

# functionInk: An efficient method to detect functional groups in multidimensional networks reveals the hidden structure of ecological communities

January 16, 2020

## Short Title

Community detection in multidimensional ecological networks

Alberto Pascual-García<sup>1, 2,\*</sup> and Thomas Bell<sup>1</sup>

(1) Department of Life Sciences. Silwood Park Campus. Imperial College London, Ascot, United Kingdom

(2) Current address: Institute of Integrative Biology. ETH-Zürich, Zürich, Switzerland

(\* ) Correspondence: alberto.pascual.garcia@gmail.com.

## Abstract

1. Complex networks have been useful to link experimental data with mechanistic models, and have become widely used across many scientific disciplines. Recently, the increasing amount and complexity of data, particularly in biology, has prompted the development of multidimensional networks, where dimensions reflect the multiple qualitative properties of nodes, links, or both. As a consequence, traditional quantities computed in single dimensional networks should be adapted to incorporate this new information. A particularly important problem is the detection of communities, namely sets of nodes sharing certain properties, which reduces the complexity of the networks, hence facilitating its interpretation.

2. In this work, we propose an operative definition of “function” for the nodes in multidimensional networks, and we exploit this definition to show that it is possible to detect two types of communities: i) modules, which are communities more densely connected within their members than with nodes belonging to other communities, and ii) guilds, which are sets of nodes connected with the same neighbours, even if they are not connected themselves. We provide two quantities to optimally detect both types of communities, whose relative values reflect their importance in the network.

3. The flexibility of the method allowed us to analyze different ecological examples encompassing mutualistic, trophic and microbial networks. We showed that by considering both metrics we were able to obtain deeper ecological insights about how these different ecological communities were structured. The method mapped pools of species with properties that were known in advance, such as plants and pollinators. Other types of communities found, when contrasted with external data, turned out to be ecologically meaningful, allowing us to identify species with important functional roles or the influence of environmental variables. Furthermore, we found the method was sensitive to community-level topological properties like the nestedness.

4. In ecology there is often a need to identify groupings including trophic levels, guilds, functional groups, or ecotypes. The method is therefore important in providing an objective means of distinguishing modules and guilds. The method we developed, *functionInk* (functional linkage), is computationally efficient at handling large multidimensional networks since it does not require optimization procedures or tests of robustness. The method is available at: [HTTPS://GITHUB.COM/APASCUALGARCIA/FUNCTIONINK](https://github.com/APASCUALGARCIA/FUNCTIONINK).

**Keywords:** Multiplex networks, community detection, modules, guilds, mutualistic networks, trophic networks, microbial networks

# 1 Introduction

1 Networks have played an important role in the development of ideas in ecology, particularly in understanding food  
2 webs [1], and flows of energy and matter in ecosystems [2]. However, modern ecological datasets are becoming  
3 increasingly complex, notably within microbial ecology, where multiple types of information (taxonomy, behaviour,  
4 metabolic capacity, traits) on thousands of taxa can be gathered. A single network might therefore need to  
5 integrate different sources of information, leading to connections between nodes representing relationships of  
6 different types, and hence with different meanings. Advances in network theory have attempted to develop tools  
7 to analyse these more sophisticated networks, encompassing ideas such as multiplex, multilayer, multivariate  
8 networks, reviewed in [3]. There could therefore be much value in extending complex networks tools to ecology  
9 in order to embrace these new concepts.

10 In this paper we aim to address a particularly relevant problem in complex networks theory, namely the  
11 detection of “communities”[4] when the network contains different types of links. In addition, we are interested in  
12 finding a method with the flexibility to identify different types of communities. This is motivated by the fact that,  
13 in ecology, it is recognized that communities may have different topologies [5] and often an intrinsic multilayer  
14 structure [6]. We aim to detect two main types of communities. Firstly, the most widely adopted definition  
15 of community is the one considering sets of nodes more densely connected within the community than with  
16 respect to other communities, often called *modules* [7]. An example in which modules are expected is in networks  
17 representing significant co-occurrences or segregations between microbial species, when these relationships are  
18 driven by environmental conditions. Since large sets of species may simultaneously change their abundances in  
19 response to certain environmental variables [8, 9], this results in large groups of all-against-all co-occurring species,  
20 and between-groups segregations. Secondly, we are interested in finding nodes sharing a similar connectivity  
21 pattern even if they are not connected themselves. An example comes from networks connecting consumers and  
22 their resources, when communities are determined looking for consumers sharing similar resources preferences.  
23 This idea is aligned with the classic Eltonian definition of niche, which emphasizes the functions of a species rather  
24 than their habitat [10]. We call this second class of communities *guilds*, inspired by the ecological meaning in which  
25 species may share similar ways of exploiting resources (i.e. similar links) without necessarily sharing the same  
26 niche (not being connected themselves), emphasizing the functional role of the species [11]. Consequently, guilds  
27 may be quite different to modules, in which members of the same module are tightly connected by definition.  
28 The situation in which guilds are prevalent is known as disassortative mixing [12], and its detection has received  
29 comparatively less attention than the “assortative” situation (which results in modules) perhaps with the exception  
30 of bipartite networks [13, 14].

31 There are many different approximations for community detection in networks, summarized in [15]. However,  
32 despite numerous advances in recent years, it is difficult to find a method that can efficiently find both modules  
33 and guilds in multidimensional networks, and that is able to identify which is the more relevant type of community  
34 in the network of interest. This might be because there is no algorithm that can perform optimally for any network  
35 [16], and because each type of approximation may be suited for some networks or to address some problems but  
36 not for others, as we illustrate below.

37 Traditional strategies to detect modules explore trade-offs in quantities like the betweenness and the clustering  
38 coefficient [7], as in the celebrated Newman-Girvan algorithm [17]. Generalizing the determination of modules  
39 to multidimensional networks is challenging. Consider, for instance, that a node A is linked with a node B and  
40 this is, in turn, linked with a node C, and both links are of a certain type. If A is then linked with node C with  
41 a different type of link, should the triangle ABC be considered in the computation of the clustering coefficient?  
42 One solution proposed comes from the consideration of stochastic Laplacian dynamics running in the network  
43 [18], where the permanence of the informational fluxes in certain regions of the network reflects the existence  
44 of communities. This approximation has been extended to consider multilayer networks [19], even if there are  
45 modules defined in different layers that highly overlap, hence defining communities (combination of modules)  
46 across layers [20]. A fundamental caveat for these methods is that the links must have a clear interpretation for  
47 how their presence affects informational fluxes. Returning to the above example, if the links AB and BC represent  
48 mutualistic interactions and the link AC represents a competitive one, can a random walk follow the link AC  
49 when this interaction does not describe a flux of biomass between the species but rather a disruption in the flux  
50 of biomass of AB and BC?

51 A related approach searches for modules using an optimization function that looks for a partitioning in a  
52 multilayer network that maximizes the difference between the observed model and a null model which considers  
53 the absence of modules [13]. This strategy can be applied to multidimensional networks, but raises questions such

as which is the appropriate null model, and how to determine the coupling between the different layers defining the different modes of interaction [21]. In addition, since these approximations focus on the detection of modules, they neglect the existence of guilds or other network structures.

Regarding the search of guilds, this problem has received notable attention in social sciences following the notion of structural equivalence. Two nodes are said to be structurally equivalent if they have the same connectivity in the network [22]. The connectivity may be defined either analyzing if two nodes share the same neighbours, if two nodes are connected with neighbors of the same type even if they are not necessarily the same (following some preassigned roles for the nodes, e.g. prey are structurally equivalent because are connected to predators), or a combination of both. Social agents often have an assigned role, which is why structural equivalence is particularly important in social networks.

An approximation that has exploited the idea of structural equivalence is stochastic blockmodelling [23], which considers generative models with parameters fitted to the observed network. The approach brings greater flexibility because different models can accommodate different types of communities [24]. Therefore, this approximation could be used to search for both modules and guilds [25]. There are, however, also caveats to the approach, since it is a challenge to determine whether the underlying assumptions of a particular block model is appropriate for the data being used [26]. Moreover, even when the model brings an analytically closed form, the estimation of the parameters may be computationally intractable [27], hence requiring costly optimality procedures or tests for robustness [28].

In this work, we build on the idea of structural equivalence noting that a node belonging to either a guild or a module is, in both cases, structurally equivalent to the other nodes in its community. This observation was acknowledged in social sciences in the definition of  $\lambda$ -communities [29], which are types of communities encompassing both modules and guilds, whose relevance has also been previously recognized in the ecological literature [5, 30]. From this observation, we wondered whether it is possible to find a similarity measure between nodes that quantifies their structural equivalence, even when different types of links are considered. We could then join nodes according to this similarity measure while monitoring whether the communities that are formed are guilds or modules. A similar approach was investigated by Yodzis *et al.* to measure trophic ecological similarity [31], but they did not identify an appropriate threshold for determining community membership (which they call “trophospecies”).

We have developed an approach that builds on these results and develops a method to determine objective thresholds for identifying modules and guilds in ecological networks. We show that a modification of the community detection method developed by Ahn *et al.* [32], leads to the identification of two quantities we call internal and external partition densities. For a set of nodes joined within a community by means of their structural equivalence similarity, the partition densities quantify whether their similarities come from connections linking them with nodes outside the community (external density) or within the community (internal density). Notably, our method generates maximum values for the two partition densities along the clustering, allowing us to objectively determine thresholds for the similarity measure in which the communities correspond to the definition of modules (for the internal density) and guilds (external density). Since the elements within both types of communities are structurally equivalent, modules and guilds can be understood as different kinds of *functional groups* —in the Eltonian sense— and this is the name we adopt here. We reserve the term “community” for a more generic use, because other types of communities beyond functional groups may exist, such as core-periphery structures [33].

We call our method functionInk (functional linkage), emphasizing how the number and types of links of a node determine its functional role in the network. We illustrate its use by considering complex ecological examples, for which we believe the notion of functional role is particularly relevant. We show in the examples that, by combining the external and internal partition densities, we are able to identify the underlying dominant structures of the network (either towards modules or towards guilds). Moreover, selecting the most appropriate community definition in each situation provides results that are comparable to state-of-the-art methods. This versatility in a single algorithm, together with its low computational cost to handle large networks, makes our method suitable for any type of complex, multidimensional network.

## 2 Methods

### Structural equivalence similarity in multidimensional networks

Our method starts by considering a similarity measure between all pairs of nodes that quantifies the fraction of neighbours connected with links of the same type that they share (Fig. 1). This is a natural definition

106 of structural equivalence for multidimensional networks, which is agnostic to the specific information that the  
 107 interaction carries. For simplicity, we present a derivation for a network that contains two types of links. We  
 108 use undirected positive (+) interactions (e.g. a mutualism) and negative (-) interactions (e.g. competition) to  
 109 illustrate the method, but these could be replaced by any two link types. Extending the method to an arbitrary  
 110 number of link types is presented in the Suppl. Material. We call  $\{i\}$  the set of  $N$  nodes and  $\{e_{ij}\}$  the set of  $M$   
 111 links in a network. We call  $n(i)$  the set of neighbours of  $i$ , that can be split into different subsets according to  
 112 the types of links present in the network.

113 For two types of links, we split the set of neighbours linked with the node  $i$  into those linked through positive  
 114 relationships,  $n_+(i)$ , or through negative relationships,  $n_-(i)$ ; we follow a notation similar to the one presented  
 115 in [32], but note that  $n(i)$  there denotes neighbours irrespective of the type of links. Distinguishing link types  
 116 induces a division in the set of neighbours of a given node into subsets sharing the same link type, shown in Fig.  
 117 2A. More specifically, in the absence of link types we define the Jaccard similarity between two nodes  $i$  and  $j$  as:

$$S^{(J)}(i, j) = \frac{|n(i) \cap n(j)|}{|n(i) \cup n(j)|} \quad (1)$$

118 where  $|\cdot|$  is the cardinality of the set (the number of elements it contains). This metric was shown to lead to  
 119 clusters of species that are more consistent with cophenetic clustering than other alternatives [31], and generalizing  
 120 this expression to multiple attributes is achieved simply by differentiating the type of neighbours depending on  
 121 the types of connections. For two attributes (see Suppl. Material for an arbitrary number of attributes) this leads  
 122 to

$$S^{(J)}(i, j) = \frac{|n_+(i) \cap n_+(j)| + |n_-(i) \cap n_-(j)|}{|n_+(i) \cup n_+(j) \cup n_-(i) \cup n_-(j)|}. \quad (2)$$

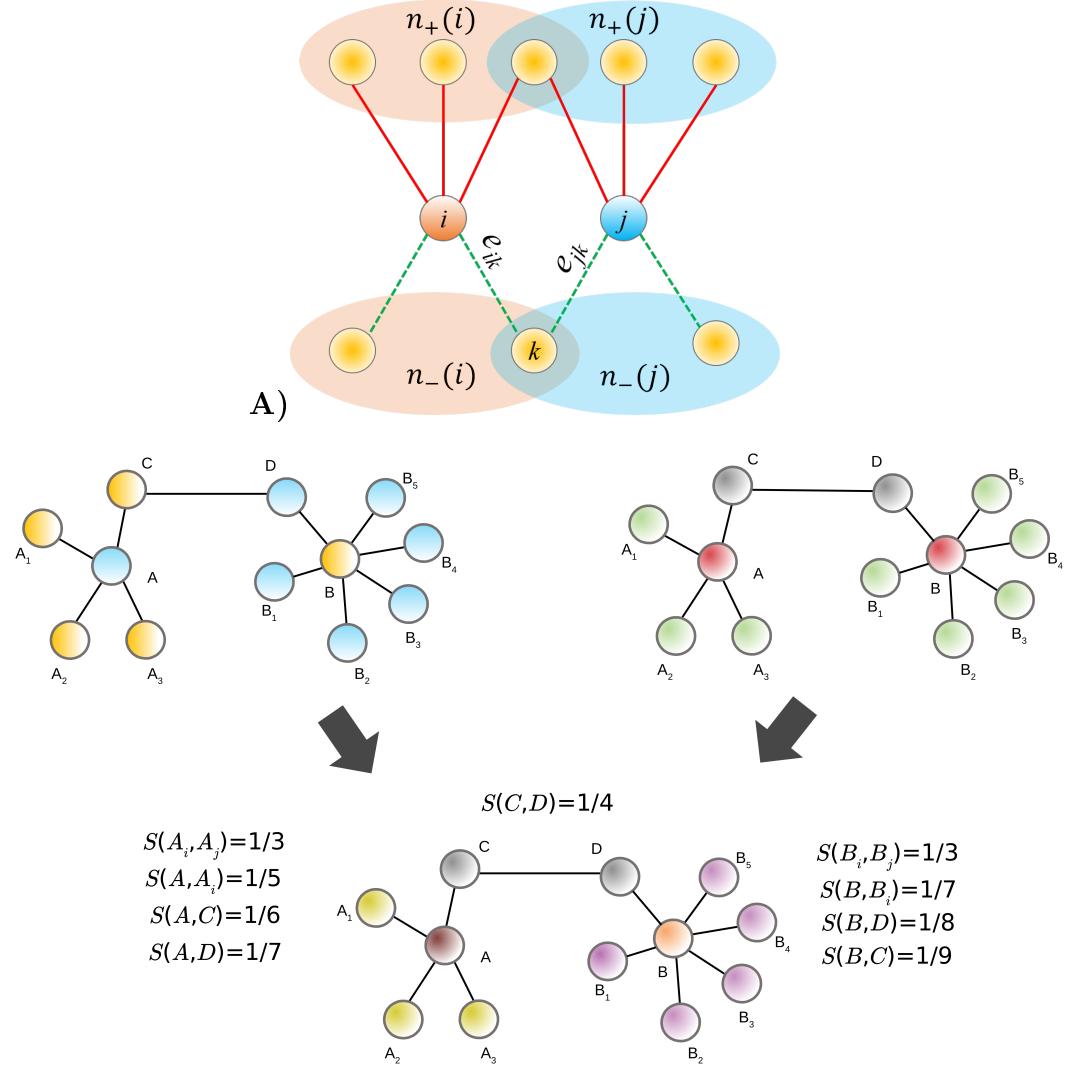
123 Accounting for the weight of the links can be made with the generalization of the Jaccard index provided by  
 124 the Tanimoto coefficient [35],  $S^{(T)}(i, j)$ , presented in Suppl. Material.

