

# TOOLS OF BIOINFORMATICS FOR COVID-19 RESEARCH

by

**Dr. Shasank Sekhar Swain**

**ICMR-Young Scientist Fellow**

ICMR-Regional Medical Research Centre,  
Bhubaneswar-23, India

Email: [swain.shasanksekhar86@gmail.com](mailto:swain.shasanksekhar86@gmail.com)





# WHAT IS TOOL..?



**Definition:** A device/ technique/ instrument that solves a problem by providing extra advantage in order to do some useful work.

**Examples:** Hammer, Screwdriver, drilling machine, Xerox machine, Google, etc.,

## MEDICAL/ BIOLOGY



Stethoscope



PCR machine

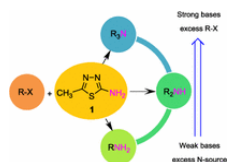


Glucometer



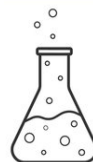
Genetic engineering

## CHEMISTRY/ PHARMACY

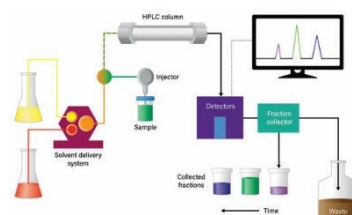


Strong bases  
excess R-X

Weak bases  
excess N-source



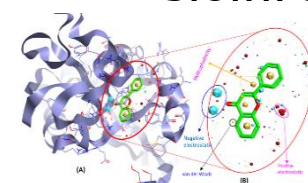
New synthesis/ drug coating protocol



New instruments analysis



## COMP. BIOLOGY/ BIOINFORMATICS

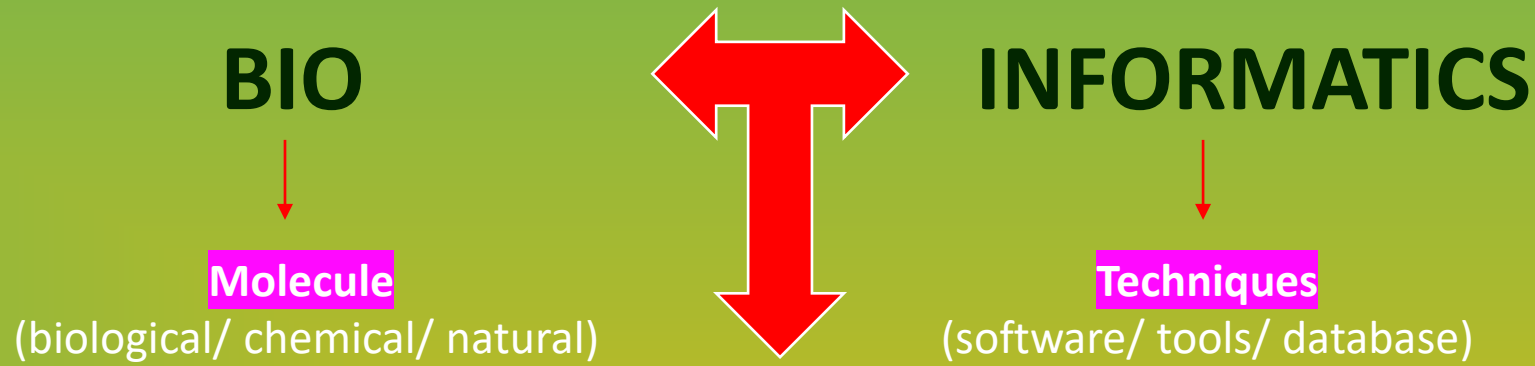


Molecular docking , networking analysis

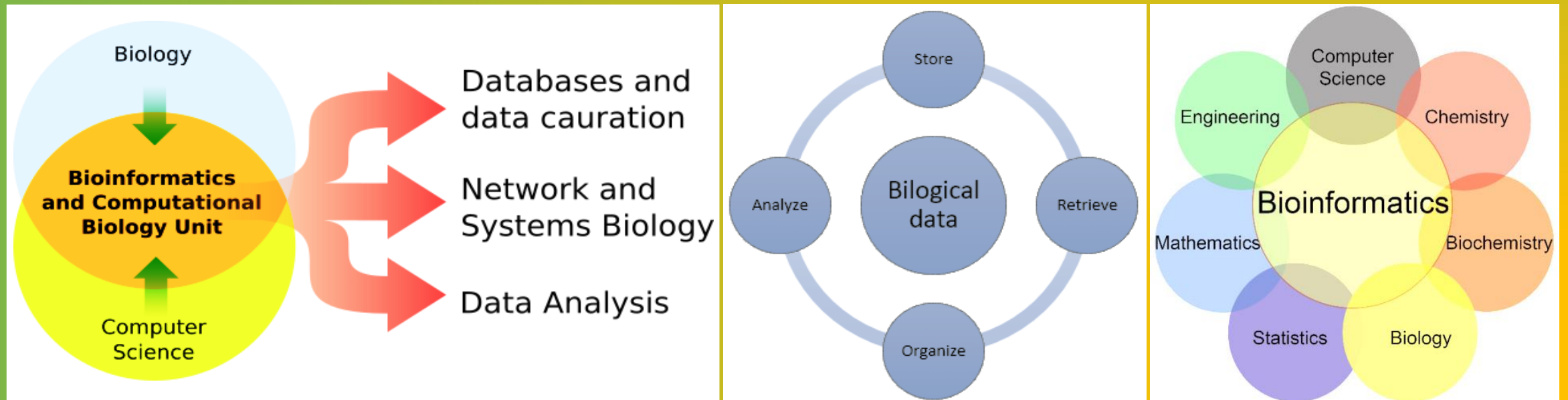


Software, program, server

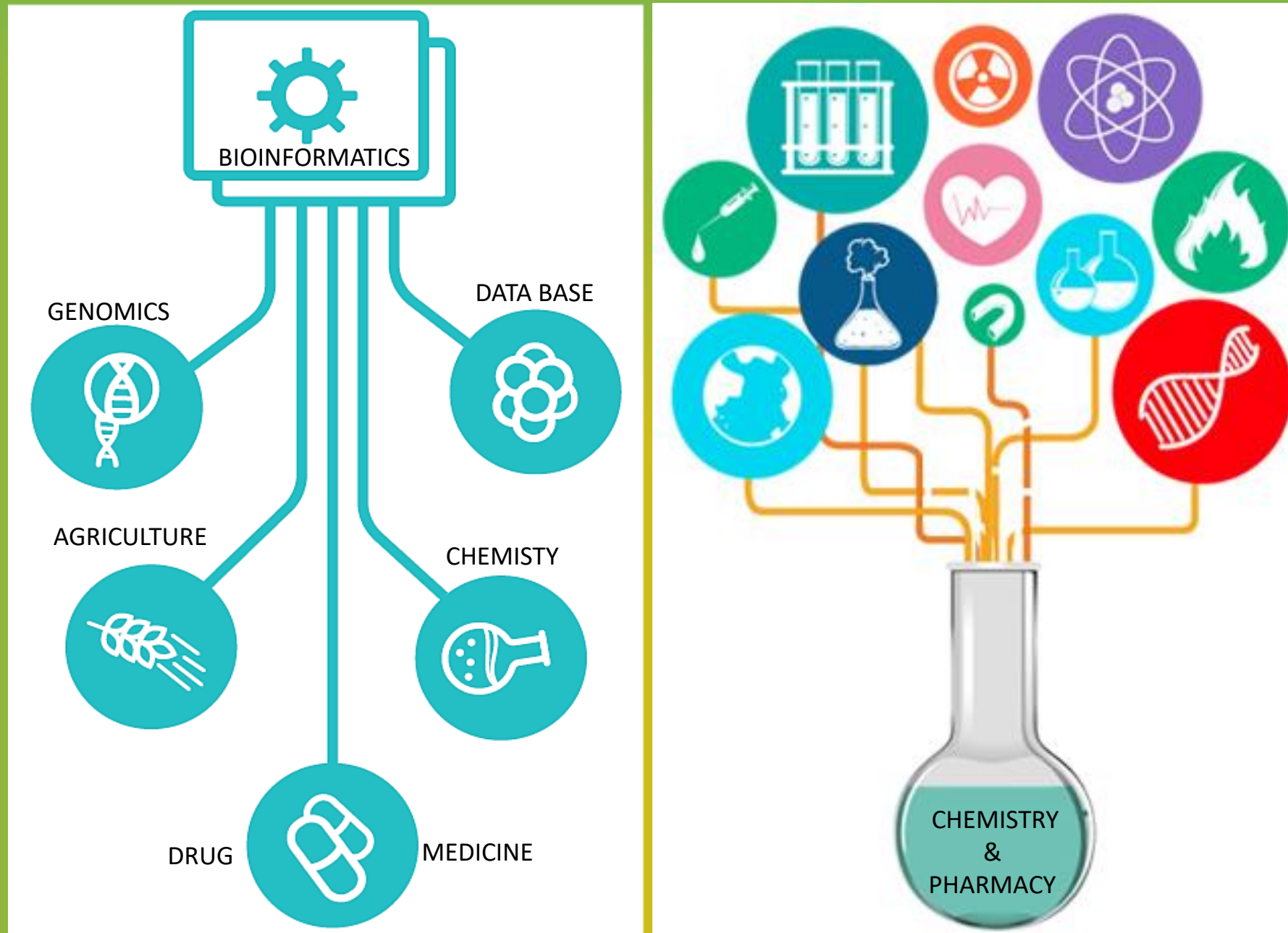
# DEFINITION OF BIOINFORMATICS



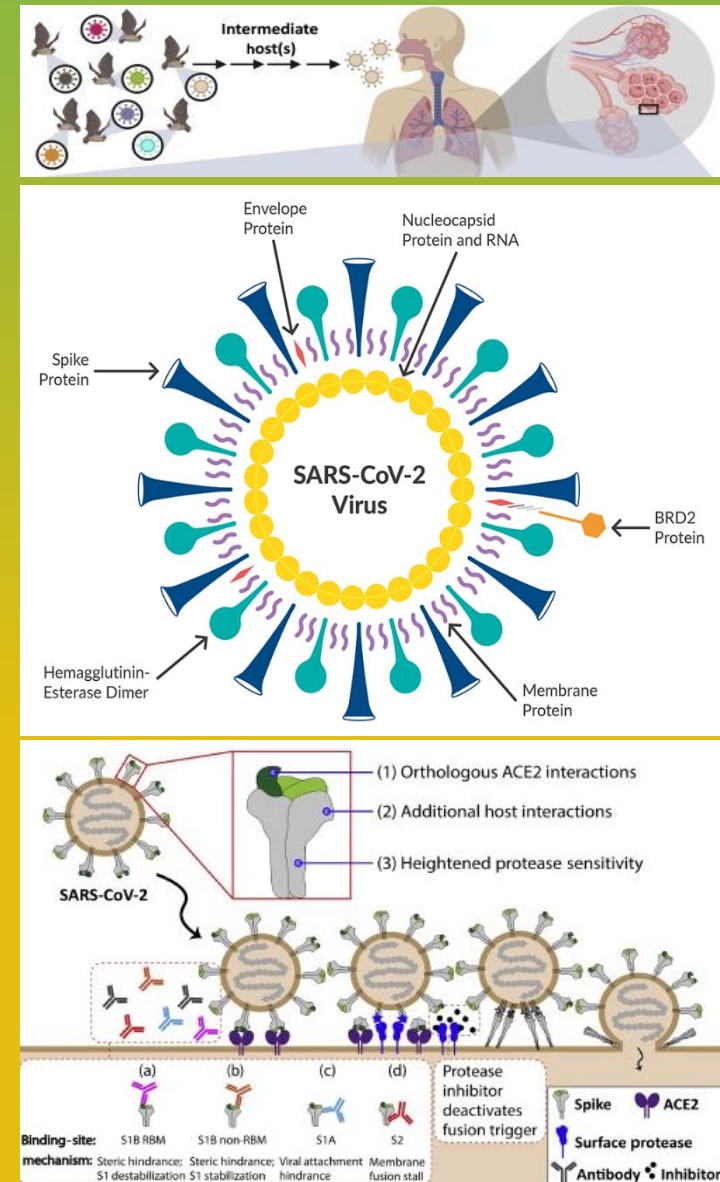
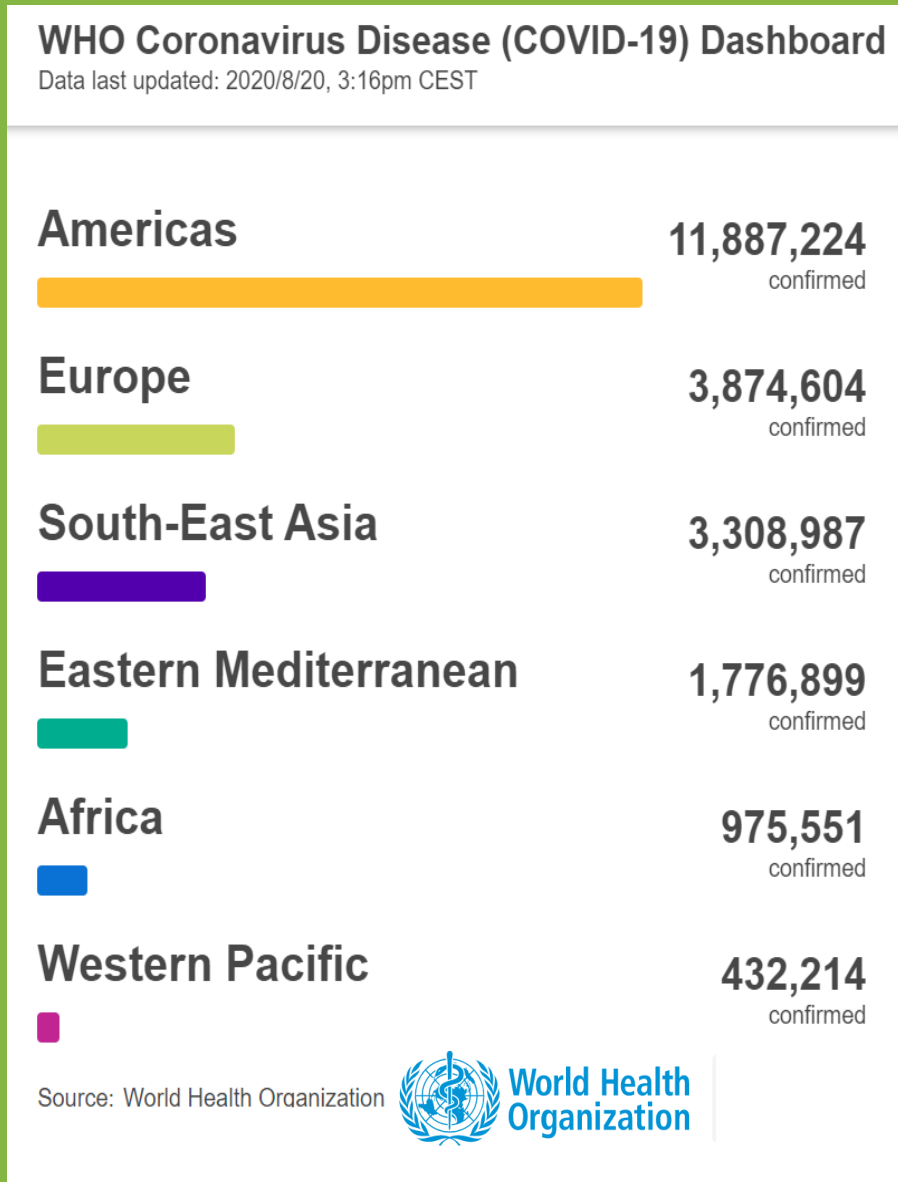
Application of information technology including statistics, mathematics to simplifying the storage, retrieval, analysis of biological dates



# BIOINFORMATICS AREA OF APPLICATION



# Coronavirus disease 2019 (COVID-19)



A novel **coronavirus** (nCoV) is a new strain that has not been previously identified in humans

