

# Inferring microbial co-occurrence networks from amplicon data: a systematic evaluation

Dileep Kishore<sup>1,2</sup>, Gabriel Birzu<sup>3,6</sup>, Zhenjun Hu<sup>1</sup>, Charles DeLisi<sup>1,3</sup>, Kirill S. Korolev<sup>†1,2,3</sup>,  
Daniel Segre<sup>†1,2,4,5</sup>

<sup>1</sup>Bioinformatics Program, Boston University, Boston, Massachusetts, USA

<sup>2</sup>Biological Design Center, Boston University, Boston, Massachusetts, USA

<sup>3</sup>Department of Physics, Boston University, Boston, Massachusetts, USA

<sup>4</sup>Department of Biomedical Engineering, Boston University, Boston, Massachusetts, USA

<sup>5</sup>Department of Biology, Boston University, Boston, Massachusetts, USA

<sup>6</sup>Department of Applied Physics, Stanford University, Stanford, California, USA

†Correspondence should be sent to korolev@bu.edu or dsegre@bu.edu

1

## Abstract

2 *Microbes tend to organize into communities consisting of hundreds of species involved in complex  
3 interactions with each other. 16S ribosomal RNA (16S rRNA) amplicon profiling provides  
4 snapshots that reveal the phylogenies and abundance profiles of these microbial communities.  
5 These snapshots, when collected from multiple samples, have the potential to reveal which  
6 microbes co-occur, providing a glimpse into the network of associations in these communities.  
7 The inference of networks from 16S data is prone to statistical artifacts. There are many tools  
8 for performing each step of the 16S analysis workflow, but the extent to which these steps affect  
9 the final network is still unclear. In this study, we perform a meticulous analysis of each step  
10 of a pipeline that can convert 16S sequencing data into a network of microbial associations.  
11 Through this process, we map how different choices of algorithms and parameters affect the  
12 co-occurrence network and estimate steps that contribute most significantly to the variance.  
13 We further determine the tools and parameters that generate the most accurate and robust  
14 co-occurrence networks based on comparison with mock and synthetic datasets. Ultimately,  
15 we develop a standardized pipeline (available at <https://github.com/segralab/MiCoNE>) that  
16 follows these default tools and parameters, but that can also help explore the outcome of any  
17 other combination of choices. We envisage that this pipeline could be used for integrating  
18 multiple data-sets, and for generating comparative analyses and consensus networks that can  
19 help understand and control microbial community assembly in different biomes.*

20 **Keywords**— Microbiome, 16S rRNA, Pipeline, Interaction, Denoising, Taxonomy, Network  
21 Inference, Correlations, Qiime, Co-occurrence, Networks

## Importance

22 To understand and control the mechanisms that determine the structure and function of microbial  
23 communities, it is important to map the interrelationships between its constituent microbial species.  
24 The surge in the high-throughput sequencing of microbial communities has led to the creation of  
25 thousands of datasets containing information about microbial abundances. These abundances can be  
26 transformed into networks of co-occurrences across multiple samples, providing a glimpse into the  
27 structure of microbiomes. However, processing these datasets to obtain co-occurrence information  
28 relies on several complex steps, each of which involves multiple choices of tools and corresponding  
29 parameters. These multiple options pose questions about the accuracy and uniqueness of the inferred  
30 networks. In this study, we address this workflow and provide a systematic analysis of how these  
31 choices of tools and parameters affect the final network, and on how to select those that are most  
32 appropriate for a particular dataset.

---

## 34      **Introduction**

35      Microbial communities are ubiquitous and play an important role in marine and terrestrial environments, urban ecosystems, metabolic engineering, and human health [1, 2]. These microbial communities, or microbiomes, often comprise several hundreds of different microbial strains interacting with each other and their environment, often through intricate metabolic and signaling relationships. Understanding how these interconnections shape community structure and functionalities is a fundamental challenge in microbial ecology, with applications in the study of microbial ecosystems across different biomes. With the advancement in DNA sequencing technologies [3] and data processing methods, more information can be extracted from these microbial community samples than ever before. In particular, high-throughput sequencing, including community metagenomic sequencing and sequencing of 16S rRNA gene amplicons, has the potential to help detect, identify and quantify a large portion of the constitutive microorganisms of a microbiome [4, 5]. These advances have led to large-scale data collection efforts involving environmental (Earth Microbiome Project) [2], marine (Tara Oceans Project) [6] and human-associated microbiota (Human Microbiome Project) [7].

49      This wealth of information on the composition and functions of a community at different times and under different environmental conditions has the potential to help us understand how 50     communities assemble and operate. A powerful tool for translating microbiome data into knowledge 51     is the construction of possible inter-dependence networks across species. The importance of these 52     networks of relationships is two fold: first, such networks can serve as maps that help identify hubs of 53     keystone species [8, 9], or basic microbiome changes that occur as a consequence of environmental 54     perturbations or underlying host conditions [10]; second, networks of inter-dependencies can serve as 55     a key bridge towards building mechanistic models of microbial communities, greatly enhancing our 56

---

57 capacity to understand and control them. For example, multiple studies have shown the importance  
58 of specific microbial interactions in the healthy microbiome [5] and others have shown how changes  
59 in these interactions can lead to dysbiosis [11, 10, 12]. In the context of terrestrial bio-geochemistry,  
60 co-occurrence networks have been proposed as a valuable approach towards reconstructing the  
61 processes leading to microbiome assembly [13], and understanding the response of microbial  
62 communities to environmental perturbations [14].

63 Direct high-throughput measurement of interactions, e.g. through co-culture micro-droplet  
64 experiments [15, 16], or spatial visualization of natural communities [17] is possible, but it requires  
65 specific technological capabilities, and has yet to be extensively used. In parallel, sequencing data  
66 across multiple samples can be used for estimating co-occurrence relationships between taxa. While  
67 the relationship between directly measured interactions and statistically inferred co-occurrence is  
68 still poorly understood [18], a significant amount of effort has gone into estimating correlations from  
69 large microbiome sequence datasets. Co-occurrence networks have microbial taxa as nodes, and  
70 edges that represent the frequent co-occurrence (or negative correlations) across different datasets.

71 One of the most frequently used avenues for inferring co-occurrence networks is the parsing and  
72 analysis of 16S sequencing data [9, 19]. A large number of software tools and pipelines have been  
73 developed to analyze 16S sequencing data, often focused on addressing the many known limitations  
74 of this methodology, including resolution, sequencing depth, compositional nature, sequencing  
75 errors and copy number variations [20, 21]. Popular methods for different phases of the analysis of  
76 16S data include tools for: (i) denoising and clustering sequencing reads [22, 23]; (ii) assigning  
77 taxonomy to the reads [24, 25]; (iii) processing and transforming the taxonomy count matrices  
78 [26]; and (iv) inferring the co-occurrence network [27, 28]. Different specific algorithms are often  
79 aggregated into popular platforms (like MG-RAST [29], Qiita [30]) and packages (such as QIIME  
80 [22]) that provide pipelines for 16S data analysis. The different methods and tools developed to solve

---

81 issues arising in 16S analysis can lead to vastly different inferences of community compositions and  
82 co-occurrence networks [31, 32], making it difficult to reliably compare networks across different  
83 publications and studies. This is partially due to the fact that existing platforms are typically focused  
84 on Operational Taxonomic Unit (OTU) generation and not on the effects of upstream statistical  
85 methods on the inferred co-occurrence networks. Furthermore, no organized framework currently  
86 exist to systematically analyze and compare existing components of the data analysis from amplicons  
87 to networks. More broadly, given the lack of comprehensive comparisons between directly observed  
88 microbial interactions (e.g. from co-culture experiments) and co-occurrence networks, there is no  
89 straightforward way to determine which set of tools or methods generate the most accurate networks.

90 In this study, we present a standardized 16S data analysis pipeline called Microbial Co-occurrence  
91 Network Explorer (MiCoNE) that produces robust and reproducible co-occurrence networks from  
92 community 16S sequence data, and allow users to interactively explore how the network would  
93 change upon using different alternative tools and parameters at each step. Our pipeline is coupled to  
94 an online integrative tool for the organization, visualization and analysis of inter-microbial networks.  
95 In addition to making this tool freely available, we implemented a systematic comparative analysis  
96 to determine which steps of the pipeline have the largest influence on the final network, and what  
97 choice seems to provide best agreement with the tested mock and synthetic datasets. We believe  
98 that these steps will ensure better reproducibility and easier comparison of co-occurrence networks  
99 across datasets. We expect that our tool will also be useful for benchmarking future alternative  
100 methods, and for ensuring a transparent evaluation of the possible biases introduced by the use of  
101 specific tools.

---

102 **Results**

103 **Microbial Co-occurrence Network Explorer (MiCoNE)**

104 We have developed MiCoNE, a flexible and modular pipeline for 16S amplicon sequencing rRNA  
105 data (hereafter mentioned simply as 16S data) analysis, that allows us to infer microbial co-occurrence  
106 networks. It incorporates various popular, publicly available tools as well as custom Python modules  
107 and scripts to facilitate inference of co-occurrence networks from 16S data (see Methods). Using  
108 MiCoNE one can obtain co-occurrence networks by applying to 16S data (or to already processed  
109 taxonomic count matrices) any combination of the available tools. The effects of changing any of  
110 the intermediate step can be monitored and evaluated in terms of its final network outcome, as well  
111 as on any of the intermediate metrics and data outputs. The MiCoNE pipeline workflow is shown in  
112 Figure 1. The different steps for going from 16S data to co-occurrence networks can be grouped  
113 into four major modules; (i) the denoising and clustering (DC) step, which handles denoising of the  
114 raw 16S sequencing data into representative sequences; (ii) the taxonomy assignment (TA) step  
115 that assigns taxonomic labels to the representative sequences; (iii) the OTU processing (OP) step  
116 that filters and transforms the taxonomy abundance table; and finally (iv) the network inferences  
117 (NI) step which infers the microbial co-occurrence network. Each process in the pipeline supports  
118 alternate tools for performing the same task (see Methods and Figure 1). A centralized configuration  
119 file contains all the specifications for what modules are used in the pipeline, and can be modified  
120 by the user to choose the desired set of tools. In what follows, we perform a systematic analysis  
121 of each step of the pipeline to estimate how much the final co-occurrence network depends on the  
122 possible choices at each step. We also evaluate a large number of tool combinations to determine a  
123 set of recommended default options for the pipeline and provide the users with a set of guidelines to

---

124 facilitate tool selection as appropriate for their data.

125 Our analysis involves two types of data: The first type consists of sets of 16S sequencing data  
126 from real communities sampled from human Stool and Oral microbiomes. The second are datasets  
127 synthetically or artificially created for the specific goal of helping evaluate computational analysis  
128 tools (see Methods). In particular, in order to objectively compare, to the extent possible, how well  
129 each step in MiCoNE best captures the underlying data, we use both mock data (labelled mock4,  
130 mock12 and mock16) from mockrobiota [33] as well as, synthetically generated reads from an  
131 Illumina read simulator called ART [34]. These mock datasets consist of fake sequencing reads  
132 generated from reads obtained from synthetic microbial isolates mixed in known proportions. They  
133 contain the expected compositions along with the reference sequences for the organisms in the  
134 mock community. The synthetic reads were simulated using three different taxonomy distribution  
135 profiles, namely soil and water microbiomes obtained Earth Microbiome Project (EMP) [2] and  
136 Stool microbiome that is used in our real community analysis [35]. Reference sequences were  
137 generated using National Center for Biotechnology Information (NCBI) and the Decard package [31]  
138 for these taxonomy profiles. Detailed information on the mock communities and the settings used to  
139 generate the synthetic data are provided in the Methods section.

140 **The choice of reference database has the biggest impact on inferred networks**

141 In order to analyze the effect of different statistical methods on the inferred co-occurrence networks,  
142 we generated co-occurrence networks using all possible combinations of methods and estimated  
143 the variability in the networks due to each choice (Figure 1). This analysis is performed while  
144 keeping the network inference algorithm (NI step) the same throughout the analysis. The effects  
145 of various steps on the final co-occurrence network is estimated by building a linear model of the  
146 edges of the network as a function of the various steps in the analysis pipeline (see Methods). Figure

---

147 2B, shows the fraction of total variation among the co-occurrence networks due to the first three  
148 steps of the pipeline. In other words, each point corresponds to a different combination of tools,  
149 and captures how much the final network is affected by such choice. The 16S reference database  
150 contributes the most (~ 25%) to variation in the networks. This is also reflected in the fact that  
151 the networks can be clearly separated based on the database used (Figure 2B). This indicates that  
152 the taxonomy assigned to the reference sequences drastically alters the co-occurrence network. In  
153 fact the variability induced by taxonomy assignment is much more significant than that due to the  
154 variability induced based on how the reference sequences themselves are identified (in the DC step).  
155 The grouping of the networks by taxonomy assignment into clusters (Figure 2B) seems to derive  
156 from the mislabelling of constitutive taxa that are present in high abundance in the community,  
157 which drastically alter the nodes and hence the underlying network topology. The residual variation  
158 (Figure 2A) can be seen as an artifact that arises when multiple steps are changed at the same time.  
159 Another interesting observation (elaborated in detail in the denoising and clustering section) is  
160 that the dissimilarity between the networks decreases when the low abundance OTUs are removed  
161 from the network. These results suggest that the most important criterion for accurate comparative  
162 analyses of co-occurrence networks is the taxonomy reference database.

163 **Denoising and clustering methods differ in their identification of less common**  
164 **reference sequences**

165 Denoising and clustering are commonly carried out to generate representative sequences from the  
166 raw 16S sequencing data and to obtain the OTU/Exact Sequence Variant (ESV) tables (counts of  
167 these representative sequences for each sample). In order to compare the OTU tables generated  
168 by different tools we processed the same 16S sequencing reads (healthy samples from a fecal  
169 microbiome transplant study [35]) using 5 different methods: open-reference clustering, closed-

---

170 reference clustering, denovo clustering, Divisive Amplicon Denoising Algorithm 2 (DADA2) [23]  
171 and Deblur [36]. The first three methods are from the Quantitative Insights Into Microbial Ecology  
172 1 (QIIME1) [22] package. We find that there is good agreement in the OTU/ESV tables when  
173 different combinations of methods are used to generate them (Supplementary Figure S1).

