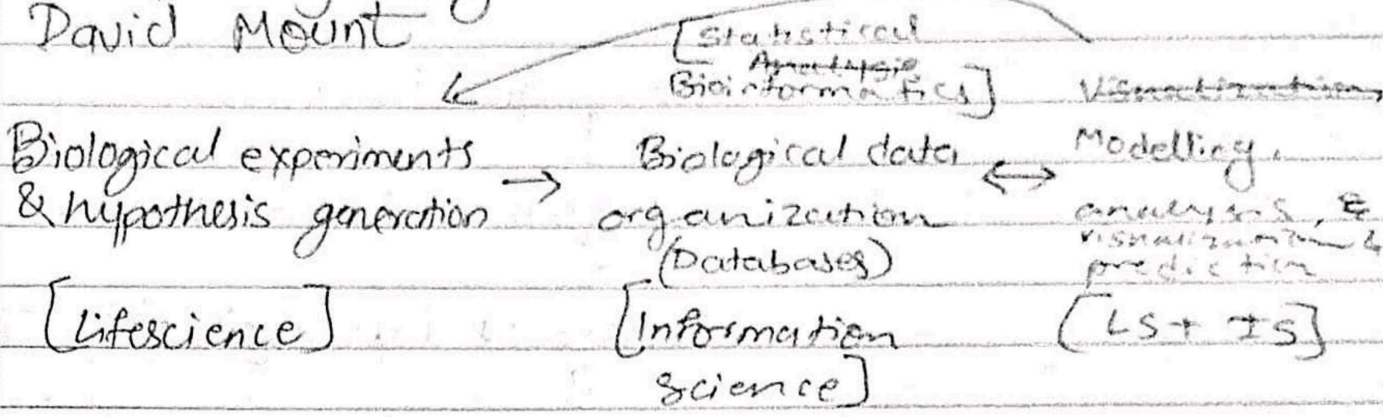


BOOK → Pavlauer  
Essentials of Bioinformatics  
↳ by Zing Zeng  
↳ David Mount



- Term Bioinformatics was invented by Paulien Hogeweg and Ben Hesper in 1970 as Study of informatic processes in biotic systems  
↳ Paulien Hogeweg is a dutch-theoretical biologist & complex systems researchers study

- Bioinformatics is a science that is not only used to also develop algorithms, store, retrieve, organize & analyze biological data but to curate data.