# TOOLS OF BIOINFORMATICS FOR COVID-19 RESEARCH

## SEQUENCE LEVEL

1. Physicochemical properties
2. Secondary structure analysis
3. Conserve and mutation analysis
4. Phylogenetic tree analysis
5. 3-D structure prediction

## STRUCTURE LEVEL

1. Structural composition
2. Therapeutic agent identification
3. Binding / active site prediction
4. Structural similarity analysis
5. Structural stability with drug

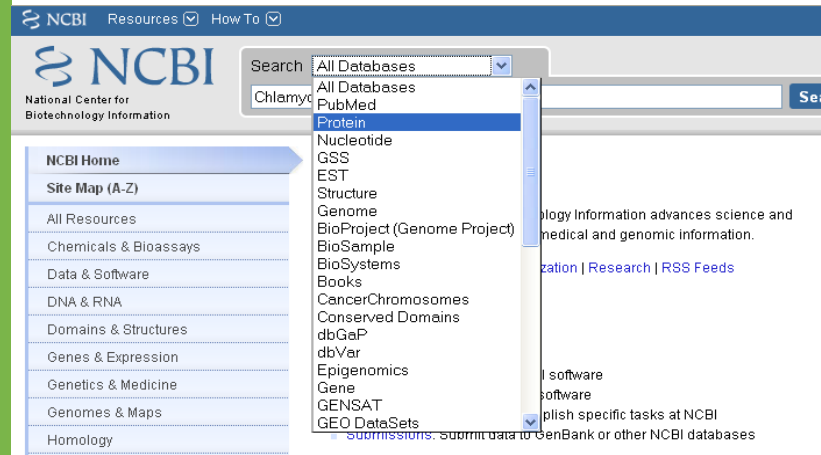
# PHYSICO-CHEMICAL PROPERTY PREDICTION

DNA

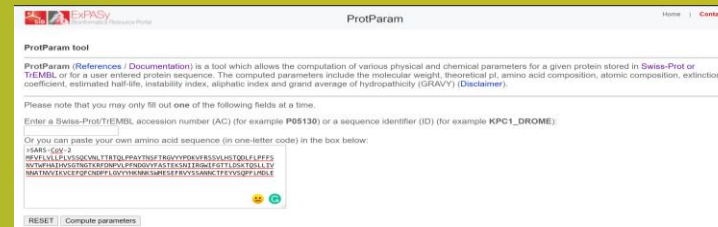
PROTEIN

```
Severe acute respiratory syndrome coronavirus 2 isolate SARS-CoV-2/human/TUN/COV0425/2020, complete genome
GenBank: MT499219.1
>MT499219.1 Severe acute respiratory syndrome coronavirus 2 isolate SARS-CoV-2/human/TUN/COV0425/2020, complete genome
ACTTCGATCTCTGTAGATCTGTTCTCTAAACGAACCTTTAAAATCTGTGTGGCTGTCACCTCGGCTGCA
CTTAGTGCACCTCACGCAGTATAATTAATAACTAATTACTGTCGTTGACAGGACACGAGTAACTCGTCT
CTTCTGCAGGCTGCTTACGGTTTCGTCGGTGTGCAGCCGATCATCAGCACATCTAGGTTTTGTCCGG
GTGACCGAAAGGTAAGATGGAGAGCCTTGTCCCTGGTTTCAACGAGAAAACACACGTCCAACCTCAGTT
CCTGTTTTACAGTTTCGCGACGTGCTCGTAGTGGCTTTGGAGACTCCGTGGAGGAGGTCTTATCAGA
CACGTCAACATCTTAAAGATGGCACTTGTGGCTTAGTAGAAGTTGAAAAAGGCGTTTTGCCTCAACTT
ACAGCCCTATGTGTTTCATCAACCGTTCGGATGCTCGAAGTGCACCTCATGGTCATGTTATGGTTGAGC
GTAGCAGAATCGAAGGCATTCAGTACGGTCGTAGTGGTGGAGACTTGGTGTCTTGTCCCTCATGT
CGGAAATACCACTGGCTTACCGCAAGGTTCTTCTTCGTAAGAACGGTAATAAAGGAGCTGGTGGCCAT
```

