INRAe

**EBAii Assemblage & Annotation 2024**

Construction and analysis of a prokaryotic genomic dataset

V. Loux, H. Chiapello, G. Gautreau

# Construction and analysis of prokaryotic genomic dataset

Outline

> Constructing a genome dataset from public ressources

> Analyzing the genome dataset : intrinsec metrics & distances

> Comparing and dereplicating the dataset

Many slides from the "***Bioinformatique par la pratique***" migale training cycle "Comparison of microbial genomes" module
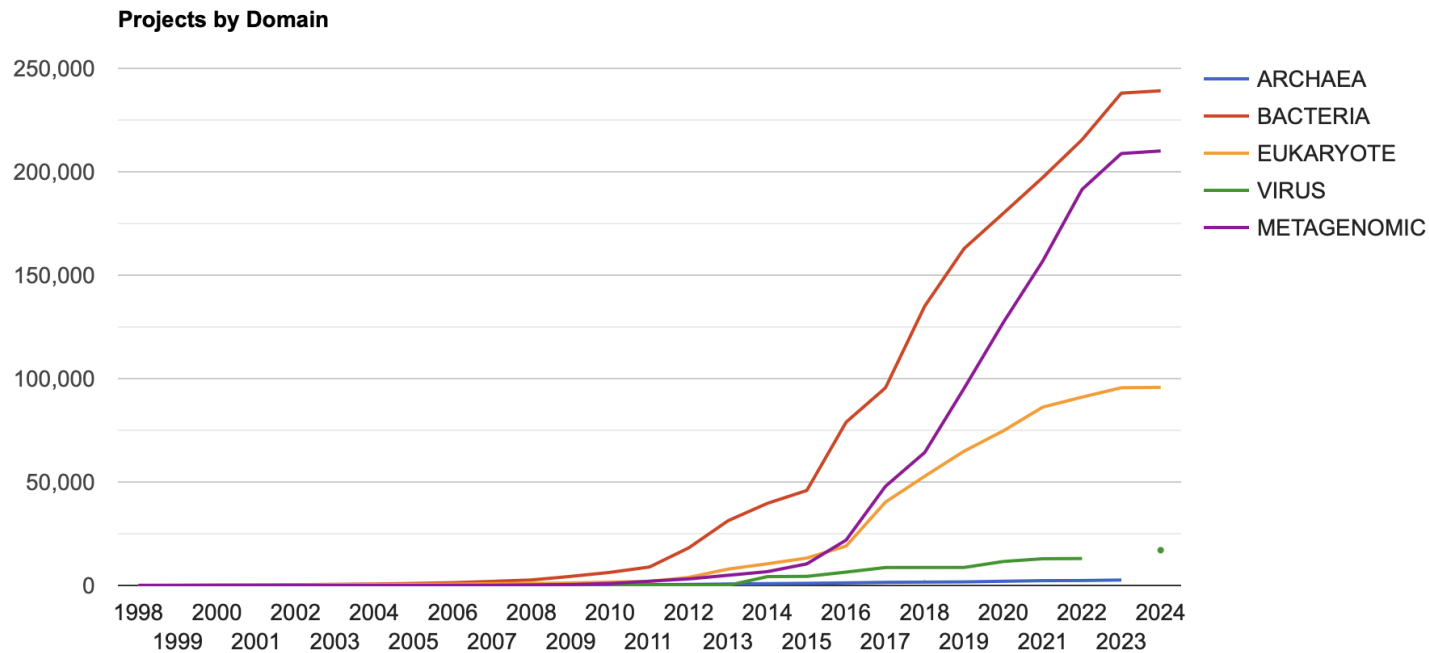
https://migale.inrae.fr/trainings

