

# genomesizeR: An R package for genome size prediction

Celine Mercier<sup>1\*</sup> and Joane Elleouet<sup>1\*</sup>

<sup>1</sup> Scion, New Zealand Forest Research Institute, New Zealand \* These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright  
and release the work under a  
Creative Commons Attribution 4.0  
International License ([CC BY 4.0](#)).

## Summary

The genome size of organisms present in an environment can provide many insights into evolutionary and ecological processes at play in that environment. The genomic revolution has enabled a rapid expansion of our knowledge of genomes in many living organisms, and most of that knowledge is classified and readily available in the databases of the National Center for Biotechnology Information (NCBI). The genomesizeR tool leverages the wealth of taxonomic and genomic information present in NCBI databases to infer the genome size of Archaea, Bacteria, or Eukaryote organisms identified at any taxonomic level. This R package uses statistical modelling on data from the most up-to-date NCBI databases and provides three statistical methods for genome size prediction of a given taxon, or group of taxa. A straightforward 'weighted mean' method identifies the closest taxa with available genome size information in the taxonomic tree, and averages their genome sizes using weights based on taxonomic distance. A frequentist random effect model uses nested genus and family information to output genome size estimates. Finally a third option provides predictions from a distributional Bayesian multilevel model which uses taxonomic information from genus all the way to superkingdom, therefore providing estimates and uncertainty bounds even for under-represented taxa.

All three methods use:

- A list of queries; a query being a taxon or a list of several taxa. The package was designed to make it easy to use with data coming from environmental DNA experiments, but works with any table of taxa.
- A reference database containing all the known genome sizes, built from the NCBI databases, with associated taxa, provided in an archive to download.
- A taxonomic tree structure as built by the NCBI, provided in the same archive.

genomesizeR retrieves the taxonomic classification of input queries, estimates the genome size of each query, and provides 95% confidence intervals for each estimate.

## Statement of need

The size of microbial genomes and its evolution can provide important insights into evolutionary and ecological processes influencing both microbial species and the environments in which they inhabit. The shedding of unnecessary genetic elements and their associated biosynthetic pathways, for example, is a common phenomenon observed in organisms with a high degree of host symbiosis (Brader et al., 2014; Moran, 2002; Vandenkoornhuys et al., 2007). Genome size reduction has also been observed in organisms experiencing arid environments (Liu et al., 2023), or a narrow range of substrates or metabolic options (Tyson et al., 2004). Among many others, these findings demonstrate the opportunities associated with including genome size as a key trait in microbial communities to provide insights spanning niche size, co-evolution, adaption, and metabolic flexibility of the microbiomes present, but also stability, and ecophysiological and functional complexity of abiotic and biotic environments.

42 However, characterizing genome size for all organisms in a microbiome remains challenging.  
 43 Methods in the past have included the use of DNA staining through to cell enumeration, flow  
 44 cytometry, culturing, gel electrophoresis, and DNA renaturation kinetics. All have merits and  
 45 limitations (Raes et al., 2007). The widespread availability and use of microbiome related  
 46 tag-amplicon DNA sequencing has tremendously increased our scientific knowledge of microbial  
 47 genomics. There is also an opportunity to explore the rapidly expanding archives of short read  
 48 DNA libraries (i.e. extant 16S and ITS amplicon sequences).

49 Alternatively, when well documented and archived in user-friendly and publicly available  
 50 databases, the exponentially growing genomic knowledge of micro-organisms is an inexpensive  
 51 resource unlocking a myriad of research opportunities in all fields of environmental sciences,  
 52 from human and agricultural microbiomes through to aquatic, soil, and atmospheric ecosystems.  
 53 Combining available genome size information to data and metadata on community structure  
 54 from existing projects can add further scientific value to investment already made in these  
 55 projects, at no added cost.

