



OPEN

Assessment of phylo-functional coherence along the bacterial phylogeny and taxonomy

Marcos Parras-Moltó & Daniel Aguirre de Cárcer✉

In this report we use available curated phylogenies, taxonomy, and genome annotations to assess the phylogenetic and gene content similarity associated with each different taxon and taxonomic rank. Subsequently, we employ the same data to assess the frontiers of functional coherence along the bacterial phylogeny. Our results show that within-group phylogenetic and gene content similarity of taxa in the same rank are not homogenous, and that these values show extensive overlap between ranks. Functional coherence along the 16S rRNA gene-based phylogeny was limited to 44 particular nodes presenting large variations in phylogenetic depth. For instance, the deep subtree affiliated to class Actinobacteria presented functional coherence, while the shallower family Enterobacteriaceae-affiliated subtree did not. On the other hand, functional coherence along the genome-based phylogeny delimited deep subtrees affiliated to phyla Actinobacteriota, Deinococcota, Chloroflexota, Firmicutes, and a subtree containing the rest of the bacterial phyla. The results presented here can be used to guide the exploration of results in many microbial ecology and evolution research scenarios. Moreover, we provide dedicated scripts and files that can be used to continue the exploration of functional coherence along the bacterial phylogeny employing different parameters or input data (<https://git.io/Jec5U>).

Our knowledge of microbial communities' composition has greatly expanded over the last decade thanks to the advent of high-throughput sequencing-based metagenomic approaches. The use of strategies based on the high-throughput sequencing of 16S rRNA gene amplicons is nowadays widespread in studies seeking to assess community structure. The use of such phylogenetic marker is often preferred over shot-gun metagenomic sequencing due to its associated reduced sequencing and computational costs. Most commonly, these studies aim not only to provide snapshots of community composition on different samples, but also to assess how these communities change along environmental gradients or sample groups, and significantly, why the observed changes take place in relation to the different putative functional roles of bacterial groups in the ecosystem.

Within these approaches, the taxonomic binning of reads based on the use of dedicated training sets remains the most common strategy in the description of community composition. The use of taxonomic ranks to organize the bacterial tree of life on the basis of evolutionary relationships has a colossal value in scientific communication, and it seems difficult to imagine the progress of microbial ecology without its use. However, the link between taxonomy, phylogeny, and ecological function is not clear-cut¹. Another common strategy in the analysis of 16S rRNA gene data is the use of sequence clusters obtained at predefined similarity thresholds (OTUs) or, less frequently yet more meaningfully, the search for study-oriented patterns along the complete 16S rRNA gene phylogeny (e.g.²).

It is nowadays most often accepted that, notwithstanding the widespread horizontal gene transfer and convergent evolution phenomena in the bacterial kingdom, many traits show conservatism along its phylogeny³. Over the last decade, different genome analyses have shown that closely related populations tend to share a higher percentage of traits than expected by chance⁴, a pattern also observed at high phylogenetic distances⁵. This link between phylogeny and shared function has also been substantiated by metabolic network analyses⁶ and in-depth literature reviews^{3,7}. More recently, Royalty and Steen⁸ studied the relative abundances of the 25 clusters of orthologous group functional categories⁹ within Park et al.'s genome-based improved taxonomy¹⁰ and showed that taxonomic rank explained 3.2%, 14.6%, 4.1%, 9.2%, 4.8% and 5.5% of the variance observed (for phylum, class, order, family and genus; respectively). The present study takes one step forward and analyzes the phylo-functional coherence of all taxa and taxonomic ranks. More importantly, we explore the frontiers along the bacterial phylogeny where gene content similarity and shared phylogeny cease to correlate; we consider that functional coherence persists if within-node gene content similarity is significantly higher than what could be

Departamento de Biología, Universidad Autónoma de Madrid, 28049 Madrid, Spain. ✉email: daniel.aguirre@uam.es

