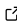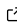# PANKEGG: Integrative Visualisation and Comparison of Metagenome-Assembled Genomes Annotation, Taxonomy, and Quality

**Renaud Van Damme** [1,4], **Arnaud Vanbelle** [1], **Tomas Klingström** [1], **Juliette Hayer** [2,3], **Amrei Binzer-Panchal** [1,4], and **Erik Bongcam-Rudloff** [1]

**1** Department of Animal Biosciences, Faculty of Center for Veterinary Medicine and Animal Science, Swedish University of Agricultural Sciences, Uppsala, Sweden **2** MIVEGEC, University of Montpellier, IRD, CNRS, Montpellier, France **3** Laboratoire Mixte International Drug Resistance in Southeast Asia **4** SLU Bioinformatics Infrastructure, Swedish University of Agricultural Sciences, Uppsala, Sweden

## Summary

To study microorganisms present in an environment, many tools are needed. Each of these tools generates a large number of results files, providing different insights. PANKEGG is an all-in-one tool that centralises these results files and helps researchers navigate and visualise them for comparison, interpretation, and understanding of microbial communities in the environment. PANKEGG supports researchers by parsing and visualising metagenome-assembled genomes (MAGs) and exploring their metabolic capabilities. It integrates quality metrics, annotation, and taxonomic classification in one interactive central database.

PANKEGG enables researchers to explore, compare, and interpret their data through a modern browser-based interface, streamlining the analysis of large and complex metagenomic datasets. The software supports output from widely used tools such as CheckM2 (Chklovski et al., 2022), GTDB-TK (Chaumeil et al., 2020), Sourmash (Brown & Irber, 2016), and EggNOG (Cantalapiedra et al., 2021), making it a flexible solution for a wide range of microbiome and environmental genomics studies. By interconnecting the different analysis results, researchers can investigate different interactions in the data more visually and conveniently.

PANKEGG answers questions such as:

- How many of my bins pass the GTDB (Parks et al., 2022) quality threshold?
- What is their taxonomic classification?
- Which Kyoto Encyclopedia of Genes and Genomes (KEGG) orthologs are present?
- Which proportion of their respective metabolic pathways is covered by the KEGG orthologs identified?

Answering these questions typically requires inspecting multiple result files and cross-referencing the information. In contrast, PANKEGG provides an all-in-one platform to explore, visualise, interpret, and save the results.

## Statement of need

The ever-growing progress of sequencing technologies has made it possible to recover thousands of draft and high-quality genomes directly from a plethora of environmental samples, accelerating our understanding of microbial diversity across ecosystems. Shotgun metagenomics with

assembly-based approaches recovers MAGs, giving access to taxonomic and functional profiles of uncultured microorganisms.

PANKEGG integrates taxonomic analysis with KEGG pathway annotations, distinguishing itself from tools like MAGFlow/BIgMAG (Yepes-García & Falquet, 2024) and Anvi'o (Eren et al., 2021). While MAGFlow/BIgMAG and Anvi'o metagenomics provide comprehensive metagenomic workflows and their visualisation, PANKEGG focuses on the visualisation and the functional interpretation of orthologs' variations across multiple samples and bins within the context of KEGG pathways (Kanehisa et al., 2023). Our focused approach on the KEGG orthologs facilitates a more direct and efficient analysis of metabolic capabilities and variations in microbial communities, even with growing datasets.

As the volume and complexity of metagenomic data increase, so do the challenges of efficiently comparing and visualising results from diverse annotation, classification, and quality assessment tools. In just one year (April 2024 to April 2025), over 135,000 new genomes were added to the Genome Taxonomy Database (GTDB). Tools like CheckM2, Sourmash, GTDB-TK, and EggNOG provide key outputs for quality, taxonomy, and functional annotation, but downstream integration and visualisation remain non-trivial.

PANKEGG enables users to merge results from any pipeline, workflow, or manual analysis that provides annotation, classification, and quality information into a standardised structured query language (SQL) database. The database allows users to explore the data through an interactive local web application. The tool is designed to analyse finalised MAGs and critically evaluate bins obtained during the binning stage of assembly-based metagenomic analysis. By integrating CheckM2 quality metrics, annotation, and taxonomic classification, PANKEGG helps users determine which bins meet the GTDB standards to be classified and reported as MAGs and which bins should be excluded due to low quality or inconsistency. PANKEGG allows the user to explore and compare the metabolic capabilities of microbial communities.

PANKEGG relies on widely used coding languages (Python, JavaScript, and HTML), SQLite as the SQL database engine (Hipp, 2000--2024), and libraries:

- flask (Pallets, 2024b)
- jinja2 (Pallets, 2024c)
- pandas (The pandas development team, 2024)
- numpy (Harris et al., 2020)
- SciKit-Learn (Pedregosa et al., 2011)
- SciPy (Virtanen et al., 2020)
- click (Pallets, 2024a)
- Python SQLite3 (Python Software Foundation, 2024)

Making its installation straightforward in most systems through pip (The pip developers, 2008), conda (Conda contributors, 2012), and pixi (Arts et al., 2023), see the PANKEGG installation chapter in our documentation. The software installation was tested on Ubuntu, Windows Subsystem for Linux (Ubuntu 16 to 22), Windows 10 & 11, macOS, and HPE Cray EX supercomputer systems. PANKEGG is based on common, reliable metagenomic annotation tools with standardised, stable output, ensuring compatibility over a long time span. However, should the tools and their output evolve, we plan to update PANKEGG to match new standards and keep it current. This unified approach reduces the barriers to integrative metagenomic analysis, enabling both specialists and non-specialists to make informed decisions based on large-scale, genome-resolved metagenomic data.

