

UCS2135 BIOMINFORMATICS TECHNOLOGIES

Unit 1

* Bioinformatics - a field of science in which biology, computer science, and information technology merge into a single discipline to analyze biological information using computers and statistical techniques.

Major Aspects: (i) well organized DBs

- (i) computationally-derived hypothesis
 - (ii) Web servers (online tools) / applications
 - (iv) Big data analysis
 - (v) Virtual screening of compounds for drug development

* Central Paradigm of Bioinformatics

Genetic Information → Molecular Structure → Biochemical Function → Phenotype
Structure → Biochemical Function → Phenotype (symptoms)

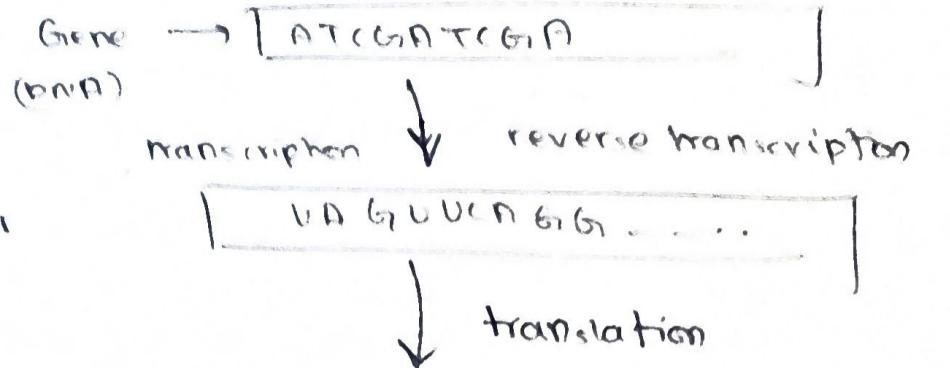
* Challenges understanding Genetic Information

- (i) Genetic info is redundant
 - (ii) Structural info is redundant
 - (iii) Genes and proteins are meta-stable
 - (iv) Genes and proteins are 1D but their function depends upon 3D structure

* Central Dogma of Molecular Biology

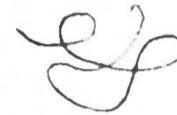
Genotype → Phenotype
(AA) (pink)