56 However, genome size estimates for many taxa found in environmental samples are missing  
 57 from public databases, or fully unknown. The evolutionary rule that phylogenetically related  
 58 organisms share genetic similarities can be exploited and genome size for taxa with unknown  
 59 genome size can be statistically inferred from related taxa with known genome size, using  
 60 taxonomy as a proxy for phylogeny. Another challenge is the precision of identification: some  
 61 taxa can only be identified at high taxonomic levels. Statistical methods can also be used to  
 62 infer their genome size range from databases. To our knowledge, there is no convenient and  
 63 fast way to obtain genome size estimates with uncertainty bounds for all organisms identified  
 64 or partially identified in an environmental sample.

65 Using the increased prevalence of whole-genome information for all organisms, we have  
 66 therefore developed genomesizeR, allowing the inference of genome size of many queries at  
 67 once, based on taxonomic information and available genome data from the National Center  
 68 for Biotechnology Information (NCBI).

## 69 Methods

### 70 NCBI database filtering and processing

71 The reference database used is built by querying all genome metadata information from the  
 72 curated NCBI RefSeq database (O'Leary et al., 2016). Filters are applied to only keep full  
 73 genomes, and discard data that the NCBI has tagged as anomalous, or abnormally large or  
 74 small.

75 This raw database is then prepared to include more pre-computed information to be used by  
 76 the package. Genome sizes are aggregated to the species level by iteratively averaging all  
 77 entries below, hence the package can only provide estimates at the level of species and above.  
 78 Average genome sizes and their associated standard error values are also pre-computed, to be  
 79 used by the weighted mean method.

### 80 Bayesian method

81 The NCBI database of species with known genome sizes was split by superkingdom (Bacteria,  
 82 Archaeae, Eukaryotes). A distributional Bayesian linear hierarchical model using the brm function  
 83 from the brms package was fitted to each superkingdom dataset. The general model structure  
 84 is outlined below.

$$\log(G_i) \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

85 where  $G_i$  is the genome size of species  $i$  in the units of 10 Mbp. The model uses predictors  
86 for both mean and standard deviation. The mean is modelled as follows:

$$\begin{aligned}\mu_i &= \alpha_0 + \alpha_{genus_{g[i]}} + \alpha_{family_{f[i]}} + \alpha_{order_{o[i]}} + \alpha_{class_{c[i]}} + \alpha_{phylum_{p[i]}} \\ \alpha_{genus_{g[i]}} &\sim \mathcal{N}(0, \sigma_{genus}^2) \\ \alpha_{family_{f[i]}} &\sim \mathcal{N}(0, \sigma_{family}^2) \\ \alpha_{order_{o[i]}} &\sim \mathcal{N}(0, \sigma_{order}^2) \\ \alpha_{class_{c[i]}} &\sim \mathcal{N}(0, \sigma_{class}^2) \\ \alpha_{phylum_{p[i]}} &\sim \mathcal{N}(0, \sigma_{phylum}^2)\end{aligned}$$

87 The following prior distributions are used:

$$\begin{aligned}\alpha_0 &\sim \mathcal{N}(0, 5) \\ (\sigma_{genus}, \sigma_{family}, \sigma_{order}, \sigma_{class}, \sigma_{phylum}, s_{class}, s_{phylum}) &\sim \mathcal{N}^+(0, 1)\end{aligned}$$

88 Differences in genome size variability was observed among taxa, therefore the model also adds  
89 predictors to the standard deviation of the response. The standard deviation is modelled as  
90 follows:

$$\begin{aligned}\log(\sigma_i) &= \lambda_0 + \lambda_{class_{c[i]}} + \lambda_{phylum_{p[i]}} \\ \lambda_{class_{c[i]}} &\sim \mathcal{N}(0, s_{class}^2) \\ \lambda_{phylum_{p[i]}} &\sim \mathcal{N}(0, s_{phylum}^2)\end{aligned}$$

91 with priors

$$\begin{aligned}\lambda_0 &\sim \mathcal{N}(0, 1) \\ (s_{class}, s_{phylum}) &\sim \mathcal{N}^+(0, 1)\end{aligned}$$

92  $\mathcal{N}^+$  is the normal distribution truncated to positive values.  $g[i], f[i], o[i], c[i]$  and  $p[i]$  are  
93 respectively the index for the genus, family, order, class, and phylum of entry  $i$  in the species-  
94 level database. Note that taxonomic groups are naturally nested within each other and the  
95 indices are designed to be unique to the particular taxonomic group it represents.