174 To compare the representative sequences generated by these methods we employ both the  
175 weighted [37] (Figure 3A) and unweighted UniFrac method [38] (Figure 3B). The weighted UniFrac  
176 distance metric takes into account the counts of the representative sequences, whereas the unweighted  
177 UniFrac distance metric does not and hence gives equal weights to each sequence. From Figure 3A  
178 one can see that the representative sequences generated by the different methods are similar to  
179 each other when weighted by their abundance. Figure 3B on the other hand shows an increase in  
180 dissimilarity between each pair of methods suggesting that the methods might differ in the treatment  
181 of sequences of low abundance. In order to verify this claim, for each of these methods we use the  
182 Greengenes (GG) taxonomy database to assign taxonomies to the representative sequences. We then  
183 correlate the abundances of matching taxonomies between a pair of DC methods (Figure S1A and B).  
184 The ESV tables generated by methods that perform denoising are very similar to each other ( $\sim 0.91$ )  
185 and the OTU tables generated by the clustering methods are very similar to each other ( $\sim 0.9$ ), but  
186 results of denoising and clustering are highly uncorrelated with each other ( $\sim 0.4$ ) (Figure S1C).

187 These comparisons only elucidate the pairwise similarity or dissimilarity of a pair of methods.  
188 In order to determine the tool that most accurately recapitulates the reference sequences in the  
189 samples, we used the 16S sequences from the mock datasets. In particular, we used the pipeline  
190 to process mock community datasets using each of the possible methods included for this step.  
191 We next compared predicted representative sequences with expected representative sequences and  
192 their distribution. The results (Figure 3C and D) show that, for the mock datasets, the different  
193 methods perform similar to each other, exactly as observed in the case of the real dataset. However,

---

194 the mock predicted sequence distributions are substantially different from the expected sequence  
195 distribution. This result is more exaggerated in the case of the unweighted UniFrac metric, where  
196 some of the datasets show a very high deviation from the expected sequences. These high deviations  
197 are primarily in two of the three datasets that were analyzed and show that the datasets themselves  
198 play a big role in the performance of these methods. This can be clearly seen in the performance  
199 (weighted UniFrac distance) of DADA2 and Deblur on mock12 and mock16 datasets, where, Deblur  
200 outperforms DADA2 on mock12 but the under-performs on mock16.

201 There is no method that clearly outperforms the rest in all datasets. Based on their slightly  
202 better performance on the mock datasets, their *de novo* error correcting nature and other previous  
203 studies [39], DADA2 and Deblur seem to be in general the most reliable. Given the unexpected  
204 poor performance of Deblur on the synthetic data, the default algorithm in the pipeline was chosen  
205 to be DADA2 (Supplementary Figure S3).

## 206 **Taxonomy databases vary widely in taxonomy hierarchy and update frequency**

207 Taxonomy databases are used to assign taxonomic identities to the representative sequences obtained  
208 after the DC step. In order to compare the assigned taxonomies from different databases, we use  
209 the same reference sequences and assign taxonomies to them using different taxonomy reference  
210 databases. The three 16S taxonomic reference databases used in this study are SILVA [25],  
211 GG [24] and NCBI RefSeq [40]. SILVA and GG are two popular 16S databases used for taxonomy  
212 identification. The NCBI RefSeq nucleotide database contains 16S rRNA sequences as a part of two  
213 BioProjects - 33175 and 33317. The three databases vastly differ in terms of their last update status -  
214 GG was last updated on May 2013, SILVA was last updated on December 2017 at the time of writing  
215 and NCBI is updated as new sequences are curated. Since updates to taxonomic classifications  
216 are frequent, these databases vary significantly in terms of taxonomy hierarchies including species

---

217 names and phylogenetic relationships [41].

218 The representative sequences obtained from the DADA2 method in DC step were used for  
219 taxonomic assignment using the three reference databases. Figure 4A depicts a flow diagram  
220 that shows how the top 50 representative sequences (sorted by abundance) are assigned a Genus  
221 according to the three different databases. We observe that not only does the assigned Genus  
222 composition vary significantly, but the percentage of unassigned representative sequences (gray)  
223 also differ. Even the most abundant representative sequence is assigned to an "unknown" Genus  
224 in two of the three databases. A representative sequence might be assigned an "unknown" Genus  
225 for one of two reasons: the first is if the taxonomy identifier associated with the sequence in the  
226 database did not contain a Genus; the second (more likely) reason is that the database contains  
227 multiple sequences that are very similar to the query (representative) sequence and the consensus  
228 algorithm (from Quantitative Insights Into Microbial Ecology 2 (QIIME2)) is unable to assign one  
229 particular Genus at the required confidence. After assigning all the representative sequences to  
230 taxonomies we perform a pairwise comparison of the similarity between assignments from different  
231 databases at every taxonomic level (Figure 4B). The assignments beyond Family level (Family,  
232 Genus and Species) are very dissimilar with < 70% similarity between any pair of databases. There  
233 are no two reference databases that are more similar than the other pairs, with GG and SILVA  
234 producing only marginally similar assignments compared to NCBI. This implies that the taxonomy  
235 assignments from each reference database are fairly unique and are largely responsible for the  
236 differences observed in the co-occurrence networks generated from different taxonomy databases.

237 Supplementary Figure S4 shows that the top 20 most abundant genera in the three resulting  
238 taxonomy composition tables are different. For example, the most abundant genus in the GG  
239 taxonomy table was *Escherichia* whereas in the SILVA taxonomy table it was *Escherichia-Shigella*.  
240 Although these are minor differences, when comparing a large number of taxonomy composition

---

241 tables these problems are hard to diagnose.

242 As in the previous section, these comparisons only indicate similarity or dissimilarity between  
243 methods. In order to obtain an absolute measure of accuracy of the taxonomic assignments we use  
244 the expected reference sequences from the mock datasets as the query sequences for the databases  
245 and the expected taxonomic composition as the standard to compare against (Figure 4C). Again, we  
246 observe that none of the databases perform better than the others in absolute terms.

247 Given that no database performs better than others against mock datasets, and that databases are  
248 almost equally distant from each other in terms of final output, the choice of which database to use  
249 should be driven by other reason. One user-specific way to choose, would be based on the known  
250 representation of taxa for the microbiome of interest (see also Discussion). Another reason could be  
251 the frequency of updates and the potential for future growth, which prompted us to set NCBI as the  
252 MiCoNE standard for taxonomy assignment. In addition to being regularly maintained and updated  
253 the NCBI database already has the advantage that its accuracy of assignments is still comparable to  
254 the SILVA and GG reference databases that are routinely used as reference databases.

255 **Networks generated using different network inference methods show notable**  
256 **difference in edge-density and connectivity**

257 The six different network inference methods used in this study are Microbial Association  
258 Graphical Model Analysis (MAGMA) [27], metagenomic Lognormal-Dirichlet-Multinomial  
259 (mLDM) [42], Sparse InversE Covariance estimation for Ecological Association and Statisti-  
260 cal Inference (SpiecEasi) [28], Sparse Correlations for Compositional data (SparCC) [19], Spearman  
261 and Pearson. These network inference methods fall into two groups, the first set of methods (Pear-  
262 son, Spearman, SparCC) infer pairwise correlations while the second set infer direct associations  
263 (SpiecEasi, mLDM, MAGMA). Pairwise correlation methods involve calculating the correlation

---

264 coefficient between every pair of OTU/ESVs leading to the detection of spurious indirect connections.  
265 On the other hand, direct association methods use conditional independence to avoid the detection  
266 of correlated but indirectly connected OTUs [28, 8].

267 For the analysis presented in this section, we used the taxonomy composition table obtained  
268 using the NCBI reference database as the input for algorithms that infer co-occurrence associations  
269 between the microbes. Figure 5A shows the networks inferred from this dataset using the different  
270 inference algorithms. The different networks differ vastly in their edge-density and connectivity;  
271 even some of the edges in common to these networks have their signs inverted. Note, however,  
272 that some of these comparisons depend on the threshold that has to be applied to the pairwise  
273 correlations methods (currently 0.3, based on [19]). To get a more quantitative picture of the  
274 differences between the inferred networks, we checked the distribution of common nodes and edges  
275 (Figure 5B) using UpSet plots [43] (only MAGMA, mLDM, SpiecEasi, SparCC are used in the  
276 comparison since Pearson and Spearman add a large number of spurious edges since they are not  
277 intended for compositional datasets). The results for the node intersections show that the networks  
278 have a large number of nodes in common (63 out of 67 nodes in the smallest network - MAGMA)  
279 and no network possesses any unique node. The edge intersections in contrast show that only  
280 19 edges (out of 98 edges in the smallest network - MAGMA) are in common between all the  
281 methods and each network has a large number of unique edges. These results indicate that there is a  
282 substantial rewiring of connections in the inferred networks.

283 Unlike the previous steps of the pipeline, where we evaluated the performance of methods on  
284 mock datasets, there is no equivalent dataset that contains a set of known interactions for the evaluation  
285 of the network inference algorithms. Therefore, we propose the construction of a consensus network  
286 (Figure 5C) involving MAGMA, mLDM, SpiecEasi and SparCC. This consensus network is built  
287 by merging the p-values generated from bootstraps of the original taxonomy composition table

---

288 using the Browns p-value combining method [44] (see Methods section). Based on this approach,  
289 MiCoNE reports as default output the consensus network, annotated with weights (correlations for  
290 SparCC and direct associations for the other methods) for all four methods.

291 **The default pipeline**

292 The systematic analyses performed in the previous sections clearly show that the choice of tools and  
293 parameters can have a big impact on the final co-occurrence network. For some of these choices (e.g.  
294 DADA2 vs. deblur) there is no clear metric to establish a best protocol. For other choices, the mock  
295 communities provide an opportunity to select combination of parameters that yield more accurate  
296 and robust results. Despite this partial degree of assessment, we wish to suggest a combination  
297 of tools and parameters that produce networks that are derived from the combination of tools  
298 which performed best on the mock communities, and displayed highest robustness to switching to  
299 alternative methods. These tools and parameters are chosen as the defaults for the pipeline and are  
300 given in Table 1.

301 The recommended tool for the Denoising and Clustering (DC) step (DADA2 or Deblur) were  
302 chosen based on their accuracy in recapitulating the reference sequences in mock communities and  
303 synthetic data. The choice of the taxonomy reference database in the Taxonomy Assignment (TA)  
304 step is dictated largely by the species expected to be present in the sample as well the database used  
305 in similar studies if comparison is a goal. Nevertheless, we suggest NCBI RefSeq along with blast+  
306 as the query tool since the database is updated regularly and has a broad collection of taxonomies.  
307 The abundance threshold at the OTU Processing (OP) step is determined automatically based on the  
308 number of samples and the required statistical power. Finally, we use the Browns p-value combining  
309 method on the networks generated using MAGMA, mLDM, SpiecEasi and SparCC to obtain a final  
310 consensus network in the Network Inference (NI) step.

---

311      Figure 6A shows the default network compared against networks generated by altering one of the  
312      steps of the pipeline from the default. These results indicate that the biggest differences in networks  
313      occur when the reference database or the network inference algorithm are changed. Furthermore, the  
314      L1 distance of networks generated by altering one of the steps of the pipeline from the default against  
315      the default network (Figure 6B) shows that the biggest deviations from the default network occur  
316      when the TA and NI steps are changed, reinforcing the same results observed in Figure 2. Figure 7  
317      shows the co-occurrence networks inferred for the hard palate for healthy subjects in a periodontal  
318      disease study [45] and the healthy stool microbiome in fecal microbial transplant study [35]. These  
319      consensus networks were generated using the default tools and parameters from Table 1.

320      **Discussion**

321      Co-occurrence associations in microbial communities help identify important interactions that drive  
322      microbial community structure and organization. Our analysis shows that networks generated using  
323      different combinations of tools and approaches can look significantly different from each other,  
324      highlighting the importance of a clear assessment of the source of variability and of tools that provide  
325      the most robust and accurate results. Our newly developed integrated software for the inference  
326      of co-occurrence networks from 16S rRNA data, MiCoNE, constitutes a freely customizable and  
327      user friendly pipeline that allows users to easily test combinations of tools and to compare networks  
328      generated by multiple possible choices (see Methods). Importantly, in addition to revisiting the test  
329      cases presented in this work, users will be able to explore the effect of various tool combinations on  
330      their own datasets of interest. The MiCoNE pipeline is built in a modular fashion. Its plug-and-play  
331      architecture will make it possible for users to add new tools and steps, either from existing packages,  
332      or from packages that were not examined in the present work, as well as future ones.

---

333     The main outcome of this work is thus two-fold: on one hand we transparently reveal the  
334     dependence of co-occurrence networks on tool and parameter choices, making it possible to more  
335     rigorously assess and compare existing networks. On the other hand, we take advantage of our  
336     spectrum of computational options and the availability of mock and synthetic datasets, to suggest a  
337     default standard setting, and a consensus approach, likely to yield networks that are robust across  
338     multiple tool/parameter choices.

339     An important caveat related to this last point is the fact that our conclusions are based on the  
340     specific datasets used in our analysis. While our datasets cover a relatively broad spectrum of  
341     biomes and sequencing pipelines, datasets that have drastically different distributions may require a  
342     re-assessment of the best settings through our pipeline.

343     It is worth pointing out some additional more specific conclusions stemming from the individual  
344     steps of our analysis.

345     The different denoising/clustering methods differ mostly in their identification of sequences that  
346     are in low abundances. Hence, they do not have much of an impact on the inferred co-occurrence  
347     networks when the sequences of low abundance are removed. However, comparison of inferred and  
348     expected reference sequences and their abundances in mock community datasets has allowed us to  
349     identify DADA2 as the method which best recapitulates the expected sequence composition. For  
350     the current work we have decided to focus on the tools most widely used at the time of the analysis.  
351     Some tools that we recently published (e.g. dbOTU3 [46]) as well as older popular methods like  
352     mothur [47] have not been included in the study, but could be added into the pipelines in future  
353     updated analyses.

