

## eLife's transparent reporting form

We encourage authors to provide detailed information *within their submission* to facilitate the interpretation and replication of experiments. Authors can upload supporting documentation to indicate the use of appropriate reporting guidelines for health-related research (see [EQUATOR Network](#)), life science research (see the [BioSharing Information Resource](#)), or the [ARRIVE guidelines](#) for reporting work involving animal research. Where applicable, authors should refer to any relevant reporting standards documents in this form.

If you have any questions, please consult our Journal Policies and/or contact us: [editorial@elifesciences.org](mailto:editorial@elifesciences.org).

### Sample-size estimation

- You should state whether an appropriate sample size was computed when the study was being designed
- You should state the statistical method of sample size computation and any required assumptions
- If no explicit power analysis was used, you should describe how you decided what sample (replicate) size (number) to use

Please outline where this information can be found within the submission (e.g., sections or figure legends), or explain why this information doesn't apply to your submission:

The submitted work does not include a sampling approach. The collection of data used for the analyses is described in the methods section and Appendix 1 – Table 5.

### Replicates

- You should report how often each experiment was performed
- You should include a definition of biological versus technical replication
- The data obtained should be provided and sufficient information should be provided to indicate the number of independent biological and/or technical replicates
- If you encountered any outliers, you should describe how these were handled
- Criteria for exclusion/inclusion of data should be clearly stated
- High-throughput sequence data should be uploaded before submission, with a private link for reviewers provided (these are available from both GEO and ArrayExpress)

Please outline where this information can be found within the submission (e.g., sections or figure legends), or explain why this information doesn't apply to your submission:

The manuscript computational analyses are based on a collection of metagenomes from 5 different metagenomic surveys and genomes from a public database (Appendix 1 – Table 5). Biological replicates are not present. The results of the analyses performed on the genomic and metagenomic datasets are presented in the figures and figure supplements. Some analyses were performed on a subset of samples or genomes. The sample selection criteria are described in the methods section referring to each analysis.

## Statistical reporting

- Statistical analysis methods should be described and justified
- Raw data should be presented in figures whenever informative to do so (typically when N per group is less than 10)
- For each experiment, you should identify the statistical tests used, exact values of N, definitions of center, methods of multiple test correction, and dispersion and precision measures (e.g., mean, median, SD, SEM, confidence intervals; and, for the major substantive results, a measure of effect size (e.g., Pearson's r, Cohen's d)
- Report exact p-values wherever possible alongside the summary statistics and 95% confidence intervals. These should be reported for all key questions and not only when the p-value is less than 0.05.

Please outline where this information can be found within the submission (e.g., sections or figure legends), or explain why this information doesn't apply to your submission:

The statistical tests applied in the manuscript are described in the methods section corresponding to each analysis and in the figures and figure supplement legends. The number of samples or genomes used for the analysis is also specified in the main text, methods and figure legends.

(For large

datasets, or papers with a very large number of statistical tests, you may upload a single table file with tests, Ns, etc., with reference to sections in the manuscript.)

## Group allocation

- Indicate how samples were allocated into experimental groups (in the case of clinical studies, please specify allocation to treatment method); if randomization was used, please also state if restricted randomization was applied
- Indicate if masking was used during group allocation, data collection and/or data analysis

Please outline where this information can be found within the submission (e.g., sections or figure legends), or explain why this information doesn't apply to your submission:

The submitted work does not include group allocation of the samples.

## Additional data files ("source data")

- We encourage you to upload relevant additional data files, such as numerical data that are represented as a graph in a figure, or as a summary table
- Where provided, these should be in the most useful format, and they can be uploaded as "Source data" files linked to a main figure or table
- Include model definition files including the full list of parameters used
- Include code used for data analysis (e.g., R, MatLab)
- Avoid stating that data files are "available upon request"

Please indicate the figures or tables for which source data files have been provided:

The dataset generated is publicly available on Figshare at <https://doi.org/10.6084/m9.figshare.12459056>

The code used for the analyses in the manuscript is available at <https://github.com/functional-dark-side/functional-dark-side.github.io/tree/master/scripts>. A list with the program versions can be found in [https://github.com/functional-dark-side/functional-dark-side.github.io/blob/master/programs\\_and\\_versions.txt](https://github.com/functional-dark-side/functional-dark-side.github.io/blob/master/programs_and_versions.txt). The code to create the figures is available at <https://zenodo.org/badge/latestdoi/276864152>

A reproducible version of the workflow is available at <https://zenodo.org/badge/latestdoi/251011742>