96 The estimation process uses Stan's Hamiltonian Monte Carlo algorithm with the U-turn  
97 sampler.

98 Posterior predictions are obtained using the predict function from the brms package, and 95%  
99 credible intervals are obtained using 2.5% and 97.5% quantiles from the posterior distribution.

100 Queries corresponding to identified species with an available genome size estimate in the NCBI  
101 database get allocated the genome size value of the database (averaged at the species level)  
102 and 95% confidence intervals are calculated as XXX

### 103 Frequentist method

104 A frequentist linear mixed-effects model using the `lmer` function from the `lme4` package (Bates  
105 et al., 2015) was fitted to the NCBI database of species with known genome sizes. The model  
106 is as follows:

$$\log(G_i) = \alpha_0 + \alpha_{\text{genus}_{g[i]}} + \alpha_{\text{family}_{f[i]}} + e_i$$

107 where  $\alpha_0$  is the overall mean,  $\alpha_{\text{genus}_{g[i]}}$  and  $\alpha_{\text{family}_{f[i]}}$  are random effect of genus and family  
108 for genus  $g[i]$  and family  $f[i]$  and  $e_i$  is the residual error of observation  $i$ .

109 The estimation process using the restricted maximum likelihood method (REML). A prediction  
110 interval is computed using the `predictInterval` function from the `merTools` package (Knowles  
111 & Frederick, 2024). As higher nested levels (order, class) are not taken into account in the  
112 model, predictions are not produced for queries above the family level.

### 113 Weighted mean method

114 The weighted mean method computes the genome size of a query by averaging the known  
115 genome sizes of surrounding taxa in the taxonomic tree, with a weighted system where further  
116 neighbours have less weight in the computed mean. The identification of related taxa is limited  
117 to levels below and including order.

118 In the calculation of the prediction for a given query, the weight for a related taxon in the  
119 database is calculated as

$$w_i = \frac{1}{d_i + 1}$$

120 where  $d_i$  is the distance of the related taxon to the query in number of nodes in the taxonomy.

121 The confidence interval is calculated as:

$$CI = 1.96 \times \sqrt{\hat{\mu}} \quad (1)$$

122 where  $\hat{\mu}$  is the computed weighted mean.

123 For queries relating to well-characterised species where many genetic studies have been  
124 performed, such as model organisms, this might lead to more precise predictions than the two  
125 other methods. This method can also perform better than the others if your queries consist of  
126 lists of taxa (for example, an output of *blastn* where several matches can be obtained for each  
127 query). Otherwise, we suggest using one of the other methods, as the confidence intervals  
128 calculated are less reliable for the weighted mean method.

### 129 Implementation

130 The main steps of all methods are multithreaded on POSIX systems using the packages `parallel`  
131 (`R Core Team`, 2024) and `doParallel` (`Corporation & Weston`, 2022).

132 The packages `ncbitax` (Machne, 2022) and `CHNOSZ` (Dick, 2019) are used to read the  
133 taxonomy data, `dplyr` (Wickham et al., 2023) and `biomformat` (McMurdie & Paulson, 2024)  
134 are used for some of the formatting, and `pbapply` (Solymos & Zawadzki, 2023) is used to  
135 display the progress bar.

