

## Supplementary information for

### Comparative genomics of the genus *Pseudomonas* reveals host- and environment-specific evolution

Zaki Saati Santamaría<sup>1,2,3\*</sup>, Riccardo Baroncelli<sup>4</sup>, Raúl Rivas<sup>1,2,5</sup>, Paula García-Fraile<sup>1,2,5</sup>

<sup>1</sup>Departamento de Microbiología y Genética, Universidad de Salamanca, 37007 Salamanca, Spain.

<sup>2</sup>Institute for Agribiotechnology Research (CIALE), Villamayor, 37185 Salamanca, Spain.

<sup>3</sup>Institute of Microbiology of the Czech Academy of Sciences, Vídeňská 1083, 142 20 Prague, Czech Republic.

<sup>4</sup>Department of Agricultural and Food Sciences (DISTAL), University of Bologna, 40126 Bologna, Italy.

<sup>5</sup>Associated Research Unit of Plant-Microorganism Interaction, USAL-CSIC (IRNASA), 37008 Salamanca, Spain.

\*Corresponding author: [zakisaati@usal.es](mailto:zakisaati@usal.es)

#### External data:

There are files hosted at Zenodo (<https://doi.org/10.5281/zenodo.7105218>). This repository includes: (I) the pangenome rarefaction curve in html format; (II) the representative sequences of the protein clusters of the *Pseudomonas* pangenome; (III) A folder with the 3,274 genomes of our study and (IV) two large files which are the output from executing FastANI on the genome collection.

We also included bioinformatic codes and source data in the GitHub repository created for this article ([https://github.com/zakisaati/Pseudomonas\\_pangenome](https://github.com/zakisaati/Pseudomonas_pangenome)).

#### Content of this PDF:

Captions for supplementary information: page 2

Figure S1: page 3

Figure S2: page 4

Figure S3: page 5

Figure S4: page 6

Figure S5: page 7

## Captions

### Supplemental Figures

**Figure S1.** Boxplots of the gene content per genome and isolation source. Different letters represent groups with significant differences ( $p \text{ adj} < 0.05$ ).

**Figure S2.** Principal component analysis (PCA) based on the COG-term content of each of the 3,274 *Pseudomonas* genomes of this study. Colors/symbols represents the isolation source: the red circles indicate strains isolated from hosts (living beings) and the blue triangles, strains not isolated from hosts.

**Figure S3.** Boxplots of the content (per genome and isolation source) of encoded proteins related to some types of environmental stress resistance. Different letters represent groups with significant differences ( $p \text{ adj} < 0.05$ ). The table summarizes the number of each specific type of resistance in the *Pseudomonas* pangenome.

**Figure S4.** Intersection plot showing the number of significantly enriched proteins in each isolation category and those enriched in two or more categories (connected nodes).

**Figure S5.** Sequence similarity network (SSN) of the proteins significantly associated (green nodes) and not associated (red nodes) with hosts. Each node represents a cluster of proteins of the *Pseudomonas* pangenome. Those clusters of proteins that are connected with edges share high similarity. Thus, each group of connected nodes represents closely related clusters of proteins. In this case, the main groups encompass both HA and NHA proteins, so there are not big clusters of closely related HA or NHA proteins.

### Supplemental Files

**Supplemental File 1.** Datasheet that includes the following 6 supplementary tables:

**Table S1. Metadata, source and features of the 3,274 genomes of our collection.**

This table includes the name for each downloaded genome, its isolation source, the database containing this genome and its reference ID, and genomic features (n° contigs, CDSs, genome length and G+C% content).

**Table S2. Summary and characteristics of the *Pseudomonas* pangenome.**

Number of proteins included in each pangenome PPanGGOLiN partition.

**Table S3. Traits (isolation source) of the *Pseudomonas* genomes.**

Each isolation category represents a column in which "1" indicates that the genome is characterised by that trait and "0" indicates that it is not.

