

Databases and ontologies

PICRUSt2-SC: an update to the reference database used for functional prediction within PICRUSt2

Robyn J. Wright^{1,*} , Morgan G.I. Langille¹ 

¹Department of Pharmacology, Faculty of Medicine, Dalhousie University, Halifax, Nova Scotia, B3H 4R2, Canada

*Corresponding author: Department of Pharmacology, Faculty of Medicine, Dalhousie University, 5849 University Avenue, Halifax, Nova Scotia, B3H 4R2, Canada. Email: robyn.wright@dal.ca

Associate Editor: Peter Robinson

Abstract

Summary: PICRUSt2 is a bioinformatic tool that predicts microbial functions in amplicon sequencing data using a database of annotated reference genomes. We have constructed an updated database for PICRUSt2 that has substantially increased the number of bacterial (19,493 to 26,868) and archaeal (406 to 1,002) genomes as well as the number of functional annotations present. The previous PICRUSt2 database relied on many timely and computationally intensive manual processes that made it difficult to update. We constructed a new streamlined process to allow regular upgrades to the PICRUSt2 database on an ongoing basis, and used this process to create a new database, PICRUSt2-SC (Sugar-Coated). Additionally, we have shown that this updated database contains genomes that more closely match study sequences from a range of different environments. The genomes contained in the database therefore better represent these environments and this leads to an improvement in the predicted functional annotations obtained from PICRUSt2.

Availability and implementation: PICRUSt2 source code is freely available at <https://github.com/picrust/picrust2> and at <https://anaconda.org/bioconda/picrust2>. The latest version of PICRUSt2 at the time of writing is also archived: <https://doi.org/10.5281/zenodo.15119781>. The PICRUSt2-SC database comes pre-installed with PICRUSt2 from version 2.6.0 onwards. Step-by-step instructions for making the updated database are at <https://github.com/picrust/picrust2/wiki/Updating-the-PICRUSt2-database>. All code used for the analyses and figures in this manuscript is at https://github.com/R-Wright-1/PICRUSt2-SC_application_note and <https://doi.org/10.5281/zenodo.15119770>.

1 Introduction

PICRUSt2 (Phylogenetic Investigation of Communities by Reconstruction of Unobserved States) is a widely used bioinformatic tool that allows the prediction of microbial community functions based on amplicon sequencing data. The first version of PICRUSt (PICRUSt1) was published in 2013 (Langille *et al.* 2013), and the second version (PICRUSt2) was published in 2020 (Douglas *et al.* 2020). In PICRUSt2, 16S rRNA gene sequences collapsed to amplicon sequence variants (ASVs) are placed into a reference tree using multiple-sequence alignment and then hidden-state prediction is used to infer the gene content of study sequences based on the gene content of the taxa used to construct the reference tree. Information on their predicted 16S rRNA gene copy number and their abundance in the original samples is used to weight the predicted genome annotations to generate a predicted metagenome.

One of the largest limitations to PICRUSt2 currently is that the default database uses functional annotations acquired from the Integrated Microbial Genomes database (Markowitz *et al.* 2012) in 2017. The number of genomes and functions has continued to expand during this time, and timely updates to PICRUSt2 are essential for its accuracy and functional comprehensiveness. However, the previous approach for updating the default database and functions was difficult due to a reference tree needing to be re-built with

newly added genomes. We have constructed a new database for PICRUSt2, PICRUSt2-SC (Sugar-Coated), that uses Genome Taxonomy Database (GTDB) r214 genomes (Parks *et al.* 2018, 2020, 2022; Rinke *et al.* 2021) and the bacterial and archaeal phylogenetic trees released with GTDB. We have annotated these genomes using EggNOG (Huerta-Cepas *et al.* 2019; Cantalapiedra *et al.* 2021), and as the GTDB genomes are readily available for download, we have provided instructions for PICRUSt2 users to add annotations to this updated database.