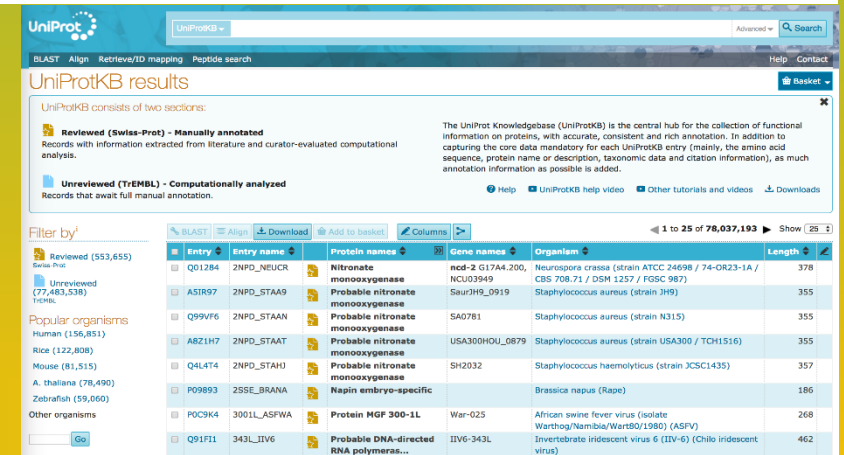
```
>sp|P0DTC2|SPIKE_SARS2 Spike glycoprotein OS=Severe acute respiratory syndrome coronavirus 2 OX=2697049 GN=S PE=1 SV=1
MFVFLVLLPLVSSQCVNLTTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFSNVTFWFH
AIHVSGTNGTKRFDNPNVLPFNDGVYFASTEKSNIRGWI FGTTLD SKTQSL LIVN NATNVV I K VCE
FQCNDFPLGVVYHKNNKSWMESEFRVYSSANNCTFEYVSQPF LMDLE GKQGNFKNLRE FVFKNID
GYFKIY SKHTPINLVRDLPQGFSALEPLVDLP IGINITRFQTL LALHRSYLT PGDSSSGWTAGAAA
YYVGYLQPRFTLLKYNENGTITDAVDCALDPLSETKCTLKSFTVEKGIYQTSNFRVQPTESIVRFP
NITNLCPFGEVFNATRFASVYAWNRRKRSNCVADY SVLYNSASFSTFKCYGVSP TKLNDLCFTNVY
ADSFVIRGDEVQR IAPGQTGKIADYNYKLPDDFTGCVI AWNSNLD SKVGVN YNYLYR LFRKSNLK
PFERDI STEIYQAGSTPCNGVEGFNCYFP LQS YGFQPTNGVGYQFVRVVLS FELLHAPATVCGPK
KSTNLVKNKCVNFNGLTGTGVLTE SNKKFLPFQQFGRDIADTTPDAVRDPQ TLEILDITPCSFGG
VSVITPGTNTSNQVAVLYQDVNCTEVPVAIHADQLTPTWRVYSTGNSV FQTRAGCLIGAEHVNNSY
NGVEGFTE SNKKFL
```



<https://www.ncbi.nlm.nih.gov/>



<https://web.expasy.org/protparam/>



<https://www.uniprot.org/>



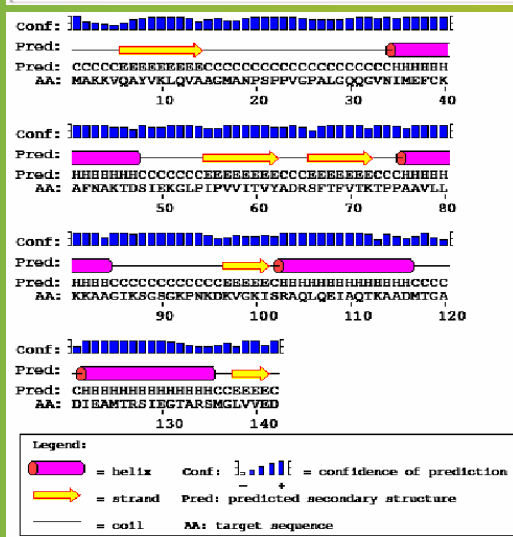
- Molecular weight
- Total amino acids composition
- Total number of atoms
- Negatively charged residues
- Ext. coefficient
- Positively charged residues
- Grand average of hydropathicity
- Aliphatic index
- Protein as stable profile
- Estimated half-life
- Molecular formula
- Theoretical isoelectric point

# SECONDARY STRUCTURE PREDICTION

- Accurate prediction of the exact elements of protein 3D structure is essential for any research targeting a protein.
- Predicting the formation of protein structures such as alpha helices and beta strands, while for nucleic acids, it means predicting the formation of nucleic acid structures like helices and stem-loop structures.

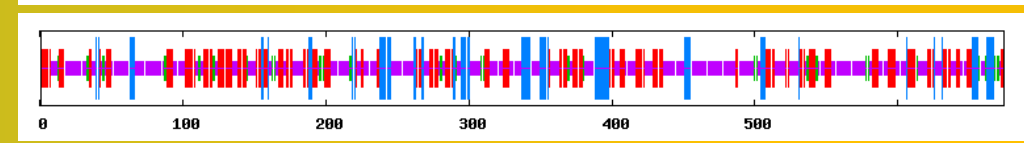
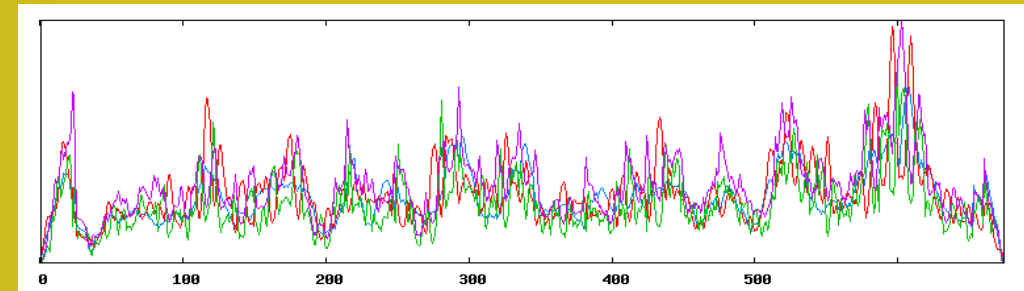
## PSIPRED 4.0 (Predict Secondary Structure)

<http://bioinf.cs.ucl.ac.uk/psipred/>



## SOPMA secondary structure prediction

[https://npsa-prabi.ibcp.fr/cgi-bin/npsa\\_automat.pl?page=/NPSA/npsa\\_sopma.html](https://npsa-prabi.ibcp.fr/cgi-bin/npsa_automat.pl?page=/NPSA/npsa_sopma.html)



Alpha helix (Hh) : 84 is 12.46%  $3_{10}$  helix (Gg) : 0 is 0.00% Pi helix (Ii) : 0 is 0.00% Beta bridge (Bb) : 0 is 0.00% Extended strand (Ee) : 196 is 29.08% Beta turn (Tt) : 41 is 6.08% Bend region (Ss) : 0 is 0.00% Random coil (Cc) : 353 is 52.37% Ambiguous states (?) : 0 is 0.00% Other states : 0 is 0.00%

[https://List\\_of\\_protein\\_secondary\\_structure\\_prediction\\_programs](https://List_of_protein_secondary_structure_prediction_programs)



# SEQUENCE ANALYSIS (MUTATION OR CONSERVED)

- A gene mutation is a permanent alteration in the DNA, protein sequence that makes up a gene, such that the sequence differs from what is found in most people/ region; In simply, sometimes our DNA sequence gets altered; this is called a mutation.

EMBL-EBI Services Research Training Industry About us

## Clustal Omega

Tools > Multiple Sequence Alignment > Clustal Omega

EMBL-EBI to be HTTPS by default from 1st October

On the 1st October the majority of services hosted on www.ebi.ac.uk will be served over HTTPS by default. Services that are becoming HTTPS by default will automatically redirect users accessing the site on insecure HTTP URLs to secure HTTPS URLs. Users of EMBL-EBI services may wish to update links, bookmarks or API clients to use the HTTPS URLs.

### Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between three or more sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

Important note: This tool can align up to 4000 sequences or a maximum file size of 4 MB.

STEP 1 - Enter your input sequences

Enter or paste a set of PROTEIN sequences in any supported format:

<https://www.ebi.ac.uk/Tools/msa/clustalo/>

aa POSITIO N	Kolkata
271	Q
408	R
614	G
723	I (1) (4T)
930	A
1124	V (2) (3G)

aa POSITIO N	STATES			
	Kolkata	Gujarat	Kerala 29	Kerala 166
271	Q	R	Q	Q
408	R	R	I	R
614	G	G	D	D
930	A	A	A	V
1124	V	G	G	G

Figure 2: Mutation analysis of isolates from Kolkata, Gujarat and Kerala

- Multiple sequence alignment of Spike protein sequence of Kolkata isolate with sequences obtained from other parts of India. Sites of mutation are showed in Red
- Tabulation of amino acid mutations among isolates from Kolkata. Mutations are shown in red. Number/s in parenthesis show number of isolates that showed the amino acid type.
- Tabulation of amino acid present at the points of mutation for isolates from different parts of India.