125 Finally, we introduce a modification to the above definition of  $S^{(J)}$  to account for the particular case in which  
 126  $i$  and  $j$  are only connected between themselves, i.e. they do not share any neighbours according to the above  
 127 definition. This is problematic because we want to distinguish this situation from the one in which they do not  
 128 share any nodes, for which we get  $S(i, j) = 0$ . We resolve this situation by considering that a node is its own  
 129 neighbour, in which case two nodes only connected between themselves would yield  $S(i, j) = 1$ . However, we note  
 130 that this would also be the value between two nodes that are connected and that also share all neighbours (a  
 131 motif known as a clique), irrespective of the number of neighbours they share, because the similarity measure  
 132 saturates. We argue that this situation is unsatisfactory because there is stronger evidence that two nodes are  
 133 structurally equivalent when they share connections and creating transitive motifs, since transitivity is a key  
 134 property in the definition of equivalence classes [36]. The situation can be resolved by using the convention that,  
 135 for two connected nodes, the intersection set is reduced by one, i.e.  $|n(i) \cap n(j)| \rightarrow |n(i) \cap n(j)| - 1$ . This  
 136 convention has the interesting property that, for cliques, increasing the number of nodes involved also increases  
 137 the similarity between its members, resulting in an upper bound of one and a lower bound of  $1/2$  (a 2-node clique).  
 138 In addition, two connected nodes that share neighbors but are not connected themselves have a smaller difference  
 139 in the similarity compared to nodes within cliques that share the same number of neighbors, thus facilitating the  
 140 identification of guilds. In Fig. 1 we illustrate the computation of this similarity with a simple example.

## 141 Identification of communities through clustering and similarity cut-offs

142 Once the similarity between nodes is computed, the next objective is to define and identify structurally equivalent  
 143 communities. As explained in the Introduction, there are different possible definitions of structural equivalence,  
 144 illustrated in Fig. 1B. In the figure is shown how the similarity metric proposed together with an agglomerative  
 145 clustering to join nodes in communities, encapsulates these different notions of structural equivalence. A critical  
 146 question, however, is how to objectively determine the threshold to stop the clustering ([37])?

147 This question is often addressed by iteratively “partitioning” the network into the distinct communities, and  
 148 monitoring each partition with a function having a well defined maximum or minimum that determines the  
 149 threshold of the optimal partition. In [32], the authors proposed to join links of a network according to a similarity  
 150 measure between the links with an agglomerative clustering, and to monitor the clustering with a quantity called  
 151 the partition density. The partition density is the weighted average across communities of the number of links  
 152 within a community out of the total possible number of links (which depends on the number of nodes in the



**Figure 1: Illustration of the method.** (A) The similarity between nodes  $i$  and  $j$  is computed considering the neighbours of each node and the types of interactions that link them. In this example, two types of link are shown: positive (+) interactions are solid links connecting the sets of neighbours  $n_+(i)$  and  $n_+(j)$ . Negative (−) links are shown as dotted links connecting the sets of neighbours  $n_-(i)$  and  $n_-(j)$ . Following Eq. 2,  $|n(i) \cap n(j)| = 2$  and  $|n(i) \cup n(j)| = 8$ , which yields  $S^{(J)}(i, j) = 2/8$ . If, for instance,  $e_{ik}$  changes from being − to +, the node  $k$  would no longer belong to the set  $n(i) \cap n(j)$ , being the new similarity:  $S^{(J)}(i, j) = 1/8$ . In Ref. [32] the similarity computed in this way is assigned to the links  $e_{ik}$  and  $e_{jk}$ . (B) Structural equivalence can be defined in different ways. In the top-left network we considered that blue and yellow colors encode *a priori* information describing the roles of the nodes. Identifying sets of nodes connected similarly to nodes with equivalent roles (i.e. the emphasis is on the roles and not on the specific neighbours, a situation called regular equivalence [34]) leads to two communities (the yellow and blue sets of nodes themselves), because every blue node is connected to a yellow one. The method of Guimerá and Amaral [33], determines communities focusing on their topological role (top-right network) by identifying central (A and B), peripheral (A1-A3 and B1-B5) and connector nodes (C and D). functionInk (bottom network) defines communities by joining nodes with approximately the same neighbours and, if there are roles for the nodes, these can be incorporated defining link types (one type for each pair of roles connected, in the example only one type is needed). All non-zero Jaccard similarities of the example are shown. Clustering these similarities will lead to different partitions and, stopping at  $S^{(J)} = 1/4$ , communities being the intersection of those found in the above networks are obtained, highlighting the potential to identify communities considering both the roles and topological features. Figure adapted from [33].

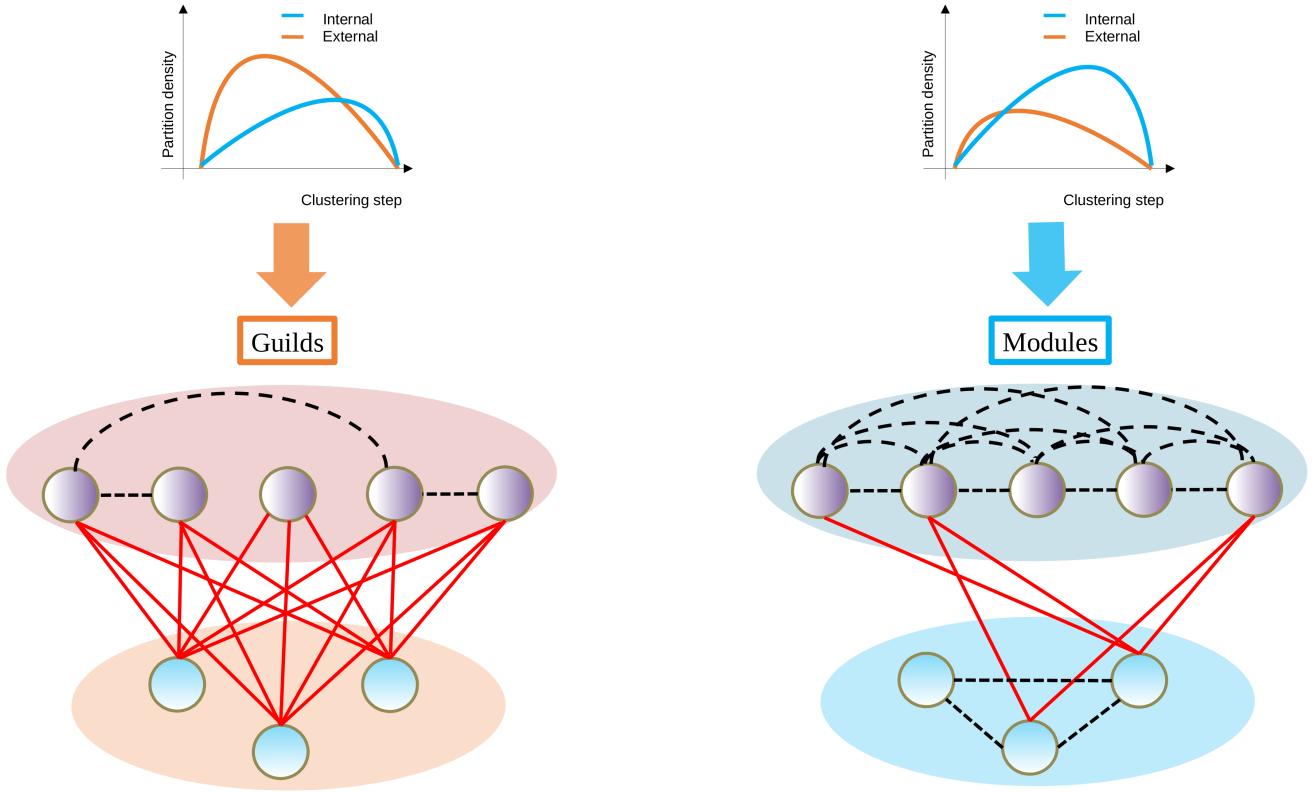


Figure 2: **Definition of guilds and modules.** For each set of nodes  $n_c^{\text{int}}$  belonging to the same community  $c$  (nodes within the same shaded area) we consider the number of links within the community (black dashed lines, called  $m_c^{\text{int}}$  in the main text) out of the total number of possible internal links, to compute the internal partition density (see upper curves). We also computed the external partition density, which is the density of links connecting nodes external to the community ( $m_c^{\text{ext}}$ , solid red lines linking nodes belonging to different communities) out of the total number of possible external links. We call guilds the communities determined at the maximum of the external partition density, and modules those found at the maximum of the internal partition density. The relative value of the external and internal partition densities allow us to estimate which kind of community dominates the network. In the example, guilds dominate the network on the left, and modules dominates the network on the right.

community). We re-considered the method of Ahn et al. [32] (which was originally defined over partitions of links, see Suppl. Material), to work over partitions of nodes, and we developed two partition densities, with two distinct meanings. To develop these measures we noted that, when joining nodes into a cluster, we are concluding that these nodes share (approximately) the same neighbours connected with the same type of links, but the nodes joined may or may not be connected between them. We therefore redefined the partition density so that it distinguishes between the contribution to the link density arising from the connections *within* a community from connections shared with external nodes *between* communities.

Formally, given a node  $i$ , we differentiate neighbours that are within the same community ( $n^{\text{int}}(i)$ , where int stands for “interior”) from neighbours that are in different communities ( $n^{\text{ext}}(i)$ ), hence  $n(i) = n^{\text{int}}(i) \cup n^{\text{ext}}(i)$  (Fig. 2). For a singleton (a community of size one)  $n^{\text{int}}(i) = \{i\}$  and  $n^{\text{ext}}(i) = \emptyset$ . Similarly, the set of links  $m(i)$  connecting the node  $i$  with other nodes can also be split into two sets: the set connecting the node with neighbours within its community  $m^{\text{int}}(i)$ , and those connecting it with external nodes  $m^{\text{ext}}(i)$ . This distinction was also considered in other context (called the problem of coloring nodes [38]).

Therefore, for each partition of nodes into  $T$  communities our method identifies, for each community  $c$ , the total number of nodes it contains,  $n_c^{\text{int}}$ , and the total number of links connecting these nodes  $m_c^{\text{int}}$ . In addition, it computes the total number of nodes in other communities that have connections to the nodes in the community,  $n_c^{\text{ext}}$ , through a number of links  $m_c^{\text{ext}}$ . Clearly, to identify  $n_c^{\text{ext}}$  neighbours, at least  $n_c^{\text{ext}}$  links are required and thus an increasing number of links in excess,  $m_c^{\text{ext}} - n_c^{\text{ext}}$  are necessary to obtain an increasing contribution to the

similarity of the nodes in the community through external links (however, this is not a sufficient condition, see Suppl. Material). In this way, a relevant quantity to characterize a community is the fraction of links in excess out of the total possible number  $(m_c^{\text{ext}} - n_c^{\text{ext}})/n_c^{\text{ext}}(n_c^{\text{int}} - 1)$ . We note this calculation does not take into account multiple link types. The weighted average of this quantity through all communities leads to the definition of external partition density:

$$D^{\text{ext}} = \frac{1}{M} \sum_c \frac{m_c^{\text{ext}}}{2} \frac{(m_c^{\text{ext}} - n_c^{\text{ext}})}{n_c^{\text{ext}}(n_c^{\text{int}} - 1)}, \quad (3)$$

where  $M$  is the total number of links. We now follow a similar reasoning to consider a necessary condition to obtain an increasing contribution to the similarity of the nodes through the internal links (see Suppl. Material). We acknowledge that in a community created by joining nodes through the similarity measure we propose, it may happen that  $n_c^{\text{int}} > 0$  even if  $m_c^{\text{int}} = 0$ . Therefore, any link is considered a link in excess, leading to the following expression for the internal partition density, which quantifies the fraction of internal links in excess out of the total:

$$D^{\text{int}} = \frac{1}{M} \sum_c m_c^{\text{int}} \frac{2m_c^{\text{int}}}{n_c^{\text{int}}(n_c^{\text{int}} - 1)}. \quad (4)$$

Finally, we define the total partition density as the sum of both internal and external partition densities:

$$D^{\text{total}} = D^{\text{int}} + D^{\text{ext}},$$

and hence, if all the fractions in  $D^{\text{int}}$  and  $D^{\text{ext}}$  are equal to one, i.e. all possible links in excess are realized,  $D^{\text{total}}$  equals to one. Since at the beginning of the clustering the communities have a low number of members, most of the contribution towards  $D^{\text{total}}$  comes from  $D^{\text{ext}}$  while, in the last steps, where the communities become large,  $D^{\text{int}}$  will dominate. All three quantities will reach a maximum value along the clustering (for the internal it could be at the last step) and, if one of them clearly achieves a higher value, it will be indicative that one type of functional group is dominant in the network. If that is the case, the maximum of  $D^{\text{total}}$  —which is always larger or equal to  $\max(\max(D^{\text{int}}), \max(D^{\text{ext}}))$ — will be at a clustering step close to the step in which the dominating quantity peaks. If neither  $D^{\text{ext}}$  nor  $D^{\text{int}}$  clearly dominates,  $D^{\text{total}}$  will peak at an intermediate step between the two partial partition densities maxima, suggesting that this intermediate step is the best candidate of the optimal partition for the network. Communities determined at this intermediate point where they can be both guilds and modules will be called, generically, functional groups.

### 3 Results

#### 195 Plant-pollinator networks

To illustrate the use of the method we start analyzing a synthetic example. In ecological systems, species are often classified into communities according to their ecological interactions, such as in mutualistic networks of flowering plants and their animal pollinators. These networks are characterized by intra and interspecific competition within both the pool of plants and the pool of animals, and by mutualistic relationships between plants and animals, leading to a bipartite network.

