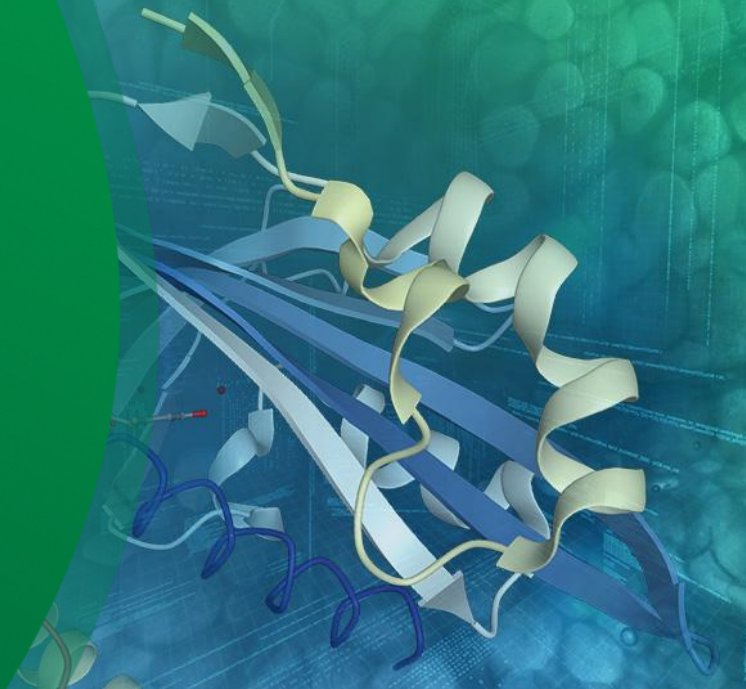


# Metagenomics



**Germana Baldi, Varsha Kale**

Bioinformaticians at MGnify

# Who are we?

# MGnify



Rob Finn  
Team Leader  
Section Head



Lorna Richardson  
Team Coordinator



Varsha Kale  
Bioinformatician



Germana Baldi  
Bioinformatician

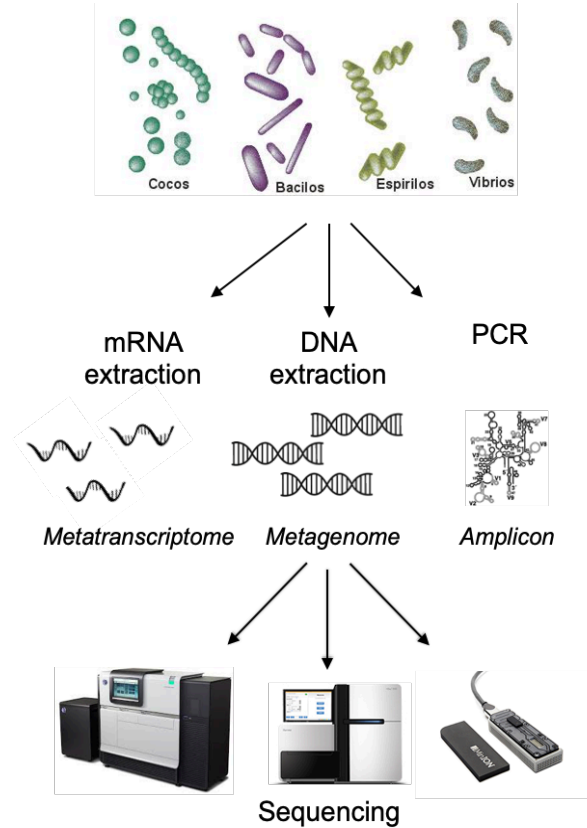
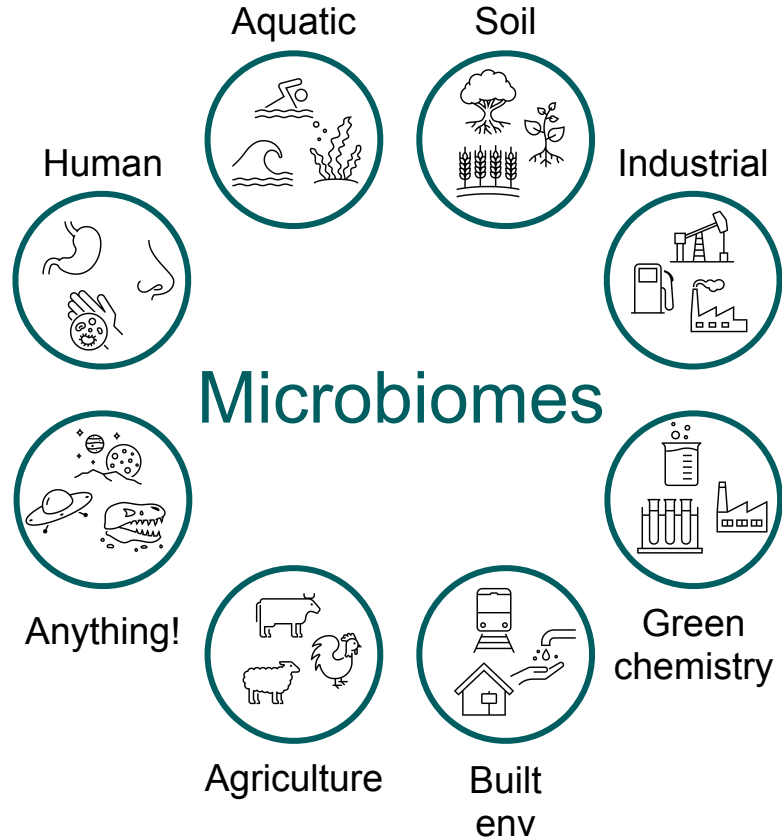


