# INRAE

# EBAii Assemblage & Annotation

Part 2: construction and analysis of a prokaryotic genomic dataset

H. Chiapello & V. Loux

Helene.chiapello@inrae.fr
https://orcid.org/0000-0001-5102-0632

Valentin.loux@inrae.fr
https://orcid.org/0000-0002-8268-915X

# 2. Construction and analysis of prokaryotic genomic dataset

Outline

> ## 2.1 Constructing a genome dataset

> 2.2 Analyzing the genome dataset

> 2.3 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

Hélène Chiapello
Training

Valentin Loux
Technical coordinator

# Dataset construction

# 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 Prokka [*See module 9*]

- It's always better if annotation is homogeneous

# Quick reminder on format

# 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.

*Example*

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

# 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 annotaed with thier position , strand and qualifiers that contains the n 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 **"//"**.

- NCBI, ENA (European Nucleotide Archive) et DDBJ (Japan) entries are synchronized each day.
- Those three bank agree on the list of feature / qualifier that one can use to annotate sequence.

# 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
```

# 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=BioP
NZ_LHTK01000001    GenBank    pseudogene    1    1014    .    -    1    ID=LFZ49_RS22320.pseudogene;
NZ_LHTK01000001    GenBank    gene    1011    1634    .    -    1    ID=LFZ49_RS00010;Name=LFZ49_RS0
NZ_LHTK01000001    GenBank    mRNA    1011    1634    .    -    1    ID=LFZ49_RS00010.t01;Parent=LFZ
```

# Practical : public genomes

# How to filter and download publicly available genomes ?

- list all publicly available genomes with
  - metadata
  - quality metrics (size, completeness,…)
- filter according to above criteria
- download genomes in various formats

# A solution : NCBI 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.

# 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_ GenBank assembly | GCF_ RefSeq assembly |
|---|---|---|
| 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



NCBI Datasets Genome Page

Genome Table || Figure Source

- 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**

# NCBI Datasets : Command Line



COMMAND EXAMPLES

datasets **download** genome accession GCF_000001405.40
datasets *summary* genome taxon ailuridae

datasets **download** ortholog symbol brca1
datasets *summary* ortholog gene-id 672

datasets **download** gene accession NP_000483.3
datasets *summary* gene gene-id 40650

datasets **download** virus genome taxon sars-cov-2
datasets *summary* virus genome accession NC_045521.2
datasets **download** virus protein ORF10

datasets rehydrate --directory my_genomes

NCBI Datasets Command Line

**genome** options :

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

# NCBI Datasets : Aplication Programmatic Interface



NCBI Datasets Python API

Jupyter Notebook

# NCBI Datasets : Galaxy Integration



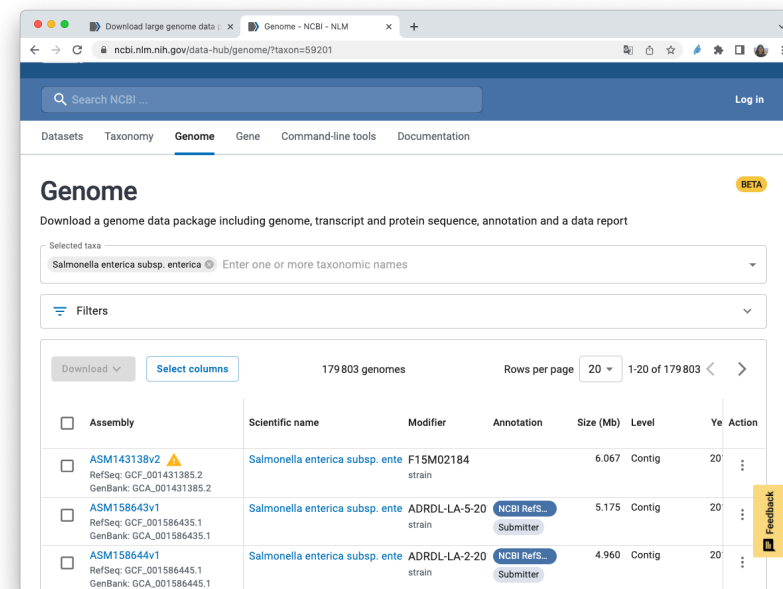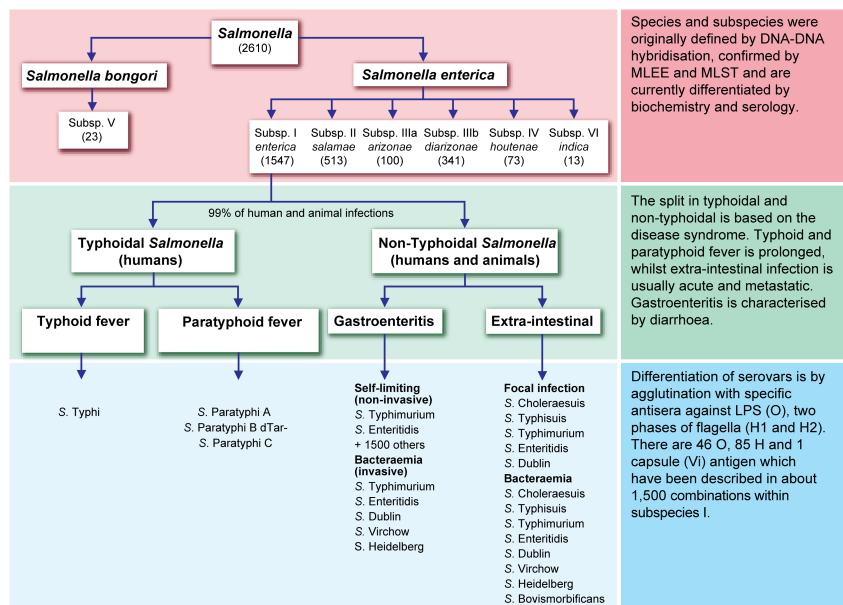A wrapper of the command line tool