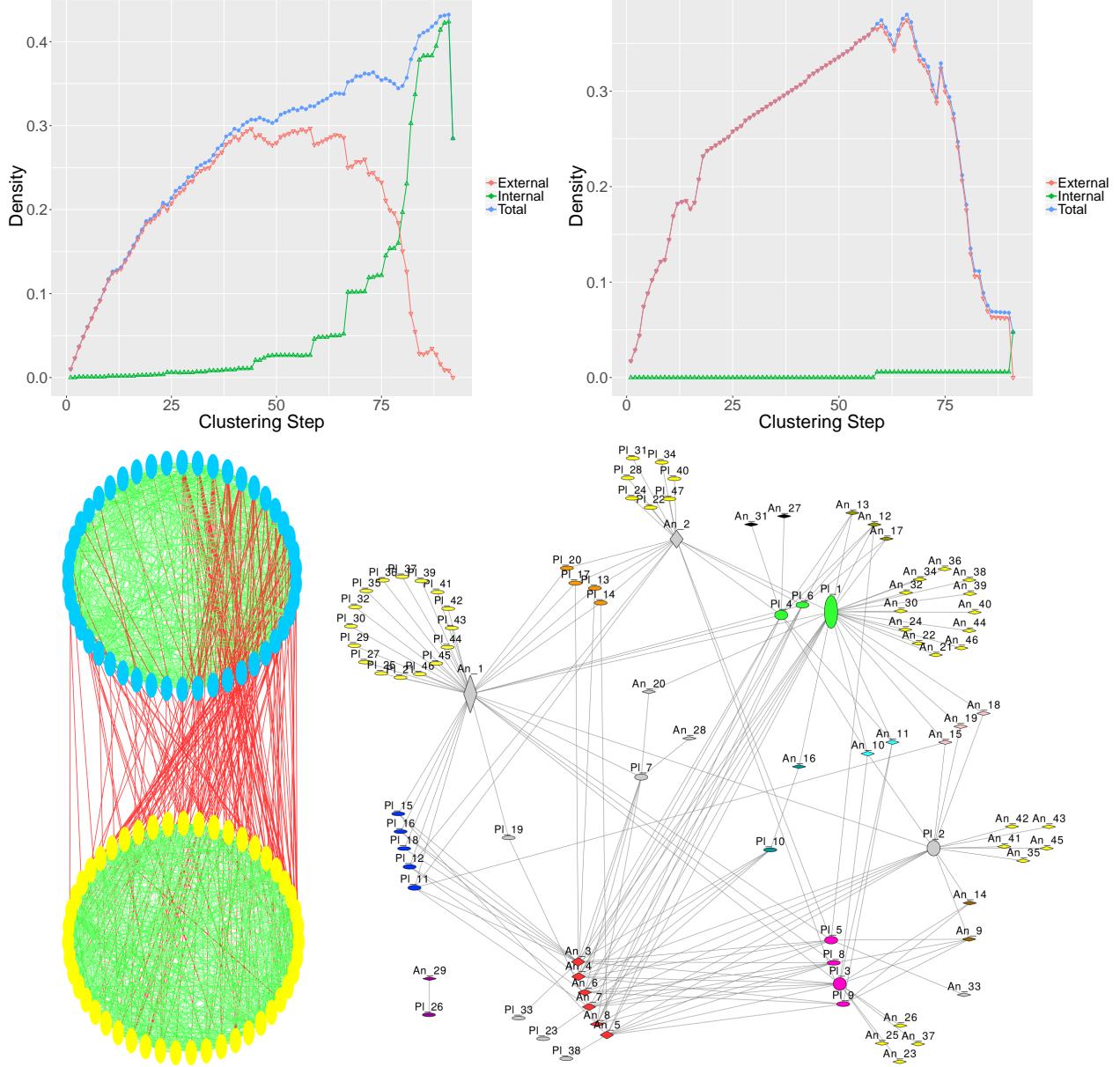
To investigate the performance of our method and, in particular, the influence of the topological properties into the partition density measures, we generated a set of artificial mutualistic networks with diverse topological properties, following the method presented in [39]. For the mutualistic interactions, we focused on two properties: the connectance  $\kappa_{\text{mut}}$ , which is the fraction of observed interactions out of the total number of possible interactions, and the nestedness  $\nu$  as defined in [40] (see Methods), which codifies the fraction of interactions that are shared between two species of the same pool, averaged over all pairs of species. We selected these measures for their importance in the stability-complexity debate in mutualistic systems [39], and the similarity between the nestedness (which, in the definition we adopt here, represents the mean ecological overlap between species) and the notion of structural equivalence we considered (see Suppl. Material). For the competition matrices, we considered random matrices with different connectances,  $\kappa_{\text{comp}}$ , since it is difficult to estimate direct pairwise competitive interactions experimentally, and they are frequently modeled with a mean field competition matrix.

We verified that in all networks the set of plants and animals are joined in the very last step of the clustering irrespective of the clustering method used, a result that must follow construction. As expected, the curves

monitoring the external and internal partition densities depends on the properties of the networks. We illustrate this finding in Fig. 3, where we have selected two networks with contrasting topological properties. One of the networks has high connectance within the pools and low connectance and nestedness between the pools. The internal partition density peaks at the last step minus one (i.e. where the two pools are perfectly separated) consistent with the definition of modules, where the intra-modules link density is higher than the inter-modules link density. On the other hand, the second network has intra-pool connectance equals to zero, and very high connectance and nestedness between the pools (see Fig. 3). We selected a  $\kappa_{\text{comp}} = 0$  for simplicity in the network representation, but similar results are obtained for low values of  $\kappa_{\text{comp}}$ , see for instance Suppl. Fig. 10. In this second network (see Fig. 3, right panel), only the external partition density peaks and, at the maximum, the communities that we identified clearly reflect the structural equivalence of the nodes members in terms of their connectance with nodes external to the group, as we expect for the definition of guilds. The ecological information retrieved for guilds is clearly distinct from the information retrieved for the modules, the former being related to the topology of the network connecting plants and animals. We observe that guilds identify specialist species clustered together, which are then linked to generalists species of the other pool: a structure typical of networks with high nestedness.

The method identified several interesting guilds and connections between them. For instance, generalists Plant 1, Animal 1 and Animal 2 (and to a lesser extent Plant 2) have a low connectivity between them but, being connected to many specialists, determine a region of high vulnerability, in the sense that a directed perturbation over these species would have consequences for many other species. This is confirmed by the high betweenness of these nodes (proportional to the size of the node in the network). In addition, the algorithm is able to identify more complex partitions of nodes into communities. As an example of this, Animal 16 (turquoise) is split from Animals 10 and 11 (cyan), which form a second community, and from Animals 15, 18 and 19 (light pink) that are joined into a third community, despite of the subtle connectivity differences between these six nodes. Finally, it also detects communities of three or more species that have complex connectivity patterns which, in this context, may be indicative of functionally redundant species (e.g. red and blue communities).

Examples with other intermediate properties are analyzed in the Suppl. Figs. 8 and 9. Broadly speaking, either the internal or the total partition density maximum peaks at the last step minus one, allowing for detection of the two pools of species. Nevertheless, the method fails to find these pools in situations in which the similarity between members of distinct pools is comparable to the similarity of members belonging to the same pool. This may be the case if the connectances are small (see Suppl. Fig. 10). The relative magnitude of the external vs. internal partition density depends on the connectance between the pools of plants and animals and on the connectance within the pools, respectively (see Suppl. Fig. 8). Interestingly, networks for which the nestedness is increased keeping the remaining properties the same generated an increase in the external partition density (see Suppl. Fig. 9). These examples illustrate how the external partition density is sensitive to complex topological properties, in particular to an increase in the dissasortativity of the network, as expected when guilds are dominant.



**Figure 3: Analysis of synthetic mutualistic networks.** (Top left) Partition densities for a network with  $\kappa_{\text{comp}} = 0.5$ , nestedness  $\nu = 0.15$  and  $\kappa_{\text{mut}} = 0.08$  and (top right) for a network with  $\kappa_{\text{comp}} = 0$ , nestedness  $\nu = 0.6$  and  $\kappa_{\text{mut}} = 0.08$ . The high density of competitive links in the first network makes the internal partition density dominate, leading to two modules representing the plant-pollinator pools (bottom left network), while reducing the density of competitive links to zero in the second network makes the external partition density to dominate, finding guilds (bottom right, with plants labeled “Pl” and animals labeled “An”). The small increase in the internal partition density for this network at step 59 is due to two specialist species joined at that step (animal 29 and plant 56, shown at the bottom left of the network). Nodes are colored according to their functional group in both networks although, in the network finding guilds (bottom right), specialist species are yellow, single species communities are gray, and the size of the nodes is proportional to their betweenness.

## 249 Trophic networks

250 We tested our method in a comprehensive multidimensional ecological network of 106 species distributed in trophic  
 251 layers with approximately 4500 interactions, comprising trophic and non-trophic interactions (approximately 1/3  
 252 of the interactions are trophic) [41]. This network was analyzed looking for communities extending a stochastic

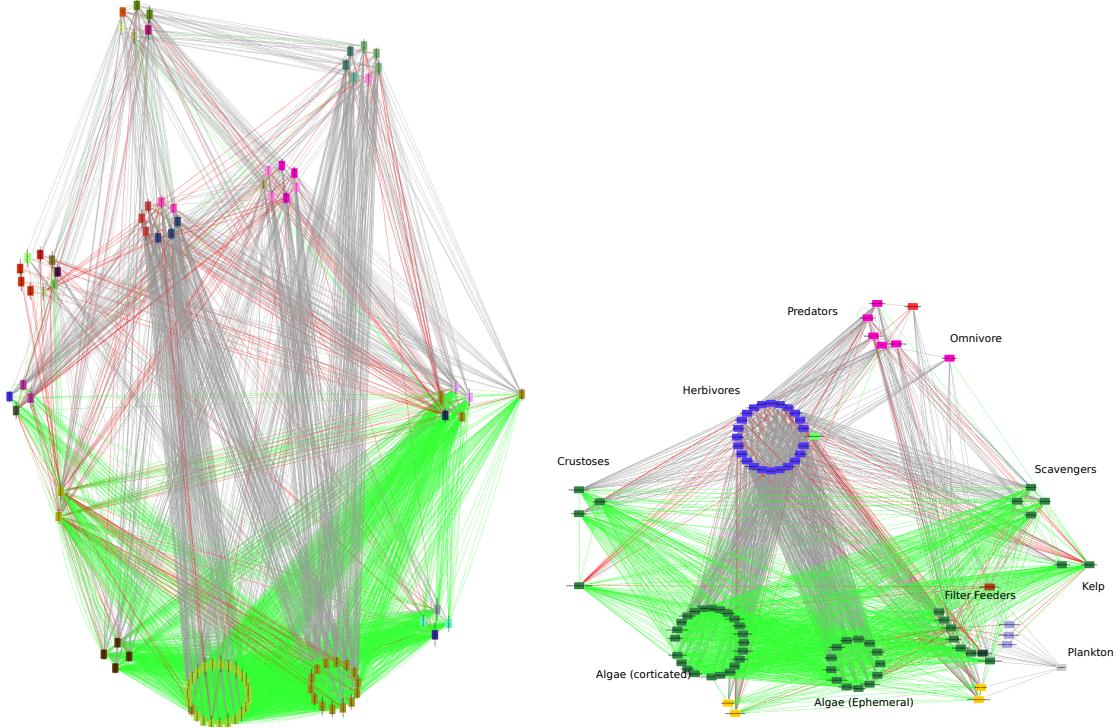


Figure 4: **Determination of guilds and modules in a large trophic network.** Trophic networks with links representing trophic (gray), non-trophic positive (red), and negative (green) interactions. (Left) Nodes are grouped according to the classification found in [41] (reference classification), and colored by the guilds found with functionInk at the maximum of the external partition density. (Right) Nodes are grouped according to the trophic levels and colored by the modules found by functionInk (see Main Text for details). The modules separate the three main trophic levels: predators, herbivores and basal species, further separating some of them into subgroups, such as filter feeders and plankton, which is an orphan module.

blockmodelling method [12] to deal with different types of interactions [41]. The estimation of the parameters of the model through an Expectation-Maximization algorithm requires controlling the influence of random starting conditions since each initial condition may lead to a different result, and hence is needed to test the robustness of the results. Here we show that, in this example, our method is comparable with this approximation, and it has the advantage of being deterministic. Moreover, the simplicity of the method allows us to handle large networks with arbitrary number of types of links and to evaluate and interpret the results, as we show in the following.

Our method finds a maximum for the internal density when there are only three communities. Previous descriptions of the network identified three trophic levels in the network (Predators, Herbivores and Basal species). The latter are further subdivided into subgroups like (e.g. Kelps, Filter feeders), and there are some isolated groups like one Omnivore and Plankton. To match these subgroups we observed that the total partition density reaches a maximum close to the maximum of the external partition density (step 69) and maintains this value along a plateau until step 95 (see Suppl. Fig. 11). We analyzed results at both clustering thresholds finding that, at step 95, we obtain modules with a good agreement with the trophic levels, shown in Fig. 4. On the other hand, at step 69 we find a larger number of communities, some of which fit the definition of modules and others the definition of guilds (see Fig. 4).

To shed some light on the information obtained from this second network, we compared the classification obtained by Kefi *et al.* [41] (in the following reference classification) and our method. We computed several similarity metrics comparing the classification we obtained at each step of the agglomerative clustering with functionInk and the reference classification (see Methods). In Fig. 5, we show that the similarity between both classifications is highly significant ( $Z$ -score  $> 2.5$ ) and is maximized when the external partition density is also maximized, i.e. at step 69. This is particularly apparent for the Wallace 01, Wallace 10 and Rand indexes (see Fig. 6 and Suppl. Fig. 12). Notably, communities in the reference classification were also interpreted as functional groups in the same sense proposed here [41].

Nevertheless, there are some discrepancies between both classifications. In particular, although there is a complete correspondence between the two largest communities in both classifications, there are a number of intermediate communities in the reference classification whose members are classified differently in our method. To illustrate these discrepancies, we plotted a heatmap of the Tanimoto coefficients of members of four communities of intermediate sizes containing discrepancies, showing their membership in both the reference and the functionInk classification with different colors (see Fig. 5). The dendrograms cluster rows and columns computing the Euclidean distance between their values. Therefore, these dendrograms are very similar to the method encoded in functionInk, and the communities must be consistent, representing a powerful way to visually inspect results. Indeed, the dendrograms are in correspondence with both functionInk and reference communities, but we observe some discrepancies. For instance, the community found by the reference classification containing several *Petrolisthes* species, joins species that have low similarity regarding the number and type of interactions as measured by the Tanimoto coefficients, while functionInk joins together the three species with high similarity, leaving aside the remainder species. Therefore, despite the methodological differences between both methods, the different classifications produce similar outcomes, but result in different sized clusters (a different cut-off in the dendrograms, with functionInk finding finer clusters). The advantage of functionInk is then apparent in the simplicity of the method, which permits validation through visual inspection of the consistency of the classification.