## Tool Overview

PANKEGG consists of two primary tools:

## PANGEGG MAKE DB

This script ingests a comma-separated values file specifying the locations of EggNOG annotations, classification (Sourmash or GTDB-TK), and quality metrics (Checkm2) files for each sample. It then constructs an SQL database aggregating all relevant results.
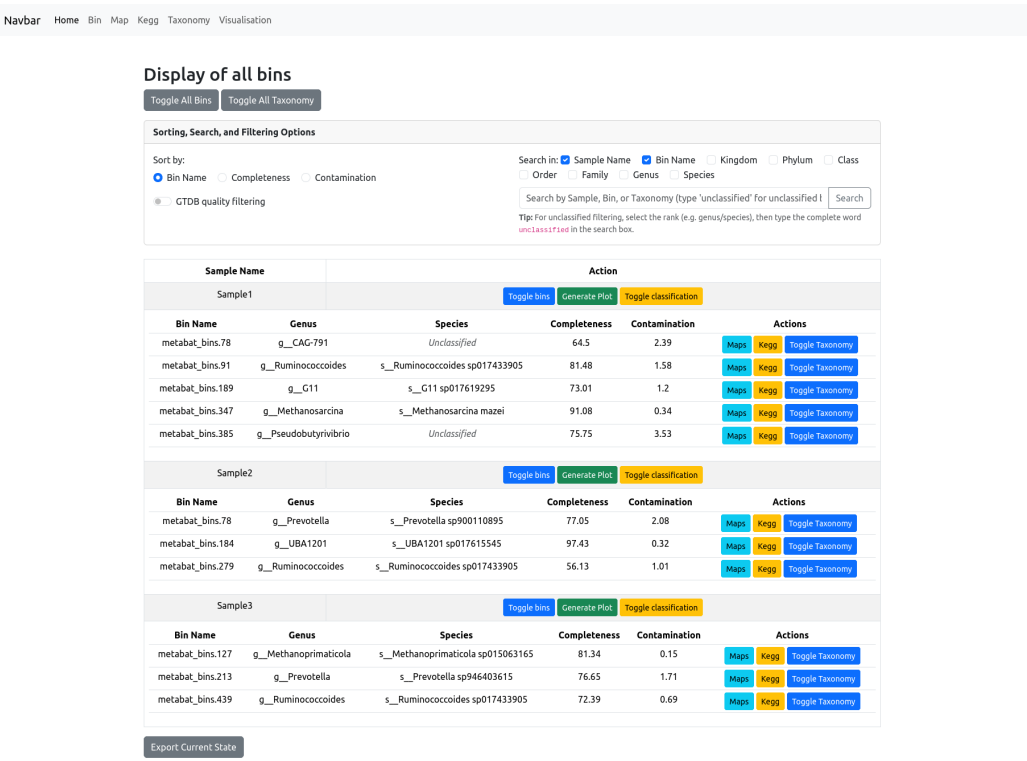
## PANKEGG APP

The app is a Flask-based web server that connects to the SQL database generated by PANKEGG_make_db.py, providing a simple interface for interfacing, cross-referencing, and filtering the data. The interface offers many different views of the data, and each page can be filtered by information from the other pages. A more detailed explanation is in our documentation's PANKEGG Web Page chapter. Here is a concise summary:

1. Bin page: Outlines all bins/MAGs in the database. Figure 1
2. Pathway page, also called Map page: Lists the pathways in the database. Each pathway contains a "completion value". This indicator is calculated by dividing the number of KEGG orthologs in the user's database by the total number of orthologs in the pathway. Figure 2
3. KEGG page: Lists the KEGG orthologs in the data. The information for each ortholog is expandable, and the view then includes the corresponding EggNOG entries (bin ID, sample ID, GO terms, KEGG orthologs associated, EggNOG description).
4. Taxonomy page: Presents the taxonomic composition of the input database.
5. Sample vs. Sample, Bins vs. Bins pages: Allow users to compare different samples or bins, respectively, with regard to pathway presence, completeness, and quality metrics.
6. PCA page: Visualises a principal component analysis (PCA) based on functional or taxonomic profiles. Beware that PCA interpretation is only valid with enough data; we recommend at least 40 bins (Shaukat et al., 2016).

PANKEGG's app is designed to make exploration intuitive. It features sortable and filterable tables, interactive plots, and external links to KEGG and other databases. These features collectively support users in hypothesis generation, genome curation, and discovery of ecological and functional trends in complex datasets.

