

Metaphor: facilitating the large-scale recovery of genomes from metagenomes



PRESENTER

Vinícius W. Salazar¹, Babak Shaban², Maria Quiroga², Robert Turnbull², Edoardo Tescari², Vanessa Rossetto Marcelino³, Heroen Verbruggen⁴, Kim-Anh Lê Cao¹

¹ School of Mathematics & Statistics, University of Melbourne, ² Melbourne Data Analytics Platform, ³ Hudson Institute of Biomedical Research, ⁴ School of Biosciences, University of Melbourne

@vinisalazar_



@vinisalazar

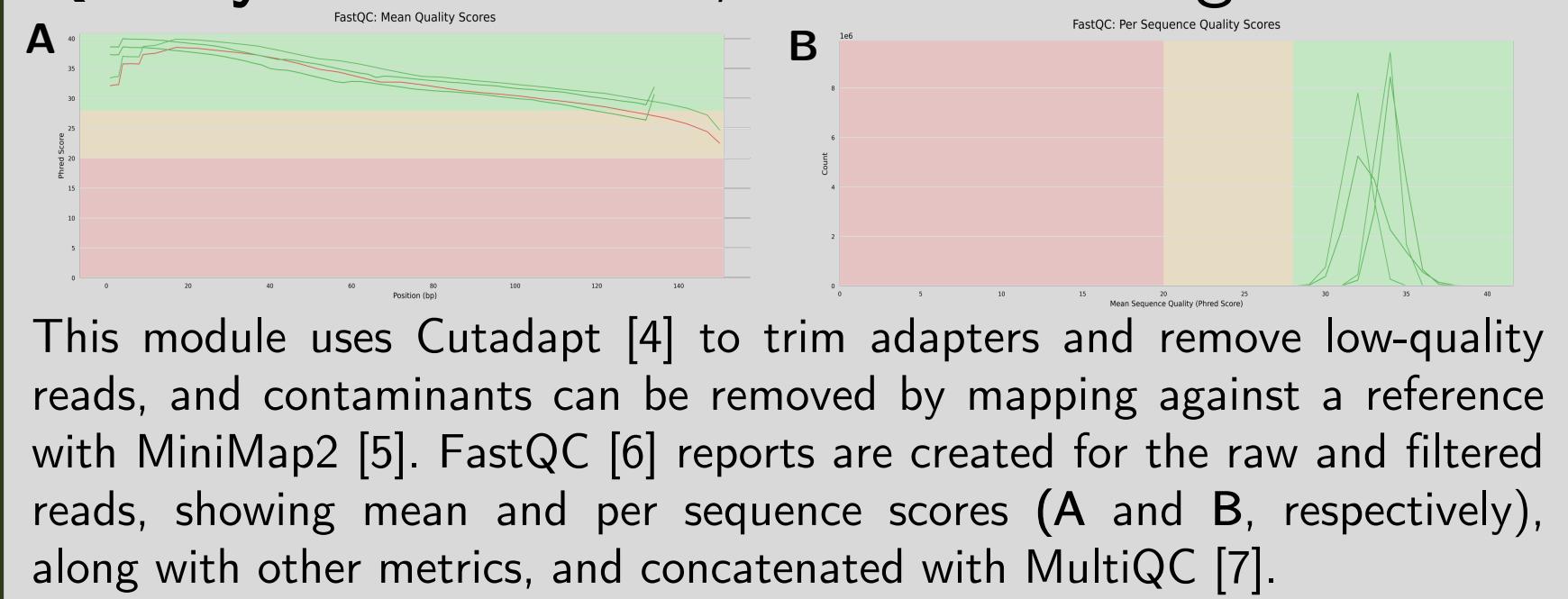
Metaphor is a general-purpose workflow for the assembly and binning of metagenomes. It advances existing metagenomic binning pipelines by the combination of two core features: the usage of multiple binning software along with a binning refinement step, and the possibility of defining groups for assembly and binning of samples. This latter feature allows users to pool samples together for the processes of assembly and binning (also called as coassembly or cobinning). We recommend doing this for samples which are biologically similar, *i.e.* that were sourced from the same community or environment, and as computational capacity allows. Metaphor is implemented with Snakemake [1], following best practices for workflow development, and the selection of software tools is informed by the Critical Assessment for Metagenome Interpretation [2]. To benchmark Metaphor, we ran it on simulated datasets with defined “ground-truth” reference genomes, allowing the comparison of different assembly and binning strategies. Metaphor is easily installable from Bioconda [3], compatible with various computing environments, and designed to be user-friendly, scalable, and automated.

