

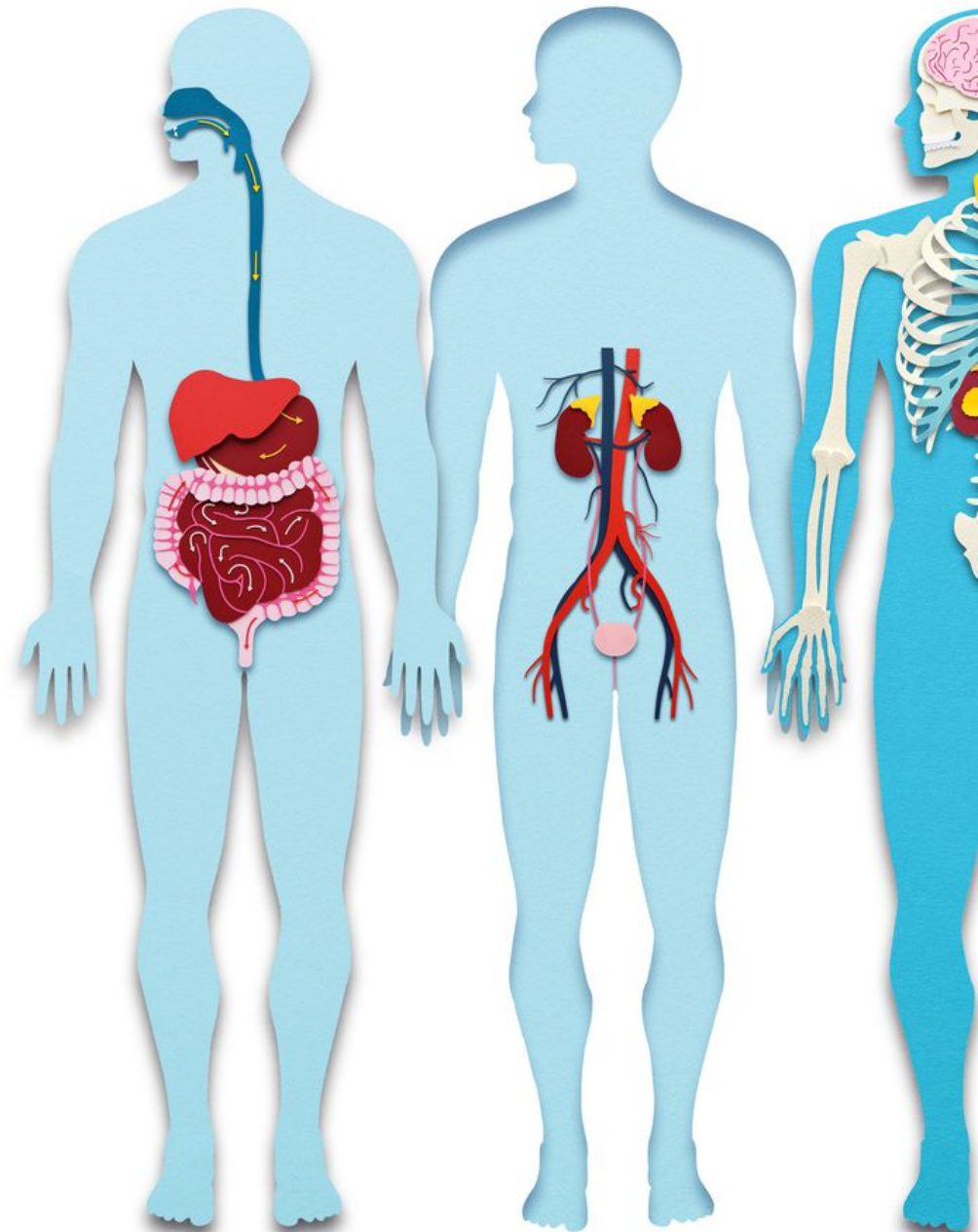
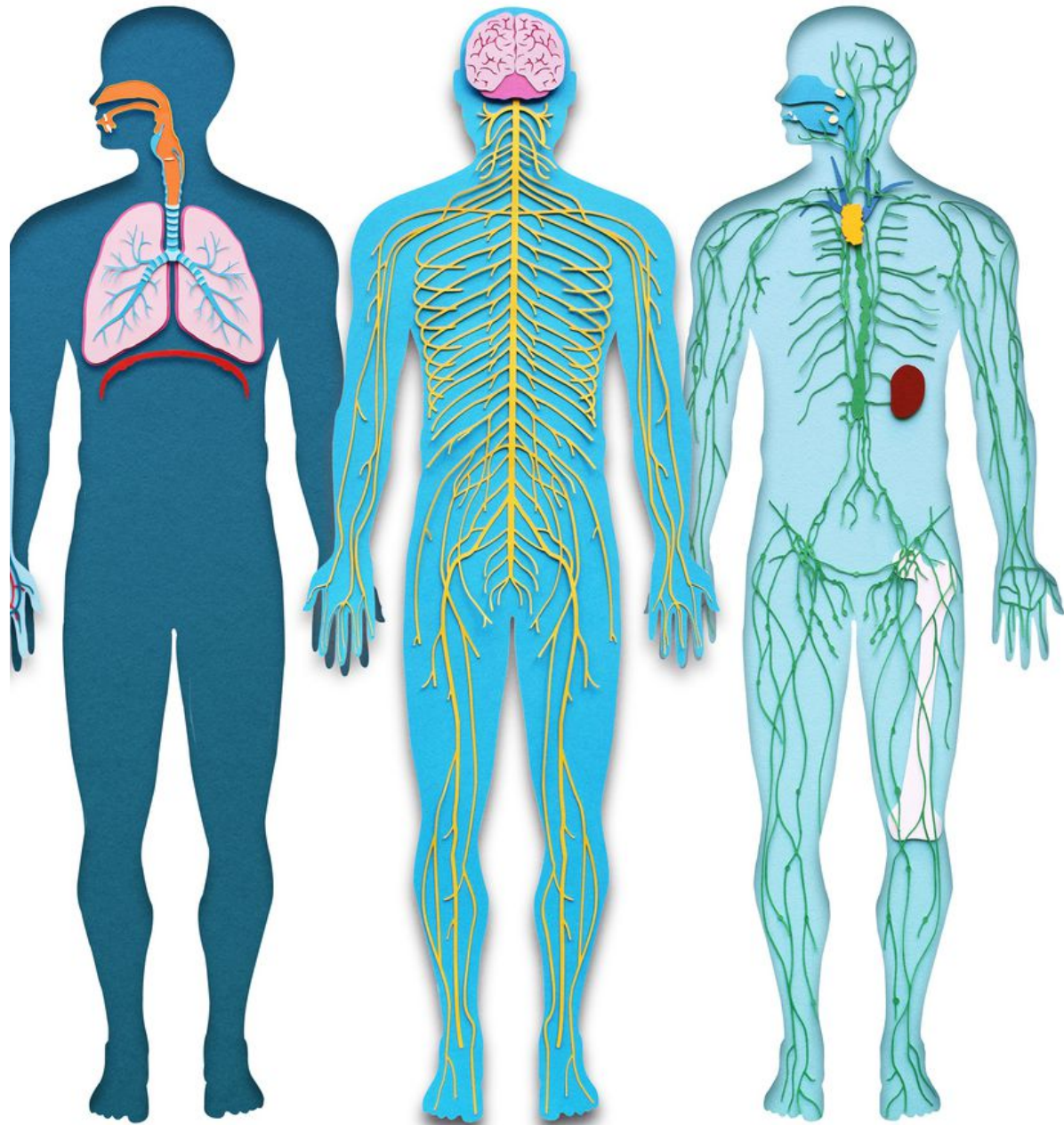
# Deep Learning in Structural Biology and Protein Design: Where, How, and Why

Deep Learning Theory Workshop and Summer School

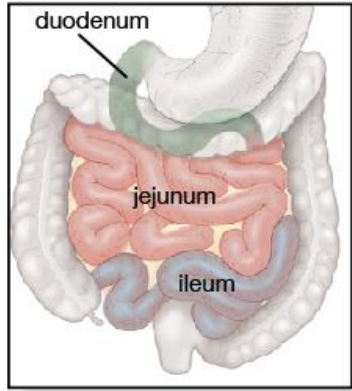
Aug 5, 2022 @ Simons Institute

By Chloe Hsu (UC Berkeley)

Intro: A view of biology from living beings to molecules



### Regions of the small intestine



large intestine

small intestine

esophagus

stomach

mucosa

submucosa

plicae circulares  
(valves of Kerckring)