Q271R

R408I

D614G

A930V

G1124V

1 SARS-CoV-2/29/human/2020/IND 100.0% 100.0%

2 hCoV-19/India/S2/2020|EPI\_ISL\_430468 100.0% 99.7%

3 QJC19491.1 100.0% 99.7%

4 SARS-CoV-2/166/human/2020/IND 100.0% 99.8%

consensus/100%

consensus/90%

consensus/80%

consensus/70%

1 10 20 30 40 50 60 70 80 90 100 110 120 130

2DUC SGFRKHAFFPSGKVEGCHYQVTCGTTLNGLALDDVYCPRHVICTAEDNLNPNYEDLLIRKSNHFLVQAGNVQLRVIGHSHQNCILLRLKYDTANPKTPKYKFYRIQPGQTFSVLACYNGSPSGVYQCAN

6Y2E SGFRKHAFFPSGKVEGCHYQVTCGTTLNGLALDDVYCPRHVICTAEDNLNPNYEDLLIRKSNHFLVQAGNVQLRVIGHSHQNCVLLRLKYDTANPKTPKYKFYRIQPGQTFSVLACYNGSPSGVYQCAN

Consensus SGFRKHAFFPSGKVEGCHYQVTCGTTLNGLALDDVYCPRHVICTAEDNLNPNYEDLLIRKSNHFLVQAGNVQLRVIGHSHQNCILLRLKYDTANPKTPKYKFYRIQPGQTFSVLACYNGSPSGVYQCAN

131 140 150 160 170 180 190 200 210 220 230 240 250 260

2DUC RPNHTIKGSFLNGSCGSYGFNIIDYDCVSYFCYHMHMELPTGVHAGTDLEGNFYGPYVDRQTAQAAGTDTTITLVNLAWL YARYINGDRWFLNRF T T L NDFNLVANKYNYEPLTQDHVDILGPLSAQTGIA

6Y2E RPNHTIKGSFLNGSCGSYGFNIIDYDCVSYFCYHMHMELPTGVHAGTDLEGNFYGPYVDRQTAQAAGTDTTITLVNLAWL YARYINGDRWFLNRF T T L NDFNLVANKYNYEPLTQDHVDILGPLSAQTGIA

Consensus RPNHTIKGSFLNGSCGSYGFNIIDYDCVSYFCYHMHMELPTGVHAGTDLEnFYGPYVDRQTAQAAGTDTTITLVNLAWL YARYINGDRWFLNRF T T L NDFNLVANKYNYEPLTQDHVDILGPLSAQTGIA

261 270 280 290 300 306

2DUC VLDMCAALKELLQNGMNGRTILGSTILEDEFTPFDVYRQC SGVTFQ

6Y2E VLDMCASLKELLQNGMNGRTILGSAALLEDEFTPFDVYRQC SGVTFQ

Consensus VLDMCAALKELLQNGMNGRTILGsaILEDEFTPFDVYRQC SGVTFQ

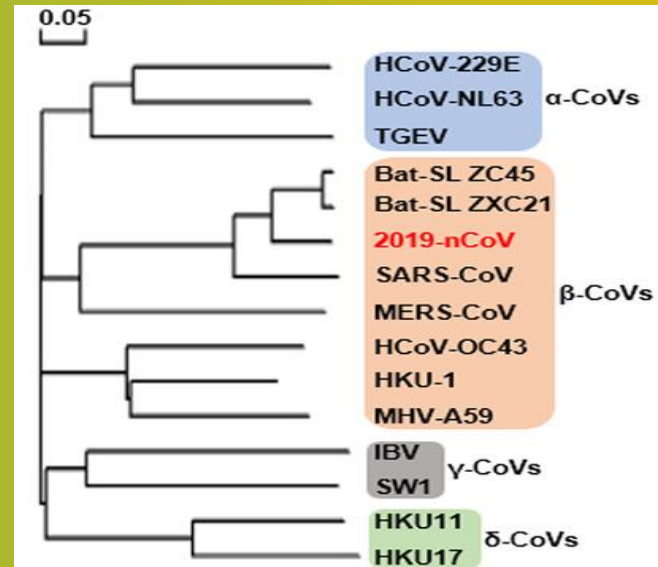
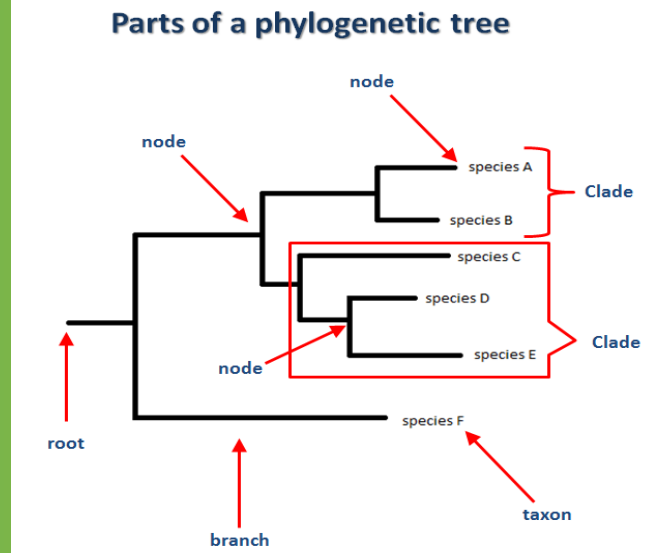
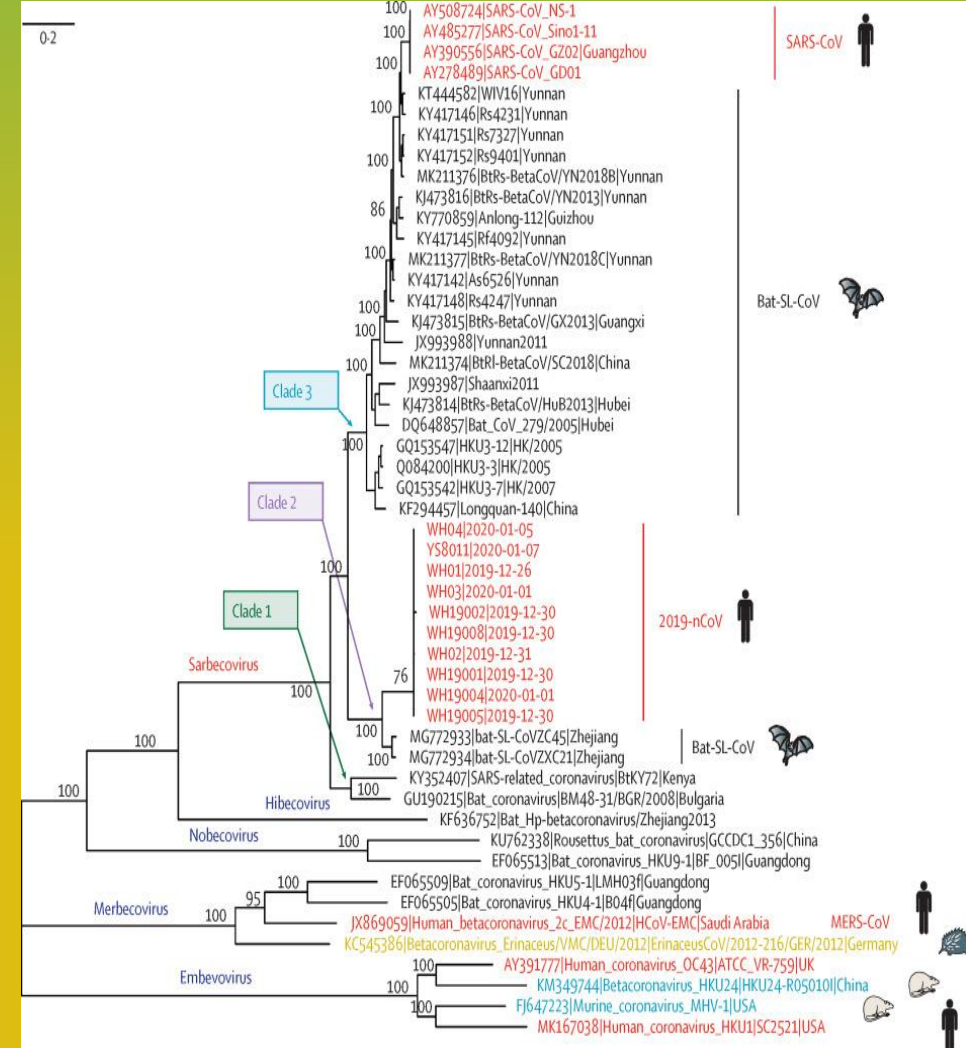
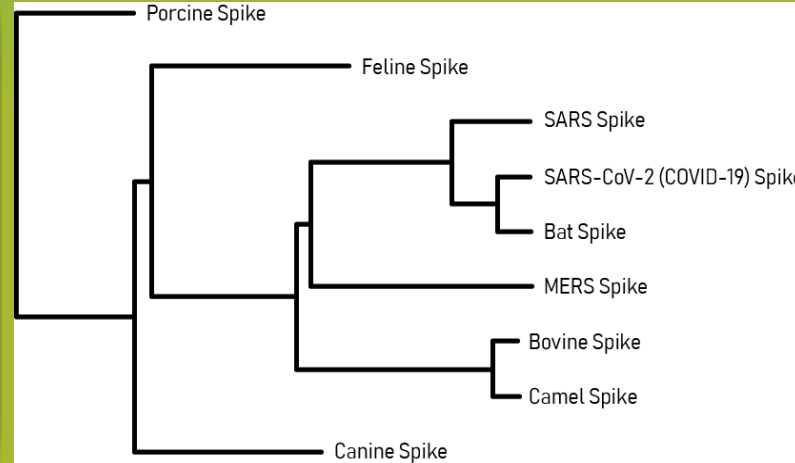
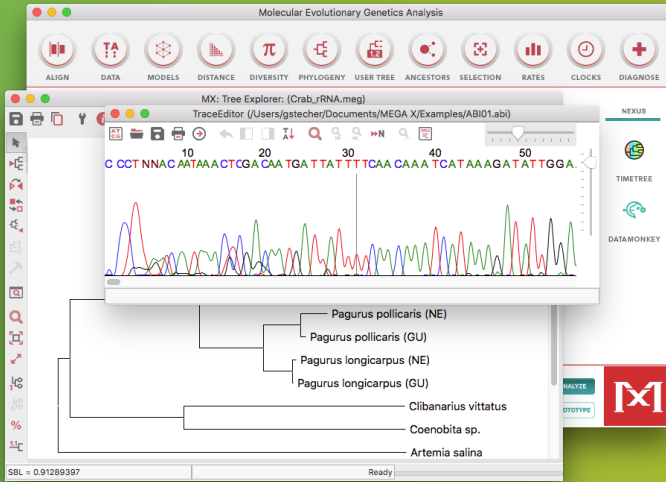
<https://www.ebi.ac.uk/Tools/msa/>