General implementation: automate and scale your analysis with Snakemake.

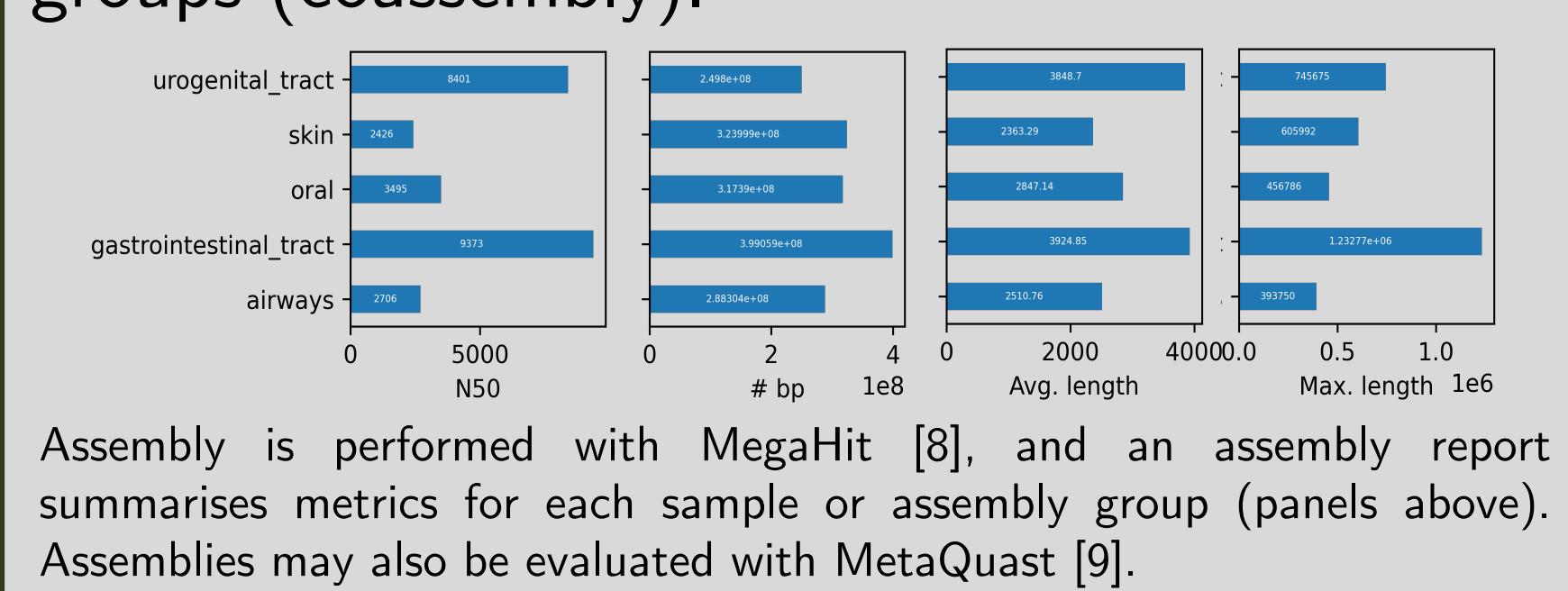
The workflow uses Snakemake, a scientific workflow management system, for calculation of a graph of jobs that can then be parallelized for faster execution. This can be done in multi-core machines, or combined with job schedulers such as SLURM and scaled to several computing nodes in high performance computing clusters. Metaphor is installed through the Bioconda software distribution channel, and Snakemake manages all dependencies automatically.