**Table S4.** Number of different COGs encoded by each of the 3,274 *Pseudomonas* genomes.

**Table S5.** Carbohydrate Active EnZYmes (CAZYs) annotations of the representative proteins of the *Pseudomonas* pangenome.

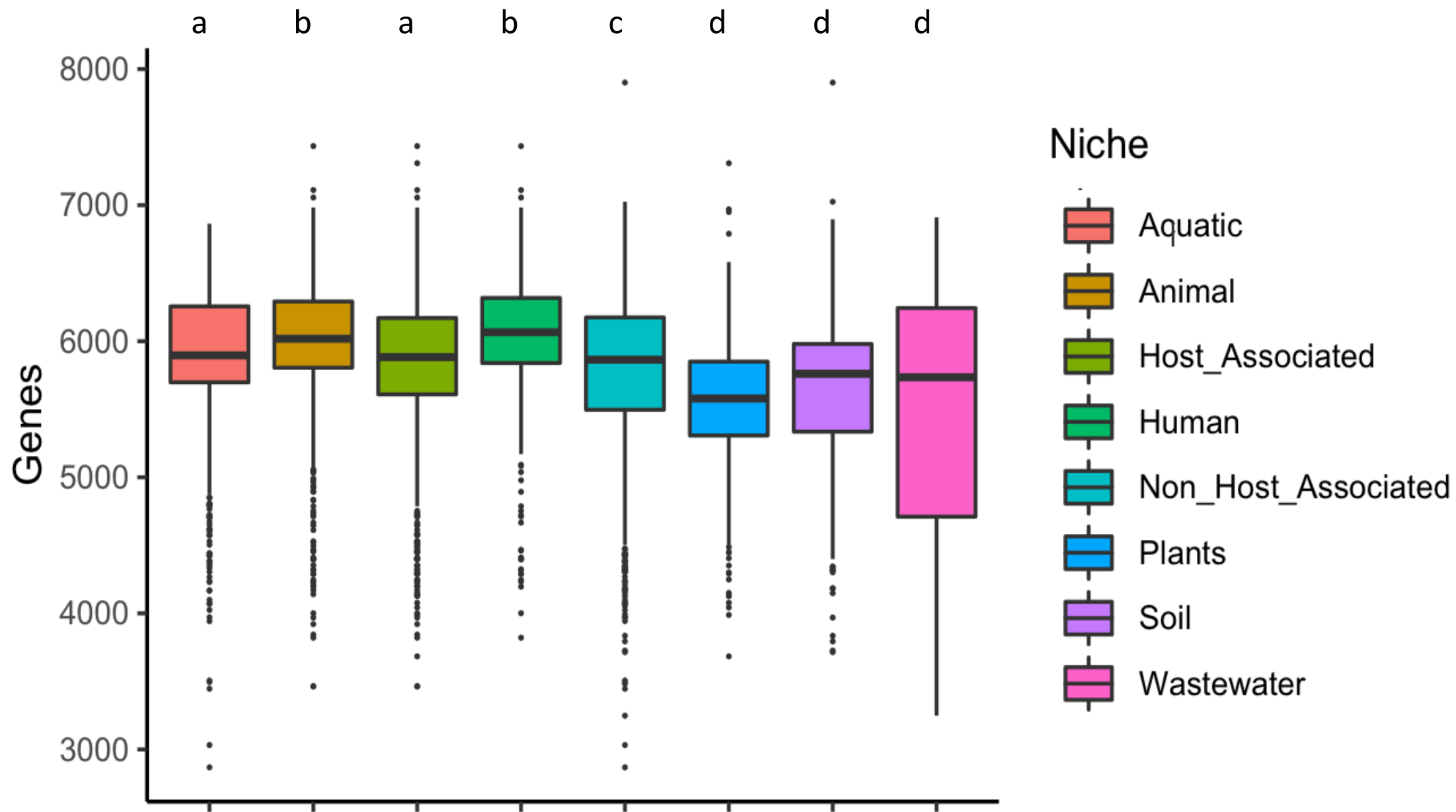
**Table S6.** Annotations of the stress-resistance-related proteins of the representative sequences of the *Pseudomonas* pangenome.