<https://bioinformaticshome.com/tools/msa/msa.html>

# PHYLOGENETIC TREE ANALYSIS/ BUILDING

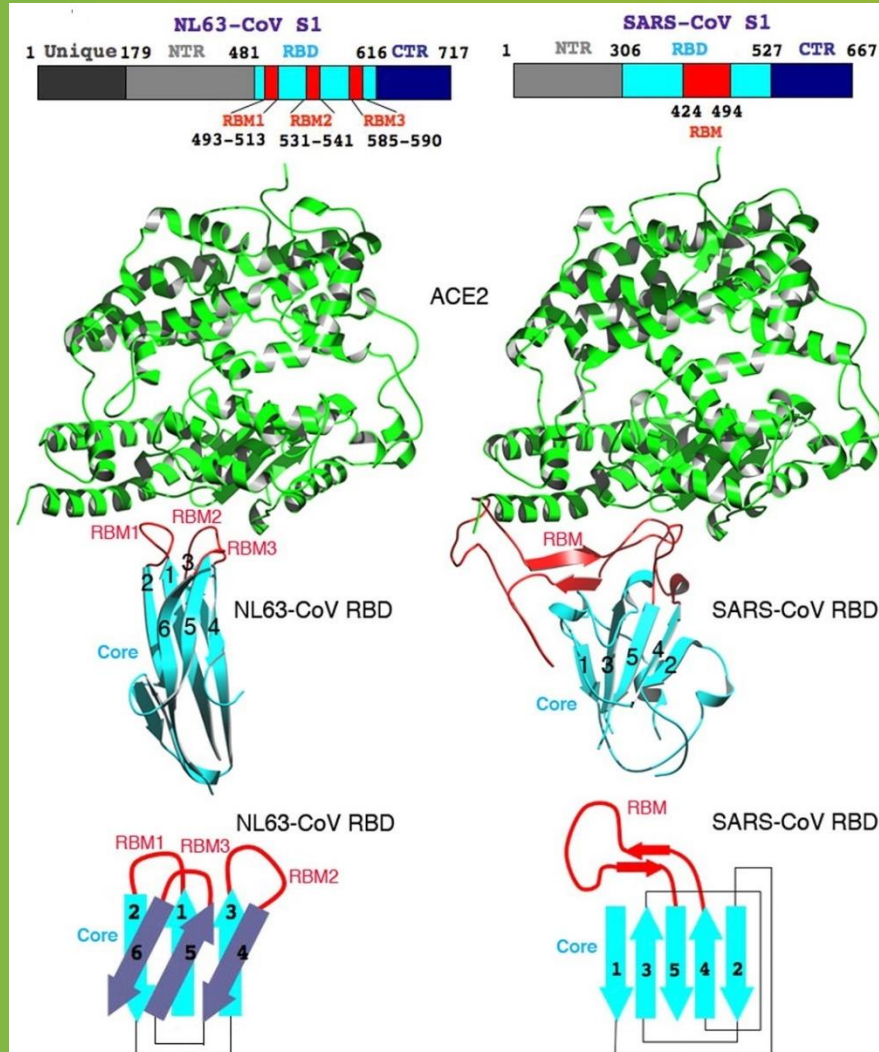
Phylogenetics is the study of evolutionary relationships among biological entities

Most usable software's: MEGA, Dendroscope, FigTree, Phylotree, ggtree



# STRUCTURAL SIMILARITY ANALYSIS

- An image quality metric that assesses the visual impact of a protein structure characteristics/ building blocks



SOFTWARE'S

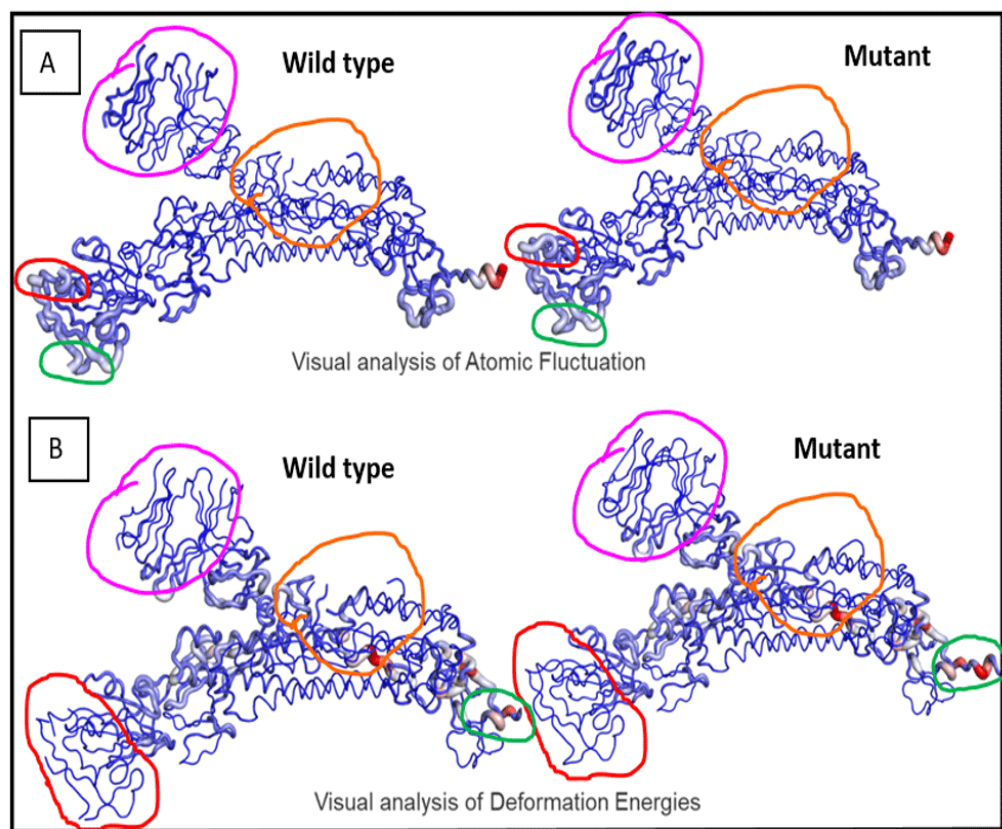


SARS-CoV-2	SGFRKMAFP	SGKVEGCMVQVTCGTTT	LNGLWLD	DDVYVYCP	PRHVICT	SE	MLNP	NYEDLL	LR	60						
SARS-CoV	-GFRKMAFP	SGKVEGCMVQVTCGTTT	LNGLWLD	DDTVYVYCP	PRHVICT	AED	MLNP	NYEDLL	LR	59						
*****																
SARS-CoV-2	KSNHN	FLVQAGNVQLRV	I	GHSMQNCV	LKLVDTAN	P	KPKYK	FVRIQP	GGT	F	SVL	ACYNG	120			
SARS-CoV	KSNHS	FLVQAGNVQLRV	I	GHSMQNCV	LKLVDTSN	P	KPKYK	FVRIQP	GGT	F	SVL	ACYNG	119			
****																
*****																
SARS-CoV-2	SPSGVYQ	CAMRPNFTIK	G	SFLNGSC	GVGFNI	D	YDCV	FCYMH	MELPT	GVH	AGT	DLEGN	180			
SARS-CoV	SPSGVYQ	CAMRPNHTIK	G	SFLNGSC	GVGFNI	D	YDCV	FCYMH	MELPT	GVH	AGT	DLEGN	179			
*****																
*****																
SARS-CoV-2	FYGFVDR	QTAQAAGT	DTTITV	NVLAWL	YAAVING	DRW	FLNR	FTTLL	NDFNL	VAMK	YNYE	240				
SARS-CoV	FYGFVDR	QTAQAAGT	DTTITL	NVLAWL	YAAVING	DRW	FLNR	FTTLL	NDFNL	VAMK	YNYE	239				
*****																
*****																
SARS-CoV-2	PLTQD	HVDILG	PLSAQT	GI	AVLDM	CA	SKELL	QNGM	NGRTIL	GS	ALLE	DEFT	PF	DVVR	QC	300
SARS-CoV	PLTQD	HVDILG	PLSAQT	GI	AVLDM	CA	SKELL	QNGM	NGRTIL	GS	ALLE	DEFT	PF	DVVR	QC	299
*****																
*****																
SARS-CoV-2	SGVTFQ	306														
SARS-CoV	S-----	300														
*****																