2 PICRUSt2-SC database creation

All genomes ($n = 402,709$) from the GTDB r214.1 were downloaded from the GTDB repository. Only representative genomes that were present in the GTDB bacterial or archaeal phylogenetic trees, had a high quality 16S rRNA gene, and that were $\leq 10\%$ contamination and $\geq 90\%$ completion were used, giving $n = 27,870$ total genomes ($n = 26,868$ and $n = 1,002$ genomes for bacteria and archaea, respectively). Full details of the methods used for 16S rRNA gene identification and processing are in [Supplementary Information section 2.1](#). This has increased the number of taxa present at every phylogenetic rank from phylum to species (Fig. 1), with approximately three-fold or four-fold more species present in the PICRUSt2-SC database for bacteria or archaea, respectively, compared with the previous PICRUSt2 database

Received: 21 February 2025; Revised: 8 April 2025; Editorial Decision: 23 April 2025; Accepted: 25 April 2025

© The Author(s) 2025. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

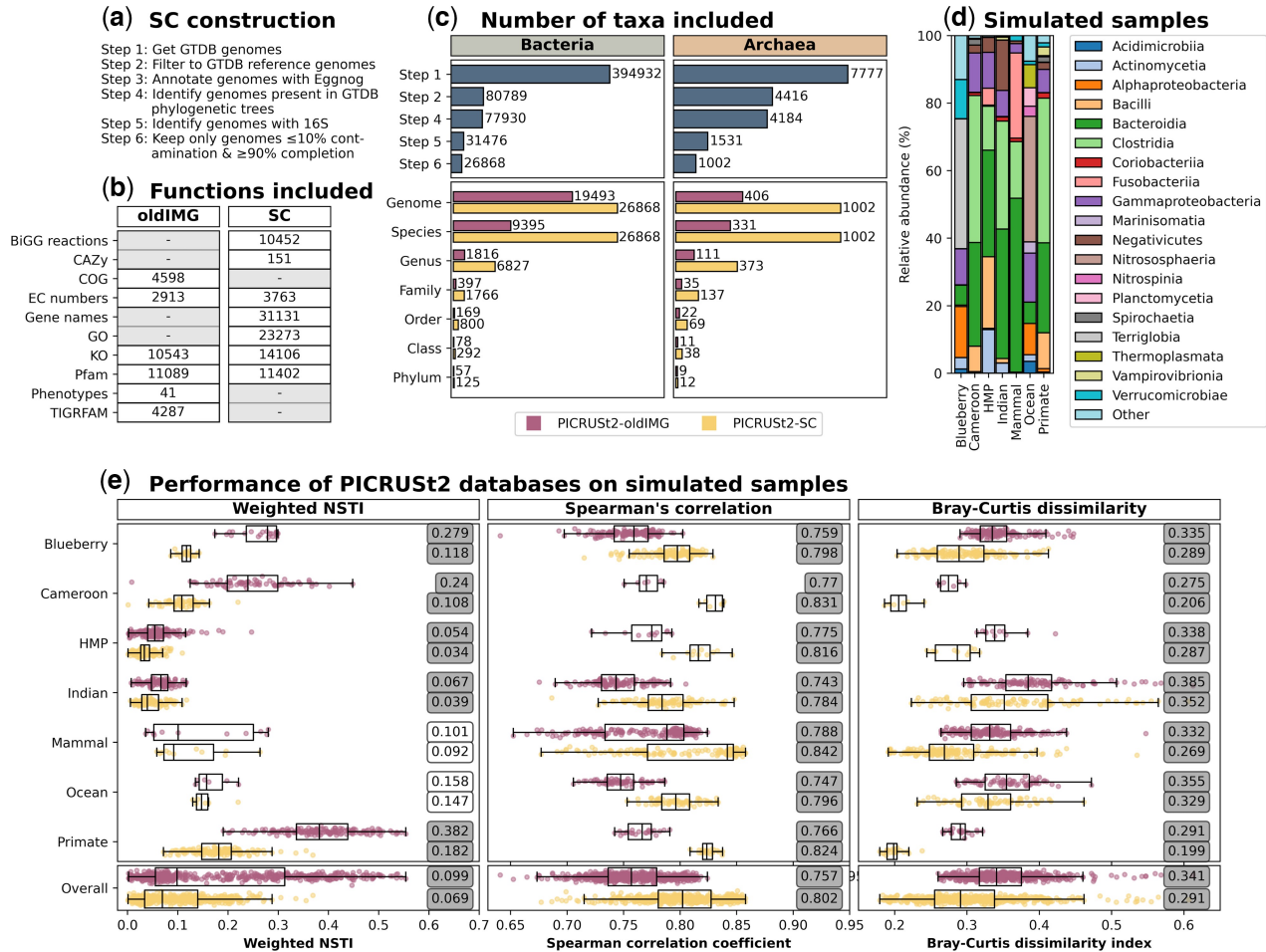


Figure 1. Comparison of the PICRUSt2-oldIMG and PICRUSt2-SC databases showing: (a) the steps in the construction of the PICRUSt2-SC database; (b) the number of functions annotated within different frameworks for the default and updated databases (note that not all frameworks were included in both databases); (c) the number of taxa included for each step of the PICRUSt2-SC construction (top) and for each phylogenetic rank (bottom) for bacteria and archaea; (d) composition at the class level for the simulated samples (the mean relative abundance is shown for each dataset); and (e) the performance of the PICRUSt2-oldIMG and PICRUSt2-SC databases on the simulated samples from each dataset and overall (bottom). Spearman's correlation and Bray-Curtis dissimilarity is shown for KEGG orthologs (EC numbers are in Fig. S5). Individual points are shown for each sample with points being coloured pink for the PICRUSt2-oldIMG database and yellow for the PICRUSt2-SC database. Boxplots represent the median, upper, and lower quartiles, and whiskers show the range of the data (1.5 times the Interquartile Range) and values in boxes are medians. The results for t-tests between the PICRUSt2-oldIMG and PICRUSt2-SC are shown with grey shading for significant ($P \leq .05$) tests.

(hereafter referred to as PICRUSt2-oldIMG). Included bacterial genomes have a median completeness of 99.35% and contamination of 0.66%, while archaeal genomes have a median completeness of 99.03% and contamination of 0.67%.