354     The choice of taxonomy database was found to be the most important factor in the inference of a  
355     microbial co-occurrence network, contributing ~ 20% of the total variance. The frequent changes  
356     in the taxonomy nomenclature coupled with the frequency of updates to the various 16S reference

---

357 databases create inherent differences [41] in taxonomy hierarchies in these databases. Our analysis  
358 revealed that no particular reference database performs better than the others across all scenarios.  
359 We suggest that that choice of the database should be made based on possible reported or inferred  
360 biases in the representation of given biomes in a specific databases [41]. The default reference  
361 database in the pipeline is the NCBI 16S RefSeq database as it is more frequently updated and is  
362 most compatible with the blast+ query tool. We also enable users to use custom databases [48] with  
363 the blast+ and naive bayes classifiers that are incorporated into the pipeline (from QIIME2).

364 Filtering out taxa that are present in low abundances in all samples did not increase (in most  
365 datasets tested) the proportion of taxa in common between taxonomy tables generated using different  
366 reference databases. However, we do observe that the reduction in the number of taxa leads to better  
367 agreement in the networks inferred through different methods. Moreover, filtering is necessary in  
368 order to increase the power in tests of significance when the number of taxa is much greater than the  
369 number of samples.

370 The networks generated by different network inference methods show considerable differences in  
371 edge-density and connectivity. One reason for this is the underlying assumptions regarding sparsity,  
372 distribution and compositionality that the algorithms make. The consensus network created by  
373 merging the networks inferred using the different network inference methods enables the creation of  
374 a network whose links have evidence based on multiple inference algorithms.

375 Exploring the effects of these combinations of methods on the resultant networks is difficult and  
376 inconvenient since different tools differ in their input and output formats and require inter-converting  
377 between the various formats. The pipeline facilitates this comparative exploration by providing a  
378 variety of modules for inter-conversion between various formats, and by allowing easy incorporation  
379 of new tools as modules.

380 We envision that MiCoNE, and the underlying tools and databases that help process amplicon

---

381 sequencing data into co-occurrence networks, will be increasingly useful towards building large  
382 comparative analyses across studies. By having a unified transparent tool to compute networks, it  
383 will be possible to reprocess available 16S datasets to obtain networks that are directly comparable  
384 to each other. Furthermore, even in the analysis of published networks across studies and processing  
385 methods, MiCoNE could help understand underlying biases of each network, which could in turn be  
386 taken into account upon making cross-study comparisons.

387 **Materials and Methods**

388 **Datasets**

389 The study uses three kinds of 16S rRNA sequencing datasets: real datasets, mock datasets and  
390 synthetic datasets. Real datasets are collections of sequencing reads obtained from naturally  
391 occurring microbial community samples. The current study used healthy stool samples from a fecal  
392 microbiome transplant study [35] and healthy saliva samples from a periodontal disease study [45]  
393 as real datasets for analysis. The mock community 16S datasets are real sequencing data obtained  
394 for artificially assembled collections of species in known proportions. The mock datasets used  
395 for this study, obtained from mockrobiota [33], are labelled mock4, mock12 and mock16. The  
396 mock4 community is composed of 21 bacterial strains. Two replicate samples from mock4 contain  
397 all species in equal abundances, and two additional replicate samples contain the same species in  
398 unequal abundances. The mock12 community is composed of 27 bacterial strains that include  
399 closely related taxa with some pairs having only one to two nucleotide difference from another. The  
400 mock16 community is composed of 49 bacteria and 10 archaea, all represented in equal amount.  
401 The synthetic datasets were generated using an artificial read simulator called ART [34]. Three  
402 different microbial composition profiles were used as input; reads were generated using a soil and

---

403 water microbiome composition profiles from the EMP [2] and healthy gut microbiome project  
404 from the fecal microbiome transplant study [35]. The reads are simulated using the NCBI RefSeq  
405 database as the reference sequence pool and the "art\_illumina" sequence profile with a mutation  
406 rate of 2%. The scripts used to generate the synthetic data are in the scripts folder of the repository  
407 (<https://github.com/segrelab/MiCoNE-pipeline-paper>).

408 **MiCoNE**

409 The flowchart describing the workflow of MiCoNE (Microbial Co-occurrence Network Explorer),  
410 our complete 16S data-analysis pipeline, is shown in Figure 1. The pipeline integrates many  
411 publicly available tools as well as custom R or Python modules and scripts to extract co-occurrence  
412 associations from 16S sequence data. Each of these tools corresponds to a distinct R or python  
413 module that recapitulates the relevant analyses. All such individual modules are available as part  
414 of the MiCoNE package. The inputs to the pipeline by default are the raw community 16S rRNA  
415 sequence reads, but the software can be alternatively configured to use trimmed sequences, OTU  
416 tables and other types of intermediate data. The final output of the pipeline is the inferred network  
417 of co-occurrence relationships among the microbes present in the samples.

418 The MiCoNE pipeline provides both a Python API as well as a command-line interface and  
419 only requires a single configuration file. The configuration file lists the inputs, output and the steps  
420 to be performed during runtime, along with the parameters to be used (if different from defaults)  
421 for the various steps. Since the entire pipeline run-through is stored in the form of a text file (the  
422 configuration file), subsequent runs are highly reproducible and changes can be easily tracked using  
423 version control. It uses the nextflow workflow manager [49] under the hood, making it readily usable  
424 on local machines, cluster or cloud with minimal configuration change. It also allows for automatic  
425 parallelization of all possible processes, both within and across samples. The pipeline is designed to

---

426 be modular: each tool or method is organized into modules which can be easily modified or replaced.  
427 This modular architecture simplifies the process of adding new tools (refer to modules section in  
428 the MiCoNE documentation). The main components of the pipeline are detailed in the subsequent  
429 sections.

### 430 **Denoising and Clustering (DC)**

431 This module deals with processing the raw 16S sequence data into OTU or ESV count tables. It  
432 consists of the following processes: quality control, denoising (or clustering) and chimera checking.  
433 The quality control process handles the demultiplexing and quality control steps such as trimming  
434 adapters and trimming low-quality nucleotide stretches from the sequences. The denoise/cluster  
435 process handles the conversion of the demultiplexed, trimmed sequences into OTU or ESV count  
436 tables (some methods, like closed reference and open reference clustering, perform clustering and  
437 taxonomy assignment in the same step). The chimera checking process handles the removal of  
438 chimeric sequences created during the Polymerase Chain Reaction (PCR) step. The output of this  
439 module is a matrix of counts, that describes the number of reads of a particular OTU or ESV (rows  
440 of the matrix) present in each sample (columns of the matrix). The options currently available in  
441 the pipeline for denoising and clustering are: open reference clustering, closed reference clustering  
442 and de novo clustering methods from QIIME1 v1.9.1 [22] and denoising methods from DADA2  
443 v1.14 [23] and Deblur v1.1.0 [36]. The quality filtering and chimera checking tools are derived  
444 from those used in QIIME2 v2019.10.0 and DADA2.

### 445 **Taxonomy Assignment (TA)**

446 This module deals with assigning taxonomies to either the representative sequences of the OTUs or  
447 directly to the ESVs. In order to assign taxonomies to a particular sequence we need a taxonomy  
448 database and a query tool. The taxonomy database contains the collection of 16S sequences of

---

449 micro-organisms of interest and the query tool allows one to compare a sequence of interest to all  
450 the sequences in the database to identify the best matches. Finally, a consensus method is used  
451 to identify the most probable match from the list of best matches. The pipeline incorporates GG  
452 13\_8 [24], SILVA 132 [25] and the NCBI (16S RefSeq as of Oct 2019) [40] databases for taxonomy  
453 assignment and the Naive Bayes classifier from QIIME2 and NCBI blast as the query tools (from  
454 QIIME2). The consensus algorithm used is the default method used by the classifiers in QIIME2.

## 455 **OTU and ESV Processing (OP)**

456 This module deals with normalization, filtering and applying transformations to the OTU or ESV  
457 counts matrix. Rarefaction is a normalization technique used to overcome the bias that might arise  
458 due to variable sampling depth in different samples. This is performed either by sub-sampling or  
459 by normalization of the matrix to the lowest sampling depth [26]. Rarefaction is usually followed  
460 by filtering, which is performed to remove samples or features (OTUs or ESVs) from the count  
461 matrix that are sparse. In order to determine the filtering threshold we fix the number of samples  
462 and correlation detection power needed and determine the number of features to be used. Finally,  
463 transformations are performed in order to correct for and overcome the compositional bias that is  
464 inherent in a counts matrix (in most cases this is handled by the network inference algorithm).

## 465 **Network Inference (NI)**

466 This module deals with the inference of co-occurrence associations from the OTU or ESV counts  
467 matrix. These associations can be represented as a network, with nodes representing taxonomies of  
468 the micro-organisms and edges representing the association between them. A null model is created  
469 by re-sampling and bootstrapping the correlation/interaction matrix and is used to calculate the  
470 significance of the inferred associations by calculating the p-values against this null model [50]. The  
471 pipeline includes Pearson, Spearman and FastSpar v0.0.10 (a faster implementation of SparCC) [50]

---

472 as the pairwise correlation metrics, and SpiecEasi v1.0.7 [28], mLDM v1.1 [42] and MAGMA [27]  
473 as the direct association metrics. The empirical Browns method [44] is used for combining p-values  
474 from the various methods to obtain a consensus p-value, which is used to create the consensus  
475 network.

476 **Network Variability**

477 In order to compare across different networks, and analyze the degree of variability induced by  
478 the choice of different modules and parameters, we organized multiple networks into a single  
479 mathematical structure that we could use for linear regression. In particular, we transformed the  
480 adjacency matrix of each co-occurrence network into a vector. We then merged the networks  
481 generated from all possible combinations of tools into a table ( $N$ , see below) in which each column  
482 represents one network.

$$N = \begin{bmatrix} edge_{1,1} & edge_{2,1} & \cdots & edge_{n,1} \\ edge_{1,2} & edge_{2,2} & \cdots & edge_{n,2} \\ \vdots & \vdots & \vdots & \vdots \\ edge_{1,n} & edge_{2,n} & \cdots & edge_{n,n} \end{bmatrix}$$

483 In other words,  $N$  is the merged table, each column  $N_i$  is the vector representation of one of the  
484 networks, and each row  $L_i$  represents the one particular edge in all networks (assigned 0 if the edge  
485 does not exist in the network).

486 We use linear regression to express each link  $L_i$  as a linear function of categorical variables that  
487 describe the possible options in each of the first three steps of the pipeline.

---

In particular, we infer parameters  $\alpha_i$  such that:

$$L_i = \sum_{j=1}^5 (\alpha_i^{DC(j)} \cdot \delta_i^{DC(j)}) + \sum_{j=1}^3 (\alpha_i^{TA(j)} \cdot \delta_i^{TA(j)}) + \sum_{j=1}^2 (\alpha_i^{OP(j)} \cdot \delta_i^{OP(j)}) + \epsilon_i$$

488 where,  $\alpha_i$  are the coefficients of the regression,  $\epsilon_i$  are the residuals and  $\delta_i$  are the indicator  
489 variables that correspond to the processes utilized in the pipeline used to create the network  $N_i$ ;  
490 for example,  $\delta_i^{DC(1)} = 1$  if the DC(1) process was used in the generation of the network  $N_i$ . Here,  
491 (i) DC(1) = "closed reference", DC(2) = "open reference", DC(3) = "de novo", DC(4) = "dada2",  
492 DC(5) = "deblur"; (ii) TA(1) = "GreenGenes", TA(2) = "SILVA", TA(3) = "NCBI"; (iii) OP(1) =  
493 "no filtering", OP(2) = "filtering".

494 The variance contributed by each step of the pipeline is calculated for every connection in the  
495 merged table through ANOVA using the Python statsmodels package and is shown in Figure 2B.  
496 The total variance for the network is calculated by adding the variances for each connection. The  
497 PCA analysis is also performed on the merged table to generate Figure 2C.

498 **Code and Data Availability**

499 Pipeline: <https://github.com/segrelab/MiCoNE>  
500 Data and scripts: <https://github.com/segrelab/MiCoNE-pipeline-paper>

501 **Acknowledgments**

502 We are grateful to members of the Segrè lab for helpful discussions and for feedback on the  
503 manuscript. This work was partially funded by grants from the National Institutes of Health  
504 (National Institute of General Medical Sciences, award R01GM121950; National Institute of Dental

---

505 and Craniofacial Research, award number R01DE024468; and National Institute on Aging, award  
506 number UH2AG064704), the U.S. Department of Energy, Office of Science, Office of Biological &  
507 Environmental Research through the Microbial Community Analysis and Functional Evaluation in  
508 Soils SFA Program (m-CAFEs) under contract number DE-AC02-05CH11231 to Lawrence Berkeley  
509 National Laboratory, the National Science Foundation (grants 1457695 and NSFOCE-BSF 1635070),  
510 the Human Frontiers Science Program (RGP0020/2016), and the Boston University Interdisciplinary  
511 Biomedical Research Office. KSK was supported by Simons Foundation Grant #409704, by the  
512 Research Corporation for Science Advancement through Cottrell Scholar Award #24010, by the  
513 Scialog grant #26119, and by the Gordon and Betty Moore Foundation grant #6790.08.

514 **Contributions**

515 Designed the research project: DK, KK, DS, ZH, CDL. Performed analysis: DK, GB. Wrote the  
516 first draft of the manuscript: DK. Revised and wrote final version of the manuscript: DK, DS, KK.