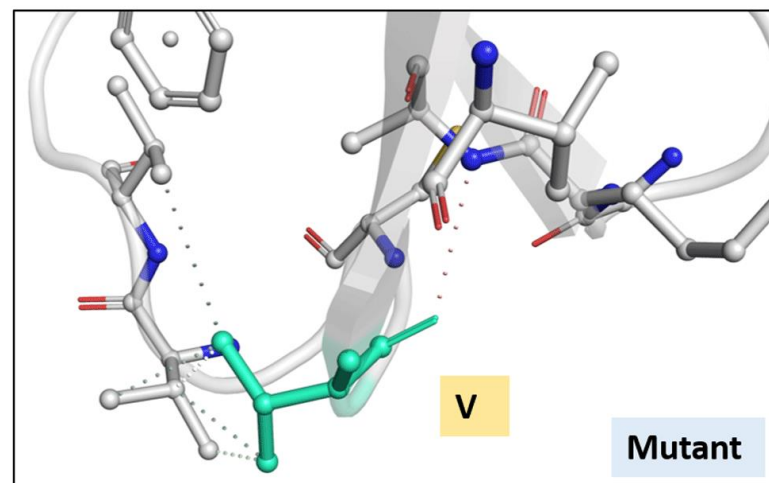
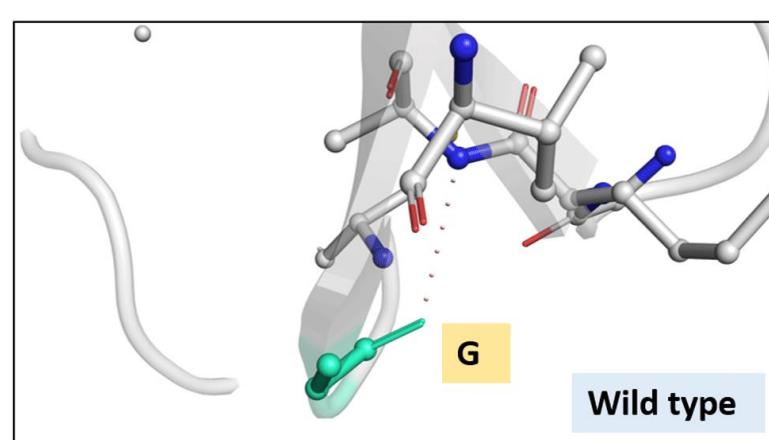
● indication of mutation, (\*) or star for a conservative sequence (:) or colon for a conservative mutation and (.) single dot for a semi-conserved region of amino acids

# STRUCTURAL STABILITY ANALYSIS IN MUTANT PROTEIN

- Proteins are highly dynamic molecules, whose function is intrinsically linked to their molecular motions (analysis carried out by the tool, DynaMut).
- Despite the pivotal role of protein dynamics has led to most structure-based approaches for assessing the impact of mutations on protein structure and function relying upon static structures.



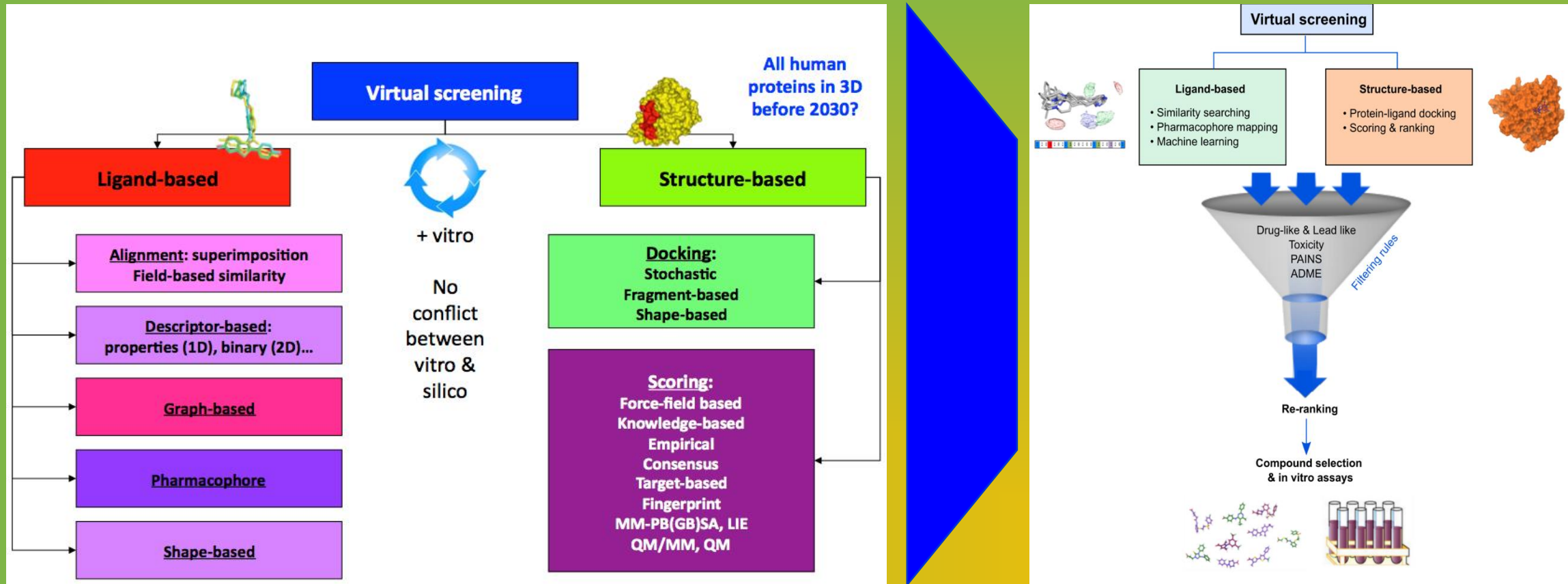
Magnitude of (A) atomic fluctuation and (B) deformation has been shown using thin to thick lines coloured blue (low), white (moderate) and red (high).



Wild-type and mutant residues are coloured in light-green and are also represented as sticks alongside with the surrounding residues which are involved on any type of interactions.

To correlate if changes in secondary structure are also reflected in the dynamics of the protein in its tertiary structure, performed normal mode analyses and studied protein stability and flexibility. Change in vibrational entropy energy ( $\Delta\Delta S_{Vib}^{ENCoM}$ ) between the wild type Wuhan isolate and the West Bengal isolate was  $-4.445 \text{ kcal}\cdot\text{mol}^{-1}\cdot\text{K}^{-1}$ . The  $\Delta\Delta G$  was  $0.905 \text{ kcal/mol}$  and the  $\Delta\Delta G^{ENCoM}$  was  $4.756 \text{ kcal/mol}$ . All these suggested a stabilizing mutation in this type of spike.

# DRUG DISCOVERY USING BIOINFORMATICS TOOLS

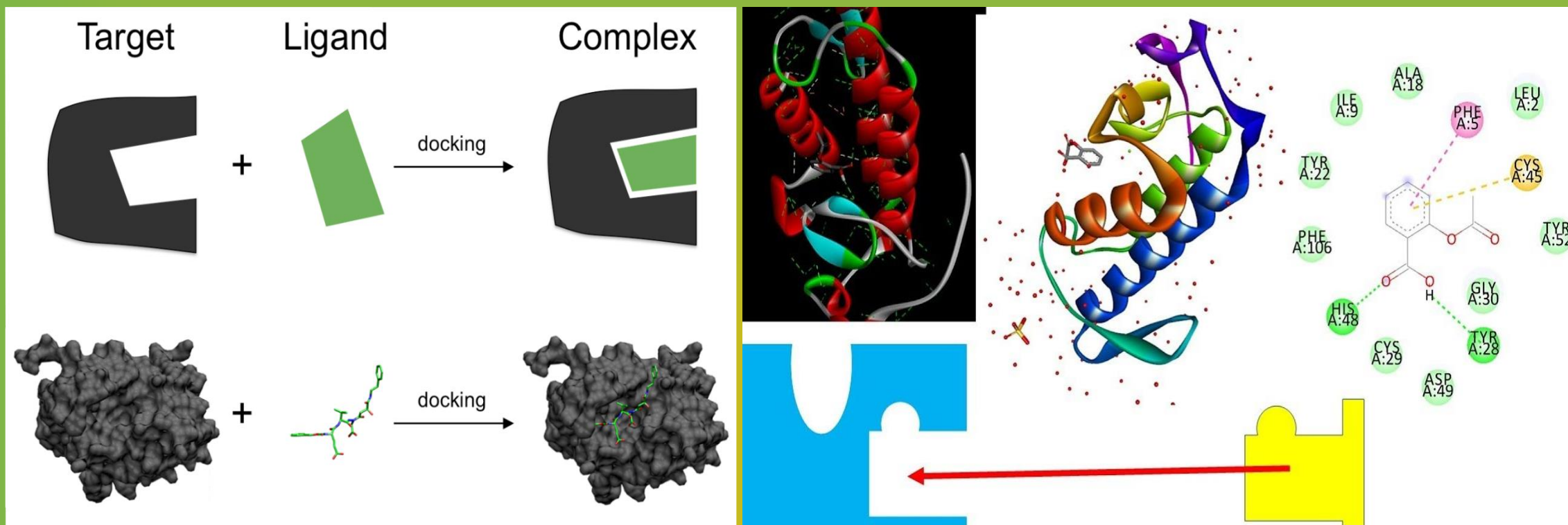


**LIGAND BASED DRUG DESIGN:** It is otherwise known as indirect drug design. It trusts on the awareness of different new ligand molecules that bind with the target protein molecule. (**Known ligand with unknown receptor**).

**STRUCTURE BASED DRUG DESIGN:** It depends on the wisdom of three-dimensional structure of the protein molecule. Practically the structure was initially identified by X-ray crystallography which improves the aptitude to produce new drugs that fight against diseases. (**Known receptor with unknown ligand**).