## Microbial networks

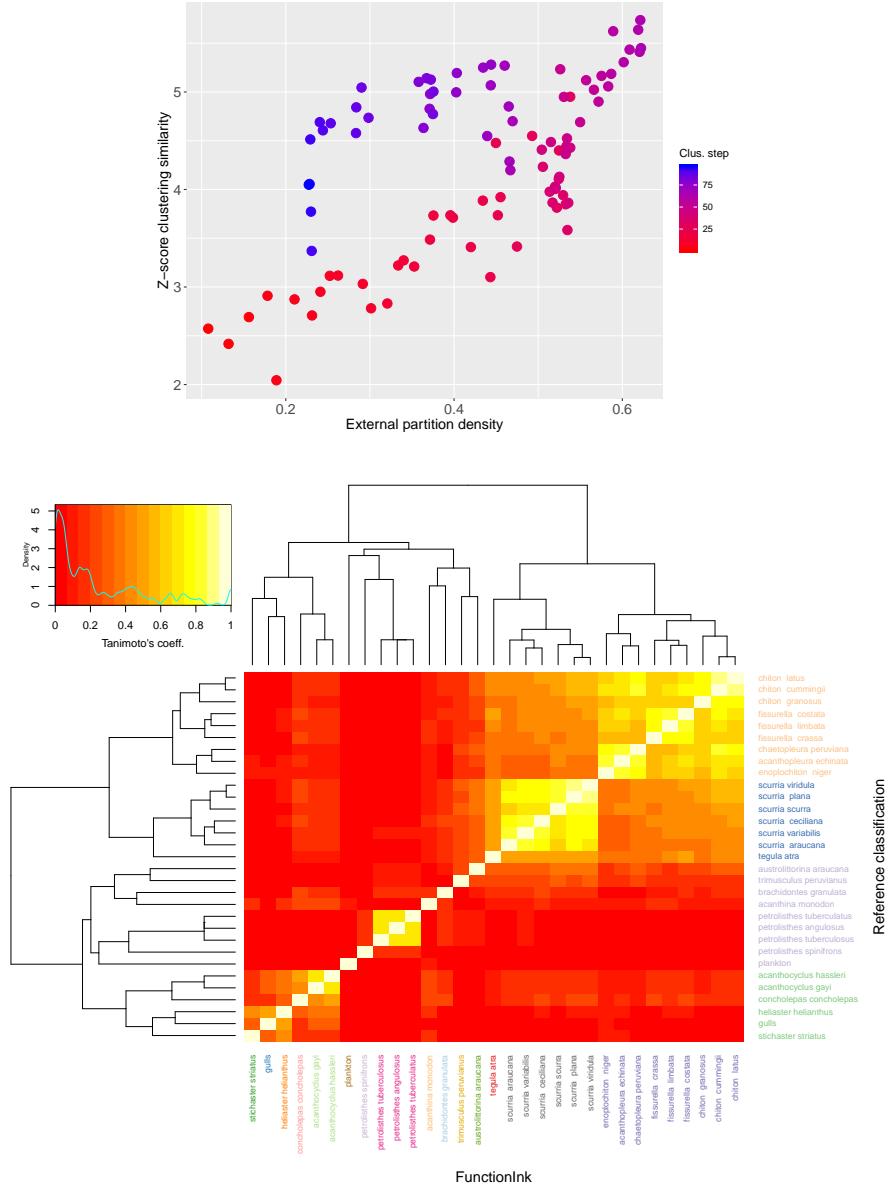
We discuss a last example of increasing importance in current ecological research, which is the inference of interactions among microbes sampled from natural environments. We considered a large matrix with more than 700 samples of 16S rRNA operative taxonomic units (OTUs) collected from rain pools (water-filled tree-holes) in the UK [43, 44] (see Suppl. Material). We analyzed  $\beta$ -diversity similarity of the samples contained in the matrix with the Jensen-Shannon divergence metric [45], further classifying the samples automatically, leading to 6 disjoint clusters we call  $\beta$ -diversity-classes (i.e. clusters of samples, see Methods). Next, we inferred a network of significant positive (co-occurrences) or negative (segregations) correlations between OTUs using SparCC [46] (see Methods), represented in Suppl. Fig. 13. Applying functionInk to the network of inferred correlations, we aimed to understand the consistency between the results of functionInk (modules and guilds) and the  $\beta$ -diversity-classes. The rationale is that, by symmetry, communities determined from significant co-occurrences and segregations between OTUs should reflect the similarity and dissimilarity between the samples, hence validating the method.

Contrasting with the trophic network analyzed in the previous example, the external partition density brings a poor reduction of the complexity of the network (peaking after only 22 clustering steps), and the internal partition density is higher, hence suggesting a more relevant role for modules (see Suppl. Fig. 14). Differences in the three stopping criteria are shown in Suppl. Fig. 13, where two large modules are apparent, with a large number of intracluster co-occurrences (continuous links) and interclusters segregations (dotted links). Note that this is quite different to what is found in macroscopic trophic networks, where pools of species (e.g. prey) have within module competitive (segregating) interactions, while between-modules interactions can be positive (for predators) or negative (for prey).

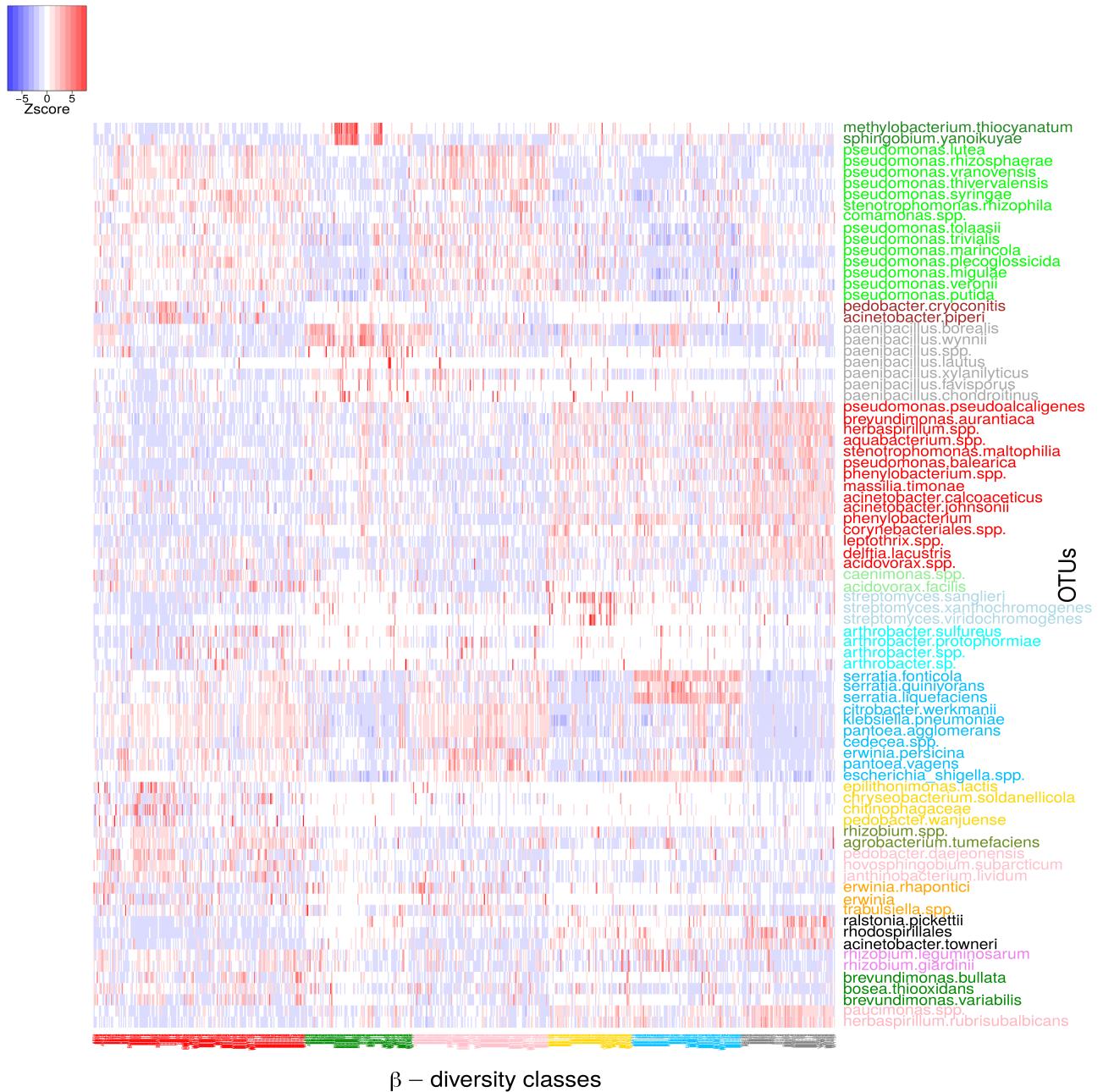
There is reasonable agreement between the functional groups found at the maximum of the total partition density and the  $\beta$ -diversity-classes, shown in Fig. 6. Moreover, the detection of networks complements the information that  $\beta$ -diversity-classes provides, since it is possible to individuate the key players of these classes (see Suppl. Material). Notably, it was shown in [43] that the  $\beta$ -diversity-classes might be related to a process of ecological succession driven by environmental variation, the functional groups are likely driven by environmental preferences rather than by ecological interactions, likely explaining the large number of positive co-occurrences. This speaks against a naive interpretation of correlation networks in microbial samples as ecological interactions unless environmental preferences are under control [9, 8].

## Discussion

We presented a novel method for the analysis of multidimensional networks, with nodes with an arbitrary number of link types. We implemented the method adopting the definition of structural equivalence, which underlies both the similarity measure definition and the rationale behind both the clustering and our definition of partition densities. We selected a set-theoretic similarity measure quantifying the number of nodes that are shared with the same type of interaction, which we believe is a natural definition of structural equivalence for multidimensional networks, and that has the advantage that it does not make assumptions on how the information flows in the



**Figure 5: Comparison between the reference classification in the trophic network and functionInk.** (Top) Z-score of the Wallace 10 index [42], measuring the similarity between the reference classification and the functionInk method at each clustering step. The similarity with the reference classification (see Main Text) is maximized around the maximum of the external partition density. (Bottom) Comparison of communities 1, 4, 7 and 9 in the reference classification, whose members were classified differently by functionInk. Colors in the names of species in rows (columns) represent community membership in the reference (functionInk) classifications. The heatmap represents the values of the Tanimoto coefficients, and the dendrograms are computed using Euclidean distance and clustered with complete linkage. Both classifications are generally consistent with the dendrograms, but with functionInk finding finer clusters.



**Figure 6: Comparison between  $\beta$ -diversity classes and functional groups in a microbial network.** Heatmap representing the z-score of the log-transformed abundances of the OTUs (see Methods). Species are colored according to their functional group membership obtained at the maximum of the total partition density. Samples are colored according to one of the six community classes found in [43] after optimal clustering with a  $\beta$ -diversity distance. Orphan clusters were excluded except for 5 *Paenibacillus* species (characteristic of the green class) that were added to the functional group formed by *Paenibacillus borealis* and *Paenibacillus wynii*. The heatmap blocks show segregation and co-occurrence between modules, further mapping the  $\beta$ -diversity classes.

327 network, typical of approximations based on Laplacian dynamics (see e.g. [18, 19]). This allow us to join  
328 nodes simply by their similarity, with no need for specific assumptions about the network structure. Moreover,  
329 this similarity can also be naturally linked to two measures of nodes' partitioning that allowed us to propose a  
330 clear differentiation between modules (determined by the maximum of the internal partition density) and guilds  
331 (determined by the maximum of the external partition density).

332 Beyond these technical advantages, we illustrated the versatility of functionInk using several ecological ex-  
333 amples. The relative value between the internal and external partition density immediately yields information  
334 on whether the network is dominated by modules, guilds, or intermediate structures. This allows for increasing  
335 flexibility in the analysis of the networks, and for a more nuanced interpretation of network structure and species'  
336 roles in the ecosystem. For both mutualistic and trophic networks, the internal partition density correctly finds  
337 the trophic layers, justifying the success of the original method [32]. Our extension recovered the functional groups  
338 as determined by Kefi et al. [41] through the external partition density, and the visual inspection reflects a good  
339 consistency with the definition we proposed for functional groups in terms of structural equivalence. Moreover,  
340 in the mutualistic networks, we showed that the functional groups discovered in this way was sensitive to changes  
341 to high-order topological properties such as the nestedness.

342 The analysis of the microbial network was dominated by modules rather than guilds. Interestingly, these  
343 modules had intra-cluster positive correlations, contrary to what would be expected in a macroscopic trophic  
344 network, where competitive interactions would be dominant between members of the same trophic layer. We  
345 selected in this example for further exploration the functional communities found at the maximum of the total  
346 partition density, with some groups having properties closer to those of guilds and others closer to modules. The  
347 communities that we identified were in good agreement with the functional communities found using  $\beta$ -diversity  
348 similarity [43], supporting the consistency of the method.

349 To finish, we highlight some limitations of the method. Firstly, it may have problems if the communities  
350 are highly overlapping [20, 32]. In these cases, it would be convenient to inspect the partition at the three  
351 classifications given by the different partition densities, since it is likely that overlapping communities are split  
352 in an earlier classification and then joined at later steps of the clustering. Another possibility is to combine it  
353 with the approximation proposed in Ref. [32], that has both compatible and, at the same time, complementary  
354 results (see Suppl. Material). To continue with, our approximation does not consider yet the case in which there  
355 are multiedges in the network, although real networks are typically very sparse and the probability of finding  
356 multiedges is small [26]. Finally, although the method might not be able to achieve the generality of other  
357 approximations aiming to find any arbitrary structure in the network [12, 24, 28], such approximations require  
358 either heuristics to find a solution for the parameters –and hence a unique optimal solution is not guaranteed–,  
359 or a computationally costly sampling of the parameter space. Our method relies on a deterministic method  
360 whose results are easily inspected, and its computational cost for a network with  $N$  nodes scales as  $N^2$  for  
361 the similarity metric computation plus the clustering, which is order  $N$ . The method is freely available in the  
362 address ([HTTPS://GITHUB.COM/APASCUALGARCIA/FUNCTIONINK](https://github.com/APASCUALGARCIA/FUNCTIONINK)) and, importantly, although we developed it  
363 with ecological networks in mind, it can be applied to any kind of network.

## 364 Acknowledgments

365 We thank Michael Schaub for critical comments on an earlier version of the manuscript, and to Sebastian Bonhoeffer  
366 for his support. We also thank to two anonymous referees whose comments help us to improve the manuscript.  
367 The research was funded by a European Research Council starting grant (311399-Redundancy) awarded to T.B.  
368 T.B. was also funded by a Royal Society University Research Fellowship. APG was also funded by the Simons  
369 Collaboration: Principles of Microbial Ecosystems (PriME), award number 542381.

## 370 Authors' Contributions

371 APG conceived the project and designed methodology; APG performed the analysis; TB contributed data; APG  
372 and TB analysed the data; APG led the writing of the manuscript. Both authors contributed critically to the  
373 drafts and gave final approval for publication.

## 374 Supplementary Methods and Results

### 375 Generalization of the Jaccard and Tanimoto coefficients to an arbitrary number of 376 link types

377 Consider a network with a set  $\{i\}$  of  $N$  nodes and a set  $\{e_{ij}\}$  of  $M$  links. These links are classified into  $\Omega$  types  
378 labeled with the index  $\alpha = (1, \dots, \Omega)$ . These types would typically account for differential qualitative responses  
379 of the nodes properties due to the interactions. For example, if we consider that the nodes are species and  
380 the property of interest is the species abundances, the effect of cooperative or competitive interactions on the  
381 abundances can be codified using two different types of links: positive and negative. If these relations are inferred  
382 through correlations between abundances, we could use a quantitative threshold (for instance a correlation equal  
383 to zero) to split the links into positive and negative correlations. In general, we may use a number of qualitative  
384 attributes or quantitative thresholds in the weights of the links to determine different types of links.