Sandy Rogers  
Website Developer

# What is metagenomics?

|          |                                      |  |
|----------|--------------------------------------|--|
| META     | “transcending”, “more comprehensive” | Transcends the individual organisms to focus on the community more comprehensively |
| GENOMICS | “the study of genomes”               |  |

# What is metagenomics?



# Sequencing Technology

## Illumina

- Short-reads (50-250 bp)
- 4-20k million reads per run
- Error rate <0.1%
- Sample preparation \$50 - \$100
- Run cost \$1k - \$4k

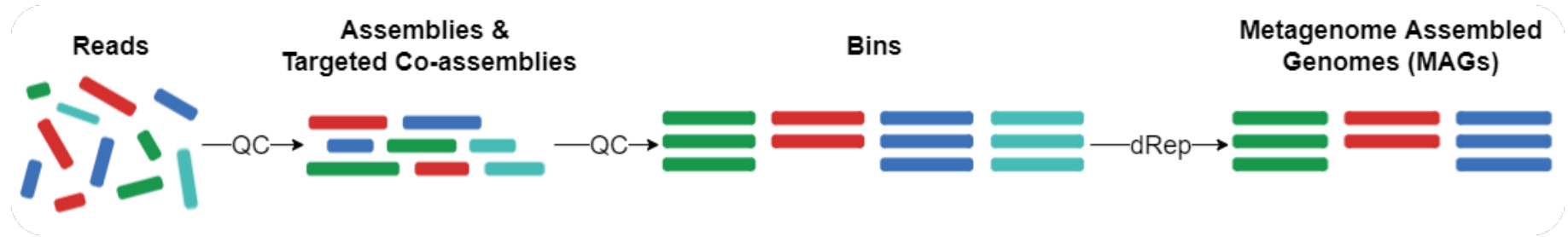
## Oxford Nanopore/PacBio SMRT

- Long-reads (>1 kbp)
- 20-100k reads per run
- Error rate 10-15%
- Sample preparation \$100 - \$500
- Run cost \$1k - \$2k