## References

- [1] Melanie Ghoul and Sara Mitri. “The Ecology and Evolution of Microbial Competition.” In: *Trends in microbiology* 24.10 (Oct. 2016), pp. 833–845. ISSN: 1878-4380. doi: 10.1016/j.tim.2016.06.011. URL: <http://www.ncbi.nlm.nih.gov/pubmed/27546832>.
- [2] Luke R. Thompson et al. “A communal catalogue reveals Earth’s multiscale microbial diversity”. In: *Nature* 551.7681 (Nov. 2017), p. 457. ISSN: 0028-0836. doi: 10.1038/nature24621. URL: <http://www.nature.com/doifinder/10.1038/nature24621>.
- [3] Takashi Narihiro and Yoichi Kamagata. “Genomics and Metagenomics in Microbial Ecology: Recent Advances and Challenges.” In: *Microbes and environments* 32.1 (2017), pp. 1–4. ISSN: 1347-4405. doi: 10.1264/jsme2.ME3201rh. URL: [http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC5371069](http://www.ncbi.nlm.nih.gov/pubmed/28367917%20http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC5371069).
- [4] Juan Jovel et al. “Characterization of the Gut Microbiome Using 16S or Shotgun Metagenomics.” In: *Frontiers in microbiology* 7 (2016), p. 459. ISSN: 1664-302X. doi: 10.3389/fmicb.2016.00459. URL: [http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC4837688](http://www.ncbi.nlm.nih.gov/pubmed/27148170%20http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC4837688).
- [5] Jason Lloyd-Price, Galeb Abu-Ali, and Curtis Huttenhower. “The healthy human microbiome.” In: *Genome medicine* 8.1 (2016), p. 51. ISSN: 1756-994X. doi: 10.1186/s13073-016-0307-y. URL: [http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC4848870](http://www.ncbi.nlm.nih.gov/pubmed/27122046%20http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC4848870).
- [6] Houjin Zhang and Kang Ning. “The Tara Oceans Project: New Opportunities and Greater Challenges Ahead.” In: *Genomics, proteomics & bioinformatics* 13.5 (Oct. 2015), pp. 275–7. ISSN: 2210-3244. doi: 10.1016/j.gpb.2015.08.003. URL: <http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC4678785>.
- [7] Barbara A. Human Microbiome Project Consortium et al. “A framework for human microbiome research.” In: *Nature* 486.7402 (June 2012), pp. 215–21. ISSN: 1476-4687. doi: 10.1038/nature11209. URL: <http://www.ncbi.nlm.nih.gov/pmcarticlerender.fcgi?artid=PMC3377744>.
- [8] Rajita Menon, Vivek Ramanan, and Kirill S. Korolev. “Interactions between species introduce spurious associations in microbiome studies”. In: *PLOS Computational Biology* 14.1 (Jan. 2018). Ed. by Stefano Allesina, e1005939. ISSN: 1553-7358. doi: 10.1371/journal.pcbi.1005939. URL: <http://dx.plos.org/10.1371/journal.pcbi.1005939>.

- 
- 552 [9] Lisa Röttjers and Karoline Faust. “From hairballs to hypotheses—biological insights from  
553 microbial networks”. In: *FEMS Microbiology Reviews* 42.6 (Nov. 2018), pp. 761–780.  
554 ISSN: 1574-6976. doi: 10.1093/femsre/fuy030. URL: <https://academic.oup.com/femsre/article/42/6/761/5061627>.
- 555 [10] Jack A. Gilbert et al. “Microbiome-wide association studies link dynamic microbial consortia  
556 to disease”. In: *Nature* 535.7610 (2016), pp. 94–103. ISSN: 14764687. doi: 10.1038/nature18850. arXiv: NIHMS150003.
- 557 [11] Baohong Wang et al. “The Human Microbiota in Health and Disease”. In: *Engineering*  
558 3.1 (Feb. 2017), pp. 71–82. ISSN: 20958099. doi: 10.1016/j.eng.2017.01.008. URL:  
559 <http://linkinghub.elsevier.com/retrieve/pii/S2095809917301492>.
- 560 [12] José E Belizário and Mauro Napolitano. “Human microbiomes and their roles in dysbiosis,  
561 common diseases, and novel therapeutic approaches.” In: *Frontiers in microbiology* 6  
562 (2015), p. 1050. ISSN: 1664-302X. doi: 10.3389/fmicb.2015.01050. URL: <http://www.ncbi.nlm.nih.gov/pubmed/26500616%20http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC4594012>.
- 563 [13] Noah Fierer. *Embracing the unknown: Disentangling the complexities of the soil microbiome*.  
564 Oct. 2017. doi: 10.1038/nrmicro.2017.87. URL: <https://pubmed.ncbi.nlm.nih.gov/28824177/>.
- 565 [14] Shuo Jiao, Weimin Chen, and Gehong Wei. “Resilience and assemblage of soil microbiome  
566 in response to chemical contamination combined with plant growth”. In: *Applied and  
567 Environmental Microbiology* 85.6 (Mar. 2019). ISSN: 10985336. doi: 10.1128/AEM.02523-  
568 18. URL: <http://aem.asm.org/>.
- 569 [15] Ryan H. Hsu et al. “Microbial Interaction Network Inference in Microfluidic Droplets”. In: *Cell  
570 Systems* 9.3 (Sept. 2019), 229–242.e4. ISSN: 24054720. doi: 10.1016/j.cels.2019.06.  
571 008. URL: [http://www.cell.com/article/S2405471219302315/fulltext%20http://www.cell.com/article/S2405471219302315/abstract%20https://www.cell.com/cell-systems/abstract/S2405-4712\(19\)30231-5](http://www.cell.com/article/S2405471219302315/fulltext%20http://www.cell.com/article/S2405471219302315/abstract%20https://www.cell.com/cell-systems/abstract/S2405-4712(19)30231-5).
- 572 [16] Xingjin Jian et al. “Microbial microdroplet culture system (MMC): An integrated platform for  
573 automated, high-throughput microbial cultivation and adaptive evolution”. In: *Biotechnology  
574 and Bioengineering* 117.6 (June 2020), pp. 1724–1737. ISSN: 0006-3592. doi: 10.1002/bit.  
575 27327. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/bit.27327>.
- 576 [17] Steven A Wilbert, Jessica L Mark Welch, and Gary G Borisy. “Spatial Ecology of the Human  
577 Tongue Dorsum Microbiome”. In: *Cell Reports* 30 (2020), 4003–4015.e3. doi: 10.1016/j.  
578 celrep.2020.02.097. URL: <https://doi.org/10.1016/j.celrep.2020.02.097>.
- 579 [18] Cristal Zuñiga, Livia Zaramela, and Karsten Zengler. “Elucidation of complexity and  
580 prediction of interactions in microbial communities”. In: *Microbial Biotechnology* 10.6  
581 (2017), pp. 1500–1522. doi: 10.1111/1751-7915.12855.
- 582
- 583
- 584
- 585
- 586
- 587
- 588

- 
- 589 [19] Jonathan Friedman and Eric J. Alm. “Inferring Correlation Networks from Genomic Survey  
590 Data”. In: *PLoS Computational Biology* 8.9 (Sept. 2012). Ed. by Christian von Mering,  
591 e1002687. ISSN: 1553-7358. doi: 10.1371/journal.pcbi.1002687. URL: <http://dx.plos.org/10.1371/journal.pcbi.1002687>.
- 593 [20] Richa Bharti and Dominik G Grimm. “Current challenges and best-practice protocols for  
594 microbiome analysis”. In: *Briefings in Bioinformatics* 2019.00 (Dec. 2019), pp. 1–16. ISSN:  
595 1477-4054. doi: 10.1093/bib/bbz155. URL: <https://academic.oup.com/bib/advance-article/doi/10.1093/bib/bbz155/5678919>.
- 597 [21] Jolinda Pollock et al. *The madness of microbiome: Attempting to find consensus "best  
598 practice" for 16S microbiome studies*. Apr. 2018. doi: 10.1128/AEM.02627-17. URL:  
599 <https://doi.org/10.1128/AEM.02627-17..>
- 600 [22] J Gregory Caporaso et al. “QIIME allows analysis of high-throughput community sequencing  
601 data”. In: *Nature Methods* 7.5 (May 2010), pp. 335–336. ISSN: 1548-7091. doi: 10.1038/nmeth.f.303. URL: <http://www.nature.com/articles/nmeth.f.303>.
- 603 [23] Benjamin J Callahan et al. “DADA2: High-resolution sample inference from Illumina  
604 amplicon data”. In: *Nature Methods* 13.7 (July 2016), pp. 581–583. ISSN: 1548-7091. doi:  
605 10.1038/nmeth.3869. URL: <http://www.nature.com/articles/nmeth.3869>.
- 606 [24] T Z DeSantis et al. “Greengenes, a chimera-checked 16S rRNA gene database and workbench  
607 compatible with ARB.” In: *Applied and environmental microbiology* 72.7 (July 2006),  
608 pp. 5069–72. ISSN: 0099-2240. doi: 10.1128/AEM.03006-05. URL: <http://www.ncbi.nlm.nih.gov/pubmed/16820507%20http://www.ncbi.nlm.nih.gov/articlerender.fcgi?artid=PMC1489311>.
- 611 [25] Christian Quast et al. “The SILVA ribosomal RNA gene database project: improved data  
612 processing and web-based tools”. In: *Nucleic Acids Research* 41.D1 (Nov. 2012), pp. D590–  
613 D596. ISSN: 0305-1048. doi: 10.1093/nar/gks1219. URL: <http://academic.oup.com/nar/article/41/D1/D590/1069277/The-SILVA-ribosomal-RNA-gene-database-project>.
- 616 [26] Sophie J Weiss et al. “Effects of library size variance, sparsity, and compositionality on the  
617 analysis of microbiome data”. In: *PeerJ PrePrints* 3 (2015), e1408. ISSN: 2167-9843. doi:  
618 10.7287/peerj.preprints.1157v1. arXiv: peerj.preprints.270v1 [10.7287].  
619 URL: <https://doi.org/10.7287/peerj.preprints.1157v1%7B%5C%7D5Cnhttps://peerj.com/preprints/1157v1/%7B%5C%7Dsupp-8>.
- 621 [27] Arnaud Cougoul, Xavier Bailly, and Ernst C Wit. “MAGMA: inference of sparse microbial  
622 association networks”. In: (2019). doi: 10.1101/538579. URL: <https://doi.org/10.1101/538579>.
- 623

- [28] Zachary D. Kurtz et al. “Sparse and Compositionally Robust Inference of Microbial Ecological Networks”. In: *PLOS Computational Biology* 11.5 (May 2015). Ed. by Christian von Mering, e1004226. ISSN: 1553-7358. DOI: 10.1371/journal.pcbi.1004226. URL: <http://dx.plos.org/10.1371/journal.pcbi.1004226>.

[29] Kevin P. Keegan, Elizabeth M. Glass, and Folker Meyer. “MG-RAST, a Metagenomics Service for Analysis of Microbial Community Structure and Function”. In: Humana Press, New York, NY, 2016, pp. 207–233. DOI: 10.1007/978-1-4939-3369-3\_13. URL: [http://link.springer.com/10.1007/978-1-4939-3369-3%7B%5C\\_%7D13](http://link.springer.com/10.1007/978-1-4939-3369-3%7B%5C_%7D13).

[30] *Qiita - open-source microbial study management platform*. URL: <https://qiita.ucsd.edu/> (visited on 05/22/2018).

[31] Jonathan L. Golob et al. “Evaluating the accuracy of amplicon-based microbiome computational pipelines on simulated human gut microbial communities”. In: *BMC Bioinformatics* 18.1 (2017), p. 283. ISSN: 1471-2105. DOI: 10.1186/s12859-017-1690-0. URL: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-017-1690-0>.

[32] Sophie Weiss et al. “Correlation detection strategies in microbial data sets vary widely in sensitivity and precision”. In: *ISME J* 10.7 (2016), pp. 1–13. ISSN: 1751-7362. DOI: 10.1038/ismej.2015.235. URL: <http://dx.doi.org/10.1038/ismej.2015.235>.

[33] Nicholas A Bokulich et al. “mockrobiota: a Public Resource for Microbiome Bioinformatics Benchmarking.” In: *mSystems* 1.5 (2016). ISSN: 2379-5077. DOI: 10.1128/mSystems.00062-16. URL: <http://www.ncbi.nlm.nih.gov/pubmed/27822553%20http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5080401>.

[34] Weichun Huang et al. “ART: a next-generation sequencing read simulator.” In: *Bioinformatics (Oxford, England)* 28.4 (Feb. 2012), pp. 593–4. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btr708. URL: <http://www.ncbi.nlm.nih.gov/pubmed/22199392%20http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3278762>.

[35] Dae-Wook Kang et al. “Microbiota Transfer Therapy alters gut ecosystem and improves gastrointestinal and autism symptoms: an open-label study”. In: *Microbiome* 5.1 (Dec. 2017), p. 10. ISSN: 2049-2618. DOI: 10.1186/s40168-016-0225-7. URL: <http://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-016-0225-7>.

[36] Amnon Amir et al. “Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns.” In: *mSystems* 2.2 (Apr. 2017), e00191–16. ISSN: 2379-5077. DOI: 10.1128/mSystems.00191-16. URL: <http://www.ncbi.nlm.nih.gov/pubmed/28289731%20http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5340863>.