Annotation of all representative genomes within the new PICRUSt2-SC database was carried out using Egnog v2.1.12 with Prodigal v2.6.3 (Hyatt *et al.* 2010) for gene prediction and the Egnog diamond database v5.0.2. This gave Biochemical, Genetic and Genomic (BiGG) model reactions (Schellenberger *et al.* 2010), Carbohydrate Active Enzymes (CAZy) (Cantarel *et al.* 2009), Enzyme Commission (EC) number, Gene Ontology (GO) (Ashburner *et al.* 2000; Aleksander *et al.* 2023), Kyoto Encyclopedia of Genes and Genomes (KEGG) orthologs (KO) (Kanehisa and Goto 2000), Protein family (Pfam) (Finn *et al.* 2014), and gene name annotations for all genomes. The default PICRUSt2 database contained annotations for KO, EC, Pfam, Clusters of Orthologous Genes (COG) (Tatusov *et al.* 1997; Galperin *et al.* 2021), Phenotype (Markowitz *et al.* 2012), and TIGRFAM (Li *et al.* 2021). By default, only the KO and EC annotations are used, as in the previous version. In the

PICRUSt2-SC database, we have 1.3-fold more annotations than the PICRUSt2-oldIMG database for both KOs (14,106 versus 10,543) and EC numbers (3,763 versus 2,913), while the number of Pfam annotations remains similar (Fig. 1). See Supplementary Information section 2.2 for further details. Users may also add their own functions to the new database, and we provide instructions for doing so on the PICRUSt2 Github.

3 Placement of sequences into phylogenetic trees

With the PICRUSt2-oldIMG database, study sequences are placed in a reference phylogenetic tree constructed using 16S rRNA gene sequences that contains both bacterial and archaeal sequences. In our PICRUSt2-SC database, we now use separate (unrooted) trees for bacteria and archaea constructed using multiple marker genes (from GTDB), and we therefore initially place sequences in both trees and perform 16S copy number prediction and Nearest Sequenced Taxon Index (NSTI) calculations for both trees. We then compare

the NSTI obtained for each and choose the most appropriate (lowest NSTI) for functional predictions. These updates to the code used for running PICRUSt2 are described in [Supplementary Information section 3.1](#). To verify that this tree placement is appropriate, we used the same 16S datasets used in [Douglas et al. \(2020\)](#): (i) Blueberry, $n = 22$, bulk soil and blueberry rhizosphere samples ([Yurgel et al. 2018, 2017](#)); (ii) Cameroon, $n = 57$, stool samples from Cameroonian individuals ([Morton et al. 2015; Lokmer et al. 2019](#)); (iii) HMP, $n = 137$, samples spanning the human body from the Human Microbiome Project ([Huttenhower et al. 2012](#)); (iv) Indian, $n = 91$, stool samples from Indian individuals ([Dhakan et al. 2019](#)); (v) Mammal, $n = 8$, mammalian stool samples ([Finlayson-Trick et al. 2017](#)); (vi) Primate, $n = 77$, non-human primate stool samples ([Amato et al. 2019](#)); and (vii) Ocean, $n = 6$, ocean samples ([Gillies et al. 2015](#)). We took the processed feature tables and representative sequences used by [Douglas et al. \(2020\)](#) and classified these taxonomically using the naïve Bayes Scikit-Learn classifier trained on full-length 16S rRNA gene sequences from the Greengenes2 database ([McDonald et al. 2024](#)) within QIIME2 (v2024.5) ([Bolyen et al. 2019](#)). We ran PICRUSt2 using both the PICRUSt2-oldIMG and the PICRUSt2-SC databases for each of the datasets.