385 We call  $n(i)$  the set of neighbours of  $i$ , and we split these neighbours into (at most)  $\Omega$  different subsets  
386 according to the types of links present in the network. The Jaccard coefficient defined in Eq. 2 can be extended,  
387 considering similarities between nodes, as:

$$S^{(J)}(n_i, n_j) = \frac{\sum_{\alpha=1}^{\Omega} |n_{\alpha}(i) \cap n_{\alpha}(j)|}{|\bigcup_{\alpha=1}^{\Omega} n_{\alpha}(i) \cup n_{\alpha}(j)|}. \quad (5)$$

388 Accounting for the weight of the links can be made with the generalization of the Jaccard index provided by  
389 the Tanimoto coefficient [35]. We first introduce the method without differentiating between different types of  
390 neighbours. Consider the vector  $\mathbf{a}_i = (\tilde{A}_{i1}, \dots, \tilde{A}_{iN})$  with

$$\tilde{A}_{ij} = \frac{1}{k_i} \sum_{i' \in n(i)} w_{ii'} \delta_{ij} + w_{ij} \quad (6)$$

391 where  $w_{ij}$  is the weight of the link connecting the nodes  $i$  and  $j$ ,  $k_i = |n(i)|$  and  $\delta_{ij}$  is the Kronecker's delta  
392 ( $\delta_{ij} = 1$  if  $i = j$  and zero otherwise). Determining the quantity  $W_{ij} = \mathbf{a}_i \cdot \mathbf{a}_j = \sum_k \tilde{A}_{ik} \tilde{A}_{kj}$ , the Tanimoto similarity  
393 is defined as

$$S^{(T)}(e_{ik}, e_{jk}) = \frac{W_{ij}}{W_{ii} + W_{jj} - W_{ij}}. \quad (7)$$

394 Working with link types requires a generalization of the above expression. Consider for the moment two  
395 types related with a positive  $w_{ij} > 0$  or a negative  $w_{ij} < 0$  weight of the links. The term  $\tilde{A}_{ii} = 1/k_i \sum_{i'} w_{ii'}$   
396 is the average of the strengths of the links connected with node  $i$ , and it is desirable to keep this meaning when  
397 considering two types to properly normalize the Tanimoto similarity. This is simply achieved redefining  $\tilde{A}_{ij}$  as

$$\tilde{A}_{ij} = \frac{1}{k_i} \sum_{i' \in n(i)} \text{abs}(w_{ii'}) \delta_{ij} + w_{ij}. \quad (8)$$

398 On the other hand, the similarity is essentially codified in the term  $W_{ij}$  that we now want to redefine to  
399 account for two types of interactions in such a way that only products  $\tilde{A}_{ik} \tilde{A}_{kj}$  between terms with the same  
400 sign contribute to the similarity. This is achieved with the following definition, which generalizes the Tanimoto  
401 coefficient

$$W_{ij} = \sum_k \tilde{A}_{ik} \tilde{A}_{kj} \delta(\text{sgn}(\tilde{A}_{ik}) - \text{sgn}(\tilde{A}_{kj})) \quad (9)$$

402 where  $\text{sgn}(\cdot)$  is the sign function and  $\delta(a - b)$  is the Dirac delta function ( $\delta(a - b) = 0$  if  $a \neq b$ ). Generalizing  
403 to an arbitrary number of types can be achieved by defining a variable  $\mu_{ij}$  that returns the type of the link, i.e.  
404  $\mu_{ij} = \alpha$  with  $\alpha$  being a factor variable which, for the example of positive and negative links, is codified by the  
405 sign of the links' weight. We generalize the expression in Eq. 9 as follows

$$W_{ij} = \sum_k \tilde{A}_{ik} \tilde{A}_{kj} \delta(\mu_{ik} - \mu_{kj}). \quad (10)$$

406 Finally, the generalization of the external and internal partition densities to consider multiple types of links  
 407 simply requires us to correctly classify the neighbours of each node accounting for the different types  $n(i) =$   
 408  $\bigcup_{\alpha=1}^{\Omega} n_{\alpha}^{\text{int}}(i) \cup n_{\alpha}^{\text{ext}}(i)$ . Similarly, the set of links  $m(i)$  connecting the node  $i$  with other nodes must be also split  
 409 into sets according to the different types  $m(i) = \bigcup_{\alpha=1}^{\Omega} m_{\alpha}^{\text{int}}(i) \cup m_{\alpha}^{\text{ext}}(i)$ . The expressions for the internal and  
 410 external partition densities remain otherwise the same.

## 411 Additional notes on the relation between the Jaccard similarity and the partition 412 densities

413 In this section we aim to provide a more explicit relation between the Jaccard similarity and the partition density  
 414 and, in particular, the notion of links in excess. We will show that the existence of links in excess is a necessary  
 415 condition to compute the similarity between the nodes in a cluster, and increases if the number of links in  
 416 excess results in a higher similarity. Since the similarity can also be partitioned between external and internal  
 417 components, this condition can be independently applied for the external and internal links in excess. We will  
 418 finally discuss why these conditions are not sufficient.

419 Consider that the network has a single type of link (and we will make some precision below for the general  
 420 case in which there are different types of links). In the computation of the partition densities, we inspect each  
 421 community and we average across communities. Therefore, let us start considering a generic community  $c$ , and  
 422 note that the neighbours of a node  $i$  are partitioned into those belonging to the same community,  $n^{\text{int}}(i)$ , and  
 423 those belonging to other communities,  $n^{\text{ext}}(i)$ . The Jaccard similarity of two nodes within the same community  
 424 can be expressed as

$$S_c^{(\text{J})}(i, j) = \frac{|n^{\text{int}}(i) \cap n^{\text{int}}(j)| + |n^{\text{ext}}(i) \cap n^{\text{ext}}(j)|}{|n^{\text{int}}(i) \cup n^{\text{int}}(j) \cup n^{\text{ext}}(i) \cup n^{\text{ext}}(j)|} = S_c^{(\text{J}), \text{int}}(i, j) + S_c^{(\text{J}), \text{ext}}(i, j), \quad (11)$$

425 where we used a subindex  $c$  to stress the fact that  $i, j \in c$ . Hence, the mean similarity of the nodes in the  
 426 community can also be split in two components

$$\overline{S_c^{(\text{J})}(i, j)} = \frac{2}{n^{\text{int}}(n^{\text{int}} - 1)} \left( \sum_{i < j} S_c^{(\text{J}), \text{int}}(i, j) + S_c^{(\text{J}), \text{ext}}(i, j) \right) = \overline{S_c^{(\text{J}), \text{int}}(i, j)} + \overline{S_c^{(\text{J}), \text{ext}}(i, j)}.$$

427 The sums  $\hat{s}_c^{\text{int}} = \sum_{i < j} |n^{\text{int}}(i) \cap n^{\text{int}}(j)|$  and  $\hat{s}_c^{\text{ext}} = \sum_{i < j} |n^{\text{ext}}(i) \cap n^{\text{ext}}(j)|$  (where  $i, j \in c$ ) encode the  
 428 contributions to the mean similarity of the community given by external and internal nodes, respectively. We aim  
 429 to relate these terms to the links in excess used to compute the partition densities.

430 To establish this relationship, we note that the degree of any node in the network,  $|n(i)|$ , can also be partitioned  
 431 between the links connecting to other nodes in the community it belongs to,  $|n^{\text{int}}(i)|$ , and the links connecting  
 432 to nodes from other communities,  $|n^{\text{ext}}(i)|$ . Again, we will add a subindex, e.g.  $|n_c(i)|$ , if there is any possible  
 433 ambiguity about the identity of the community to which  $i$  belongs. In addition, we would like to specifically  
 434 identify how the degree  $|n^{\text{ext}}(i)|$  is partitioned with respect to each of the communities the node  $i$  is connected to,  
 435 i.e.  $|n_c^{\text{ext}}(i)| = \sum_{c' \neq c} |n_c^{\text{ext}}(i|c')|$ , where  $n_c^{\text{ext}}(i|c')$  stands for the set neighbours of  $i \in c$  belonging to the communities  
 436  $c' \neq c$ . Similarly, we can now look at the degrees of the neighbours with respect to the community  $c$ , and we  
 437 write  $n_{c'}^{\text{ext}}(k|c)$  for the set of neighbours of  $k \in c'$  belonging to the community  $c$ . Therefore, we can write the total  
 438 number of external links of  $c$  as:

$$m^{\text{ext}} = \sum_{\substack{c' \neq c \\ k \in c'}} |n_{c'}^{\text{ext}}(k|c)|,$$

439 where the summatory runs for every element  $k \in c'$  and every community  $c' \neq c$ .

440 In addition, it is evident that a node  $k$  in an external community  $c'$  will contribute to the mean similarity of  
 441 the nodes in community  $c$  if and only if its degree with respect to that community is at least 2, i.e.  $n_{c'}^{\text{ext}}(k|c) \geq 2$ ,  
 442 and thus it is immediate to relate  $\hat{s}_c^{\text{ext}}$  and  $m^{\text{ext}}$ :

$$\begin{aligned}
\hat{s}_c^{\text{ext}} &= \sum_{i < j} |n_c^{\text{ext}}(i) \cap n_c^{\text{ext}}(j)| = \sum_{\substack{c' \neq c \\ k \in c'}} \frac{|n_{c'}^{\text{ext}}(k|c)|(|n_{c'}^{\text{ext}}(k|c)| - 1)}{2} = \\
&= \frac{1}{2} \left( \sum_{\substack{c' \neq c \\ k \in c'}} |n_{c'}^{\text{ext}}(k|c)|^2 - \sum_{\substack{c' \neq c \\ k \in c'}} |n_{c'}^{\text{ext}}(k|c)| \right) \\
&= \frac{1}{2} \left( \sum_{\substack{c' \neq c \\ k \in c'}} |n_{c'}^{\text{ext}}(k|c)|^2 - m_c^{\text{ext}} \right),
\end{aligned}$$

and therefore, if  $m_c^{\text{ext}}$  increases  $\hat{s}_c^{\text{ext}}$  will increase (note that the argument of  $m_c^{\text{ext}}$  is the same that the quadratic term, but the latter will dominate the sum). We should still relate  $\hat{s}_c^{\text{ext}}$  with the links in excess, which we defined as  $m^{\text{ext}} - n_c^{\text{ext}}$  because we require at least one link with each neighbor to identify it as such. Therefore, we subtract from each node one link, reducing its degree in one, and we again compute all the possible remaining pathways, that we subtract from  $\hat{s}_c$ :

$$\hat{s}_c^{\text{ext}} - \sum_{\substack{c' \neq c \\ k \in c'}} \frac{(|n_{c'}^{\text{ext}}(k|c)| - 1)(|n_{c'}^{\text{ext}}(k|c)| - 2)}{2} = \sum_{\substack{c' \neq c \\ k \in c'}} (|n_{c'}^{\text{ext}}(k|c)| - 1) = m_c^{\text{ext}} - n_c^{\text{ext}}.$$

Note that the subtracted term has the same arguments that we used for the computation of  $\hat{s}_c^{\text{ext}}$  and that it must be smaller than  $\hat{s}_c^{\text{ext}}$ . A similar reasoning can be followed for the internal partition density, although now we should note that all links contribute to the similarity because we followed a convention in which  $|n^{\text{int}}(i) \cap n^{\text{int}}(j)| = 1$  if  $i$  and  $j$  are linked, what leads to

$$\hat{s}_c^{\text{int}} - \sum_i \frac{|n^{\text{int}}(i)|(|n^{\text{int}}(i)| - 1)}{2} = m^{\text{int}}.$$

The main reason why an increase in the number of links in excess is not a sufficient condition for the similarity to increase is that, if the links in excess are of different types, they may not contribute to the similarity. This is a weakness of our approximation that would require reinspecting the types of connections between neighbours during the clustering procedure, and developing a procedure to solve ambiguities if the types are different. This is not required in the current implementation, in which the method is separated in two distinct problems, namely the computation of the similarities (where there is inspection of the neighbours and the types of connections) and the clustering, hence improving the computational efficiency. Since the clustering is performed using the  $N$  highest similarities (out of  $N(N - 1)/2$  possible similarities, being  $N$  the number of nodes in the network), our choice is designed in a way in which nodes joined in a community are mostly connected to the same neighbours through the same type of links (otherwise the similarities should be low and the nodes will not be joined), which we confirmed in the examples shown in the article text. Hence, a more accurate partition density will likely not change the qualitative behaviour of the current functions, which have the advantage of being computationally efficient.

In addition, the relative contribution of  $\hat{s}_c^{\text{ext}}$  and  $\hat{s}_c^{\text{int}}$  may not be strictly translated into the same relative contribution to the similarities because we are omitting in the above computation the normalization terms. For instance, it may happen that two nodes  $i$  and  $j$  are joined in a cluster and that, when an external node  $k$  joins the cluster, we have  $s^{\text{ext}}(i, k) > s^{\text{ext}}(j, k)$  and  $S^{\text{int}}(i, k) > S^{\text{int}}(j, k)$  (implying  $|n(i)| \ll |n(j)|$ ). This means that the neighbours shared by  $i$  and  $k$  and those shared between  $j$  and  $k$  substantially differ. However, since  $k$  has few neighbours and  $i$  and  $j$  were joined first, the inequality  $S(i, j) > \max(S(i, k), S(j, k))$  would be unlikely to be fulfilled ( $n(i)$  is large and  $|n(i) \cap n(j)|$  is small since a fraction of the few neighbors that  $j$  has are shared with  $k$ ). Therefore, again the hierarchical clustering helps avoid “pathological” situations, hence its importance in our approximation.

**474 Original definition of partition density and relation with functionInk**

475 For completeness, we present the definition of partition density presented in [32] and a comparison with the new  
 476 method. In short, Ahn *et al.* method starts building a similarity measure between any pair of links sharing  
 477 one node in common. Two links will be similar if the nodes that these two links do not share have, in turn,  
 478 similar relationships with any other node, shown in Fig. 1. From this similarity measure, links are clustered and  
 479 an optimal cut-off for the clustering is found monitoring a measure called *partition density* (which in this paper  
 480 relates to the *internal partition density*). The optimal classification found at the cut-off, determines groups of links  
 481 that are similar because they connect nodes that are themselves similar in terms of their connectivity. Therefore,  
 482 the nodes are classified indirectly, according to the groups that their respective links belong, and a node may not  
 483 belong to a single community but to several communities if its links belong to different clusters. This is claimed  
 484 to be an advantage with respect to other methods (in particular for high density networks) as membership to a  
 485 single cluster is not enforced. At every step of the clustering it is obtained a partition  $P = P_1, \dots, P_C$  of the links  
 486 into  $C$  subsets. For every subset, the number of links is  $m_c = |P_c|$  and the number of nodes that these links are  
 487 connecting is  $n_c = |\cup_{e_{ij} \in P_c} \{i, j\}|$ . The density of links for the cluster  $C$  is then

$$D_c = \frac{m_c - (n_c - 1)}{n_c(n_c - 1)/2 - (n_c - 1)} \quad (12)$$

488 where the normalization considers the minimum ( $n_c - 1$ ) and maximum ( $n_c(n_c - 1)/2$ ) number of links that  
 489 can be found in the partition. The difference with respect to Eq. 4, is that a term ( $n_c^{int} - 1$ ) is now subtracted.  
 490 The reason is that, in Ahn *et al.* method, clustering with links implies that two nodes in the same cluster must  
 491 share links. But, according to our definition of function, two nodes may be structurally equivalent even if there  
 492 is no interaction between them.

493 The partition density  $D$  is then given by the average of the density of links for all the partitions, weighted by  
 494 the number of links

$$D = \frac{2}{M} \sum_c m_c \frac{m_c - (n_c - 1)}{(n_c - 2)(n_c - 1)} \quad (13)$$

495 where  $M$  is the total number of links. It was shown that when using agglomerative clustering, this function  
 496 achieves a maximum which determines the optimal partition [32].