https://documents.migale.inrae.fr/#training-materials

Hélène Chiapello          Valentin Loux

**INRAe**

# Sequenced genomes by kingdom

**Project Totals in GOLD (by year and Domain Group)**

**Projects by Domain**



Legend:
- ARCHAEA
- BACTERIA
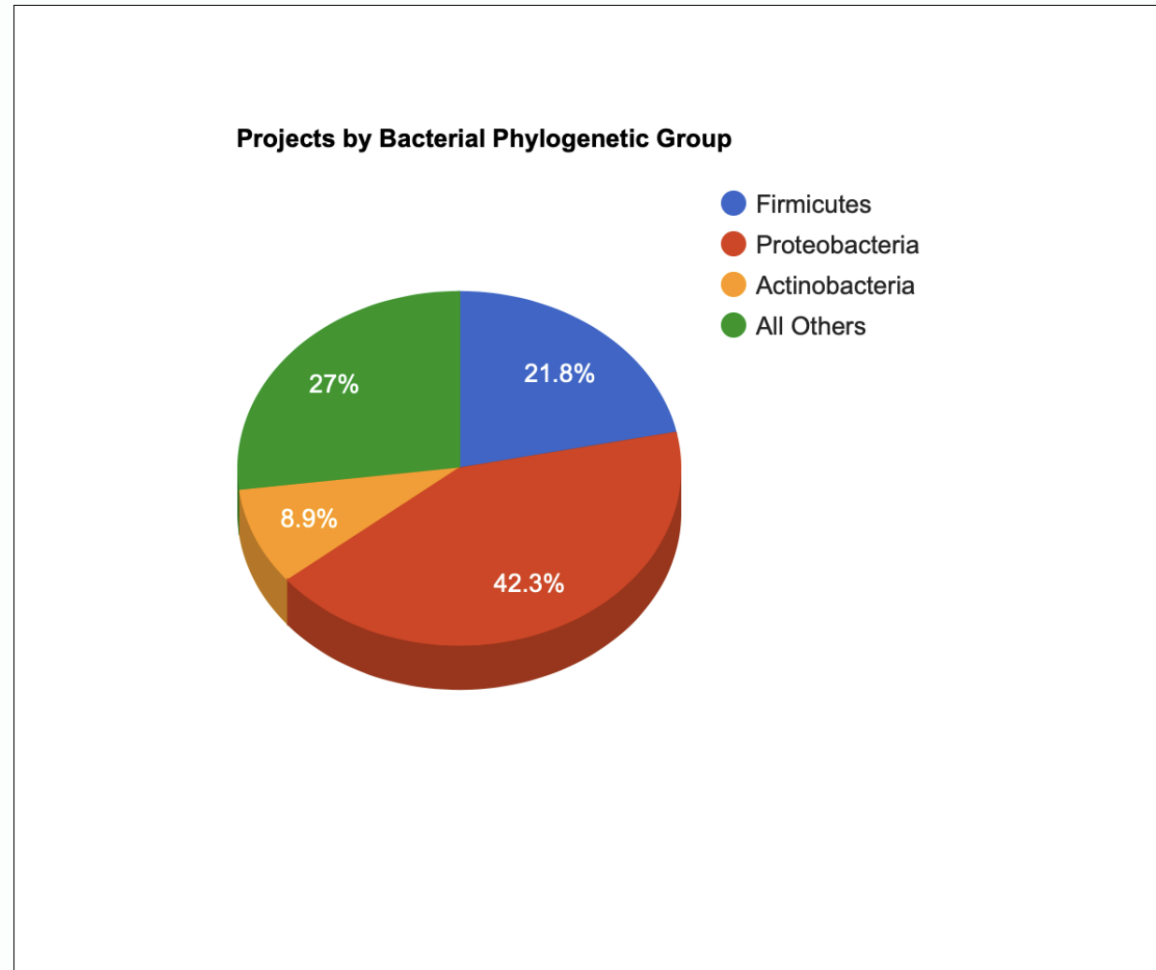- EUKARYOTE
- VIRUS
- METAGENOMIC

**June 24**

- Archea : 2 629
- Eukaryote : 95 593
- Virus : 13 009
- Metagenomics : 210 060
- Bacteria : 239 146

⚠ genomes in GOLD ≠ assembly in Genbank ( > 2 M bacteria)

Source : Gold ( https://gold.jgi.doe.gov/statistics )

INRAe

migale

# Sequenced genomes by kingdom



Phylogenetic distribution of Bacterial Genome Projects

Projects by Bacterial Phylogenetic Group

- Firmicutes
- Proteobacteria
- Actinobacteria
- All Others

21.8% — Firmicutes
42.3% — Proteobacteria
8.9% — Actinobacteria
27% — All Others

Source : Gold ( https://gold.jgi.doe.gov/statistics )

INRAE

EBAii Assemblage & Annotation
27/09/22/ MaIAGE-Migale/ H. Chiapello, V. Loux

migale

# Why using public data and do comparative genomics ?

- Answer to (not so simple) questions like :
  - What is the genomic diversity into a species / genus ?
  - Is the genome structure conserved into a species / genus ?
  - How does the gene repertory evolves into a species / genus ?
  - Is the gene repertory corelated to a given habitat ?

- Does this diversity could explain a given phenotype :
  - metabolism
  - probiotics (anti-inflamatory)
  - pathogenicity

- …

# Dataset building

**Genomes of interest** could be

- already **published** and **available** at public databanks (ENA, NCBI, …)
- private, not yet published.

At least, **we need** :

- [As much as possible] **complete genome assemblies** (contigs / scaffolds in Fasta format)
- Syntactic and functional annotation :
  - Genbank or GFF format
  - For private genomes, you could/should use Bakta [remember what GG told you yesterday afternoon ?]

It's always better, if not mandatory, if [syntaxic] **annotation is homogeneous**

# Quick reminder on formats !

INRAᏋ

# FASTA format

The **FASTA format** is used to represent **sequence information**.

The format is very simple:

A > symbol on the FASTA header line indicates a fasta record start.

A string of letters called the sequence id may follow the > symbol.

The header line may contain an arbitrary amount of text (including spaces) on the same line.

Subsequent lines contain the sequence.

```
>foo
ATGCC
>bar other optional text could go here
CCGTA
>bidou
ACTGCAGT
TTCGN
>repeatmasker
ATGTGTcgggggggATTTT
>prot2; my_favourite_prot
MTSRRSVKSGPREVPRDEYEDLYYTP
SSGMASP
```

# Genbank Format

The Genbank format is used to represent sequence **and** annotation information together.

The start of the annotation section is marked by a line beginning with the word "LOCUS".

Features (CDS, genes) are annotated with their position , strand, and qualifiers that contain the  annotation.

The **start of sequence** section is marked by a line beginning with the word "ORIGIN" and the end of the section is marked by a line with only "//".

Those three bank agree on the list of **feature / qualifier** that one can use to annotate sequence. (Cf https://www.ebi.ac.uk/ena/WebFeat/ )

NCBI, ENA (European Nucleotide Archive) et DDBJ (Japan) entries are synchronized each day.

