sequence similarity threshold. They assigned taxonomy to OTUs

218     using RDP Classifier (Wang *et al.* 2007) against the Greengenes database (DeSantis *et al.*

219     2006). To obtain the abundance matrix, they computed the number of sequences that

220     mapped to each OTU at every sample site. To control for sample contamination and

221     potential sequencing errors, they discarded OTUs with fewer than 5 sequences across

222     all locations, which reduced the number of OTUs from 4,087 to 1,577.

223

224     Like many environmental DNA surveys, the resulting soil microbiome abundance

225     matrix is very sparse. Abundance values of zero comprise fully 85% of the matrix. Most

226     sites contained 150-300 OTUs, but only 1% of matrix entries have more than 10

227    sequences for a given OTU at a given site. In other words, although there were on the

228    order of 1000 sequences from each location, most OTUs at a site were phylogenetically

229    distinct.

230

231    In order to calculate the correlation of abundance patterns between a pair of OTUs, we

232    must choose a similarity score function. The most common choices in past studies are

233    Pearson and Spearman correlations, which exhibit good statistical sensitivity and

234    specificity under standard conditions (Berry & Widder 2014). However, the Pearson

235    correlation assumes that variables are normally distributed and linearly correlated, and

236    it behaves poorly when relationships are nonlinear, as may be the case in complex

237    microbial systems. Spearman's rank correlation, which measures the degree to which

238    two variables monotonically co-vary, does not suffer from this problem and is the more

239    common choice in microbiome studies (Lozupone *et al.* 2012; see also Weiss *et al.* 2016

240    for a review of correlation methods).

241

242    **A correction for matrix sparsity in Spearman ranks**

243    In this setting, Spearman will overestimate correlations when nearly all abundances are

244    either zero or some integer close to zero. As an intermediate step, Spearman assigns a

245    rank value to each location, and locations with equal abundance receive the same rank.

246    Thus, both matrix sparsity and a heavy-tailed distribution of abundances will induce a

247    very large number of multi-way ties, which will then have identical ranks. The result is

248    an inflated pairwise correlation score under Spearman. (Standard implementations of

249    Spearman's in Matlab, R, and Python all rely on the user to correct for ties in the data.)

250

251      This behavior can be corrected through breaking ties at random by adding a small

252      amount of real-valued noise to each entry in the abundance matrix. After adding these

253      minor perturbations, the set of all pairwise Spearman rank correlation coefficients ($\rho$)

254      form a smooth distribution (Figure 1A), as desired, rather than a perverse disjoint

255      distribution when ties are not broken (Figure 1B).

256

257      Crucially, the noise added to each observed value must not disturb the partial ordering

258      obtained without the noise. In practice, this is easily accomplished by using Monte Carlo

259      to sample from the many total orderings that are consistent with the original partial

260      ordering. Under a particular choice of significance threshold, this procedure will

261      generate a set of equally plausible networks, which are free from the statistical artifacts
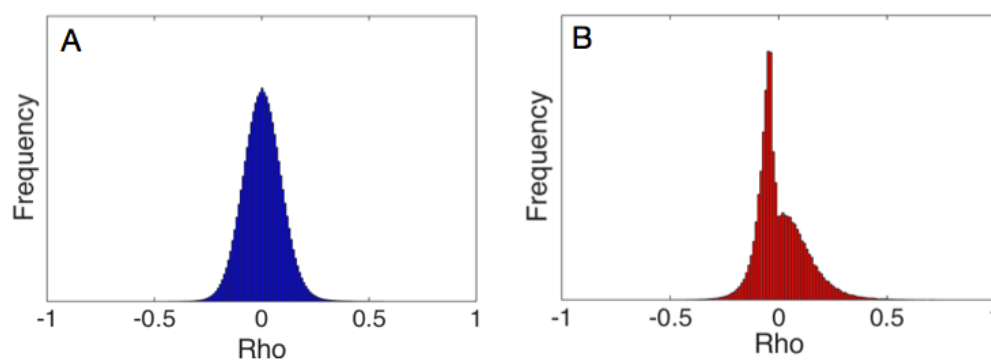
262      of tied ranks.

263

264

265



266      *Fig 1: Null distributions of Spearman rank correlation coefficients across sites for the Barberan*

267      ***et al. soil microbiome data.*** *(A) Coefficients under Monte Carlo sampling, using noise to break ties*

268      *randomly. (B) Coefficients without correcting for tied ranks between locations.*

269

270    This correction prevents the spurious conclusion that two taxa are ecologically related

271    because they are both absent from many of the same locations.  There are many reasons

272    why a taxon could have zero abundance at a given location, including habitat filtering,

273    local extinction due to ecological drift, dispersal limitation, or competitive exclusion. Or,

274    it may indicate that the taxon's DNA failed to bind to the 16S primer during

275    amplification, was undetected due to sequencing depth, or was absent by chance from

276    the soil sample.  In short, an abundance of zero is highly ambiguous, and a conservative

277    approach is to avoid inferring the presence of an interaction based primarily on shared

278    absences.

279

280    **Converting the sparsity-corrected data into a network**

281    To convert the abundance matrix into a network, we must apply a threshold to the

282    similarity scores. In this way, only OTU pairs for which the absolute value of their score

283    is above the threshold are connected in the network. It follows that a node with no

284    scores above the threshold will have a degree of zero in the network, and by convention

285    we omit such singletons from subsequent analysis (Barberán *et al.* 2012). As a result,

286    the number of nodes $n$ in the inferred network will typically be less than the number of

287    OTUs $N$ in the abundance matrix.

288

289    **Picking a threshold for significance**

290    Choosing an appropriate threshold of significance for similarity scores is an open

291    question, particularly for sparse data sets like OTU abundance matrices (Thomas &

292    Blitzstein, 2011). The goal of this choice is to eliminate pairwise interactions that are

293    likely due to statistical fluctuations or sampling noise, without excluding interactions

294    due to biological processes. Furthermore, we would like the scientific conclusions that

295    we draw from the resulting data to be robust to reasonable variations in threshold

296    choice (Thomas & Blitzstein, 2011). Currently, however, there is no generally reliable

297    method for balancing these two conflicting goals in OTU abundance matrices. Some

298    studies have used random permutations of the abundance matrix to compute a null

299    distribution of similarity scores, and then selected as a threshold the similarity value

300    corresponding to a conventional $p$-value choice of 0.01 or 0.05 (Faust & Raes 2012).

301    However, this procedure tends to select very low thresholds, and this may potentially

302    result in a high false positive rate for interactions. Other studies have used arbitrarily

303    chosen thresholds (Friedman & Alm 2012, Qin 2010).

304

305    Here, we use a repeated element-wise random permutation of the noise-added

306    abundance matrix to first compute a null distribution of similarity scores. We then

307    compute the size of the largest component -- the largest set of nodes for which any pair

308    is connected by some sequence of edges -- in the induced network for a wide range of

309    threshold values. Because the permutations break any ecologically-driven correlations

310    in the abundance matrix, this curve has a characteristic sigmoidal shape (Figure 2). The

311    location of the curve's transition to less than 1% of OTUs in the largest component

312    serves as a reasonable choice for the lower bound on the threshold. Networks derived

313    from this permuted data treatment are composed of all spurious links, so a threshold

314    below that transition, which would include these links, is overly inclusive. In practice, a

315 conservative choice of threshold will be a value slightly above this transition point.

316 Including the sparsity correction from above within this procedure serves to correct the

317 substantial distributional bias in similarity scores that would otherwise occur (see

318 Figure 1) as a result of multiple tied ranks and the heavy-tailed distribution of

319 abundance values.

320

321 We subject the OTU abundance data to three different treatments and systematically

322 vary the threshold to illustrate its impact on each. The three treatments are (i) the

323 original data, (ii) the original data with the Spearman correction, and (iii) the original

324 data with both Spearman correction and permutation null distribution. To illustrate the

325 effect of threshold choice on each treatment, we measure the fraction of OTUs $N$

326 contained in the largest component of the network across similarity thresholds (Figure

327 2). The size of this component provides a simple quantitative measure of overall graph

328 connectivity, and is a monotonically decreasing function of the threshold. That is, higher

329 thresholds will tend to produce smaller, less connected graphs, and lower thresholds

330 will tend to produce larger, more densely connected graphs.

