

1 **Detection of phylogenetic core groups in diverse microbial**
2 **ecosystems**

3 Marcos Parras-Moltó¹, Daniel Aguirre de Cárcer^{1*}

4 ¹Departamento de Biología, Universidad Autónoma de Madrid, Madrid, Spain.

5

6 **Correspondence:**

7 Dr. Daniel Aguirre de Cárcer

8 daniel.aguirre@uam.es

9

10 **Keywords:**

11 Community assembly, microbiome, Phylogenetic clustering, 16S rRNA gene

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29 ABSTRACT

30 The detection and subsequent analysis of phylogenetic core groups (PCGs) in a
31 microbial ecosystem has been recently proposed as a potentially important analytical
32 framework with which to increase our understanding of its structure and function.
33 However, it was still unclear whether PCGs represented an infrequent phenomenon in
34 nature. Here we provide evidence of PCGs in a large and diverse array of environments,
35 which seems to indicate that their existence is indeed a predominant feature of microbial
36 ecosystems. Moreover, we offer dedicated scripts to examine the presence and
37 characteristics of PCGs in other microbial community datasets.

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60 **Background.** It is nowadays commonly believed that microbial communities assemble
61 on the basis of function alone. This idea is supported by the predominant observation
62 that different community compositions can translate into functionally-equivalent
63 microbial ecosystems. In this model, multiple unrelated populations would be
64 functionally redundant [1] in a particular microbial ecosystem type. However, this idea
65 is somewhat challenged by the extended phenomenon of phylogenetic clustering in
66 microbial communities; the tendency of bacteria to co-occur with phylogenetically
67 related populations more often than expected by chance alone [2, 3].

68 Phylogenetic clustering in a microbial ecosystem can be studied in terms of
69 phylogenetic core groups (PCGs), representing discrete portions of the bacterial
70 phylogeny present in all instances of a given ecosystem type. PCGs have been detected
71 so far in the rice rhizosphere [4] and human gut (fecal) [5] environments. The existence
72 of a PCG in a particular microbial ecosystem type has been theorized to be linked to
73 selection based on a combination of biotic and abiotic factors characteristic of that
74 ecosystem, and the existence in populations belonging to that PCG of a
75 phylogenetically-conserved set of traits allowing them to surpass such selection [4].
76 Thus, the study of PCGs in a given ecosystem could help understand the selection
77 forces at play in the ecosystem, and thus illuminate overall community assembly and
78 function.

79 It is still unclear whether PCGs are a predominant feature of microbial ecosystems or a
80 rare phenomenon. Thus, to test these possibilities we evaluate here the existence of
81 PCGs in a wide array of diverse microbial ecosystems. Also, so far PCGs had been
82 detected in terms of 16S sequence clusters of varying depth, which represents a
83 reasonable proxy. However, sequence clustering lacks true transitivity, which, jointly
84 with differential initial seeding between clustering runs, may translate into slightly
85 different clusters for the same input dataset generated by different runs or clustering
86 algorithms. Thus, here we analyze PCGs also on the basis of nodes in a phylogenetic
87 tree detected in all instances of the ecosystem type, an approach that also provides
88 increased phylogenetic resolution.

89 **Methods.** Here we analyze the existence of PCGs in nine different datasets from the
90 literature presenting a comparatively high number of ecosystem replicates and
91 sequencing depth (Table 1); The human microbiome is represented by datasets
92 *FlemishGut* [6] (fecal), *TwinsUK* [7] (fecal), *Illeum* [8] (mucosa), *Rectum* [8] (mucosa),
93 and *Vagina* [9] (mucosa). Plant-associated environments are represented by *Rice* (root
94 samples) [10] and *Leaf* [11], animal microbiomes by *Sponge* (*Carteriospongia*
95 *foliascens*) [12] and *Mice* [13], and environmental communities by *Wastewater* [14].
96 *Rice* was further subdivided by root environment (rhizosphere, rhizoplane, and
97 endosphere), *Mice* by origin (wild or lab), and *Vagina* in terms of previously reported
98 community types [9].

99 For each dataset, samples presenting very low sequence depths were removed, then all
100 samples were subsampled to a (minimum) common depth. Finally, the normalized
101 datasets were analyzed with *BacterialCore.py* (<https://git.io/Je5V3>). The script uses
102 various QIIME processes [15] and *R* libraries to detect PCGs and produce associated
103 analyses and statistics. It employs the clustering-based core detection approach
104 previously described [5], and a new approach based on a 16S rRNA gene phylogeny.
105 Here, the algorithm traverses the tree from leaves to root; if a leaf/node is present in a

106 (selected) percentage of samples it is flagged as “core”, and its abundance values
107 removed from all parental nodes before continuing, so that reported core groups are
108 non-overlapping. Additionally, *BacterialCore.py* provides per core-group information,
109 statistics, and consensus taxonomies.

