# Introduction to AlphaFold for 3D Protein Structure Prediction on Grace

Introduction to AF2 - 05/25/2022
Devon J. Boland
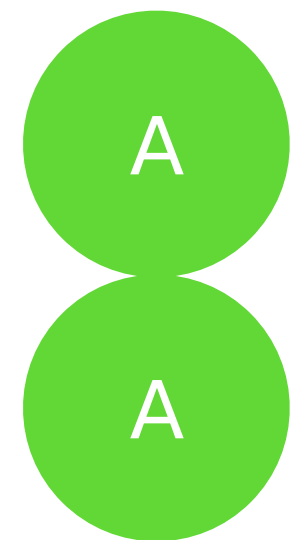Devarenne Laboratory

TEXAS A&M
AGRILIFE
RESEARCH

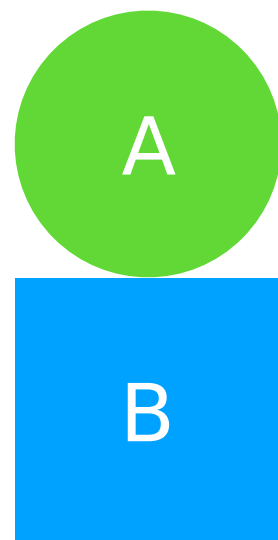# Basics of 3D Protein Structure

## Four Tiers of Protein Structure

- 1° - Sequence of AA's (polypeptides)
- 2° - Interactions of the carbon backbone of 1°
- 3° - Folding of 2° onto itself
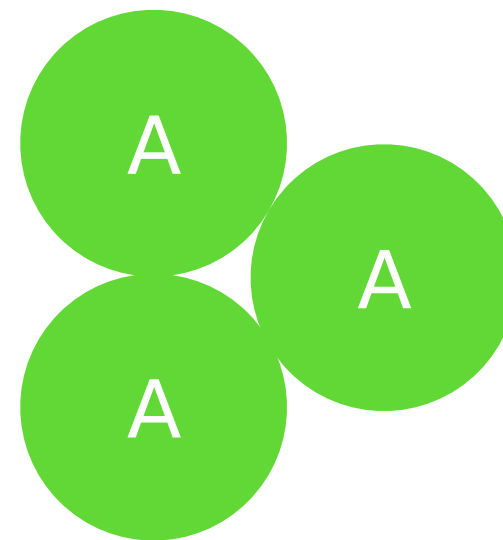- 4° - Multiple 3° units (Monomers) assembling together
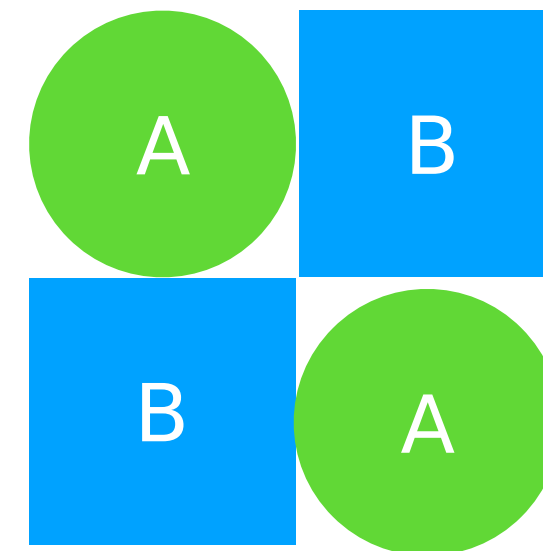
## Different Types of Quaternary Structure
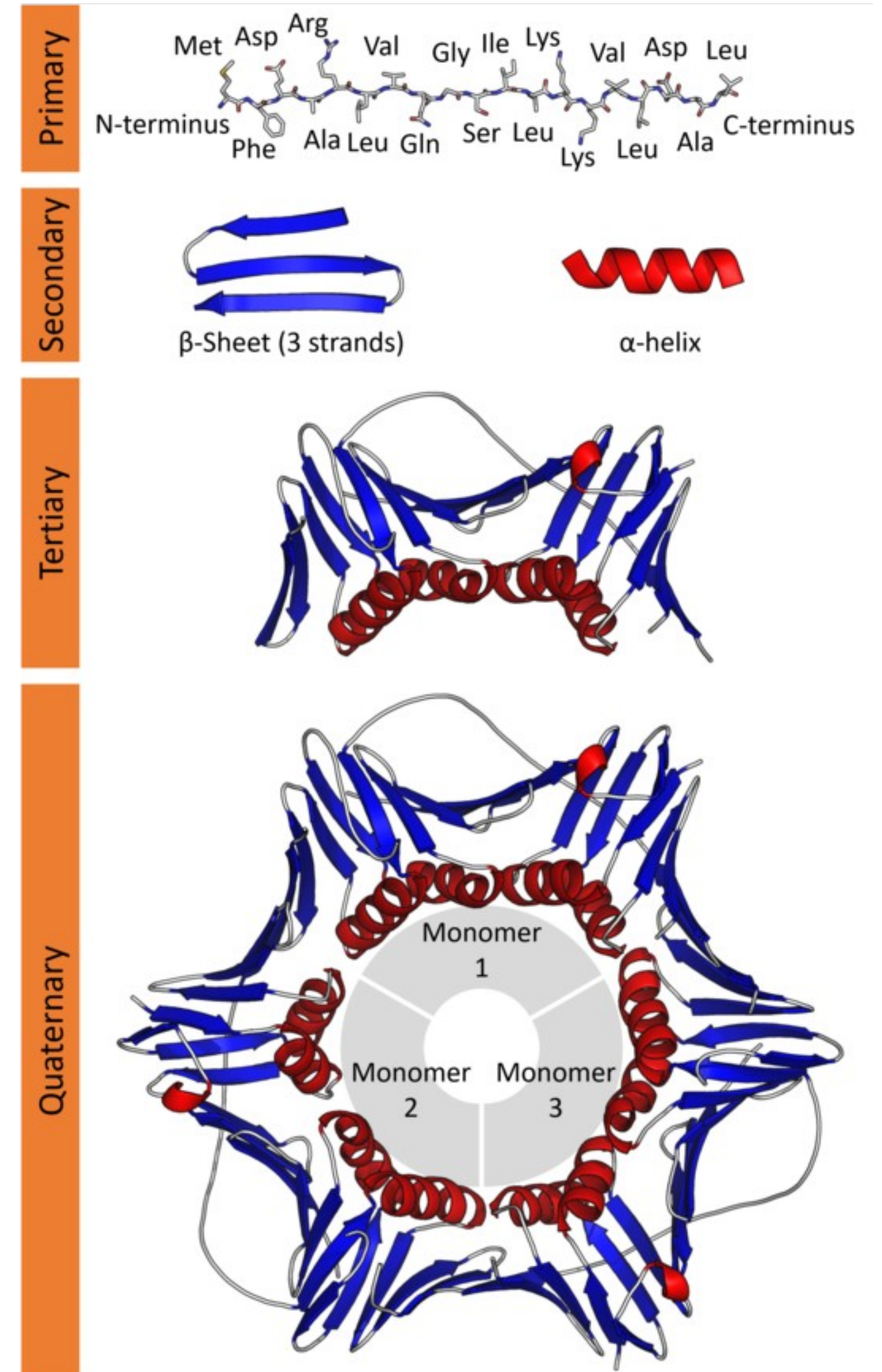


