

1 **Straightforward and reproducible analysis of bacterial pangenomes using Pagoo**

2

3 Ignacio Ferrés^{1,*}, Gregorio Iraola^{1,2,3,*}

4

5

6

7 ¹ Microbial Genomics Laboratory, Institut Pasteur Montevideo, Montevideo, Uruguay.

8 ² Center for Integrative Biology, Universidad Mayor, Santiago de Chile, Chile.

9 ³ Wellcome Sanger Institute, Hinxton, United Kingdom.

10

11

12 * To whom correspondence should be addressed: Gregorio Iraola. Tel: +5985220910; Fax:

13 +5985220911; Email: giraola@pasteur.edu.uy. Correspondence may also be addressed to:

14 Ignacio Ferrés. Tel: +5985220910; Fax: +5985220911; Email: iferres@pasteur.edu.uy.

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37 **Pangenome analysis is fundamental to explore evolutionary processes occurring in**
38 **bacterial populations. However, the lack of standardized methods for handling diverse**
39 **pangenomic datasets and complex metadata hinders more straightforward and**
40 **reproducible downstream analyses. To fill this gap, we introduce Pagoo, a new**
41 **framework that integrates pangenome data, analytical methods and visualization tools**
42 **in a single object that can be easily stored, shared and responsively queried for**
43 **improved biological interpretation of bacterial evolution.**

44 The exponentially growing number of diverse bacterial genomes has prompted
45 pangenome reconstruction as a gold standard to explore genetic diversity of bacterial
46 populations^{1,2}. Pangenome comparisons reveal genome evolutionary dynamics associated
47 with important biological processes such as speciation, host-adaptation, pathogenicity or the
48 acquisition of antimicrobial resistance. Pangenome reconstruction is typically performed
49 from genes annotated in a set of whole-genome sequences. In general, coding sequences of
50 different strains are grouped in orthologous clusters based on different similarity criteria.
51 Then, pangenome data informs about the belonging of each gene encoded in each genome to
52 a certain orthologous cluster. In the recent years, several software tools have been developed
53 to reconstruct bacterial pangenomes, such as Roary³, panX⁴ or PanOCT⁵. These tools focus
54 on automation of steps and optimization of computational costs to cluster thousands of
55 sequences of increasingly large genomic datasets. However, there is a lack of tools that can
56 take the output of pangenome reconstruction softwares and provide standardized and
57 straightforward methods for data integration, storage, analysis and visualization.

58 Here, we introduce Pagoo, the first pangenome post-processing tool that can take the
59 output of pangenome reconstruction softwares providing a standardized framework for its
60 analysis. Pagoo is based on an object-oriented design built on a novel class system in R
61 which implements: i) an integrative data structure for standardized storage of pangenome
62 information such as orthologous clusters, sequences, annotations and metadata in a single
63 object; ii) a set of straightforward methods for responsive querying, handling and subsetting
64 of this data structure; and iii) a set of standard statistics and active visualizations leveraging
65 flexible downstream comparative analyses. Along with extensive documentation, we show
66 how Pagoo interacts with other widely used microbial genomics tools and the R ecosystem
67 for improved analysis of bacterial populations.

68 A pangenome can be represented as individual genes which belongs to organisms
69 (genomes) and that are also assigned to a cluster of orthologous genes. Pagoo stores this as a
70 three-column matrix, with one column identifying an individual gene, the next one
71 identifying the organism that this gene belongs to, and the last one identifying the
72 orthologous cluster that the gene was assigned by the pangenome reconstruction method.
73 Optionally, this matrix can contain additional columns as gene-specific metadata like
74 annotations or functional assignments. Orthologous clusters and organisms can also take
75 metadata represented as two different matrices, with the condition that each one must contain
76 a column that correctly maps each observation (cluster or organism) into the former matrix.
77 Gene sequences can also be added to this structure, with the condition that their names must
78 also map to rows in the first matrix (Fig. 1A). This relational structure optimizes data storage
79 avoiding duplication, enables flexibility for working with different data types and facilitates
80 complex querying and analysis.