# Genbank entry example

```
LOCUS       SCU49845      5028 bp     DNA                      PLN        21-JUN-1999
DEFINITION  Saccharomyces cerevisiae partial genes.
ACCESSION   U49845
VERSION     U49845.1  GI:1293613
KEYWORDS    .
SOURCE      Saccharomyces cerevisiae (baker's yeast)
  ORGANISM  Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomyces.
REFERENCE   1  (bases 1 to 5028)
  AUTHORS   Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.
  TITLE     Cloning and sequence of REV7, a gene whose function is required for
            DNA damage-induced mutagenesis in Saccharomyces cerevisiae
  JOURNAL   Yeast 10 (11), 1503-1509 (1994)
  PUBMED    7871890
FEATURES             Location/Qualifiers
     source          1..5028
                     /organism="Saccharomyces cerevisiae"
                     /db_xref="taxon:4932"
                     /chromosome="IX"
                     /map="9"
     CDS             <1..206
                     /codon_start=3
                     /product="TCP1-beta"
                     /protein_id="AAA98665.1"
                     /db_xref="GI:1293614"
                     /translation="SSIYNGISTSGLDLNNGTIADMRQLGIVESYKLKRAVVSSASEA
                     AFVLLRVDNTTRARPRTANROHM"
```

# GFF Format

The **General Feature Format** contains **annotation** and (optionally) **sequence**. It consists of one line per feature, each containing 9 columns of data, plus optional track definition line.

```
##gff-version 3
##sequence-region NZ_LHTK01000001 1 688985
# organism Salmonella enterica subsp. arizonae serovar 62:z36:- str. 5335/86
# date 17-JAN-2020
NZ_LHTK01000001    GenBank    contig     1     688985    .    +    1    ID=NZ_LHTK01000001;Dbxref=BioProject:PRJNA224116,taxon
NZ_LHTK01000001    GenBank    pseudogene    1    1014    .    -    1    ID=LFZ49_RS22320.pseudogene;Alias=LFZ49_RS22320;Name
NZ_LHTK01000001    GenBank    gene    1011    1634    .    -    1    ID=LFZ49_RS00010;Name=LFZ49_RS00010;old_locus_tag=LFZ49
NZ_LHTK01000001    GenBank    mRNA    1011    1634    .    -    1    ID=LFZ49_RS00010.t01;Parent=LFZ49_RS00010
```

# Practical : public genomes

How to gather a list of public genomes of interest ?

List **all** available public genomes of interest with :

- Associated metadata

- Metrics quality (size, N50, completnesse…)

- Filter according to above criteria

- Download genomes in **various formats**

Work from the complete list of prokaryotic public genomes available at NCBI [https://www.ncbi.nlm.nih.gov/datasets/genome/]

# A solution : NSBI datasets



NCBI Datasets components

NCBI Datasets is a new resource that lets you easily gather data from across NCBI databases. You have the choice of getting the data through three interfaces:

- NCBI Datasets website Command-line tools

- API (Application programming Interface)

NCBI Datasets delivers data and metadata as a **cohesive data package** contained in a zip archive. i.e., for an **assembly : sequences, annotation (CDS, transcripts, genome…) and metadata.**

INRAe

# Source for genome assemblies

- **A GenBank** (**GCA**) genome assembly contains assembled genome sequences submitted by investigators to GenBank or another member of the International Nucleotide Sequence Database Collaboration (INSDC)

- A **RefSeq** (**GCF**) genome assembly represents an NCBI-derived copy of a submitted GenBank (GCA) assembly. In the majority of cases, the annotation is generated by the NCBI prokaryotic or eukaryotic genome annotation pipelines

| | **GCA_** | **GCF_** |
|---|---|---|
| Also known as | GenBank assembly | RefSeq assembly |
| Submitter-owned assembly archive | ✔ | ✘ |
| NCBI-maintained assembly copy | ✘ | ✔ |
| Always includes annotation | ✘ | ✔ |
| NCBI may add sequences (e.g. mitochondrial genomes) | ✘ | ✔ |
| NCBI may remove sequences (e.g. contamination) | ✔ * | ✔ |

\* following submitter request or agreement

**NCBI Datasets website genome sources**

Source : Dataset documentation

# NCBI Datasets : Datasets Genome Table



- Find **all current genomes**, including metagenomes
  View **multiple taxa** such as birds and bees, or polyphyletic groups like fish

- Easily find genomes with **NCBI RefSeq** annotations
  Get more accurate genome counts, since **each row now represents a single genome with GenBank and RefSeq accessions** for that genome in the same row

- **Customize your downloads** to include either GenBank or RefSeq files, or both
  Download **tables** or **data packages**

Figure Source

# NCBI Datasets : Command Line



**genome** options :

- summary according to accession or taxid
- filter according to quality criteria & metadata
- download packages (or rehydrate) in various formats

**INRAℓ**

# NCBI Datasets : Aplication Programmatic Interface



[Jupyter Notebook](#)

# NCBI Datasets : Galaxy Integration



A wrapper of the command line tool

Parameters to define packages files

# NCBI Datasets : Galaxy Integration

**A few caveats of the wrapper** :

- Some (not so easy) errors when select / filter fails
- Often broken 😩
- Impossible to just download a list of **genomes as a file** and "**rehydrate**" it after

**What we recomend to do** [for now] :

- use the NCBI dataset genome page to **browse / filter a list of genomes** of interest, and choose one of these solution
  - **A few genomes** :
    - Download the genomes as package un re-upload into Galaxy 😩
  - **A lot of genomes :**
    - Extract the URL of interest of the files from the TSV, feed it into Galaxy

**INRAQ**

# The training datasets

We will work on 3 datasets of public Salmonella genomes



- *217 949 Salmonella enterica enterica* public assemblies at NCBI!!