**Dimer (homo)**

**Dimer (hetero)**

**Trimer (homo)**

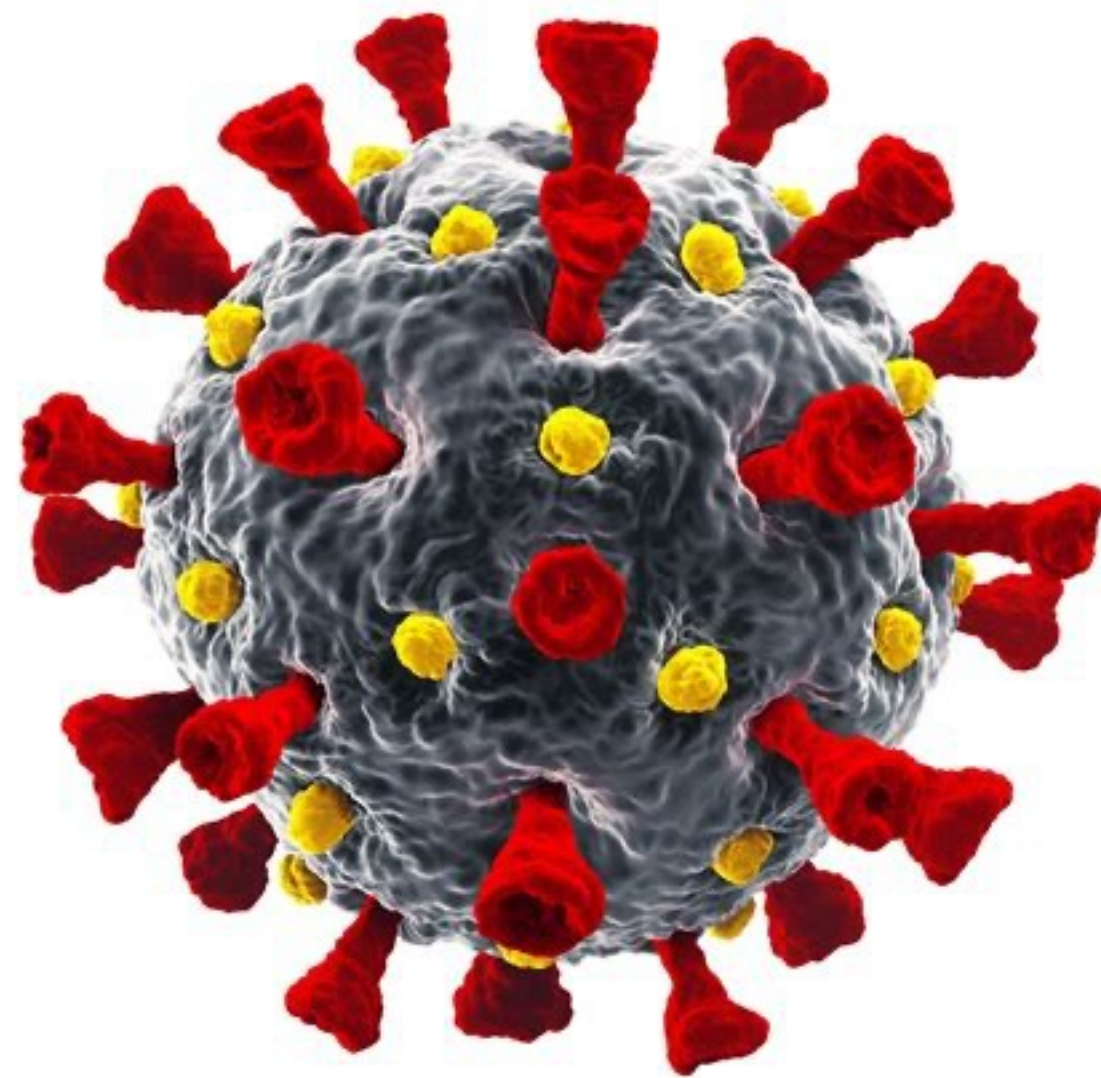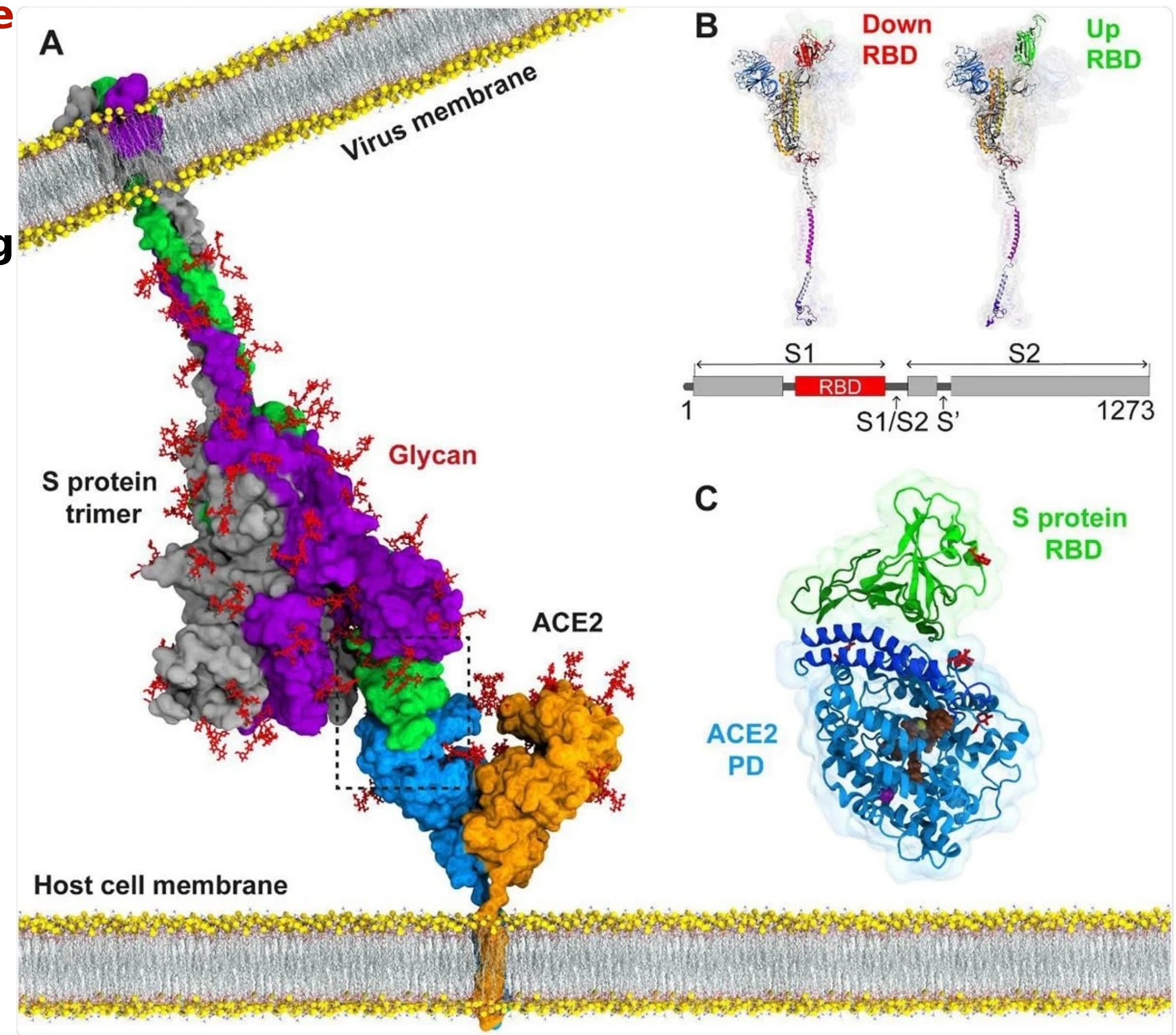**Tetramer (hetero)**



