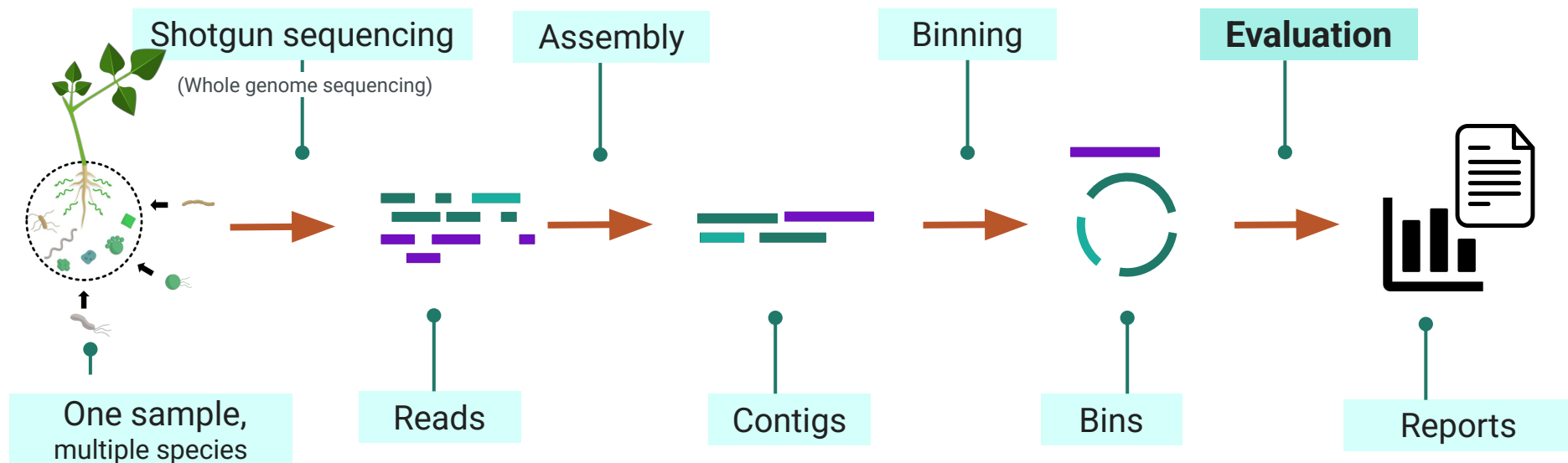




Mapler, a metagenome assembly and evaluation pipeline

Nicolas Maurice, Claire Lemaitre, Riccardo Vicedomini and Clémence Frioux

Metagenome assembly



Highly complex ecosystems

More species

- ➡ Lower coverage per species
- ➡ More shared sequences between species
- ➡ Lower quality assembly, less MAGs (Metagenome Assembled Genomes)

Highly complex ecosystems

More species

- ➡ Lower coverage per species
- ➡ More shared sequences between species
- ➡ Lower quality assembly, less MAGs (Metagenome Assembled Genomes)

- ◆ Mock community¹ : 18,0 Gbp, 17 species ➡ $\approx 265X$ / species⁵
- ◆ Human gut² : 18,8 Gbp, $>100^6$ species ➡ $< 50X$ / species
- ◆ Dead wood³ : 16,1 Gbp, $>1\ 000^7$ species ➡ $< 4X$ / species
- ◆ Soil⁴ : 81,2 Gbp, $>10\ 000^8$ species ➡ $< 2X$ / species

1 : ZymoD6331: SRR13128014

2 : Biocollective 139369 : SRR15275211

3 : Richy et al. (2024) : SRR28211698 to SRR28211701, coassembled

4 : Belliaro, Maurice et al. (2025)

5 : Assuming average genome ≈ 4 Mb

6 : Rowland et al. (2017)

7 : Pioli et al. (2023)

8 : Roesch et al. (2007)

Metagenome assembly evaluation

Can its bin be exploitable as MAGs ?



Completeness, contamination, contiguity...



Near complete, high / medium / low quality

- ◆ Single bin assessment
- ◆ Single-sample multiple assembly comparison

Many evaluation tools already exist^{1,2}
Poorly suits complex ecosystems with few MAGs

Metagenome assembly evaluation

Can its bin be exploitable as MAGs ?