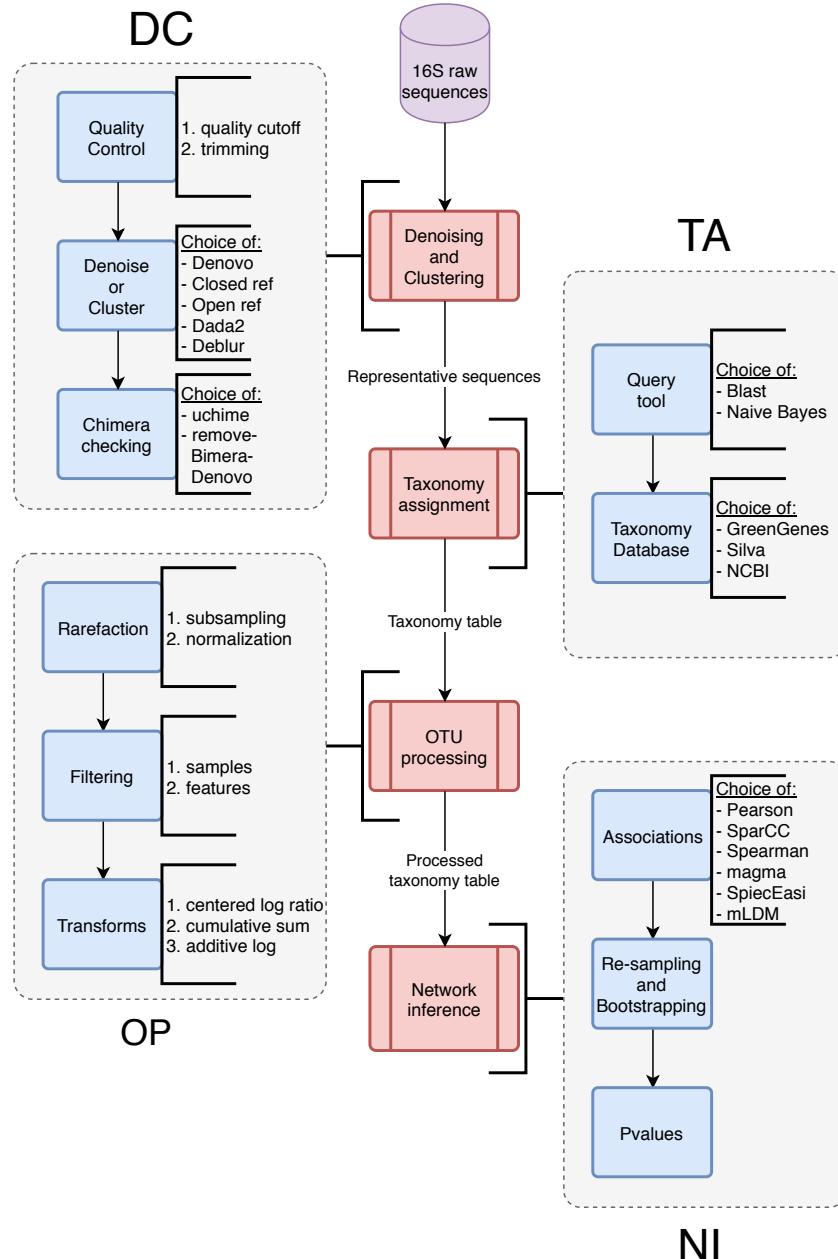
- 
- 661 [37] Catherine A Lozupone et al. “Quantitative and qualitative beta diversity measures lead  
662 to different insights into factors that structure microbial communities.” In: *Applied and*  
663 *environmental microbiology* 73.5 (Mar. 2007), pp. 1576–85. ISSN: 0099-2240. DOI: 10.1128/  
664 AEM.01996-06. URL: <http://www.ncbi.nlm.nih.gov/pubmed/17220268%20http://www.ncbi.nlm.nih.gov/pmcmlerender.fcgi?artid=PMC1828774>.
- 666 [38] Catherine Lozupone and Rob Knight. “UniFrac: a new phylogenetic method for comparing  
667 microbial communities.” In: *Applied and environmental microbiology* 71.12 (Dec. 2005),  
668 pp. 8228–35. ISSN: 0099-2240. DOI: 10.1128/AEM.71.12.8228-8235.2005. URL: <http://www.ncbi.nlm.nih.gov/pubmed/16332807%20http://www.ncbi.nlm.nih.gov/pmcmlerender.fcgi?artid=PMC1317376>.
- 671 [39] Jacob T. Nearing et al. “Denoising the Denoisers: an independent evaluation of microbiome  
672 sequence error-correction approaches”. In: *PeerJ* 6 (Aug. 2018), e5364. ISSN: 2167-8359.  
673 DOI: 10.7717/peerj.5364. URL: <https://peerj.com/articles/5364>.
- 674 [40] Eric W Sayers et al. “Database resources of the National Center for Biotechnology Information.” In: *Nucleic acids research* 37.Database issue (Jan. 2009), pp. D5–15. ISSN: 1362-  
675 4962. DOI: 10.1093/nar/gkn741. URL: <http://www.ncbi.nlm.nih.gov/pubmed/18940862%20http://www.ncbi.nlm.nih.gov/pmcmlerender.fcgi?artid=PMC2686545>.
- 679 [41] Monika Balvočiūtė and Daniel H. Huson. “SILVA, RDP, Greengenes, NCBI and OTT  
680 — how do these taxonomies compare?” In: *BMC Genomics* 18.S2 (Mar. 2017), p. 114.  
681 ISSN: 1471-2164. DOI: 10.1186/s12864-017-3501-4. URL: <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12864-017-3501-4>.
- 683 [42] Yuqing Yang, Ning Chen, and Ting Chen. “Inference of Environmental Factor-Microbe  
684 and Microbe-Microbe Associations from Metagenomic Data Using a Hierarchical Bayesian  
685 Statistical Model.” In: *Cell systems* 4.1 (Jan. 2017), 129–137.e5. ISSN: 2405-4712. DOI:  
686 10.1016/j.cels.2016.12.012. URL: <http://www.ncbi.nlm.nih.gov/pubmed/28125788>.
- 688 [43] Alexander Lex et al. *UpSet: Visualization of Intersecting Sets*. Tech. rep. URL: <http://grouplens.org/datasets/movielens/>.
- 690 [44] William Poole et al. “Combining dependent P-values with an empirical adaptation of  
691 Brown’s method”. In: (). DOI: 10.1093/bioinformatics/btw438. URL: <https://www.bioconduc>.
- 693 [45] Casey Chen et al. “Oral microbiota of periodontal health and disease and their changes after  
694 nonsurgical periodontal therapy”. In: *The ISME Journal* 12.5 (May 2018), pp. 1210–1224.  
695 ISSN: 1751-7362. DOI: 10.1038/s41396-017-0037-1. URL: <http://www.nature.com/articles/s41396-017-0037-1>.

- 
- 697 [46] Scott W. Olesen, Claire Duvallet, and Eric J. Alm. “DbOTU3: A new implementation of  
698 distribution-based OTU calling”. In: *PLoS ONE* 12.5 (2017), pp. 1–13. ISSN: 19326203. DOI:  
699 [10.1371/journal.pone.0176335](https://doi.org/10.1371/journal.pone.0176335).
- 700 [47] Patrick D Schloss et al. “Introducing mothur: open-source, platform-independent, community-  
701 supported software for describing and comparing microbial communities.” In: *Applied and*  
702 *environmental microbiology* 75.23 (Dec. 2009), pp. 7537–41. ISSN: 1098-5336. DOI: [10.1128/AEM.01541-09](https://doi.org/10.1128/AEM.01541-09). URL: [http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2786419](http://www.ncbi.nlm.nih.gov/pubmed/19801464%20http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2786419).
- 705 [48] Jarmo Ritari et al. In: *BMC Genomics* 6.1 (Dec. 2011), p. 1056. ISSN: 1471-2164. DOI:  
706 [10.1186/s12864-015-2265-y](https://doi.org/10.1186/s12864-015-2265-y).
- 707 [49] Paolo Di Tommaso et al. “Nextflow: A tool for deploying reproducible computational  
708 pipelines”. In: *F1000Research* 4 (July 2015). DOI: [10.7490/F1000RESEARCH.1110183.1](https://doi.org/10.7490/F1000RESEARCH.1110183.1).  
709 URL: <https://f1000research.com/posters/4-430>.
- 710 [50] Stephen C Watts et al. “FastSpar: rapid and scalable correlation estimation for compositional  
711 data”. In: *Bioinformatics* (Aug. 2018). Ed. by Oliver Stegle. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/bty734](https://doi.org/10.1093/bioinformatics/bty734). URL: <https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty734/5086389>.

## Tables and Figures

Process	Tool	Parameters
Denoising and Clustering	Dada2/Deblur	default
Taxonomy Assignment	NCBI with Blast	RefSeq database
OTU Processing	Based on statistical power	Dynamic cutoff
Network Inference	Consensus method	-

Table 1: Default tools and parameters for the pipeline



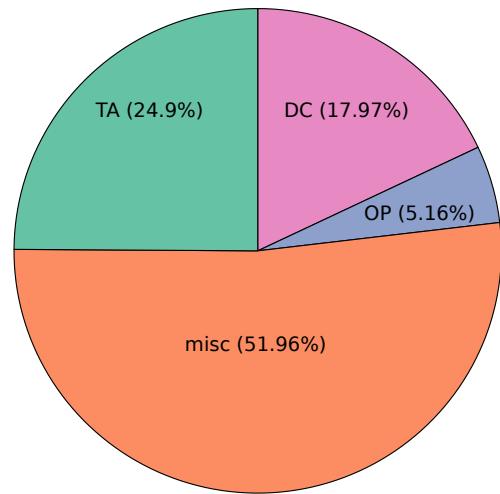
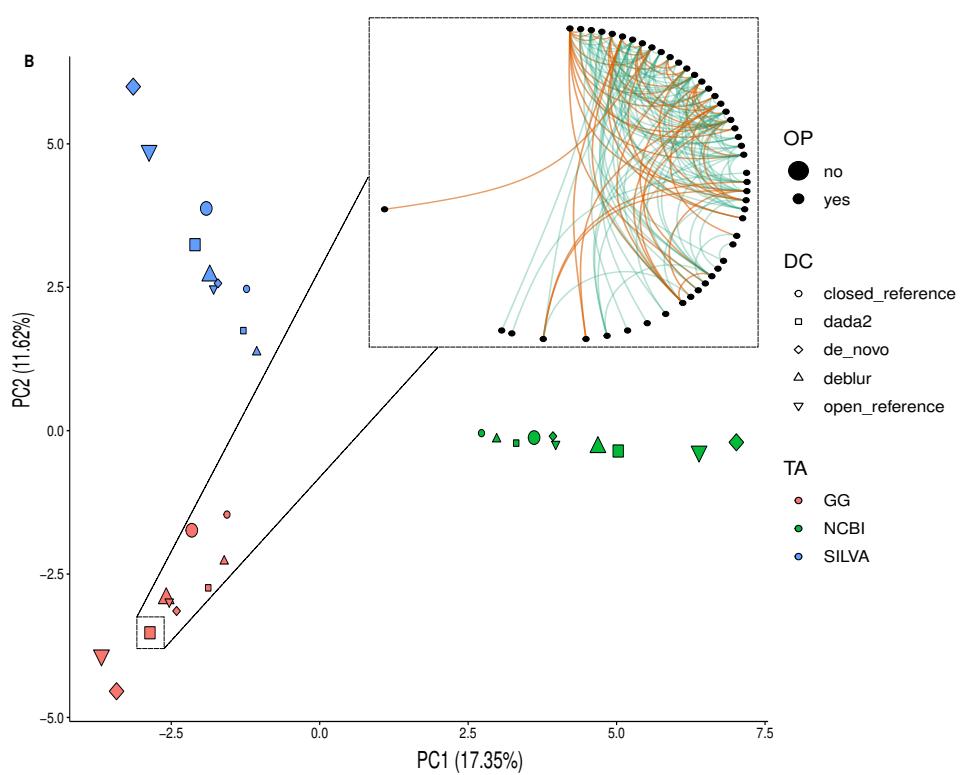
**Figure 1: The workflow of the microbial co-occurrence analysis pipeline.** The steps can be grouped into four major groups: **(DC)** Denoising and Clustering, **(TA)** Taxonomy Assignment, **(OP)** OTU or ESV Processing, and **(NI)** Network Inference. Each step incorporates several processes, each of which in turn have several alternate algorithms for the same task (indicated by the text to the right of the blue boxes). The text along the arrows describes the data that is being passed from one step to another. For details on each process and data types, see Methods.

---

**A**

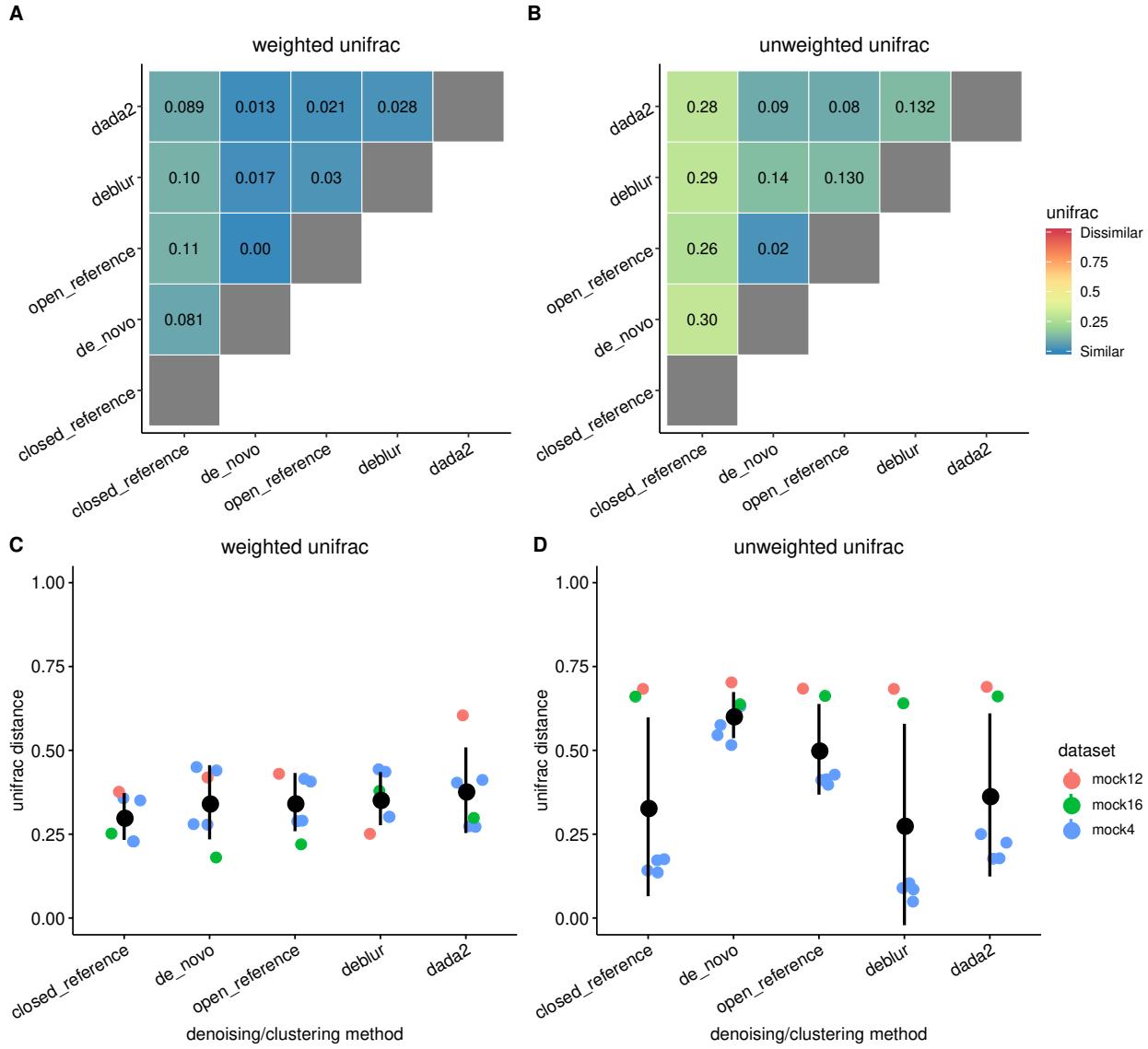
$$L_i = \begin{bmatrix} \text{edge}_{i,1} \\ \text{edge}_{i,2} \\ \vdots \\ \text{edge}_{i,n} \end{bmatrix}$$