https://en.wikipedia.org/wiki/Protein_structure

2

**Structure = Function → Function = Structure**

- **visualizing binding interface**
- **location of allosteric inhibition**
- **conformational changes**
- **molecular scaffold for molecular docking**
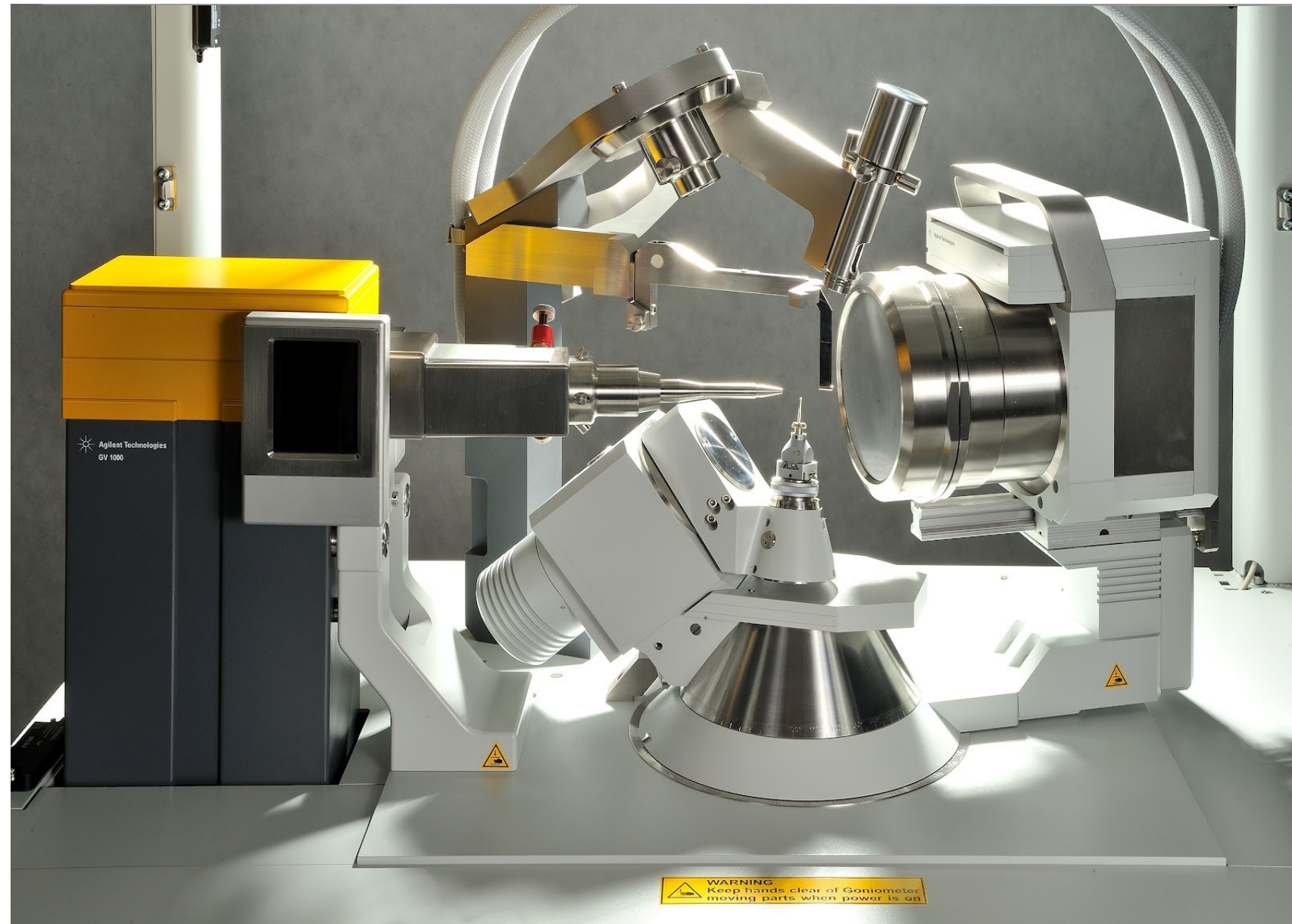- **mutational design**
- **MANY OTHER FACETS...**