# Figures



**Figure 1:** Bin page: This page shows the bins for each sample in the database. Three samples, with eleven bins, are displayed, including their classification, CheckM2 completeness, and contamination.



**Figure 2:** Map page, On this page, we see a list of maps from Samples 1 and 3 filtered for pathways containing the word 'metabolic.' Only one pathway is visible here, which is 15.43% complete. Below, a list of KEGG orthologs detected in the samples for this pathway is displayed.

# Author contribution

# Acknowledgements

# References

Arts, R., Zalmstra, B., Vollprecht, W., de Jager, T., Morcotilo, N., & Hofer, J. (2023). *Pixi: A modern package and environment manager for python and conda*. https://prefix.dev/docs/pixi/

Brown, C. T., & Irber, L. (2016). Sourmash: A library for MinHash sketching of DNA. *The Journal of Open Source Software*, *1*(5), 27. https://doi.org/10.21105/joss.00027

Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P., & Huerta-Cepas, J. (2021). eggNOG-mapper v2: Functional annotation, orthology assignments, and domain prediction at the metagenomic scale. *Molecular Biology and Evolution*, *38*(12), 5825–5829. https://doi.org/10.1093/molbev/msab293

Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P., & Parks, D. H. (2020). GTDB-tk: A toolkit to classify genomes with the genome taxonomy database. *Bioinformatics*, *36*(6), 1925–1927. https://doi.org/10.1093/bioinformatics/btz848

Chklovski, A., Parks, D. H., Woodcroft, B. J., Tyson, G. W., & Hugenholtz, P. (2022). CheckM2: A rapid, scalable and accurate tool for assessing microbial genome quality using machine learning. *Genome Biology*, *23*(1), 260. https://doi.org/10.1038/s41592-023-01940-w

Conda contributors. (2012). *Conda: A system-level, binary package and environment manager running on all major operating systems and platforms*. https://docs.conda.io/projects/conda/

Eren, M. A., Kiefl, E., Shaiber, A., Veseli, I., Miller, S. E., Schechter, M. S., Fink, I., Pan, J. N., Yousef, M., Fogarty, E. C., Trigodet, F., Watson, A. R., Esen, O. C., Moore, R. M., Clayssen, Q., Lee, M. D., Kivenson, V., Graham, E. D., Merrill, B. D., … Willis, A. D. (2021). Community-led, integrated, reproducible multi-omics with anvi'o. *Nature Microbiology*, *6*, 3–6. https://doi.org/10.1038/s41564-020-00834-3

Harris, C. R., Millman, K. J., Walt, S. J. van der, Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., & others. (2020). Array programming with NumPy. *Nature*, *585*(7825), 357–362. https://doi.org/10.1038/s41586-020-2649-2

Hipp, R. D. (2000--2024). *SQLite*. SQLite Consortium. https://sqlite.org/consortium.html

Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., & Morishima, K. (2023). KEGG for taxonomy-based analysis of pathways and genomes. *Nucleic Acids Research*, *51*(D1), D587–D592. https://doi.org/10.1093/nar/gkac963

Pallets. (2024a). *Click (version 8.2.1)*. https://palletsprojects.com/p/click/.

Pallets. (2024b). *Flask (version 3.1.1)*. https://flask.palletsprojects.com/.

Pallets. (2024c). *Jinja2 (version 3.1.6)*. https://palletsprojects.com/p/jinja/.

Parks, D. H., Chuvochina, M., Rinke, C., Mussig, A. J., Chaumeil, P.-A., & Hugenholtz, P. (2022). GTDB: An ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Research*, *50*(D1), D785–D794. https://doi.org/10.1093/nar/gkab776

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, *12*, 2825–2830. https://scikit-learn.org/

Python Software Foundation. (2024). *sqlite3 — DB-API 2.0 interface for SQLite databases*.

Shaukat, S. S., Rao, T. A., & Khan, M. A. (2016). Impact of sample size on principal component analysis ordination of an environmental data set: Effects on eigenstructure. *Ekológia (Bratislava)*, *35*(2), 173–190. https://doi.org/10.1515/eko-2016-0014

The pandas development team. (2024). Pandas-dev/pandas: Pandas (version 2.2.3). *Zenodo*. https://doi.org/10.5281/zenodo.3509134

The pip developers. (2008). *Pip - the python package installer*. https://pip.pypa.io/

Van Damme, R., Hölzer, M., Viehweger, A., Müller, B., Bongcam-Rudloff, E., & Brandt, C. (2021). Metagenomics workflow for hybrid assembly, differential coverage binning, metatranscriptomics and pathway analysis (MUFFIN). *PLOS Computational Biology*, *17*(2), e1008716. https://doi.org/10.1371/journal.pcbi.1008716

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., Walt, S. J. van der, Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., … Contributors, S. 1.0. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in python. *Nature Methods*, *17*, 261–272. https://doi.org/10.1038/s41592-019-0686-2

Yepes-García, J., & Falquet, L. (2024). Metagenome quality metrics and taxonomical annotation visualization through the integration of MAGFlow and BIgMAG. *F1000Research*, *13*, 640. https://doi.org/10.12688/f1000research.152290.2