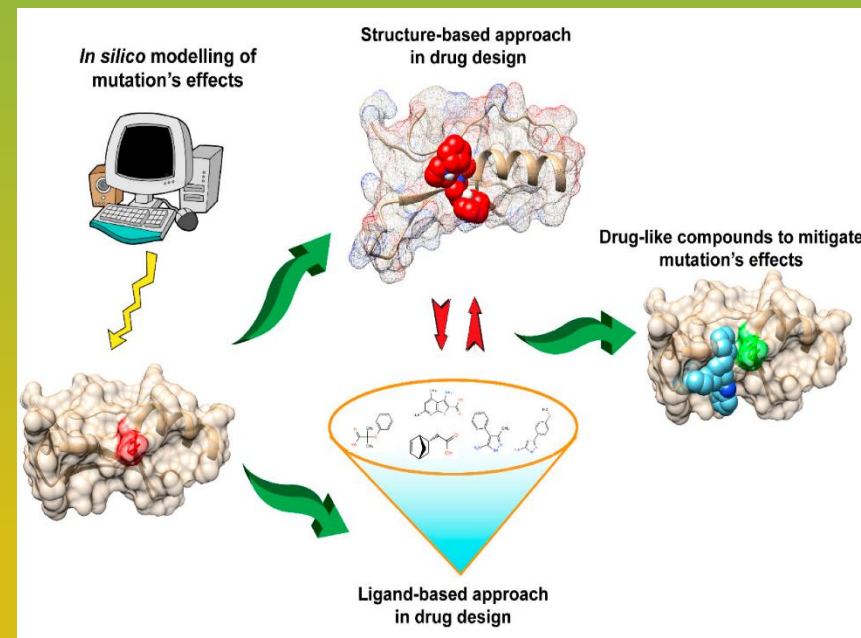
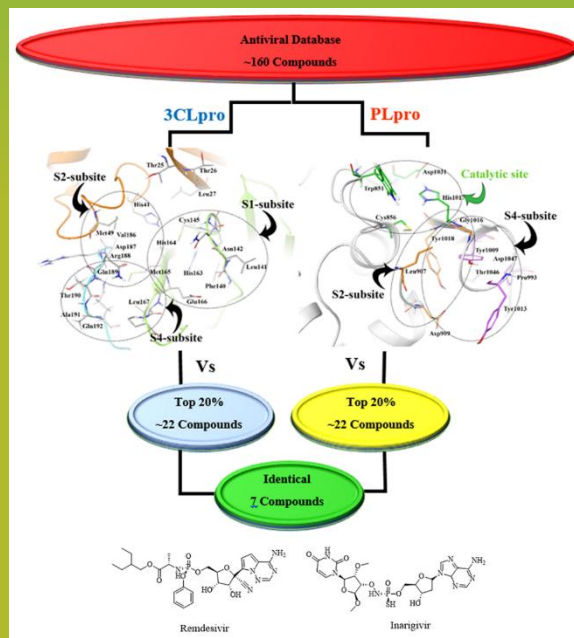
# CONCEPT OF MOLECULAR DOCKING



- ❖ Molecular docking, which predicts interaction patterns based on scoring function between proteins and small molecules as well as proteins and proteins, to evaluate the binding between two molecules is widely used in the field of drug screening and design.
- ❖ It is currently used as a standard computational tool in drug design for lead compound optimisation and in virtual screening studies to find novel biologically active molecules.

# HIGH-THROUGHPUT VIRTUAL SCREENING

Virtual screening is an important part of computer-aided drug design methods. It may be the cheapest way to identify potential lead compounds, and many successful cases have proven successful using this technology



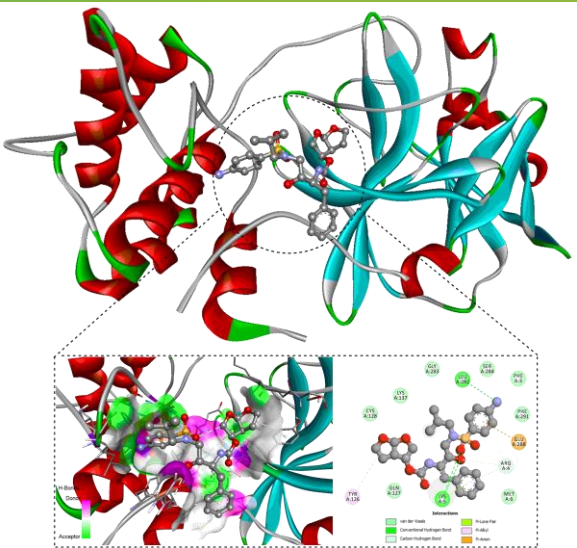
Sections	Type of modeling	License	Reference
GOLD	Protein-ligand	Commercial	Verdonk et al <sup>5</sup>
Gaussian (ONIOM)	QM/MM	Commercial	Gaussian <sup>6</sup>
AUTODOCK	Protein-ligand	Open	Morris et al <sup>7</sup>
GLIDE	Protein-ligand	Commercial	Friesner et al <sup>8</sup>
RosettaDock	Protein-protein	Open	Lyskov and Gray <sup>9</sup>
pyDOCK	Protein-protein	Open	Jimenez-Garcia et al <sup>10</sup>
AquaSol	Solvent effect	Open	Koehl and Delarue <sup>11</sup>

Abbreviations: ONIOM, our own N-layered integrated molecular orbital and molecular mechanics; QM, quantum mechanics; MM, molecular mechanics.

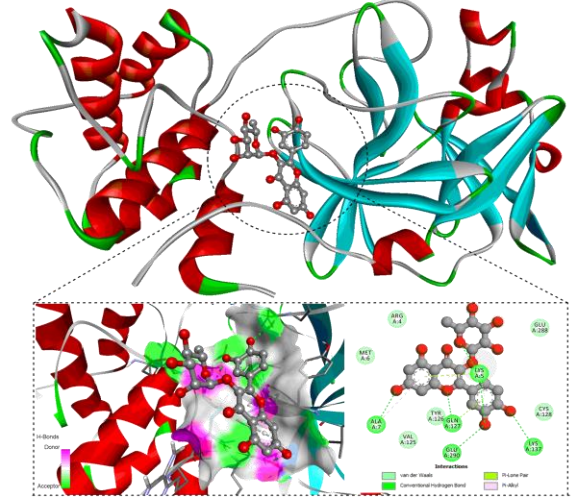
Docking approach	Examples
Matching of descriptors	DOCK, QSDOCK, SLIDE
Incremental construction	FlexX, Hammerhead
Monte Carlo Simulated Annealing	AutoDock, MCDOCK
Monte Carlo Minimization	ICM, QXP
Molecular Dynamics	MDD
Genetic Algorithms	GOLD, AutoDock3



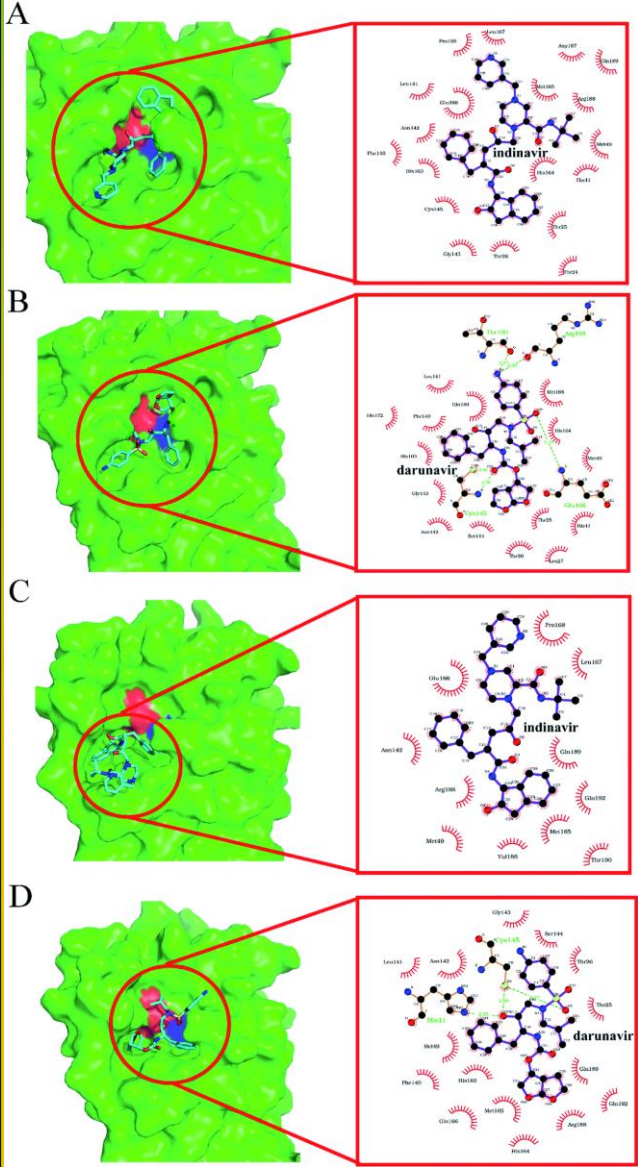
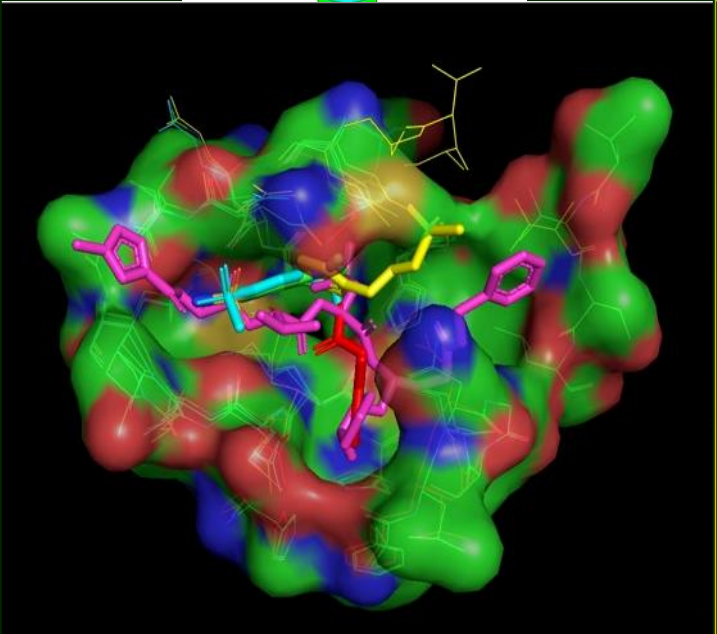
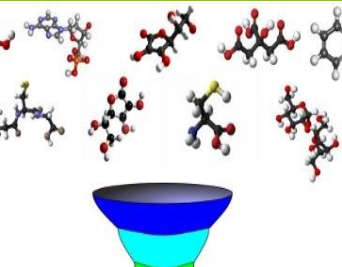
# PROTEIN-LIGAND INTERACTION