SARS-COV-2

# Methods for Elucidating Protein Structure

## X-Ray Crystallography

## NMR Spectroscopy

## Cryo-EM



**pros:**
- **widely used**
- **high-throughput**
- **~2.5-1Å resolution**

**cons:**
- **crystallization process**
- **high variation**

**pros:**
- **native state**
- **non-destructive**
- **real-time**

**cons:**
- **only works on "small" proteins**
  - **<100kDa**
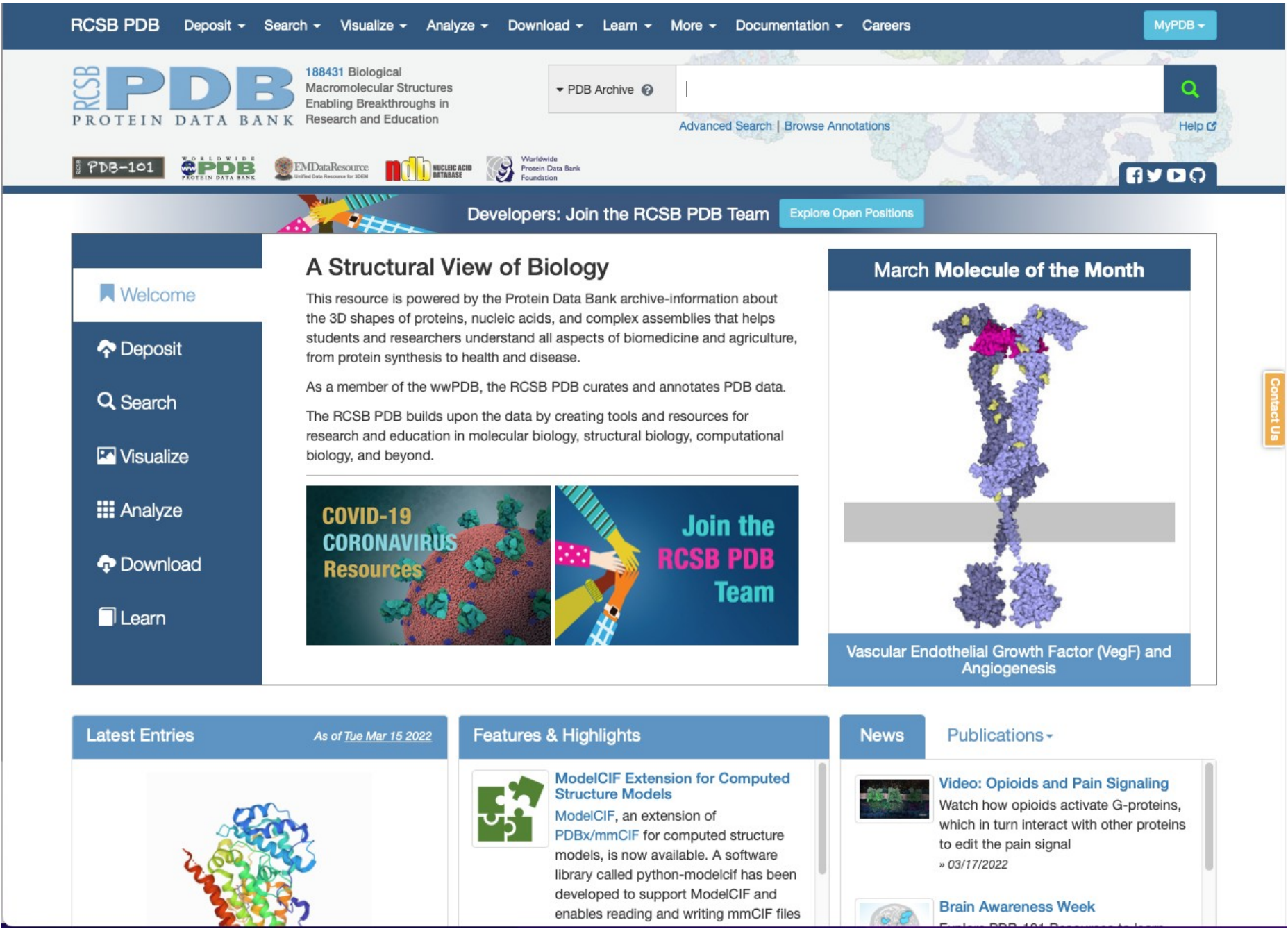- **high [protein]**
- **$$$ maintenance**