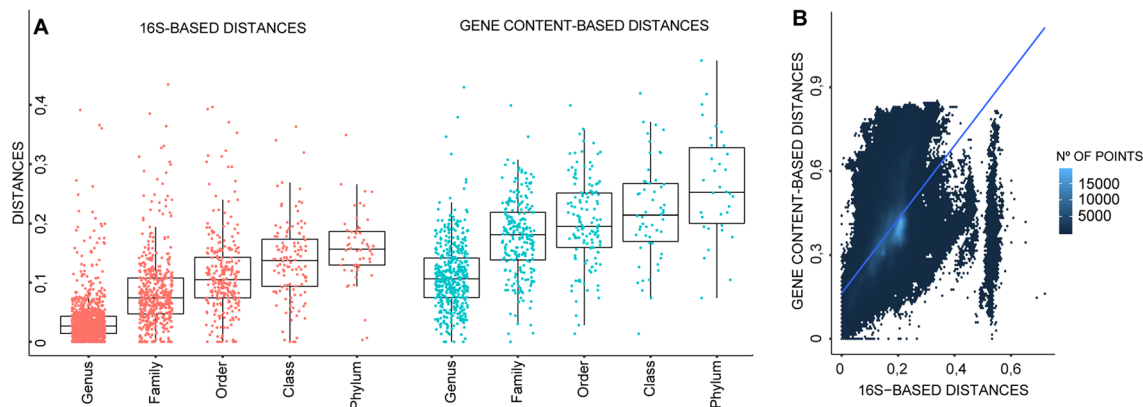


Figure 1. Panel (A) Box plots describing average within-group phylogenetic (16S rRNA gene sequence distance) and gene-content (pairwise Jaccard distance) distance values for each of the proposed taxa in all considered ranks. Panel (B) Correlation between 16S rRNA gene and gene content-based distances. Each point represents both distance values for each of the possible pairwise comparisons among the 6,989 dereplicated genomes with genome annotations ($y = 1.32x + 0.16$, $R^2 = 0.37$).

expected from a random draw of genomes along the phylogeny. Our main goal is to provide adequate guidance for future analyses employing 16S gene surveys aiming to improve our understanding of the role that Bacteria play in the ecosystem.

Results

Our results are based on available curated phylogenetic trees and 16S rRNA gene sequence data from the Genome Taxonomy Database (<http://gtdb.ecogenomic.org/downloads>), as well as genome annotations from proGenomes (<http://progenomes.embl.de/>). Due to uneven sampling, the number of available members per taxa varied wildly (Suppl. Figure 1, Suppl. Mat. 1). In all five ranks considered (genus, family, class, order, phylum) a large proportion of groups contained few members, while a small proportion presented a large number of members. This fact will likely impact the overall accuracy of classifiers using the available data as training set¹¹.

Phylogenetic similarity, assessed using 16S rRNA gene distances as proxy, was also quite variable among groups within the same taxonomic rank (Fig. 1A, Suppl. Mat. 1). As expected, the per-rank averages of intra-taxa distance averages followed an ascending trend from genus to phylum (Fig. 1A). While correlated, the observed 16S rRNA gene distance values are markedly smaller than the controversial cut-offs often employed in the field, 0.05, 0.10 and 0.20, for genus, family, and phylum, respectively¹². Another interesting trend relates to the coefficients of variation associated with the per-rank distributions of intra-taxa distance averages (from genus to phylum; 1.08, 0.73, 0.59, 0.48 and 0.36). The values seem to indicate that phylogenetic depth of taxa in a rank becomes more stable from genus to phylum. Equally anticipated, per-rank distributions of intra-taxa gene content distance averages also followed an ascending trend from genus to phylum (Fig. 1A, Suppl. Mat. 1).