497 functionInk shifts the attention from links back to nodes (see Fig. 1A). There are a number of reasons justifying  
 498 this shift:

- 499 • Working with nodes is more natural, thus favoring the interpretation of the communities. From a biological  
 500 perspective, we are interested in the function of nodes. For instance, the definition of “niche” is straightfor-  
 501 ward for nodes, justifying both the similarity measure definition and the rational behind both the clustering  
 502 and the definition of partition densities we proposed. In particular, focusing on nodes allowed us to identify  
 503 both modules and guilds.
- 504 • Inspecting and interpreting results is more convenient working with node partitions. For instance, most of  
 505 the network visualization programs can easily separate communities defined over nodes partitions but not  
 506 over links partitions.
- 507 • Monitoring changes in communities defined over node partitions becomes easier (even more if there are  
 508 changes in the network such as addition/removal of nodes or links), because the number of links scales as  
 509  $N^2$  while the number of nodes scales  $N$  (being  $N$  the number of nodes), so it may be expected a more  
 510 profound change in the links partitioning when adding/removing nodes.

511 Nevertheless, there are some types of networks for which a sharp classification of nodes into communities may be  
 512 elusive, such as when communities overlap (see section 11 in [4]), a situation motivating the development of the  
 513 method of Ahn *et al.* [32]. To illustrate this point we perform an explicit comparison between the method of Ahn  
 514 *et al.* and functionInk with a specific example. The open question to compare both methods is how to determine  
 515 partitions over sets of nodes from partitions over sets of links. A possibility comes from the identification of sets  
 516 of nodes whose links belong to the same set of link partition(s), since this is the idea behind the development of  
 517 functionInk. Therefore, we should find similar communities (clusters in the nodes’ partition) if we use the method  
 518 of Ahn *et al.* to identify a partition of links and i) we then identify communities as those nodes sharing links  
 519 classified in the same clusters in the partition of links or ii) we find these communities directly with functionInk.

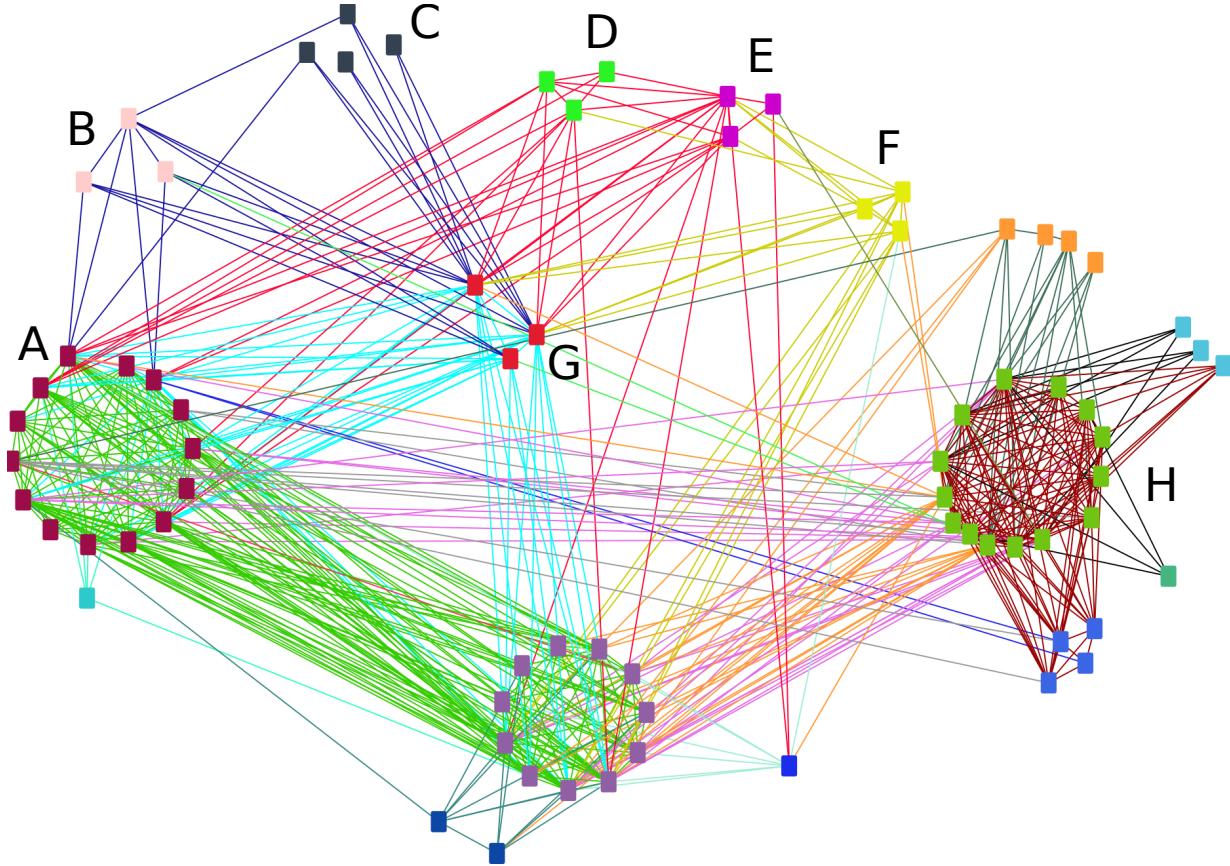


Figure 7: **Comparison between functionInk and Ahn *et al.* method [32].** Microbial network analyzed in Fig. 13 now comparing the partition over the set of nodes found with the total partition density with functionInk (leading to communities of nodes sharing the same color and close in space) and the partition over the set of links found with Ahn *et al.* method (leading to groups of links sharing the same color and connecting similar community nodes). The network is reduced in size and rearranged for the sake of a more clear comparison. Analysis of the communities labeled A-G is found in the text. Note that there is not full correspondence with the communities found in Fig. 13, since we do not use the information from the types of links here.

To explore this question, we re-analyze the microbial network shown in Fig. 13 with functionInk considering a single type of link and the communities found at the maximum of the total partition density, together with the partition in links determined with the method of Ahn *et al.*, using average linkage as a clustering algorithm in both cases. In Fig. 7 we show the network in which, for clarity, we removed communities with less than three members (although we kept some nodes with a high density of links).

Partitions defined by Ahn *et al.* method are identified by different colors for the links, and those defined by functionInk by different colors for the nodes (which are further located close in space). It is immediately apparent the good match between both methods, as expected since functionInk departs from the method of Ahn *et al.* For instance, the community labeled G can be defined as the one having cyan, olive, red and blue links. We also observe, however, some differences. Communities B and C are those having only blue links. functionInk splits these nodes in two communities because all three members of community B are linked to the same member in community G, but this is a connection that none of the C-members have, and similarly with respect to two members of community A (although, in this case, one member of C has one such connection). The same applies for communities D and E, which could be defined as those preferentially having red links, but community E is split from D because it has connections with other nodes, in particular with community F. Nevertheless, functionInk is not able to split community H into sub-communities, which is suggested for instance by subsets of nodes having links to several link partitions, e.g. pink and orange links. Even if these communities would be split stopping at an earlier step (for instance at the maximum of the external partition density) thanks to the flexibility provided

538 by the availability of several partition function definitions, it is also clear that functionInk may have difficulties  
539 with highly overlapping communities, a situation in which it may be useful to combine both methods.

## 540 Clustering algorithm

541 After computing the similarity between nodes with the method presented in the Results, the algorithm clusters  
542 nodes using one of three hierarchical clustering algorithms: average linkage [47], single linkage and complete  
543 linkage. Starting from each node being a separate cluster, at each step  $t$  all algorithms join the two most similar  
544 clusters  $A$  and  $B$ , and compute the similarity between the new combined cluster and all other clusters  $C$  in a way  
545 that depends on the clustering algorithm.

546 Single linkage is the most permissive algorithm, because the similarity it assigns to the new cluster is the  
547 maximum similarity between the two clusters joined and clusters  $C$ :

$$S^{t+1}(AB, C) = \max(S(A, C), S(B, C)).$$

548 where  $t$  labels the step of the algorithm,  $A$  and  $B$  are the clusters that are joined,  $AB$  denotes the new  
549 composite cluster, and  $C$  is any other cluster. On the other hand, complete linkage is the most restrictive,  
550 assigning the minimum similarity

$$S^{t+1}(AB, C) = \min(S(A, C), S(B, C)).$$

551 Finally, average linkage assigns an intermediate value computed as the weighted average similarity with the  
552 two joined clusters

$$S^{t+1}(AB, C) = \frac{n_A S^t(A, C) + n_B S^t(B, C)}{n_A + n_B}$$

553 being  $n_A$  and  $n_B$  the number of elements that  $A$  and  $B$  contain, respectively. Identification of the two pools  
554 of plants and animals is independent of the clustering method used, but the maximum of the external partition  
555 density is achieved earlier for single linkage and later for complete linkage; we found a good compromise between  
556 the number and the size of the clusters working with average linkage, in agreement with previously reported  
557 results [31], but the clustering method could be selected according to information known from the links. In our  
558 experience, single linkage is easily dominated by the giant cluster in high density networks in which modules are  
559 prevalent (rather than for guilds). The appropriate clustering method should be guided by the research question.  
560 For instance, if gene homology is explored, it is probably more appropriate to use single linkage (as a relative of  
561 one gene's relative is also its relative, i.e. transitivity is automatically fulfilled [36]). On the other hand, if we  
562 analyse well-differentiated functional similarity, it might be more appropriate to be conservative and use complete  
563 linkage.

## 564 Plant-pollinator networks and topological properties

565 We selected six plant-pollinator networks artificially generated in [39] with known topological properties, sum-  
566 marized in Table. We consider as topological properties the connectance (fraction of links) of the mutualistic  
567 matrix, the connectance of the competition matrices, and the definition of nestedness provided in [40]. Given a  
568 mutualistic matrix  $A_{ik}^{(P)}$  representing presence-absence of interaction between the set of plants, indexed by  $i$ , and  
569 the set of animal species, indexed by  $k$ , we compute the degree of a species as  $n_i^{(P)} = \sum_k A_{ik}^{(P)}$  (see Ref. [40]  
570 in Supplementary Material). A similar definition would apply for animals  $n_k^{(A)} = \sum_i A_{ik}^{(P)}$ . Next we define the  
571 ecological overlap between two species of plants  $i$  and  $j$  as the number of insects that pollinate both plants:

$$n_{ij}^{(P)} = \sum_{k \in A} A_{ik}^{(P)} A_{jk}^{(P)}, \quad (14)$$

572 a definition that is equivalent to the numerator Jaccard similarity used in this work. Summing over every pair  
573 of plants and normalizing leads to the definition of nestedness:

$$\nu^{(P)} = \frac{\sum_{i < j} n_{ij}^{(P)}}{\sum_{i < j} \min(n_i^{(P)}, n_j^{(P)})}. \quad (15)$$

$\nu$	$\kappa_{\text{mut}}$	$\kappa_{\text{comp}}$
0.15	0.08	0
0.35	0.16	0.15
0.6	0.28	0.5

Table 1: **Topological properties of the bipartite networks analyzed.** Different combinations of nestedness ( $\nu$ ), intra-pools connectance  $\kappa_{\text{comp}}$  and inter-pools connectance  $\kappa_{\text{mut}}$  were analyzed.

A symmetric definition applies for animals, so we take as final definition of nestedness  $\nu = \max(\nu^{(P)}, \nu^{(A)})$ . Note that in Eq. 14,  $k$  indexes animal species and, since the two pools of plants and animals are separated until the very last steps, changes in the nestedness will have an effect only on the external partition density.

## 577 Trophic networks

We downloaded the network and metadata provided in [41] and compared the clusters found with those obtained by functionInk. After computing the Tanimoto coefficients as explained above, we cluster the nodes and retrieve the classification found at each step. We then computed five indexes (Rand, Fowlkes and Mallows, Wallace 10, Wallace 01 and Jaccard), implemented in the R PCI function of the PROFDP package [42]. In order to assign a significance value for the different indexes we obtained, for each index  $x$ , a bootstrapped distribution with mean  $\bar{x}_{(B)}$  and standard deviation  $\sigma_{(B)}$ , re-sampling with replacement the samples and recomputing the indexes  $10^3$  times. Next we computed  $10^3$  completely random classifications, obtained by shuffling the identifiers relating each sample with one of the classifications, and retrieving the maximum  $x_{(R)}$ . We finally verified that the random value was significantly different from the bootstrapped distribution by computing the z-scores:

$$z = \frac{\text{abs}(x_{(R)} - \bar{x}_{(B)})}{\sigma_{(B)}},$$

which we considered significant if it was higher than 2.5. Heatmaps were generated with the HEATMAP.2 function in R package GPLOTS.

## 590 Bacterial networks

We considered a public dataset of 753 bacterial communities sampled from rainwater-filled beech tree-holes (*Fagus* spp.) [44], leading to 2874 Operative Taxonomic Units (OTUs) at the 97% of 16 rRNA sequence similarity. These communities were compared with Jensen-Shannon divergence [45], and automatically clustered following the method proposed in Ref. [48] to identify enterotypes. The clusters found with this method in [43] were used to color the community labels in Fig. 13.

The inference of the OTU network started quantifying correlations between OTUs abundances with SparCC [46]. To perform this computation, from the original OTUs we reduced the data set removing rare taxa with less than 100 reads or occurring in less than 10 samples, leading to 619 OTUs. Then, the significance of the correlations was evaluated bootstrapping the samples 100 times the data and estimating pseudo p-values for each of the  $N(N - 1)/2$  pairs. A relationship between two OTUs was considered significant and represented as a link in the network if the correlation was larger than 0.2 in absolute value and the pseudo p-value lower than 0.01. The network obtained in this way was analyzed using Cytoscape [49].

In Suppl. Fig. 13 we show the network obtained with the partitions found at the maximum of the external, total and internal partition densities. The combination of the three partitions allow us to individuate different key players in the network, and its relation with the definition of structural equivalence we adopted. For instance, only one OTU from the green functional group in the partitioning generated with the total partition density, have an important number of co-occurrences with members of the red functional group. In addition, it is the only one having a significant segregation with respect to a highly abundant member of the blue functional groups. These two highly abundant segregating OTUs are *Pseudomonas putida* (green functional group) and *Serratia fonticola* (blue functional group), both of which were shown to dominate two of the  $\beta$ -diversity-classes [43]. In the partitioning found with the external partition density, *Pseudomonas putida* appears as an orphan node and the functional group of *Serratia fonticola* is split in three. This example illustrates how the different stopping criteria can be used to individuate larger or smaller differences in the connectivity.