**pros:**
- **large proteins/complexes**
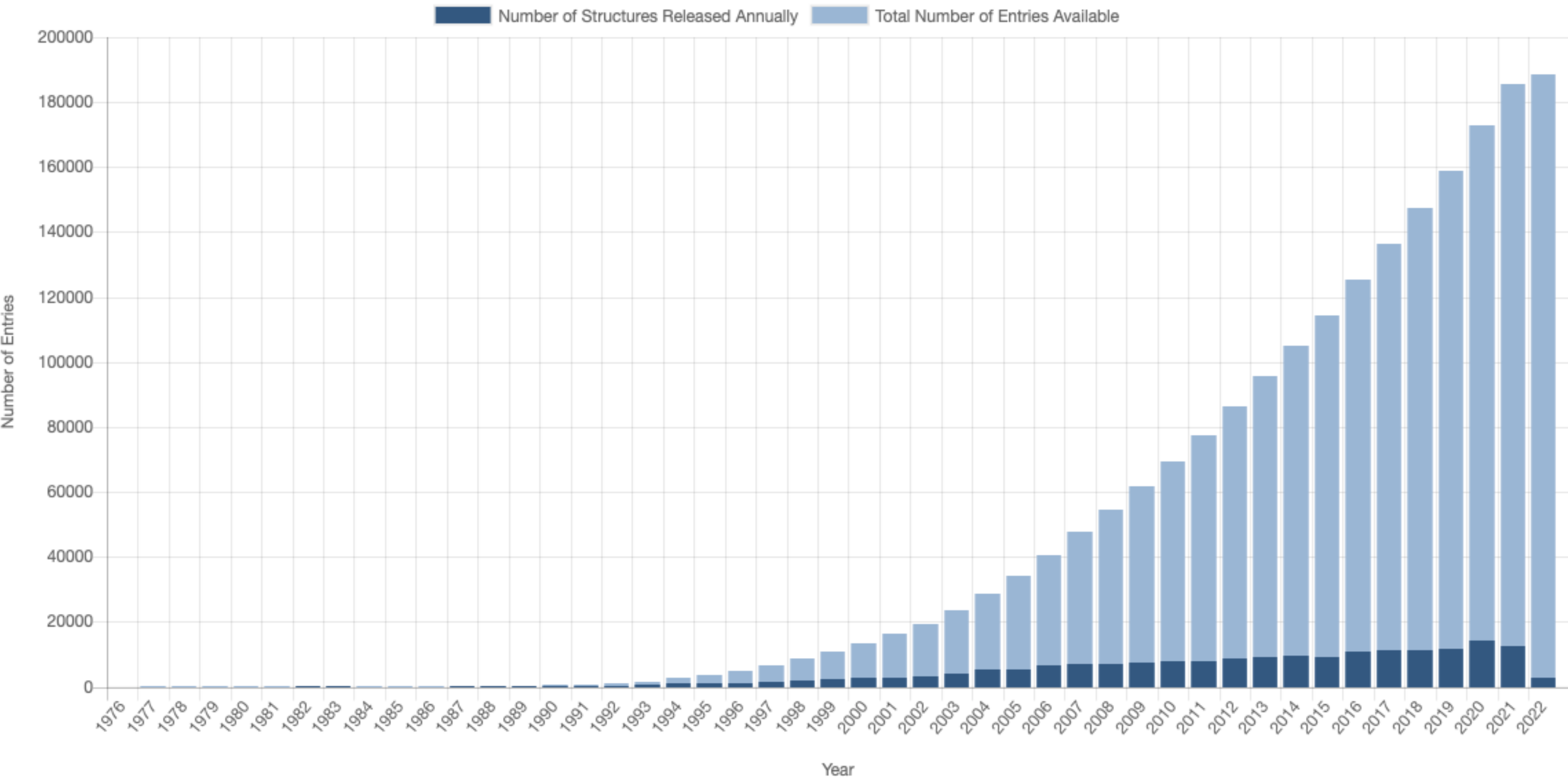- **near-native state**

**cons:**
- **only works on "large" proteins**
  - **>200kDa**
- **freezing samples**
- **computationally intensive**
- **$$$ maintenance**

# Protein Data Bank (PDB)

**After we elucidate a structure where does it go?**



**Most publishing journals require a structure be deposited to the PDB prior to publication of a study!**

**Parameters th<u>at MUST</u> be considered:**
1. **Primary Structure (bond angles)**
2. **Secondary Structure (α-helix, β-sheet, loops)**
3. **Tertiary Structure (folding of secondary structure)**
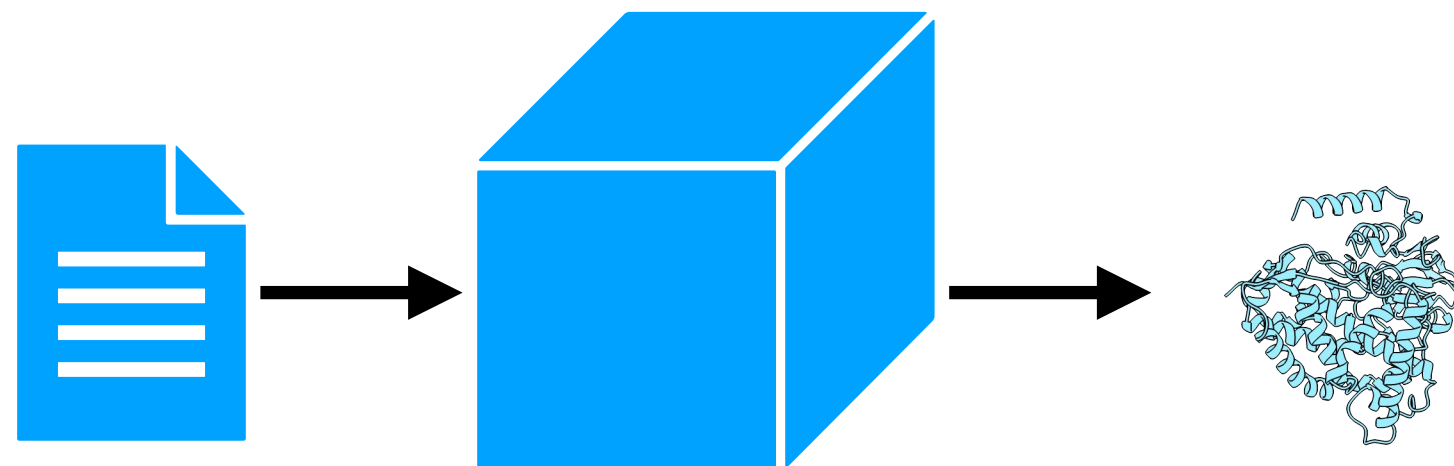4. **Quaternary Structure (Optional)**
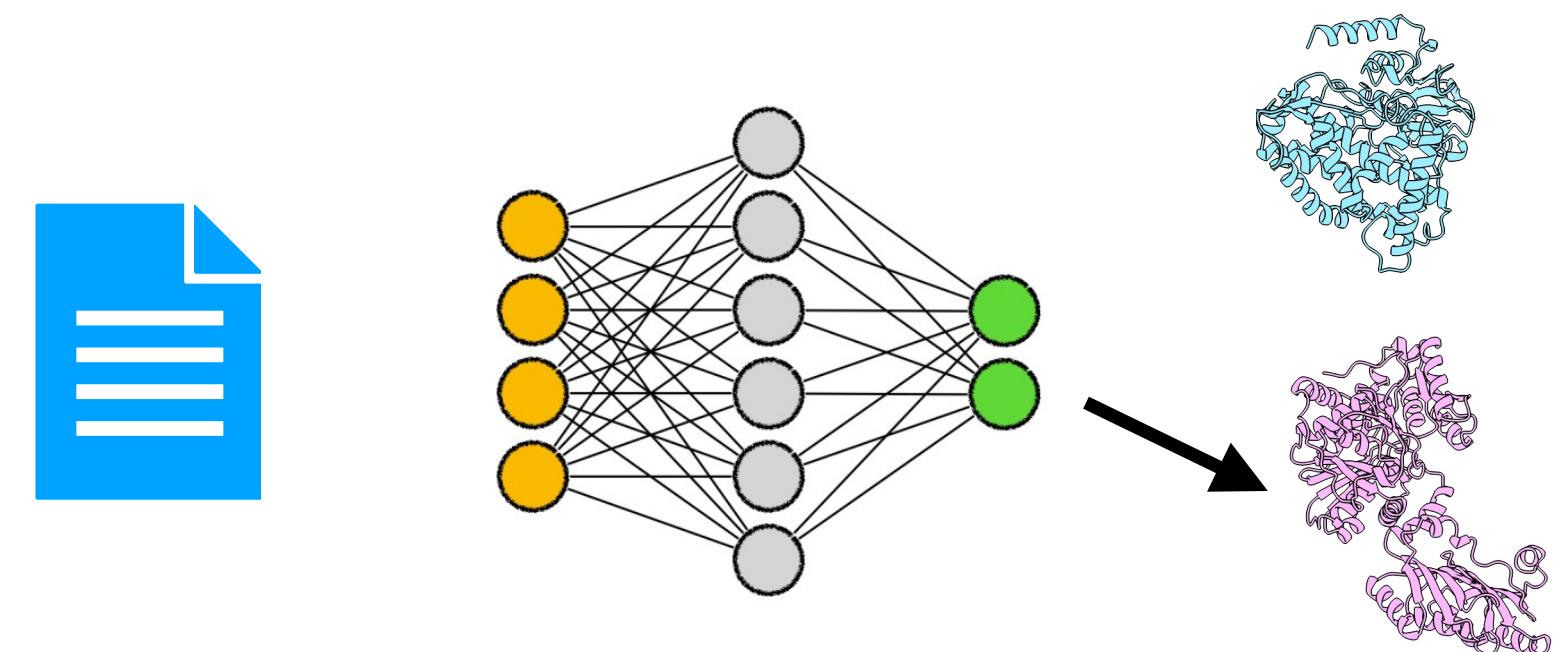
**3°/4° Structure Have Additional Parameters:**
1. **Conserved protein domains**
2. **Protein families/superfamilies/clans**