81 Indeed, a salient and unique feature of Pagoo is that this data structure is stored and
82 managed in an encapsulated, object-oriented fashion using the R6 package as backend. In
83 contrast with traditional R programming, the R6 paradigm considers that methods belong to
84 objects rather than to generic functions, so an object contains both the data and embedded
85 methods to analyze it. In this context, the Pagoo object is built on three novel R6 classes.
86 PgR6 is the most basic class that contains methods and functions for data handling and
87 subsetting. Then, PgR6M inherits all the methods and fields from PgR6 and incorporates
88 statistical methods and visualization tools based on the ggplot2 package⁶. PgR6MS inherits
89 all capabilities from the others and adds methods for manipulation of DNA sequences using
90 the Biostrings package⁷ (Fig. 1B). These classes support the main data types that typically
91 represent a pangenome, providing a novel and synergistic framework to manage both the raw
92 data and methods to perform operations and explore results with customized visualizations.
93 Moreover, any of these classes could be further inherited and easily extended by third party
94 applications.

95 Another remarkable feature of Pagoo is that raw data stored in the pangenome object is
96 kept unaltered in the background, while users can query, mutate or subset the object using
97 active bindings. This allows changing the state of the object without altering the original data.
98 For example, users can temporarily hide certain organisms from the dataset, actively set
99 thresholds that change the definition of core genes, or extract specific information from
100 organisms, genes, clusters or sequences. Class-specific methods for generic subset operators
101 are also implemented enabling seemingly extraction of relevant field subsets straight from the

102 object by using widely known R subset notation. Also, Pagoo provides specific methods to
103 automatically generate the pangenome object from output files produced by standard
104 pangenome reconstruction tools like Roary, and to save any changes to the object along with
105 the unaltered original data as a single file. Importantly, Pagoo lacks of external dependencies
106 and is built and tested in all three major operative systems (Linux, Windows and Mac). A
107 detailed explanation of each method and operator for data input, saving and loading the
108 pangenome object, and for specific data handling and subsetting is provided in the online user
109 manual (<https://iferres.github.io/pagoo/>). Together, this implementation represents a new
110 concept for pangenome data handling, facilitating reproducibility and enabling multiple and
111 flexible analyses.

112 Pagoo also includes statistical and visualization methods. Customized plots and
113 statistical analyses can be generated directly from the pangenome object using active
114 bindings on the console or by deploying a built-in R-Shiny application. This interactive
115 application is divided in two main components: (i) a general dashboard that interactively
116 displays summary statistics including number of organisms, orthologous clusters and genes,
117 core and accessory genome sizes, gene frequency barplots, pangenome curves and scrollable
118 information about core genome clusters and genes (i.e. annotation or any other metadata);
119 and (ii) a specific dashboard showing clustering of genomes according to accessory gene
120 distances and Principal Components Analysis, genome-specific accessory genome sizes,
121 visualization of gene presence/absence matrix with associated metadata and information
122 about accessory gene clusters (Supplementary Information; Fig. S1). This interactive
123 application allows responsive exploration of evolutionary trends in bacterial populations to
124 guide downstream analyses, leveraging the interaction of Pagoo with other tools.

125 Remarkably, more complex comparative pangenome analyses can be performed by
126 applying concise code recipes. We define recipes as relatively short snippets that pipe
127 pangenome information extracted from the object as input to other R tools. We have
128 developed example recipes (available in the online user manual at
129 <https://iferres.github.io/pagoo/articles/6-Recipes.html>) to build core genome phylogenies,
130 identify population structure, explore genome-wide selective pressures acting over the core
131 genes and compare individual gene sequences against specific databases. Importantly, the
132 development and implementation of recipes enable full reproducibility of publication-quality
133 figures generated directly from the pangenome object (Fig. 2).

134 As a working example we used Pagoo to reanalyze a previously published study on the
135 evolution of *Campylobacter fetus* pangenome. This species has a strong population structure

136 with different lineages adapted to livestock or humans⁸. Briefly, we reconstructed a
137 pangenome from 69 selected *C. fetus* genomes with Roary³ using default parameters and used
138 its output to build a Pagoo object. Then, we performed a comparative analysis between
139 livestock- and human-derived *C. fetus* genomes. The dynamic exploration of results using the
140 Pagoo Shiny application (https://microgenlab.shinyapps.io/pagoo_campylobacter/) allowed
141 us to recover main diversity patterns reported for this species, such as a marked difference
142 between accessory genome size and gene presence/absence patterns between livestock- and
143 human-adapted strains.