SARS-CoV-Mpro (PDB ID: 6Y2E)-Darunavir



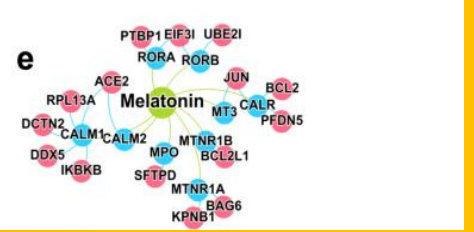
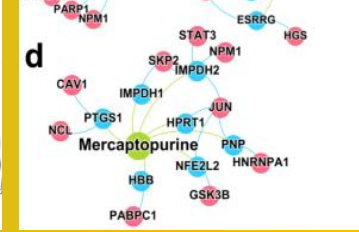
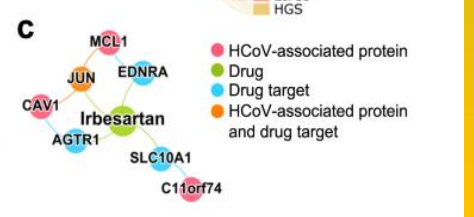
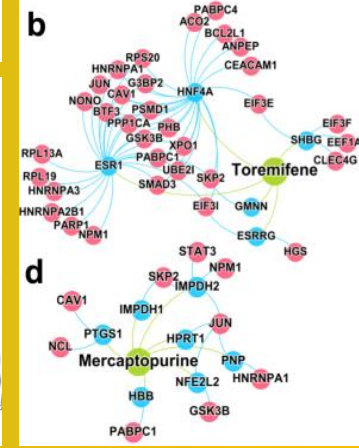
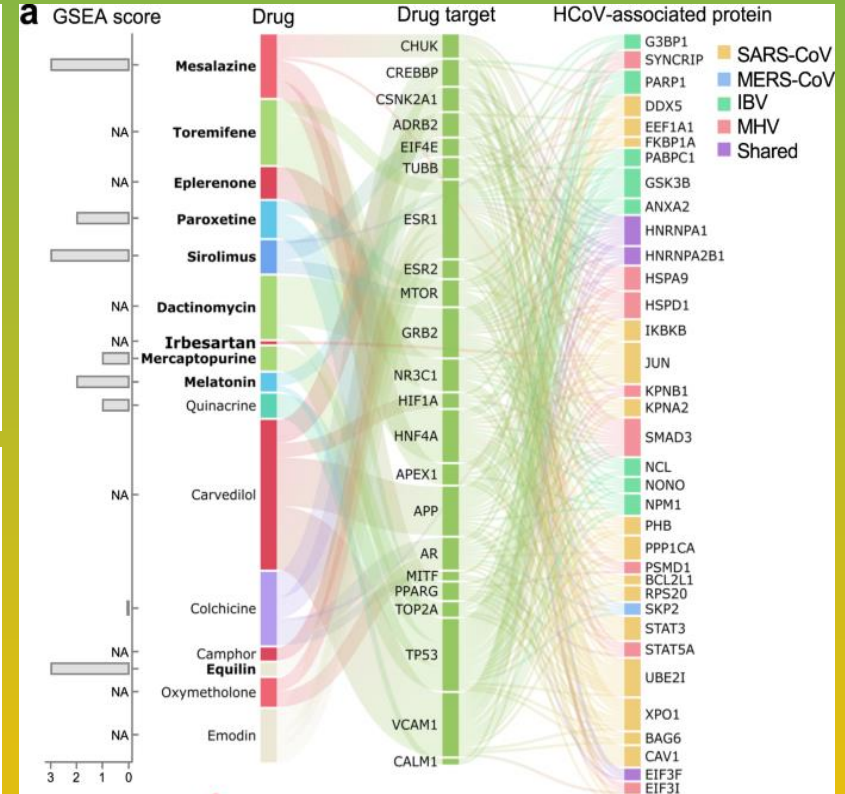
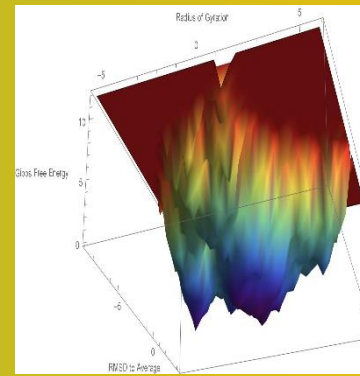
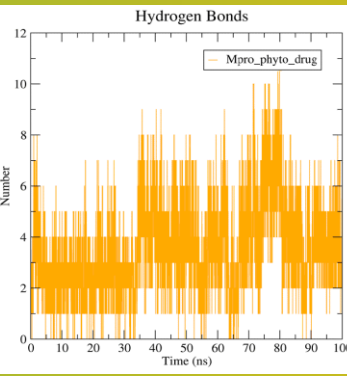
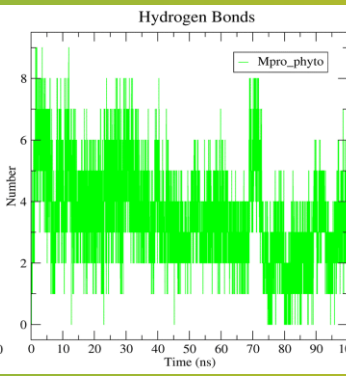
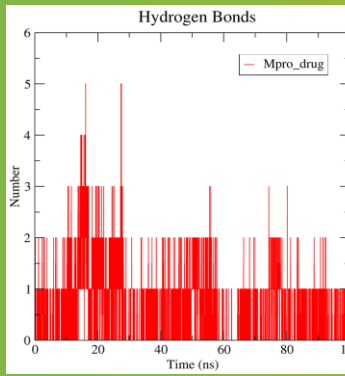
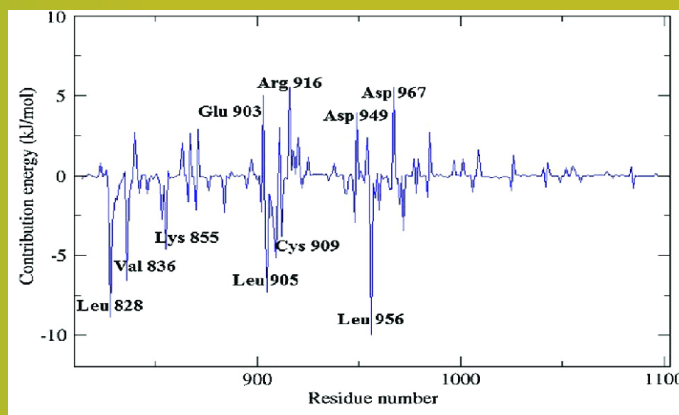
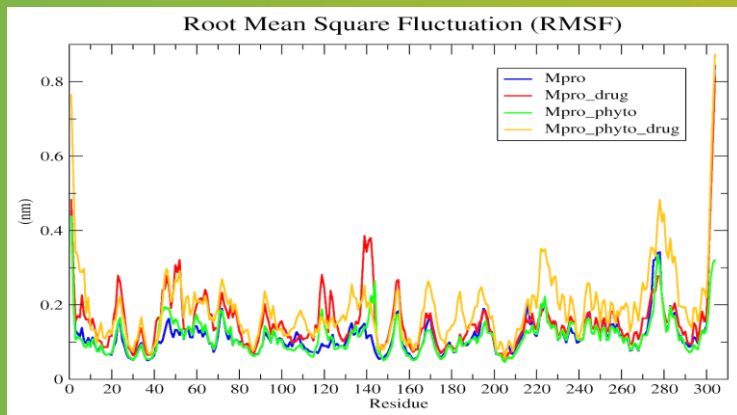
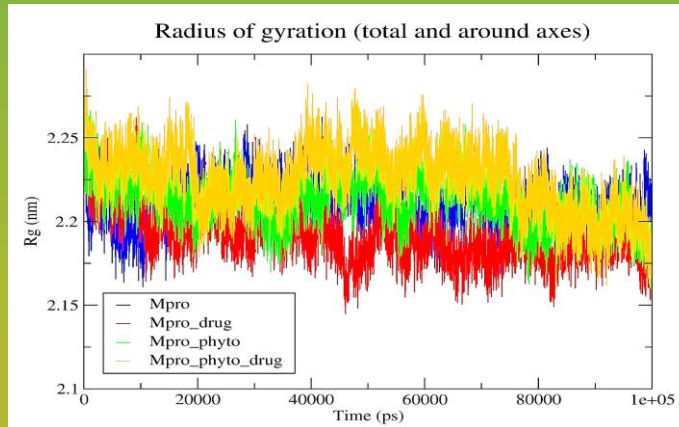
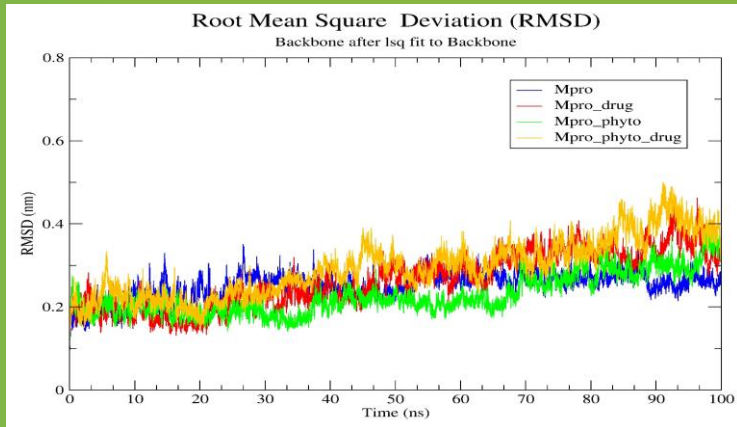
SARS-CoV-Mpro -Quercetin-3-rhamnoside



# ADVANCED BIOINFORMATICS TOOLS IN DRUG DISCOVERY

M D S I M U L A T I O N

G E N E N E T W O R K I N G



## CONCLUDING REMARKS

### Various common advantages of computational method drug design as follows

- To reduce the complexity
  - Time consuming
  - Accurate results
  - Reproducibility
  - Lower cost
- Novel target identification

### Major advantages of computation in the drug design process as follows

- Virtual screening and de novo drug design
- In silico* pharmacokinetic properties prediction
- Improved methods for to determine protein-ligand binding.

*Thank you for your patience*



*STAY HOME ..STAY HEALTHY FROM*

