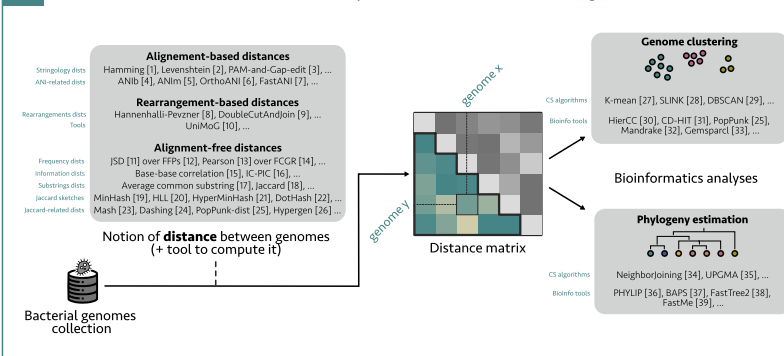


Subquadratic Storage of Pairwise Distance Matrices of Large Microbial Genome Collections

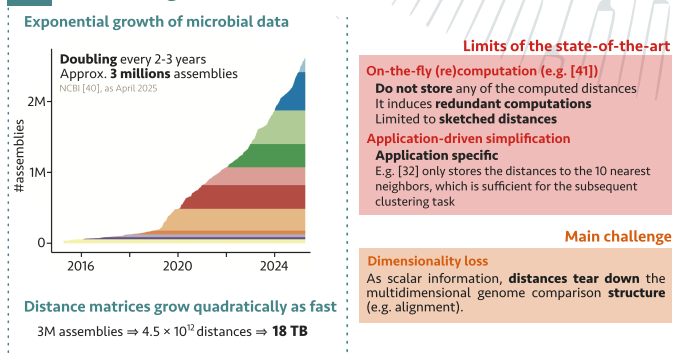
Léo Ackermann*, Pierre Peterlongo*, Karel Břinda*

*Univ. Rennes, Inria, CNRS, IRISA - UMR 6074, Rennes, F-350

C.1 Distance matrices underlie major bioinformatics analyses



C.2 Their storage has become infeasible

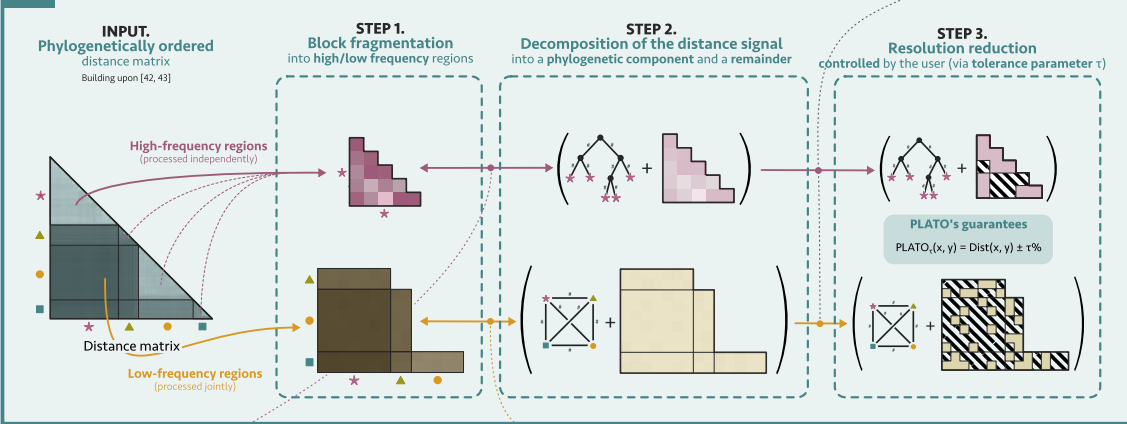


PLATO: Compression of microbial distance matrices with constant-time random access and controlled loss