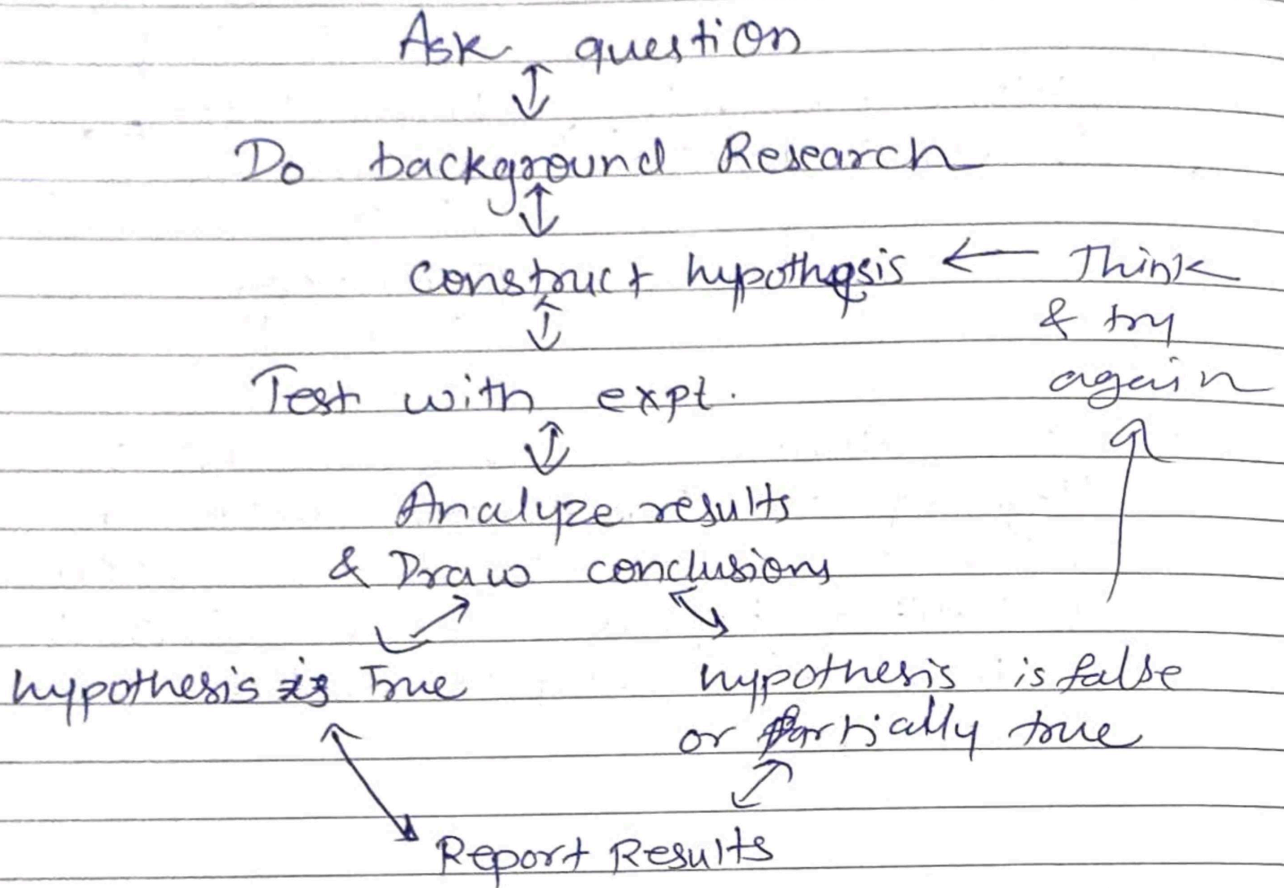
- Components of Bioinformatics
  - Biological data
  - Algorithms
  - Computers