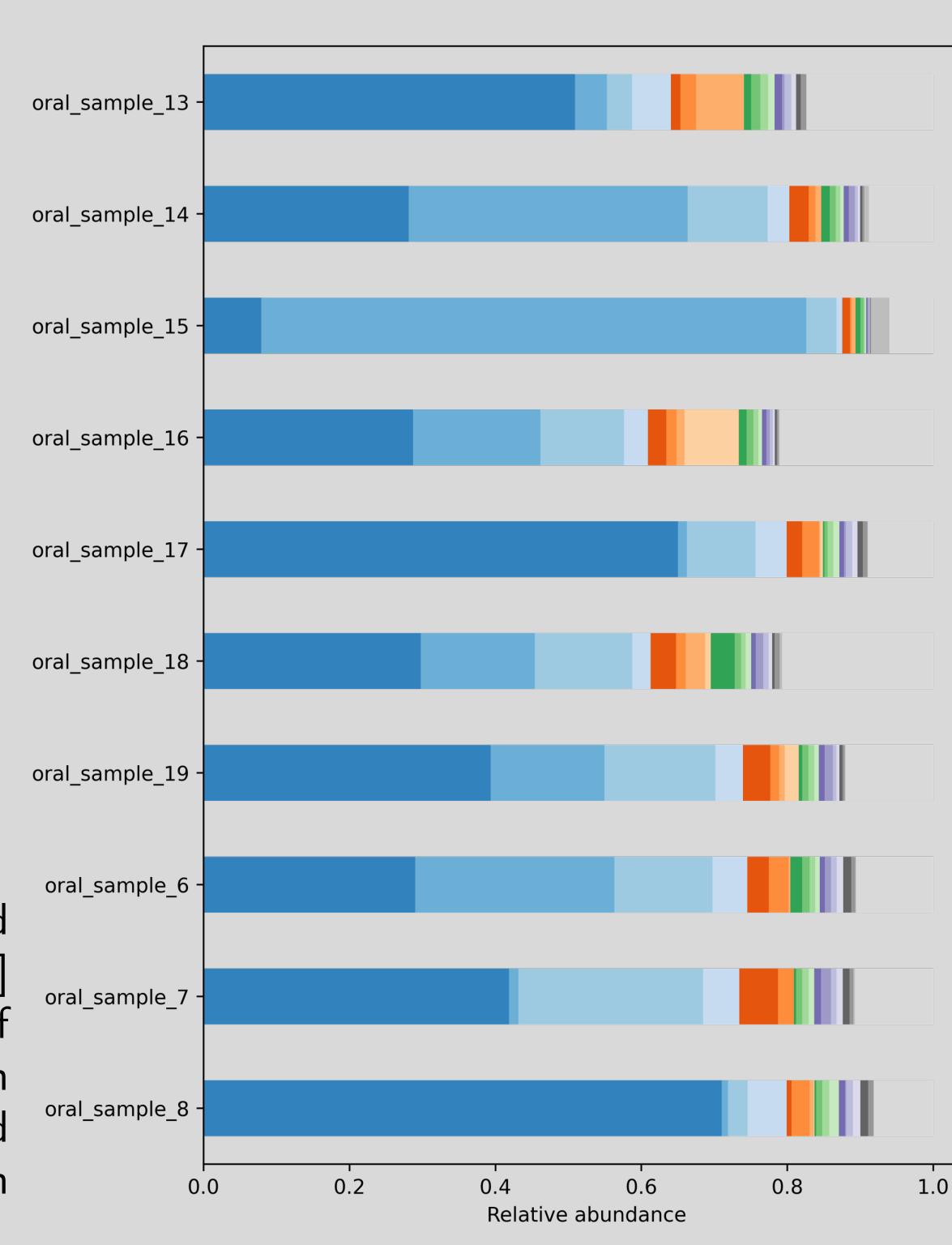
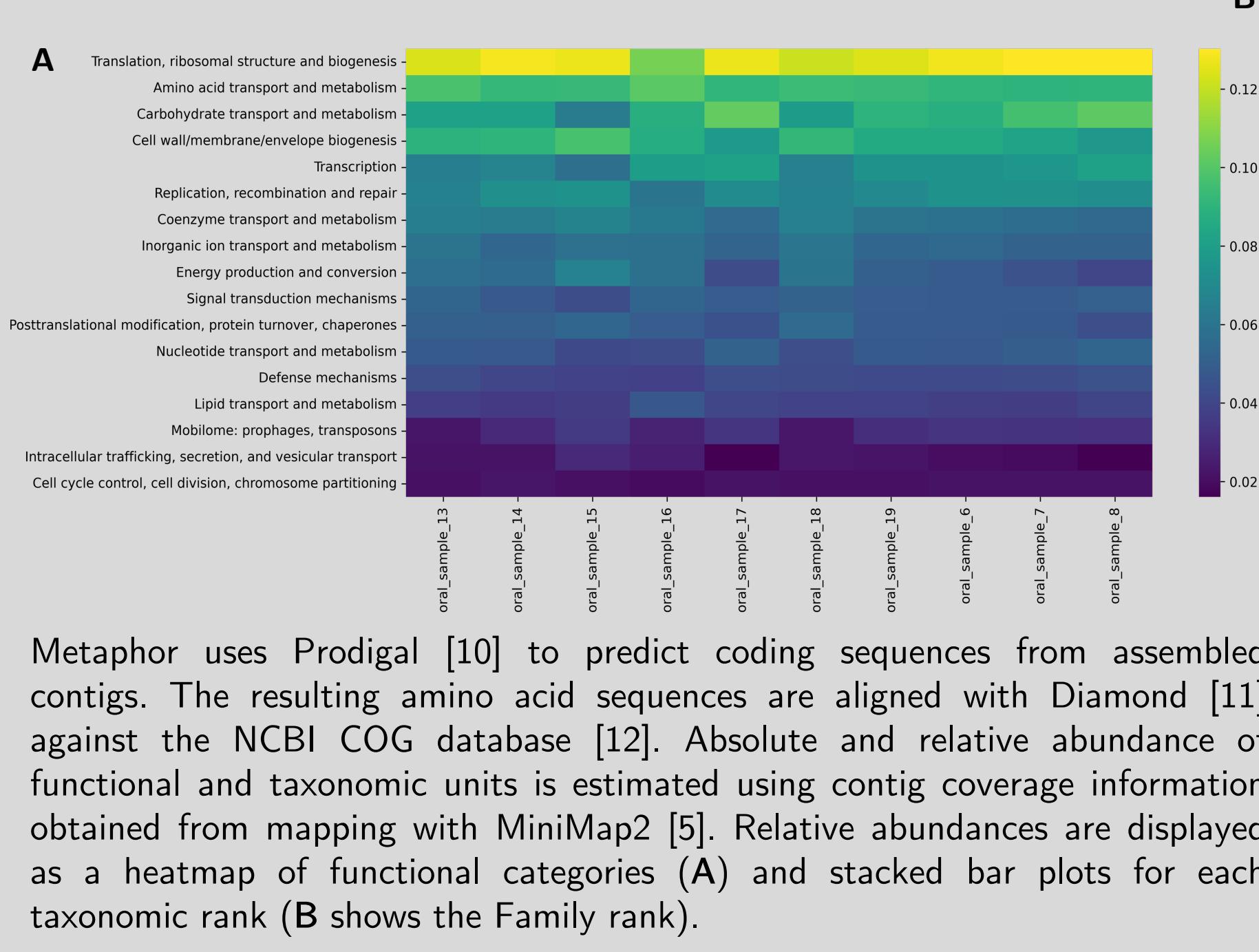
Quality Control: trim, filter and merge reads.