INRAe

# Training

We will construct 3 datasets of public *Salmonella* genomes

- **Dataset 1**: list all Salmonella enterica subsp. enterica assemblies using their taxon id and assembly level (Chromosome)

- **Dataset 2**: list all the **Salmonella bongori** assemblies to choose and download the best outgroup of a salmonella enterica dataset

- **Dataset 3**: download 16 Salmonella enterica public assemblies (2 sub-species, 4 serotypes) from their accession numbers

# Data set 1 : Taxonomy browser



List all *Salmonella enterica subsp. enterica* assemblies using their *taxon id* and *assembly level (Chromosome)*

INRA@

# Data set 1 : genome table



- Notice :
  - Filter parameters
  - Select columns button
  - Download button (table or package)

INRAØ

# •Data set 2: download a genome



**List all** the *Salmonella bongori assemblies* to **choose the best outgroup** of a *Salmonella enterica* dataset

**Download** the genome in Genbank, nuclotide and fasta and GFF format

# Data set 3 : 16 S. enterica public genomes (part 1)

| Assembly_accession | Subspecies | Serotype | Strain | assembly_level |
|---|---|---|---|---|
| GCF_001951465.1 | arizonae | 18:z4,z23 | CVM N27 | Scaffold |
| GCF_001448925.1 | arizonae | 62:z36 | 5335/86 | Contig |
| GCF_000756465.1 | arizonae | 62:z36 | RKS2983 | Complete Genome |
| GCF_000018625.1 | arizonae | 62:z4 | z23 | Complete Genome |
| GCF_000983595.1 | enterica | ParatyphiA | na | Scaffold |
| GCF_000026565.1 | enterica | ParatyphiA | AKU_12601 | Complete Genome |
| GCF_000011885.1 | enterica | ParatyphiA | ATCC 9150 | Complete Genome |
| GCF_000484015.1 | enterica | ParatyphiB | SARA61 | Contig |

# Data set 3 : 16 S. enterica public genomes (part 1)

| Assembly_accession | Subspecies | Serotype | Strain | assembly_level |
|---|---|---|---|---|
| GCF_001951465.1 | arizonae | 18:z4,z23 | CVM N27 | Scaffold |
| GCF_900002585.1 | enterica | Typhi | na | Scaffold |
| GCF_000256015.1 | enterica | Typhi | BL196 | Contig |
| GCF_000195995.1 | enterica | Typhi | CT18 | Complete Genome |
| GCF_000007545.1 | enterica | Typhi | Ty2 | Complete Genome |
| GCF_001120665.1 | enterica | Typhimurium | DT104 | Scaffold |
| GCF_000006945.2 | enterica | Typhimurium | LT2 | Complete Genome |
| GCF_000210855.2 | enterica | Typhimurium | SL1344 | Complete Genome |
| GCF_000312745.2 | enterica | Typhimurium | STm6 | Contig |

# Data set 3 : from a tabular file

**Download** 16 Salmonella enterica public assemblies (2 sub-species, 4 serotypes) from their *accession numbers*.

**Input : [Filtered]** List of assembly accession in a tabular file downloaded from **Dataset genome Table**

- Import ncbi_dataset_salmonella_genome_table.tsv from Shared Data / Data Library / EBAII A&A 2022 / Prokaryotic Annotation / NCBI Dataset

- **Filter lines concerning Refseq assemblies** ( starts with "GCF_") using **Select lines that match an expression** tool

- Select the first column of the file ( Assembly Accession) using Cut columns from a table

- Feed **NCBI Datasets Genomes** download genome sequence, annotation and metadata with the list of accession

- BROKEN 😣 Retrieve all file format of interest **including** genbank annotated files BROKEN 😣

  - [also in Shared Data / Data Library/ Libraries /EBAII A&A 2022 Prokaryotic annotation/ Salmonella dataset]

# Construction and analysis of prokaryotic genomic dataset

Outline

> Constructing a genome dataset

> Analyzing the genome dataset

> Comparing and dereplicating the dataset

INRAe

# Analyzing a genome dataset

Why?

- Frequent problems in genome analysis and comparison
  - Heterogenous quality of sequencing and assembly
  - Presence of huge number or public genomes OR absence of any close genomes of the same species in public databases
  - Difficulties regarding microbial taxonomy (classification) and nomenclature (naming of genus, species and strain naming) for many non-model organisms

# Introduction

- What is a species?



Ernst Mayr (1942) :
*"Species are groups of actually or potentially interbreeding natural populations, which are reproductively isolated from other such groups"*
⇨ **Not relevant for bacteria**

# What is a bacterial species?

No consensual definition for procaryotes

▶ **No universal criteria**
▶ **Several approaches used to classify bacterial**
  - Phenotypes and morphological criteria
  - DNA-DNA hybridization
▶ **Universal markers**
  - 16S rRNA
  - MLST (Multi Locus Sequence Typing)
▶ **Genomic-based taxonomy are now becoming a gold-standard**

| % ADN-ADN hybridation | >70% |
| --- | --- |

| % rRNA 16S identity | >98,7% |
| --- | --- |

# Example: the Genome-based taxonomy for prokaryotic genomes

- Objective: a standardized microbial taxonomy based on genome phylogeny

- Taxonomy inferred from concatenated single copy marker proteins

Parks et al. 2018, 2021

https://gtdb.ecogenomic.org/

# Evaluating genome diversity in a dataset

- Why ?
  - Identify outlier genomes
  - Identify groups of (very) similar genomes and de-replicate datasets
  - Estimate genome similarity in a dataset and design an adapted comparative strategy