Completeness, contamination, contiguity...

Near complete, high / medium / low quality

- ◆ Single bin assessment
- ◆ Single-sample multiple assembly comparison

Many evaluation tools already exist^{1,2}
Poorly suits complex ecosystems with few MAGs

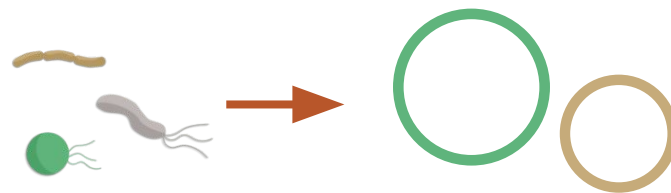
1 : PacBio HiFi-MAG-Pipeline ; 2 : CheckM

How representative of the sample is the assembly ?

How much of the diversity is captured by the assembly ?

What are the characteristics of that uncaptured diversity ?

Can such insight help us assemble better metagenomes ?



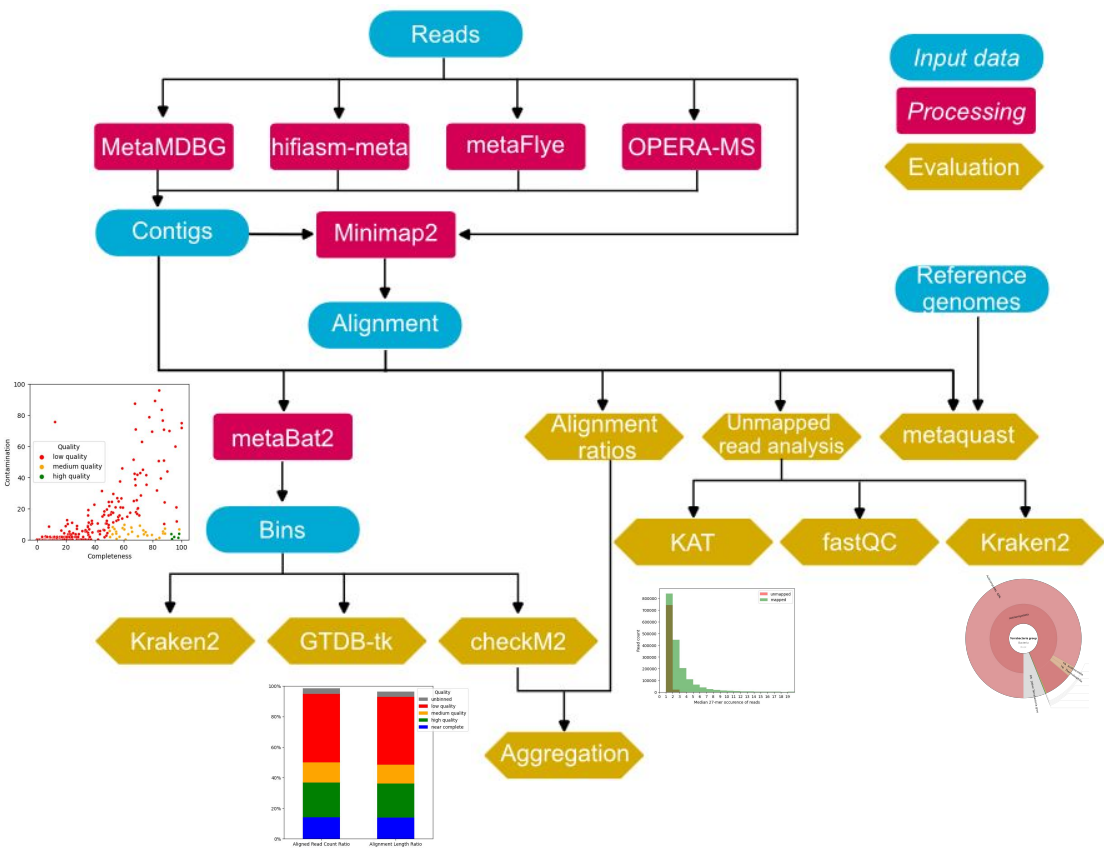
- ◆ Overall metagenome assessment i
- ◆ Multiple samples comparison

