

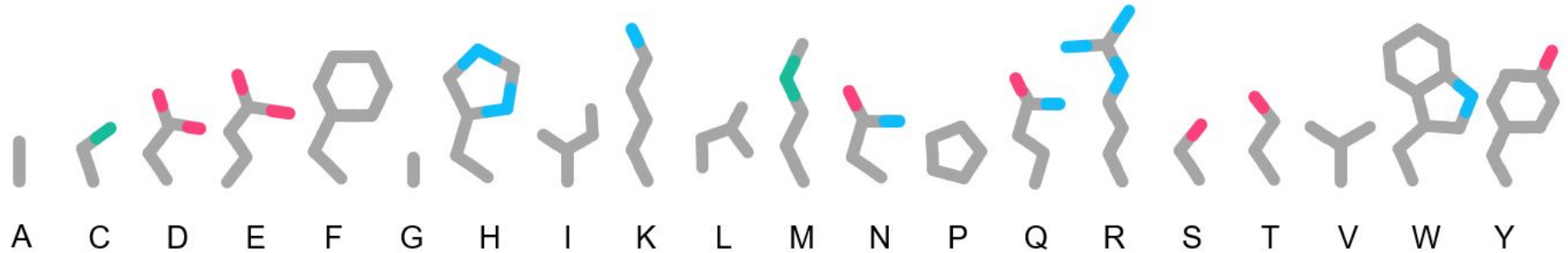
# CSE7850/CX4803 Machine Learning in Computational Biology



## Lecture 15: Learning from Structure Data

Yunan Luo

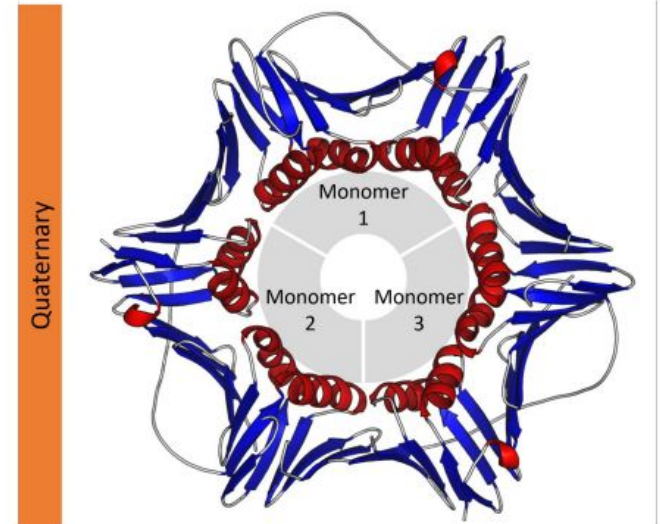
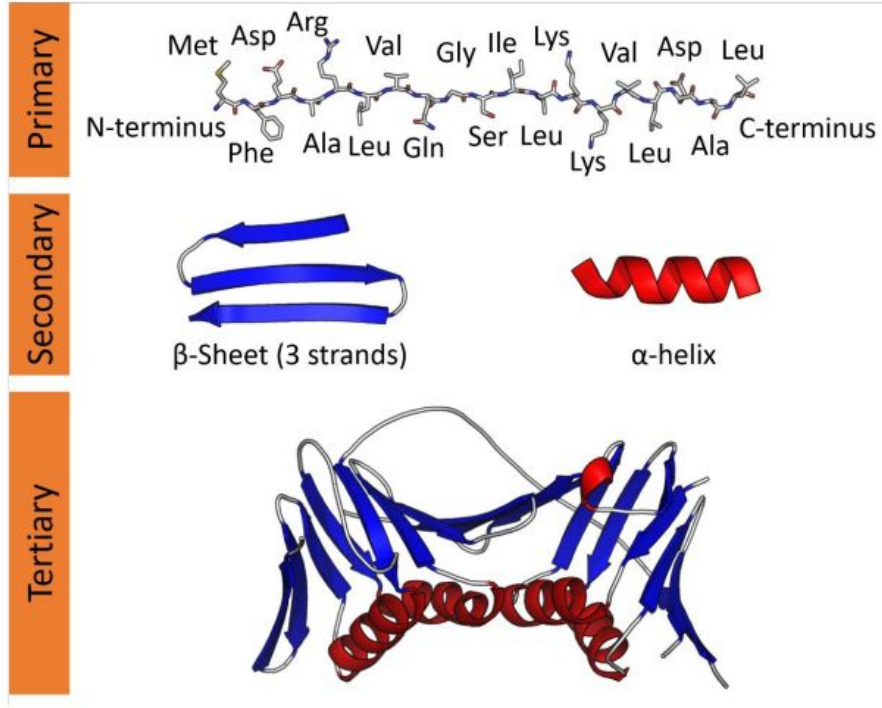
# Amino acids are the building blocks of proteins