When analyzing intra-taxa distance averages (Fig. 1A) and standard deviations (Suppl. Figure 2) it becomes clearer that taxonomic ranks are not consistent in terms of phylogenetic depth and gene content similarity. That is, for all taxonomic ranks, the different taxa within the same rank do not present a homogeneous level of either phylogenetic depth or gene content similarity. In this regard, the phylogenetic and gene content similarity values of many taxa were lower than those for taxa in superior ranks, and vice versa. This circumstance most likely stems from the fact that available genomes are not evenly distributed throughout the bacterial tree of life, that for operational and communication purposes the number of ranks is small, and that the main focus of modern taxonomy is related to the production of monophyletic taxa¹⁰.

Our initial results showed the expected overall correlation between gene content (pairwise Jaccard distances) and 16S-based distances (Fig. 1B). However, a noticeable discontinuity appeared in the pairwise 16S-based distance values at *ca.* 0.5 which, in turn, seems to mark the limit of the observed correlation. To better explore the frontiers of functional coherence along the bacterial phylogeny, we developed *FunCongr.R*; a script that compares within-node average gene content distance against what could be expected from a random draw of genomes along the phylogeny. If the average gene content distance between the genomes of a node is significantly smaller than expected from the null model, the node is deemed as functionally-coherent (see “Methods” section for a complete description of the procedure).

The analysis of the 16S rRNA gene phylogeny (Fig. 2, Suppl. Figure 3, Suppl. Mat. 2) returned 44 functionally-coherent nodes. These nodes present large variations in depth, showing within-node average 16S rRNA gene distance values ranging from 0.03 to 0.25 (average 0.12 ± 0.05). Without the goal of being exhaustive, a subsequent analysis based on the consensus (80%) taxonomy of leaves in each node, with the sole criterion of whether a particular taxon appears assigned to a single coherent node, revealed various deep-branching nodes affiliated to the Acidobacteriota (phylum), Alphaproteobacteria (class), Actinobacteria (class), Bacilli (class), Bacteroidia (class), Cyanobacteria (class), Fusobacteriales (order) Campylobacteriales (order), and a node (node2739) dominated by sequences affiliated to Bacilli (class) and Clostridia (class) (61% and 22%, respectively). On the other hand, the Gammaproteobacteria-affiliated region of the tree was partitioned among many different coherent