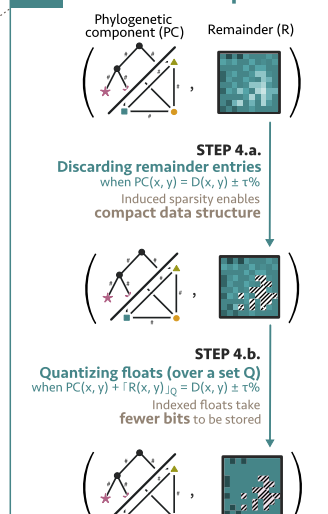
M.1 Key ideas

A **significant part** of the distance signal is **explained** by restricted classes of **phylogenies** that can be **efficiently estimated and stored** (steps 1 and 2)
Resolution needed by comparison-based downstream analyses is **relative**, allowing **efficient compression of the remaining part** of the signal (step 3)

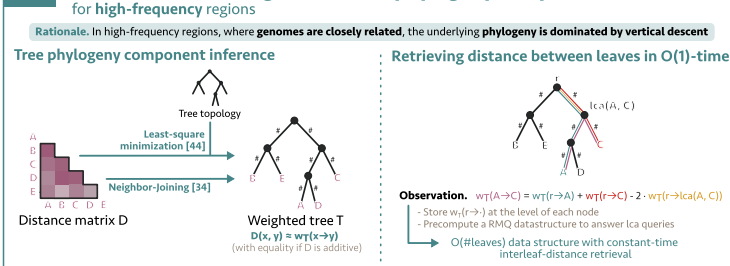
M.2 Workflow overview



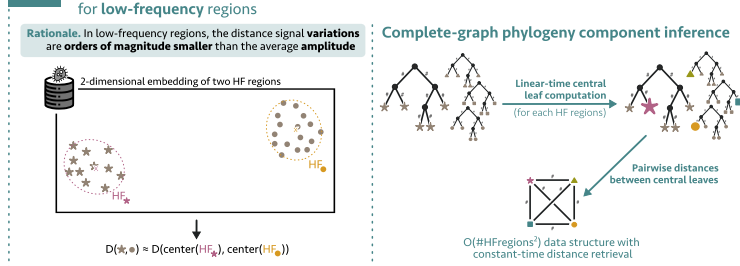
M.3 Remainder compression



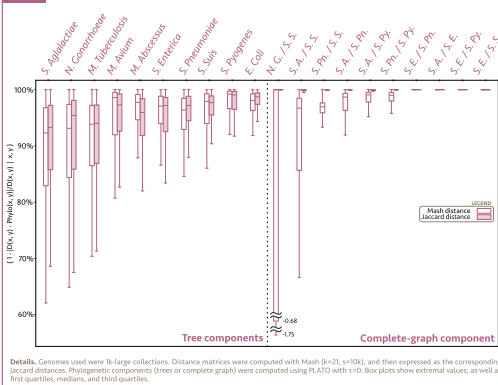
M.4 Estimation and storage of the tree phylogeny component for high-frequency regions



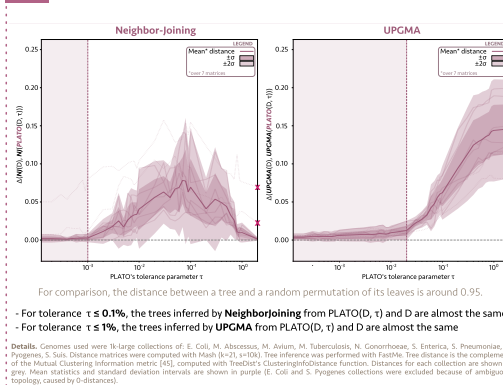
M.5 Estimation and storage of the complete-graph phylogeny component for low-frequency regions