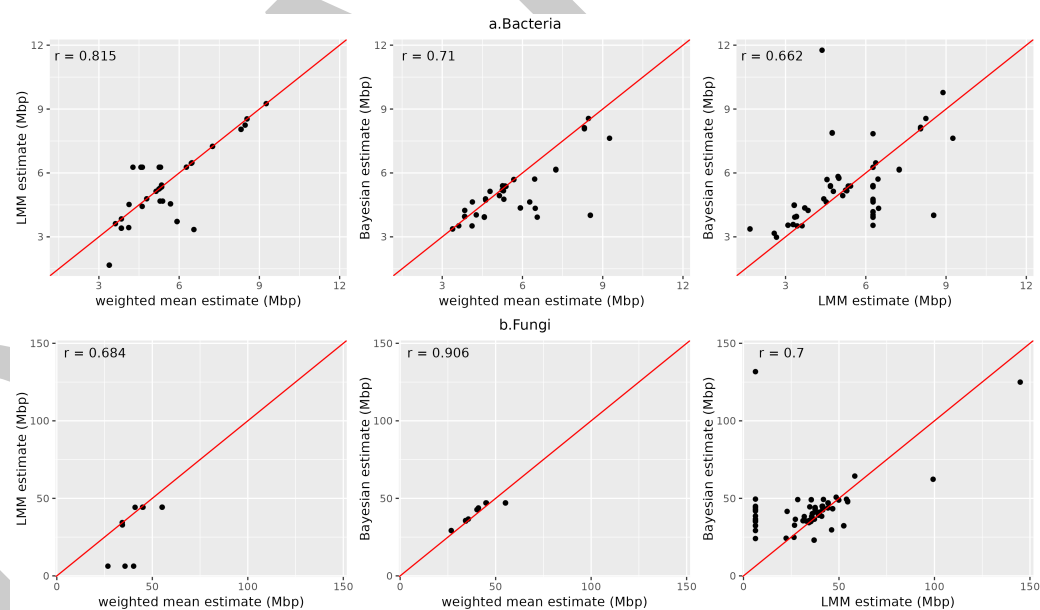
The R package accepts as input formats the common 'taxonomy table' format used by popular packages such as phyloseq (McMurdie & Holmes, 2013) and mothur (Schloss et al., 2009), and any file or data frame with a column containing either NCBI taxids or taxon names. The output format is a data frame with the same columns as the input, with some added columns providing information about the estimation and the quality of the estimation. The user can also choose a simple output format only containing the estimation information.

Several plotting functions using the ggplot2 (Wickham, 2016) and ggtree (Yu, 2020) packages are also provided to visualise the results.

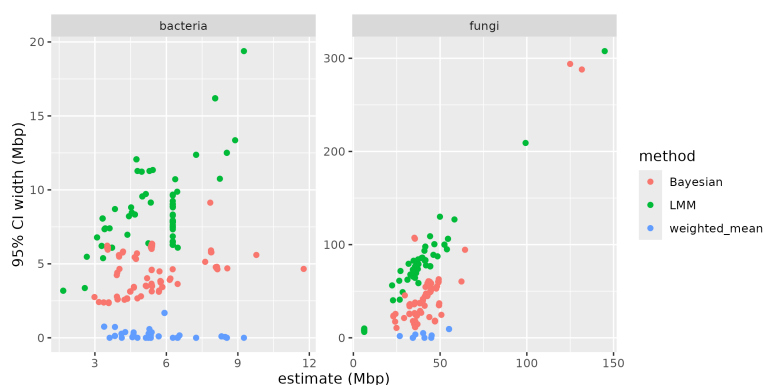
## Method comparison

The applicability of each method varies. The Bayesian method outputs results for any taxon that is recognised in the NCBI taxonomy. The frequentist random effects model method only outputs results for queries that have a match at the species, genus, or family level. The weighted mean method only performs an estimation for queries that have at least two matches at the species, genus, family, or order level. Below is a comparison of estimates for an example set of bacteria and fungi queries where the highest level of match with the database is the family level. Note that there are fewer successful estimations with the weighted mean method than with the two model-based methods.