110 **Results and discussion.** The microbial ecosystems analyzed presented a considerable
111 number of PCGs detected at different phylogenetic depths along the bacterial phylogeny
112 (Table 1, Figure 1, Suppl. Mat. 1, Suppl. Mat. 2). The exceptions to this pattern were
113 the mucosa environments analyzed (*Illeum*, *Rectum*, and *Vagina*) as well as the *Leaf*
114 ecosystem, featuring the presence of very few PCGs. This phenomenon could be
115 hypothesized to relate to the more homogeneous abiotic conditions of these
116 environments translating to less diverse communities. However, the *Rice rhizoplane* and
117 *endosphere* ecosystems, which could also be *a priori* considered as presenting more
118 homogeneous abiotic conditions, presented a large number of PCGs. The low number of
119 PCGs detected in the mucosal ecosystems could be related to their comparatively low
120 sequencing depth. Nevertheless, the leaf environment presents a substantial sequencing
121 depth, but only two PCGs.

122 Overall, the detected PCGs represented a preeminent fraction of the total community
123 (Table 1), with the lowest pooled abundance values being 18.5% (*Leaf*) and 34.9%
124 (*Illeum*), and the largest 77.6% (*Sponge*) and 93.4% (*Vagina*). In general, there was a
125 good correspondence between the clustering and tree-based approaches (Figure 1,
126 Supplementary Material 2), both of which produced correlated results in terms of
127 number of PCGs and their phylogenetic depth. Commonly, results for the clustering
128 approach represented a subset of those from the tree-based approach (Supplementary
129 Material 1; Venn diagrams)

130 In this brief report we have detected PCGs in terms of 16S sequence clusters and nodes
131 in a phylogenetic tree of different depths present in all samples from the same
132 ecosystem type. While this is a useful heuristic, other criteria such as a Poisson
133 distribution [16], a competitive lottery schema [17], invariance metrics [18], or the use
134 of neutral models [19, 20], could be employed and implemented within
135 *BacterialCore.py*.

136 **Conclusion.** The use of observed phylogenetic clustering patterns of community
137 assembly may represent an important clue to understand the assembly and function of a
138 microbial ecosystem. Here we provide evidence of PCGs in a large and diverse array of
139 environments, which seems to indicate that their existence is indeed a predominant
140 feature of microbial ecosystems. Moreover, we offer dedicated scripts to examine the
141 presence and characteristics of PCGs in other microbial community datasets.

142

143 **Availability of data:**

144 The datasets analyzed during the current study are available from their original sources.
145 Additional result files and scripts are available from the corresponding author upon
146 request.