614 **Supplementary Figures**

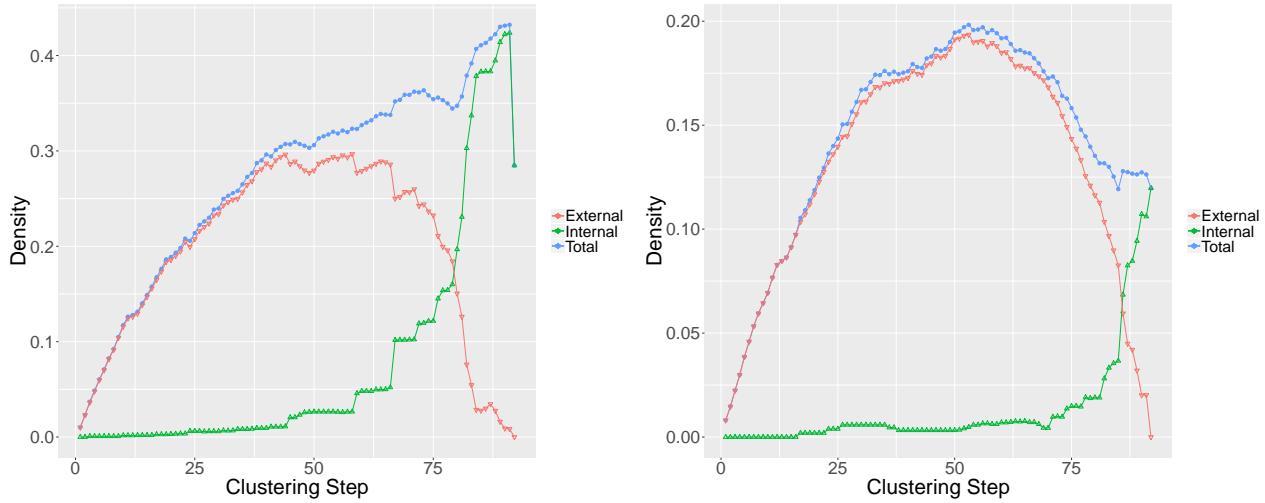


Figure 8: **Partition densities of synthetic mutualistic networks.** Networks with nestedness  $\nu = 0.15$ ,  $\kappa_{\text{mut}} = 0.08$ , and  $\kappa_{\text{comp}} = 0.5$  (left) or  $\kappa_{\text{comp}} = 0.15$  (right). Changing the connectance change the relative value between the external and internal partition densities.

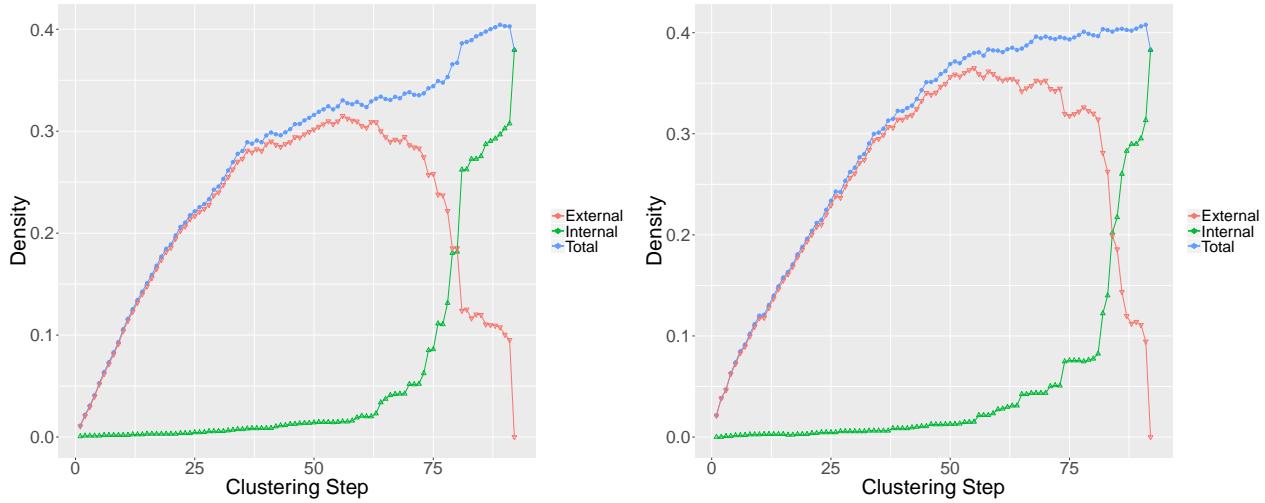


Figure 9: **Partition densities for synthetic mutualistic networks.** Networks with  $\kappa_{\text{comp}} = 0.5$ ,  $\kappa_{\text{mut}} = 0.28$  and  $\nu = 0.35$  (left) or  $\nu = 0.6$  (right). The high connectance of both networks make the internal partition density dominant, and two pools are detected through the total partition density. Nevertheless, the increase of the nestedness is detected through an increase in the external partition density, which makes the second network more disassortative.

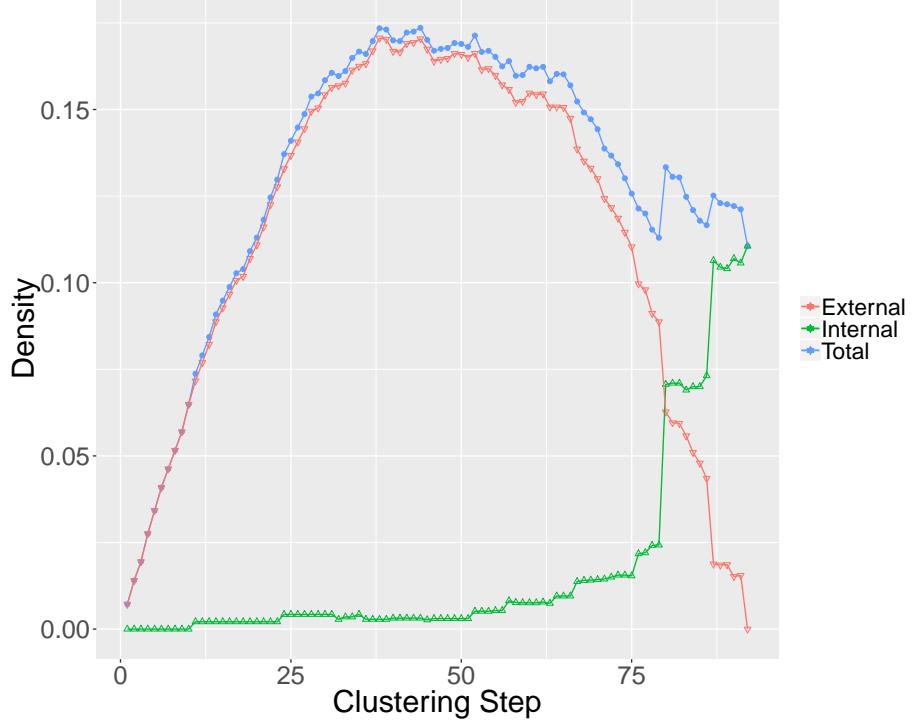


Figure 10: **Partition densities of synthetic mutualistic networks.** Network with nestedness  $\nu = 0.05$ ,  $\kappa_{\text{mut}} = 0.065$  and  $\kappa_{\text{comp}} = 0.15$ . The low connectance hinders the detection of the two pools of plants and pollinators.

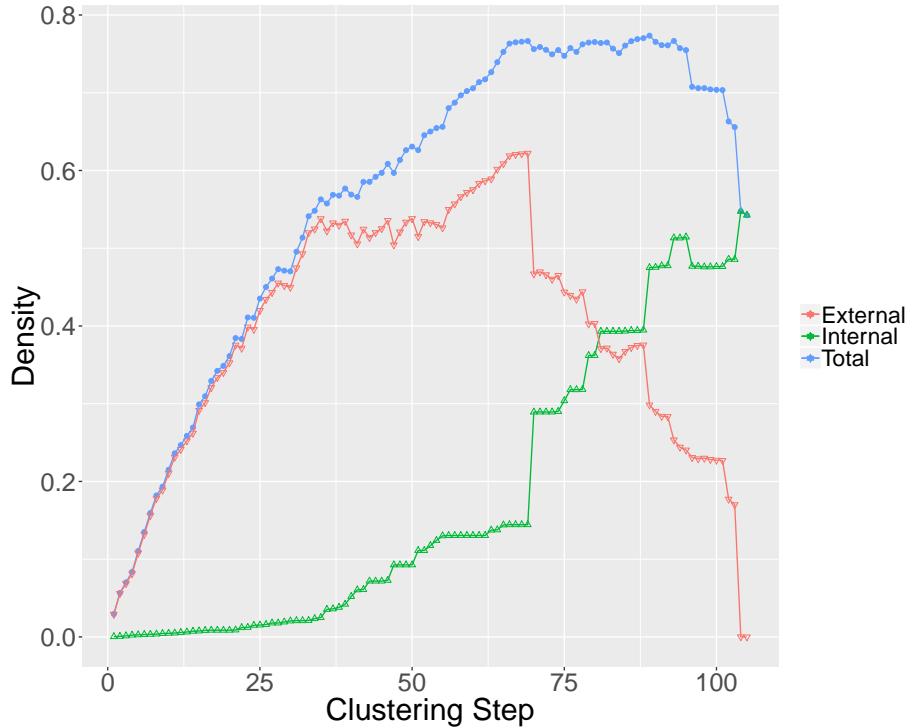


Figure 11: **Partition density of the trophic network.** The internal partition density peaks when there are three clusters, consistent with the existence of three trophic layers. The external partition density has a maximum at step 69, which is analyzed in detail with respect to the reference classification found in [41].

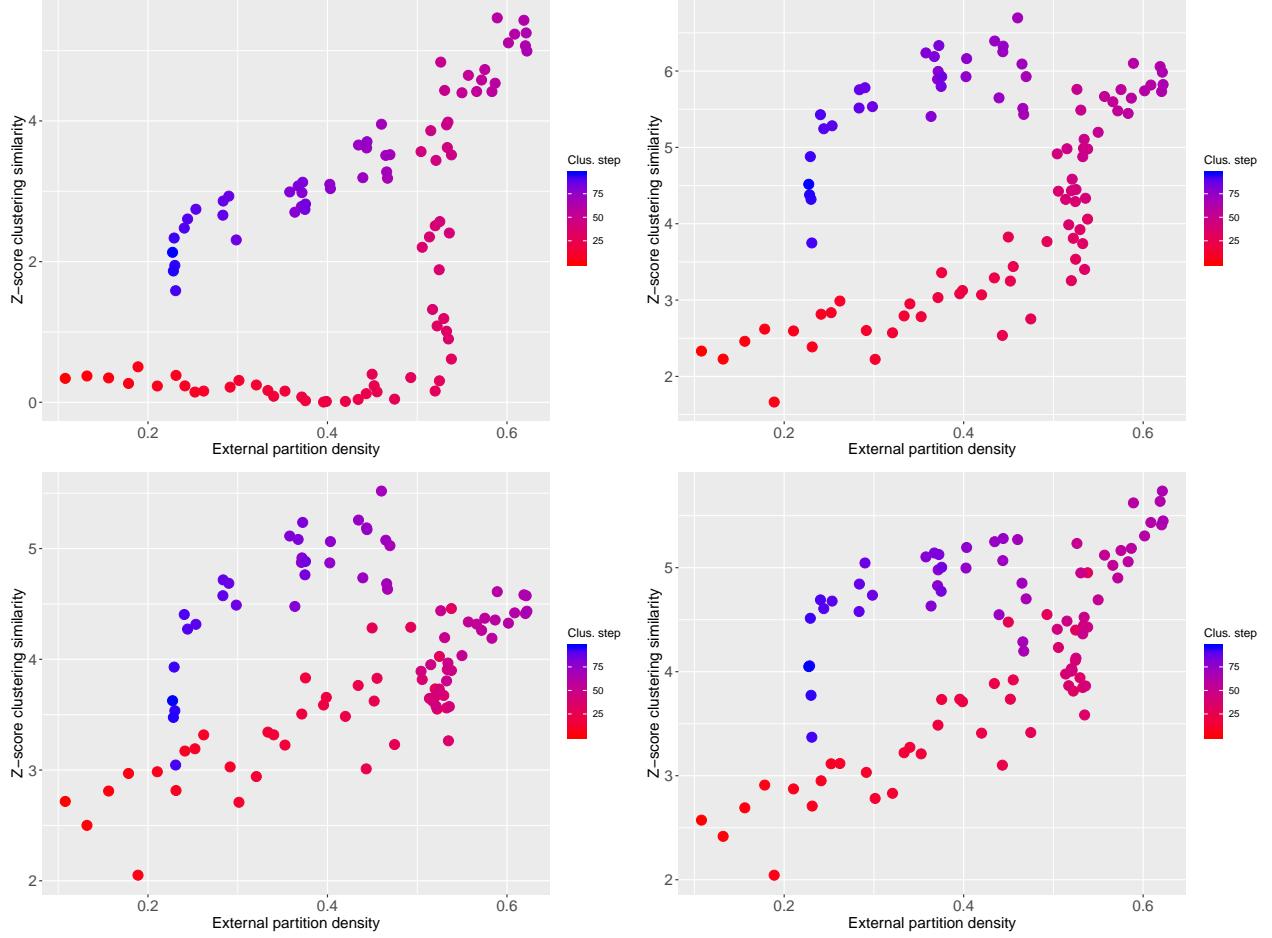
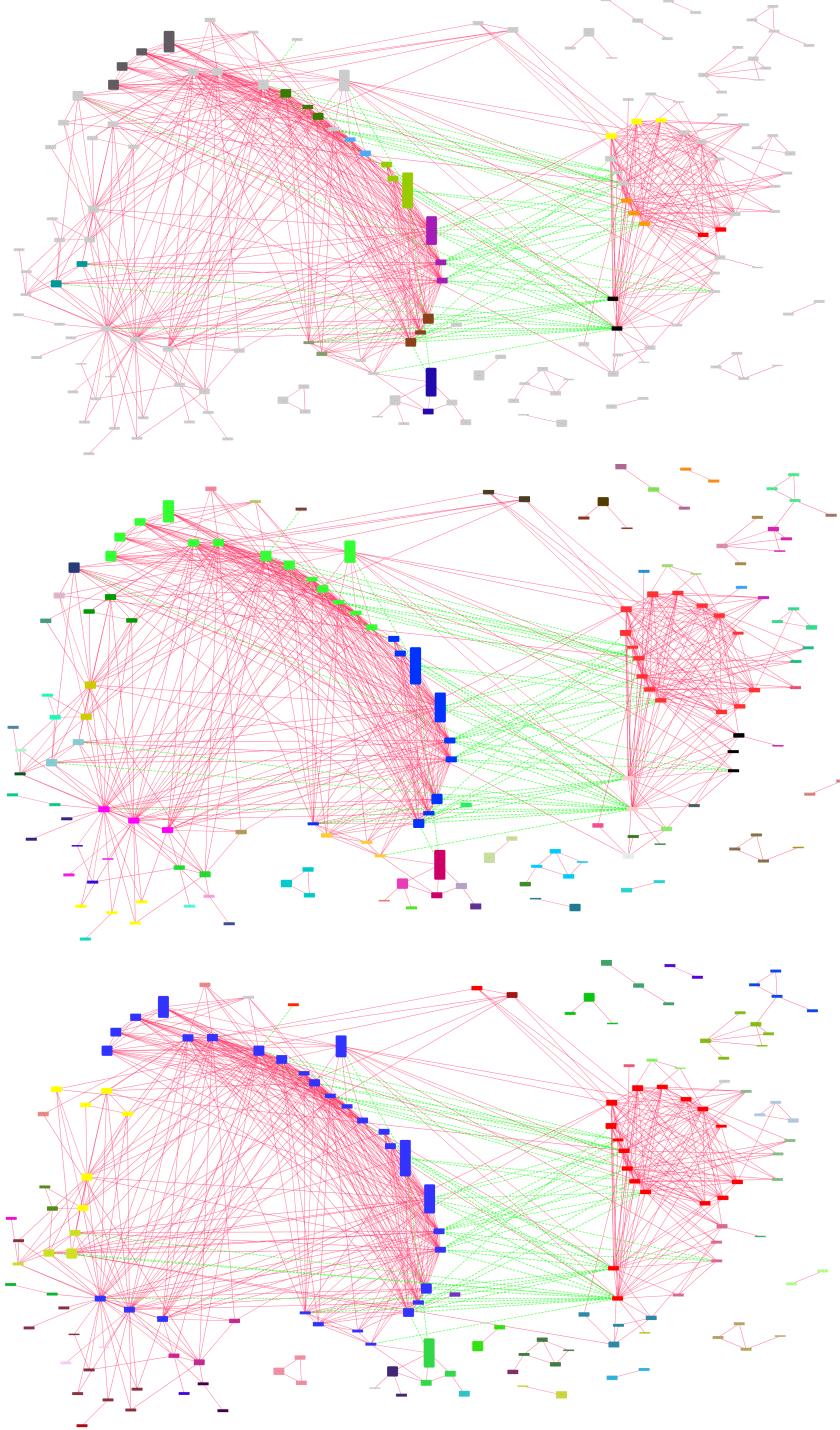


Figure 12: **Comparison between classifications of the trophic network.** Similarity between the reference classification found in [41] and the one found with functionInk is performed with the Z-score of a different indexes: Wallace 01 (Top left), Fowlkes and Mallows (Top right), Jaccard (Bottom left) and Rand (Bottom right). All indexes bring significant values and the maximum similarity is close to the maximum of the external partition density.