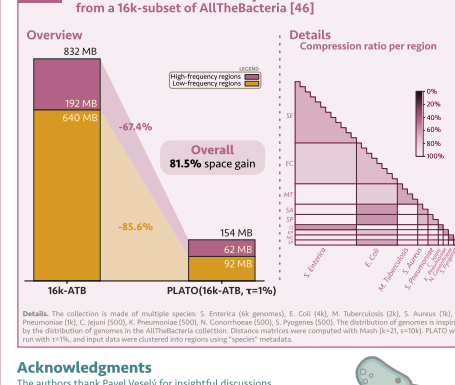
R.1 Phylogenetic components capture most of the distance signal



R.2 Tree inference is marginally impacted by small relative errors



R.3 APPLICATION: Compression of distances from a 16k-subset of AllTheBacteria [46]



Conclusion

We introduce **PLATO**, a method to **compress any distance matrix** up to a user-chosen **relative precision** while keeping **constant-time random access to entries**.

PLATO's core idea is to **estimate and store phylogenies** from which the distances can be recovered, **instead of the distances themselves**.

We use it to **reduce by 81.5% the storage** needed by a typical small collection of genomes, as a proof of concept.

Perspectives

Full characterization of the compression capability of the workflow, e.g. as a function of the diameter (CS) or the mutation rate (Bio) of the region.

Scaling to larger microbial collections (e.g. 661k [47]). This requires new subsampling and clustering strategies to stay efficient, e.g. during the tree phylogeny inference step.

Acknowledgments

The authors thank Pavel Vesely for insightful discussions. This work is supported by the ANR project REALL (ANR-24-CE45-1226).

References

- Hannenhall, Bill System Technical Journal, 1992
- Levenshtein, Soviet Physic Doklady, 1966
- Hoodman et al., Journal of Molecular Biology, 1970
- Coris et al., IJSEM, 2007
- Richter et al., PNAS, 2009
- Lee et al., IJSEM, 2016
- Jain et al., Nature Communications, 2018
- Hannenhall et al., FOCs, 1995
- Friedberg et al., Bioinformatics, 2008
- Holler et al., Bioinformatics, 2012
- Lin, IEEE Transactions on Information Theory, 1991
- Sims et al., PNAS, 2011
- Francis, Proceedings of the Royal Society of London, 1895
- Jeffrey, Nuclear Acids Research, 1990
- Li et al., BMC, 2008
- Guo et al., Gene, 2012
- Ullah et al., Journal of Computational Biology, 2006
- Jaccard, Bulletin de la Société Vaudoise des sciences naturelles, 1901
- Broder, Compression and Complexity of Sequences, 1997
- Faloutsos et al., Discrete Mathematics and Theoretical Computer Science, 2007
- Yu et al., IEEE Transactions on Knowledge and Data Engineering, 2012
- Nunes et al., Knowledge Discovery and Data Mining, 2013
- Chen et al., Genome Biology, 2016
- Baker et al., Genome Biology, 2019
- Levy et al., Genome Research, 2019
- Xu et al., Bioinformatics, 2014
- McQueen, Symposium on Mathematical Statistics and Probability, 1967
- Sibson, The Computer Journal, 1973
- Ester et al., Knowledge Discovery and Data Mining, 1996
- Zhou et al., Bioinformatics, 2013
- Fu et al., Bioinformatics, 2013
- Levy et al., Philos Trans R Soc Lond B Biol Sci, 2012
- von Kloten et al., Preprint, 2015
- Saitou et al., Molecular Biology and Evolution, 1987
- Sokal et al., University of Kansas Science Bulletin, 1958
- Felsenstein, Journal of Molecular Evolution, 1981
- Corander et al., Bulletin of Mathematical Biology, 2007
- Prie et al., PLoS ONE, 2010
- Larkin et al., Molecular Biology and Evolution, 2015
- Colclough et al., Nuclear Acids Research, 2015
- Inglis et al., RECOMB, 2016
- Břinda et al., Nature methods, 2015
- Hendrychová, MSc Thesis, 2025
- Legendre, 1805
- Smith, Bioinformatics, 2006
- Hunt et al., Preprint, 2015
- Blackwell et al., PLoS Biology, 2011