How ?
  - Alignment based approaches (ANI)
  - k-mer based approaches (MASH)

# Average Nucleotide Identity (ANI)

- Meet the need for a robust measure of genomic relatedness and a systematic and scalable species assignation technique

- Mean identity percent of aligned regions of a pair of genomes

- Rely on pairwise alignments from
  - aligned core genes
  - genomic alignments

- Can easily be used to build phylogenetics tree using distance methods

- Is implemented in several bioinformatics tools: ANIn (nucmer based, Richter 2009) gANI (coding regions, Varghese 2015), fastANI (computing efficiency, Jain 2018), …



Genetic diversity within five important bacterial groups. Konstantidinis et al. 2006. The bacterial species definition in the genomic era
DOI: 10.1098/rstb.2006.1920

# Average Nucleotide Identity (ANI)

- ANI strongly correlates (R = 0.79 for logarithmic correlation) with the 16S rRNA gene sequence identity and can resolve areas where the 16S rRNA gene is inadequate (intra-species level)

- The average rate of synonymous substitutions shows a tight correspondence to ANI, suggesting that ANI may also be a useful descriptor of the evolutionary distance

- ANI shows a strong linear correlation to DNA–DNA reassociation values, and the 70% DNA–DNA reassociation standard corresponds to ≈93–94% ANI i.e. strains that show >94% ANI should belong to the same species



ANI vs 16 rRNA %id          ANI vs average syn. Mutation rate          ANI vs DNA-DNA reassociation value

Konstantidinis et al. 2005. Genomic insights that advance the species definition for prokaryotes
https://doi.org/10.1073/pnas.0409727102

INRAɛ

# MASH: fast (meta)genome distance estimation using MinHash

- Mash allows to compute a pairwise mutation distance without alignment using k-mer counts
- Mash provides two basic functions for sequence comparisons:
  - **sketch**: converts a sequence or collection of sequences into a MinHash sketch
  - **dist**: compares two sketches and returns an estimate of the Jaccard index (i.e. the fraction of shared k- mers), a P value, and the Mash distance



$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

Ondov, B.D., Treangen, T.J., Melsted, P. et al. Mash: fast genome and metagenome distance estimation using MinHash. Genome Biol 17, 132 (2016). https://doi.org/10.1186/s13059-016-0997-x

# MASH distances correlate well with ANI

- Dataset: 500 complete E. coli genomes
  - Each plot column shows a different sketch size
  - Each plot row a different k-mer size k.
  - Gray lines: model relationship D = 1– ANI

- Increasing the sketch size improves the accuracy of the MASH distance, especially for more divergent sequences.

- Limit on how well the MASH distance can approximate ANI, especially for more divergent genomes (e.g. ANI considers only the core genome)

Ondov, B.D., Treangen, T.J., Melsted, P. et al. Mash: fast genome and metagenome distance estimation using MinHash. Genome Biol 17, 132 (2016). https://doi.org/10.1186/s13059-016-0997-x

# Back to procaryote taxonomy

| % ADN-ADN hybridation | >70% |
|---|---|



| % rRNA 16S identity | >98,7% |
|---|---|



| % genome identity (ANI) | >94% |
|---|---|



| K-mer distance | <0,06 |
|---|---|

**MASH**

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

# Quality Control & filter

- Already presented: evaluate Quality using the 3Cs

1. **Contiguity**. Produce the longest possible contigs.

2. **Correctness**. Assemble contigs with few/no errors.

3. **Completeness**. Cover the entire original sequence and minimize missing regions

- An additional key point for microbes: evaluate **Contamination**
  - From genomic fragments of divergent taxa
  - From genomic fragments of multiple strains (i.e. strain **heterogeneity**)
  - Be sure that the "announced" taxa is good

# CheckM

- a set of tools for assessing the quality of genomes recovered from isolates, single cells, or metagenomes
- provides robust estimates of genome **completeness** and **contamination**
  - use collocated sets of genes that are ubiquitous and single-copy within a phylogenetic lineage
  - propose a fixed vocabulary for defining genome quality based on estimates of completeness and contamination
- Evaluate by simulations the accuracy of quality estimates



**INRAΘ**

**CheckM consists of a workflow for precomputing lineage-specific marker genes for each branch within a reference genome tree (top box) and an online workflow for inferring the quality of putative genomes (bottom box).**



**Donovan H. Parks et al. Genome Res. 2015;25:1043-1055** © 2015 Parks et al.; Published by Cold Spring Harbor Laboratory Press

INRAe

migale

# CheckM relies on several other tools and data

- *prodigal* to predict genes

- A reference genome tree based on 43 phylogenetically informative marker genes and 5656 trusted reference genomes
  - Marker genes are identified in assemblies using HMMER
  - The resulting genes are used to placed the genome into the tree using *pplacer*

- Lineage-specific marker sets determined for all nodes within the reference genome tree by identifying single-copy genes present in ≥97% of all descendant genomes.