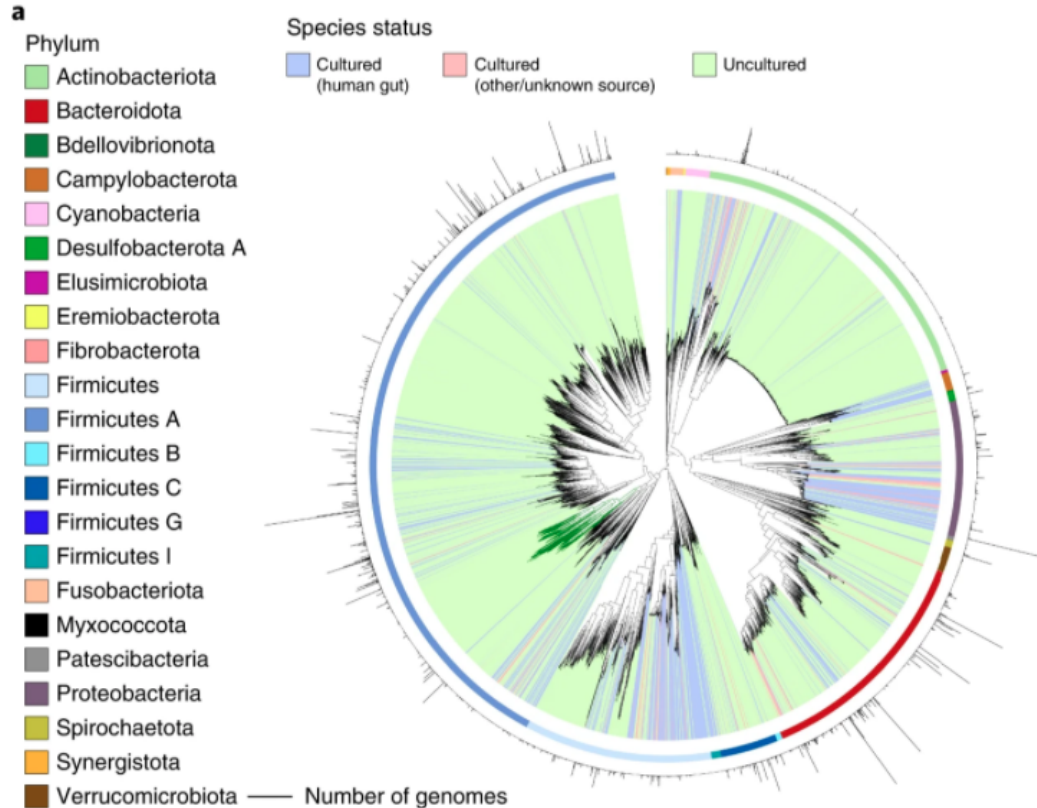


Sequencing

# Overview of MAG generation workflow



# Most gut species lack isolate genomes



# Things to note

- The majority of what we will cover is prokaryotic-focused (also for historical reasons!)
- Your sample may contain a lot of variety
  - Viruses
  - Eukaryotes
  - Prophages
  - Plasmids
  - ...
- There is no one correct answer to “how should I analyse my microbiome data?” (sorry!)

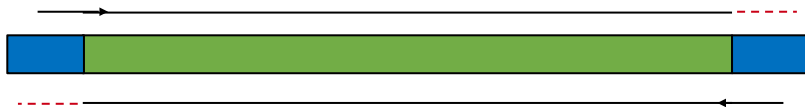


# Methods of quality control

Cleaning raw reads reduces the risk of contamination in downstream analyses.

You might want to consider:

- GC content
- Duplicates
- Trim reads/regions by quality score
- Sequence length
- Contamination from:
  - Human/host: map samples against reference
  - Reagents: use negative controls
  - Sequencing process: clipping known primers/adaptor sequences



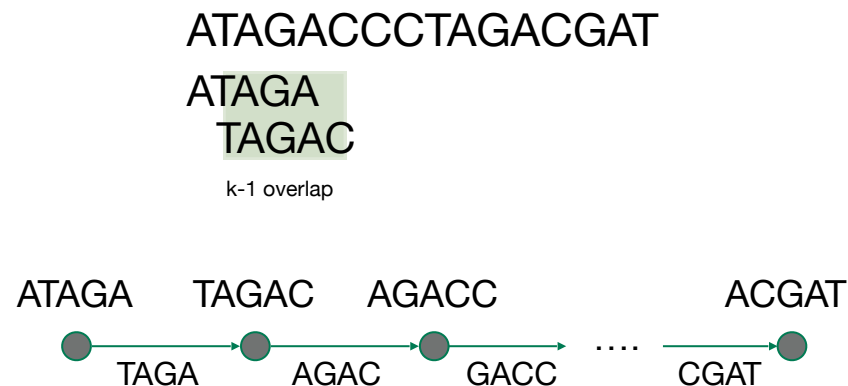
It is important to keep high quality controls throughout the whole workflow