**Supplemental File 2.** Datasheet including the *Pseudomonas* protein relationships with the isolation source based on Scoary analyses.

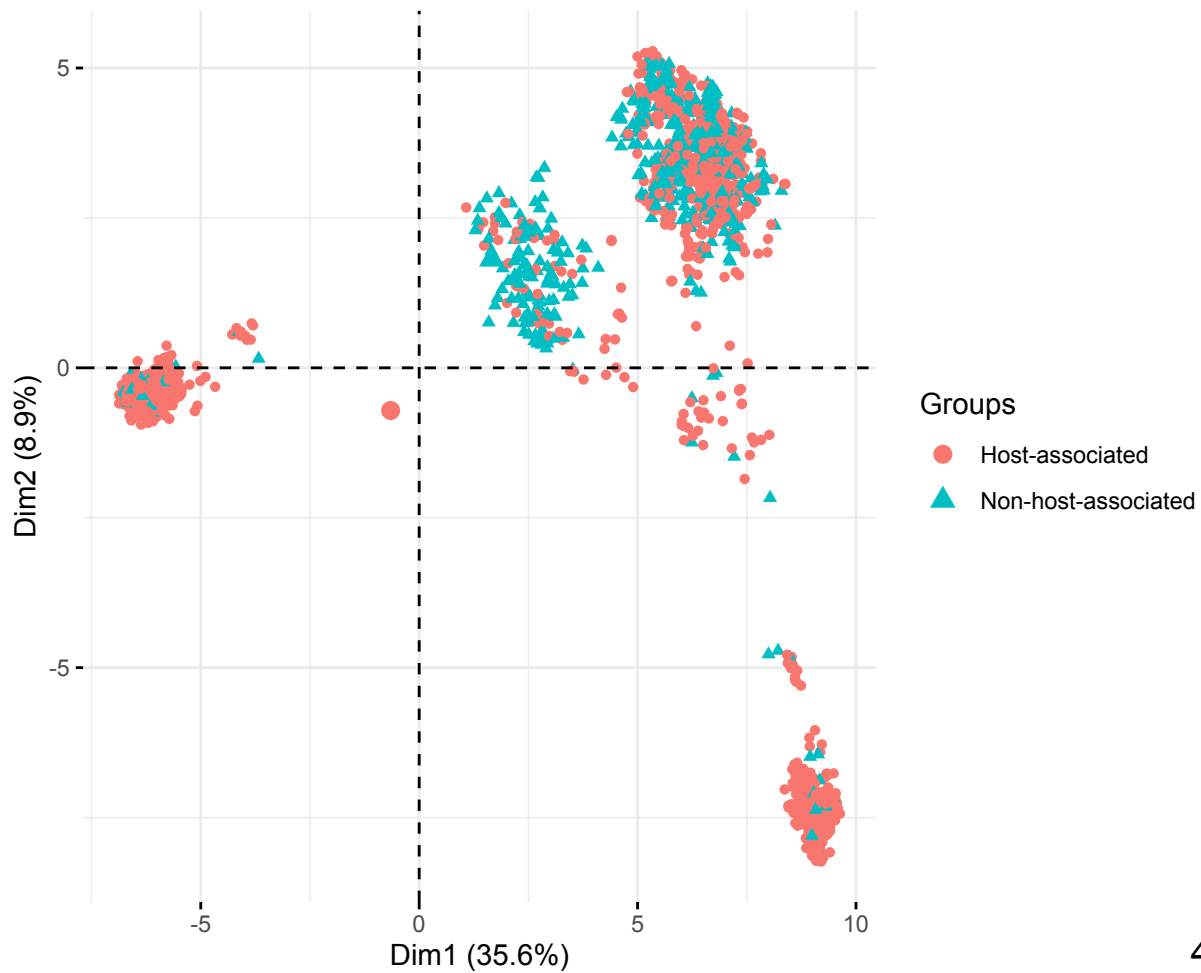
**Supplemental File 3.** Datasheet including the *Pseudomonas* Clusters of Orthologous Groups (COGs) relationships with the isolation source based on Scoary analyses.