INRAe

# CheckM report

Provides classic quality metrics and plots, including:

- Results of binning

    >Marker lineage, #genomes, #markers, #marker sets

- CheckM metrics

    > Completeness, Contamination, Strain heterogeneity

- Classical Quality metrics

    > #ambiguous bases, #scaffolds, #contigs, N50 (scaffolds), N50 (contigs), Mean scaffold length (bp),Mean contig length (bp), Longest scaffold (bp),Longest contig (bp), GC, GC std (scaffolds > 1kbp)

# CheckM report – binning part

Marker lineage:  indicates the taxonomic rank of the lineage-specific marker set used to estimated genome completeness, contamination, and strain heterogeneity.

#genomes:  number of reference genomes used to infer the lineage-specific marker set

#markers:  number of marker genes within the inferred lineage-specific marker set

#marker sets: number of co-located marker sets within the inferred lineage-specific marker set

0-5+: number of times each marker gene is identified

# CheckM report

- **Completeness:** estimated completeness of genome as determined from the presence/absence of marker genes and the expected colocalization of these genes

- **Contamination:** estimated contamination of genome as determined by the presence of multi-copy marker

- **Strain heterogeneity:** % determined from the number of multi-copy marker pairs which exceed a specified **amino acid identity threshold** (default = 90%).
  - High strain heterogeneity suggests the majority of reported contamination is from one or more closely related organisms (i.e. potentially the same species),
  - Low strain heterogeneity suggests the majority of contamination is from more phylogenetically diverse sources

INRAE

# CheckM: proposed genome quality classification scheme

- **Finished genomes**: genomes assembled into a single contiguous sequence containing no gaps or ambiguities and where extensive efforts have been made to identify errors

- **Noncontiguous finished:** genomes assembled into multiple sequences as a result of repetitive regions, but otherwise of a finished quality

- **Draft genomes:** all other genomes

**Table 3.** Controlled vocabulary of draft genome quality based on estimated genome completeness and contamination

| Completeness | Classification | Contamination | Classification |
|---|---|---|---|
| ≥90% | Near | ≤5% | Low* |
| ≥70% to 90% | Substantial | 5% to ≤10% | Medium |
| ≥50% to 70% | Moderate | 10% to ≤15% | High |
| <50% | Partial | >15% | Very high |

(*) Genomes estimated to have 0% contamination can be designated as having "no detectable contamination".

**Donovan H. Parks et al. Genome Res. 2015;25:1043-1055**

⚠️ **Evolution of the thresholds with the evolution of sequencing technologies** :

- Those threshold are now more for **MAGs** (see Pasolli et al. 2019, Bowers et al., 2017 )
- For **isolates genome**, for instance , Refseq defines high quality as >98

**INRAe**

# CheckM result interpretation limits

- CheckM is dedicated to eubacterial and archeal genomes
    - Eukaryotic or phage genomes will be reported as highly incomplete
    - The quality of plasmids must also be assessed independently of CheckM

- The novelty of a genome will also influence the accuracy of CheckM estimates
    - Estimates for bacterial and archaeal genomes from deep basal lineages with few reference genomes are generally based on domain-level marker sets
    - Quality estimates may be not reliable for genomes of novel lineages
    - Gene loss or duplication may be an issue

Conclusion : use CheckM as a tool to detect outliers and further investigate!

# Construction and analysis of prokaryotic genomic dataset

Outline

> Constructing a genome dataset

> Analyzing the genome dataset

> **Comparing and dereplicating the dataset**

# Comparing and dereplicating a genome dataset

Why?

> ## To deal with

- The huge number or public genomes for some taxonomical groups including very similar or identical ones
  - Ex : E. coli, S. enterica
- The heterogenous quality of sequencing and assembly of these data

> ## To design a relevant comparative strategy adapted to the dataset

- **Back to genome diversity evaluation**

Two main methods
- Alignment based approaches (ANI)
  - slow (need pairwise comparisons)
  - Robust to genome incompleteness
- k-mer based approaches (MASH)
  - Rapid (hash technics)
  - Not robust to genome incompleteness
  - Only provides an estimate of ANI
    - > Become very approximative for very divergent genomes

- ## **Comparing and dereplicating a genome dataset**

The dRep tool

> dRep is a python program which performs **rapid pairwise genome comparisons** using genomic distances

> it can be used for genome **dereplication**: identification of the 'same' genomes from a large set + determination of the highest quality genome in each replicate set