**Ideally, we would like to predict a protein's 3-D structure given only its AA sequence (1° stucture)**
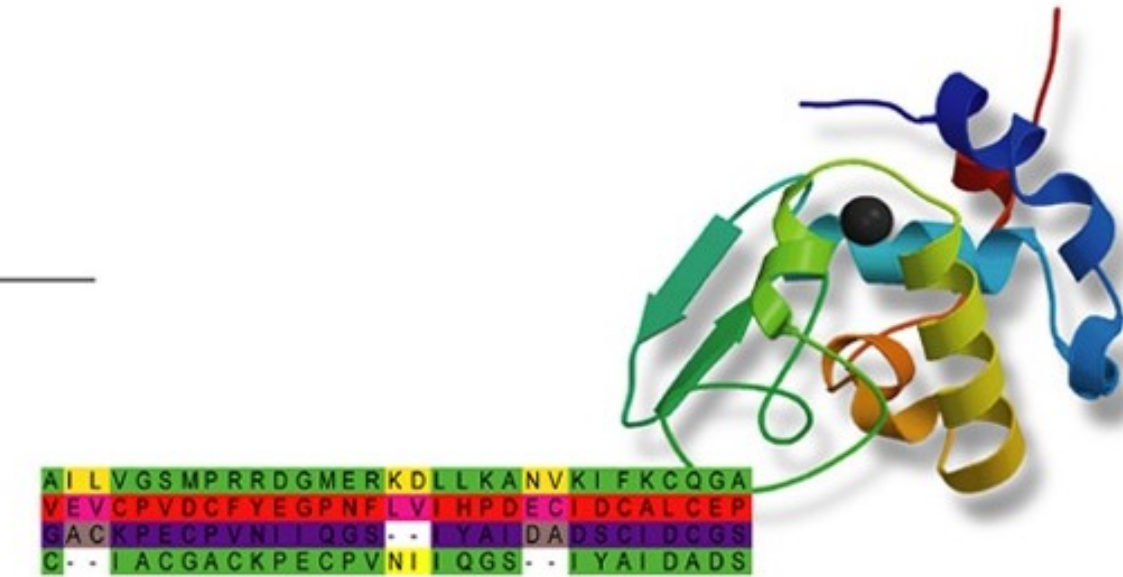
**Algorithm-Based**

**Machine Learning (AI)**

**I-TASSER**
Protein Structure & Function Predictions

**Modeller**
Program for Comparative Protein
Structure Modelling by Satisfaction
of Spatial Restraints

- **Widely used molecular replacement modeling system**
- **Often paired with X-ray crystallography during the refinement step**
- **Best case use when a high sequence homolog exists**

- **CASP rated highest for blind protein search/modeling**
- **Completely open-source**

BIOZENTRUM
University of Basel
The Center for Molecular Life Sciences    SWISS-MODEL

- **Fully-integrated with ExPASy Suite**
- **Extensive manually curated databases**