331

332 **Section 2: Nonlinear effects of the threshold choice**

333 Figure 2 shows the percentage of nodes in the largest component as a function of the

334 choice of threshold, for each of the three treatments. To facilitate comparison with past

335 work on this data set (Barberán *et al.* 2012), we include a dashed vertical line at a

336 threshold of 0.36. This yields a network from the noise-added data of comparable size

337 to this past work ($n=300$). The location of the noise transition in the green line ($\Delta$), near

338    a threshold of 0.30 represents a lower bound on reasonable choices of a threshold.

339    Across thresholds, the original data shows a relatively slow decline in the size of this

340    largest component. Compared to the other treatments, which better eliminate spurious

341    connections, this slow decline is clearly an artifact of the presence of many false

342    positives in the network. By applying the Spearman correction or that correction and

343    the null distribution from permutations, the largest component shrinks much more

344    quickly. The difference between the treated lines and the original data illustrates the

345    dramatic extent to which not controlling for these statistical artifacts can alter the

346    extracted structure of the species interaction network.  A further observation is that the

347    smooth variation of the noise-added data treatment indicates that there is no obviously

348    best choice for a threshold, except somewhere close to but slightly above the noise

349    transition.

350

351    This finding illustrates the complexities that arise when using a threshold to extract a

352    network from a correlation matrix, and suggests that a particular choice requires some

353    justification or at least a robustness analysis to demonstrate that scientific conclusions

354    do not depend sensitively on that choice. From a data analysis perspective, we would

355    preserve the most ecological signal by not applying a threshold and instead using the

356    correlation scores as weights for edges in a fully connected or complete graph (Thomas

357    & Blitzstein 2011). However, many common network analysis techniques do not

358    generalize to weighted complete networks, or such methods have not yet been

359    developed. As a result, thresholding may be necessary to address certain classes of
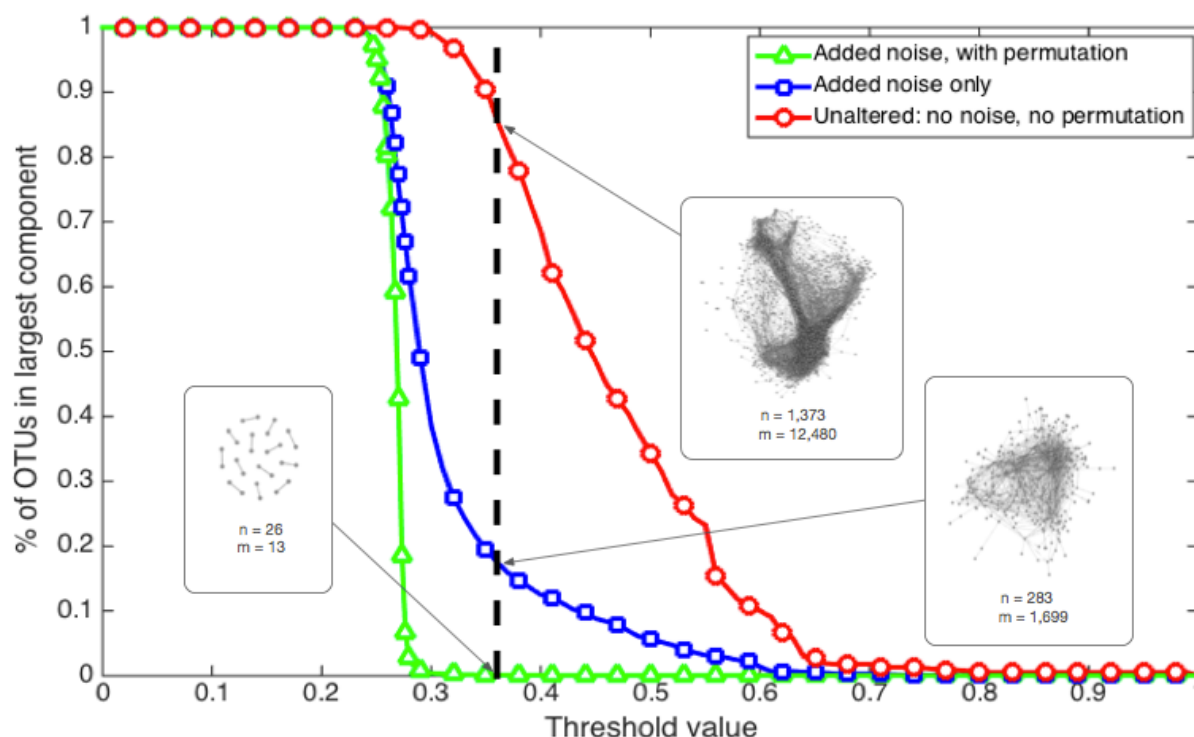
360    ecological questions.

**Fig 2: Fraction of all OTUs in the largest component, as a function of correlation threshold.**

*When the pairwise correlation threshold is 0, all edges are included and thus all nodes are in the largest component. When the threshold is 1, all edges are excluded and all singletons are discarded, so all of the OTUs are excluded from the analysis. The inset networks result from applying a threshold of 0.36, shown by the bold dashed line, to each of the treatments. The 0.36 threshold corresponds to 86% of OTUs in the largest component for the unaltered data, but just 18% of the OTUs in the noise-added treatment. For the permuted treatment, with noise added, the threshold intersects after the phase transition, yielding <1% of OTUs in the largest component.*

To further illustrate the impact of threshold choice on the structure of the induced network, we measured five standard network summary statistics as a function of threshold choice. These summary statistics are (i) the average degree, (ii) the average path length, (iii) the diameter, which is the maximal-length shortest path among any

376    pair of nodes, (iv) the modularity, which quantifies the extent to which nodes cluster

377    into groups, with more edges occurring inside groups than expected at random, and (v)

378    the clustering coefficient.

379

380    If the functional relationship between threshold and network statistic were constant or

381    linear, the particular choice of threshold is less likely to impact scientific conclusions

382    that depend on its particular value. For all five of these measures, however, we find a

383    nonlinear relationship between the measure and the choice of threshold. That is, the

384    structure of the network does not change smoothly, and different threshold choices can

385    lead to very different patterns of connectivity within the network (Figure 3).

386

387    For instance, even the average degree of this network exhibits a surprisingly nonlinear

388    pattern across thresholds (Figure 3A). The non-monotonicity, illustrated by the bump

389    around a threshold of 0.35, results from the convention of discarding nodes with no

390    connections. Thus, as the threshold increases, more of these nodes are created and then

391    excluded, which allows the average degree to increase again as the giant component

392    shrinks but the connectivity of its nodes stays relatively steady. (The average degree

393    touches the x-axis at a threshold of 0.75; when singletons are included, this transition

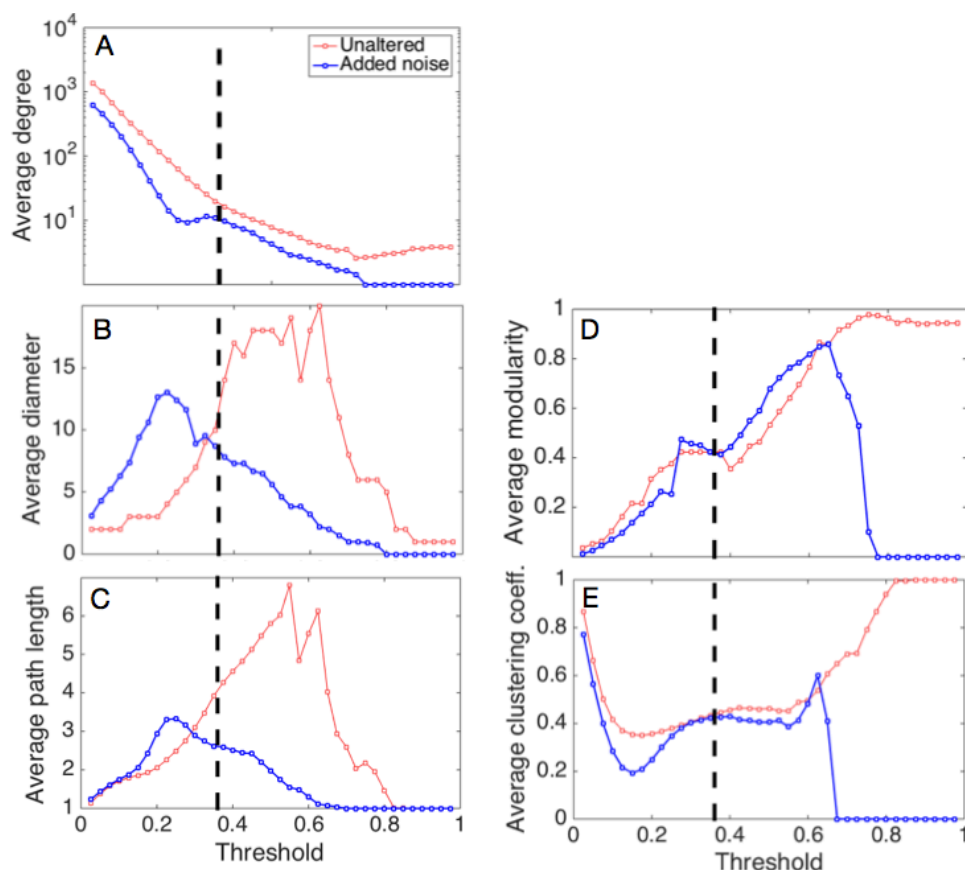394    occurs around a threshold of 0.40 (Supp. Fig.1).)