# De novo assembly

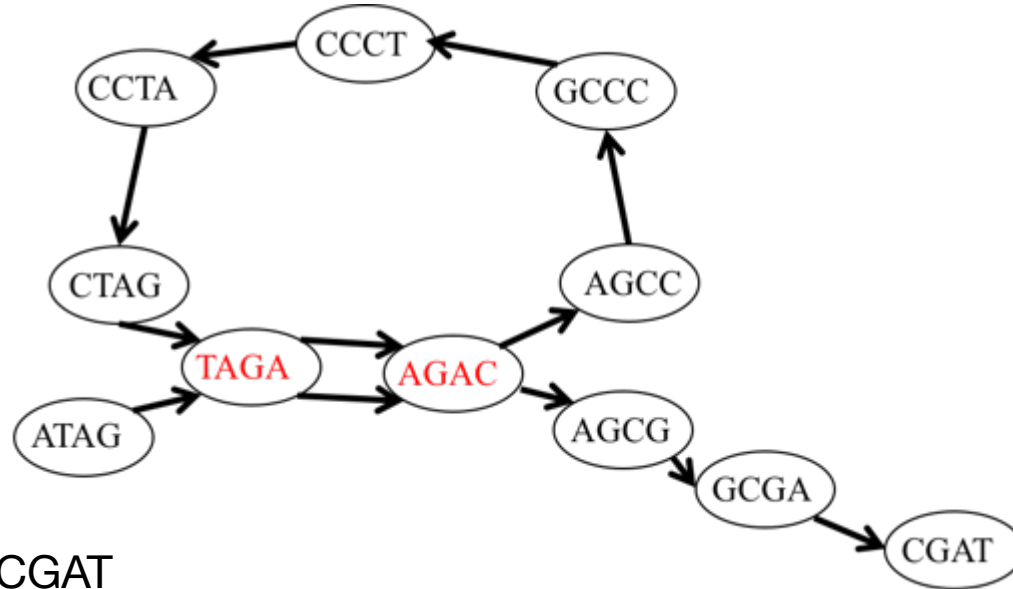
Method for reconstructing genomes from DNA/RNA fragments, with no prior knowledge of the original sequence or the order of those fragments.

## De Bruijn graph

- Extract all substrings of length  $k$  from input reads (i.e. **k-mers**)
- Model relationship in a de Bruijn graph
  - nodes: k-mers
  - edges: adjacent k-mers overlapping by  $k-1$  letters
- Visit each node exactly once through the graph (i.e. identify Eulerian path)



# De novo assembly



ATAGACCCTAGACGAT

ATAGA  
TAGAC

k-1 overlap



# Co-assembly

Merging (appending) two or more samples to be assembled together

## PROS

- More data, better/longer assemblies
- Access to lower abundant organisms

## CONS

- Higher computational overhead
- Risk of shattering the assembly graph by strain variations
- Risk of increased contamination

## When to co-assemble?

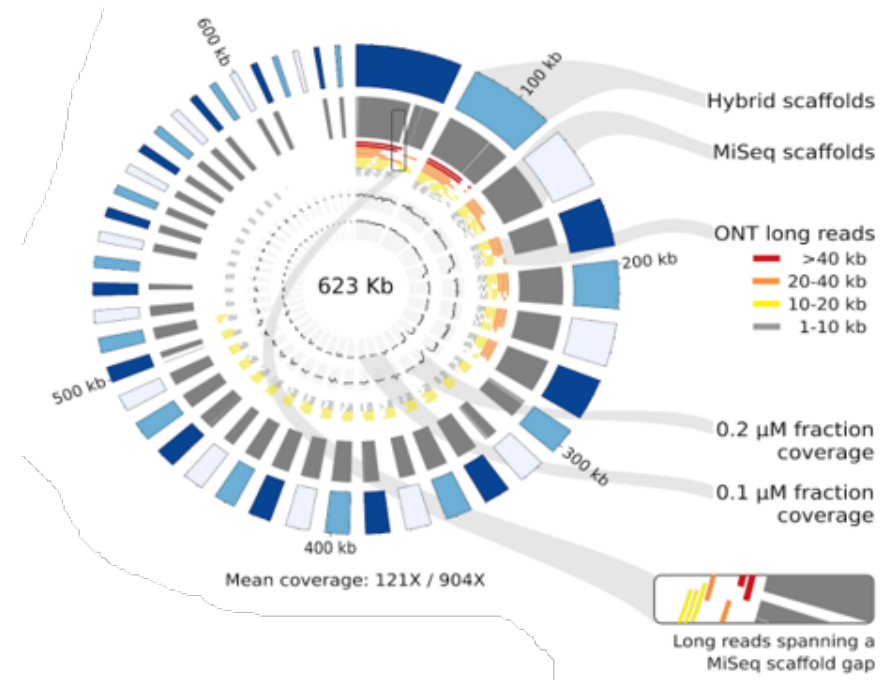
- Same sample
- Same sampling event
- Longitudinal sampling of the same individual
- Related samples

# Combining sequencing technologies

Short-reads + Long-reads = better assembly

2 main strategies

- Assemble short reads, extend contigs and resolve repetitive regions with long reads (**hybrid assembly**)
- Assemble long reads, polish them with short reads