Parameters to define packages files

# NCBI Datasets : Galaxy Integration

- A few caveats of the wrapper :

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

- What recomend to :

    - use the NCBI dataset genome page to browse / filter a list of genomes of interest
    - download the list as a `tsv` file
    - feed NCBI datset with the list to donwload the genomes in diverse formats

# The training datasets

We will work on 3 datasets of public *Salmonella* genomes



179.803 salmonella enterica enterica public assemblies at NCBI!!

# The training datasets

We will work on 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 from their *taxon id*
- dataset 3: download 16 Salmonella enterica public assemblies (2 sub-species, 4 serotypes) from their *accession numbers*
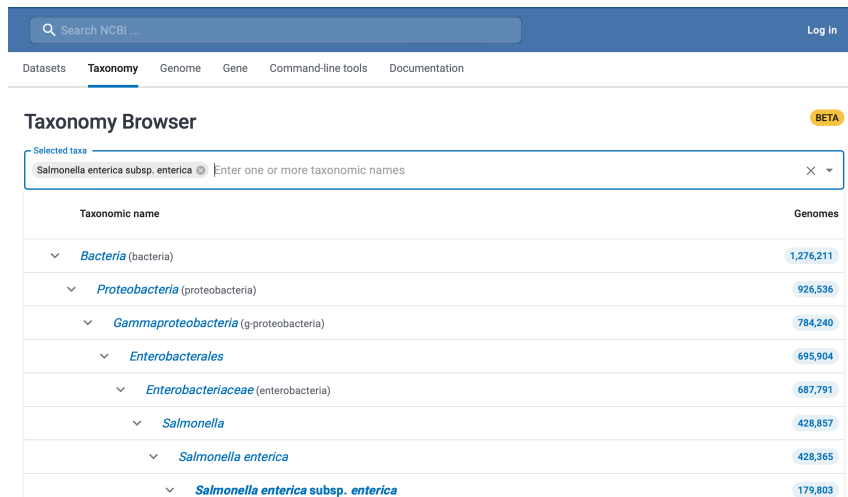
# Use case 1 : Taxonomy browser

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

# Use case 1 : genome table

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



Notice :

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

# Use case 2 : genome table

List all the *Salmonella bongori assemblies* to choose and download the best outgroup of a salmonella enterica dataset from their *taxon id*



Notice :

- Filter parameters :
  - "reference genome"
  - annotated by RefSeq
- Refseq accession for Galaxy

# Use case 2 : Galaxy dataset

- **List all** the *Salmonella bongori* assemblies using their species name or taxon id
- **Choose and Download** the best genome for being an outgroup of a *Salmonella enterica* dataset

# Training dataset 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 |

# Training dataset 3: 16 *S. enterica* public genomes (part 2)

| 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 |

# Use case 3 : from a tabular file

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

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