# Les databases

Pratique #1

# Genome DataBases

- All published genome sequence have to be available in a public DB

  - Consortia made by 3 big DBs

    - EMBL (European Molecular Biology Laboratory nucleotide sequence database at EBI, Hinxton, UK)

    - GenBank (at National Center for Biotechnology information, NCBI, Bethesda, MD, USA)

    - DDBJ (DNA Data Bank Japan at CIB , Mishima, Japan)

    - Ces grandes banques généralistes s'échangent systématiquement leur contenu depuis 1987 et adoptent un système de conventions communes (The DDBJ/EMBL/GenBank Feature Table Definition).

  - GOLD ("Genomes OnLine Database") : base de données qui recense les milliers de génomes séquencés ou en voie de séquençage.
  - PATRIC ("the Pathosystems Resource Integration Center")
  - MAGE (Genoscope)

# Retrieve genomes

- "Assemblies" from NCBI /!\

Species in repositories of WGS data are often mislabelled! **Trust no one!**



**☰ NCBI**  Resources ⌄  How To ⌄

**Genome**        [ Genome ⌄ ] [                    ]
                        Limits   Advanced

Organism Overview ; **Genome Assembly and Annotation report [401]** ; Genome Tree report [401] ; Plasmid Annotation Report [31]

## Staphylococcus epidermidis

[                                                                    ] [ Search ] [ Clear ]

**Anomalous: All  Levels:** ☑ All  ☑ ⬤Complete [8]  ☑ ◖Chromosome [1]  ☑ ◐Scaffold [139]  ☑ ◔Contig [253]

| Organism/Name | Strain | CladeID | BioSample | BioProject | Assembly | Level | Size (Mb) | GC% | Replicons | |
|---|---|---|---|---|---|---|---|---|---|---|
| Staphylococcus epidermidis | 1312_SEPI | 19672 | SAMN03197289 | PRJNA267549 | GCA_001070555.1 contaminated | ◐ | 7.72405 | 49.90 | - | |
| Staphylococcus epidermidis | 6_SEPI | 19510 | SAMN03197799 | PRJNA267549 | GCA_001073525.1 contaminated | ◐ | 6.00094 | 36.20 | - | |
| Staphylococcus epidermidis | 110_SEPI | 19993 | SAMN03197063 | PRJNA267549 | GCA_001070175.1 contaminated | ◐ | 4.79355 | 31.60 | - | NO !!! |
| Staphylococcus epidermidis | 926_SEPI | 20104 | SAMN03198144 | PRJNA267549 | GCA_001075745.1 contaminated | ◔ | 4.38758 | 36.00 | - | |
| Staphylococcus epidermidis | 1068_SEPI | 19993 | SAMN03197030 | PRJNA267549 | GCA_001068615.1 contaminated | ◔ | 3.5814 | 37.50 | - | |
| Staphylococcus epidermidis | 114_SEPI | 19993 | SAMN03197095 | PRJNA267549 | GCA_001068755.1 | ◐ | 2.91249 | 31.70 | - | YES |
| Staphylococcus epidermidis | ABKUX | 19993 | SAMN04286996 | PRJNA302961 | GCA_002221835.1 | ◔ | 2.85146 | 31.90 | - | |