147

148 **Acknowledgements**

149 This work was funded by the Spanish Ministry of Science and Innovation grant
150 BIO2016-80101-R.

151

152

153

154 REFERENCES

- 155 1. Adair KL, Douglas AE. Making a microbiome: the many determinants of host-
156 associated microbial community composition. *Curr Opin Microbiol.* 2017;35:23-9.
- 157 2. Stegen JC, Lin X, Konopka AE, Fredrickson JK. Stochastic and deterministic assembly
158 processes in subsurface microbial communities. *The Isme Journal.* 2012;6:1653.
- 159 3. Horner-Devine MC, Bohannan BJ. Phylogenetic clustering and overdispersion in
160 bacterial communities. *Ecology.* 2006;87:S100-8.
- 161 4. Aguirre de Cárcer D. A conceptual framework for the phylogenetically constrained
162 assembly of microbial communities. *Microbiome.* 2019;7:142.
- 163 5. Aguirre de Cárcer D. The human gut pan-microbiome presents a compositional core
164 formed by discrete phylogenetic units. *Scientific Reports.* 2018;8:14069.
- 165 6. Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K *et al.* Population-level
166 analysis of gut microbiome variation. *Science.* 2016;352:560-4.
- 167 7. Goodrich JK, Davenport ER, Beaumont M, Jackson MA, Knight R, Ober C *et al.* Genetic
168 Determinants of the Gut Microbiome in UK Twins. *Cell Host Microbe.* 2016;19:731-
169 43.
- 170 8. Gevers D, Kugathasan S, Denson LA, Vazquez-Baeza Y, Van Treuren W, Ren B *et al.*
171 The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe.*
172 2014;15:382-92.
- 173 9. Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SSK, McCulle SL *et al.* Vaginal
174 microbiome of reproductive-age women. *Proceedings of the National Academy of
175 Sciences.* 2011;108:4680-7.
- 176 10. Edwards J, Johnson C, Santos-Medellín C, Lurie E, Podishetty NK, Bhatnagar S *et al.*
177 Structure, variation, and assembly of the root-associated microbiomes of rice.
178 *Proceedings of the National Academy of Sciences.* 2015;112:E911.
- 179 11. Wagner MR, Lundberg DS, Del Rio TG, Tringe SG, Dangl JL, Mitchell-Olds T. Host
180 genotype and age shape the leaf and root microbiomes of a wild perennial plant. *Nat
181 Commun.* 2016;7.
- 182 12. Moitinho-Silva L, Nielsen S, Amir A, Gonzalez A, Ackermann GL, Cerrano C *et al.* The
183 sponge microbiome project. *Gigascience.* 2017;6:1-7.
- 184 13. Rosshart SP, Vassallo BG, Angeletti D, Hutchinson DS, Morgan AP, Takeda K *et al.*
185 Wild Mouse Gut Microbiota Promotes Host Fitness and Improves Disease Resistance.
186 *Cell.* 2017;171:1015-28.e13.
- 187 14. Saunders AM, Albertsen M, Vollertsen J, Nielsen PH. The activated sludge ecosystem
188 contains a core community of abundant organisms. *ISME J.* 2016;10:11-20.
- 189 15. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK *et al.*
190 QIIME allows analysis of high-throughput community sequencing data. *Nat Meth.*
191 2010;7:335-6.
- 192 16. Gumiere T, Meyer K, Burns A, Gumiere S, Bohannan B, Andreote F: A probabilistic
193 model to identify the core microbial community; 2018.
- 194 17. Verster AJ, Borenstein E. Competitive lottery-based assembly of selected clades in
195 the human gut microbiome. *Microbiome.* 2018;6:186.
- 196 18. Bradley PH, Pollard KS. Proteobacteria explain significant functional variability in the
197 human gut microbiome. *Microbiome.* 2017;5:017-0244.
- 198 19. Harris K, Parsons TL, Ijaz UZ, Lahti L, Holmes I, Quince C. Linking Statistical and
199 Ecological Theory: Hubbell's Unified Neutral Theory of Biodiversity as a
200 Hierarchical Dirichlet Process. *Proceedings of the IEEE.* 2017;105:516-29.
- 201 20. Adair KL, Wilson M, Bost A, Douglas AE. Microbial community assembly in wild
202 populations of the fruit fly *Drosophila melanogaster*. *The Isme Journal.* 2018;12:959-
203 72.