Very good documentation:
https://drep.readthedocs.io/en/latest/



dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication
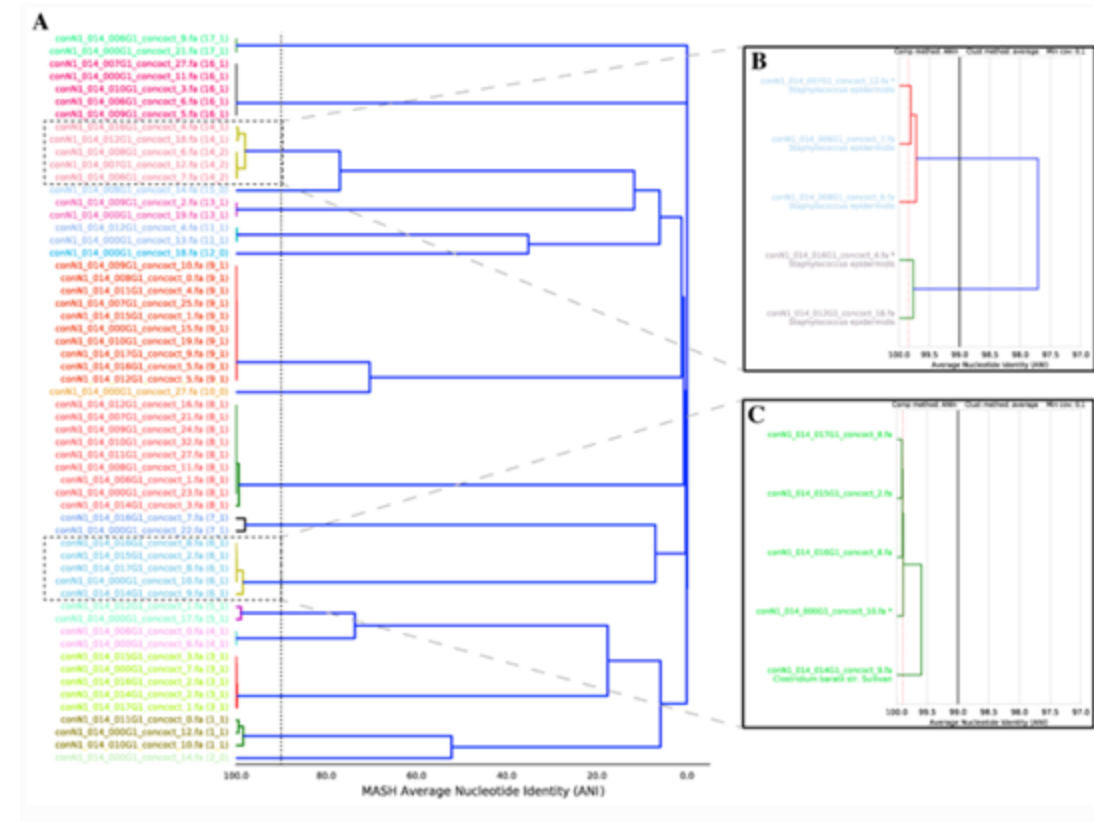
# Comparing and dereplicating a genome dataset

The dRep tool

dREP uses 2 main steps:

1. a first (rapid) clustering of genomes using MASH similarity (90% by default)

2. a second more sensitive step based on ANI on pairs of genomes that have at least a minimum level of "MASH" similarity (99% by default)

Very good documentation:
https://drep.readthedocs.io/en/latest/

INRAE

- **dRep important concepts**

**1. dRep primary clustering use a greedy algorithm,** i.e. an algorithm that take shortcuts to run faster and generally produces "quasi-optimal" solutions. Genomes that are not on the same MASH primary clustering will never be compared with ANI

**2. Importance of genome completeness:** MASH is very sensitive to genome completness. the more incomplete of genomes you allow into your genome list, the more you must decrease the primary cluster threshold.

**3. The secondary ANI threshold (default value: 99%, limit: 99.99%) indicates how similar genomes need to be to be considered the "same".** Depending on the application,you may modify this parameter, i.e.: 95% ANI for species-level de- replication or 98% ANI to generate a set of genomes that are distinct when mapping short reads.

**4. A score is used to pick representative genomes takes into account several parameters such as Completeness, Contamination, strain heterogeneity and centrality** (a measure of how similar a genome is to all other genomes in it's cluster).

INRAⒺ

- **dRep commands and parameters**

**1. dREp compare:** compare and cluster a set of genomes using one or two clustering steps.

**2. dREp dereplicate**: compare, cluster and dereplicate a set of genomes. During dereplication the first step is identifying groups of similar genomes, and the second step is picking a Representative Genome (RG) for each cluster

Parameters of primary and secondary clustering may have to be adjusted depending on the diversity of the dataset and on the objective of the comparison/dereplication

Default values of dRep clustering parameters:

```
-pa P_ANI, --P_ani P_ANI

                ANI threshold to form primary (MASH) clusters

                (default: 0.9)

-sa S_ANI, --S_ani S_ANI

                ANI threshold to form secondary clusters (default:

                0.99)
```

# • **dRep practice**

**use dREP-dreplicate** to explore the Salmonella genome dataset diversity and completenes and dereplicate the dataset

> input : 16 Salmonella genome fasta files

➢ Default parameters

➢ **All outputs**

**explore and interpret results**

**INRAE**

EBAii Assemblage & Annotation

27/09/22/ MaIAGE-Migale/ H. Chiapello, V. Loux

# • **dRep tools and result files**

**dRep rely on several other programs:**

**1. Mash**: to build the primary clusters

**2. Mummer**: to perform the ANI computation on pairwise genome alignements (used by default but **fastANI** or **gANI** may also be used)