Assembly: assemble samples individually or in groups (coassembly).

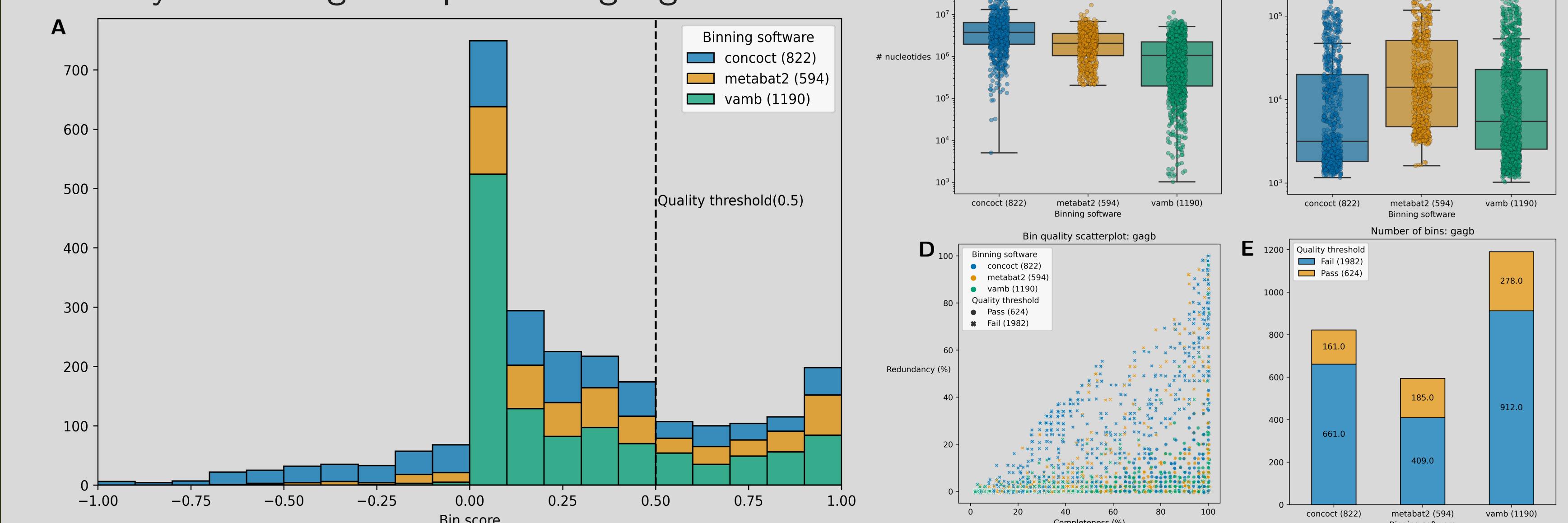


Annotation: obtain contig-level taxonomic and functional profiles.