**Figure 13: Comparison of functional groups in the microbial network.** Network of significant co-occurrences (continuous links) and segregations (dotted links) at the species level (nodes). Colors indicate functional group membership, which was determined by the maximum of the external (top), total (middle) and internal partition densities (bottom). Orphan nodes are colored gray in the top figure for clarity. The higher value of the internal partition density (see Suppl. Fig. 14) suggests that a modular structure is the more appropriate to describe the functional groups. This is confirmed by the low number of guilds (top figure) and the good agreement between the global topological structure and the modules (bottom figure). Communities were automatically located close in space according to the partition found with the total partition density (middle) and blue and green communities rearranged manually to make more clear their connections, in particular we separated one node in the green community being the only one with co-occurrences with other communities on the right-hand side.

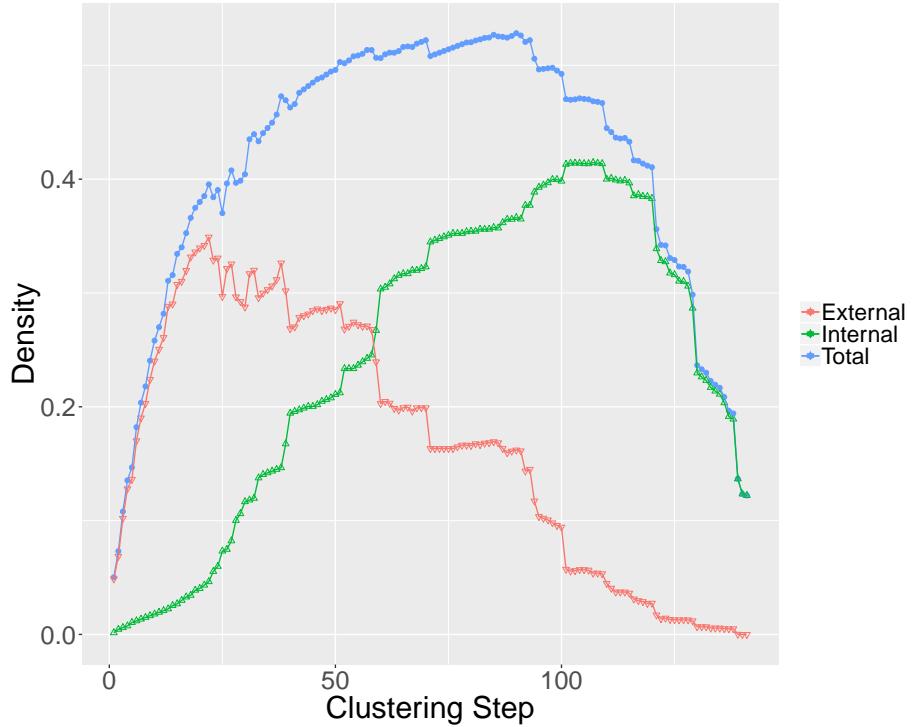


Figure 14: **Partition density of the microbial network.** The external partition density brings a poor reduction in the complexity of the network, with only 22 elements joined, while the internal partition density achieves a higher value and still a good number of clusters. Results suggest that modules are more relevant in this network given the high number of intra-cluster co-occurrences, later confirmed by visual inspection in the Main Text.

## 616 References

- 617 [1] J. E. Cohen and D. W. Stephens, *Food webs and niche space*. No. 11, Princeton University Press, 1978.
- 618 [2] R. MacArthur, "Fluctuations of animal populations and a measure of community stability," *Ecology*, vol. 36,  
619 p. 533, July 1955.
- 620 [3] M. Kivelä, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter, "Multilayer networks,"  
621 *Journal of Complex Networks*, vol. 2, no. 3, pp. 203–271, 2014.
- 622 [4] S. Fortunato, "Community detection in graphs," *Physics Reports*, vol. 486, no. 3–5, pp. 75–174, 2010.
- 623 [5] S. Allesina and M. Pascual, "Food web models: a plea for groups," *Ecology Letters*, vol. 12, no. 7, pp. 652–662,  
624 2009.
- 625 [6] S. Pilosof, M. A. Porter, M. Pascual, and S. Kéfi, "The multilayer nature of ecological networks," *Nature  
626 Ecology & Evolution*, vol. 1, no. 4, p. 0101, 2017.
- 627 [7] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang, "Complex networks: Structure and  
628 dynamics," *Physics reports*, vol. 424, no. 4, pp. 175–308, 2006.
- 629 [8] A. Carr, C. Diener, N. S. Baliga, and S. M. Gibbons, "Use and abuse of correlation analyses in microbial  
630 ecology," *The ISME journal*, p. 1, 2019.
- 631 [9] A. Pascual-García, J. Tamames, and U. Bastolla, "Bacteria dialog with santa rosalia: Are aggregations of cos-  
632 mopolitan bacteria mainly explained by habitat filtering or by ecological interactions?," *BMC Microbiology*,  
633 vol. 14, no. 1, p. 284, 2014.
- 634 [10] C. S. Elton, *Animal ecology*. New York: The Macmillan Company, 1927.
- 635 [11] D. Simberloff and T. Dayan, "The guild concept and the structure of ecological communities," *Annual Review  
636 of Ecology and Systematics*, vol. 22, no. 1, pp. 115–143, 1991.
- 637 [12] M. E. Newman and E. A. Leicht, "Mixture models and exploratory analysis in networks," *Proceedings of the  
638 National Academy of Sciences*, vol. 104, no. 23, pp. 9564–9569, 2007.
- 639 [13] M. E. Newman, "Finding community structure in networks using the eigenvectors of matrices," *Physical  
640 Review E*, vol. 74, no. 3, p. 036104, 2006.
- 641 [14] E. Estrada and J. A. Rodríguez-Velázquez, "Spectral measures of bipartivity in complex networks," *Physical  
642 Review E*, vol. 72, no. 4, p. 046105, 2005.
- 643 [15] M. T. Schaub, J.-C. Delvenne, M. Rosvall, and R. Lambiotte, "The many facets of community detection in  
644 complex networks," *Applied Network Science*, vol. 2, no. 1, p. 4, 2017.
- 645 [16] L. Peel, D. B. Larremore, and A. Clauset, "The ground truth about metadata and community detection in  
646 networks," *Science advances*, vol. 3, no. 5, p. e1602548, 2017.
- 647 [17] M. Girvan and M. E. Newman, "Community structure in social and biological networks," *Proceedings of the  
648 National Academy of Sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.
- 649 [18] R. Lambiotte, J.-C. Delvenne, and M. Barahona, "Laplacian dynamics and multiscale modular structure in  
650 networks," *arXiv preprint arXiv:0812.1770*, 2008.
- 651 [19] P. J. Mucha, T. Richardson, K. Macon, M. A. Porter, and J.-P. Onnela, "Community structure in time-  
652 dependent, multiscale, and multiplex networks," *Science*, vol. 328, no. 5980, pp. 876–878, 2010.
- 653 [20] M. De Domenico, A. Lancichinetti, A. Arenas, and M. Rosvall, "Identifying modular flows on multilayer  
654 networks reveals highly overlapping organization in interconnected systems," *Physical Review X*, vol. 5,  
655 no. 1, p. 011027, 2015.

- 656 [21] M. Bazzi, M. A. Porter, S. Williams, M. McDonald, D. J. Fenn, and S. D. Howison, "Community detection in temporal multilayer networks, with an application to correlation networks," *Multiscale Modeling & Simulation*, vol. 14, no. 1, pp. 1–41, 2016.
- 657  
658
- 659 [22] S. Wasserman and K. Faust, *Social network analysis: Methods and applications*, vol. 8. Cambridge university press, 1994.
- 660  
661
- 662 [23] P. W. Holland, K. B. Laskey, and S. Leinhardt, "Stochastic blockmodels: First steps," *Social networks*, vol. 5, no. 2, pp. 109–137, 1983.
- 663  
664
- 665 [24] C. De Bacco, E. A. Power, D. B. Larremore, and C. Moore, "Community detection, link prediction, and layer interdependence in multilayer networks," *Physical Review E*, vol. 95, no. 4, p. 042317, 2017.
- 666  
667
- 668 [25] M. E. Newman, "Equivalence between modularity optimization and maximum likelihood methods for community detection," *Physical Review E*, vol. 94, no. 5, p. 052315, 2016.
- 669  
670
- 671 [26] B. Karrer and M. E. Newman, "Stochastic blockmodels and community structure in networks," *Physical Review E*, vol. 83, no. 1, p. 016107, 2011.
- 672  
673
- 674 [27] T. Valles-Catala, F. A. Massucci, R. Guimera, and M. Sales-Pardo, "Multilayer stochastic block models reveal the multilayer structure of complex networks," *Physical Review X*, vol. 6, no. 1, p. 011036, 2016.
- 675  
676
- 677 [28] M. Ganji, J. Chan, P. J. Stuckey, J. Bailey, C. Leckie, K. Ramamohanarao, and I. Davidson, "Image constrained blockmodelling: a constraint programming approach," in *Proceedings of the 2018 SIAM International Conference on Data Mining*, pp. 19–27, SIAM, 2018.
- 678  
679
- 680 [29] S. P. Borgatti, M. G. Everett, and P. R. Shirey, "LS sets, lambda sets and other cohesive subsets," *Social networks*, vol. 12, no. 4, pp. 337–357, 1990.
- 681  
682
- 683 [30] J. J. Luczkovich, S. P. Borgatti, J. C. Johnson, and M. G. Everett, "Defining and measuring trophic role similarity in food webs using regular equivalence," *Journal of Theoretical Biology*, vol. 220, no. 3, pp. 303–321, 2003.
- 684  
685
- 686 [31] P. Yodzis and K. O. Winemiller, "In search of operational trophospecies in a tropical aquatic food web," *Oikos*, pp. 327–340, 1999.
- 687  
688
- 689 [32] Y.-Y. Ahn, J. P. Bagrow, and S. Lehmann, "Link communities reveal multiscale complexity in networks," *Nature*, vol. 466, no. 7307, pp. 761–764, 2010.
- 690  
691
- 692 [33] R. Guimera and L. A. N. Amaral, "Cartography of complex networks: modules and universal roles," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2005, no. 02, p. P02001, 2005.
- 693  
694
- 695 [34] S. P. Borgatti and M. G. Everett, "The class of all regular equivalences: Algebraic structure and computation," *Social networks*, vol. 11, no. 1, pp. 65–88, 1989.
- 696  
697
- 698 [35] T. T. Tanimoto, "elementary mathematical theory of classification and prediction," 1958.
- 699  
700
- 701 [36] A. Pascual-García, D. Abia, Á. R. Ortiz, and U. Bastolla, "Cross-over between discrete and continuous protein structure space: insights into automatic classification and networks of protein structures," *PLoS Computational Biology*, vol. 5, no. 3, p. e1000331, 2009.
- 702  
703
- 704 [37] A. Harrer and A. Schmidt, "Blockmodelling and role analysis in multi-relational networks," *Social Network Analysis and Mining*, vol. 3, no. 3, pp. 701–719, 2013.
- 705  
706
- 707 [38] M. G. Everett and S. P. Borgatti, "Exact colorations of graphs and digraphs," *Social networks*, vol. 18, no. 4, pp. 319–331, 1996.
- 708  
709
- 710 [39] A. Pascual-García and U. Bastolla, "Mutualism supports biodiversity when the direct competition is weak," *Nature Communications*, vol. 8, p. 14326, Feb. 2017.
- 711  
712
- 713 [40] U. Bastolla, M. A. Fortuna, A. Pascual-García, A. Ferrera, B. Luque, and J. Bascompte, "The architecture of mutualistic networks minimizes competition and increases biodiversity," *Nature*, vol. 458, pp. 1018–1020, Apr. 2009.

- 700 [41] S. Kéfi, V. Miele, E. A. Wieters, S. A. Navarrete, and E. L. Berlow, "How structured is the entangled  
701 bank? the surprisingly simple organization of multiplex ecological networks leads to increased persistence  
702 and resilience," *PLoS Biology*, vol. 14, no. 8, p. e1002527, 2016.
- 703 [42] M. S. Shotwell *et al.*, "profpm: An r package for map estimation in a class of conjugate product partition  
704 models," *J Stat Softw*, vol. 53, no. 8, pp. 1–18, 2013.
- 705 [43] A. Pascual-García and T. Bell, "Community-level signatures of ecological succession in natural bacterial  
706 communities," *bioRxiv*, p. 636233, 2019.
- 707 [44] D. W. Rivett and T. Bell, "Abundance determines the functional role of bacterial phylotypes in complex  
708 communities," *Nature microbiology*, p. 1, 2018.
- 709 [45] D. M. Endres and J. E. Schindelin, "A new metric for probability distributions," *IEEE Transactions on  
710 Information Theory*, vol. 49, pp. 1858–1860, July 2003.
- 711 [46] J. Friedman and E. J. Alm, "Inferring correlation networks from genomic survey data," *PLoS Computational  
712 Biology*, vol. 8, no. 9, p. e1002687, 2012.
- 713 [47] R. R. Sokal, "A statistical method for evaluating systematic relationships," *Univ Kans Sci Bull*, vol. 38,  
714 pp. 1409–1438, 1958.
- 715 [48] M. Arumugam, J. Raes, E. Pelletier, D. Le Paslier, T. Yamada, D. R. Mende, G. R. Fernandes, J. Tap,  
716 T. Bruls, J.-M. Batto, *et al.*, "Enterotypes of the human gut microbiome," *Nature*, vol. 473, no. 7346,  
717 pp. 174–180, 2011.
- 718 [49] P. Shannon, A. Markiel, O. Ozier, N. S. Baliga, J. T. Wang, D. Ramage, N. Amin, B. Schwikowski, and  
719 T. Ideker, "Cytoscape: a software environment for integrated models of biomolecular interaction networks,"  
720 *Genome research*, vol. 13, no. 11, pp. 2498–2504, 2003.