Overholt, W. A. et al., Environmental Microbiology (2020)

# Binning

## Supervised approach

- Relies on known reference genomes
- Uses homology or sequence composition similarity for binning

Multiple bidders exist:

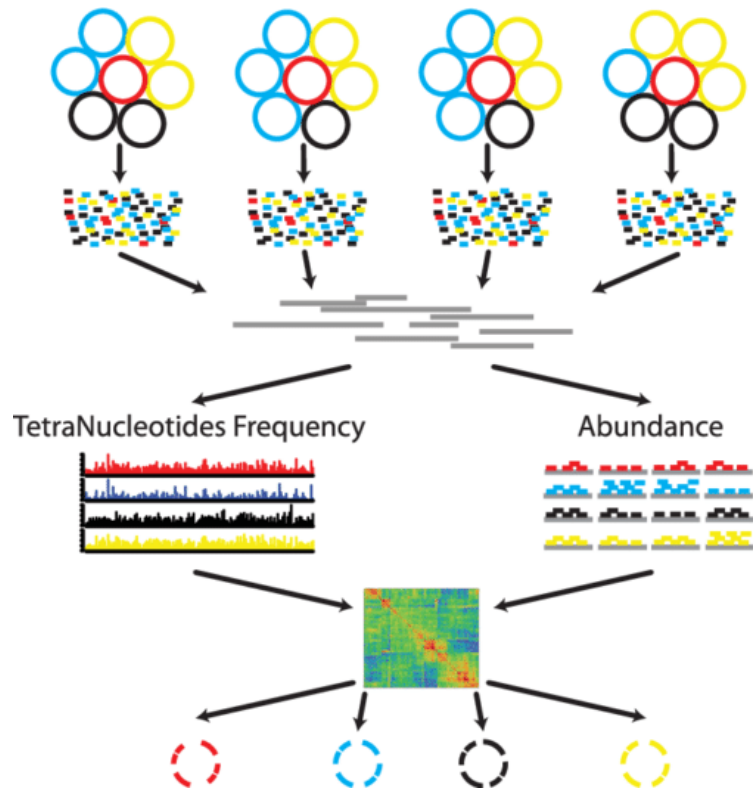
- MetaBAT
- MaxBin2
- CONCOCT
- Semibin
- ...

## Unsupervised approach

- Does not need a reference genome
- Relies on sequence composition similarity and/or species abundance for binning



# MetaBAT

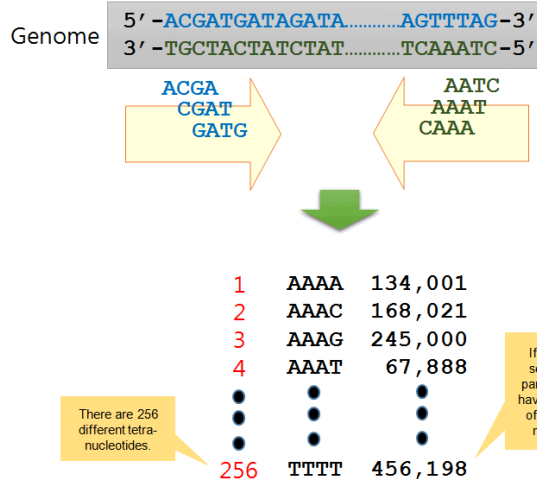


## Preprocessing

- 1 Samples from multiple sites or times
- 2 Metagenome libraries
- 3 Initial de-novo assembly using the combined library

## MetaBAT

- 4 Calculate TNF for each contig
- 5 Calculate Abundance per library for each contig
- 6 Calculate the pairwise distance matrix using pre-trained probabilistic models
- 7 Forming genome bins iteratively

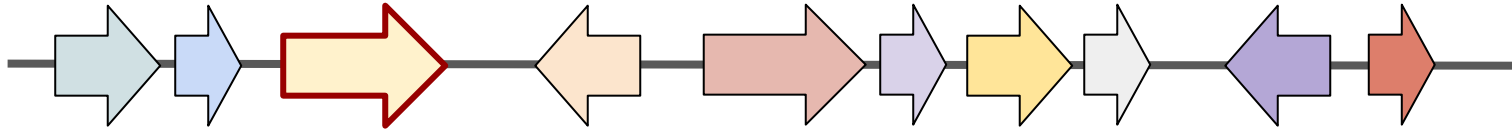




# Quality assessment: CheckM

Uses a set of lineage-specific single-copy marker genes (SCMG) - genes that are present in every genome within a lineage and are single copy.