395

**Fig 3: Network properties vary as a function of threshold.** *This figure shows the change of network properties as the similarity score threshold varies between 0 and 1. The red lines represent the unaltered abundance data; the blue lines represent the noise-added data to correct rank ties. The vertical line at 0.36 is the same threshold used in Figure 2. Panels correspond to the following properties: (A) average degree, (B) diameter, (C) average path length, (D) maximum modularity, and (E) clustering coefficient.*

Similar patterns appear for the average and maximal path length (Figs. 3B and 3C). At lower thresholds, the network is relatively dense, making short paths among nodes plentiful. As the threshold increases, edges are removed, which makes the largest component sparser and increases path lengths. Finally, both measures decline above a

408    threshold of 0.25 as the size of the largest component itself begins to shrink, which

409    shortens path lengths again.

410

411    As the threshold increases, the largest component becomes sparser and the estimated

412    maximum modularity score also increases (Figure 3D), implying the existence groups of

413    nodes with relatively high internal connectivity (Clauset *et al.* 2004).  This property

414    deviates from a simple linear increase between threshold values of 0.25 and 0.38. At

415    higher values of the threshold, the largest component breaks up into small but fully

416    connected subgraphs, which have the highest possible marginal contributions to

417    modularity. However, for very high threshold values, the average degree falls below 1

418    and the network is composed primarily of disconnected edges, which yields a

419    modularity score of 0.

420

421    Because very low thresholds produce very dense networks, the clustering coefficient

422    (Figure 3E) is initially very high, but decreases quickly. Interestingly, and unlike the

423    other network statistics on this data set, the clustering coefficient stabilizes across

424    intermediate choices of thresholds, even as other network statistics are still changing.

425    As with the behavior of modularity, the clustering coefficient rises quickly and then falls

426    to 0 as the network crosses from being composed primarily of disconnected triangles

427    and edges to being composed entirely of disconnected edges.

428

429    The nonlinear dependence of the structure of the extracted network on the threshold

430    applied to the correlation matrix demonstrates the importance of performing

431     robustness analyses in this setting. Higher thresholds tend to naturally produce

432     networks with many small components, high modularity and shorter path lengths.

433     Lower thresholds tend to produce a large component, often with lower modularity

434     scores.  The threshold at which the transition between these two regimes occurs is

435     likely to be data dependent, and thus should be quantified in order to clarify the

436     confounding role that network size and density have on other network measures.

437

438     **Choosing a threshold for significant interactions**

439     If there existed a labeled data set, such as fully-defined microbial communities where

440     every individual microbial cell had fully sequenced 16S ribosomal RNA, we could train a

441     machine learning model to choose the threshold that best balances false positive

442     (spurious) links against false negative (missing) links, when those communities are

443     sampled. However, it is typically impractical to fully characterize the taxa that make up

444     an *in vivo* microbial community. Thus, in practice, choosing an intermediate value for

445     the threshold is a reasonable strategy. The threshold should be large enough to be

446     above the noise transition (Figure 2, green line), but small enough that the network is

447     not mostly disconnected. However, because of the nonlinear relationships between

448     network structure and threshold choice, a robustness analysis should always be

449     performed in order to determine whether a particular conclusion depends sensitively

450     on which intermediate threshold is chosen.

451

452     **Section 3: Measuring non-random network structure**

453

454    Given a choice of threshold and the corresponding network derived from corrected

455    Spearman correlation scores, we can now ask whether the distribution of the network's

456    links represents non-random patterns.  We use a second class of null models to find

457    statistically significant properties of the derived network by controlling for

458    connectivity. The two models in this class will allow us to distinguish whether a

459    particular pattern in the distribution of edges across the network is likely due to

460    chance.

461

462    The first null model is the Erdős–Rényi random graph, which preserves the average

463    degree of the derived network while removing any taxonomic information from the

464    nodes (Erdős & Rényi 1960, Kara *et al*. 2013). This model is sometimes denoted *G(n,p),*

465    where *n* is the number of nodes and *p = <k>/(n-1)*, where the mean degree *<k> = 2m/n*

466    is the probability that any pair of vertices is connected and where *m* is the number of

467    edges in the derived network. Drawing a large number random graphs from this model

468    (e.g., 2000 graphs) allows us to numerically estimate a null distribution for any network

469    property, while controlling only for the average degree of a node.

470

471    The second null model in this class is a Chung-Lu random graph model (Chung & Lu

472    2002) where we prohibit self-loops (an edge *(i,i)* for some node *i*).  Like the Erdős–

473    Rényi model, a Chung-Lu model starts with the same number of nodes as the derived

474    network.  Rather than giving each edge equal probability, this model preserves the

475    expected degree sequence by making the probability of an edge between two nodes

476    proportional to the product of their expected degrees. Specifically, the probability of an

477      edge between nodes $i$ and $j$ is $P_{i,j} = (k_i * k_j) / 2m - 1$, where $k_i$ is the degree of node $i$ in the

478      derived network. This model is similar to the popular configuration model (Molloy &

479      Reed 1995), but like the Erdős–Rényi model, it only produces simple networks, i.e.,

480      those without self-loops or multiple connections between the same pair of nodes. As

481      before, drawing a large number of random graphs from this model allows us to

482      numerically estimate a null distribution for the same network properties of interest, but

483      now controlling for the average degree of a node and the degree distribution across

484      nodes.

485

486      To illustrate how these models can be used to distinguish plausible structural patterns

487      from those generated by chance, we apply them to the soil microbe network extracted

488      in the previous section from the corrected Spearman scores. The derived network has

489      about $n = 268$ nodes and $m = 1730$ edges; the precise numbers vary depending on the

490      noise addition step. We then compare the null vs. the derived network's distributions

491      for (i) mean path length, (ii) modularity, (iii) diameter, and (iv) clustering coefficient.

492      Both null models are parameterized to match the mean degree and thus the random

493      graphs match the derived network on that measure by design.

494

495      Both path length and diameter are slightly elevated in the networks derived from the

496      corrected Spearman data compared to the null models (Figures 4A-B). The average path

497      length is $2.935 \pm 0.052$ for the corrected data, compared with $2.611 \pm 0.019$ for Erdős–

498      Rényi and $2.604 \pm 0.040$ for Chung-Lu. Similarly, the average diameter for the corrected

499      data is $7.411 \pm 0.874$, compared with $4.114 \pm 0.318$ in the Erdős–Rényi and $5.582 \pm$

500    0.544 for the Chung-Lu models. These differences are statistically significant, although

501    the effect size is small. That is, the extracted microbial interaction networks are only

502    less compact than we would expect if edges were distributed at random.