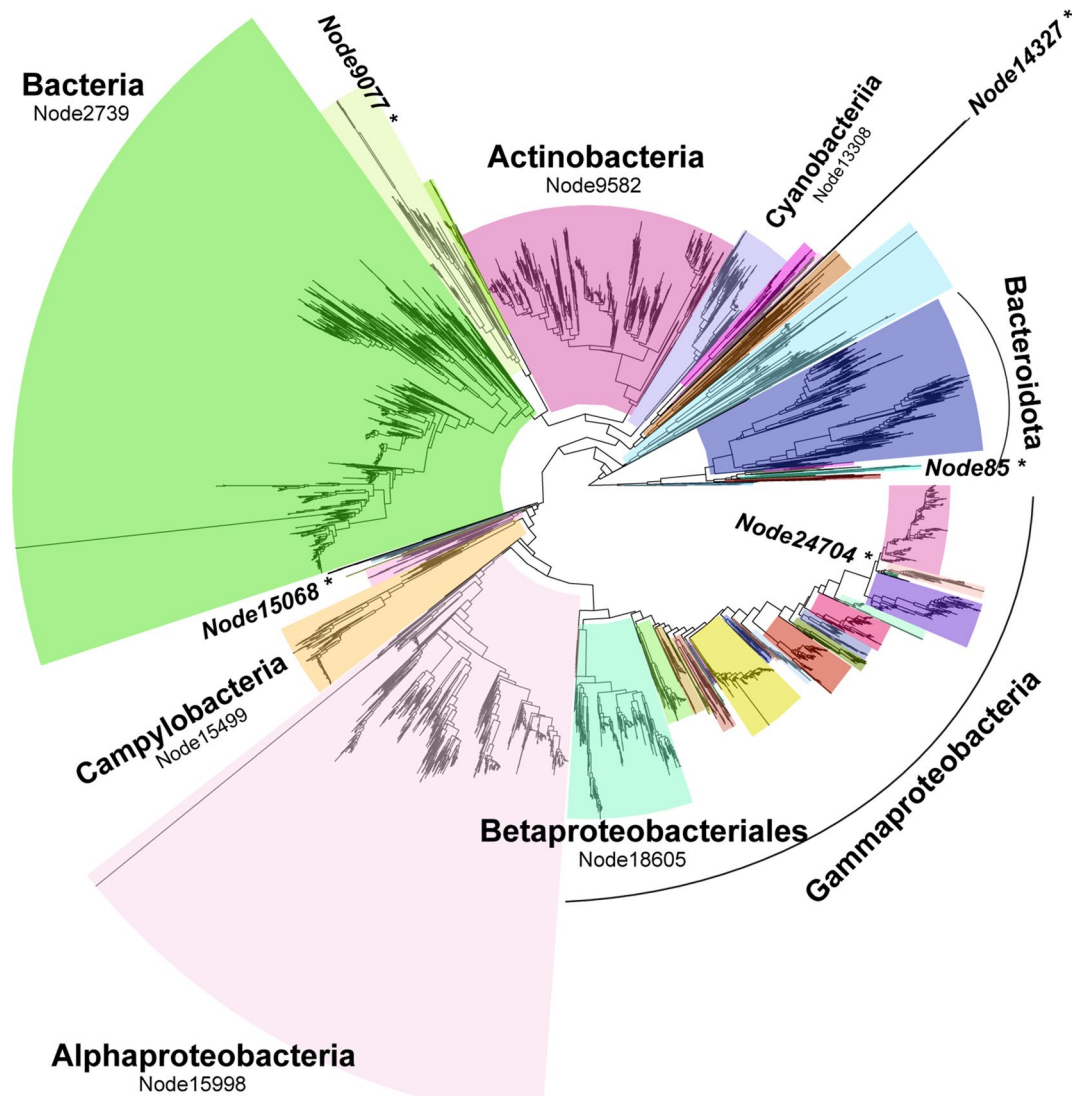


Figure 2. Frontiers of functional coherence along the 16S rRNA gene phylogeny. The tree represents Park et al.'s 16S rRNA gene-based phylogeny, pruned to contain only leaves with nearly full length 16S rRNA sequences and functional annotations. Colored boxes mark limits of phylo-functional coherence. Only the most prominent results are annotated (see Suppl. Figure 3 for an annotated and vertical version of the tree). *Nodes that failed to pass the functional coherence test.

nodes, hence indicating that such class is not functionally coherent in itself (when assessed within a 16S rRNA gene-based phylogeny). Here, the Betaproteobacteriales (order) stood out as a large and coherent node, while many different coherent nodes appeared affiliated to the Pseudomonadales (order), Enterobacteriales (order), and even Enterobacteriaceae (family). On the other hand, the analysis of the genome-based phylogeny (Suppl. Figure 4–5, Suppl. Mat. 2), inferred from a concatenated alignment of 120 ubiquitous single-copy proteins¹⁰, returned four functionally-coherent deep subtrees, respectively affiliated to phyla Actinobacteriota, Deinococcota, Chloroflexota, Firmicutes (node15643), and a fifth functionally-coherent deep subtree containing the rest of the bacterial phyla (node6). As mentioned earlier, the 16S rRNA gene provides low phylogenetic resolution at high distances, and thus phylogenies constructed on the basis of 16S rRNA gene distances may fail to recover true deep-branching evolutionary relationships, which are better gauged by whole genome-based phylogenies such as Park et al.'s employed here. Hence, it was not at all unexpected to find a deeper phylo-functional coherence in the genome-based phylogeny than in the 16S rRNA gene-based tree.