- Speeds data production & data interpretation

- Bioinformatics = Biologist + CS + IT

- Human genome  $\Rightarrow$  3.5 billion basepairs

$\Rightarrow$  How to work with bioinfo

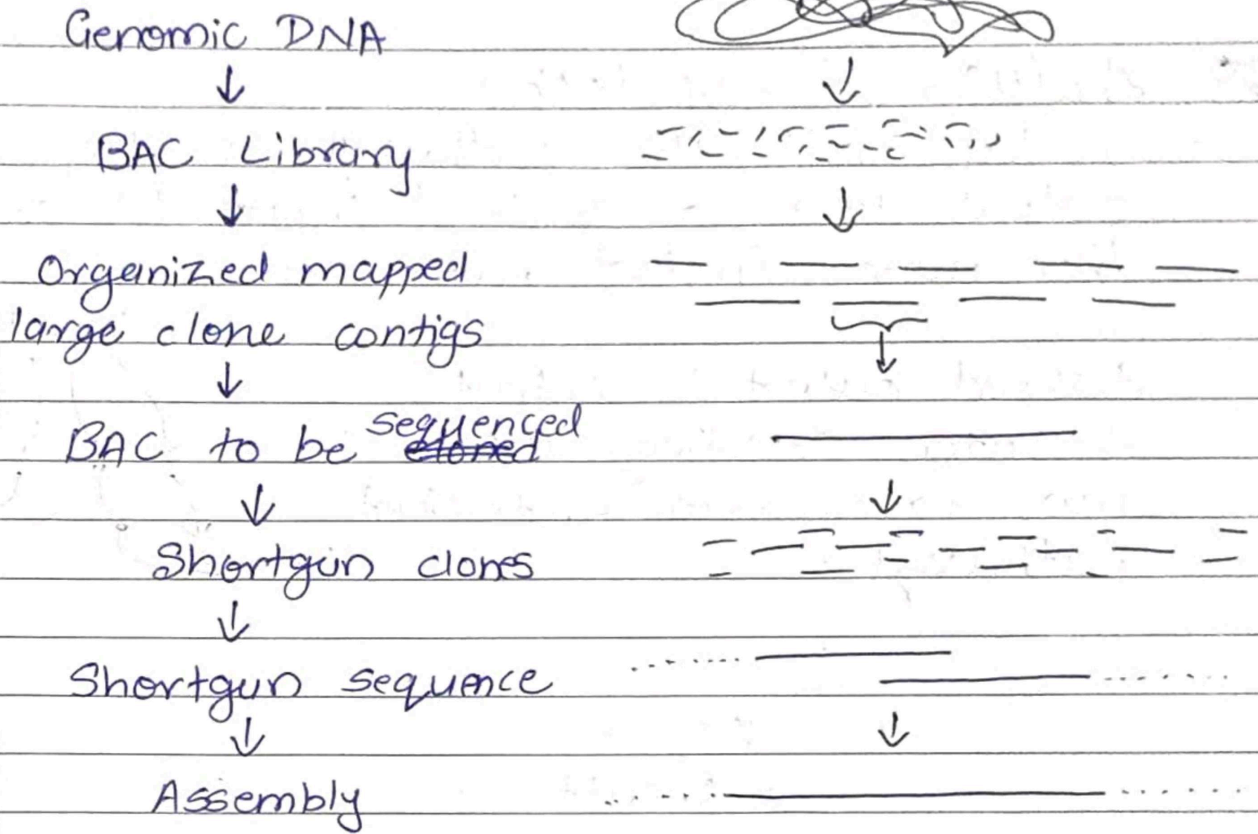




# \* Sequence Assembly

→ NGS [Shotgun]

Sequence Repository Database (SRA)



→ Illumina is the base platform for NGS

# \* Genome annotation

↳ naming/labeling of the genes/genome that comes from NGS.

- annotations can be done by researchers or NCBI. NCBI is preferred for cross-checking.
- Genome Browsers → UCSC, ENSEMBL, NCBI, etc.
- Genome tracks
- Data validation via cross-checking & backtracking

# \* Molecular evolution

↳ Eg. Phylogenetic analysis

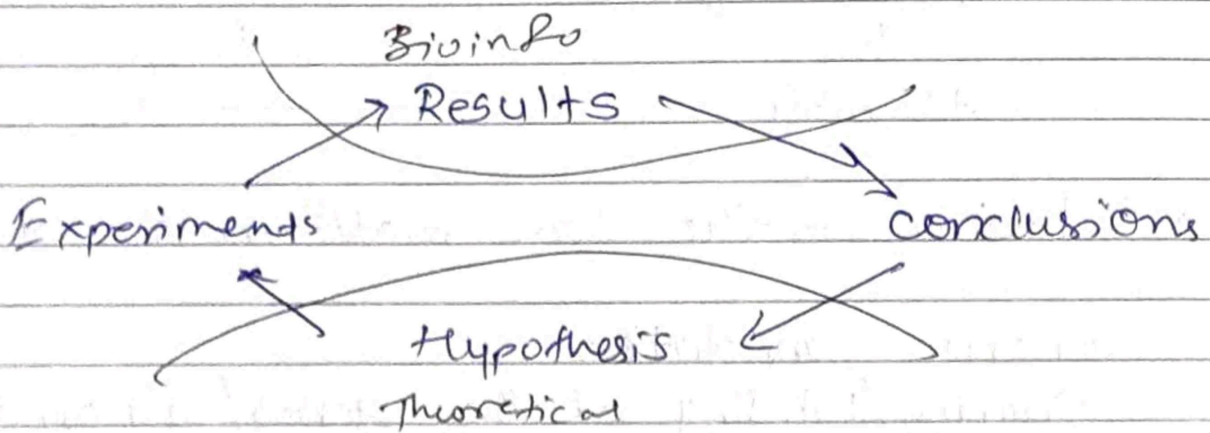
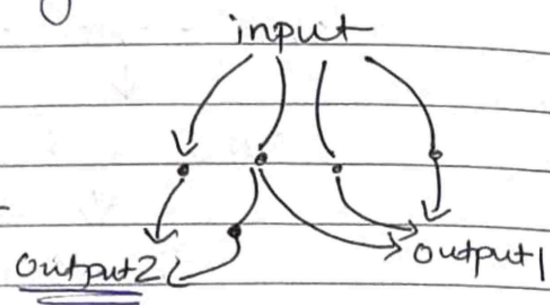
### \* Analysis of gene expression.