For PICRUSt2-SC predictions, the domain obtained from the taxonomic classifications and the domain of the tree giving the lowest NSTI value for each sequence agreed for 21,543 of 21,552 (99.9%) total sequences within these datasets. See [Supplementary Information section 3.2](#) for details on the nine ASVs that were potentially placed into the wrong tree. Comparing the NSTIs obtained with the PICRUSt2-SC versus the PICRUSt2-oldIMG database on a per sequence basis and on a per sample basis (i.e. the weighted NSTIs that account for sequence abundance within samples), we find that the NSTIs have decreased significantly (t-test $P \leq .05$) in almost all cases ([Supplementary Information section 3.3](#)).

4 Verification of the functional predictions obtained with the updated database

The functional predictions obtained with the PICRUSt2-SC database were correlated with those obtained with the PICRUSt2-oldIMG database ([Supplementary Information section 4.1](#)). In [Douglas et al. \(2020\)](#), Spearman correlation coefficients were calculated between PICRUSt2-obtained KO predictions and HUMAnN2-obtained KO annotations of the paired metagenome samples. Although we did not expect this to be a reliable standard for a database constructed at a different time using different databases for annotation, we calculated Spearman correlation coefficients and Bray-Curtis dissimilarity indices between the PICRUSt2-predictions for both databases and the HUMAnN2 annotations for comparability with the previous study ([Supplementary Information section 4.2](#)).

We instead constructed mock samples using Egg-nog-annotated genomes with simulated abundances. Briefly, we matched the taxa in the seven datasets described above to GTDB genomes that were not present in the PICRUSt2-SC database and used these to construct mock communities that matched the abundances within the seven real datasets. Full details of this mock sample construction can be found in [Supplementary Information section 4.3](#). The mock samples

that we constructed comprised a range of taxa typical of the environment that they were designed to simulate ([Fig. 1](#)), for example, the Blueberry soil samples contained genomes from the Terriglobia class, while the Ocean water samples contained genomes from the Nitrososphaeria class. We ran PICRUSt2 with both the PICRUSt2-oldIMG and the PICRUSt2-SC databases using the mock representative sequences and feature tables for each dataset. Comparing the NSTIs obtained with the PICRUSt2-SC database versus the PICRUSt2-oldIMG database on a per sample basis (i.e. the weighted NSTIs that account for sequence abundance within samples), we find that the NSTIs have decreased in all cases, most of which are statistically significant ([Fig. 1](#)). The NSTIs have decreased from a median of 0.099 (range 0.054 [HMP]—0.382 [Primate]) with PICRUSt2-oldIMG to 0.069 (range 0.034 [HMP]—0.182 [Primate]) with PICRUSt2-SC. While the median NSTI per sequence ([Supplementary Information section 4.4](#)) was slightly lower with the PICRUSt2-oldIMG database for the Cameroon and Indian datasets, the median weighted NSTIs were always lower with the PICRUSt2-SC database, suggesting that it is rarer sequences that have higher NSTIs.

We compared the functional predictions obtained using PICRUSt2 with either the PICRUSt2-oldIMG or the PICRUSt2-SC database with the gold-standard and calculated both Spearman correlation coefficients and Bray-Curtis dissimilarity indices. Median Spearman correlation coefficients are significantly higher for the PICRUSt2-SC (median 0.802; range 0.784 [Indian]—0.842 [Mammal]) versus the PICRUSt2-oldIMG (median 0.757; range 0.743 [Indian]—0.788 [Mammal]) database for both KOs ([Fig. 1](#)) and EC numbers ([Supplementary Information section 4.4](#)). Bray-Curtis dissimilarity indices were also significantly lower for the PICRUSt2-SC (median 0.291; 0.352 [Indian]—range 0.199 [Primate]) versus the PICRUSt2-oldIMG (median 0.341; 0.385 [Indian]—range 0.275 [Cameroon]) database for both KOs ([Fig. 1](#)) and EC numbers ([Supplementary Information section 4.4](#)). KOs had both higher Spearman correlation coefficients and lower Bray-Curtis dissimilarity indices between the mock and PICRUSt2-predictions than EC numbers, and the lowest correlations and highest Bray-Curtis dissimilarity indices were found consistently with the Indian dataset.