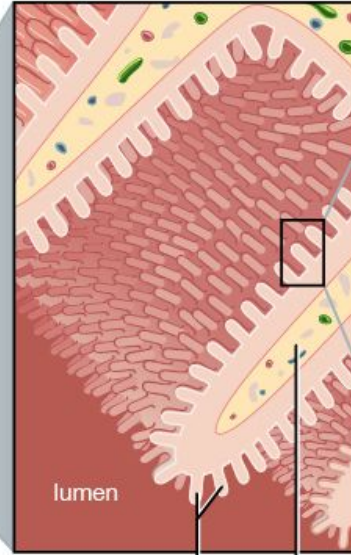
lumen

serosa

longitudinal  
muscle layer

circular  
muscle layer

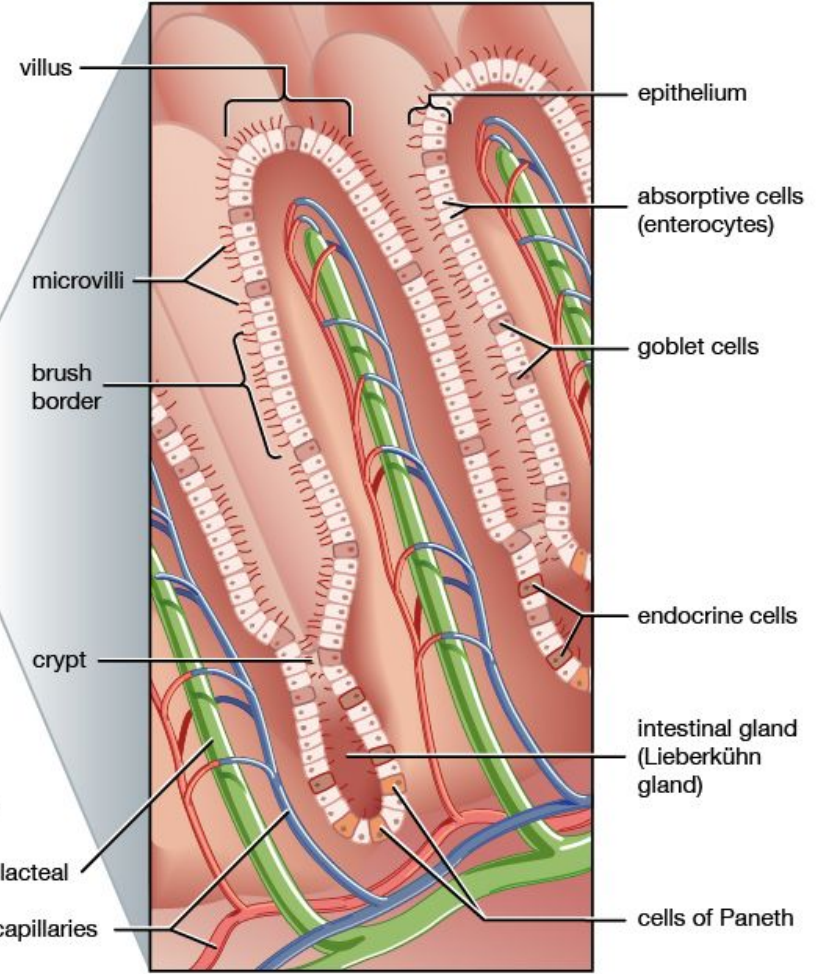
### Enlargement of plicae circulares

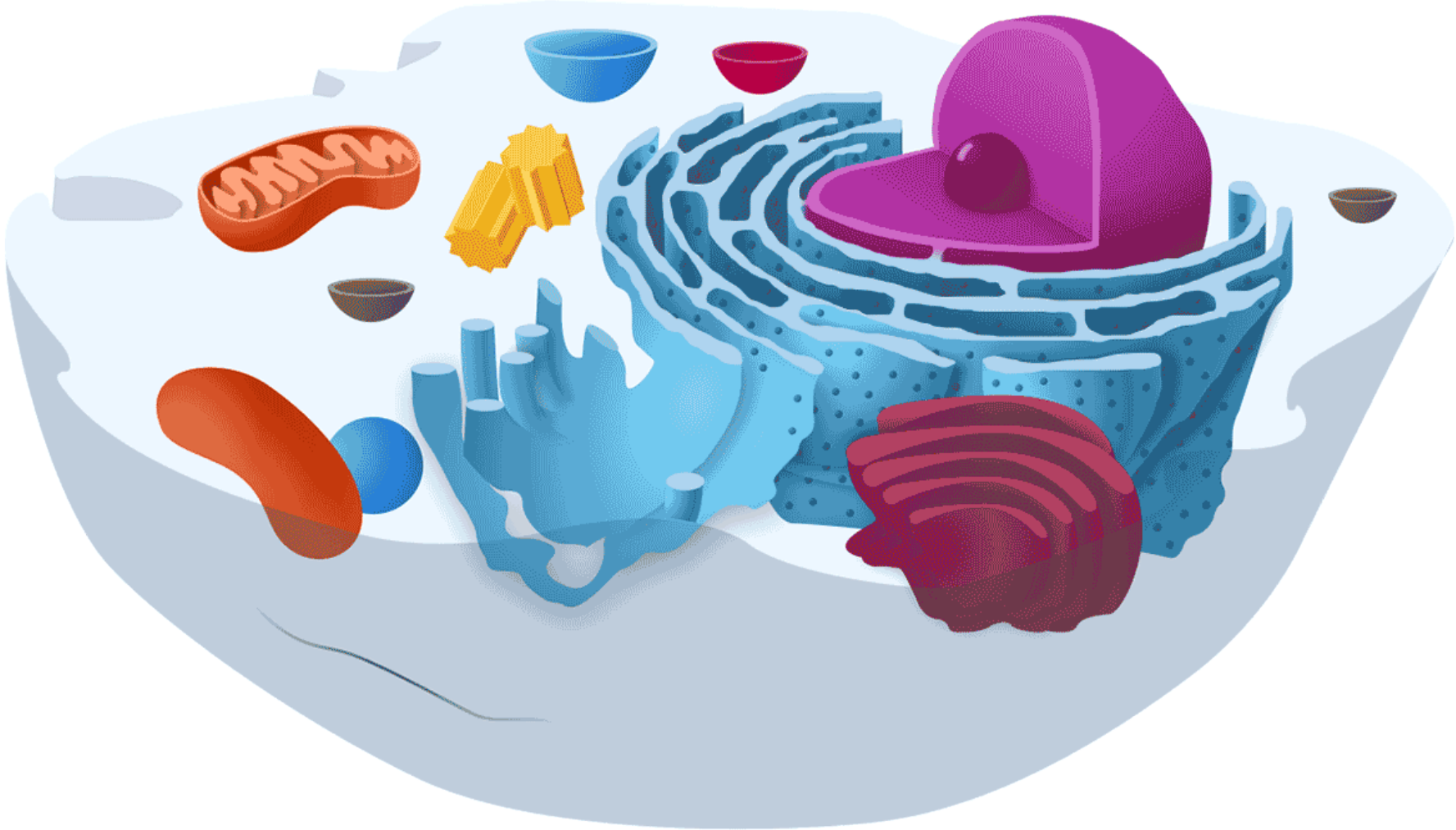


villi

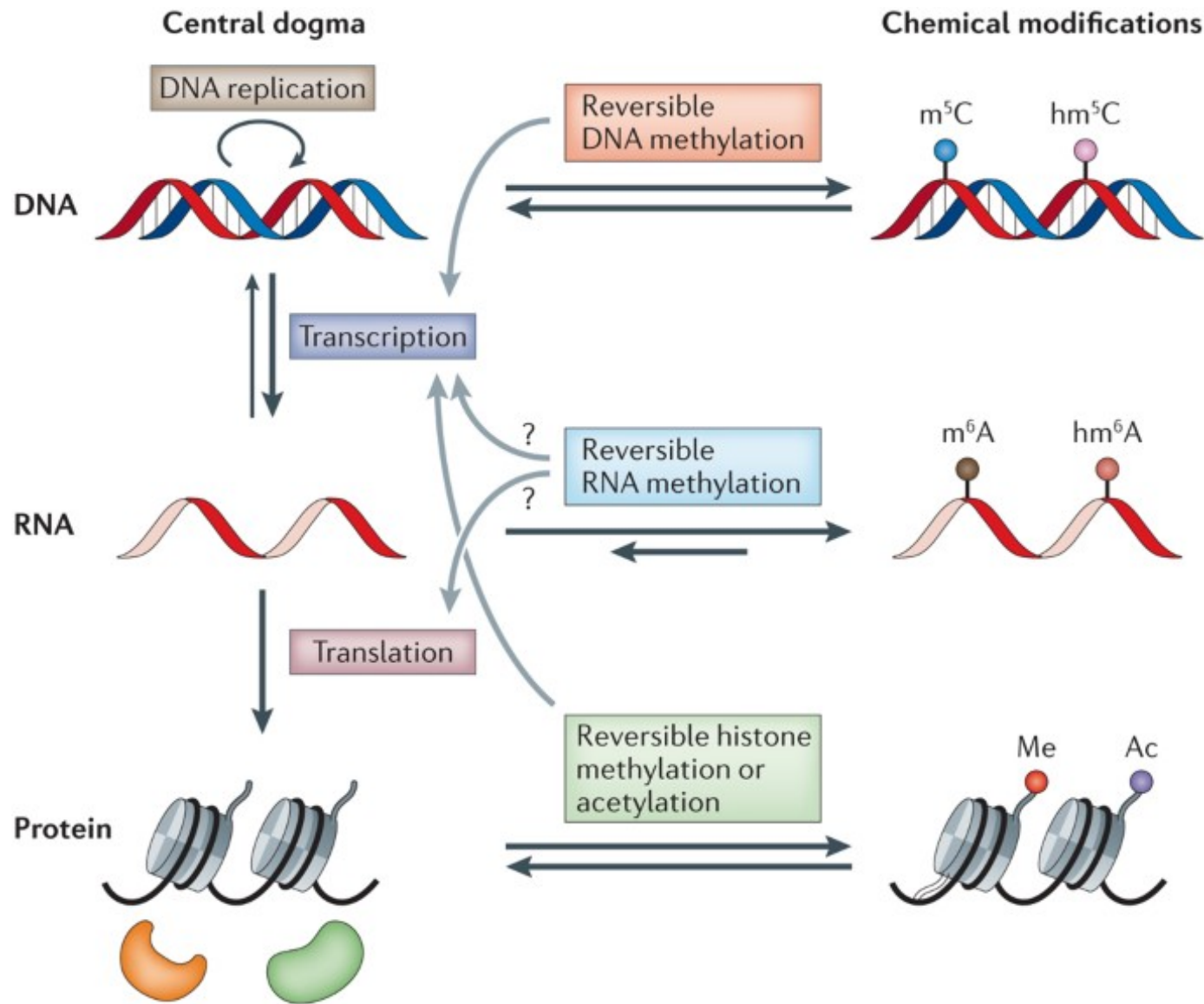
plica circularis

### Structure of a villus





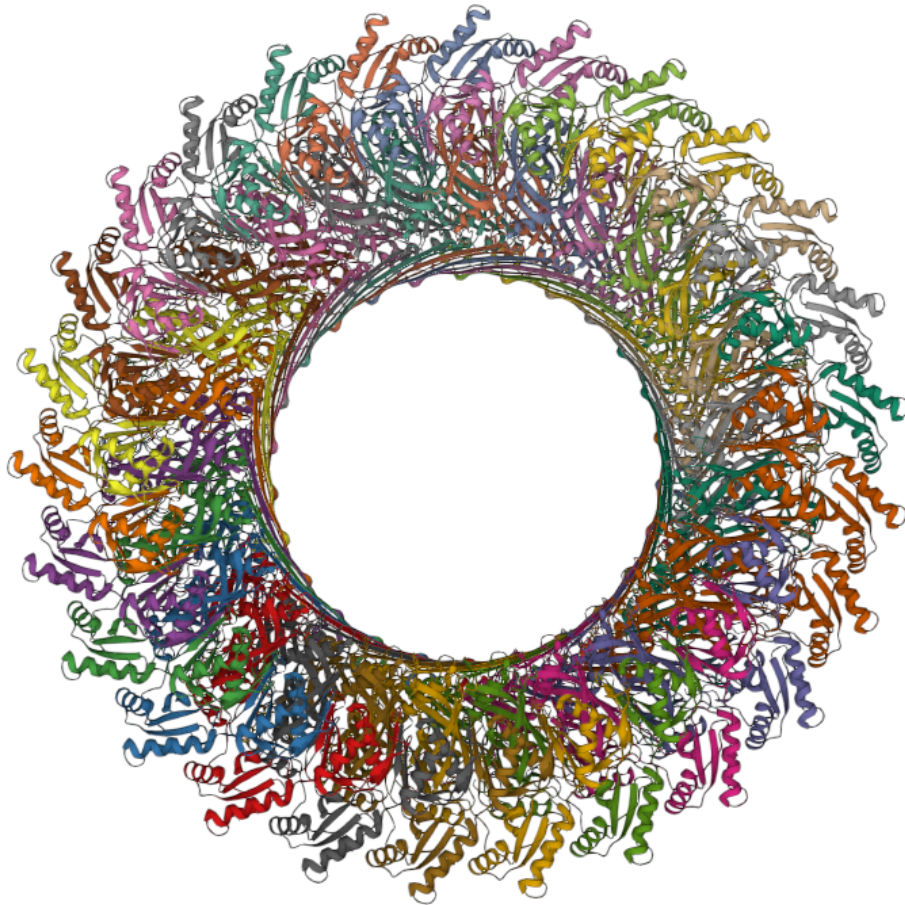




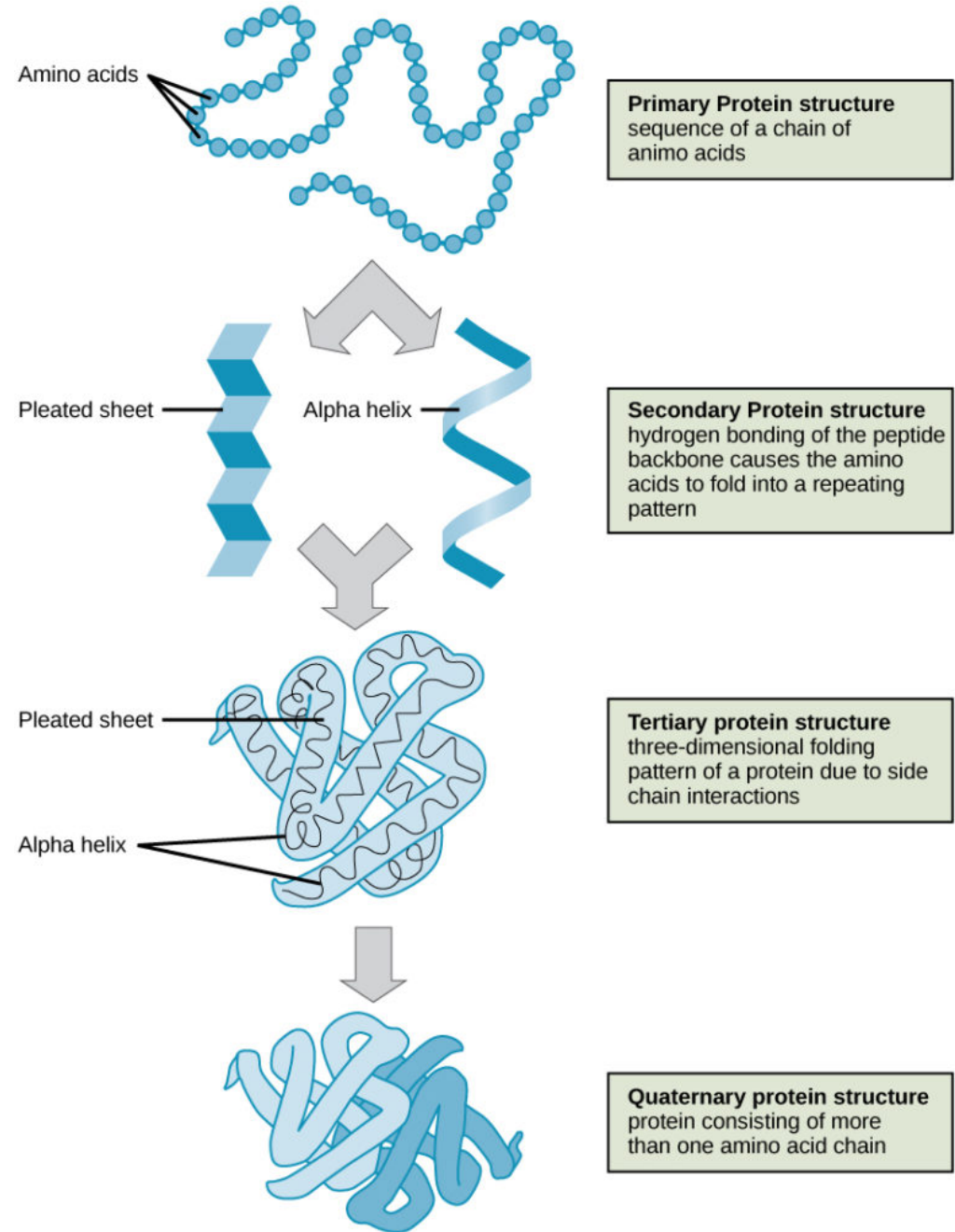
“Where”: Advances in structural biology and protein design



# Protein structures



Salmonella LP ring (It moves and swims!)