**All of these solutions heavily rely on homology to some characterized protein**
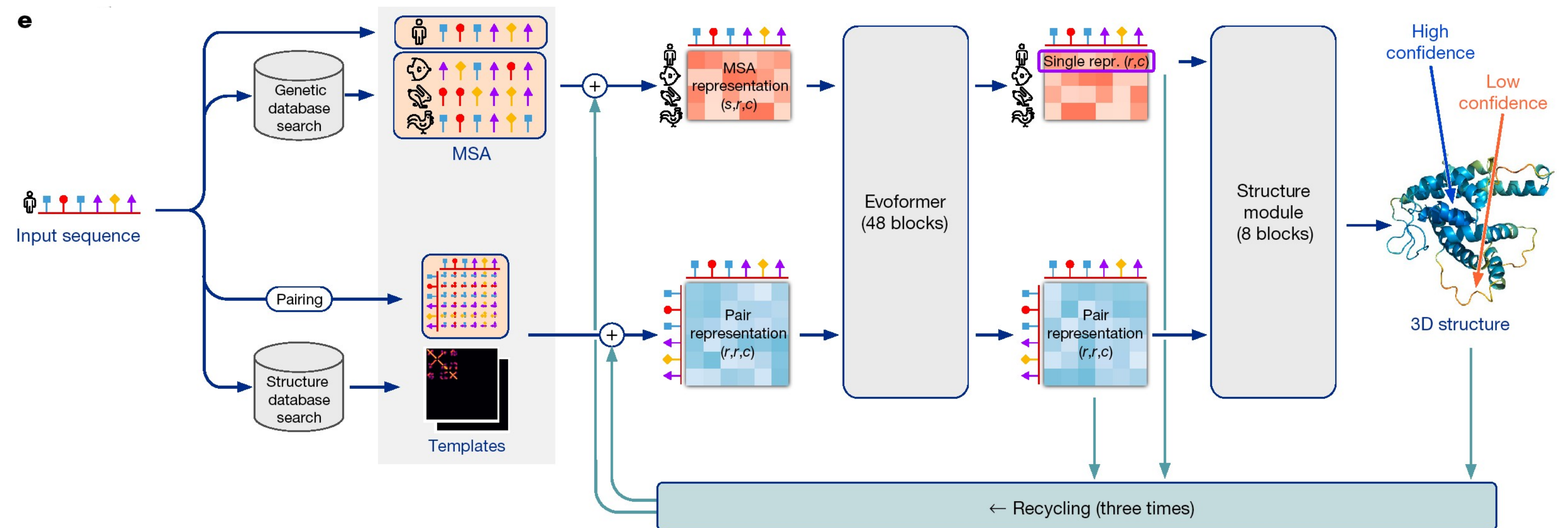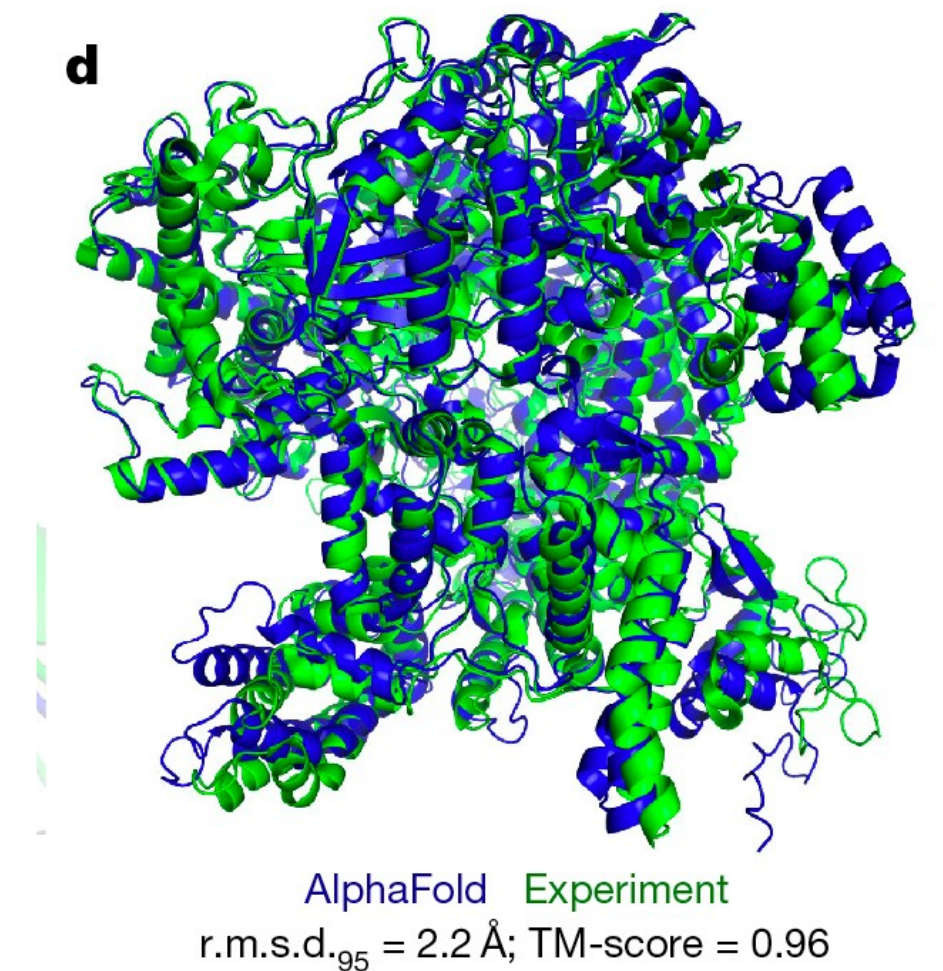
CASP8 target 512-D1
all models
(3dsm)

**CASP evaluates current methods of predicting protein structure since 1994.**

"These techniques are expensive and slow: it can take hundreds of thousands of dollars and years of trial and error for each protein. AlphaFold can find a protein's shape in a few days." - MIT Technology Review