503

504    Similarly, the modularity scores (Figure 4C) are higher in the derived network

505    compared to those of the null models. The modularity is $0.415 \pm 0.014$ for the derived

506    network, while it is $0.217 \pm 0.005$ for the Erdős–Rényi model, and $0.280 \pm 0.012$ for the

507    Chung-Lu model. For these null models, the observed modularity scores are highly

508    statistically significant, and thus may represent a true ecological signal. However, as we

509    observed in the previous section, the modularity score is highly dependent on the

510    choice of threshold. For instance, under a threshold of 0.25 instead of 0.36, the

511    difference in modularity scores between the Chung-Lu null model and the derived

512    network vanishes (both are approximately 0.299). As such, the significance of the

513    modularity score should be interpreted cautiously.
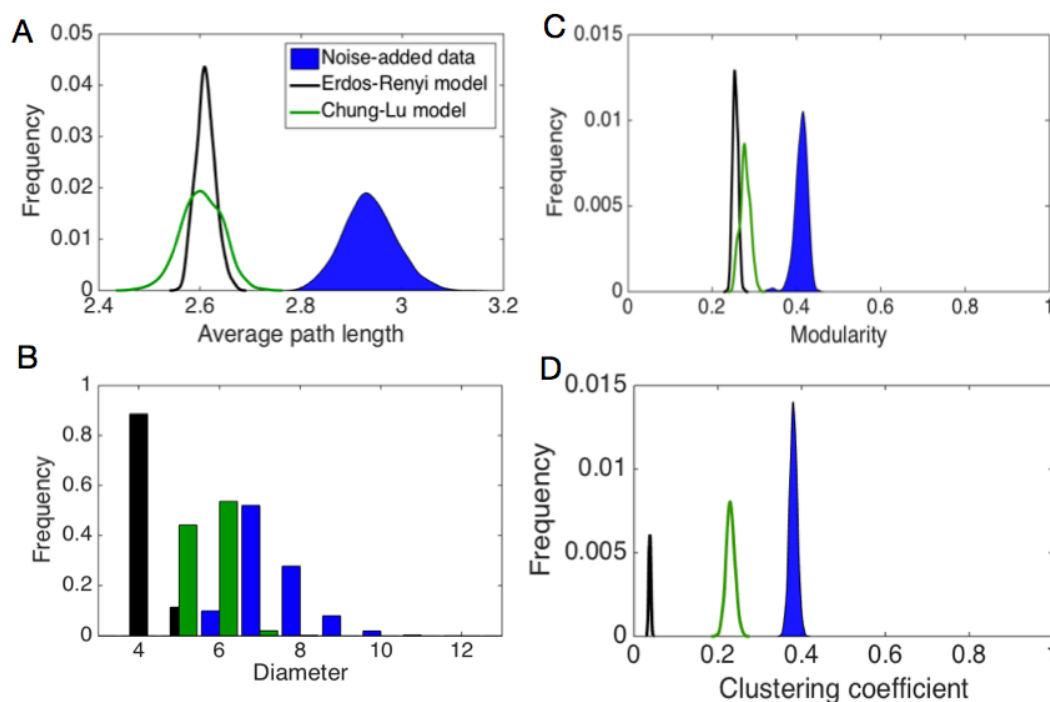
514

515

516

517

518

519

520

521

***Fig 4: Network properties compared with null network models with fixed connectivity.***

*Distributions of network properties across observed data and null models from the second class of*

*models: Erdős–Rényi and Chung-Lu. The observed data is graphed as blue in each plot. Panels show the*

*following properties: (A) average path length, (B) diameter, (C) modularity, and (D) clustering*

*coefficient.*

Compared to both null models, the derived network has a substantially higher

clustering coefficient (Figure 4D), which is similar to the scores observed in social

networks (Newman 2012; page 237).  The clustering coefficient for the derived

network is $0.380 \pm 0.009$, while it is $0.038 \pm 0.002$ for Erdős–Rényi random graphs and

$0.230 \pm 0.010$ for Chung-Lu random graphs. The difference in null distributions

indicates that about half of the value of the observed clustering coefficient can be

explained as an artifact of heterogeneous degree structure, which the Chung-Lu model

captures but the Erdős–Rényi model does not. This suggests that microbial

536    communities are enriched in three-way interactions (triangles) and these represent

537    potentially ecologically meaningful functional relationships among triplets of OTUs.
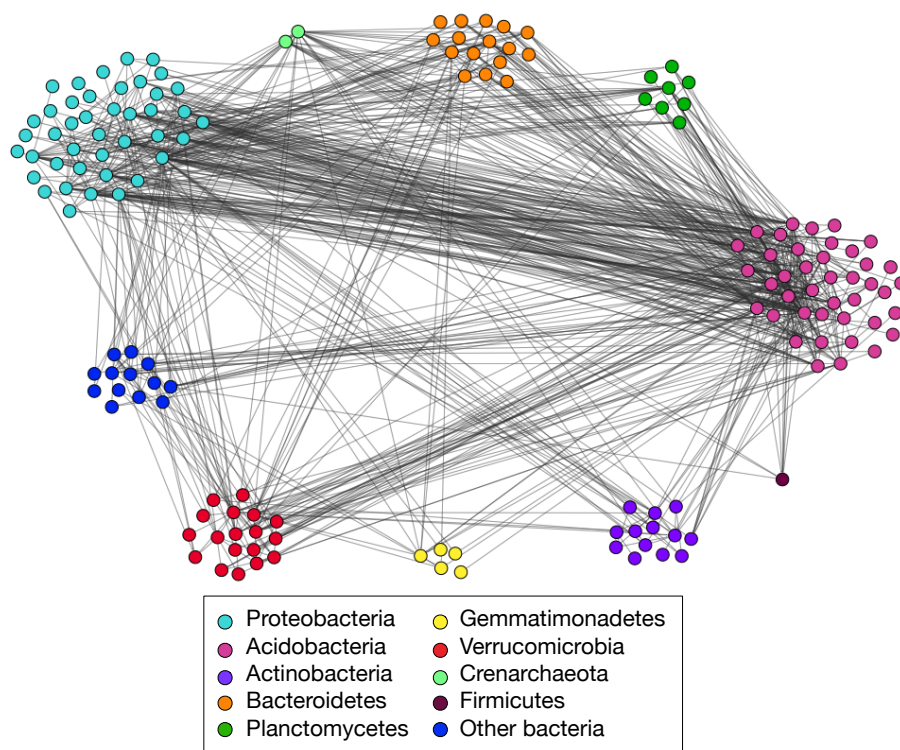
538



*Fig 5: Consensus network of edges, organized by phylum.* Edges in this figure are present in 90% of Monte Carlo simulations of noise addition.

539

540

541

542

543

544    **The consensus network**

545    Because the network properties of the derived network appear statistically significant

546    relative to our random graph null models, we can now construct and interpret a

547    "consensus network," which contains every pairwise interaction that is present in at

548    least 90% of the Monte Carlo samples. This consensus network is composed of 158

549    nodes and 787 edges. A simple but scientifically interesting question we may address

550    with this network is whether microbes tend to co-occur with others in the same

551    phylum. A positive signal of this assortative mixing pattern (Newman 2003) would

552    suggest a phylogenetic structuring for niche preferences or potential synergistic

553    relationships within phyla (Barberán *et al*. 2012).

554

555    However, we see little evidence for this hypothesis, finding instead that soil microbes

556    are not more likely to co-occur with taxa within phyla rather than across phyla (Figure

557    5).  Specifically, the number of edges between two phyla appears roughly proportional

558    to the number of taxa in both phyla, exactly as we would expect if such co-occurrences

559    were largely due to chance. As an additional check, we calculate the fraction of edges

560    that connect each phylum (Table 1). This enables us to investigate the potential

561    heterogeneous mixing of phyla. We observe that Acidobacteria and Proteobacteria have

562    the highest proportions of within-phylum edges, so these phyla are most likely to co-

563    occur with species within their respective phyla, when we enforce that clusters must

564    correspond to phyla. But the modularity of this network, which provides a quantitative

565    measure of assortativity among categorical labels on nodes (in this case, phyla), we find

566    a score of 0.0745 – much lower than the estimated maximal modularity when nodes are

567    allowed to mix independently of their phyla label (Fig. 4C). That is, non-phylogenetic

568    factors dominate the structure of OTUs interactions in this data set.