Lack of existing tools
Highly suits complex ecosystems with few MAGs

Mapler: a pipeline for assessing assembly quality in taxonomically rich metagenomes sequenced with HiFi reads



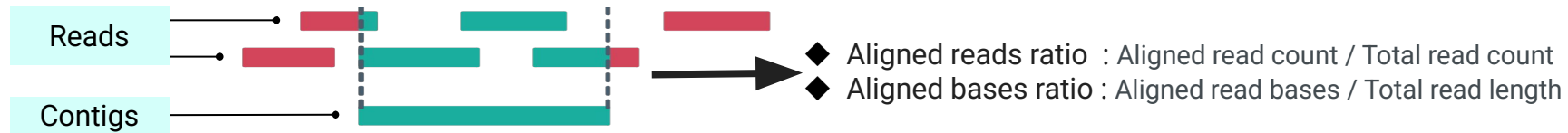
<https://doi.org/10.1093/bioinformatics/btaf334>



Bioinformatics, 2025

- ◆ Assembly & Binning
- ◆ Evaluation

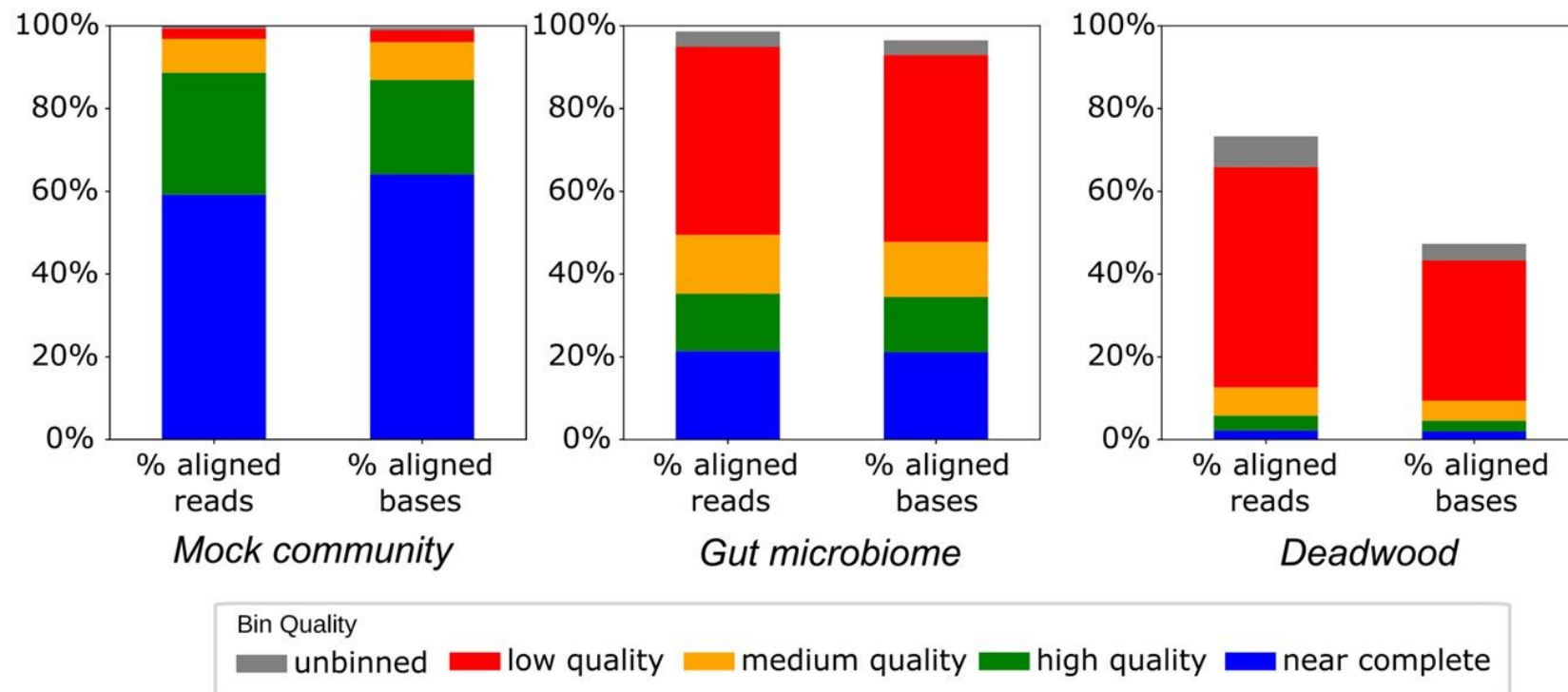