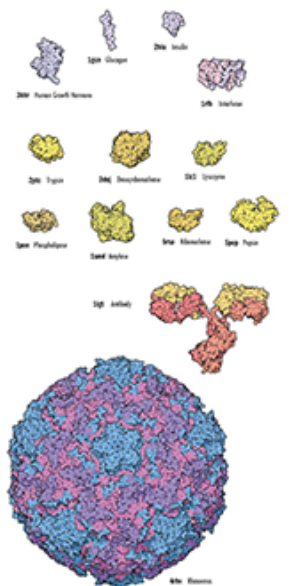


# MOLECULAR MACHINERY: A Tour of the Protein Data Bank

**1** Living cells are filled with complex molecular machinery, a million times smaller than familiar machines like computers or automobiles. Cells use these tiny molecular machines to perform all of the jobs needed for life. Some are molecular scissors that cut food into cell-sized pieces. Some build new molecules when cells grow or when damaged tissues are repaired. Some are molecular boxes and muscles that support cells and help them move and crawl. Some fight off attackers, defending against infection.

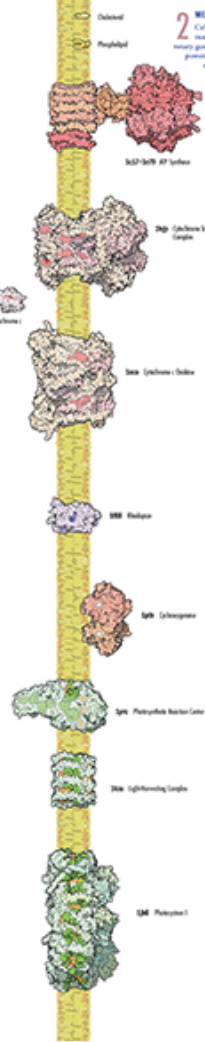
Researchers around the world are studying these molecules and determining their precise atomic structures. These structures are available on the internet through the Protein Data Bank (<http://www.pdb.org>), the central storehouse of biomolecular structures. A few of the thousands of structures held in the Protein Data Bank are shown here. In these pictures, the molecules are all drawn at a magnification of 2,000,000 times, and each atom is shown as a small sphere. Many of these structures are composed of several subunits, which are indicated by different colors. An enormous range of sizes is shown here: the water molecule at the left has only three atoms and the ribosome shown below has hundreds of thousands.

By David S. Goodsell, The Scripps Research Institute, La Jolla, California, USA  
Graphic design by Gal M. Eisenberg, Los Alamos Supercomputer Center

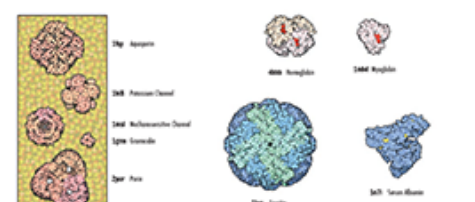


**1 OUTSIDE THE CELL**  
Some molecular machines perform their jobs outside of cells. Many are enzymes, so that they can diffuse quickly to their sites of action. This is one of the four hair combs shown at the top: insulin and albumin, which together regulate blood sugar levels; antibodies, which react against the immune system; and human growth factor. The seven digestive enzymes (to reflect on size) used in this module, so that they can survive the hostile environment in the digestive tract. Each of these enzymes has a small groove (colored towards the top) in which they bind to a different sugar molecule and digest it. At the bottom is elastin, the wire that connects the connective cells, and an antibody, one major defense against viruses. Antibodies bind to viruses and prevent them from binding to cell surfaces, thus blocking infection.

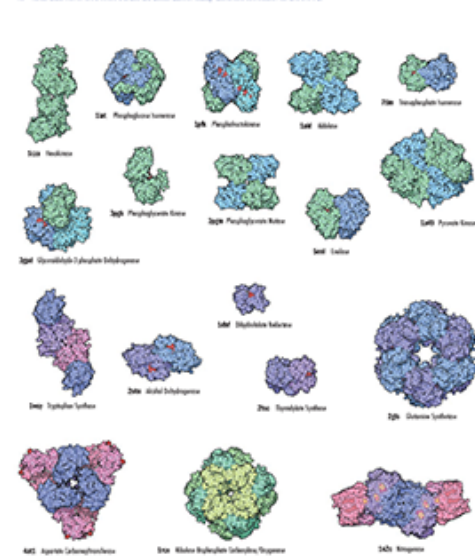
**PROTEIN DATA BANK**  
<http://www.pdb.org/> • [info@rcsb.org](mailto:info@rcsb.org)  
RESEARCH COLLABORATORY FOR STRUCTURAL BIOCHEMISTRY  
HELSING, THE STATE UNIVERSITY OF NEW YORK  
SAR DILLIG SUPERCOMPUTER CENTER  
NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY



**2 NEWERIDES**  
Cells are surrounded by a membrane made of lipids, like the phospholipid and cholesterol molecules shown at the top. Membranes keep the cellular machinery inside and unwanted material out. Many proteins are embedded in the membrane, performing a variety of essential tasks. ATP synthase is a rotary generator that produces ATP (adenosine triphosphate), the small molecule used for powering cells. The two large complexes below it charge a battery that stores ATP molecules, and the two protein synthetases a double channel between them. Ribosome is found in membranes in the active. The small vesicle molecule inside of it changes shape when illustrated, causing the surrounding protein to send a signal to the brain. Cytochrome *c* holds one of the molecules used to signal pain—also cytochrome molecule here, however, is linked by two molecules of apigenin, shown inside in white. In the bottom are two molecules involved in photosynthesis, which capture energy from light and use it to pump the electrons of sugar in plant cells.



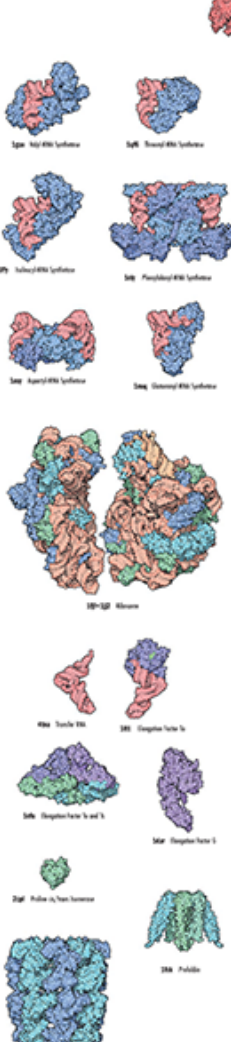
**3 TRANSPORT AND STORAGE**  
Of course, a perfectly sealed membrane would be of little use to cells, because nutrients could not get in and wastes could not get out. The box shows a membrane binding face on five proteins that form channels through the membrane as shown. To the right of the box are several soluble proteins involved in transport and storage of molecules. Hemoglobin and myoglobin carry oxygen. Ferritin stores a lot of iron that is used later. Ferritin also carries many different molecules in the blood.



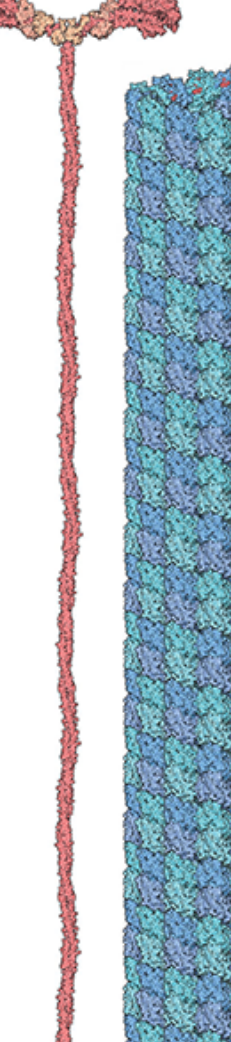
**4 CHEMICAL FACTORIES**  
Cells hold a bewildering variety of enzymes—proteins that perform chemical reactions. At the top are the two enzymes that perform glycolysis, the breakdown of sugar to form ATP. Below are two related enzymes that perform different phosphorylation reactions. Diphospholipase activates a few cofactor molecules and alcohol dehydrogenase breaks down alcohol. Ribonucleic polymerase catalyzes transcription in the most common manner on the Earth, and performs a key step in the synthesis of carbon dioxide by plants to form sugar. The three enzymes and the transferrin make different building blocks for creating new molecules. Nitrogenase performs an essential role in the ecosystem by converting nitrogen gas into a form that living cells can use.



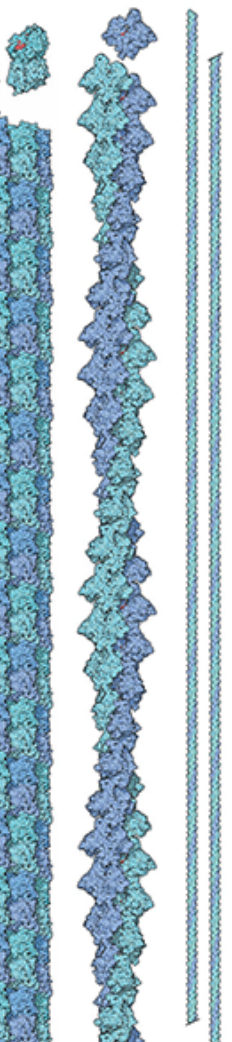
**5 DNA**  
Genetic information is stored in the DNA double helix, with varying forms top to bottom here. Many proteins are used to replicate, read, and store this information. RNA polymerase copies the information from a strand of DNA that will be used to direct the construction of new proteins. It is assisted by transcription factors, which release enzymes when the helix is unwound and unwound, and guide it to appropriate starting points by the two protein complexes below it. DNA polymerase replicates DNA strands; then, the protein is filling a gap in the double helix. Some proteins, like the top nitrogenase, guide DNA and bend it sharply, or even wrap it all the way around themselves, like the two nucleosomes at the bottom.



**6 BUILDING NEW PROTEINS**  
New proteins are built by ribosomes using the molecular factories that read the genetic code and use it to direct construction. Many molecular machines are needed to make the protein. Twenty different aminoacyl-tRNA molecules (top) are shown below the building blocks of amino acids, ready to be added to a growing protein chain. Several protein factors, shown below the ribosome, guide each tRNA into the proper spot. The three ribosome proteins shown at the bottom help each new protein fold into a proper shape.



**7 BEANS AND BIRDS**  
Cells are held up and supported by a complex infrastructure. This infrastructure is formed of sturdy filaments like actin and microtubules, composed of many subunits stacked like bricks. Myosin is a molecular motor that crawls along actin filaments, allowing the cell to move. Collagen, broken into two pieces here, is actually bound outside of cells, where it forms connective tissue between cells.