144 The advent of high-throughput sequencing technologies more than fifteen years ago
145 pushed microbiology towards the field of comparative genomics, that rapidly transitioned
146 from studies including few to thousands of genomes². This substantially increased the
147 complexity of datasets, requiring new approaches to systematically handle and track different
148 components of interrelated pangenomic data. Pagoo introduces a new framework
149 underpinned in a concept that leverages the simplicity of storing all the information in a
150 standardized and reproducible manner in a single, shareable object. Along with future
151 developments and add-ons, Pagoo aims to improve and facilitate current practices on the
152 genomic analysis of bacterial populations.

153

154 **References**

155

1. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: implications for the microbial ‘pan-genome’. *Proc. Natl. Acad. Sci. U.S.A.* **102**, 13950–13955 (2005).
2. Vernikos, G., Medini, D., Riley, D. R. & Tettelin, H. Ten years of pan-genome analyses. *Curr. Opin. Microbiol.* **23**, 148–154 (2015).
3. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691–3693 (2015).
4. Ding, W., Baumdicker, F. & Neher, R. A. panX: pan-genome analysis and exploration. *Nucleic Acids Res.* **46**, e5 (2018).
5. Fouts, D. E., Brinkac, L., Beck, E., Inman, J. & Sutton, G. PanOCT: automated clustering of orthologs using conserved gene neighborhood for pan-genomic analysis of bacterial

strains and closely related species. *Nucleic Acids Res.* **40**, e172 (2012).

6. Wickham, H. ggplot2. *WIREs Computational Statistics* **3**, 180–185 (2011).

7. Biostrings: Efficient manipulation of biological strings version 2.56.0 from Bioconductor.

<https://rdrr.io/bioc/Biostrings/>.

8. Iraola, G. *et al.* Distinct *Campylobacter fetus* lineages adapted as livestock pathogens and human pathobionts in the intestinal microbiota. *Nature Communications* **8**, 1367 (2017).

156

157

158 **Competing interests.** Nothing to declare.

159

160 **Acknowledgments.** We thank Pablo Fresia and Daniela Costa for insightful comments and
161 suggestions during testing of Pagoo. I.F. is funded by grant ANII-
162 POS_NAC_2018_1_151494 from Agencia Nacional de Investigación e Innovación (ANII),
163 Uruguay.

164

165 **Author contributions.** G.I. and I.F. conceived the idea, I.F. developed the software and
166 performed experiments, G.I. and I.F. wrote the manuscript.