Discussion

While bacterial taxonomy and phylogeny resources, as well as related bioinformatic tools, continue to improve, the question remains as to how they should best be employed in studies using 16S rRNA gene surveys to assess bacteria-ecosystem relationships, a widespread approach. With regards to the use of taxonomic binning, our results show that within-group phylogenetic and gene content similarity of taxa in the same rank are not homogeneous, and that these values show extensive overlap between ranks. Similarly, Royalty and Steen pinpointed a

lack of homogeneity in within-group genome diversity for the different bacterial phyla⁸. Thus, we believe that most taxonomic ranks should be considered when assessing bacteria-ecosystem relationships, since it is not possible to know a priori which sections and depths of the bacterial phylogeny will be affected by the studied environmental gradient or sample group. The possible exception is phylum, since the 16S rRNA gene presents low phylogenetic resolution at the deeper and shallowest taxonomic ranks¹³. While our results were produced using GTDB's taxonomy, we expect the same patterns to arise if employing other classifications, and thus we argue that our recommendation above should also be followed if using classifications other than GTDB, such as SILVA¹⁴.

With regards to the use of phylogeny-based approaches, the 16S rRNA gene-based phylogeny showed remarkable correlation with gene content similarity. However, the breaking points of such relationship appeared at different depths along each different branching path of the phylogeny. Thus, single, arbitrary tree topology or sequence distance thresholds should not be employed regarding phylo-functional coherence. For instance, while a deep Actinobacteria (class)-affiliated node passed our functional coherence test, the Enterobacteriaceae (family)-affiliated shallower subtree did not.

The strategy employed could be biased by low quality genome annotations, which prompted us to employ the curated proGenomes database. In this sense, the depth of the frontiers of functional coherence described here could be understood as conservative estimates; for instance within the Bacteroidota (phylum)-affiliated subtree of the 16S-based phylogeny, node85 failed to pass the functional coherence test, and thus the subtree was split into several coherent units (Fig. 2, Suppl. Figure 3). While reducing the p-value of the permutation test to 0.01 returned the exact same tree topology, node85 contains only five genomes, and thus missannotation could potentially have artificially reduced the limits of phylo-functional coherence in this subtree. On the other hand, within that same phylogeny the loss of functional coherence within the Enterobacteriaceae (family)-affiliated subtree, and hence that of the larger Gammaproteobacteria (class)-affiliated subtree as well (Fig. 2, Suppl. Figure 3), was supported by the analysis of 35 genome annotations (node24704).

In addition to our previous results-driven recommendation that all ranks from genus to class/phylum be employed if using taxonomic binning, we argue that our functional coherence proxy (i.e. more within-node gene content similarity than expected by chance, and no descendant nodes failing to pass such test) is intuitive and useful in the exploration of phylo-functional coherence along the bacterial phylogeny. Furthermore, the existence of limits to phylo-functional coherence, and the fact that these limits vary in depth along the bacterial phylogeny, should be considered, for instance, in the evaluation of phylogenetic factors driving observed patterns of microbial community composition. The goal here is to identify the clades that respond to experimental or environmental variables, and thus it would be important to assess whether such clades present functional coherence (e.g.^{15,16}). The same holds true for studies of microbial community assembly from a phylogenetic perspective (e.g.^{17–19}) since observed patterns of phylogenetic signal are presumed to arise from phylo-functional coherence²⁰.

Thus, the results presented here can be used by the community, together with Park et al.'s provided phylogenetic trees, to obtain more meaningful results in many microbial ecology and evolution research scenarios. Moreover, the associated scripts and files are freely available (<https://git.io/jec5U>) so that (e.g.) different phylogenies, genome annotations, or significance thresholds can be used in the assessment of phylo-functional coherence.