**8**

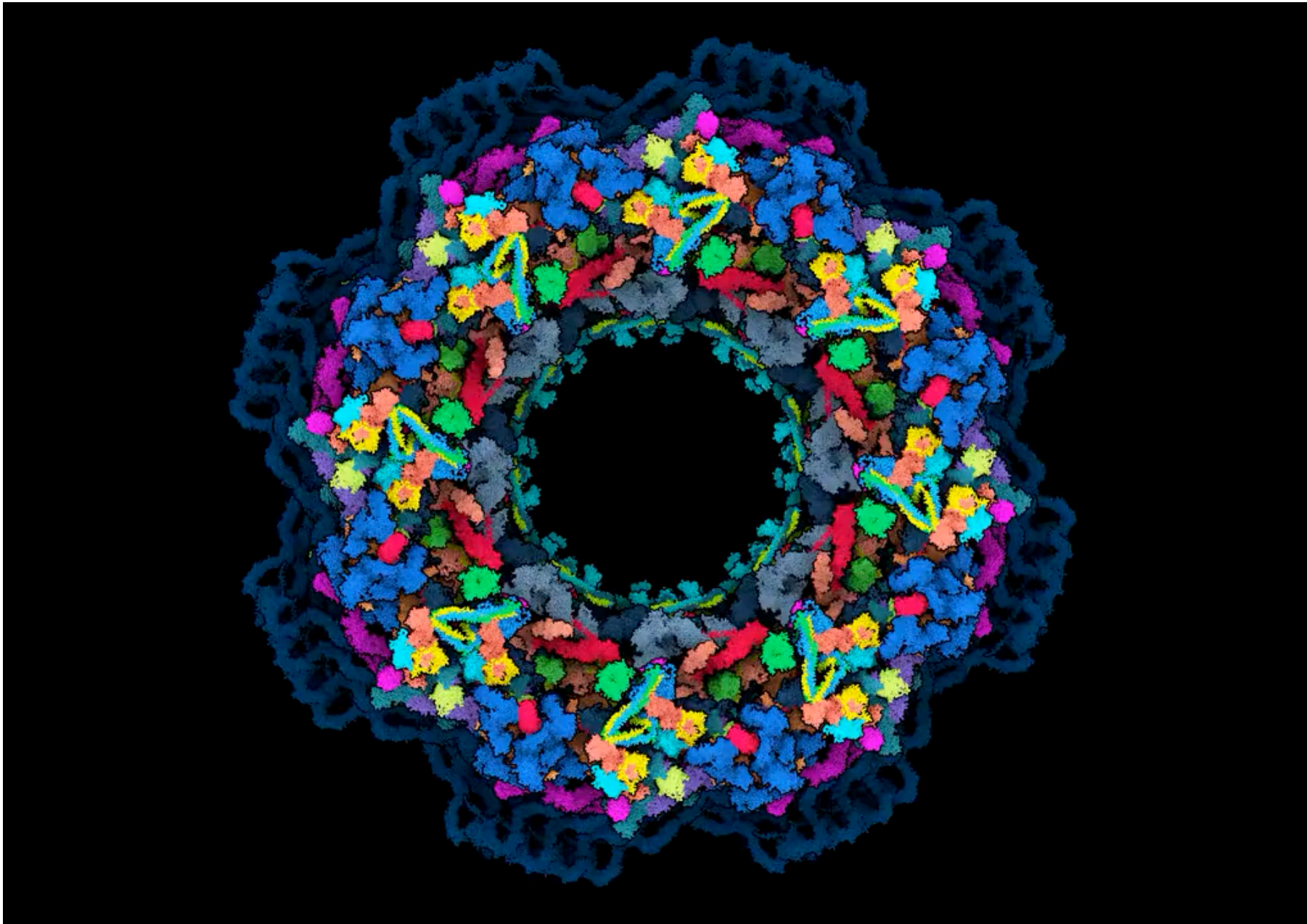
# Structure prediction

Substantial improvement in structure prediction from AlphaFold2

Also notable in new capabilities for predicting protein-protein interactions (although not perfect yet)

Applications:

- **Interprets complex experimental data in structural biology at a new speed**
- Creates an unprecedented size of predicted structures for data mining & learning
- ...



Mosalaganti, Shyamal, et al. "Artificial intelligence reveals nuclear pore complexity." BioRxiv (2021).

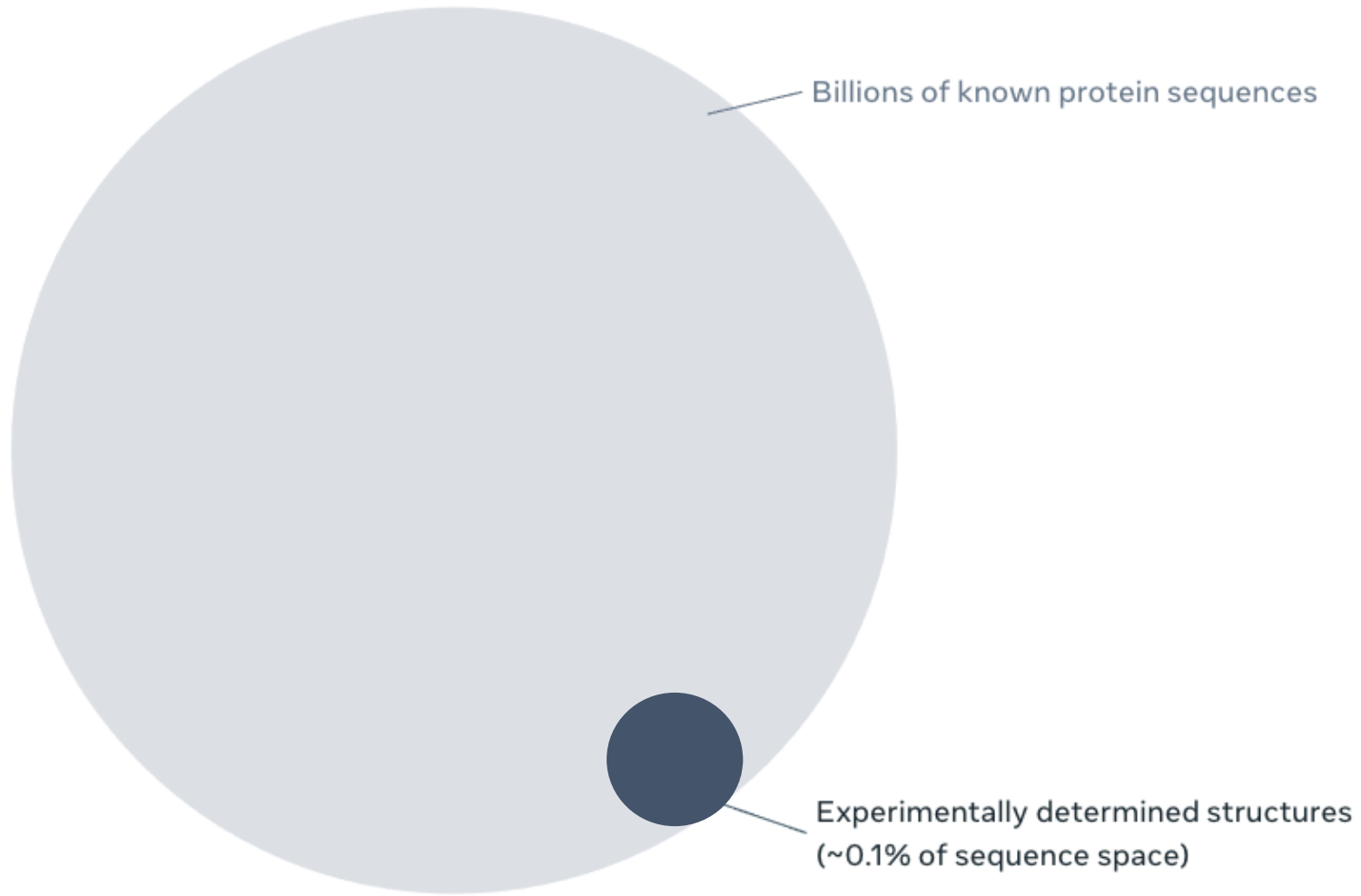
# Structure prediction

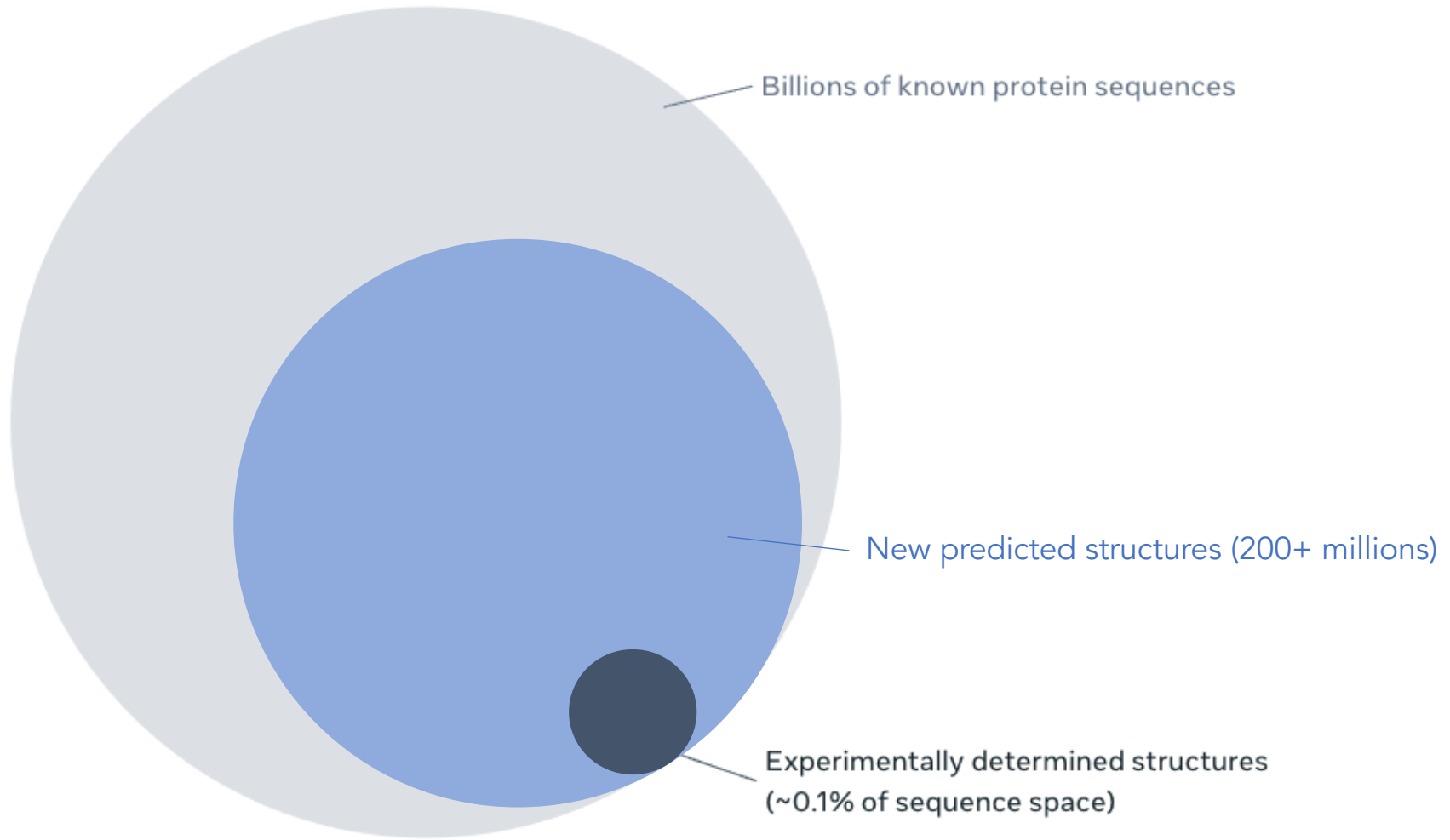
Substantial improvement in structure prediction from AlphaFold2

Also notable in new capabilities for predicting protein-protein interactions (although not perfect yet)

Applications:

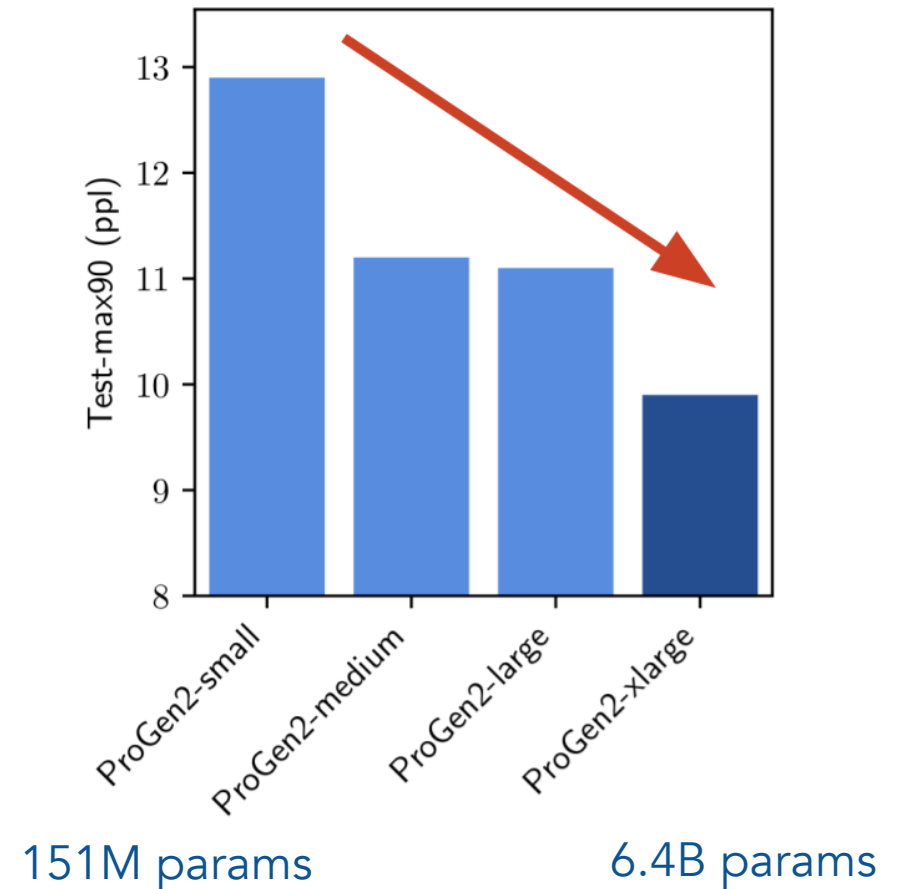
- Interprets complex experimental data in structural biology at a new speed
- **Creates an unprecedented size of predicted structures for data mining & learning**
- ...





# Protein sequence models

- Modeling the distribution of protein sequences through language models or other density models
- Bigger is better??
- Applications
  - Predicts effects of genetic mutations
  - Guides sequence choices in protein engineering





Sequence density model  
(language model)

Protein sequence  $\longrightarrow$

Sequence likelihood

VALYDYEARTED

Neutral variation



I



G

Disease variation

VALYDYEARTED

Neutral variation



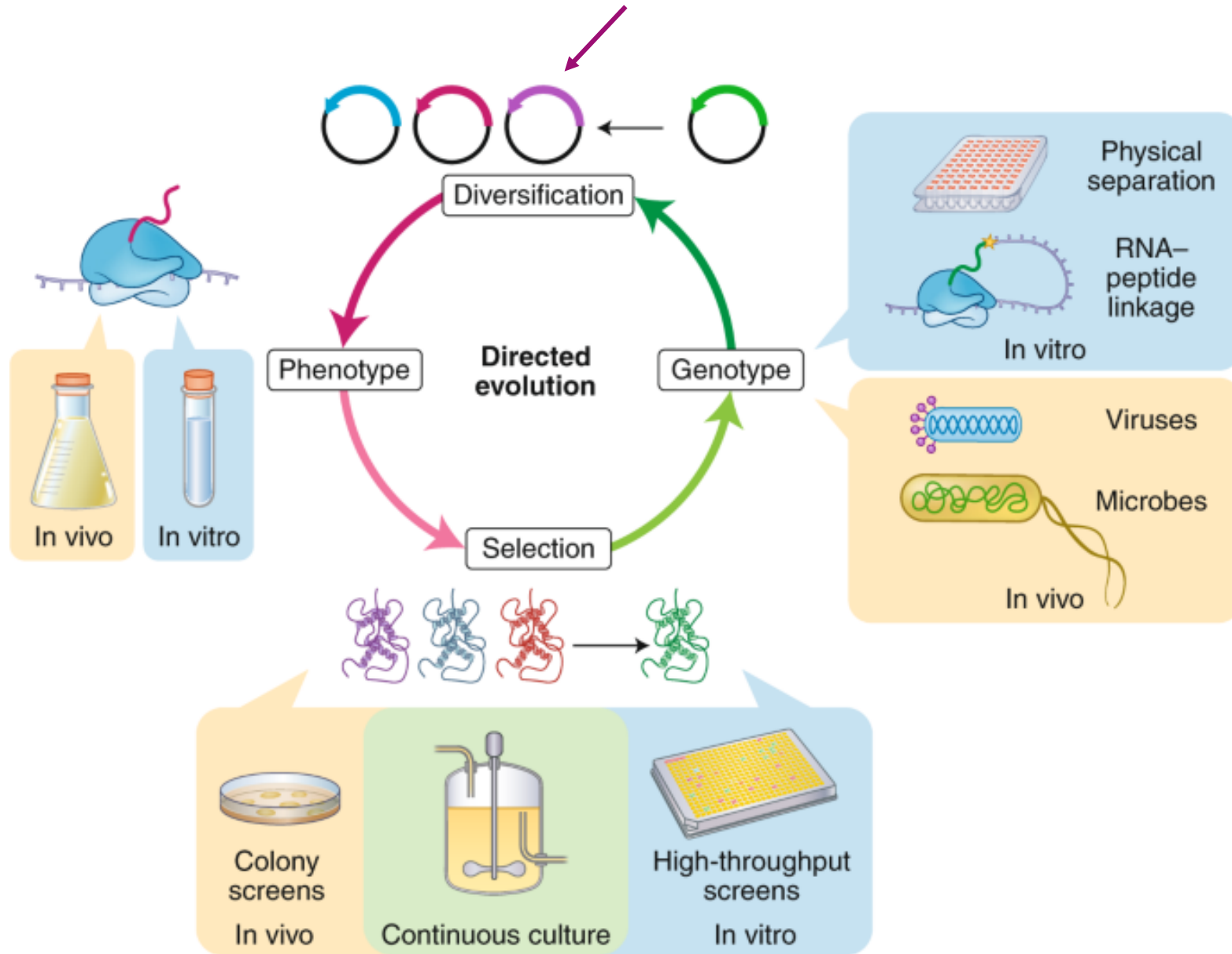
R



W

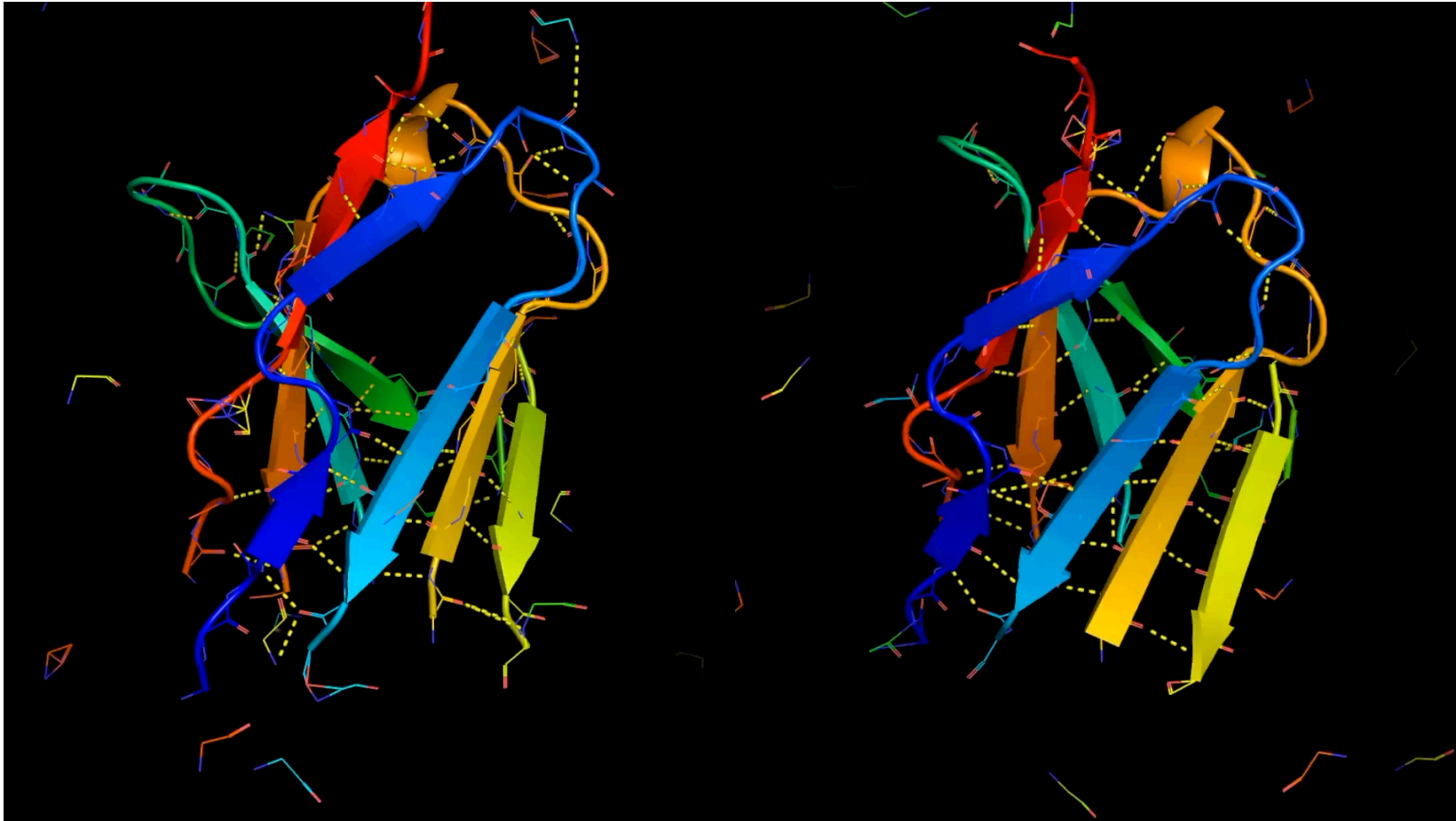
Enhanced variation

# Machine learning-based library design



# Generative models for protein design

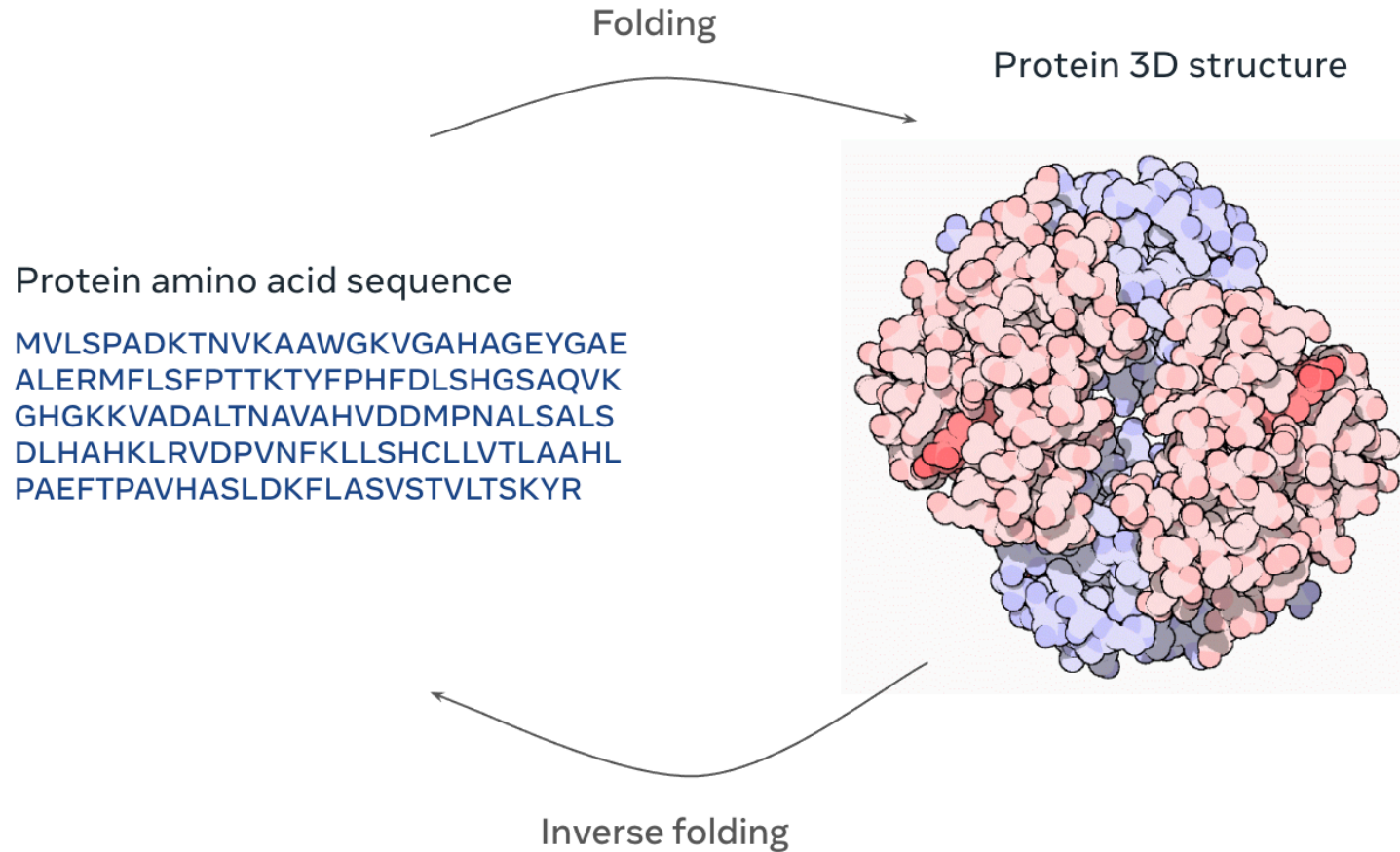
Denoising diffusion probabilistic model for protein structure and sequence