replication → ↗ ↘

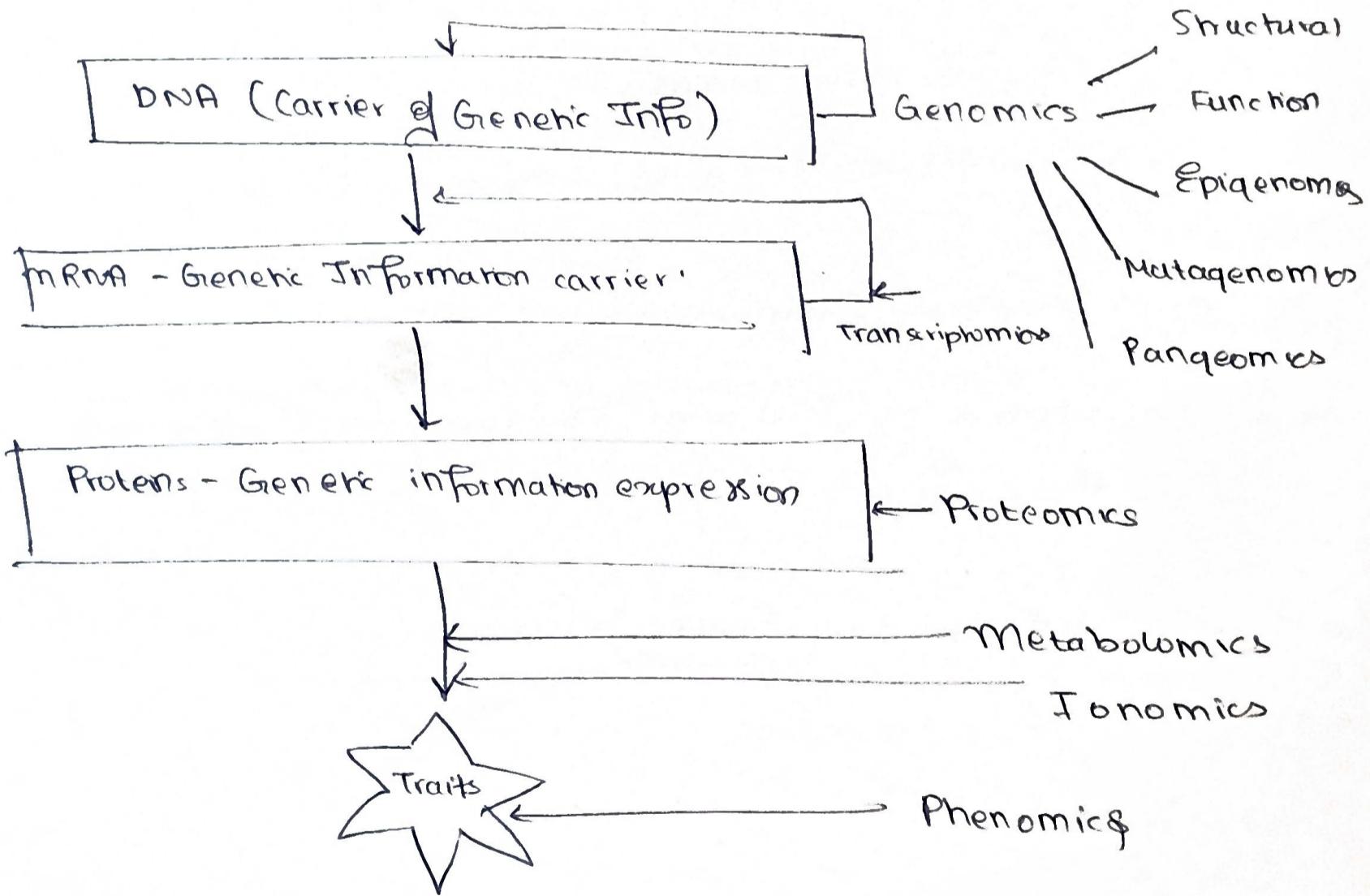
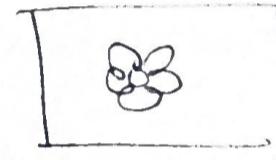