Methods

Phylo-functional coherence of taxa and ranks. First, the 16S rRNA gene sequences for the representative (dereplicated) genomes employed to produce Parks et al.'s taxonomy (bac_ssu_r86.1.fna) were obtained from their repository (<http://gtdb.ecogenomic.org/downloads>)¹⁰. Then, *Mothur*²¹ commands were used to align the sequences against SILVA reference alignment²², and trimmed to contain only the sequence delimited by universal primers 27f and 1492R²³. Sequences not spanning the region were removed, leaving a total of 15,186 sequences. Taxa not including at least three members were not considered for subsequent analyses. Distances between the remaining 16S rRNA gene sequences were obtained using the *dist.seqs* command in *Mothur*, and later analyzed on the basis of their taxonomic affiliation (bac_metadata_r86.tsv), producing summary statistics for all taxa at all ranks.

NCBI's TaxIDs for each of the representative genomes (included within bac_metadata_r86.tsv) were linked to the gene content information obtained from the proGenomes resource²⁴ (complete annotation table kindly provided by Daniel Mende), resulting in 6989 different annotations for representative genomes presenting nearly full length 16S rRNA gene sequences. Only gene annotations with COG or NOG in their identifier were employed for consistency. Pairwise Jaccard distances between these genomes were derived from the shared presence of genes using the *R* package *vegan*²⁵, again producing within-taxa summary statistics for all taxa at all ranks.

Frontiers of functional coherence within phylogenetic trees. We developed a dedicated R script, *FunCongr.R*, to explore the limits of functional coherence along the bacterial phylogeny. Our approach employs phylogenetic trees (gtb_r86.ssu.bacteria.fasttree.tree and bac120_r86.1.tree¹⁰) and pairwise Jaccard distances between genomes obtained as mentioned above, although any other suitable tree or pairwise metrics could be used. The trees were initially processed to include constant node labels, and pruned to remove leaves without functional annotation.

The script traverses the tree from leaves to root (i.e. nodes are evaluated only if its descendant nodes have already been evaluated). Only nodes with sufficient information are analyzed (we arbitrarily fixed the minimum at 5 genomes). While traversing the tree, each node is flagged as non-coherent if it either (1) fails to pass a null model test, or (2) presents descendant nodes that did not pass the test previously. The test compares the average value of pairwise Jaccard distances between the node's leaves to an empirical cumulative distribution (ECD) of 1000 average values of pairwise Jaccard distances between N random leaves along the phylogeny, where N is the

actual number of leaves of the node being evaluated. If the node's average Jaccard distance value is not within the lowest 5% of the ECD (or arbitrarily selected threshold), the node is flagged as non-coherent.

In this manner, we delimit functional coherence along the tree by pinpointing nodes that either fail to pass the test, or present descendant nodes that fail to pass the test. Thus, here functional coherence is presumed to persist if within-node average gene content distance is lower (i.e. more similar) than what could be expected from a random draw of genomes along the phylogeny (with $p < 0.05$). The results are finally processed by recording the last functionally-coherent node along each branching path, along with their summary statistics. All figures of this work were produced using *ggplot2* package in R²⁶.

The article was previously published as a preprint (<https://doi.org/10.1101/795914>).

Data availability

The datasets analyzed during the current study are available from their original source (as stated above) and at (<https://git.io/jec5U>).

