# What is a Database?

- A ***database*** is a computerized archive used to store and organize data in such a way that information can be retrieved easily via a variety of search criteria.

# Types of Databases

## 3 Major types of electronic databases:

1. **Flat files** - simple text files
   - no organization to facilitate retrieval



2. **Relational** - data organized as tables ("relations")
   - shared features among tables allows rapid search



3. **Object-oriented** - data organized as "objects"
   - objects associated hierarchically

# Biological Databases

- Currently in all 3 types
  - MANY flat files despite the obvious drawbacks of them

- What are goals of biological databases?
  - Information retrieval
  - Knowledge discovery

# Types of Biological Databases

1. **Primary**
   - Simple archives of sequences, structures, images, etc.
   - Raw data, minimal annotations, not always well curated!

2. **Secondary**
   - Enhanced with more complete annotation of sequences, structures, images, etc.
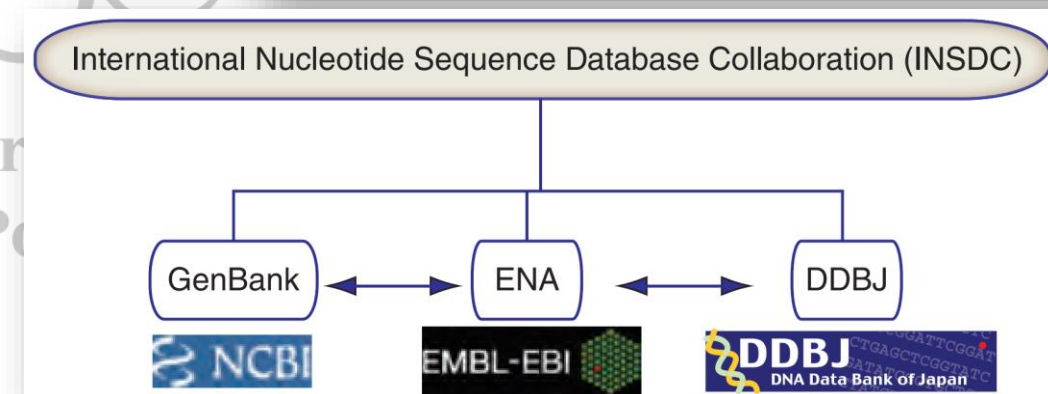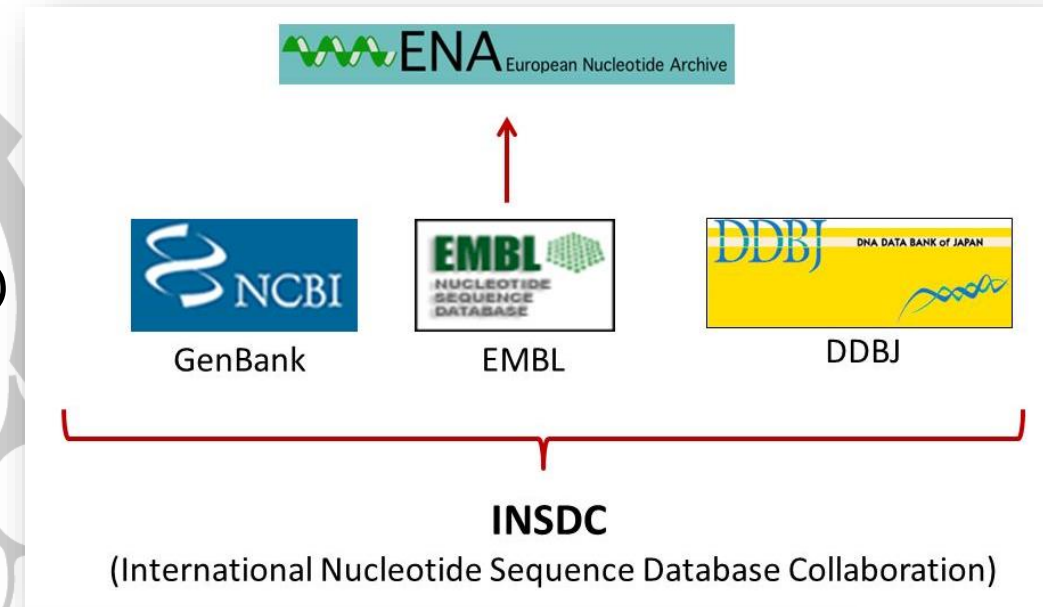   - Usually curated!

3. **Specialized**
   - Focused on a particular research interest or organism
   - Usually highly curated

# Examples of Biological Databases

**1- Primary**

- **DNA sequences**
  - GenBank - US
  - European Molecular Biology Lab (EMBL)
  - DNA Data Bank of Japan (DDBJ)

- **Structures (Protein, DNA, RNA)**
  - PDB - Protein Data Bank (3D)
    - Flat file format
    - Archives atomic coordinates
  - NDB - Nucleic Acid Databank

# Examples of Biological Databases

**2- Secondary**

- **Protein sequences**
  - Swiss-Prot, TreEMBL, PIR
  - These recently combined into UniProt
  - Pfam, Blocks, DALI

**3- Specialized**

- **Species-specific (or "taxonomic" specific)**
  - HIV sequence, Flybase, WormBase, AceDB, PlantDB, GenBank EST
- **Molecule-specific, disease-specific**

# Pitfalls of Biological Databases

- There are many errors in sequence databases
- Lack of documentation: quality or reliability of data
- Limited mechanisms for "data checking" or preventing propagation of errors  *(esp. annotation errors!!)*
- Redundancy: NCBI has now created a *nonredundant* database, called RefSeq
- Inconsistency in annotation
- Incompatibility (format, terminology, data types, etc.)

# Information Retrieval from Biological Databases

- **2 most popular retrieval systems:**

  - **ENTREZ -** developed and maintained by **NCBI**

  - **Sequence Retrieval Systems (SRS) -** maintained by **EBI**

- **Both:**

  - Provide access to multiple databases

  - Allow complex queries

# References

- Mostly used:
  - Essential bioinformatics, Chapter 2 (Introduction to Biological Databases)

- Second reference:
  - Bioinformatics and functional genomics, Chapter 2 (Access to Sequence Data and Related Information)

# Thanks for your attention