569

570

571

572  **Table 1: Fraction of edges connecting clusters based on phylum identity.**

| | Acido | Actino | Other | Bacteroid | Crenarch | Firmicutes | Gemma | Plancto | Proteo | Verruco | All |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Acido | **0.172** | 0.010 | 0.014 | 0.029 | 0.010 | 0.004 | 0.004 | 0.019 | 0.098 | 0.027 | 0.388 |
| Actino | 0.010 | **0.017** | 0.000 | 0.000 | 0.002 | 0.000 | 0.000 | 0.000 | 0.010 | 0.003 | 0.041 |
| Other | 0.014 | 0.000 | **0.004** | 0.002 | 0.002 | 0.001 | 0.001 | 0.001 | 0.011 | 0.008 | 0.043 |
| Bacteroid | 0.029 | 0.000 | 0.002 | **0.014** | 0.005 | 0.000 | 0.002 | 0.000 | 0.021 | 0.006 | 0.078 |
| Crenarch | 0.010 | 0.002 | 0.002 | 0.005 | **0.000** | 0.000 | 0.001 | 0.001 | 0.005 | 0.001 | 0.027 |
| Firmicutes | 0.004 | 0.000 | 0.001 | 0.000 | 0.000 | **0.000** | 0.000 | 0.000 | 0.001 | 0.000 | 0.006 |
| Gemma | 0.004 | 0.000 | 0.001 | 0.002 | 0.001 | 0.000 | **0.001** | 0.000 | 0.002 | 0.001 | 0.013 |
| Plancto | 0.019 | 0.000 | 0.001 | 0.000 | 0.001 | 0.000 | 0.000 | **0.000** | 0.014 | 0.001 | 0.036 |
| Proteo | 0.098 | 0.010 | 0.011 | 0.021 | 0.005 | 0.001 | 0.002 | 0.014 | **0.110** | 0.010 | 0.281 |
| Verruco | 0.027 | 0.003 | 0.008 | 0.006 | 0.001 | 0.000 | 0.001 | 0.001 | 0.010 | **0.031** | 0.087 |
| All | 0.388 | 0.041 | 0.043 | 0.078 | 0.027 | 0.006 | 0.013 | 0.036 | 0.281 | 0.087 | |

573

574  *Table 1: Fraction of edges connecting pairs of OTUs across nine phyla (or "other", for OTUs that don't map to*

575  *a known phylum). Edges connecting OTUs from the same phylum are highlighted. Note that this table is*

576  *symmetric because edges are undirected.*

577

578

579  ***Discussion***

580

581  A key step in better understanding the complex structure and function of microbial

582  ecosystems is identifying the ecologically meaningful interactions among microbes.

583  Distinguishing spurious interactions from real interactions is a key step in this process.

584  However, common approaches in this setting can contaminate the extracted network

585  with statistical artifacts that may confound ecological interpretation. Here, we have

586  developed and demonstrated simple and appropriate null models for addressing this

587  question at both the network extraction and the analysis steps, and we used them to

588  reanalyze a previously studied large soil microbiome data set.

589

590    After adding noise to the sparse OTU abundance data, we examined in detail the

591    difficulty of choosing a similarity threshold. Since network analysis depends on this

592    initial network derivation step, a conservative approach would test whether

593    conclusions about the network hold (or the same pattern appears) across a range of

594    reasonable threshold choices. In practice, we suggest choosing a threshold slightly

595    above the noise transition produced by the permutation test, and well below the point

596    where the network breaks up into small, disconnected components. An interesting line

597    of future work would examine the efficacy of supervised learning techniques from

598    machine learning to automatically choose a threshold that optimizes some downstream

599    performance measure (De Choudhury *et al.* 2010), e.g., likelihood of the extracted

600    network under a probabilistic generative model like the stochastic block model (Karrer

601    & Newman 2011).

602

603    Next, we used null models that preserve network connectivity to investigate the

604    variation in network measures. We did find slight but statistically significant elevation

605    of average path lengths and diameters in the derived network. One interpretation is

606    that microbial communities in soil are robust to environmental perturbations and have

607    evolved to recover or maintain structural stability amidst disturbances. Combining

608    future research on different microbial communities, such as the human gut microbiome,

609    with this type of network analysis would help clarify the role of average path length and

610    diameter (if any) in community robustness, e.g. after the administration of antibiotics.

611

612    We discovered that the clustering coefficient was higher in the derived network

613    compared to the network null models and that the score remained consistent across a

614    range of intermediate threshold values.  The elevated clustering coefficient may imply

615    that habitat filtering is playing an important role in the distribution and abundance of

616    OTUs in the soil. However, more research is needed to incorporate metabolic data or

617    other functional predictors into the model. Levy and Borenstein (2012) have shown

618    that in the human microbiome, co-occurrence is more often found in metabolically

619    competitive species than in metabolically complementary species -- evidence that

620    community assembly is best explained by habitat filtering in the human gut.  Similarly,

621    Goberna *et al.* (2014) also found that phylogenetic clustering was stronger in habitats

622    where competitive traits prevailed (i.e., in areas with high resource availability). Future

623    analysis of soil microbes should focus on metabolic competition and complementarity,

624    especially within OTU triads, to determine whether the elevated triad occurrence

625    corresponds to a specific community assembly mechanism (e.g., Pholchan *et al.* 2013,

626    Coyte *et al.* 2015). Future inquiry should focus on whether elevated clustering

627    coefficients are also present in networks derived from freshwater, marine, and human

628    microbiome samples.

629

630    We also discovered elevated maximum modularity scores in the derived network

631    compared to the null models. Higher modularity has been interpreted as corresponding

632    to greater niche partitioning (Faust & Raes 2012, Montoya *et al.* 2015). Further analysis

633    of metabolic functions of OTUs should investigate whether the highest-scoring

634    modularity partitions indicate true functional niches, wherein OTUs are more likely to

635     co-occur with OTUs in their own group than with OTUs in outside groups. For example,

636     gene expression data can be compared within and across the proposed functional

637     niches to identify shared or related metabolic functions (Levy & Borenstein 2013).

638     Future work may glean more from co-occurrence networks that focus on the level of

639     genes, rather than OTUs, which will become increasingly informative as more microbial

640     genomes are fully sequenced.

641

642     The consensus network was composed of 50% generalist OTUs and 50% OTUs that

643     were neither generalists nor specialists. The 79 generalist OTUs were identified based

644     on appearing in more than 80 locations. The other half of the OTUs were neither

645     generalists, nor specialists which appear in fewer than 10 sites with more than 18

646     sequences on average (Barberán *et al.* 2012). While no specialists appeared in the

647     consensus network, only 17 OTUs out of the 1577 total OTUS were identified as

648     specialists; given that about 10% of OTUs appeared in the consensus network, the

649     expected number of specialists in the consensus network would be 1.7. It is not possible

650     from this study to distinguish whether consensus networks are inherently biased

651     against specialists or whether there was simply not enough data in this sample to

652     distinguish specialists from noise. We do observe, however, that 76% of generalists (79

653     out of 104) are included in the consensus network.

654

655     The consensus network's strong modularity score may be due to the relative

656     concentration of generalists (Barberán *et al.* 2014). This might also explain why the

657     optimal partitioning of the consensus network did not correspond with phylogeny,

658      which was unexpected. The consensus network partitioning contrasts with basic

659      assumptions that ecological functions and niches are phylogenetically conserved

660      (Philippot *et al*. 2010). However, other recent work (Langille et al. 2013, Martiny et al.

661      2013) shows that while complex traits and housekeeping genes are generally deeply

662      conserved, other functional traits like assimilation of carbon sources are broadly

663      dispersed with respect to phylogeny. More work is required to identify the degree and

664      manner in which functional diversity structures real co-occurrence in the soil

665      microbiome.

666

667      The consensus network incorporates taxonomic information into the microbial

668      interaction networks, allowing us to use 16S rRNA sequence similarity to evaluate the

669      network structure. How best to incorporate that phylogenetic information is another

670      area of active research (Agler *et al.* 2016). Previous research has shown that using

671      lower binning thresholds for OTU identification does not reveal more about microbial

672      interactions, suggesting that even relatively broad binning strategies can be useful for

673      gaining ecological insight (Knights *et al*. 2011, Faust & Raes 2012).  However, other

674      authors recommend using the highest possible similarity threshold (Berry & Widder

675      2014). Future research should continue to address the phylogenetic information we

676      have about OTUs and how that data can be incorporated into identifying real ecoligal

677      interactions (O'Dwyer *et al.* 2012).