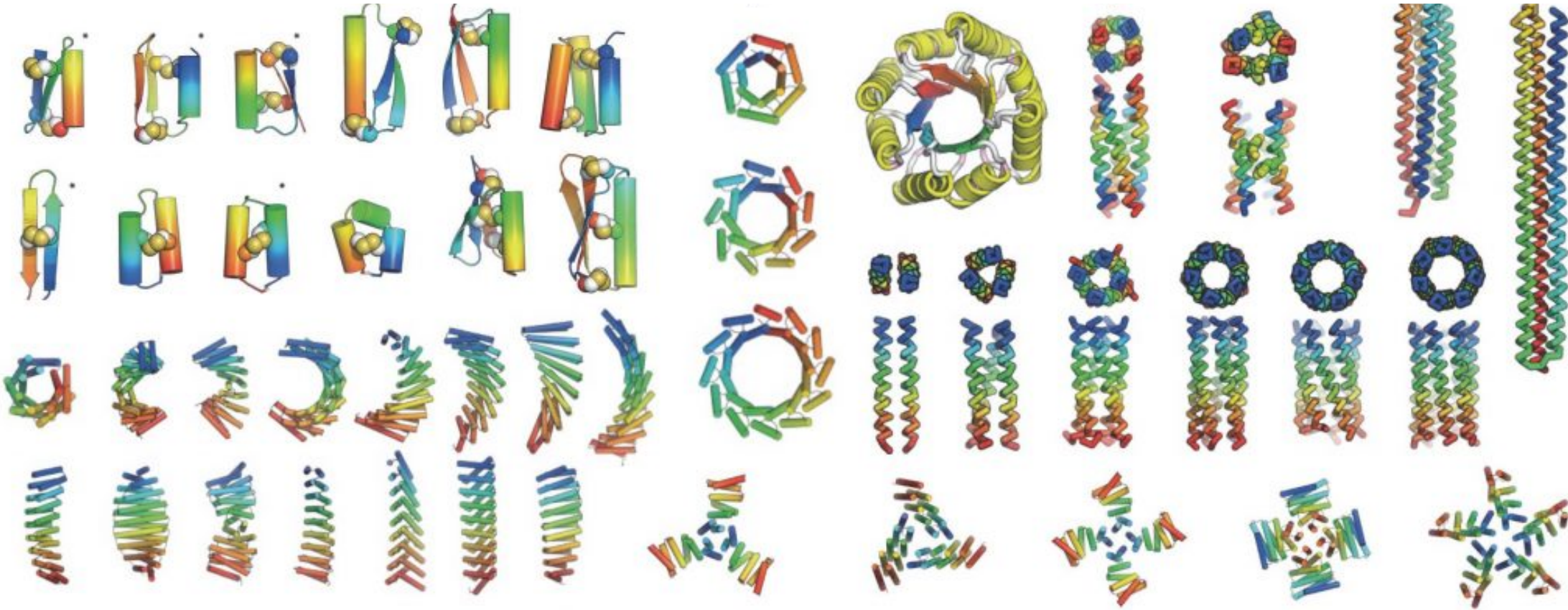
Anand, Namrata, and Tudor Achim. "Protein Structure and Sequence Generation with Equivariant Denoising Diffusion Probabilistic Models." arXiv preprint arXiv:2205.15019 (2022).

# Generative models for protein design

Conditional generative models for protein design, e.g. inverse folding models



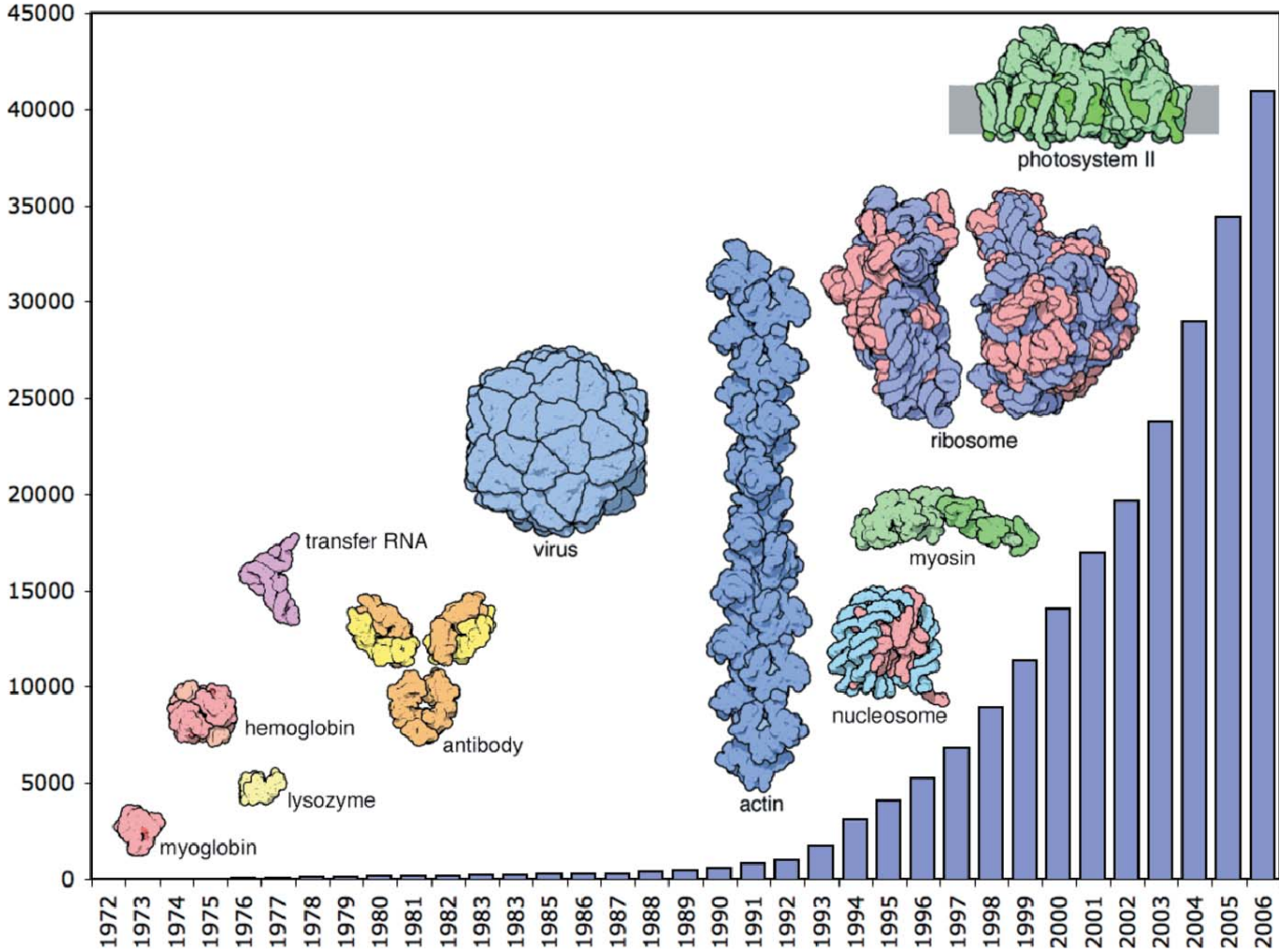
# Beyond naturally existing proteins: De novo proteins



Huang, Po-Ssu, Scott E. Boyken, and David Baker. "The coming of age of de novo protein design." *Nature* 537.7620 (2016): 320-327.

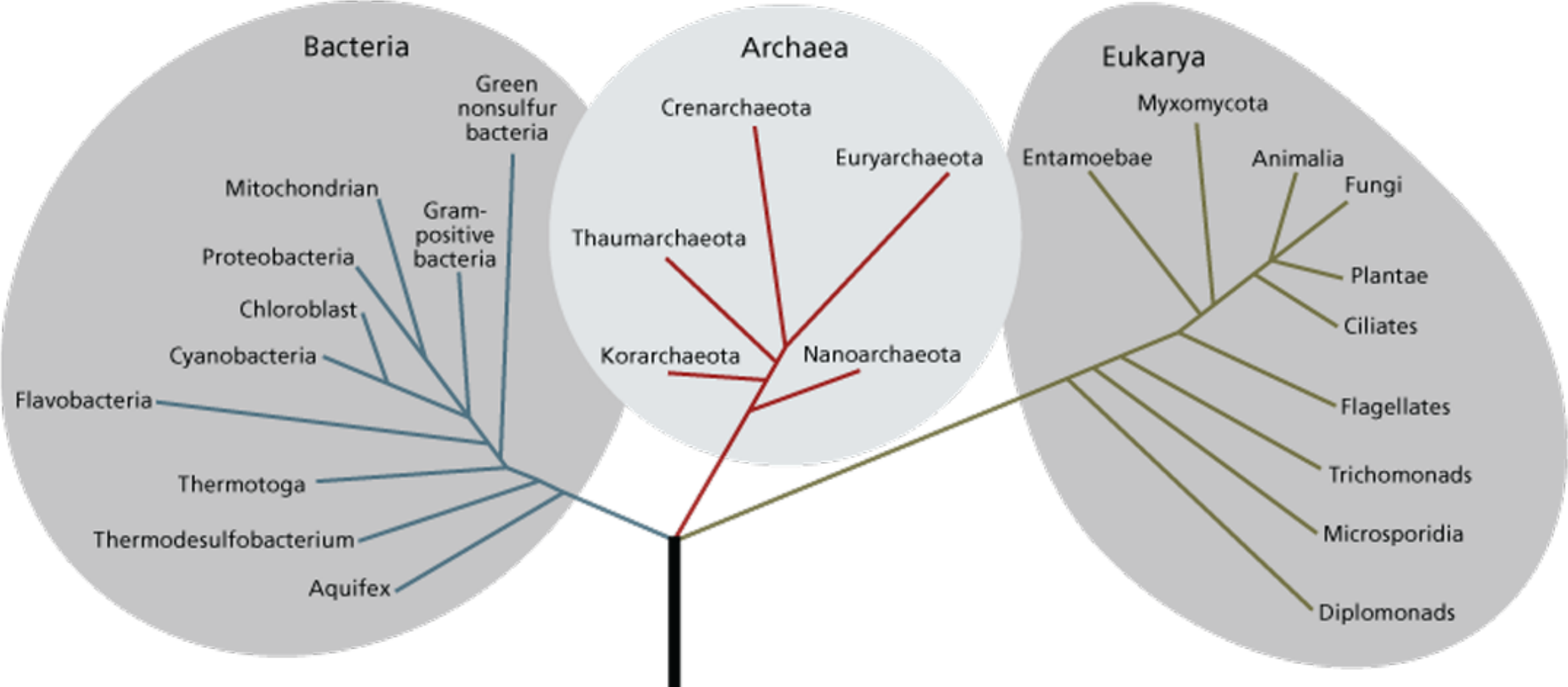
“How”: The data and the learning algorithms