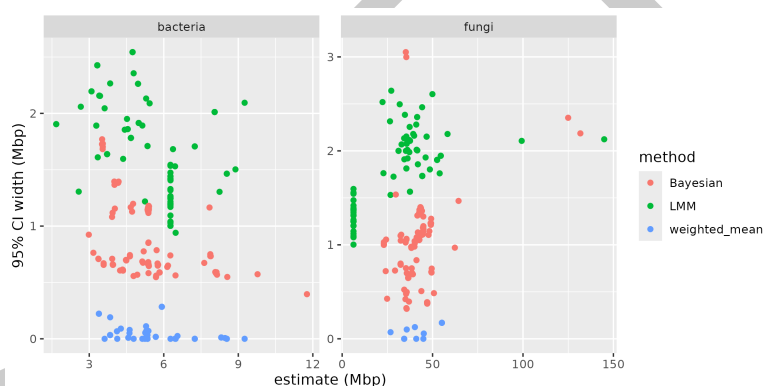
Figures below show that estimates and the width of confidence intervals differ between methods (figures Figure 1, Figure 2 and Figure 3).



**Figure 1:** Pairwise comparison of estimates from different methods for a. bacteria and b. fungi. Pearson's correlation coefficient is displayed at the top left.



**Figure 2:** Pairwise comparison of 95% confidence intervals from different methods for a. bacteria and b. fungi.



**Figure 3:** Pairwise comparison of relative 95% confidence intervals (scaled by estimated size) from different methods for a. bacteria and b. fungi.

## Example

```

155 This example data is a subset of the dataset from Labouyrie et al. (2023).
156 First, the genome sizes are predicted from the taxa:
157
158 results = estimate_genome_size(example_input_file, refdata_archive_path,
159 sep='\t', match_column='TAXID', output_format='input')
160 #####
161 # Genome size estimation summary:
162 #
163 # 32.22222 % estimations achieving required precision
164 #
165     Min.    1st Qu.    Median      Mean    3rd Qu.      Max.
166 2260662  4850009  15733757  23739845  39882934  161191581
167
168 # Estimation status:
169 Confidence interval to estimated size ratio > ci_threshold | OK
170                                                            122 | 58
171 Then, the results can be visualized using the plotting functions provided. Figure 4 shows a
172 histogram and a boxplot of the estimated genome sizes for each sample. Figure 5 shows a tree

```





## Availability

Project name: genomesizeR Project home page: <https://github.com/ScionResearch/genomesizeR>  
 Operating system(s): Platform independent Programming language: R License: GNU General Public License

## Acknowledgements

We acknowledge contributions from Sean Husheer. The authors declare that they have no conflict of interest. Funding for this research came from the Tree-Root-Microbiome programme, which is funded by MBIE's Endeavour Fund and in part by the New Zealand Forest Growers Levy Trust (C04X2002). We make no warranties regarding the accuracy or integrity of the Data. We accept no liability for any direct, indirect, special, consequential or other losses or damages of whatsoever kind arising out of access to, or the use of the Data. We are in no way to be held responsible for the use that you put the Data to. You rely on the Data entirely at your own risk.