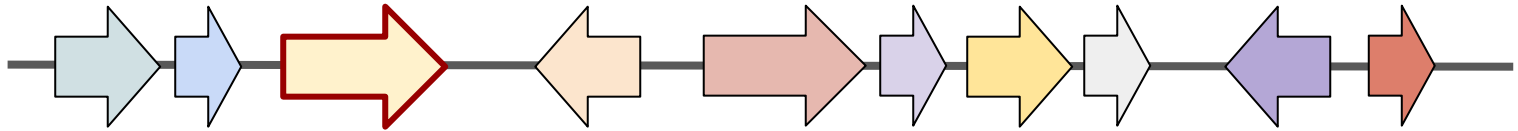
## Reference SCMG set



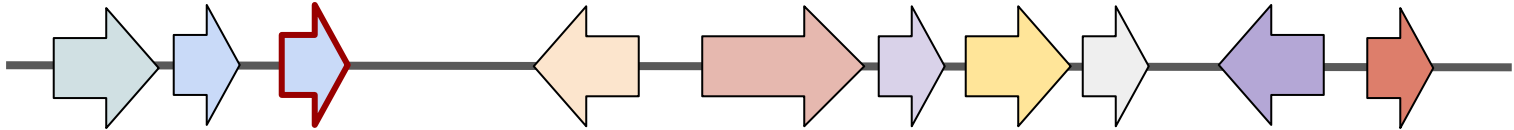
# Quality assessment: CheckM

Uses a set of lineage-specific single-copy marker genes (SCMG) - genes that are present in every genome within a lineage and are single copy.

## Reference SCMG set



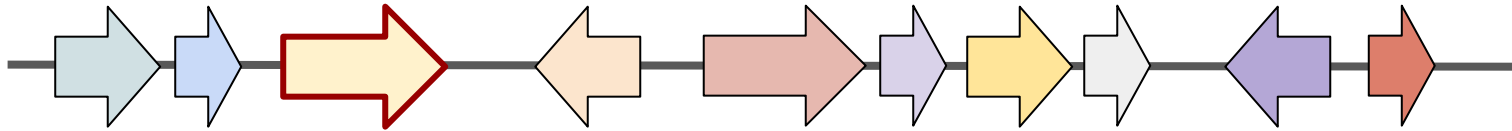
## New genome assembly to evaluate



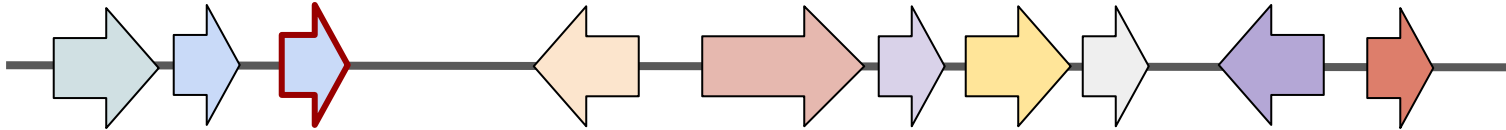
# Quality assessment: CheckM

Uses a set of lineage-specific single-copy marker genes (SCMG) - genes that are present in every genome within a lineage and are single copy.

## Reference SCMG set



## New genome assembly to evaluate



**Completeness:** 90% (9 out of 10 genes are present)

**Contamination:** 10% (1 gene occurs twice)

# Quality assessment: CheckM

**Strain heterogeneity:** indicates the source of contamination (other strains of the same species vs. more distant taxa)

CheckM output:

**Completeness:** 85%

**Contamination:** 7%

**Strain heterogeneity:** 100%



Contamination is likely to come from other strains of the same species

# Quality assessment: CheckM

**Strain heterogeneity:** indicates the source of contamination (other strains of the same species vs. more distant taxa)

CheckM output:

**Completeness:** 85%

**Contamination:** 7%

**Strain heterogeneity:** 0%

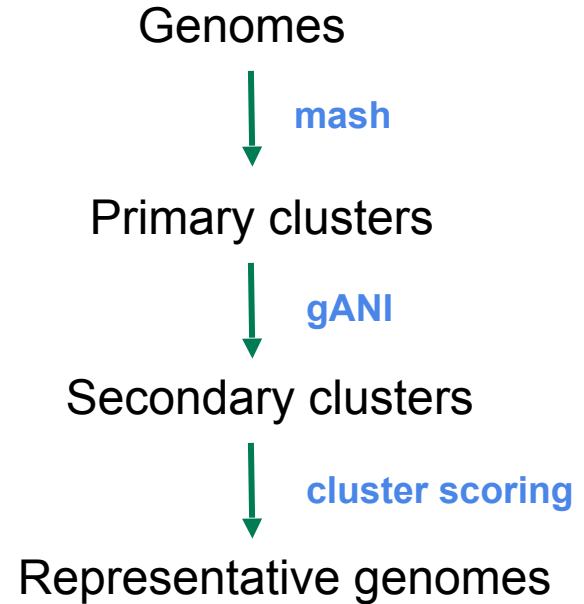
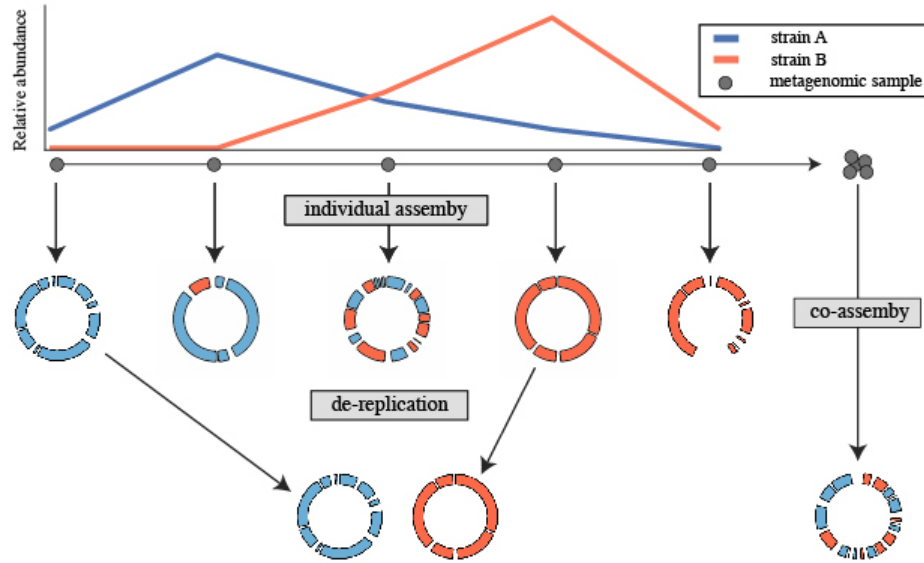


Contamination is likely to come from distant species

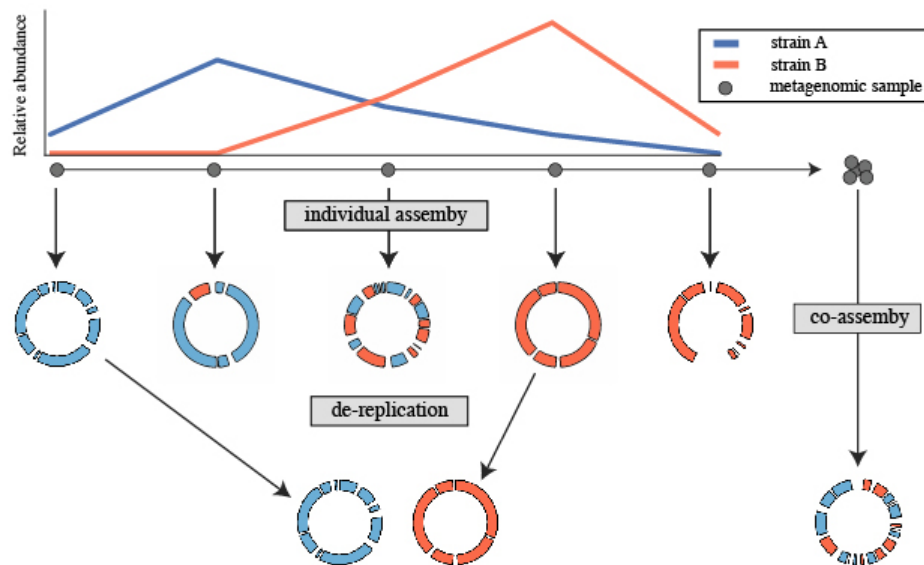
Tools to remove contamination:

- GUNC (<https://grp-bork.embl-community.io/gunc/>)
- MAGpurify (<https://github.com/snayfach/MAGpurify>)