- to find out which changes in which gene or its expression causes which changes (GEO)
- Database  $\Rightarrow$  Gene Expression

### \* Analysis of regulation.

- study all the pathways for the particular output from a specified input to select the most efficient pathway.

desired output is output 2 & later selecting the most favourable & efficient pathway.





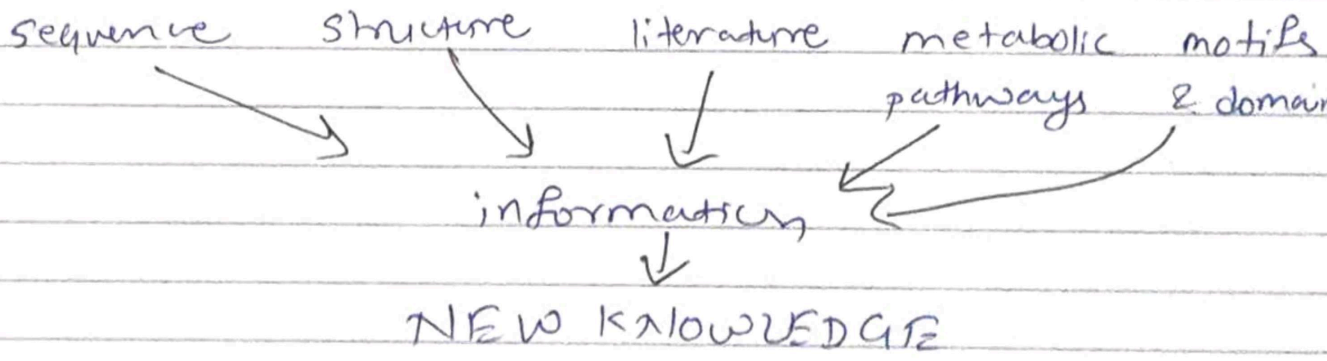
# Pairwise Sequence Alignment

## File formats for biomolecular sequences

### ⇒ Aims

- to understand the conventions regarding the presentation of DNA & Protein sequence info.
- to understand the logic underlying these conventions.
- to become familiar with commonly used sequence file formats
- to become familiar with READSEQ programme for the interconversion of file formats
- present a nucleotide or protein

### ⇒ where does data comes from



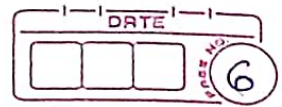
### ⇒ Database / Resource

- Collection of data in related format
  - ↳ structured
  - ↳ updated periodically (releases)

Crosslinked  
data within  
all the dbs  
at NCBI

NCBI ← Sequence

PDB ← structure



→ Includes associated tools/software necessary for accessing database, updating DB, DB information insertion, DB information deletion, etc.

→ Type & content of data.

→ Sequence, structure, nucleic acid or Protein  
→ imp. biological info such as about enzyme & their pathways, mutations, drugs, diseases, images, etc.

→ Based on source of data

- primary db → raw info shared to portal by user
- secondary db → curated info
- knowledge bases → open access to all/specialized info.
- integrated db → cross-linking of data.