**3. checkM** (Parks et al. 2015) to determine contamination and completeness of genomes

**4. Prodigal** (Hyatte et al. 2010): to predict genes (used by checkM and gANI)

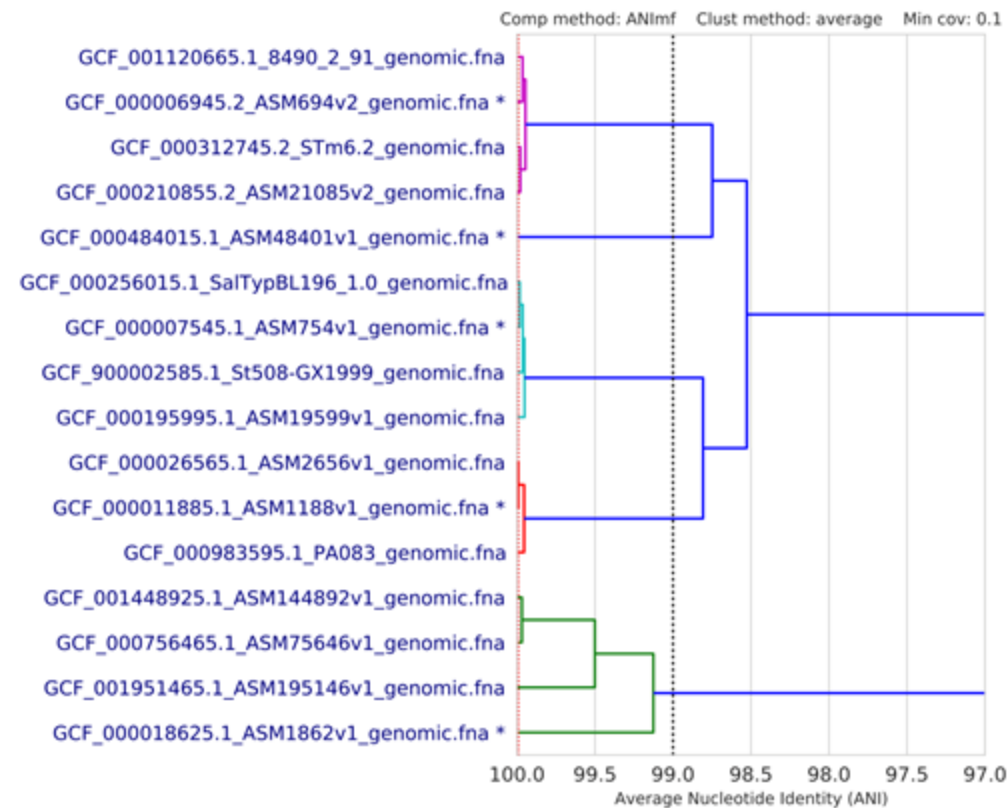**5. cipy** (Jones et al. 2001) to produce a final hierarchical clustering.

```
workDirectory
./data
...../checkM/
...../Clustering_files/
...../gANI_files/
...../MASH_files/
...../ANIn_files/
...../prodigal/
./data_tables
...../Bdb.csv   # Sequence locations and filenames
...../Cdb.csv   # Genomes and cluster designations
...../Chdb.csv  # CheckM results for Bdb
...../Mdb.csv   # Raw results of MASH comparisons
...../Ndb.csv   # Raw results of ANIn comparisons
...../Sdb.csv   # Scoring information
...../Wdb.csv   # Winning genomes
...../Widb.csv  # Winning genomes' checkM information
./dereplicated_genomes
./figures
./log
...../cluster_arguments.json
...../logger.log
...../warnings.txt
```
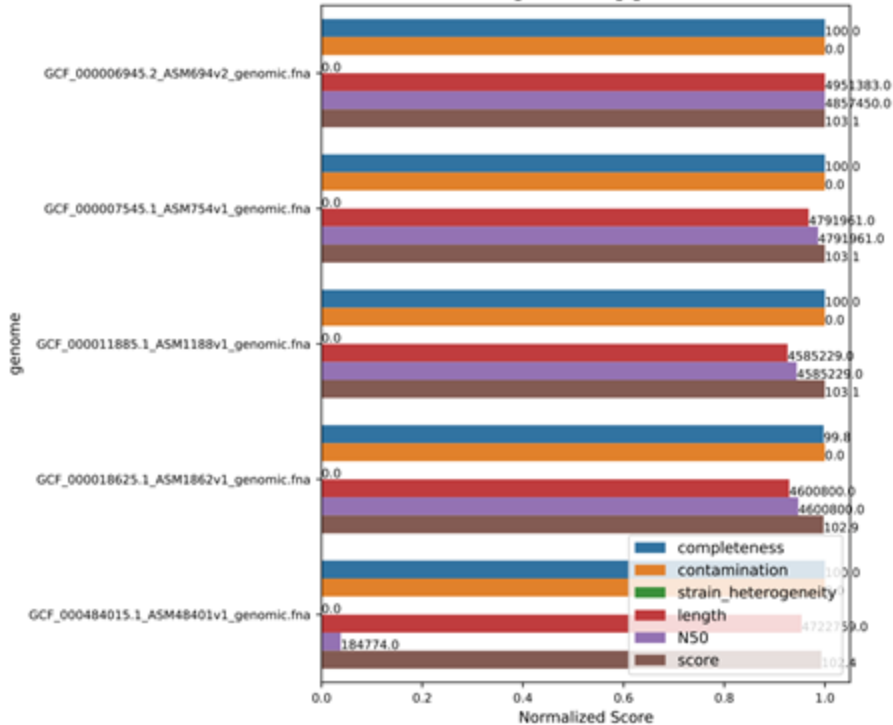
Output files of dRep

# dRep results interpretation



Secondary_clustering_dendrograms.pdf



Winning_genomes.pdf

INRAE

# •Next steps

- Linking Genomes with meta-data :
  - Meta-data extracted from « Genome » csv file
  - **Omnicrobe** : a database of habitats, phenotypes and uses of microorganisms. (Litterature , genbank, DSMZ, CIRM…)

  https://omnicrobe.migale.inrae.fr

- Pangenome analysis [Cf next talk]
  - **Roary**
    - ⚠ File format. GFF from Bakta OR GFF convetered from Gbk with Genbank to GFF3 converter (Cf https://sanger-pathogens.github.io/Roary/ )
  - **PPanGGOLiN** : Depicting microbial species diversity via a Partitioned PanGenome Graph Of Linked Neighbors