Alignment ratios

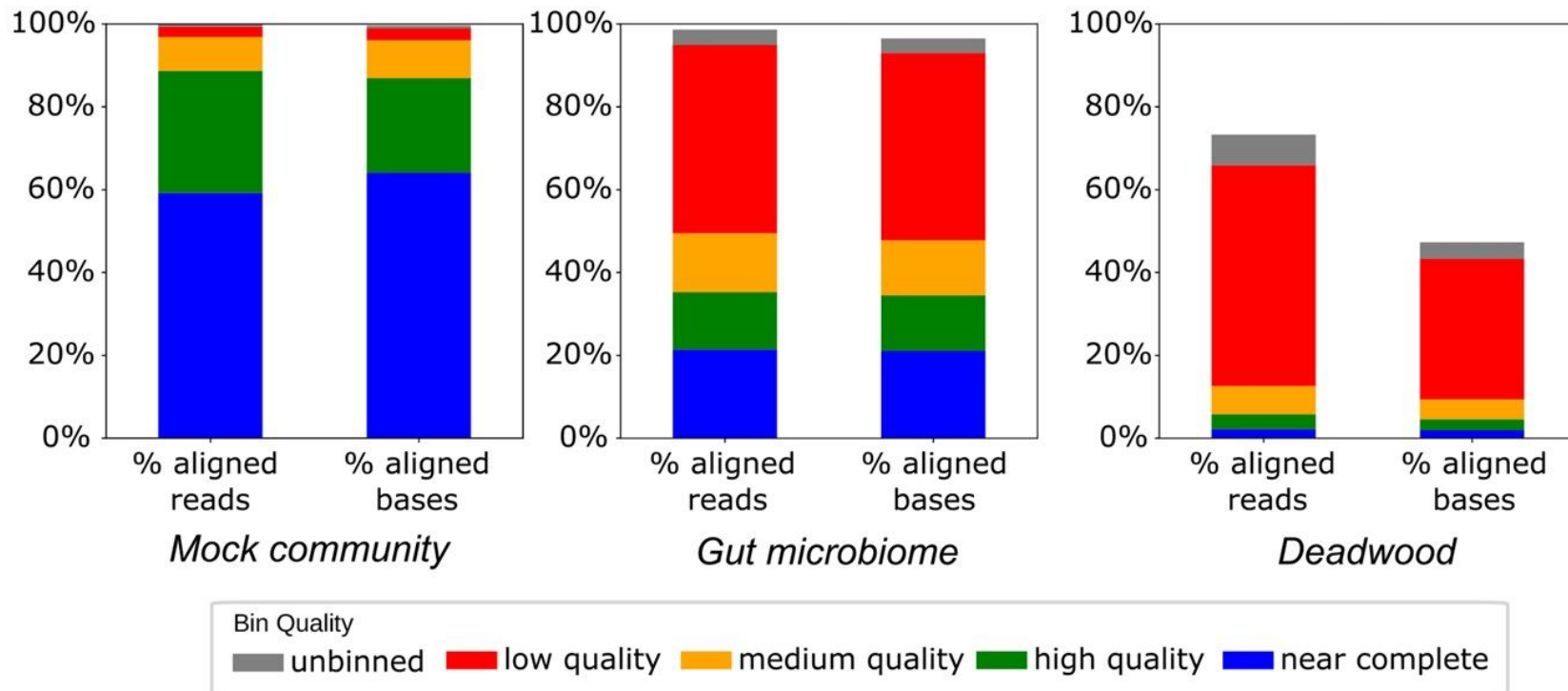


Better alignment => Better representation of the sample diversity

Aligned reads ratio >> Aligned bases ratio => Partially aligned reads

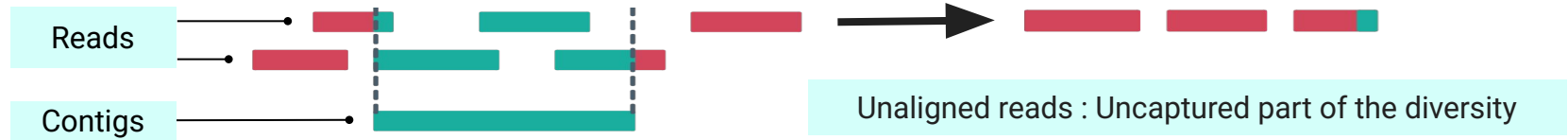
Alignment ratios of 3 metaMDBG¹ assemblies²



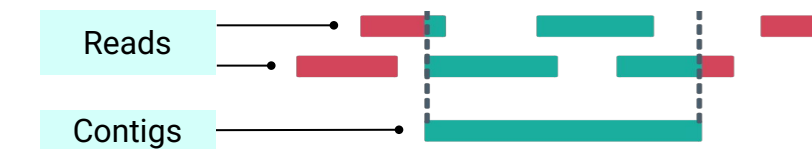
Alignment ratios of 3 metaMDBG¹ assemblies²