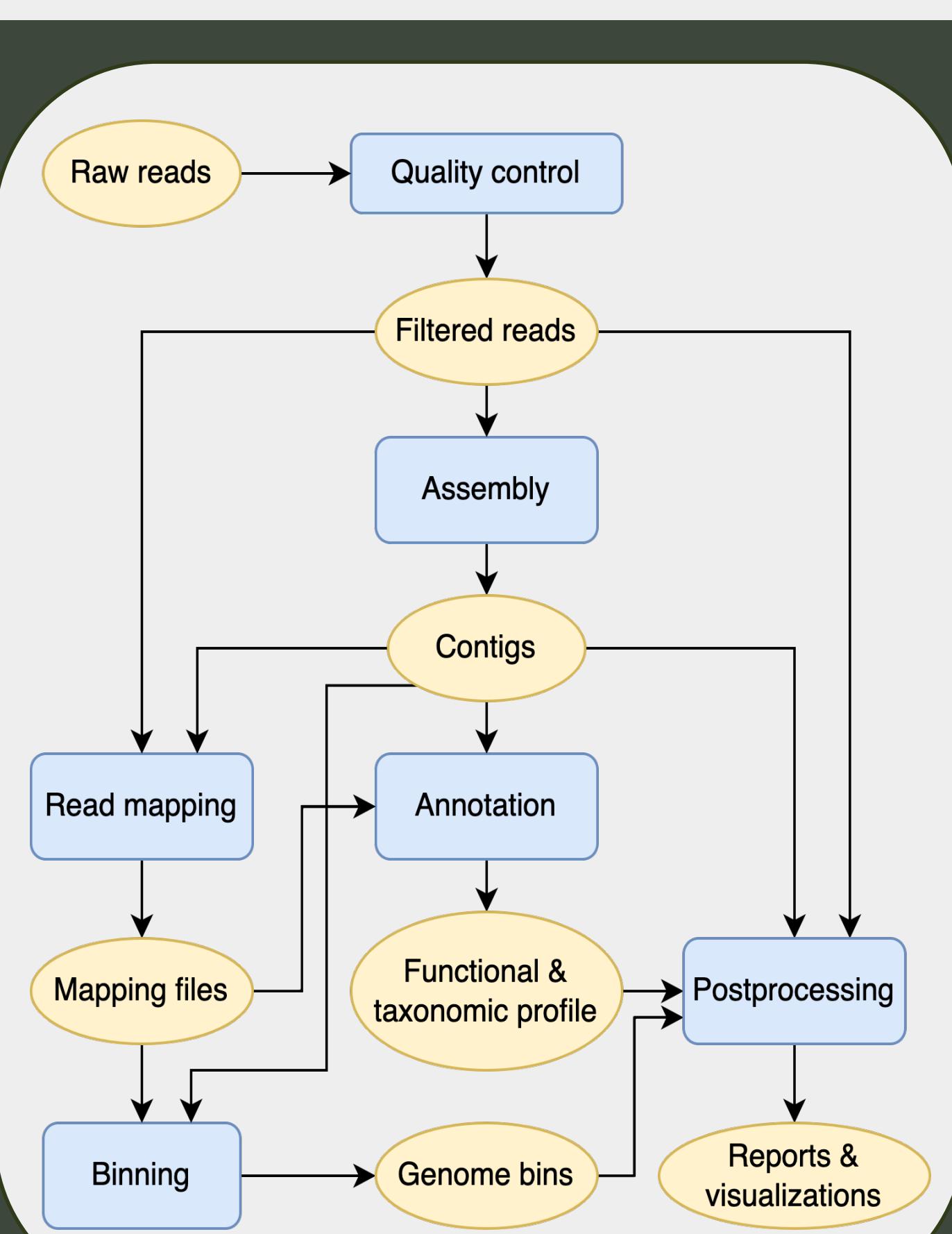
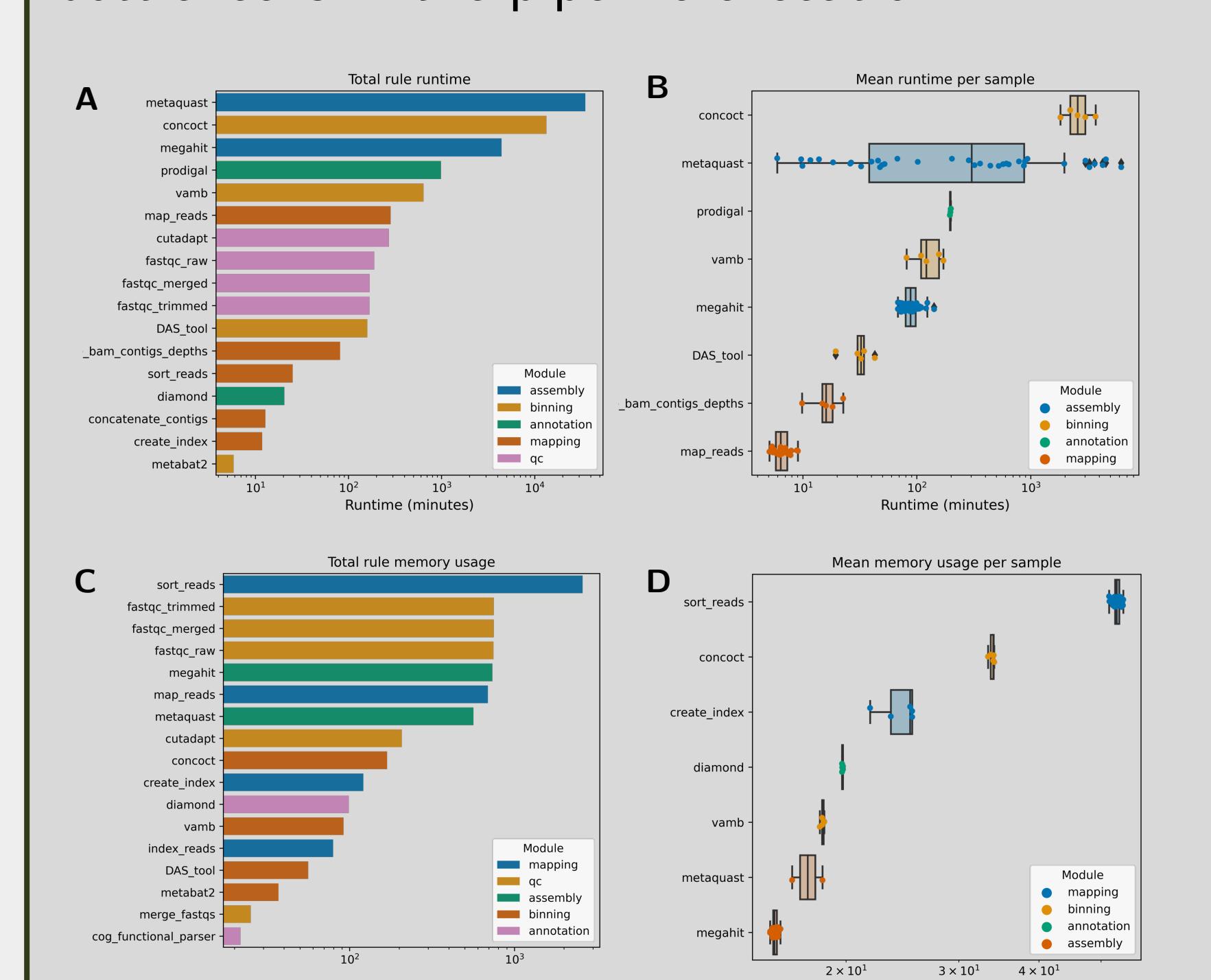