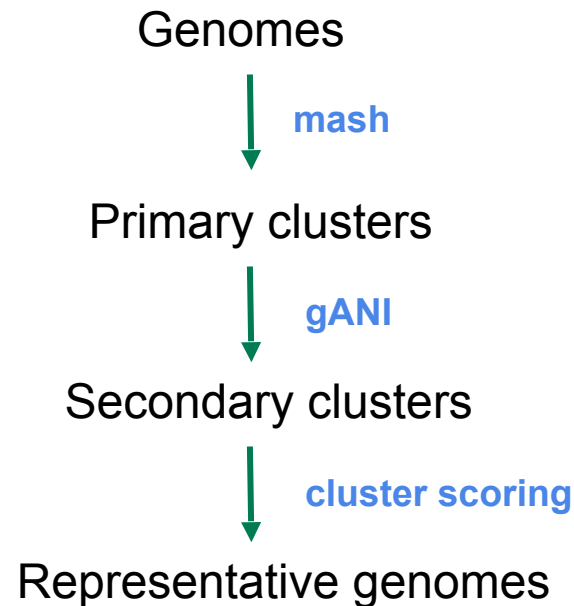
# Removing redundant genomes: dRep



# Removing redundant genomes: dRep



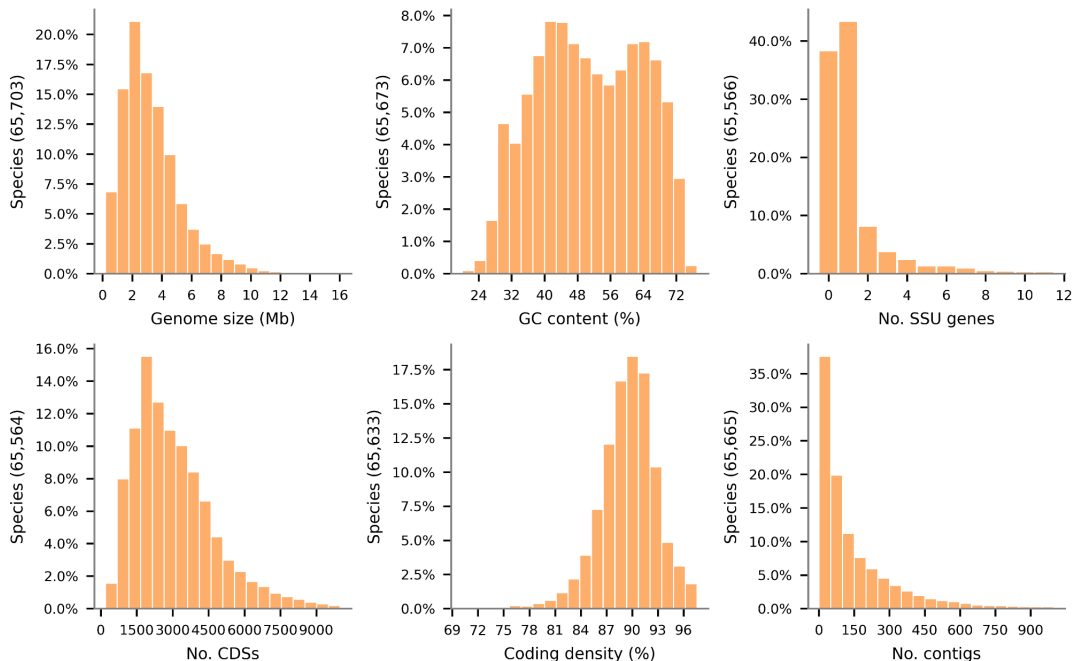
$\text{score} = A * \text{Completeness} - B * \text{Contamination} +$   
 $C * (\text{Contamination} * (\text{strain heterogeneity} / 100)) +$   
 $D * \log(N50) + E * \log(\text{size})$   
 Olm M. R. et al., ISME J (2017)



Species-level ANI distance: 95%  
 Strain-level ANI distance: 99%

# The Genome Taxonomy Database (GTDB)

- Different from NCBI taxonomy, which relies more on isolate genomes
- GTDB balances taxonomic groups based on number of organisms within taxonomic levels
- ~317,000 genomes (Bacteria and Archaea, MAGs and isolates) organised into ~65,000 species clusters
- GTDB-tk allows placement of your genomes within this framework





# Acknowledgements

Rob Finn

Lorna Richardson

Varsha Kale

Sandy Rogers

Tatiana Gurbich

Juan Caballero

MGnify team

# MGnify