## References

- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- Brader, G., Compant, S., Mitter, B., Trognitz, F., & Sessitsch, A. (2014). Metabolic potential of endophytic bacteria. *Current Opinion in Biotechnology*, 27, 30–37.
- Corporation, M., & Weston, S. (2022). *doParallel: Foreach parallel adaptor for the 'parallel' package*. <https://CRAN.R-project.org/package=doParallel>
- Dick, J. M. (2019). CHNOSZ: Thermodynamic calculations and diagrams for geochemistry. *Front. Earth Sci.*, 7.
- Knowles, J. E., & Frederick, C. (2024). *merTools: Tools for analyzing mixed effect regression models*. <https://github.com/jknowles/mertools>
- Labouyrie, M., Ballabio, C., Romero, F., Panagos, P., Jones, A., Schmid, M. W., Mikryukov, V., Dulya, O., Tedersoo, L., Bahram, M., Lugato, E., Heijden, M. G. A. van der, & Orgiazzi, A. (2023). Patterns in soil microbial diversity across europe. *Nat. Commun.*, 14(1), 3311.
- Liu, H., Zhang, H., Powell, J., Delgado-Baquerizo, M., Wang, J., & Singh, B. (2023). Warmer and drier ecosystems select for smaller bacterial genomes in global soils. *iMeta*, 2(1), e70.
- Machne, R. (2022). *Ncbi tax: Simple base r parser for NCBI taxonomy*. <https://github.com/raim/ncbitax>
- McMurdie, P. J., & Holmes, S. (2013). Phyloseq: An r package for reproducible interactive analysis and graphics of microbiome census data. *PLoS ONE*, 8(4), e61217. <http://dx.plos.org/10.1371/journal.pone.0061217>
- McMurdie, P. J., & Paulson, J. N. (2024). *Biomformat: An interface package for the BIOM file format*. <https://doi.org/10.18129/B9.bioc.biomformat>
- Moran, N. A. (2002). Microbial minimalism: Genome reduction in bacterial pathogens. *Cell*, 108(5), 583–586.
- O'Leary, N. A., Wright, M. W., Brister, J. R., Ciufu, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse, B., Smith-White, B., Ako-Adjei, D., Astashyn, A., Badretdin, A., Bao, Y., Blinkova, O., Brover, V., Chetvernin, V., Choi, J., Cox, E., Ermolaeva, O., ... Pruitt, K.



- 216 D. (2016). Reference sequence (RefSeq) database at NCBI: Current status, taxonomic  
217 expansion, and functional annotation. *Nucleic Acids Res.*, 44(D1), D733–45.
- 218 R Core Team. (2024). *R: A language and environment for statistical computing*. R Foundation  
219 for Statistical Computing. <https://www.R-project.org/>
- 220 Raes, J., Korb, J. O., Lercher, M. J., Von Mering, C., & Bork, P. (2007). Prediction of  
221 effective genome size in metagenomic samples. *Genome Biology*, 8, 1–11.
- 222 Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B.,  
223 Lesniewski, R. A., Oakley, B. B., Parks, D. H., Robinson, C. J., Sahl, J. W., Stres, B.,  
224 Thallinger, G. G., Van Horn, D. J., & Weber, C. F. (2009). Introducing mothur: Open-  
225 source, platform-independent, community-supported software for describing and comparing  
226 microbial communities. *Applied and Environmental Microbiology*, 75(23), 7537–7541.  
227 <https://doi.org/10.1128/AEM.01541-09>
- 228 Solymos, P., & Zawadzki, Z. (2023). *Pbapply: Adding progress bar to '\*apply' functions*.  
229 <https://CRAN.R-project.org/package=pbapply>
- 230 Tyson, G. W., Chapman, J., Hugenholtz, P., Allen, E. E., Ram, R. J., Richardson, P. M.,  
231 Solovyev, V. V., Rubin, E. M., Rokhsar, D. S., & Banfield, J. F. (2004). Community struc-  
232 ture and metabolism through reconstruction of microbial genomes from the environment.  
233 *Nature*, 428(6978), 37–43.
- 234 Vandenkoornhuyse, P., Mahé, S., Ineson, P., Staddon, P., Ostle, N., Cliquet, J.-B., Francez,  
235 A.-J., Fitter, A. H., & Young, J. P. W. (2007). Active root-inhabiting microbes identified  
236 by rapid incorporation of plant-derived carbon into RNA. *Proceedings of the National*  
237 *Academy of Sciences*, 104(43), 16970–16975.
- 238 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.  
239 ISBN: 978-3-319-24277-4
- 240 Wickham, H., François, R., Henry, L., Müller, K., & Vaughan, D. (2023). *Dplyr: A grammar*  
241 *of data manipulation*. <https://dplyr.tidyverse.org>
- 242 Yu, G. (2020). Using ggtree to visualize data on Tree-Like structures. *Curr. Protoc. Bioinfor-*  
243 *matics*, 69(1). <https://doi.org/10.1002/cpbi.96>