Received: 22 January 2020; Accepted: 6 April 2021

Published online: 15 April 2021

References

- Kunin, V., Ahren, D., Goldovsky, L., Janssen, P. & Ouzounis, C. A. Measuring genome conservation across taxa: Divided strains and united kingdoms. *Nucleic Acids Res.* **33**, 616–621 (2005).
- Washburne, A. D. *et al.* Phylogenetic factorization of compositional data yields lineage-level associations in microbiome datasets. *PeerJ* **9**, e2969 (2017).
- Martiny, J. B., Jones, S. E., Lennon, J. T. & Martiny, A. C. Microbiomes in light of traits: A phylogenetic perspective. *Science* **350**, 9323 (2015).
- Tamames, J., Sánchez, P. D., Nickel, P. I. & Pedrós-Alió, C. Quantifying the relative importance of phylogeny and environmental preferences as drivers of gene content in prokaryotic microorganisms. *Front. Microbiol.* **7**, 433–433 (2016).
- Philippot, L. *et al.* The ecological coherence of high bacterial taxonomic ranks. *Nat. Rev. Microbiol.* **8**, 523–529 (2010).
- Borenstein, E., Kupiec, M., Feldman, M. W. & Ruppin, E. Large-scale reconstruction and phylogenetic analysis of metabolic environments. *Proc. Natl. Acad. Sci. U S A.* **105**, 14482–14487 (2008).
- Goberna, M. & Verdu, M. Predicting microbial traits with phylogenies. *ISME J.* **10**, 959–967 (2016).
- Royalty, T. M. & Steen, A. D. Quantitatively partitioning microbial genomic traits among taxonomic ranks across the microbial tree of life. *mSphere*. **4**, e00446-e419 (2019).
- Galperin, M. Y. *et al.* COG database update: Focus on microbial diversity, model organisms, and widespread pathogens. *Nucleic Acids Res.* **49**, D274–d281 (2021).
- Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **27**, 2 (2018).
- Beiko, R. G. Microbial malaise: How can we classify the microbiome?. *Trends Microbiol.* **23**, 671–679 (2015).
- Schloss, P. D. & Handelsman, J. Status of the microbial census. *Microbiol. Mol. Biol. Rev.* **68**, 686 (2004).
- Janda, J. M. & Abbott, S. L. 16S rRNA gene sequencing for bacterial identification in the diagnostic laboratory: Pluses, perils, and pitfalls. *J. Clin. Microbiol.* **45**, 2761–2764 (2007).
- Quast, C. *et al.* The SILVA ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids Res.* **41**, D590–D596 (2013).
- Washburne, A. D. *et al.* Phylogenetic factorization of compositional data yields lineage-level associations in microbiome datasets. *PeerJ* **5**, e2969 (2017).
- Isobe, K., Allison, S. D., Khalili, B., Martiny, A. C. & Martiny, J. B. H. Phylogenetic conservation of bacterial responses to soil nitrogen addition across continents. *Nat. Commun.* **10**, 2499 (2019).
- de Aguirre Cárcer, D. The human gut pan-microbiome presents a compositional core formed by discrete phylogenetic units. *Sci. Rep.* **8**, 14069 (2018).
- Parras-Moltó, M. & de Cárcer, D. A. Detection of phylogenetic core groups in diverse microbial ecosystems. *bioRxiv*. 2020.2001.2007.896985 (2020).
- Darcy, J. L. *et al.* A phylogenetic model for the recruitment of species into microbial communities and application to studies of the human microbiome. *ISME J.* **14**, 1359–1368 (2020).
- de Aguirre Cárcer, D. A conceptual framework for the phylogenetically constrained assembly of microbial communities. *Microbiome*. **7**, 142 (2019).
- Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009).
- Yilmaz, P. *et al.* The SILVA and “all-species living tree project (LTP)” taxonomic frameworks. *Nucleic Acids Res.* **42**, D643–D648 (2014).
- Weisburg, W. G., Barns, S. M., Pelletier, D. A. & Lane, D. J. 16S ribosomal DNA amplification for phylogenetic study. *J. Bact.* **173**, 697–703 (1991).
- Mende, D. R. *et al.* proGenomes: A resource for consistent functional and taxonomic annotations of prokaryotic genomes. *Nucleic Acids Res.* **45**, D529–D534 (2017).
- Oksanen, J. *et al.* Vegan: Community Ecology Package. R package version 2.0–2. (2012).
- Wickham, H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag, 2016).

Acknowledgements

This work was funded by the Spanish Ministry of Science and Innovation grant BIO2016-80101-R and PID2019-108797RB-I00. The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

Author contributions

D.A. conceived the study, M.P. produced the scripts and carried out the analyses, and D.A. wrote the manuscript with input from M.P.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-87909-1>.

Correspondence and requests for materials should be addressed to D.A.d.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021