5 Computational resources required to run the updated database

We used the sequences within the 16S datasets described above to simulate datasets with 10, 100, 1,000 or 10,000 sequences and 10, 100 or 1,000 samples ([Supplementary Information section 5](#)). Because of the increase in the number of genomes and annotations, the updated PICRUSt2-SC database requires more time and memory to run than PICRUSt2-oldIMG; a dataset containing 1,000 sequences and 1,000 samples being run using 12 threads now takes ~23 mins to run and used 36 GB RAM (compared with ~12 mins and 27 GB RAM with PICRUSt2-oldIMG; [Table S1](#)). These memory requirements can be further reduced with modifications to some of the parameters used to run, and we believe that the resources required are still reasonable for most researchers.

6 Conclusions

We have shown that the updated PICRUSt2-SC database using GTDB genomes performs as expected and outperforms the previous PICRUSt2-oldIMG database in almost all cases, although PICRUSt2-SC requires more time and computational resources than PICRUSt2-oldIMG. This provides PICRUSt2 users with a database that contains more genomes and functions and provides a framework for the PICRUSt2 database to be updated more frequently in the future. The relative ease with which users may now download the database genomes allows users to add their own functions of interest to the database and opens the door for new and improved methods for functional annotations to be used, such as those being developed that use Artificial Intelligence. It remains important for PICRUSt2 users to remember that the output of PICRUSt2 contains functional predictions that are based on the functional annotations of the genomes that are present in the database, but large differences in the genomic content of even very closely related strains with identical 16S rRNA gene sequences may be possible. PICRUSt2-generated functional predictions are therefore recommended as a hypothesis-generating tool and not as the sole piece of evidence for associations between a function and biological phenomena.

Acknowledgements

We thank Dr Gavin Douglas for invaluable insight and assistance with troubleshooting and several PICRUSt2 users for their time testing and giving feedback on an early version of the PICRUSt2-SC database.

Author contributions

Robyn J. Wright (Conceptualization [lead], Data curation [lead], Formal analysis [lead], Investigation [lead], Methodology [equal], Project administration [equal], Validation [lead], Visualization [lead], Writing—original draft [lead], Writing—review & editing [equal]), Morgan G I Langille (Conceptualization [supporting], Funding acquisition [lead], Methodology [equal], Project administration [lead], Resources [lead], Supervision [lead], Writing—original draft [supporting], Writing—review & editing [supporting])

Supplementary data

Supplementary data are available at *Bioinformatics* online.

Conflict of interest: None declared.

Funding

This work was supported by an NSERC Discovery Grant (2016–05039) and was enabled in part by the Digital Research Alliance of Canada (alliancecan.ca).

Data availability

The datasets used for validation were previously published and are available at: NCBI SRA PRJNA389786 and PRJNA484230 (blueberry), ENA PRJEB27005 and MG RAST mgp15238 (Cameroon), <https://www.hmpdacc.org/>

HMTWGS/healthy/ (HMP), ENA PRJNA397112 (Indian), NCBI SRA SRP115632 and SRP115643 (Mammal), QIITA 11212 (Primate), and NCBI SRA SRP056891 (Ocean).