$$L \sim DC + OP + TA$$

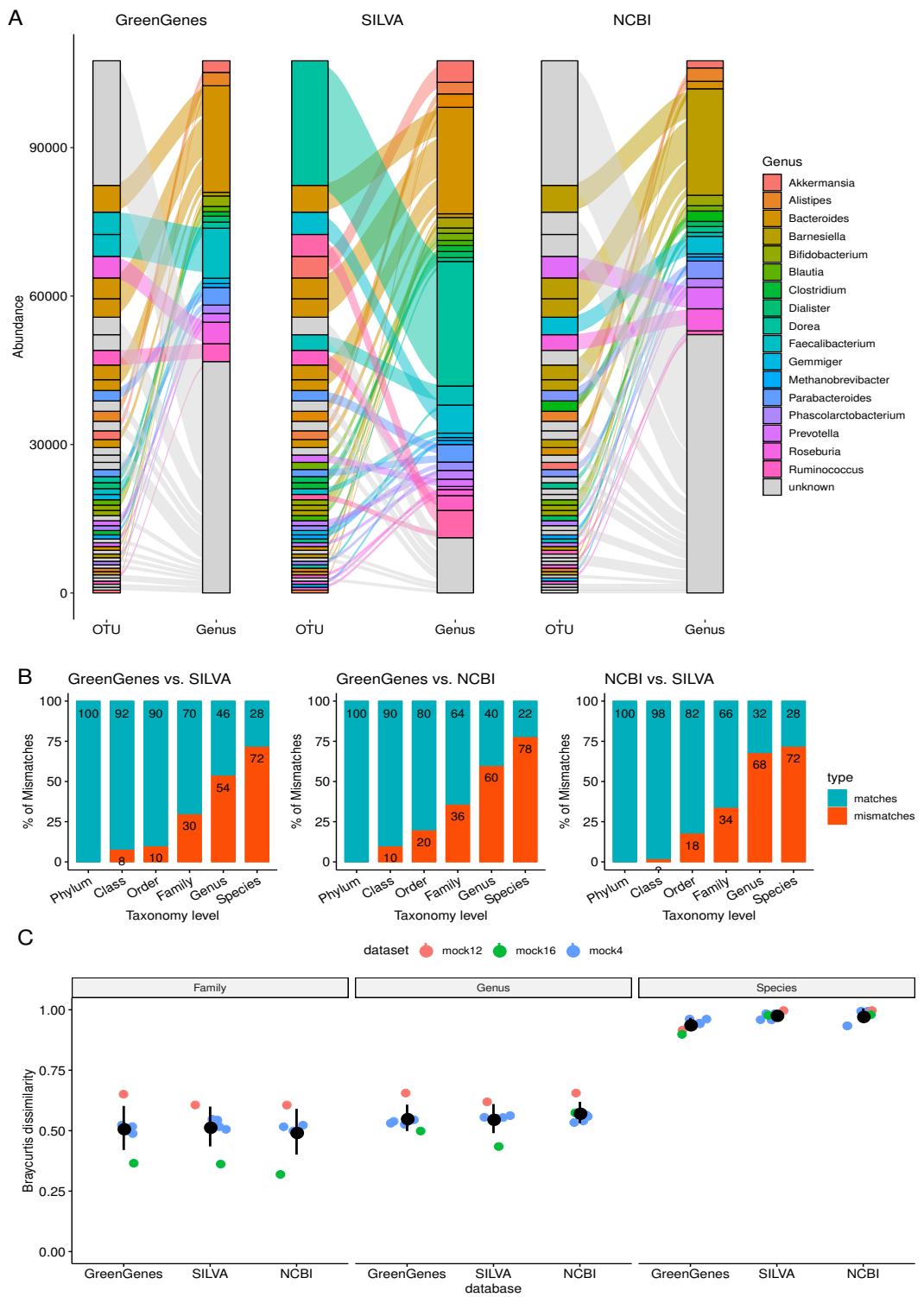
**B**

---

**Figure 2: The choice of database contributes to the most variance in the networks.** **(A)** The total relative variance in the networks contributed by the DC, TA and OP steps of the pipeline (right) and the linear model used to calculate the relative variance (left), see the Methods section for details. **(B)** All combinations of inferred networks are shown as points on a PCA plot. The color of the points corresponds to the taxonomy database, the shape corresponds to the denoising/clustering method and the size corresponds to whether low abundance OTUs were removed or not. **(B inset)** The network generated using DC=dada2, TA=GG, OP=no and NI=SPARCC and represents the particular point shown (big red square). The plot clearly shows that the points can be separated based on the TA step and that the differences due to the DC and OP steps are not as significant.



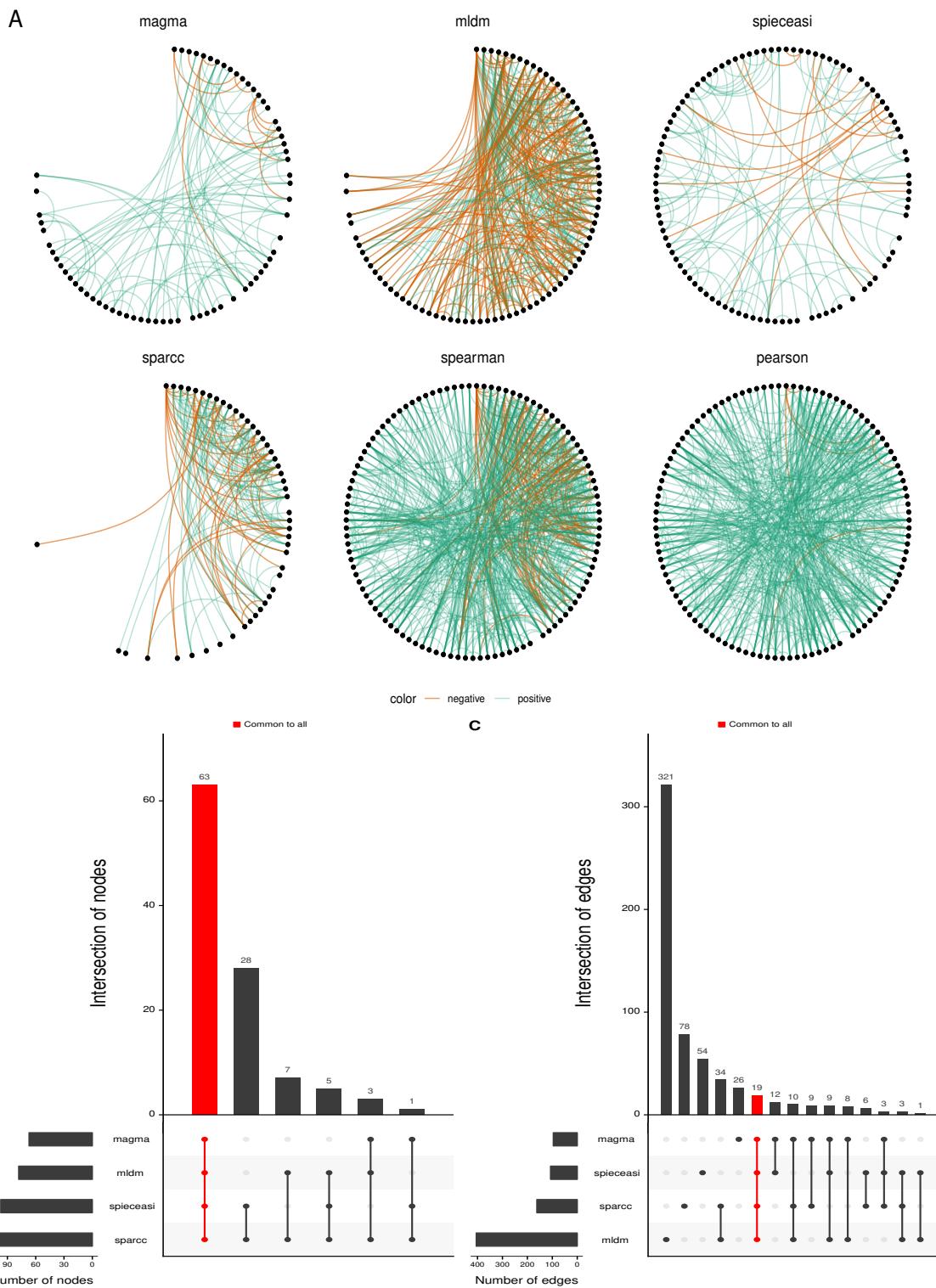
**Figure 3: The representative sequences generated by the different denoising/clustering methods are very similar but differ in the sequences that are in low abundance.** (A) The average weighted UniFrac distance between the representative sequences shows that the representative sequences and their compositions are fairly identical between the methods, (B) The relatively larger average unweighted UniFrac distance indicates that methods differ in their identification of sequences of low abundance, (C, D) The distributions of the average weighted UniFrac distance between the expected sequence profile and the calculated sequence profile in mock datasets. (D) The distributions of the average unweighted UniFrac distance show that dada2 and Debblur were the best performing methods in most of the datasets.



---

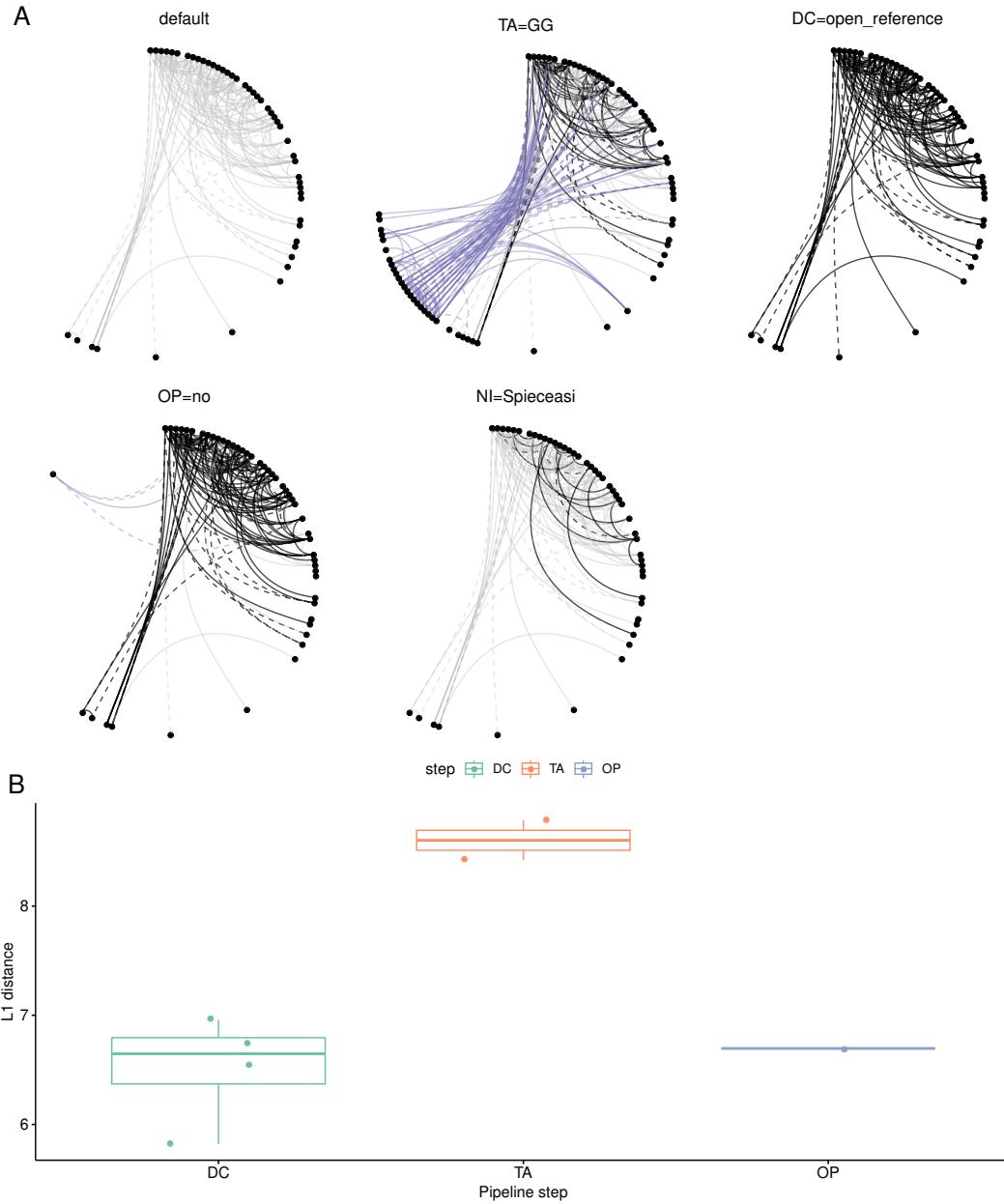
**Figure 4: Taxonomic reference databases vary widely in terms of their taxonomy assignments.**

**(A)** The assignment of the top 50 representative sequences to their respective taxonomies using the three different reference databases shows how the same sequences are assigned to different Genus. **(B)** The percentage of OTUs assigned to the same taxonomic label when using different reference databases. The percentage of mismatches decrease at higher taxonomic levels but even at the Phylum level there exists around 10% of mismatches. **(C)** The Bray-Curtis dissimilarity between the expected taxonomy profile and calculated taxonomy profile in the mock datasets shows that there is no singular best choice of database for every dataset.