204

205 TABLES

Table 1. PCGs in diverse ecosystems

Dataset	Samples	Depth	Frequency	Core groups per clust. threshold
FlemishGut	873	8,383	0.50±0.13	97 ² , 95 ¹ , 93 ¹ , 92 ¹ , 90 ¹ , 89 ³ , 88 ³ , 86 ³ , 84 ¹ , 82 ¹ , 77 ¹
TwinsUK	2,727	14,082	0.37±0.11	90 ² , 89 ¹ , 88 ¹ , 86 ¹ , 83 ² , 81 ¹ , 77 ¹
Illeum	429	3,624	0.34±0.17	84 ¹ , 81 ¹
Rectum	304	3,763	0.42±0.18	86 ¹ , 79 ¹ , 75 ¹
Rice	372	16,884	0.43±0.20	97 [□] , 95 ¹ , 90 ¹ , 87 ¹ , 86 ¹ , 85 ¹ , 84 ¹ , 83 [□] , 82 [□] , 81 ³ , 80 ³ , 79 [□] , 78 ³ , 77 ³ , 76 ² , 75 ¹
Rice (Rhizosphere)	125	16,884	0.50±0.03	97 ^{□□} , 96 [□] , 95 [□] , 94 [□] , 93 [□] , 92 [□] , 91 [□] , 90 [□] , 89 ^{1°} , 88 ^{1□} , 87 ^{1□} , 86 ^{1□} , 85 [□] , 84 ^{1□} , 83 ¹² , 82 [□] , 81 [□] , 80 [□] , 79 ¹³ , 78 [□] , 77 [□] , 76 [□] , 75 ¹
Rice (Endosphere)	133	17,735	0.70±0.15	97 ^{1□} , 96 ¹ , 95 ¹ , 92 ¹ , 91 ¹ , 90 ¹ , 88 ¹ , 87 ² , 86 ³ , 85 ² , 84 ² , 83 ¹ , 82 [□] , 81 [□] , 80 ² , 79 ² , 78 ² , 77 ² , 76 ² , 75 ³
Rice (Rhizoplane)	114	17,821	0.54±0.09	97 ³³ , 96 [□] , 95 ³ , 94 ³ , 93 ³ , 92 [□] , 91 [□] , 90 [□] , 89 ² , 88 [□] , 87 ¹¹ , 86 [□] , 85 [□] , 84 ¹¹ , 83 ^{1□} , 82 [□] , 81 [□] , 80 [□] , 79 [□] , 78 [□] , 77 [□] , 76 ¹ , 75 ²
Vagina (#1-3,5)	286	693	0.93±0.11	82 ¹
Vagina (#1)	105	693	0.80±0.19	97 ¹
Vagina (#2)	25	774	0.79±0.17	97 ¹
Vagina (#3)	135	1,271	0.85±0.13	97 ¹
Vagina (#4)	108	881	NA	NA
Vagina (#5)	21	936	0.76±0.17	97 ¹ , 95 ¹
Wastewater	43	48,668	0.53±0.07	97 ^{□□} , 96 ² , 95 ² , 94 ¹ , 93 [□] , 92 [□] , 91 [□] , 90 [□] , 89 [□] , 88 [□] , 87 [□] , 86 ^{1□} , 85 [□] , 84 [□] , 83 ¹¹ , 82 ³ , 81 [□] , 8 ³ , 79 ³ , 78 [□] , 77 ³ , 76 [□] , 75 [□]
Sponge	143	27,921	0.77±0.10	97 ^{1°} , 96 ¹ , 93 ¹ , 92 ³ , 91 ¹ , 90 ³ , 89 ¹ , 88 [□] , 87 [□] , 86 [□] , 85 ¹ , 84 [□] , 83 ² , 81 ³ , 80 ³ , 78 ² , 76 ¹
Leaf	175	10,322	0.18±0.10	97 ¹ , 91 ¹
Mice (Total)	230	10,086	0.52±0.08	95 ¹ , 92 ¹ , 91 ² , 91 ¹ , 89 ² , 88 ³ , 87 ² , 86 ² , 85 [□] , 84 [□] , 83 [□] , 82 [□] , 81 ¹ , 80 ¹ , 79 ¹ , 78 ¹
Mice (Lab)	129	15,932	0.51±0.10	97 ²¹ , 96 ³ , 95 ² , 94 ² , 93 ² , 92 [□] , 91 ^{1°} , 90 ^{1°} , 89 [□] , 88 [□] , 87 [□] , 86 [□] , 85 ³ , 84 ³ , 83 [□] , 81 [□] , 80 ³ , 79 ¹ , 78 ¹
Mice (Wild)	101	10,086	0.57±0.07	97 ³ , 96 ¹ , 93 ¹ , 92 ² , 90 ² , 89 ¹ , 88 ² , 87 [□] , 86 [□] , 85 [□] , 84 ² , 83 ¹ , 82 ¹ , 81 ¹ , 79 ¹ , 78 ² , 77 ¹ , 75 ¹

Depth; sequences per sample. **Frequency**; average pooled abundance of members of the core OTUs across the dataset. **Core groups per clust. threshold**; Numbers represent similarity clustering thresholds (x10⁻²) were core OTUs were detected, and superscript values indicate the number of such OTUs observed for each threshold.

206

207

208

209

210 FIGURE LEGENDS

211 **Figure 1. Detection of PCGs in datasets.** Results for selected datasets based on the
212 dynamic clustering of 16S rRNA gene sequences from 97% to 75% sequence identity
213 (right to left) [OTUs] and the phylogenetic tree-based approach [Tree]. For each
214 threshold, OTUs/nodes present in all samples (i.e. core) appear vertically stacked with
215 individual heights representing average relative abundance of each core OTU/node in
216 the dataset. For the tree-based approach, x-axis values represent the maximum intra-
217 node distance, not the average.
218

219 SUPPLEMENTARY MATERIALS

220 **Supplementary Material 1.** *BacterialCore.py* result files (intermediate clustering
221 files have been omitted due to their large size).

222 **Supplementary Material 2.** Detection of PCGs in datasets. Results based on the
223 dynamic clustering of 16S rRNA gene sequences from 97% to 75% sequence identity
224 (right to left) [Page 1; OTUs], and the phylogenetic tree-based approach where x-axis
225 values represent the maximum intra-node distance [Page 2; MaxS] or the average intra-
226 node distance [Page3; MeanS]. For each threshold, OTUs/nodes present in all samples
227 (i.e. core) appear vertically stacked with individual heights representing average relative
228 abundance of each core OTU/node in the dataset. Results arising from both approaches
229 are also compared [Pages 4-5].
230

231