678

679      Many recent microbial association studies have focused on problems with analyzing

680      compositional data. For instance, several studies point out that compositional effects

681    are a concern when there are big differences in component sizes (Yang *et al.* 2016) or

682    when there are relatively few components (Ban *et al.* 2015). These problems are more

683    prevalent in marine metagenomics samples or host-associated microbes, but less for

684    the soil microbiome. We argue that using a relatively simple and nonparametric

685    similarity measure such as Spearman correlation coefficient can prevent the imposition

686    of preexisting notions about how taxa are distributed and how they interact. Compared

687    to other techniques, Spearman correlation coefficients are also efficient to calculate, a

688    problem acknowledged in the mLDM algorithm by its authors (Yang *et al.* 2016). For

689    data not derived from the soil microbiome, the suggested approaches for compositional

690    data could be used in conjunction with the network derivation methods described here.

691

692    In general, analyses of OTU-location matrices have uncertain scientific value as long as

693    we lack large sets of empirically validated OTU-OTU interactions by which to evaluate

694    network extraction methods. One possible remedy for this would be to remove some

695    fraction of observed edges from the inferred network and use predictive modeling to

696    identify the most probable missing edges. This type of link prediction has been used in

697    other contexts where the observation of the network is incomplete or error-prone

698    (Goldberg & Roth 2003), or where the network is changing, as in evolving social

699    networks (Liben-Nowell & Kleinberg 2007). A generative link-prediction model allows

700    us to test the degree to which our assumptions about the underlying structure of the

701    system are correct (Clauset *et al.* 2008, Guimerá & Sales-Pardo 2009).

702

703     Future research should apply different models to recover community structure. The

704     bipartite stochastic block model (Larremore *et al*. 2014) offers a compelling alternative

705     to clustering OTUs based on their similarities across locations. That is, instead of

706     converting the abundance matrix into a similarity matrix and applying an arbitrary

707     threshold, this model operates directly on the OTU-location matrix, obtaining both a

708     clustering of OTUs, a clustering of locations, and a mixing matrix that describes how

709     OTU groups interact with location groups. By operating on the original OTU-location

710     data, this approach would reduce the number and strength of assumptions used in the

711     analysis of such data. This model would be useful for finding OTUs that co-occur and

712     thus may be ecologically interacting, though it is defined only for occurrence data rather

713     than abundance data. To include sequence abundances, the weighted stochastic block

714     model (Aicher *et al.* 2015) could be used to directly analyze the OTU abundance values,

715     without having to choose a threshold.  For the task of clustering OTUs, these community

716     detection methods are a promising set of tools.

717

718     While the approach we have outlined for testing different network derivation

719     thresholds and evaluating null model connectivity has been applied to microbial

720     abundance data, it can also be applied across other biological network analyses. Sparse

721     data sets appear in a wide variety of biological settings, from eukaryotic environmental

722     DNA surveys (e.g., Stoeck *et al.* 2010) and gene co-expression networks. The issues of

723     threshold choice and appropriate null model selection are relevant across all disciplines

724     which use network science. Utilizing the appropriate statistical approaches will allow

725    researchers to draw stronger conclusions about correlation data, while leveraging the

726    quantitative tools from network science accurately.

727

728    **Materials and Methods**
729

730    Figure 6 illustrates the data analysis procedure used in this research.
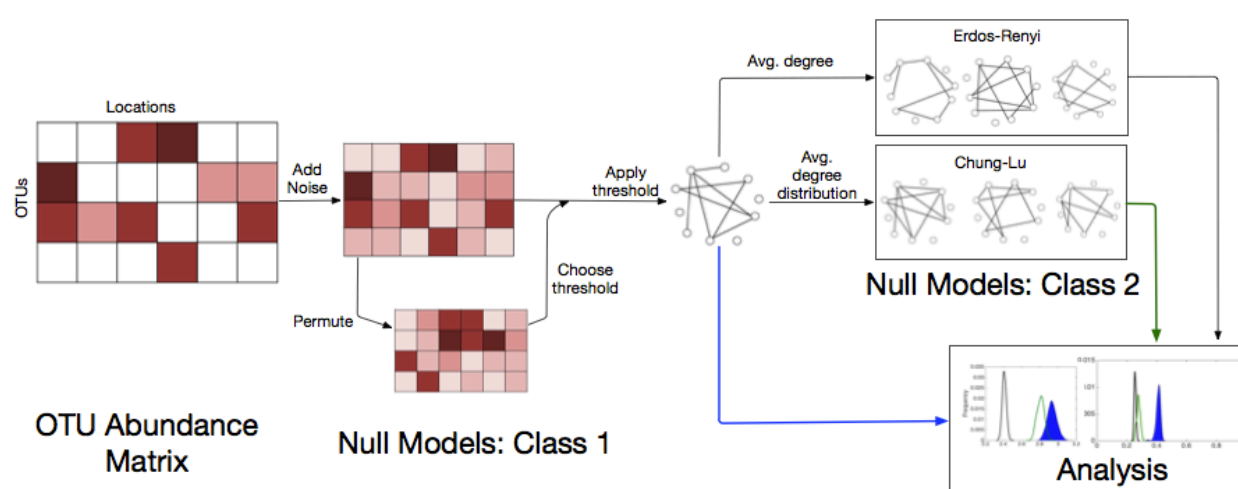
731



732

733

734    *Fig 6: Data analysis procedure for OTU abundance matrices. We start with the OTU abundance*

735    *matrix of N OTUs at L different locations. In the first class of null models, noise is added to every entry of*

736    *the matrix. Additionally, the noise-added matrix is permuted; the distribution of similarity scores in the*

737    *permuted matrix is used to set the lower bound for the threshold.  Next, the threshold is applied to*

738    *derive the observed network. This network is used to construct the second class of null models, Erdős–*

739    *Rényi, based on the average degree, and the Chung-Lu model, based on the average degree distribution.*

740    *Finally, the null network properties are compared to the observed network properties in the analysis*

741    *step.*

742

743

744 **Soil data**

745 The data set used in this experiment was acquired from previous work. Lauber *et al.*

746 (2009) acquired the bacteria and archaea data by pyrosequencing soil samples from

747 locations across North and South America. Their data set covers 151 sampling sites and

748 4088 unique OTUs, binned at 90% similarity (for explanation of the choice of 90%

749 similarity for binning, see Barberán *et al.* 2012). The data excludes OTUs with fewer

750 than 5 sequences across all sampling sites, decreasing the number of OTUs to 1577.

751 We use Spearman rank correlation coefficients to evaluate similarities between pairs of

752 OTUs based on their abundance patterns. For each OTU, Spearman converts a vector of

753 abundances into a vector of ranks, from largest to smallest. When there are identical

754 abundance values in several locations for a given OTU, the corresponding locations in

755 the rank vector are assigned the average rank for all tied entries. Given a pair of such

756 rank vectors x and y, the Spearman rank correlation coefficient is given by:

757
$$\rho = \frac{6 \, \Sigma \, r_i^2}{L(L^2 - 1)}$$

758 where $r_i = x_i - y_i$ is the difference between ranks between OTU $x$ and OTU $y$ in location $i$,

759 and where $L$ is the number of locations.

760

761 **Random noise addition**

762 Rather than allowing for ties among Spearman ranks, we correct for sparsity in the OTU

763 abundance matrix $A$ by adding noise to every OTU × location entry. We draw N × L

764 entries from a uniform distribution, $U([0,1])$, creating an N × L matrix *rand(N,L)*. To

765 ensure that we are breaking ties without reversing any true orderings, we adjust the

766   random values to be several orders of magnitude smaller than the minimum difference

767   between entries in $A$:

768
$$\Delta = \arg\min(\mathbf{A}_{i,j} - \mathbf{A}_{k,l})$$

769

770
$$b = \frac{\Delta}{1000}$$

771

772
$$\mathbf{E} = -b + [\, 2b \times \mathrm{rand}(N, L) \,]$$

773

774
$$\mathbf{A}' = \mathbf{A} + \mathbf{E}$$

775

776   To ensure that the configurations of equally likely location ranks were well sampled, we

777   repeated the noise addition steps 2000 times to generate a distribution of plausible

778   interaction networks.

779

780   **Random matrix permutations**

781   The most basic null model is the element-wise permutation of the OTU abundance

782   matrix. We chose a uniformly random permutation of the entries in the OTU abundance

783   matrix while maintaining the background distribution of abundances from which the

784   values were sampled.  The permuted data quickly transitions to having <1% of the

785   OTUs in the largest component; this is where we set the lower bound for the similarity

786   score threshold.