**Supplemental File 4.** Datasheet including the *Pseudomonas* Carbohydrate Active enZYmes (CAZys) relationships with the isolation source based on Scoary analyses.

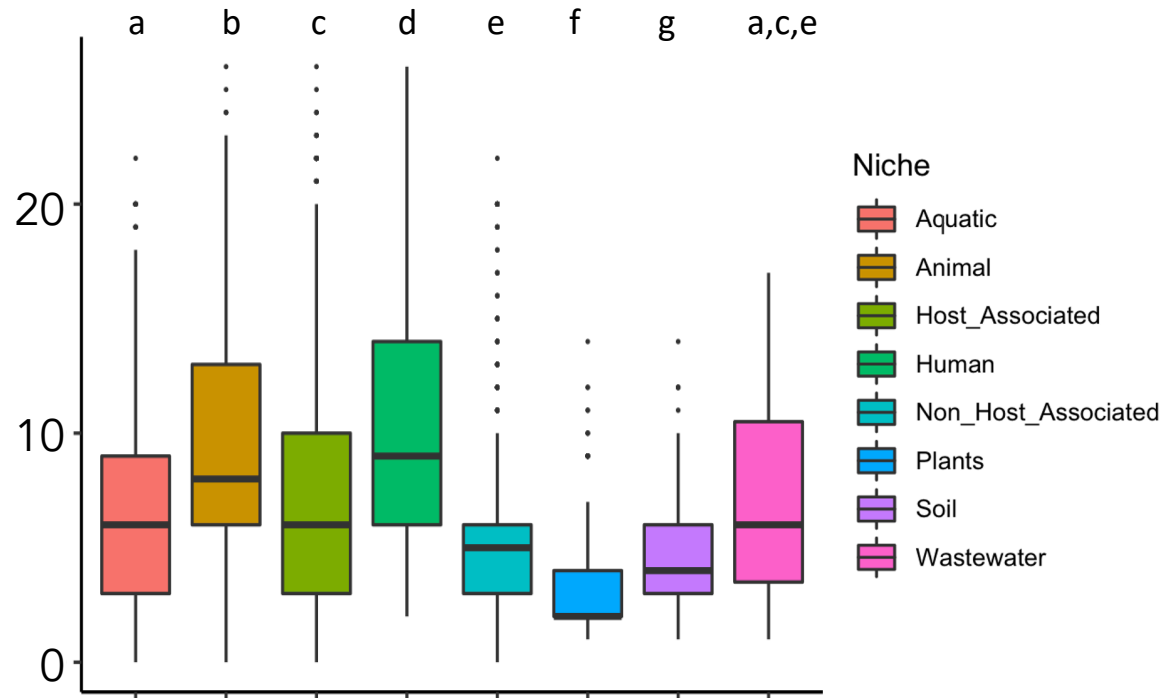
**Supplemental File 5.** Datasheet including the *Pseudomonas* resistance-related gene relationships with the isolation source based on Scoary analyses.



# Individuals - PCA



Resistance genes



Resistance mechanisms			
Antibiotics	Aminoglycoside		51
	Beta-lactam		105
	Bleomycin		
	Efflux pumps		
		Fluoroquinolones	1
		Fosfomycin	6
		Glycopeptides	1
		Macrolides	5
	Phenicols		13
	Quinolones		2
	Rifamycins		4
	Streptogramins		1
	Sulfonamides		3
	Tetracyclines		4
Metals	Sulfonamide		3
	Tetracycline		4
	Trimethoprim		1
Heat			7
Biocides			2