Amino acid  
(one & 3 letter code)

physiochem  
prop.

structure

⇒ DNA



17 8 51

DATE		

Page No. 8



⇒ File types

Myoglobin  
↓

H<sub>2</sub>N-Met - Val - Tyr - Gly - Iso - Lys - COOH

↓  
M V Y G I K

single letter protein

- <sup>software</sup> many packages have been developed for the analysis of DNA and protein sequences
  - ↳ these software stores & analyzes the DNA & protein sequence in different file formats.
- The variety of software packages will usually only accept a specific file format.
- Diff. databases holds info. in different formats.
  - ↳ This makes the situation worse.
  - ↳ in order to use this data, the info has to be modified & desired file format must be created.
- It is an essential skill to be able to recognize the diff. file formats & to be able to interconvert files between formats.

## ⇒ Primary Biological Databases

→ Nucleic acid

- ① EMBL
- ② GenBank
- ③ DDBJB

→ Protein

- ① PIR
- ② MIPS
- ③ SWISS-PROT
- ④ TrEMBL
- ⑤ NRL-3D

### → Nucleotide DB

- ① EMBL → Nucleotide seq. db
- ② Ensembl → Automatics annotation of eukarya genomes
- ③ Genome server → Overview of completed genomes at EBI
- ④ Genome MOT → Genome monitoring table
- ⑤ EMBL-Align → Multiple sequence alignment db
- ⑥ Parasites → Parasite genome db
- ⑦
- ⑧

### → EMBL/GenBank/DDBJB

- same info (diff syntax & format)
- non-confidential data is exchanged daily.



→ DB related to Genomics.

- info on genes, gene mapping (location), nomenclature & links to seq. db.
- exists for most org. imp. for life sci. research
- Eg. ① MIM

②

③

④

⑤

⇒ Plain Sequence Format.

- may contain only IUPAC characters & spaces but not numbers.
- file in (PSF) may contain one sequence, while

⇒ FASTA Format (FF)

- can contain several sequence
- begins with greater than sign (>) & single-line description, followed by the sequence on next line. & ends with 2 slashes (//)

⇒ EMBL Format (EF)

- contains several seq.

⇒ GenBank Format (GF)

### \* ENSEMBL

- contains all human genome DNA seq. currently available in public domain.
- Automated annotation
  - ↳ by using diff. software tools, features are identified in DNA sequence.
  - ①
  - ②
  - ③
  - ④
- Maintained by EMBL & Sanger's



## → Protein Databases

### → SWISS-PROT

- Annotated Sequence DB

### → TrEMBL

- DB of EMBL nucleotide translated sequences

### → InterPro

- Integrated res. for protein families, domains & functional sites.
- best for cross-referencing proteins

### → CluSTr

- offers automatic classification of SWISS-PROT & TrEMBL

### → IPI

- A non-redundant human proteome set constructed with SWISS-PROT, TrEMBL, Ensembl, RefSeq.
- no duplicated data

### → GO

- Provides assignments of gene products of Gene Ontology (GO) resource

### → Proteome

- Statistical & comparative analysis of predicted proteomes of fully

### → Protein profiles

- Tables of SWISSPROT & TrEMBL entries & alignment for protein families of protein profile

→

## IntEnz

- Integrated relational Enzyme database (IntEnz) will contain enzyme data approved by Nomenclature Committee

★

## SWISS PROT

- Annotated protein sequence db, established in 1986
- maintained collaboratively by Dept. of Medical Biochemistry, University of Geneva & EBI, since 1987.
- complete, curated, non-redundant & cross-referenced with 34 other dbs
- Highly cross-referenced.
- more than 80000 species are available in DB.
- The first 20 results give about 42%.

⇒

## SWISS PROT file format.



- ★ TrEMBL [Translation of EMBL]
- computer-annotated supplement of SWISS-PROT, as it is impossible to cope up with the flow of data. ~~using only EMBL~~
  - Well structure SWISS-PROT-like resource
  - It contains the data that isn't yet available in SWISS-PROT

- ★ Protein DB (PDB)
- imp in solving real problem in molecular biology
  - PDB, established in 1972 at Brookhaven National Laboratory (BNL)
  -

- \* - PubMed contains summarized or brief information about the specific paper or topic. Whereas, PubMed Central (PMC) contains the whole scientific or research papers or articles.
- Amount of data in PMC is much more than it is in PubMed.