messenger
RNA

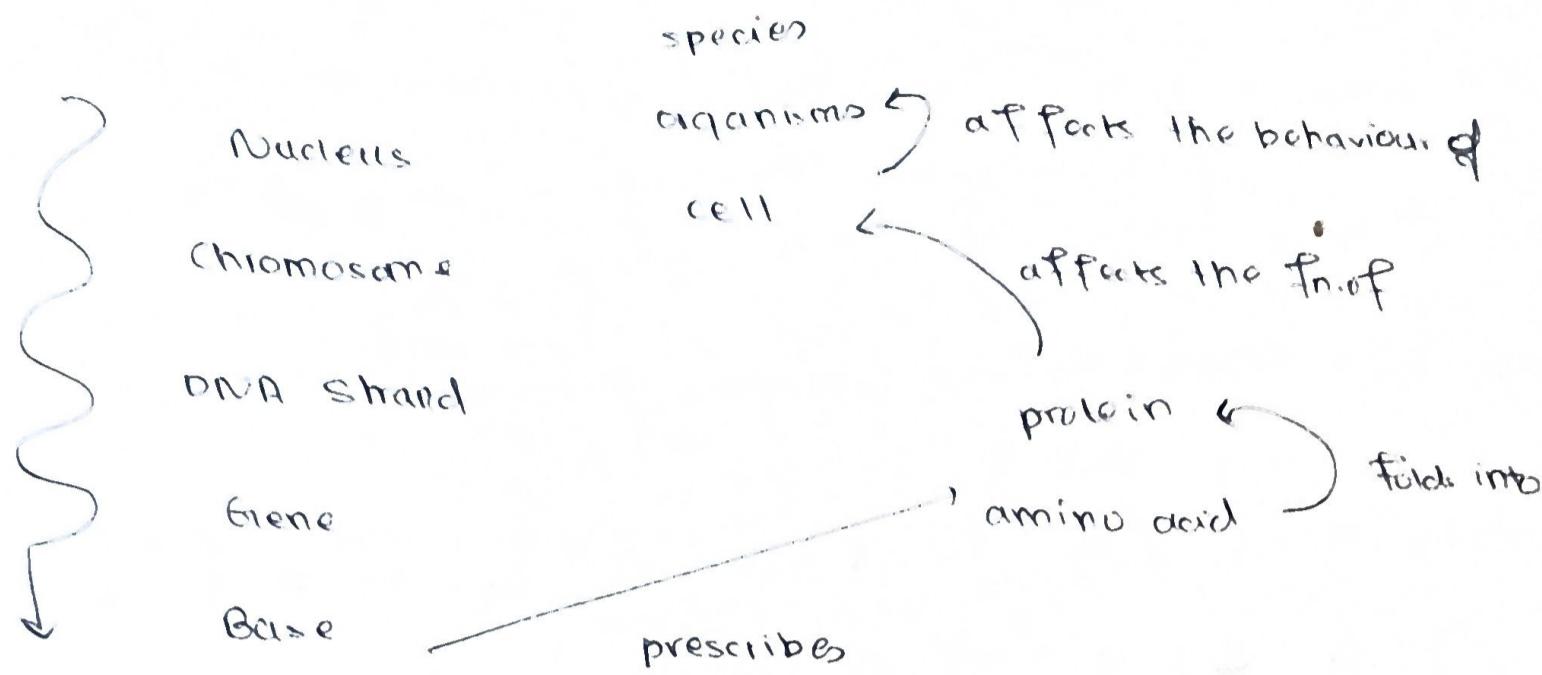
Protein



Trait:



* Substructure and Effect



* Synthesis of Proteins

- (i) DNA double helix is unzipped
- (ii) One strand is transcribed to messenger RNA
- (iii) RNA acts as a template
- (iv) Amino acid sequences fold into a 3d molecule

Transcription

→ write base counterparts

A ↔ T
G ↔ C

→ change T → U
Thymine to Uracil

For eg. G G A T G C C A A T G

C C T A C G G T T A G

C C U A C G G U U A G

- Translation
- Input triples of 3 base letters (codon)
 - Output = amino acid
 - / e.g. ACC becomes Threonine

e.g. Transcribe and Translate the following.

$$\begin{array}{l} A \leftrightarrow T \\ G \leftrightarrow C \end{array}$$

T C G G T G A A T C T G T T T G A T
A G C C A G T T A G A C A A A C T A
 Ser His Leu Asp Lys Phe
 ↓ ↓ J. ↓ ↓ K ↓ L

<u>A G C</u>	<u>C A G</u>	<u>U U A</u>	<u>G A C</u>	<u>A A A</u>	<u>C U A</u>
Ser	His	Leu	Asp	Lys	Phe
↓	↓	J.	↓	↓	↓
S	His	L	P	K	L

- * Evolution of Genes
- mutation
 - radiation from toxicity
 - Deletion, Insertion, Substitution

* Applications of Bioinformatics

① Pharmaceuticals

1. Biobank - biorepository that stores human biological samples
 - used in genomics, personalized medicine
 - helps accelerate the discovery and development of drugs

2. Veterinary Science / Animal Health - covers diagnostics,

medicines and vaccines for farmed animals and pets

- due to a growing demand for animal proteins as well as
- a strong consumer need for companion animal health care