# Protein structures: Protein Data Bank



(Today over 100,000)

# Evolutionary data: observing the products of evolution



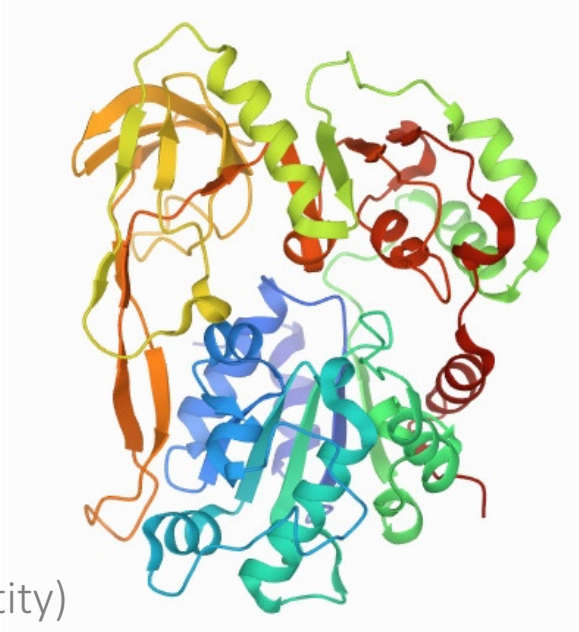


# Evolutionary data: observing the products of evolution

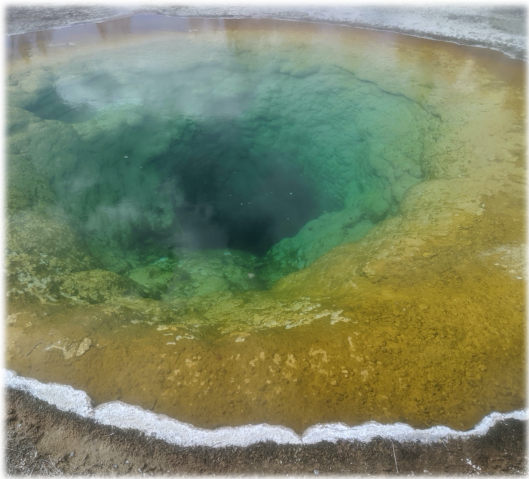
Human PIF1 helicase (PDB: 6HPH)



Thermophilic bacteria PIF1 helicase (PDB: 6S3E)













(~30% sequence identity)



# Evolutionary data: observing the products of evolution

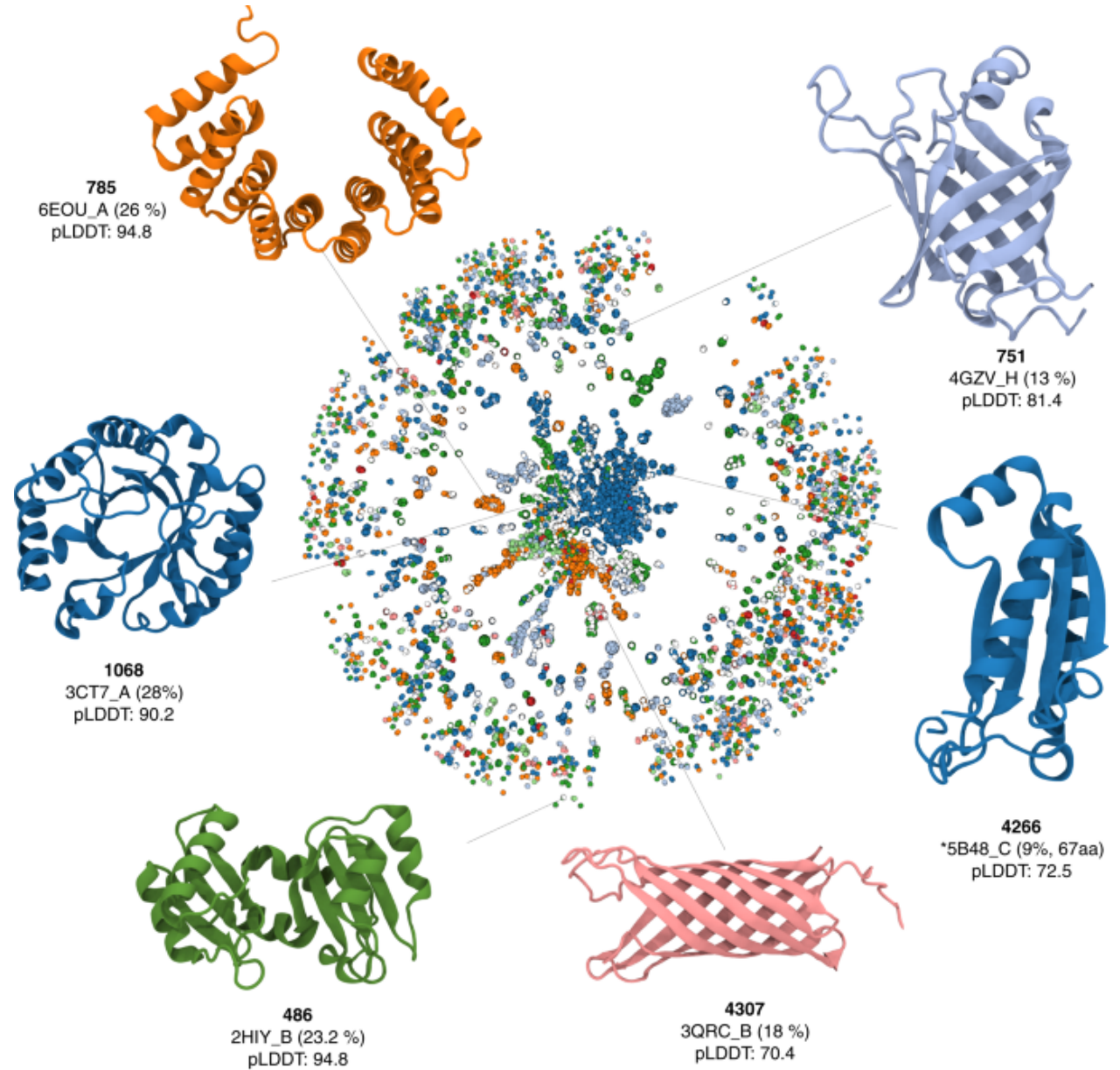
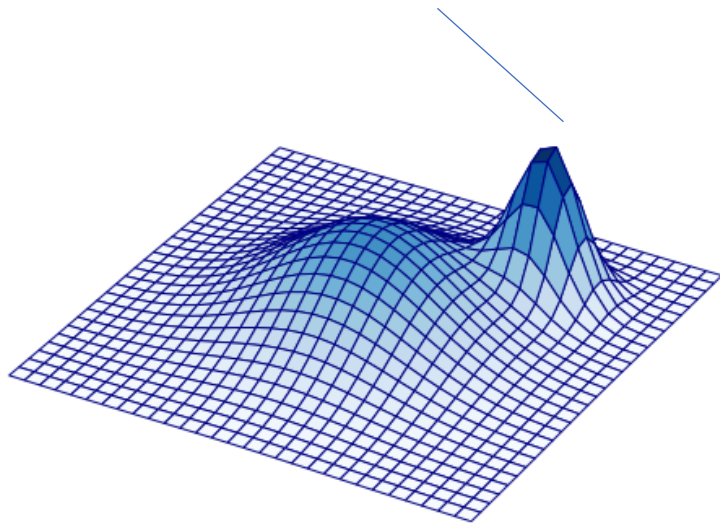
## Multiple sequence alignment (MSA)

PIF1 helicase

	WALRKTRKRLEEPFGGVKVLVLLGDTRQLEPVVPGGEEALYIARTWGGPFFFQAHVWEE--	180
	R R++ +PFGG++++ GD QL PV G + F FQ+ W+	
	AVARAVRQQ-NKPFGGIQLIICGDFLQLPPVTKGSQPP-----RFCFQSKSWKRCV	168
	-VALRVHRLWESQRQREDPLFAELLKRLRQG--DPQALETLNRAAVRPDGGEEPGLTILT	237
	V L + ++W ++ D F LL+ +R G + L A G + L	
	PVTLELTKVW----RQADQTFISLLQAVRLGRCSDEVTRQLQATASHKVGRDGIVATRLC	224
	PRRKEADALNLKRLEALPGKPLEYQAQVKG-EFAET---DFPTEAALTLLKKAQVILLRN	293
	+ + N +RL+ LPGK ++A E A T P L LK GAQV+L++N	
	THQDDVALTNERRLQELPGKVHRFEAMDSNPELASTLDAQCPVSQLLQLKLGAQVMLVKN	284
	DPLGE-YFNGDLGWVEDLEAEALAVRLKR--NGRRVVIRPFVWEKIVYTYDSEREEIKPQ	350
	+ NG G V EAE + R G VI W T + ++ +	
	LSVSRGLVNGARGVVVGFEAEGRGLPQVRFLCGVTEVIHADRW-----TVQATGGQLLSR	339
	VVGTFRQVPVRLAWALTVHKAQGLTLDKVVHLELGRGLFAHGQLYVALTRVRRLLQDL	406
	+Q+P++LAWA+++HK+QG+TLD V + LGR +FA GQ YVAL+R R LQ L	
	-----QQLPLQLAWAMSIHKSQGMTLDCVEISLGR-VFASGQAYVALSRARSLQGL	389

# Evolutionary density model

High density area = likely functional proteins



## Before deep learning: Simple models on evolutionary data



Hidden Markov Models (HMMs) only model the site-specific amino acid features.



Potts models, also known as Markov Random Fields (MRFs), model the site-specific amino acid features and the pairwise interactions between sites.

**They work quite well, too**

(when used on specific protein families in combination with search algorithms)

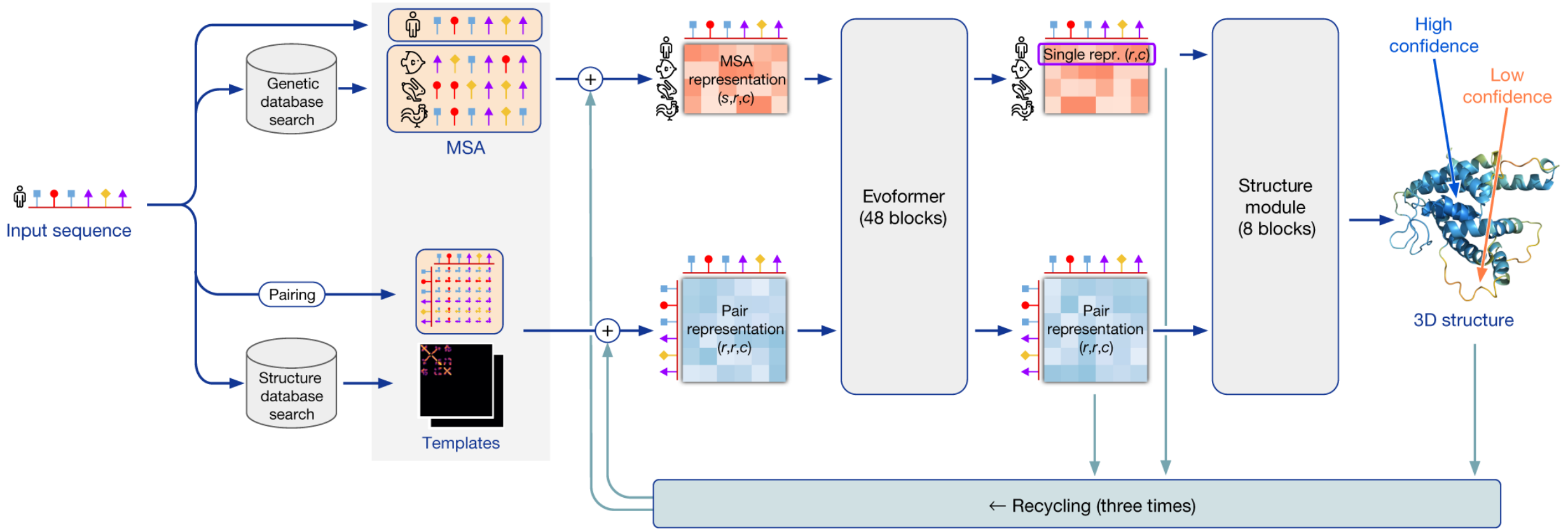
# Replacing multiple sequence alignments with deep learning models

