# Microbial sequence data

What's out there?

# Levels of sequence data

# Raw data

- Amplicon
  - PCR product, usually Sanger sequence (.ab1, .fasta)
- Locus
  - Multiple overlapping amplicons assembled (.fasta)
- Genome
  - Whole genome sequencing reads (.fastq.gz)
- Transcriptome
  - RNA (cDNA) sequencing reads (.fastq.gz)

# Derived data

- Assembled genome (.fasta)
  - Draft - multiple contigs
  - Complete - one contig per replicon

- Annotated genome (.gbk or .gff)
  - Genomic features labelled *e.g.* genes

- Protein sequences
  - Translated from predicted genes
  - Find in assembled transcripts

# Curated data

- "Curated" means
    - Assessed for quality
    - Usually some human contribution

- Orthologs
    - Protein families

- Sequence "profiles"
    - Alignments of orthologous sequences
    - DNA or Protein

# The INDSC

# INSDC International Nucleotide Sequence Database Collaboration

## International Nucleotide Sequence Database Collaboration

- The International Nucleotide Sequence Database Collaboration (INSDC) is a long-standing foundational initiative that operates between DDBJ, EMBL-EBI and NCBI. INSDC covers the spectrum of data raw reads, though alignments and assemblies to functional annotation, enriched with contextual information relating to samples and experimental configurations.

| Data type | DDBJ | EMBL-EBI | NCBI |
|---|---|---|---|
| Next generation reads | Sequence Read Archive | European Nucleotide Archive (ENA) | Sequence Read Archive |
| Capillary reads | Trace Archive | | Trace Archive |
| Annotated sequences | DDBJ | | GenBank |
| Samples | BioSample | | BioSample |
| Studies | BioProject | | BioProject |

Hosted at NCBI in Washington, USA

Hosted at EBI in Cambridge UK

Hosted at NIG in Mishima, Japan

# Sequence Read Archive (SRA)

# The FASTQ dumping ground



- NCBI → "SRA"
  Download ".sra" files
  Convert with "fastq-dump"

- ENA → "ENA"
  Download ".fastq.gz" files

- DDBJ → "DRA"
  ???

**Sample 1**



**Sample 2**



**Study**

**Sample 3**

Sample 1

Sample 2

Experiment 2

Experiment 1

Study

Sample 3    Experiment 3

**Sample 1**

**Sample 2**

**Experiment 2**

**Experiment 1**

**Study**

**Run 2**

**Sample 3**

**Experiment 3**

**Run 3**

**Run 1**

**Sample** · **Experiment** · **Run**

one ← many

| Sample | ← | Experiment |
| Experiment | ← | Run |
| Study | ← | Experiment |

# Reference genomes

# Reference genomes (ENA)

http://www.ebi.ac.uk/genomes/bacteria.html

When downloading use *"Sequence - Plain"*

List of available genomes (on 5-MAY-2015)

| | Description | Length (bp) | Sequence | | Project |
|---|---|---|---|---|---|
| | | | **Plain** | **HTML** | |
| | *Acaryochloris marina* | | | | |
| 1 | Acaryochloris marina MBIC11017 | 6,503,724 | CP000828 | CP000828 | PRJNA12997 |
| | *Acetobacter pasteurianus* | | | | |
| 2 | Acetobacter pasteurianus 386B | 2,818,679 | HF677570 | HF677570 | PRJEB1172 |
| 3 | Acetobacter pasteurianus IFO 3283-01 | 2,907,495 | AP011121 | AP011121 | PRJDA31129 |
| 4 | Acetobacter pasteurianus IFO 3283-01-42C | 2,815,241 | AP011163 | AP011163 | PRJDA31141 |

# Curated databases

## Available Databases

### Salmonella
**Strains:186080**

Assembled
- Legacy:7229
- From NGS:178851
- In Progress:5

Schemes
- rMLST:177753
- Achtman 7 Gene MLST:184953
- cgMLST V2:177394
- wgMLST:174068
- CRISPR:51158

Database Home     ⊙

### Escherichia/Shigella
**Strains:92429**

Assembled
- Legacy:9611
- From NGS:82818
- In Progress:2

Schemes
- wgMLST:81297
- Achtman 7 Gene MLST:92123
- rMLST:82304
- cgMLST V1:82096

Database Home     ⊙

### Clostridioides
**Strains:7215**

Assembled
- From NGS:7215
- In Progress:0

Schemes
- Griffiths 7 Gene:7208
- cgMLST V1:7205
- rMLST:7206
- wgMLST:7202

Database Home     ⊙

### Vibrio
**Strains:6880**

Assembled
- From NGS:6880
- In Progress:1

Schemes
- rMLST:6876

Database Home     ⊙

### Yersinia
**Strains:3596**

Assembled
- Legacy:1165
- From NGS:2433
- In Progress:1

Schemes
- Achtman 7 Gene:3229
- McNally 7 Gene:2795
- cgMLST V1:2433
- rMLST:2430
- wgMLST:2432

Database Home     ⊙

### Moraxella
**Strains:557**

Assembled
- Legacy:420
- From NGS:137
- In Progress:0

Schemes
- Achtman 7 Gene:559
- rMLST:137

Database Home     ⊙

### Helicobacter
**Strains:535**

Assembled
- From NGS:535
- In Progress:0

Schemes
- rMLST:531

Database Home     ⊙

### Dev. Sandbox
**Strains:128**

Assembled
- From NGS:128
- In Progress:3

Schemes
- Achtman 7 Gene:67
- rMLST:64

Database Home     ⊙

on

ta, please

**BBSRC**
bioscience for the future

## Search ➤

**Search** our comprehensive database for:

- Genomes
- Genes & proteins
- Immune epitopes
- 3D protein structures
- Host Factor Data
- Antiviral Drugs

**Browse All Search Types**

## Analyze ➤

**Analyze** data online:

- Sequence Alignment
- Phylogenetic Tree
- Sequence Variation (SNP)
- Metadata-driven Comparative Analysis
- BLAST

**Browse All Tools**

## Save to Workbench

Sign up for a **workbench** to:

- Store and share data
- Combine working sets
- Integrate your data with ViPR data
- Store and share analyses
- Custom search alert

**Sign In**

## ◢ Virus Families

Click on icon of family or species of interest. Click **here** to to view all families and species in list format. Don't know family of species? | Provide species name |

### Single-Stranded Positive-Sense RNA

- *Caliciviridae*
- *Hepeviridae*
- *Coronaviridae*
- *Picornaviridae*
- *Flaviviridae*
- *Togaviridae*

### Single-Stranded Negative-Sense RNA

- *Arenaviridae*
- *Paramyxoviridae*
- *Bunyaviridae*
- *Rhabdoviridae*
- *Filoviridae*

### Double-Stranded RNA

- *Reoviridae*

### Double-Stranded DNA

- *Herpesviridae*
- *Poxviridae*

http://plasmodb.org/

# Conclusions

# Conclusions

- **Be skeptical!**
  - ENA + Genbank accept anything
  - Garbage in, Garbage out

- **Curated data**
  - Refseq vs GenBank
  - Specialised sites (better QC)
  - Be wary of draft assemblies
  - Go back to primary reads

Further reading

https://www.ncbi.nlm.nih.gov/core/assets/sra/files/Factsheet_SRA.pdf

https://www.ddbj.nig.ac.jp/dra/index-e.html

https://p.ddbj.nig.ac.jp/pipeline/Login.do