Family
Streptococcaceae
Neisseriaceae
Pasteurellaceae
Enterococcaceae
Enterobacteriaceae
Carnobacteriaceae
Clostridioidesceae
Bifidobacteriaceae
Flavobacteriaceae
Peptostreptococcaceae
Staphylococcaceae
Lachnospiraceae
Planctomycetaceae
Budviciaceae
Bacillaceae
Aerococcaceae
Lactobacillaceae
Oscillospiraceae
Spirosomaceae
Umbeliferaceae
Showing the 20 most abundant
Avg of 12.04% reads were filtered as low abundance.

Binning: generate high-quality and highly complete MAGs by combining multiple binning algorithms.

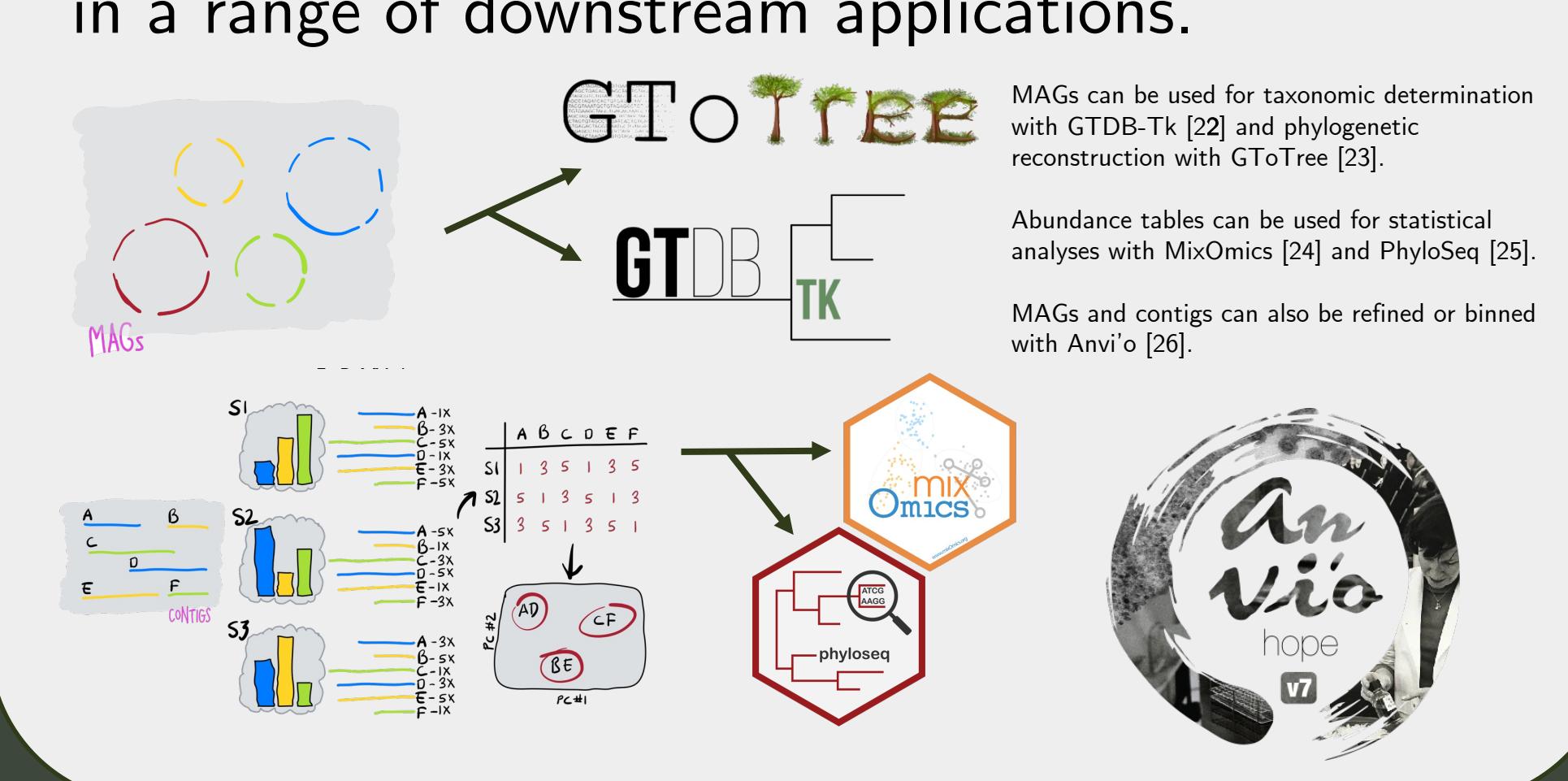


Postprocessing: identify computational bottlenecks in the pipeline execution.