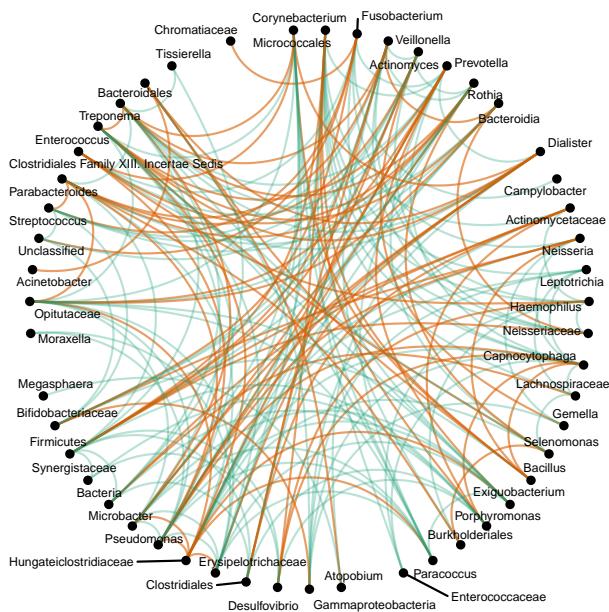
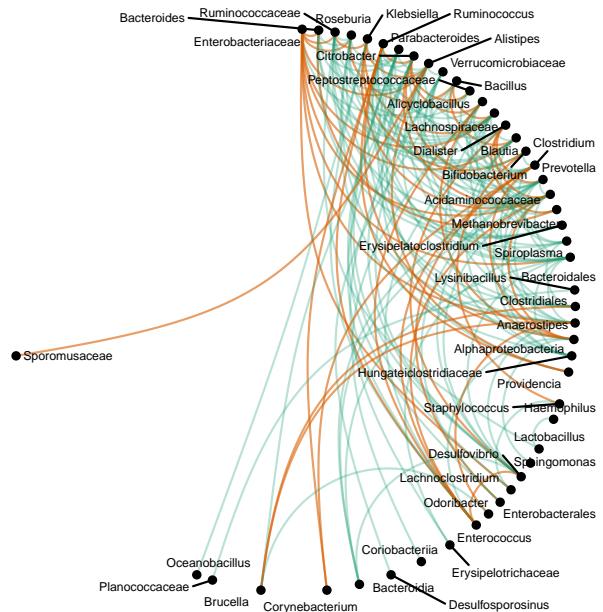
---

**Figure 5: Networks generated using different network inference methods show notable differences both in terms of edge-density and connectivity.** **(A)** The six different networks generated by the different network inference methods are very dissimilar. The green links are positive associations and the orange links are negative associations. A threshold of 0.3 was set for the methods that infer pairwise correlations (SparCC, Spearman, Pearson) and no threshold was set for the other methods. **(B)** The node overlap Upset plot [43] indicates that all the networks have a large number of common nodes involved in connections. Whereas, **(C)** The edge overlap Upset plot shows that a very small fraction of these connections are actually shared.



**Figure 6: Network inference and taxonomic assignment have the highest influence on the inferred network structures.** (A) The network constructed using the default pipeline parameters (DC=DADA2, TA=NCBI, OP=on, NI=SparCC) is compared with networks generated when one of the steps use a different tool. The common connections (common with the default network) are in black, connections unique to the network are colored purple and connections in the default network but not present in the current network are gray. (B) The L1 distance between the networks generated by changing one step of the default pipeline and the network generated using the default parameters.

---

**A****Hard Palate****B****Stool**

color — negative — positive

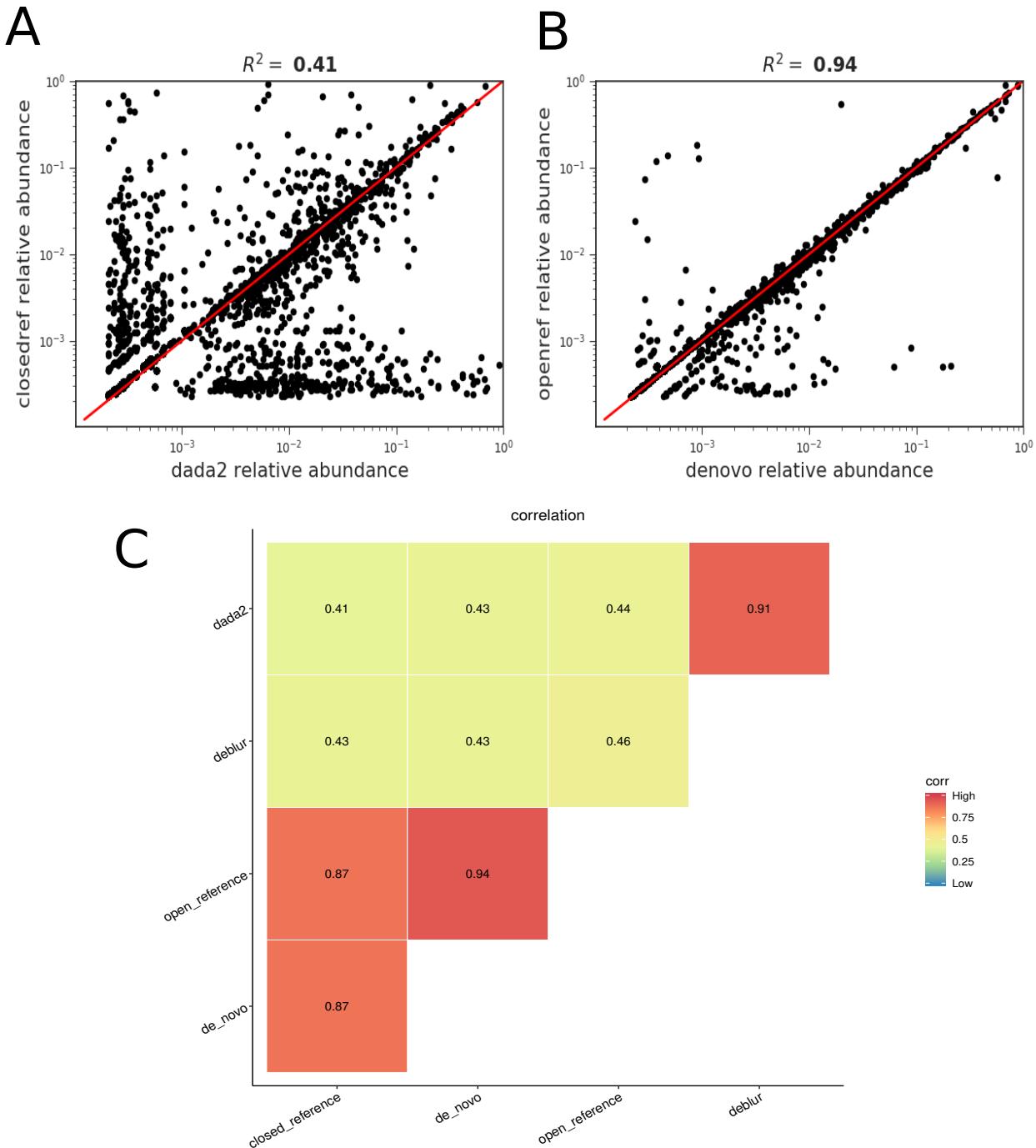
Figure 7: The consensus networks generated using the default pipeline settings. **(A)** Co-occurrence network of the Hard Palate microbiome generated from samples of healthy subjects in a periodontal diseases study. **(B)** Co-occurrence network of the Stool microbiome generated from samples of healthy subjects in a fecal microbiome transplant study.

---

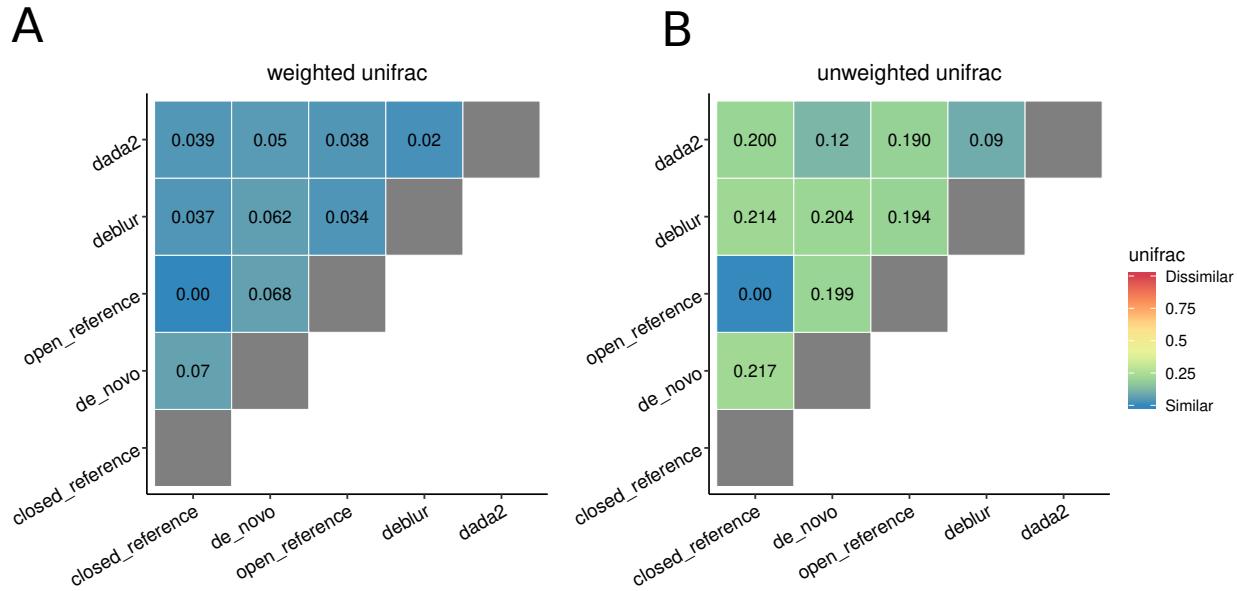
715 **Supplementary**

Step	Task	Tool	Parameter	Value
Sequence Processing	demultiplex_illumina	join_reads	min_overlap	6
			perc_max_diff	8
		demultiplex_454	rev_comp_barcodes	False
			rev_comp_mapping_barcodes	False
		trim_filter_fixed	-	-
	Chimera Checking		seq_sample_size	10,000
		uchime	ncpus	1
			trunc_q	2
			max_ee	2
		remove_bimera	-	-
Denosing and Clustering	de_novo	ncpus	-	1
		chimera_method	-	consensus
			enable_rev_strand_match	True
			suppress_de_novo_chimera_detection	True
		closed_reference	ncpus	1
	Denoise Cluster		enable_rev_strand_match	True
			suppress_de_novo_chimera_detection	True
		open_reference	ncpus	1
			reference_sequences	97_otus.fasta
			enable_rev_strand_match	True
Taxonomy Assignment	naive_bayes		suppress_de_novo_chimera_detection	True
		ncpus	-	1
		reference_sequences	-	97_otus.fasta
		picking_method	-	uclust
		dada2	ncpus	1
	blast	big_data	-	FALSE
		ncpus	-	1
		deblur	mind_reads	2
			min_size	2
			confidence	0.7
OTU/ESV Processing	Assign	mem_per_core	-	8G
		ncpus	-	1
		max_accepts	-	10
		perc_identity	-	0.8
		evalue	-	0.001
	Filter	min_consensus	-	0.51
		abundance	count_thres	500
			prevalence_thres	0.05
		group	abundance_thres	0.01
			tax_levels	[‘Phylum’, ‘Class’, ‘Order’, ‘Family’, ‘Genus’, ‘Species’]
Network Inference	partition	partition	-	-
			count_thres	500
			axis	sample
		Transform	prevalence_thres	0.05
		normalize	abundance_thres	0.01
	Export	rm_sparse_obs	rm_sparse_obs	True
		biom2tsv	rm_sparse_samples	True
			-	-
		bootstrap	bootstraps	1000
		Bootstrap	ncpus	1
Correlation	resample	filter_flag	filter_flag	True
		pvalue	ncpus	1
		sparcc	iterations	50
		pearson	ncpus	1
		spearman	-	-
	Correlation		method	mb
		spiceeasy	ncpus	1
			nreps	50
			nlambda	20
		42	lambda_min_ratio	1e-2
Network	mldm	z_mean	z_mean	1
		max_iteration	max_iteration	1500
	magma	-	-	-
	make_network	-	-	-

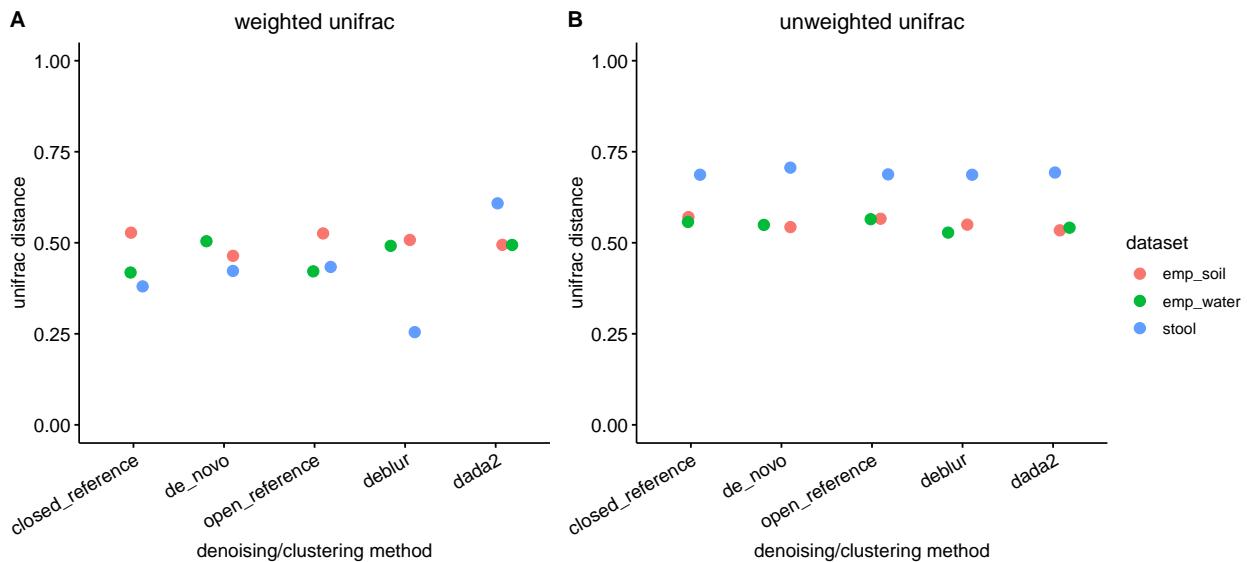
Table S1: The default parameters used in the various tools of the pipeline



**Figure S1: Comparison of various denoising and clustering algorithms used in the pipeline.** (A, B) Correlation of the abundances of the taxa that are in common between the count matrices created by two different methods. (A) The worst correlation (least similar methods) is between open-reference and dada2. (B) The best correlation (most similar methods) is between open-reference and denovo. (C) A heatmap showing the  $R^2$  of all pairwise comparisons of the methods.