References

- Aleksander SA, Balhoff J, Carbon S, Gene Ontology Consortium *et al.* The gene ontology knowledgebase in 2023. *Genetics* 2023; 224:iyad031.
- Amato KR, G Sanders J, Song SJ *et al.* Evolutionary trends in host physiology outweigh dietary niche in structuring primate gut microbiomes. *ISME J* 2019;13:576–87.
- Ashburner M, Ball CA, Blake JA *et al.* Gene ontology: tool for the unification of biology. The gene ontology consortium. *Nat Genet* 2000; 25:25–9. <https://doi.org/10.1038/75556>
- Bolyen E, Rideout JR, Dillon MR *et al.* Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2. *Nat Biotechnol* 2019;37:852–7.
- Cantalapiedra CP, Hernández-Plaza A, Letunic I *et al.* eggNOG-mapper v2: functional Annotation, orthology assignments, and domain prediction at the Metagenomic Scale. *Mol Biol Evol* 2021; 38:5825–9.
- Cantarel BL, Coutinho PM, Rancurel C *et al.* The Carbohydrate-Active EnZymes database (CAZy): an expert resource for glycogenomics. *Nucleic Acids Res* 2009;37:D233–D238.
- Dhakan DB, Maji A, Sharma AK *et al.* The unique composition of Indian gut microbiome, gene catalogue, and associated fecal metabolome deciphered using multi-omics approaches. *Gigascience* 2019; 8:giz004.
- Douglas GM, Maffei VJ, Zaneveld JR *et al.* PICRUSt2 for prediction of metagenome functions. *Nat Biotechnol* 2020;38:685–8.
- Finlayson-Trick ECL, Getz LJ, Slaine PD *et al.* Taxonomic differences of gut microbiomes drive cellulolytic enzymatic potential within hind-gut fermenting mammals. *PLoS One* 2017;12:e0189404.
- Finn RD, Bateman A, Clements J *et al.* Pfam: The protein families database. *Nucleic Acids Res* 2014;42:D222–30. <https://doi.org/10.1093/nar/gkt1223>
- Galperin MY, Wolf YI, Makarova KS *et al.* COG database update: focus on microbial diversity, model organisms, and widespread pathogens. *Nucleic Acids Res* 2021;49:D274–81.
- Gillies LE, Thrash JC, deRada S *et al.* Archaeal enrichment in the hypoxic zone in the northern Gulf of Mexico. *Environ Microbiol* 2015; 17:3847–56.
- Huerta-Cepas J, Szklarczyk D, Heller D *et al.* EggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res* 2019;47:D309–D314.
- Huttenhower C, Gevers D, Knight R *et al.* Structure, function and diversity of the healthy human microbiome. *Nature* 2012; 486:207–14.
- Hyatt D, Chen GL, LoCascio PF *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 2010;11:119.
- Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res* 2000;28:27–30.
- Langille MGI, Zaneveld J, Caporaso JG *et al.* Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences. *Nat Biotechnol* 2013;31:814–21.
- Li W, O'Neill KR, Haft DH *et al.* RefSeq: expanding the prokaryotic genome annotation pipeline reach with protein family model curation. *Nucleic Acids Res* 2021;49:D1020–8.
- Lokmer A, Cian A, Froment A *et al.* Use of shotgun metagenomics for the identification of protozoa in the gut microbiota of healthy individuals from worldwide populations with various industrialization levels. *PLoS One* 2019;14:e0211139.
- Markowitz VM, Chen I-MA, Palaniappan K *et al.* IMG: the integrated microbial genomes database and comparative analysis system. *Nucleic Acids Res* 2012;40:D115–22.

- McDonald D, Jiang Y, Balaban M *et al.* Greengenes2 unifies microbial data in a single reference tree. *Nat Biotechnol* 2024;**42**:715–8.
- Morton ER, Lynch J, Froment A *et al.* Variation in rural african gut microbiota is strongly correlated with colonization by entamoeba and subsistence. *PLoS Genet* 2015;**11**:e1005658.
- Parks DH, Chuvochina M, Chaumeil P-A *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea. *Nat Biotechnol* 2020;**38**:1079–86.
- Parks DH, Chuvochina M, Waite DW *et al.* A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol* 2018;**36**:996–1004.
- Parks DH, Chuvochina M, Rinke C *et al.* GTDB: An ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Res* 2022;**50**:D785–94.
- Rinke C, Chuvochina M, Mussig AJ *et al.* A standardized archaeal taxonomy for the genome taxonomy database. *Nat Microbiol* 2021;**6**:946–59.
- Schellenberger J, Park JO, Conrad TM *et al.* BiGG: a Biochemical Genetic and Genomic knowledgebase of large scale metabolic reconstructions. *BMC Bioinformatics* 2010;**11**:213.
- Tatusov RL, Koonin EV, Lipman DJ. A genomic perspective on protein families. *Science* 1997;**278**:631–7.
- Yurgel SN, Douglas GM, Dusault A *et al.* Dissecting community structure in wild blueberry root and soil microbiome. *Front Microbiol* 2018;**9**:1187.
- Yurgel SN, Douglas GM, Comeau AM *et al.* Variation in bacterial and eukaryotic communities associated with natural and managed wild blueberry habitats. *Phytobiomes J* 2017;**1**:102–13.