- **Contigs** = first **level**
- **Scaffolds** (supercontigs) = place several contigs in the correct order and orientation and represent sequencing gaps between the contigs with series of NNN's
- **Chromosome** = generally 1record for each chromosome (with N's)
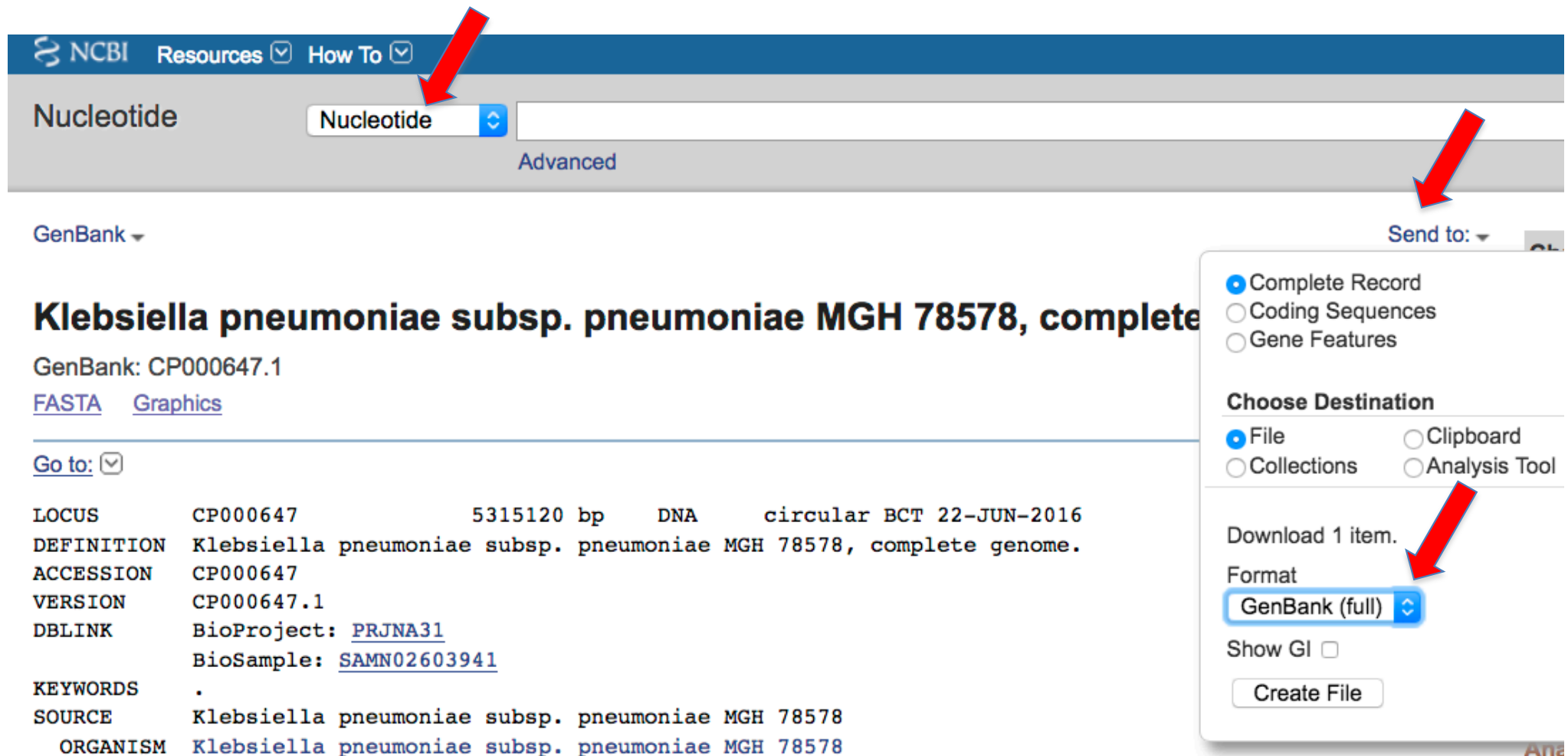- **Complete =** assemblies without sequencing gaps