787

788

789 **Thresholding**

790 The threshold that we use, 0.36, produces a network of approximately 300 nodes from

791 the sparsity-corrected Spearman score data. This threshold was chosen to improve

792 comparability between our results and those of past studies on the same data

793 (Barberán *et al.* 2012). It is also similar to the threshold used by Friedman and Alm

794 (2012) of 0.30. We did not identify any additional quantitative guidelines for threshold

795 choice in other studies.

796

797 **Network derivation**

798 The network was derived by defining edges as connections between pairs of OTUs with

799 a $\rho$ value greater than the absolute value of the chosen threshold. Nodes with no edges

800 (also known as singletons) were omitted from the network, which is conventional in

801 defining the network. Average degree, average path length, and diameter were

802 calculated following the definitions in Newman (2010). Average degree is given by:

803

804
$$\langle k \rangle = \frac{2m}{n}$$

805

806 where $m$ is the number of edges and $n$ is the number of nodes. The average path length

807 is given by:

808
$$l = \frac{1}{n(n-1)} \sum_{i \neq j} d_{i,j}$$

809

810    where $d_{i,j}$ is the shortest path between nodes $i$ and $j$ (this is different from the $d_i$ values

811    used for the Spearman rank calculation). Diameter is the maximum value of $d_{i,j}$ across

812    all pairs of nodes in the network (i.e., the longest shortest path).

813

814    The clustering coefficient is defined as the global proportion of open triangles that are

815    closed by a third edge. We find all open triangles (i.e., paths of length 2) by taking the

816    dot product of the derived network's adjacency matrix $\boldsymbol{Q}$ with itself. Since this is an

817    undirected graph, we analyze the upper triangle of the matrix only, not including the

818    diagonal. Next, to find the proportion of length-two paths traversing three nodes that

819    are also closed triangles, we multiply the upper triangle by the original matrix $\boldsymbol{Q}$,

820    element-wise. The clustering coefficient $c$ is the fraction of open triangles that contain a

821    third edge to close the triad:

822

823    $$\mathbf{R} = \mathbf{Q} \cdot \mathbf{Q}$$

824

825    $$\mathbf{U} = \begin{cases} \mathbf{R}_{i,j}, & \text{if } i < j \\ 0, & \text{otherwise} \end{cases}$$

826

827    $$c = \frac{\Sigma\, \mathbf{U}}{\Sigma\, (\mathbf{U} \times \mathbf{Q})}$$

828

829

830    The maximum modularity was calculated using the using the popular greedy

831    agglomerative algorithm of Clauset, Newman and Moore (2004). This algorithm begins

832    with all nodes in their own group and then repeatedly merges the pair of groups that

833    maximizes the marginal improvement in the modularity score until only one group

834    remains. It then reports the maximum modularity value Q and the corresponding

835    grouping of nodes D that it traversed in this sequence. We used the implementation of

836    the algorithm in the igraph package in R (Csardi & Nepusz 2006).

837

838    **Class 2: Null network models**

839    Erdős–Rényi random graphs were created based on the average degree of the derived

840    network. Given an average degree of 11.64 and 300 nodes:

841

842
$$11.64 = \frac{2m}{n}$$

843

844
$$n = 300$$

845

846
$$m = \frac{11.64 \times 300}{2} = 1746$$

847

848

849    Once we had calculated the number of edges that we wanted in order to produce

850    similar average degrees to the real data, we used the $300 \times 300$ adjacency matrix $\mathbf{T}$ to

851    determine the correct threshold $p$:

852

853
$$|\mathbf{T}| = 300 \times 300 = 90{,}000$$

41

854

$$90{,}000 \times p = 1746 \times 2$$

856

857

$$p = \frac{1746 \times 2}{90{,}000} = 0.0388$$

858

859    To derive the Erdős–Rényi random graphs, we generated a $300 \times 300$ uniform random

860    matrix. We thresholded the upper triangle of the matrix at 1 - 0.0388 = 0.9612 and

861    reflected it across the diagonal. All entries on the diagonal were set to 0. This approach

862    is consistent with the mathematical definition of Erdős–Rényi random graphs, where

863    edges are randomly chosen for all pairs. Thus, the Erdős–Rényi graph has no self-loops

864    or multi-edges, as each pair is handled once.

865

866    For the second null model on the derived network, we used a modified Chung-Lu model

867    to produce edges between nodes while preserving the expected degree distribution.

868    The probability that an edge exists between OTU $i$ and OTU $j$ is given by,

869

870    $$p_{i,j} = \frac{k_i k_j}{2m}$$

871

872    where $k_i$ is the degree of node $i$ in the derived network. To generate a single Chung-Lu

873    network we set $n$ equal to the number of nodes in the observed network. Then, for each

874    pair of nodes $(i, j)$, we picked a uniform random number between 0 and 1. If the random

875    number was between 0 and $p_{i,j}$ -- which is proportional to the product of their degrees –

876    we created an edge connecting nodes $i$ and $j$ in the Chung-Lu network. We repeated this

877    method 2000 times to generate a distribution of Chung-Lu random graphs.

878

879    **Consensus network**

880    To derive the consensus network, we applied 2000 Monte Carlo runs to the corrected

881    Spearman data at the 0.36 threshold.  We included edges that appeared in 90% of the

882    trials to produce the consensus network.  We visualized networks using the software

883    Gephi (Bastian *et al.* 2009). Nodes were color-coded by phylum.

884    **Code**

885    All code for processing the data, applying null models, deriving networks, and

886    measuring network properties are publically available on GitHub. The repository is

887    saved under nkinboulder/MicrobeCommunities. The code for this project was written

888    in Matlab and it is extensively commented for clarity.

889

895

896

897

## References

899    Agler MT, Ruhe J, Kroll S, Morhenn C, Kim S, Weigel D, Kemen EM. Microbial hub taxa

900    link host and abiotic factors to plant microbiome variation. PLoS Biol. 2016;14(1):

901    e1002352.

902    Aicher C, Jacobs AZ, Clauset A. Learning latent block structure in weighted networks. J

903    Complex Netw. 2015;3(2): 221-48.

904    Ban Y, An L, Jiang H. Investigating microbial co-occurrence patterns based on

905    metagenomic compositional data. Bioinformatics 2015;31(20): 3322-9.

906    Barberán A, Bates S, Casamayor E, Fierer N. Using network analysis to explore co-

907    occurrence patterns in soil microbial communities. ISME J. 2012;6(2): 343-51.

908    Barberán A, Ramirez KS, Leff JW, Bradford MA, Wall DH, Fierer N. Why are some

909    microbes more ubiquitous than others? Predicting the habitat breadth of soil bacteria.

910    Ecol Lett 2014;17(7): 794-802.

911    Barberán A, Ladau J, Leff JW, Pollard KS, Menninger HL, Dunn RR, Fierer N. Continental-

912    scale distributions of dust-associated bacteria and fungi. Proc Natl Acad Sci U S A.

913    2015;112(18): 5756-61.

914    Bastian M, Heymann S, Jacomy M. Gephi: an open source software for exploring and

915    manipulating networks. ICWSM. 2009 May 17;8: 361-2.

916

917   Berry D, Widder S. Deciphering microbial interactions and detecting keystone species

918   with co-occurrence networks. Front Microbiol. 2014;5(219): 1-14.

919   Csárdi G, Nepusz T. The igraph software package for complex network research.

920   InterJournal, Complex Systems. 2006 Jan 11;1695(5): 1-9.

921   Chung F, Lu L. Connected components in random graphs with given expected degree

922   sequences. Ann Comb. 2002;6: 125-45.

923   Cho I, Blaser MJ. The human microbiome: at the interface of health and disease." Nat

924   Rev Genet. 2012;13: 260-70.

925   Clauset A, Moore C, Newman MEJ. Hierarchical structure and prediction of missing links

926   in networks. Nature. 2008;453: 98-101.

927   Clauset A, Newman MEJ, Shalizi C. Finding community structure in very large networks.

928   Phys Rev **E**. 2004;70: 066111.

929   Coyte KZ, Schluter J, Foster KR. The ecology of the microbiome: networks, competition,

930   and stability. Science. 2015;350: 663-6.

931   De Choudhury M, Mason WA, Hofman JM, Watts DJ. Inferring relevant social networks

932   from interpersonal communication. In**:** Proceedings of the 19th international

933   conference on World wide web; 2010 Apr 26-30; Raleigh NC USA. New York: ACM;

934   2010. p. 301-10.