In highly complex ecosystem, most of the diversity is not captured by assembly

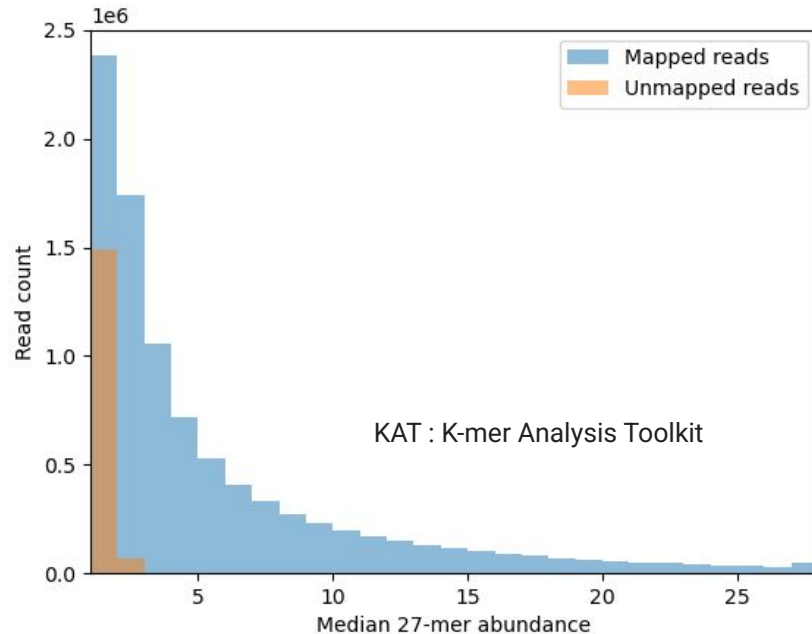
Uncaptured diversity



Uncaptured diversity in the soil sample



Unmapped reads : Uncaptured part of the diversity



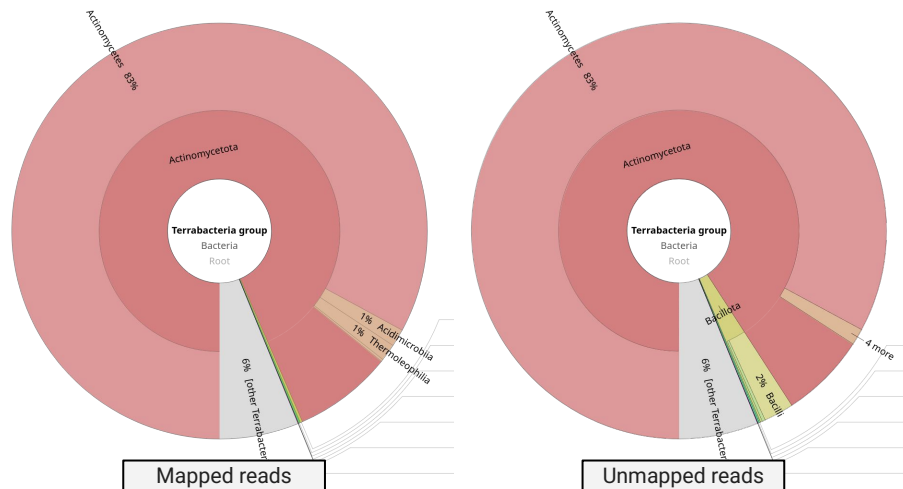
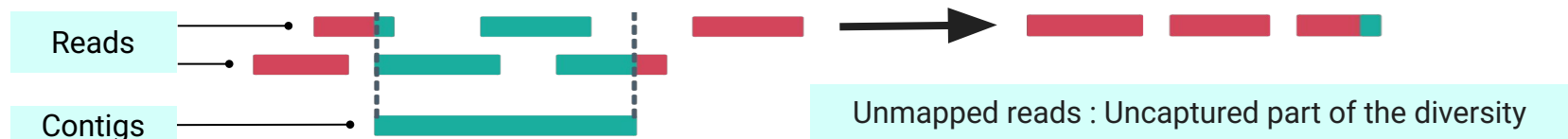
KAT : K-mer Analysis Toolkit

Belliardo, Maurice et al. (2025)

Compared to mapped reads, unmapped reads are :

- ◆ made up of rarer k-mers
- ➔ Higher sequencing depth can improve assembly
- ➔ k-mer abundance is not the only factor

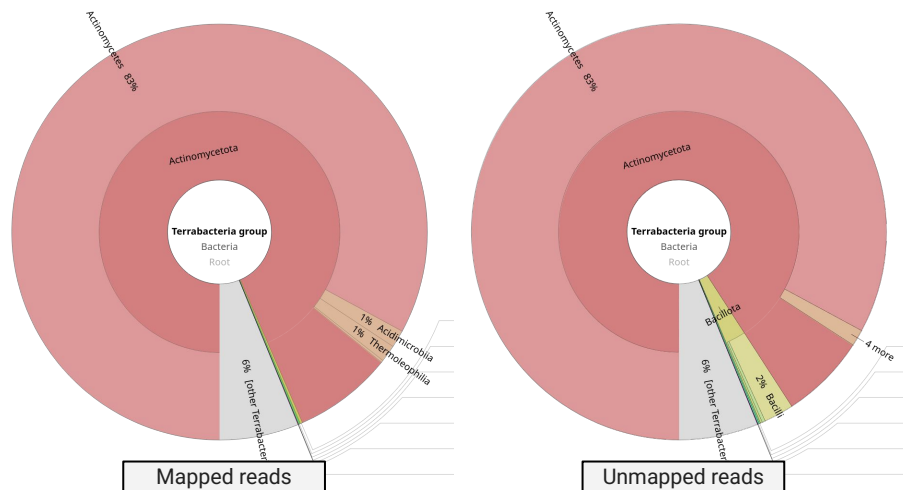
Uncaptured diversity in the soil sample



Compared to mapped reads, unmapped reads are :

- ◆ made up of rarer k-mers
 - ➔ Higher sequencing depth can improve assembly
 - ➔ k-mer abundance is not the only factor
- ◆ sometimes from under-represented taxa (Bacillota)

Uncaptured diversity in the soil sample



Compared to mapped reads, unmapped reads are :

- ◆ made up of rarer k-mers
 - ➔ Higher sequencing depth can improve assembly
 - ➔ k-mer abundance is not the only factor
- ◆ sometimes from under-represented taxa (Bacillota)
- ◆ of slightly lower length (6.3 vs 7.1 Kbp average)
 - ➔ Higher read length can improve assembly, although not the most significant factor