② Ramifications

1. Pharmacogenomics - not all drugs work on all patients, can cause death

- do gene analysis before treatment
- customized treatment

2. Gene Therapy - replace or supply the defective or missing gene

- e.g. Insulin and Factor VIII or hemophilia

③ Diagnosis of Disease

- identification of genes which cause the disease will help detect disease at an early stage - e.g. Huntington disease



weird movements

personality changes & intellectual impairment

- excessively repeated sections of

c n o

④ Drug Design

- expensive and time consuming
- computational methods can improve testing methods

⑤ Drug Discovery

- identify the molecule on which germs rely for its survival
- develop a drug that will bind to the target
- So the germ will not be able to interact with the target
- Proteins are the most common targets

e.g. HIV produces HIV protease which is a protein and which in turn eats other proteins

HIV protease has an active site where it binds to other molecules

⑥ Phylogenetic Trees

- understand our genes
- genes are homologous

⑦ Predicting Protein Structure

- proteins fold to set up an active site
- has a small but highly effective substructure
- activity sites determine the activity of the protein
- no one has found rules governing how proteins fold

* Need for Bioinformatics Technologies

intersection of fields

computer and statistical fields

handle voluminous biological data.

identify interesting patterns

* Bioinformatics Technologies

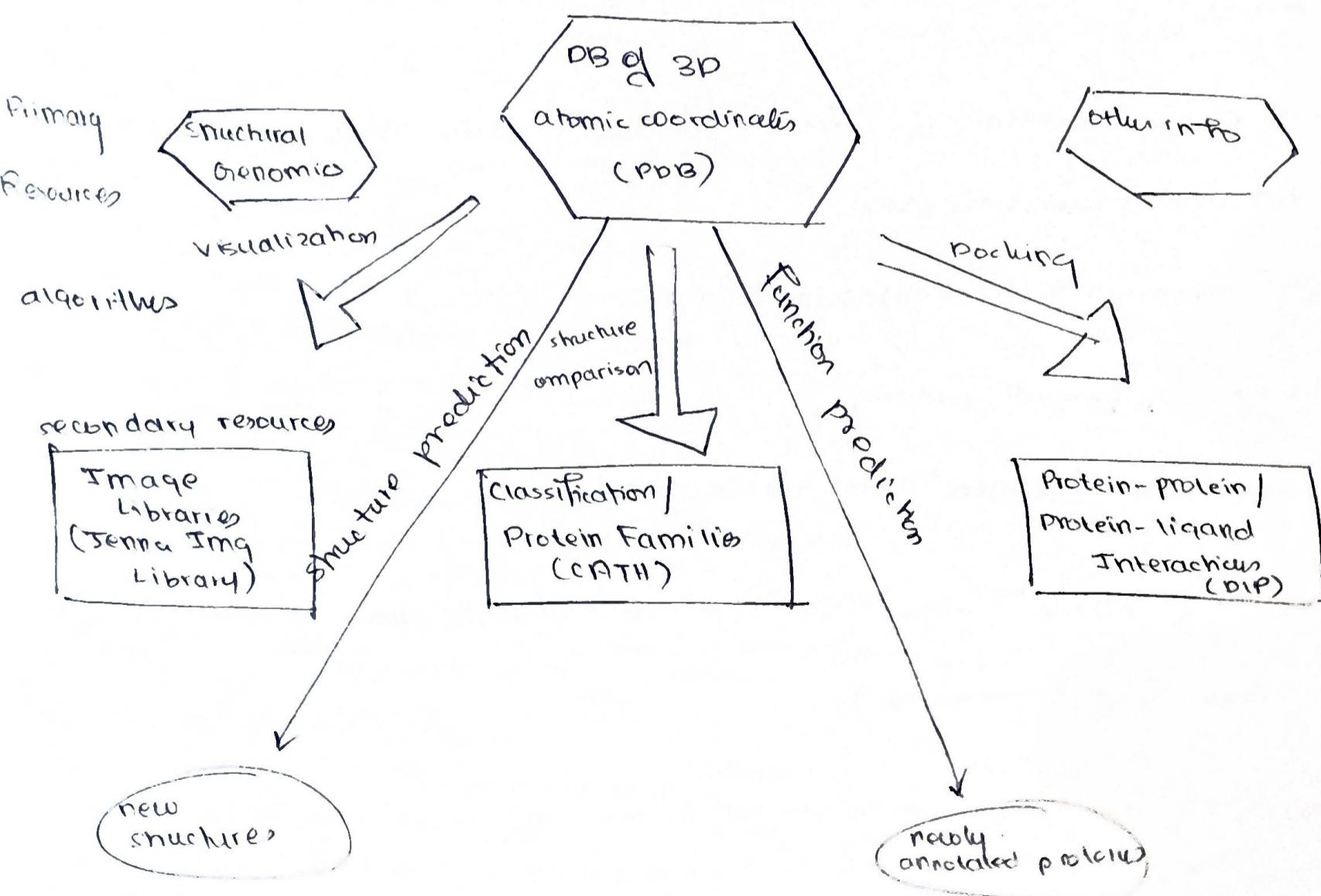
Sequence Analysis

Structure Analysis

Expression Analysis

→ grouping, profiling of genes

* Structural Bioinformatics



Structural bioinformatics represents the subset that deals directly or indirectly with the structure of macromolecules.

Primary Sources

- contains atomic coordinates unimportant/uninformative to the majority of structural bioinfo
e.g. PDB
- Protein Data Bank - stores 3D biological macromolecular structures
 - contains publicly available 3D structures of proteins, nucleic acids and a variety of other complex biomolecules determined by X-ray crystallography, NMR spectroscopy

Data Format

- consists of a collection of fixed format records that describe the
 - (i) atomic coordinates
 - (ii) chemical & biochemical features
 - (iii) experimental details
 - (iv) structural features such as H bonds
- dictionary based format is the macromolecular crystallographic information file (mmcif)
- underlying data organization is a set of relational tables
- mmcif dictionary is an ontology that describes macromolecular structure & the various experiments used to derive it.

Data Processing and Quality Control

Data deposition

Data validation - assess quality of deposited atomic models (structure validation) and assess how well they fit experimental results