167

168 **Figure legends**

169

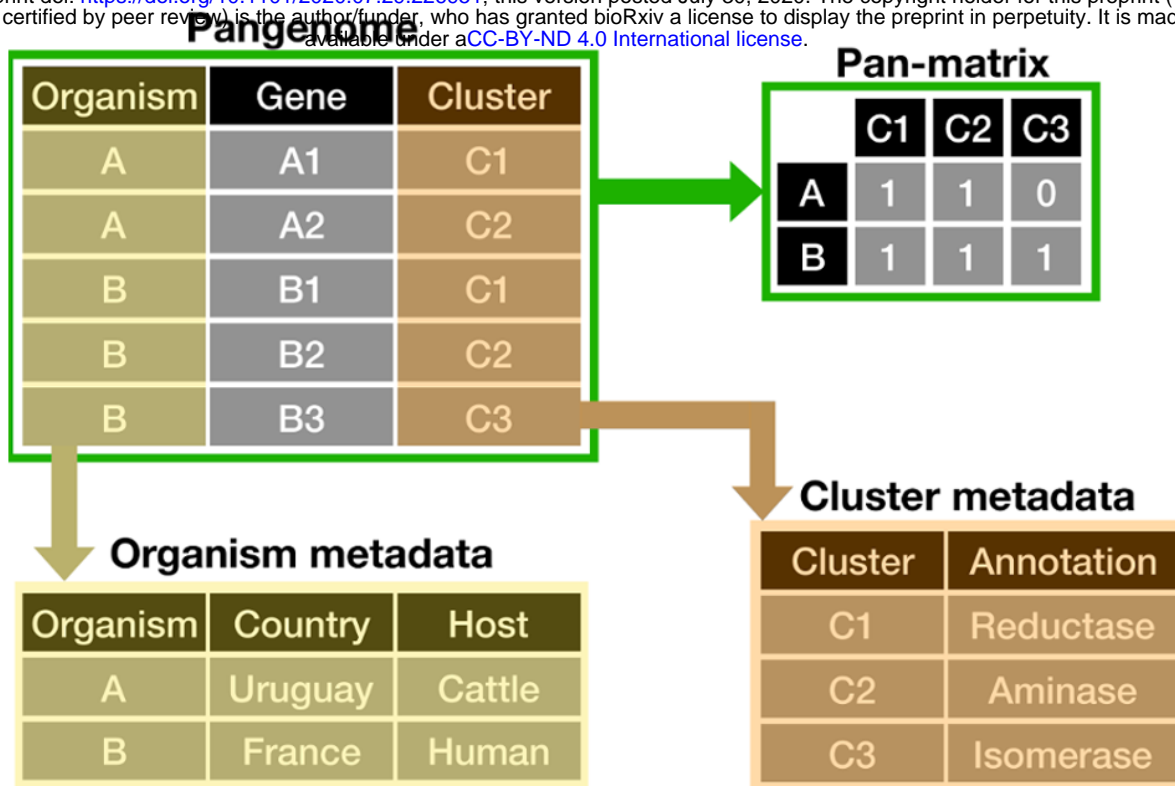
170 **Figure 1. Framework and overall design of Pagoo.** A) Example of the relational structure
171 implemented to store, link and operate over different pangenome data types. B) General
172 description of the workflow from assembled genomes to Pagoo analysis. Once pangenome
173 files are created with any available pangenome reconstruction software, these files can be
174 loaded to create the Pagoo object. The specific R6 classes store and manage different data
175 types that allow to store all the information in a single file or perform comparative analyses
176 using the R console interface or the Pagoo Shiny application.

177

178 **Figure 2. Results extracted from the pangenome object.** Exploration of the *C. fetus*
179 pangenome using information directly extracted from the pangenome object and customized
180 aesthetics. Panel (A) shows pangenome and core genome curves with grey circles
181 representing different sub-samples at increasing number of genomes; the black lines show the
182 fitting to the power law and exponential decay functions, respectively. Panel B shows the

183 distribution of genes in different subset of genomes. Panel C shows a Principal Components
184 Analysis generated from the gene presence/absence matrix that clearly two groups of
185 genomes, representing human-derived strains (red) and bovine-derived strains (green). Panel
186 D shows the distribution of the pangenome in core genes and accessory genes (shell and
187 cloud genes).

A



B

Bacterial genomes



Pangenome reconstruction

Pangenome files

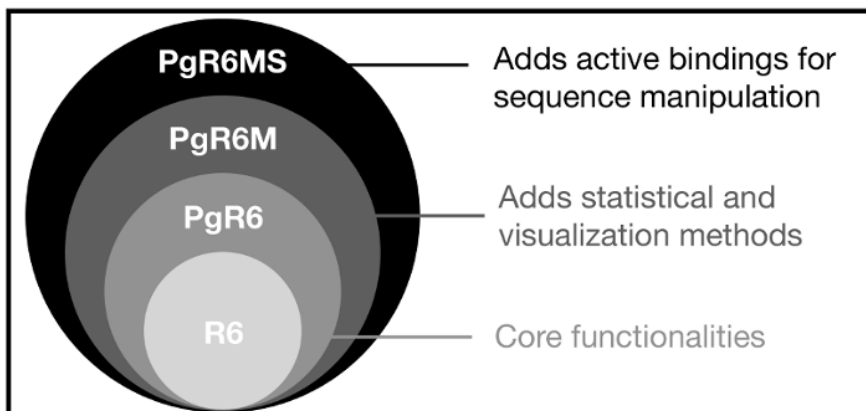


Input files to Pagoo

Pagoo object



Analysis



```
> p$summary_stats
DataFrame with 4 rows and 2 columns
  Category      Number
<character> <integer>
1 Total        3179
2 Core         1531
3 Shell        1355
4 Cloud         293
>
```

Save to single file