**Figure S2: Heatmaps showing the weighted and unweighted unifrac distances for the hard palate dataset analysis.** (A) weighted unifrac distances and (B) unweighted unifrac distances between the representative sequences generated by different denoising and clustering algorithms. These results are in agreement with the stool microbiome dataset.



**Figure S3: The distributions of the average weighted UniFrac distance between the expected sequence profile and the calculated sequence profile in the synthetic datasets.** We observe no significant difference between the various methods on the synthetic datasets used for this study.

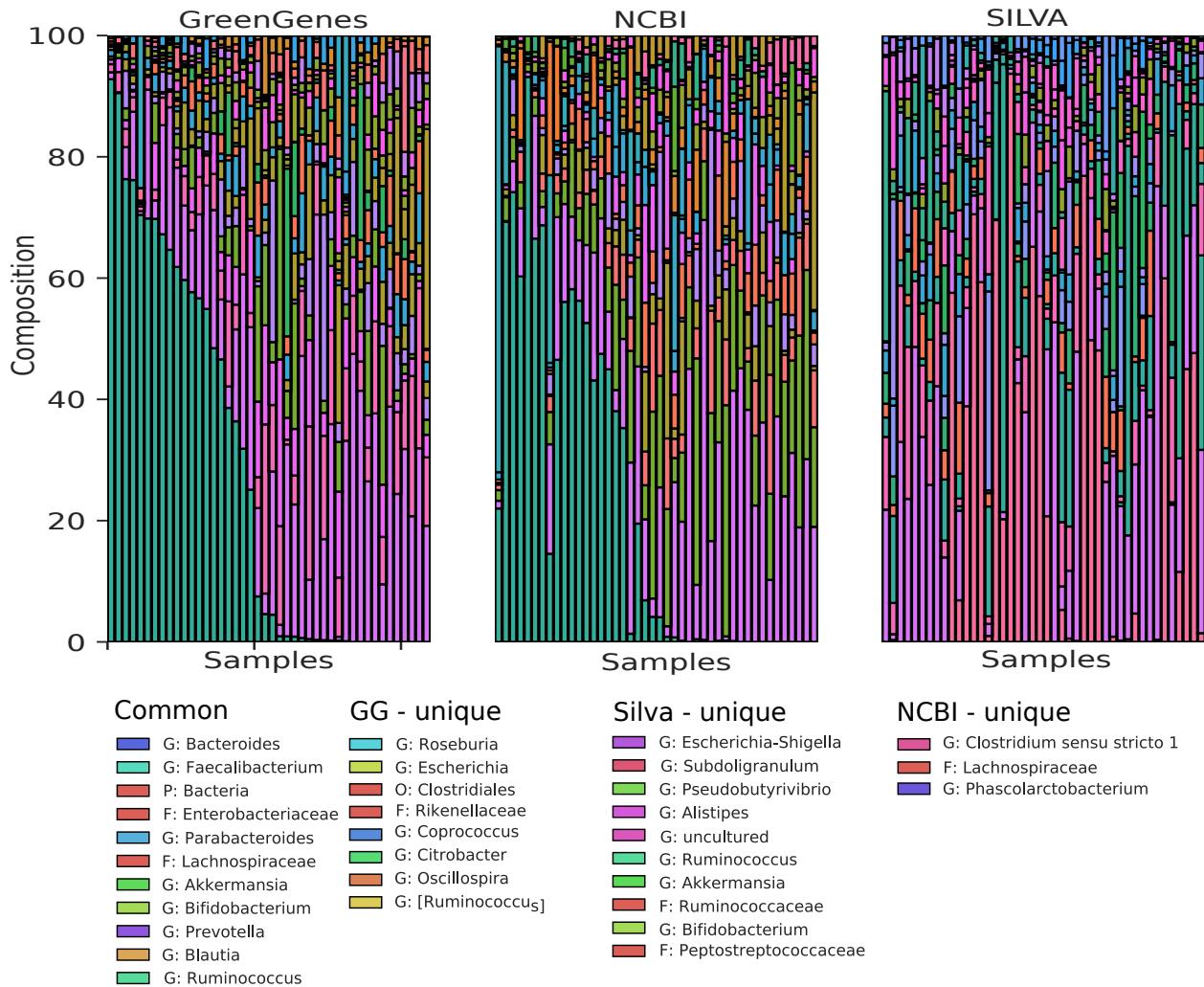


Figure S4: (A) Taxonomy composition of the 20 most abundant genera predicted for the stool microbiome dataset generated using different taxonomy references databases: Greengenes, SILVA and NCBI. The legend shows the common and the unique genera among the taxonomy assignments.

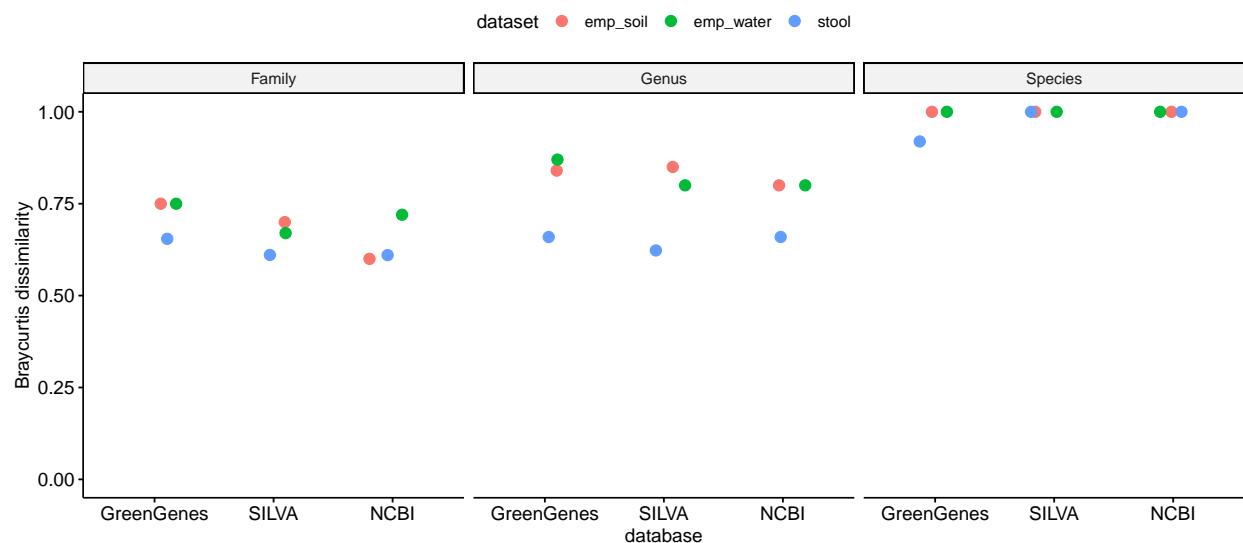


Figure S5: The bray-curtis dissimilarity between the expected taxonomic composition and generated taxonomic composition for the synthetic datasets.

---

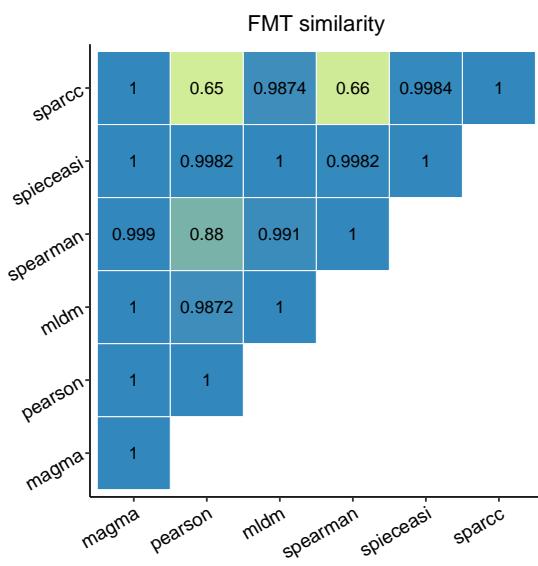
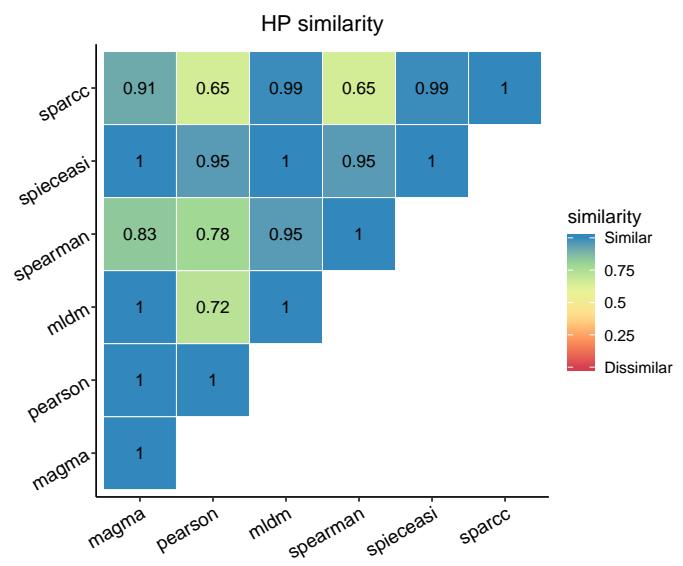
**A****B**

Figure S6: The similarity between the networks generated using the different network inference algorithms for stool dataset (A) and the hard palate dataset (B). The similarity between the various methods was found to vary with the dataset used.

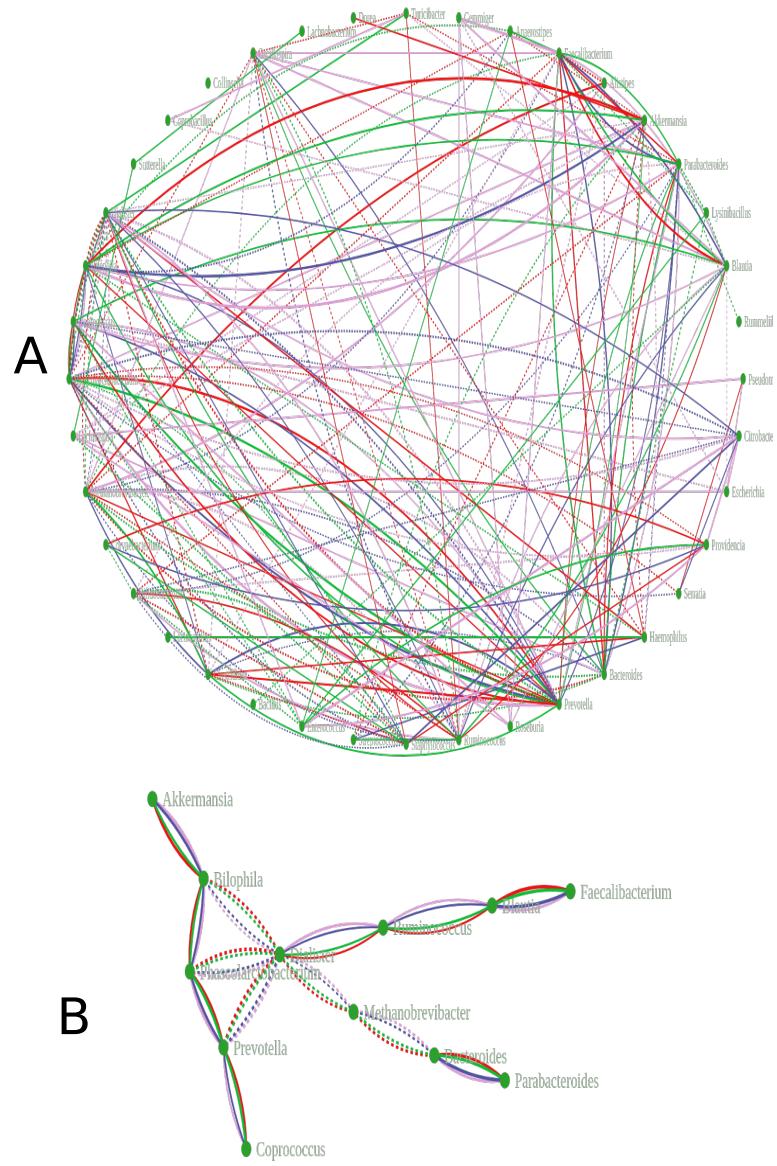


Figure S7: A network showing union (A) and intersection (B) of networks generated using different denoising and clustering tools on the Stool dataset.