Steps	Metaphor	ATLAS [18]	MetaWRAP [19]	nf-core/MAG [20]	MAGNETO [21]
Preprocessing					
Reads trimming	✓	✓	✓	✓	✓
Contamination	✓	✓	✓	✓	✓
Assembly					
Coassembly possible	✓		✓	✓	✓
Compute sets to coassemble	✓			✓	
Binning					
Cobinning possible	✓		✓	✓	✓
Multiple binning software	✓	✓	✓		
Bin refinement	✓	✓	✓		
Bin reassembly			✓		
Postprocessing					
MAGs quality check	✓	✓	✓	✓	✓
Dereplication step	✓	✓	✓	✓	✓
Genome annotation	✓	✓	✓	✓	✓
Gene catalogue	✓			✓	✓
Reproducibility					
Workflow management	✓	✓		✓	✓
Packages management	✓	✓		✓	✓

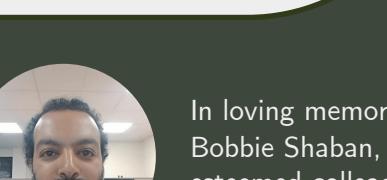
The output of Metaphor may be directly used in a range of downstream applications.



References

1. Metzger, M. et al. Sustainable data analysis. Preprint at <https://doi.org/10.1101/269> (bioRxiv 2020).
2. Meyer, L. et al. Critical Assessment of Metagenome Interpretation: the second round of challenges. *Nat Methods* 19, 429–440 (2022).
3. Grüning, B. et al. Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat Methods* 15, 471–476 (2018).
4. Kiefer, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *bioRxiv* 92, 10–11 (2014).
5. Li, H. Minimap2: fast genome remapping. *Bioinformatics* 35, 1019–1021 (2019).
6. Li, H. et al. FastQC: A Quality Control tool for High Throughput Sequence Data. *bioRxiv* 1–10 (2010).
7. Li, H. et al. MeGAHET: an ultra-fast single-node solution for large and complex metagenomic assembly via succinct de Bruijn graphs. *Bioinformatics* 31, 1674–1676 (2015).
8. Li, D. et al. MetaBAT 2: an adaptive binning strategy for robust and efficient genome reconstruction from metagenomic assemblies. *PeerJ* 10, 1–13 (2018).
9. Li, D. et al. MEGABAT 2: An adaptive binning strategy for robust and efficient genome reconstruction from metagenomic assemblies. *PeerJ* 10, 1–16 (2018).
10. Sieber, C. M. K. et al. Recovery of genomes from metagenomes via a deep learning approach. *Nature Methods* 15, 159–165 (2018).
11. Hsu, Y. et al. ATLAS: A Snakemake workflow for assembly, annotation, and genome binning of metagenomic sequence data. *Bioinformatics* 37, 1–8 (2021).
12. Lee, J. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
13. Chai, Y. et al. MAGNETO 2: An adaptive binning strategy for robust and efficient genome reconstruction from metagenomic assemblies. *bioRxiv* 1–13 (2019).
14. Albergaria, P. et al. Metabat 2: A flexible binning strategy for large and complex metagenomic datasets. *Nature Methods* 11, 141–146 (2014).
15. Kiefer, M. et al. Recovery of genomes from metagenomes via a deep learning approach and a scoring strategy. *Nature Methods* 15, 159–165 (2018).
16. Sieber, C. M. K. et al. Recovery of genomes from metagenomes via a deep learning approach and a scoring strategy. *Nature Methods* 15, 159–165 (2018).
17. Lee, J. et al. ATLAS: A Snakemake workflow for assembly, annotation, and genome binning of metagenomic sequence data. *Bioinformatics* 37, 1–8 (2021).
18. Lee, J. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
19. Ursenyk, G. V., DiRuggiero, J. B. & Taylor, J. Metabat2: a flexible pipeline for genome-retrieved metagenomic assembly. *Bioinformatics* 35, 10–12 (2019).
20. Hsu, Y. et al. ATLAS: A Snakemake workflow for assembly, annotation, and genome binning of metagenomic sequence data. *Bioinformatics* 37, 1–8 (2021).
21. Chai, Y. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
22. Chai, Y. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
23. Lee, J. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
24. Robert, F., Gauger, B., Singh, A. & Lê Cao, K. A. mixOmics: An R package for omics feature selection and multiple data integration. *PLoS Computational Biology* 17, e1012377 (2021).
25. Lee, J. et al. MAGNETO: An Automated Workflow for Genome-Rewritten Metagenomics. *bioRxiv* 44, 074322–074322 (2022).
26. Community-led, integrated, reproducible multi-omics with anvio | Nature Microbiology. <https://www.nature.com/articles/s41564-020-08834-3>.

This work was funded by a Melbourne Research Scholarship.



In loving memory of Bobbie Shaban, an esteemed colleague and friend.

Scan for poster and documentation.



THE UNIVERSITY OF MELBOURNE