# GenBank Format

suite

```
LOCUS       SCU49845      5028 bp     DNA             PLN       21-JUN-1999
DEFINITION  Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Axl2p
            (AXL2) and Rev7p (REV7) genes, complete cds.
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)
  MEDLINE   95176709
  PUBMED    7871890
REFERENCE   2  (bases 1 to 5028)
  AUTHORS   Roemer,T., Madden,K., Chang,J. and Snyder,M.
  TITLE     Selection of axial growth sites in yeast requires Axl2p, a
novel
            plasma membrane glycoprotein
  JOURNAL   Genes Dev. 10 (7), 777-793 (1996)
  MEDLINE   96194260
  PUBMED    8846915
REFERENCE   3  (bases 1 to 5028)
  AUTHORS   Roemer,T.
  TITLE     Direct Submission
  JOURNAL   Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University,
New
            Haven, CT, USA
FEATURES             Location/Qualifiers
  source          1..5028
                  /organism="Saccharomyces cerevisiae"
                  /db_xref="taxon:4932"
                  /chromosome="IX"
```

```
  gene            687..3158
                  /gene="AXL2"
  CDS             687..3158
                  /gene="AXL2"
                  /note="plasma membrane glycoprotein"
                  /codon_start=1
                  /function="required for axial budding pattern of S.
                  cerevisiae"
                  /product="Axl2p"
                  /protein_id="AAA98666.1"
                  /db_xref="GI:1293615"
                  /translation="MTQLQISLLLTATISLLHLVVATPYEAYPIGKQYPPVARVNESF
                  TFQISNDTYKSSVDKTAQITYNCFDLPSWLSFDSSSRTFSGEPSSDLLSDANTTLYFN
                  VILEGTDSADSTSLNNTYQFVVTNRPSISLSSDFNLLALLKNYGYTNGKNALKLDPNE
                  VFNVTFDRSMFTNEESIVSYYGRSQLYNAPLPNWLFFDSGELKFTGTAPVINSAIAPE
                  TSYSFVIIATDIEGFSAVEVEFELVIGAHQLTTSIQNSLIINVTDTGNVSYDLPLNYV
                  YLDDDPISSDKLGSINLLDAPDWVALDNATISGSVPDELLGKNSNPANFSVSIYDTYG
                  DVIYFNFEVVSTTDLFAISSLPNINATRGEWFSYYFLPSQFTDYVNTNVSLEFTNSSQ
                  DHDWVKFQSSNLTLAGEVPKNFDKLSLGLKANQGSQSQELYFNIIGMDSKITHSNHSA
                  NATSTRSSHHSTSTSSYTSSTYTAKISSTSAAATSSAPAALPAANKTSSHNKKAVAIA
                  CGVAIPLGVILVALICFLIFWRRRRENPDDENLPHAISGPDLNNPANKPNQENATPLN
                  NPFDDDASSYDDTSIARRLAALNTLKLDNHSATESDISSVDEKRDSLSGMNTYNDQFQ
                  SQSKEELLAKPPVQPQPESPFFDPQNRSSSVYMDSEPAVNKSWRYTGNLSPVSDIVRDS
                  YGSQKTVDTEKLFDLEAPEKEKRTSRDVTMSSLDPWNSNISPSPVRKSVTPSPYNVTK
                  HRNRHLQNIQDSQSGKNGITPTTMSTSSSDDFVPVKDGENFCWVHSMEPDRRPSKKRL
                  VDFSNKSNVNVGQVKDIHGRIPEML
BASE COUNT     1510 a   1074 c    835 g   1609 t
ORIGIN
        1 gatcctccat atacaacggt atctccacct caggtttaga tctcaacaac ggaaccattg
       61 ccgacatgag acagttaggt atcgtcgaga gttacaagct aaaacgagca gtagtcagct
      121 ctgcatctga gccgctgaa gttctactaa gggtggataa catcatccgt gcaagaccaa
      181 gaaccgccaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaaccg

        ...
//
```

# Download genbank @NCBI nucleotide

# Sequences brutes @NCBI

# Sequences brutes @ENA

# Protein and domain DataBases

- Uniprot
  - ("Universal Protein Resource") : c'est la base de données des protéines : ExPASy Proteomics Server. Consortium [EBI - SIB - PIR].

- PDB ("Protein Data Bank")

- PFAM et INTERPRO