Data annotation - process & adding info. to the entry

PDB uses accepted community standards to validate structures

Visualization - molecular graphics program

MCE

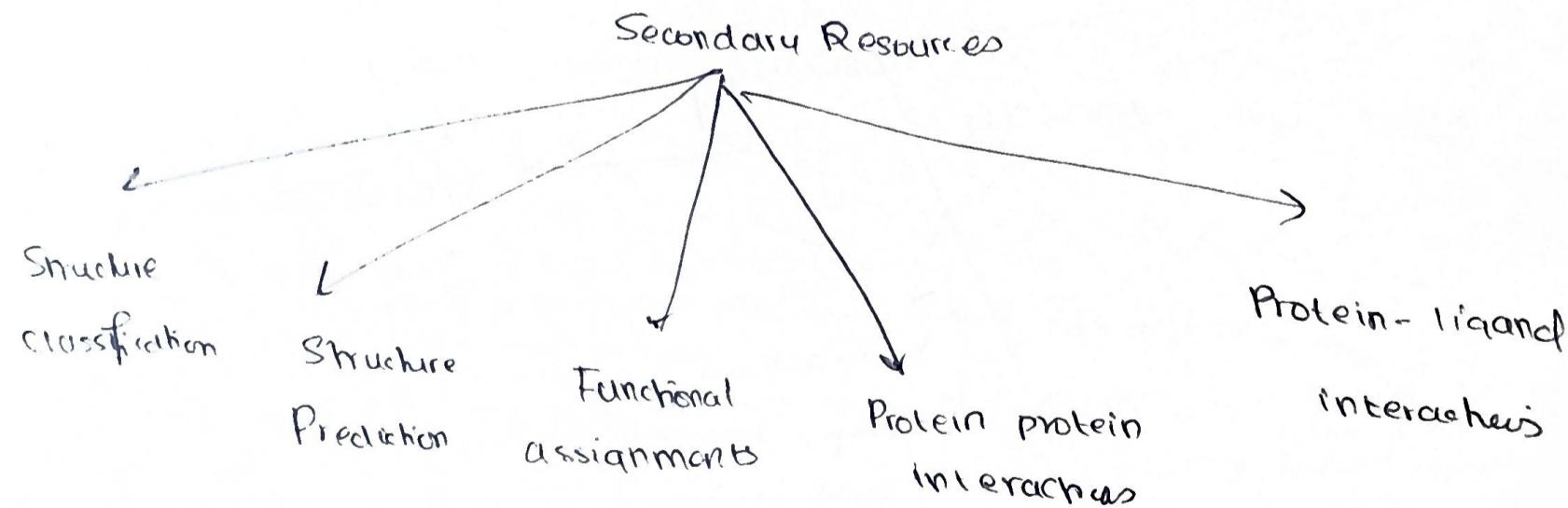
Mol Script

PyMOL

Secondary Structures

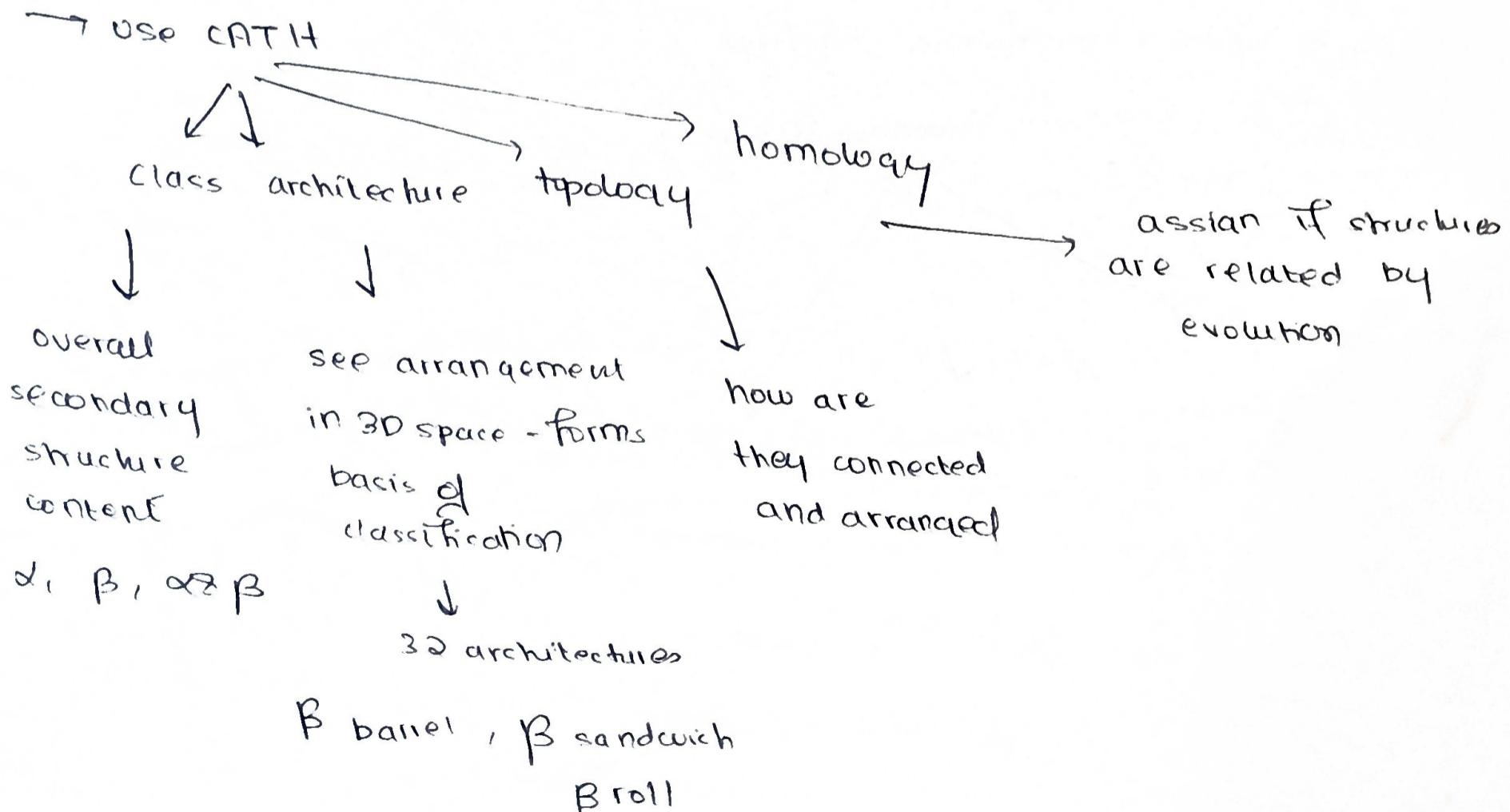
Resources

→ They are value-added structural databases - from data reducing using algorithms or human expertise

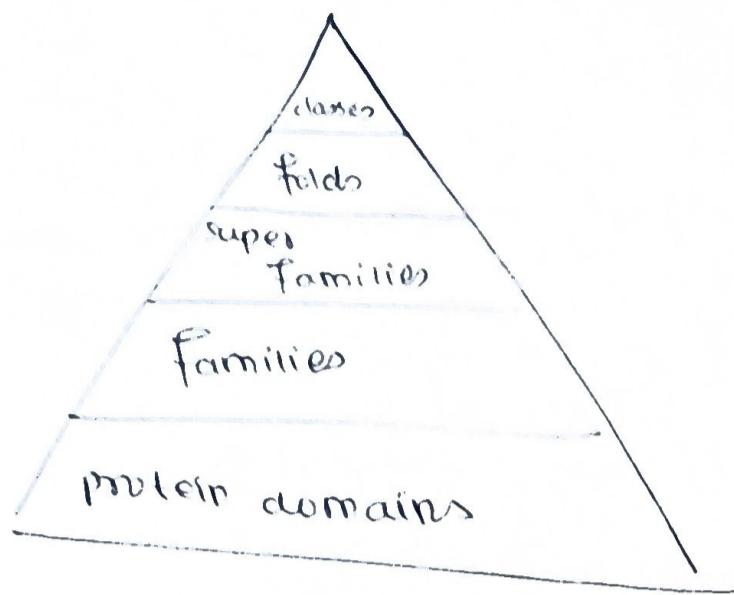


A. Structure Classification

- process of grouping proteins together by the level of 3D sequence similarity
- use structural comparators
- Divided into 3 major steps:
 - (I) representation of 2 structures in a coordinate-independent space
 - (II) comparison & optimisation
 - (III) measure statistical significance of the alignment
- Handle tradeoff between geometric alignment & its biological sequence.



SCOP - Structural classification of proteins



B. Structural Predictions

Methods are

1. Homology modeling,
2. Fold recognition
3. Ab initio

Experiments used are:

Critical Assessment of Structure Prediction (CASP)

"Fully Automated Structure Prediction
(CAFASp)

C. Functional assignment

- what exactly a protein does, when & where in the cell is it active, what are its interacting patterns.
- assign function to the protein simply from its seq. or structure
- proteins w/ similar sequences have similar functions
 - >40% → same structure, similar fn
 - <30 inconclusive
 - 20-30 twilight
 - <20 = midnight

D. Protein Protein Interactions

- network of interacting components
- infer using phylogenetic profiles, gene fusion, and gene neighborhoods
- If 2 proteins are present or both missing \Rightarrow likely to be involved in functional interactions
- Interacting genes are placed closely to each other — correlated using STRING

E. Protein-Ligand Interactions

SMILES, MARVIN

- use ligand-based (analog) target-based (target) for drug discovery
 - use pharmacophore and quantitative structure activity relationships (QSAR)
- 2 methods
clocking
building
- match from existing, DB
 - generate new ligands by connecting atoms of molecular fragments

Diva Design

2RAY, NMR