935  DeSantis TZ, Hugenholtz P, Larsen N, Rojas M, Brodie EL, Keller K, Huber T, Dalevi D, Hu

936  P, Andersen GL. Greengenes, a chimera-checked 16S rRNA gene database and

937  workbench compatible with ARB. Appl Environ Microbiol 2006;72(7): 5069-72.

938  Erdős P, Rényi A. On the evolution of random graphs. Publ. Math. Inst. Hungar. Acad. Sci.

939  1960;5: 17-61.

940  Faust K, Raes J. Microbial interactions: from networks to models. Nat Rev Microbiol.

941  2012;10: 538-50.

942  Friedman J, Alm EJ. Inferring correlation networks from genomic survey data. PLoS

943  Comput Biol. 2012;8(9): e1002687.

944  Goberna M, Navarro-Cano JA, Valiente-Banuet A, Garcia C, Verdu M. Abiotic stress

945  tolerance and competition-related traits underlie phylogenetic clustering in soil

946  bacterial communities. Ecol Lett 2014;17(10): 1191-201.

947  Goldberg DS, Roth FP. Assessing experimentally derived interactions in a small world.

948  Proc Natl Acad Sci U S A.  2003;100(8): 4372-6.

949  Guimerá R, Sales-Pardo M. Missing and spurious interactions and the reconstruction of

950  complex networks. Proc Natl Acad Sci U S A. 2009;106(52): 22073-8.

951  Kara EL, Hanson PC, Hu YH, Winslow L, McMahon KD. A decade of seasonal dynamics

952  and co-occurrences within freshwater bacterioplankton communities from eutrophic

953  Lake Medota, WI, USA. ISME J. 2013;7: 680-4.

954     Karrer B, Newman MEJ. Stochastic blockmodels and community structure in networks.

955     Phys Rev E. 2011;83(1): 016107.


956     Knights D, Costello EK, Knight R. Supervised classification of human microbiota. FEMS

957     Microbiol Rev. 2011;35: 343-59.


958     Langille MGI, Zanefeld J, Caporaso JG, McDonald D, Knights D, Reyes JA, Clemente JC,

959     Burkepile DE, Vega Thurber RL, Knight R, Beiko RG, Huttenhower C. Predictive

960     functional profiling of microbial communities using 16S rRNA marker gene sequences.

961     Nat Biotechnol 2013;31(9): 814-21.


962     Larremore DB, Clauset A, Jacobs AZ. Efficiently inferring community structure in

963     bipartite networks. Phys Rev **E**. 2014;90: 012805. doi:10.1103/PhysRevE.90.012805


964     Lauber CL, Hamady M, Knight R, Fierer N. Pyrosequencing-based assessment of soil pH

965     as a predictor of soil bacterial community structure at the continental scale. Appl

966     Environ Microbiol. 2009;75(15): 5111-20.


967     Levy R, Borenstein E. Metabolic modeling of species interaction in the human

968     microbiome elucidates community-level assembly rules. Proc Natl Acad Sci U S A.:

969     2013;110(31): 12804-9.


970     Liben-Nowell D, Kleinberg J. The link-prediction problem for social networks. J Am Soc

971     Inf Sci Technol, 2007;58(7): 1019-31.

972   Lovell D, Pawlowsky-Glahn V, Egozcue JJ, Marguerat S, Bahler J. Proportionality: A valid

973   alternative to correlation for relative data. PLoS Comput Biol. 2015;11(3). doi:

974   10.1371/journal.pcbi.1004075

975   Lozupone C, Faust K, Raes J, Faith JJ, Frank DN, Zanefeld J, Gordon JI, Knight R.

976   Identifying genomic and metabolic features that can underlie early successional and

977   opportunistic lifestyle of human gut symbionts. Genome Res. 2012;22: 1974-84.

978   Martiny AC, Treseder K, Pusch G. Phylogenetic conservatism of functional traits in

979   microorganisms. ISME J. 2013;7: 830-8.

980   Molloy M, Reed B. A critical point for random graphs with a given degree sequence.

981   Random Struct Alg. 1995;6(2-3): 161–80.

982   Montoya D, Yallop ML, Memmott J. Functional group diversity increases with

983   modularity in complex food webs. Nat Commun. 2015;6: 7379.

984   Navarrete AA, Tsai SM, Mendes LW, Faust K, de Hollander M, Cassman NA, Raes J, van

985   Veen JA, Kuramae EE. Soil microbiome responses to the short-term effects of

986   Amazonian deforestation. Mol Ecol. 2015;24(10): 2433-48.

987   Newman MEJ. Mixing patterns in networks. Phys Rev E. 2003;67(2): 026126.

988   Newman MEJ. Networks: An Introduction. 1st ed. Oxford University Press; 2010.

989   O'Dwyer JP, Kembel SW, Green JL. Phylogenetic diversity theory sheds light on the

990   structure of microbial communities. PLoS Comput Biol. 2012;8(12): e1002832.

991    Peura S, Bertilsson S, Jones RI, Eiler A. Resistant microbial co-occurrence patterns

992    inferred by network topology. Appl Environ Microbiol. 2015;81(6): 2090-7.

993    Philippot L, Andersson SGE, Battin TJ, Prosser JI, Schimel JP, Whitman WB, Hallin S. The

994    ecological coherence of high bacterial taxonomic ranks. Nat Rev Microbiol 2010;8: 523-

995    9.

996    Philippot L, Spor A, Henault C, Bru D, Bizouard F, Jones CM, Sarr A, Maron PA. Loss in

997    microbial diversity affects nitrogen cycling in soil. ISME J. 2013;7: 1609-19.

998    Pholchan MK, de C. Baptista J, Davenport RJ, Sloan WT, Curtis TP. Microbial community

999    assembly, theory and rare functions. Front Microbiol. 2013;4: 68.

1000   Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut microbial

1001   gene catalogue established by metagenomic sequencing. Nature 2010;464: 59-65.

1002   Rappe MS, Giovannoni, SJ. The uncultured microbial majority. Annu Rev Microbiol.

1003   2003;57: 369-94.

1004   Stoeck T, Bass D, Nebel M, Christen R, Jones MDM, Breiner HW, Richards TA. Multiple

1005   marker parallel tag environmental DNA sequencing reveals a highly complex eukaryotic

1006   community in marine anoxic water. Mol Ecol. 2010;19: 21-31.

1007   Thomas AC, Blitzstein JK. Valued ties tell fewer lies: Why not to dichotomize network

1008   edges with thresholds; 2011. Preprint. Available: arXiv:1101.0788v2. Accessed 27 June

1009   2016.

1010    Tucker CM, Shoemaker LG, Davies KF, Nemergut DR, Melbourne BA. Differentiating

1011    between niche and neutral assembly in metacommunities using null models of β-

1012    diversity. Oikos. 2016;125: 778-89.


1013    Van Der Heijden MGA, Bardgett RD, Van Straalen NM. The unseen majority: soil

1014    microbes as drivers of plant diversity and productivity in terrestrial ecosystems. Ecol

1015    Lett 2008;11: 296-310.


1016    Wang Q, Garrity GM, Tiedje JM, Cole JR. Naïve Bayesian classifer for rapid assignment of

1017    rRNA sequences into the new bacterial taxonomy. Appl Environ Microbiol.

1018    2007;73(16): 5261-7.


1019    Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, Xia LC, Xu ZZ, Ursell

1020    L, Alm EJ, Birmingham A, Cram JA, Fuhrman JA, Raes J, Sun F, Zhou J, Knight R.

1021    Correlation detection strategies in microbial data sets vary widely in sensitivity and

1022    precision. ISME J. 2016;10: 1669-81.


1023    Williams RJ, Howe A, Hofmockel KS. Demonstrating microbial co-occurrence pattern

1024    analyses within and between ecosystems. Front Microbiol. 2014;5: 358.


1025    Yang Y, Chen N, Chen T. mLDM: a new hierarchical Bayesian statistical model for sparse

1026    microbial association discovery; 2016. Preprint. Available: bioRxiv:042630. Accessed 8

1027    August 2016.

1028    Yano JM, Yu K, Donaldson GP, Shastri GG, Ann P, Ma L, Nagler CR, Ismagilov RF,

1029    Mazmanian SK, Hsiao EY. Indigenous bacteria from the gut microbiota regulate host

1030    serotonin biosynthesis. Cell 2015;161(2): 264-76.