- **2020 was the first year a program met a true success rate (>90%)**

  - **DeepMind's AlphaFold2**

  - **Worked on both homologous and novel proteins**

  - **Only needed to provide the AA sequence**

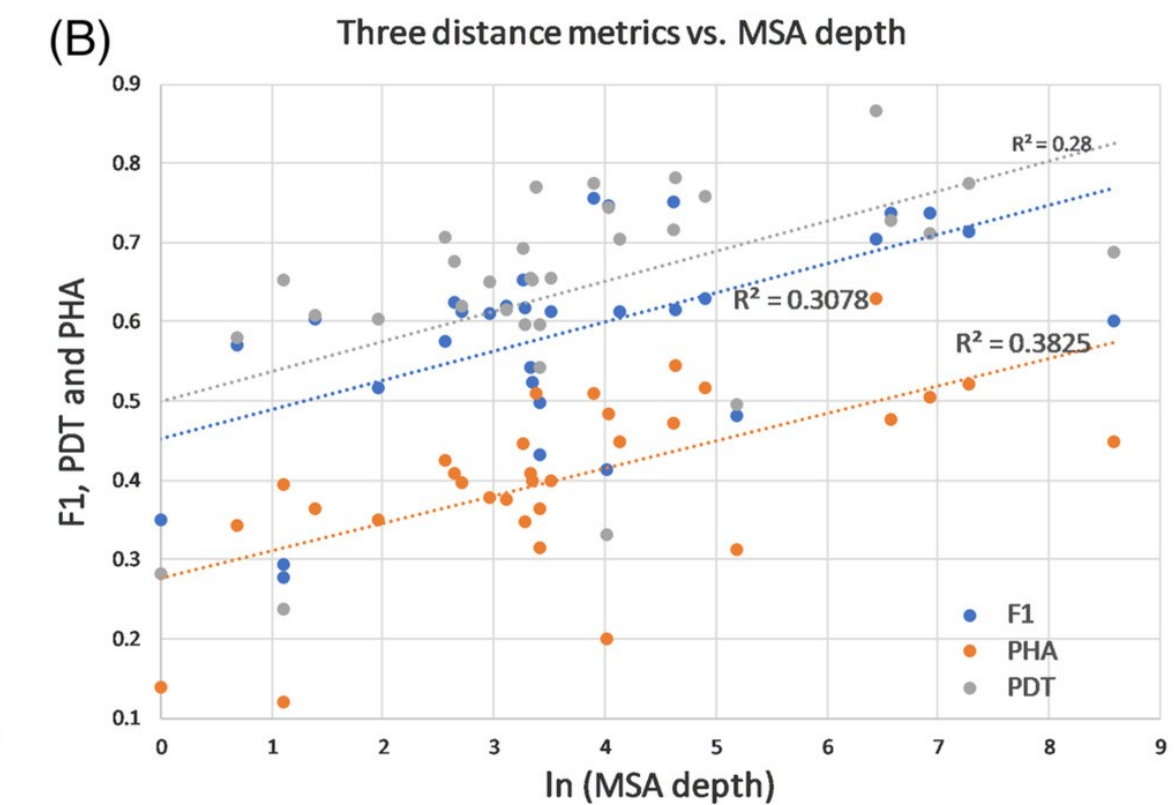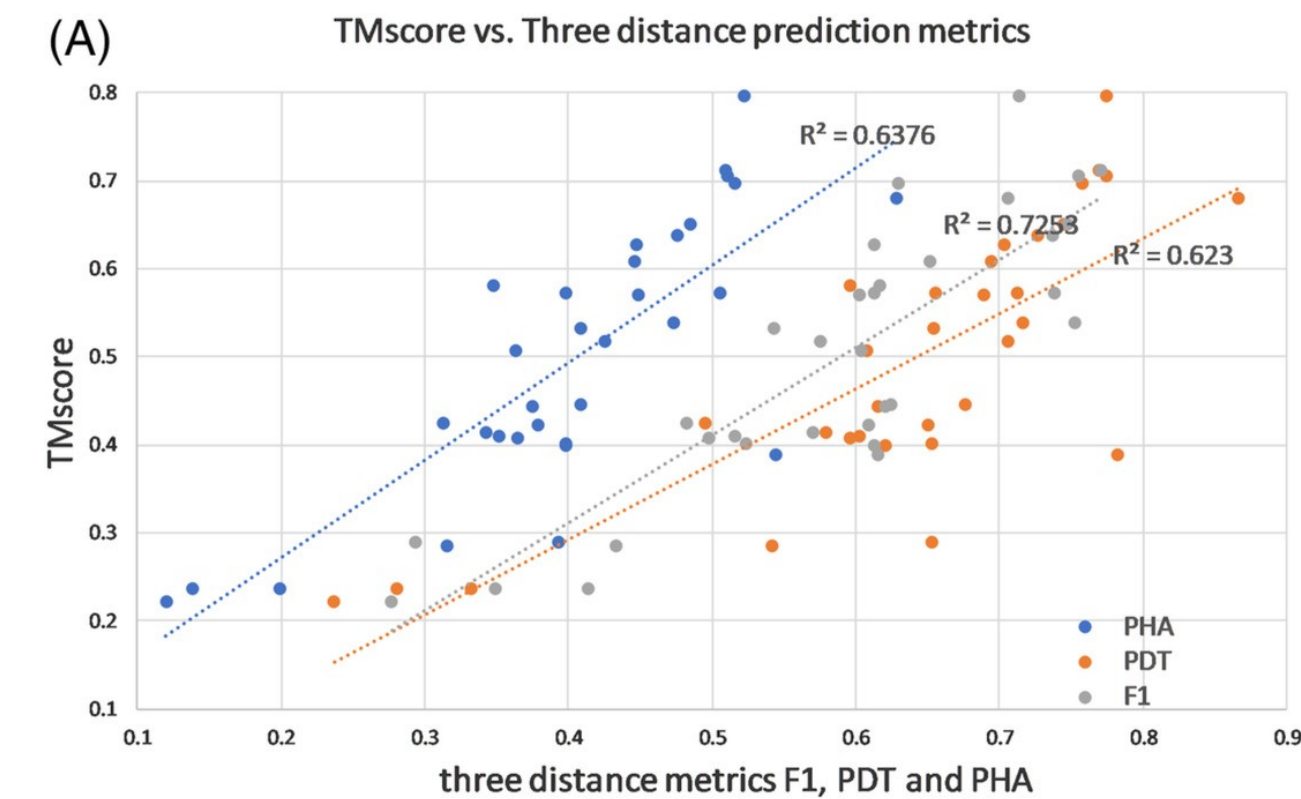  - **AlphaFold2 was released under the Apache Common Use license in 2021**

- **First program to ever meet the success rate at CASP**

- **Assessed on alignment of predicted structure to experimental**

  - **RMSD (Å)**

- **AF2 custom value to "quantify" model confidence**
  - **pLDDT**

AlphaFold    Experiment
r.m.s.d.95 = 2.2 Å; TM-score = 0.96

Structure Prediction is only as good as the input MSA

https://docs.google.com/presentation/d/

Xu, J.*et. al*. *Proteins* 2019, *87* (12), 1069-1081.

**Typically MSA depth of >30 sequences/residue produce a highly accurate and confident model**

## Structure Template Search Space



**PDB70**

**PDB_mmcif**

**PDBseqres**

**Recently, it has been observed that structure templates have little to no impact on the final predicted strucutre's accuracy.**

## Sequencing Alignment Search Space



**MGnify (bacterial proteome darkspace)**

**Uniref (dereplicated UniprotKB)**

**Uniclust30 (% clustering of UniprotKB)**



**BFD (massive cluster of bacterial proteome space)**

**Uniprot Database)**

11

**CYP153A - Bacterial cytochrome p450 (collaboration)**

*S. lycopersicum*
**Threonine Deaminase 2**

**Flagellum-22 Up-regulated**
*S. lycopersicum* **Protease**

- No "cost" to you

- Best accuracy with homolog. seqns.

- Running the same protein yeilds differences (Å)

  - CollabFold number of recycling steps

# Grace Cluster



"Grace is a 925-node Intel cluster from Dell with an InfiniBand HDR-100 interconnect, A100 GPUs, RTX 6000 GPUs and T4 GPUs. All nodes are based on the Intel Cascade Lake processor." - HPRC

# Using AlphaFold2

**Terminal via SSH (Local Client)**



**OnDemand Portal (Web Browser)**



**Google Colab Notebook (Colab Fold)**

- **AlphaFold2 was the first program to be considered a success for protein structure prediction**
- **It works on proteins with structural homologs, and those without**
- **Its major limitation is still requiring homologous sequences for accurate structure prediction**
- **Its dependent on two parts:**
  - **Finding homologous sequences**
  - **Finding homologous structures**
- **The confidence of the model is expressed as pLDDT**
- **When in doubt, run AF2 on your protein, worst that happens is you waste your time**

**Any Questions!**



**devonjboland@tamu.edu**