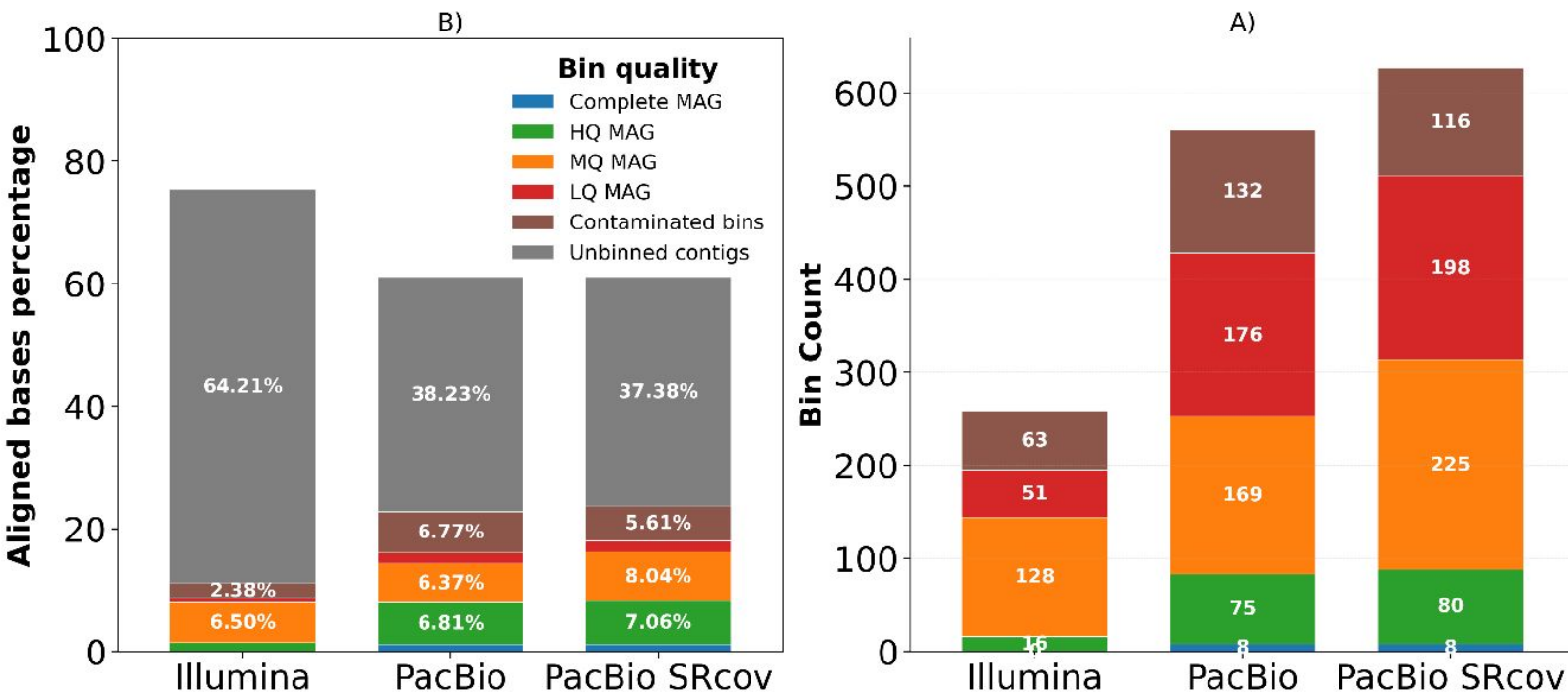
Application

Accurate MAG reconstruction from complex soil microbiome through combined short- and HiFi long-reads metagenomics



<https://doi.org/10.1101/2025.09.12.675765>

Mapler can highlight the diversity captured by Illumina assembly, even if poor contiguity prevents quality binning



In highly complex environments, most of the diversity cannot be captured

Mapler can assess the characteristics of that unassembled diversity :

- ◆ Low abundance populations
- ◆ Partially aligned reads

Which offer hindsight on how to improve assembly :

- ◆ Increase sequencing depth to reduce unique k-mers
 - Although, 18Mb of mock community reads match 20 Gb of soil reads in term of assembly quality
- ◆ Improve DNA extraction and sequencing for longer reads
- ◆ Remap reads on contigs to improve alignment
- ◆ Use short-reads for binning if insufficient long-reads coverage

Many thanks to



PEPR "Agroécologie Numérique"
MISTIC project

Clémence Frioux
Claire Lemaitre
Riccardo Vicedomini
Carole Belliaro
Etienne Danchin
Marc Bailly-Bechet