(Protein language models)

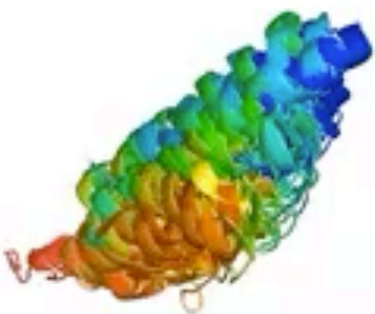


New ability to jointly model all protein families regardless of “alignments”

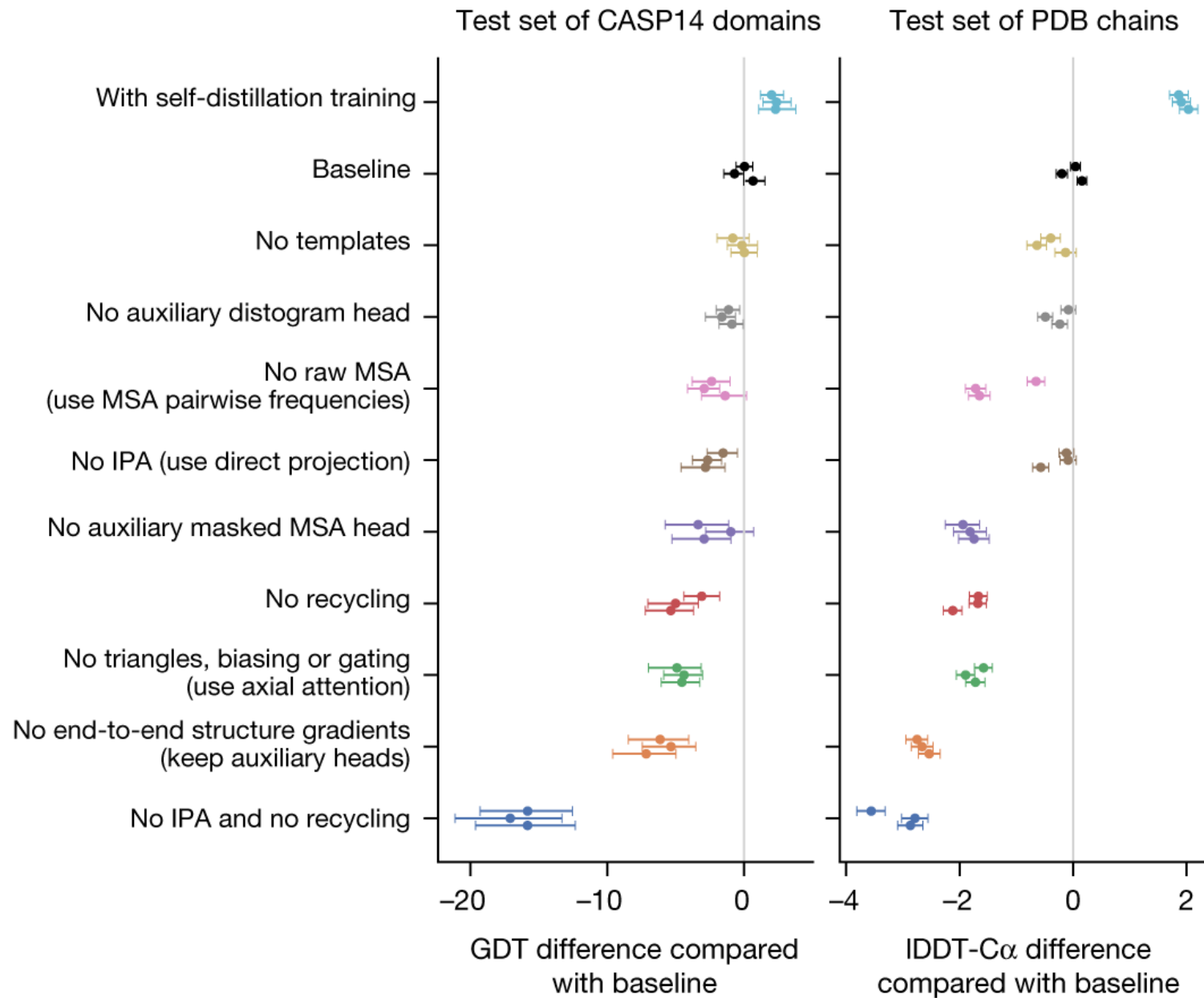
# AlphaFold2 also relies on evolutionary data (multiple sequence alignments)



Jumper, J., Evans, R., Pritzel, A. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).



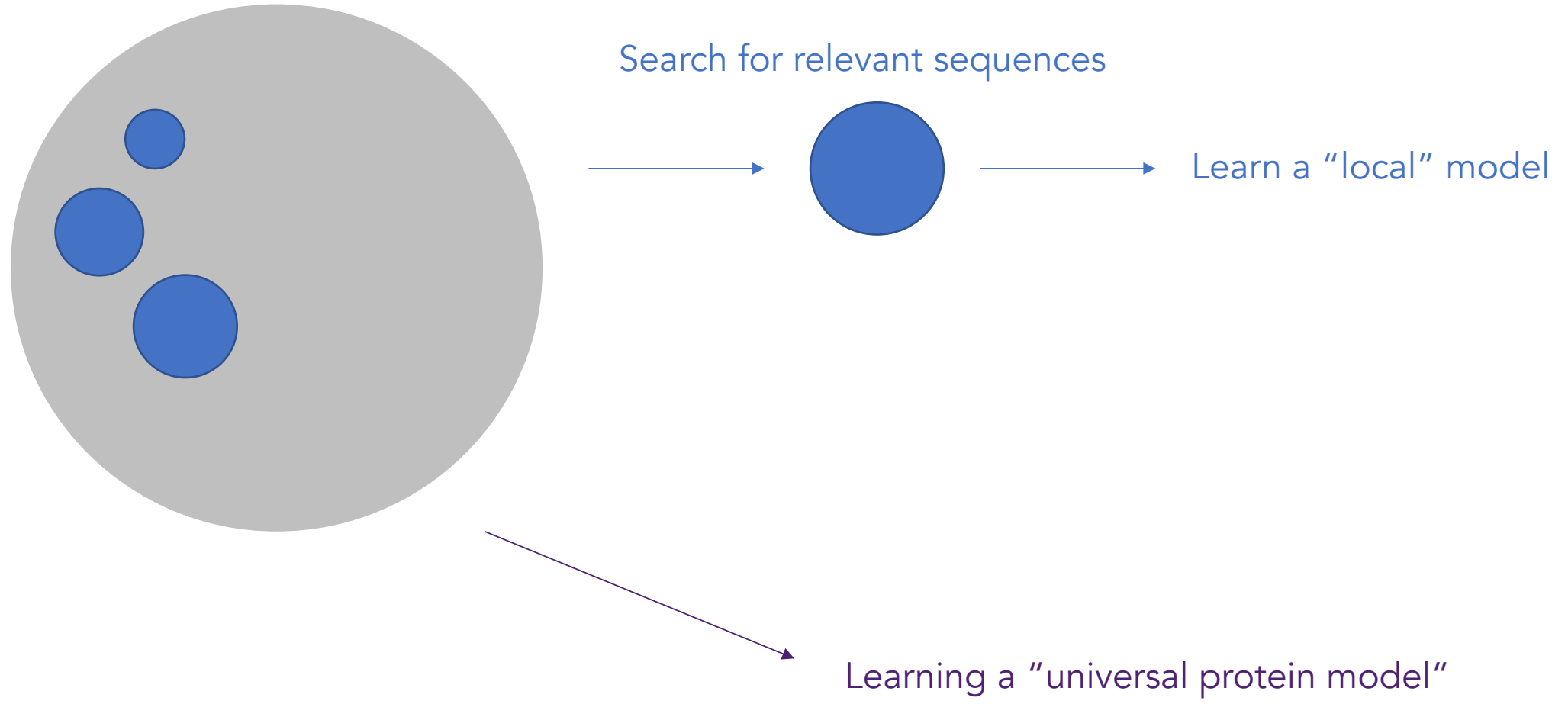
Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction





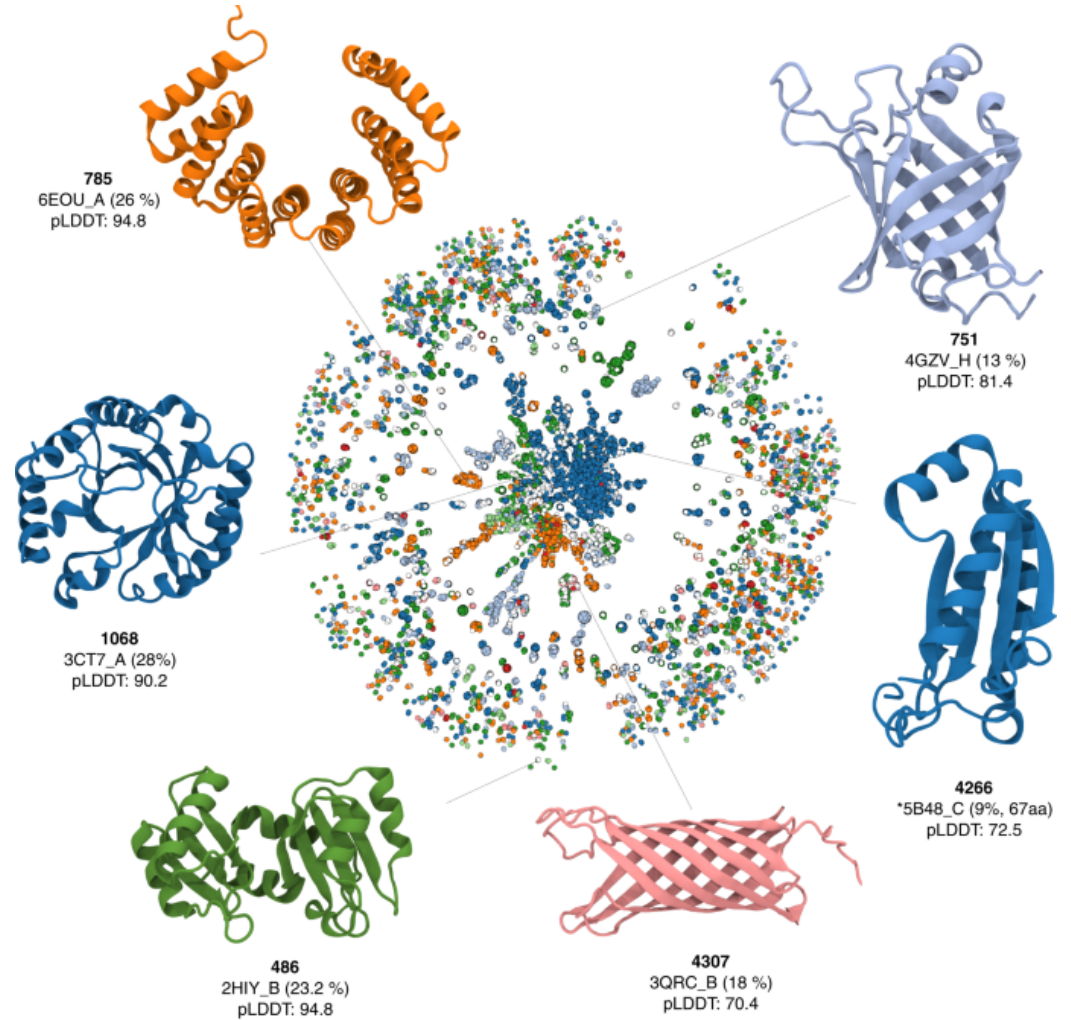
“Why”: Common themes and open questions

# "Search & learn" vs. "universal learning"



# “Generalization” vs “memorization”

Are language models / single-sequence folding models effectively memorizing the evolutionary database?



Ferruz, N., Schmidt, S. & Höcker, B. ProtGPT2 is a deep unsupervised language model for protein design. *Nat Commun* **13**, 4348 (2022).