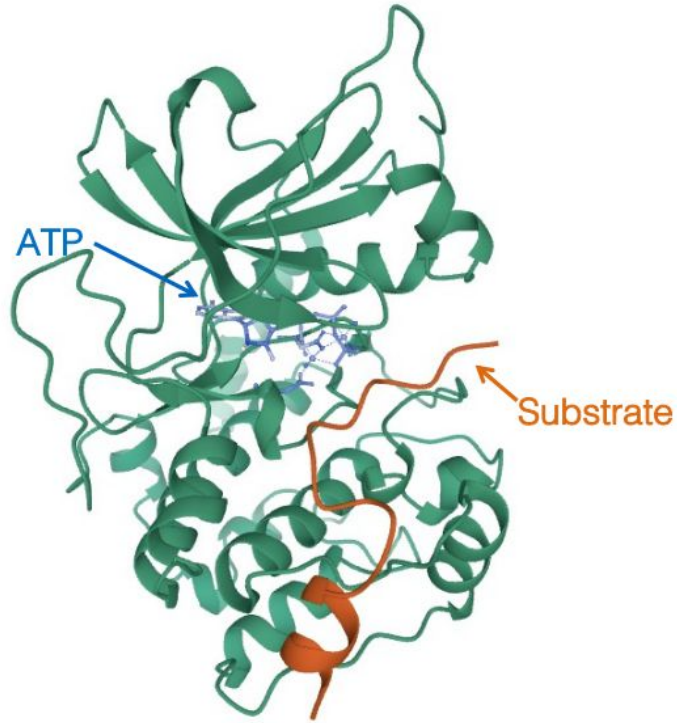
Amino acids vary in

- Size
- Shape
- Polarity
- Charge

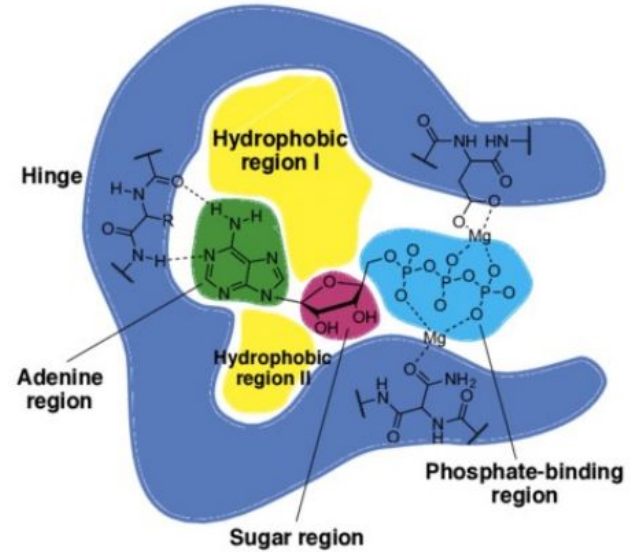
# Protein structure hierarchy



# Structure provides insight on function

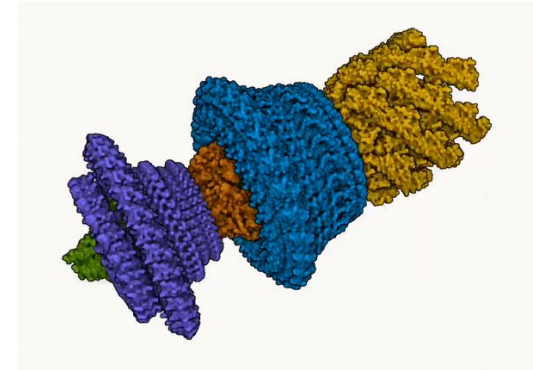
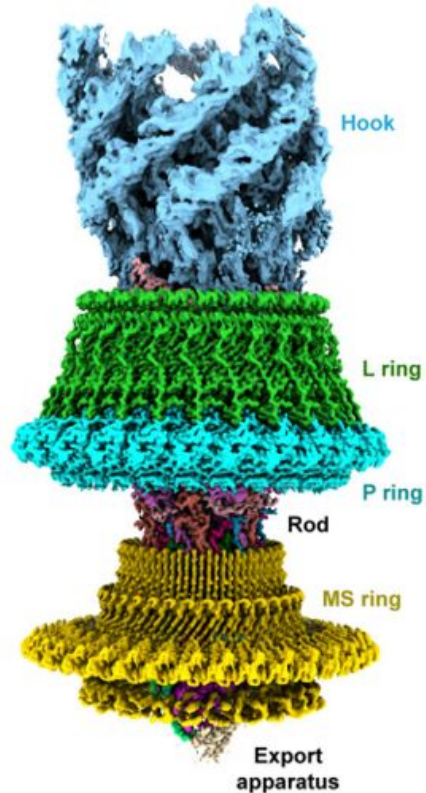


Protein Kinase



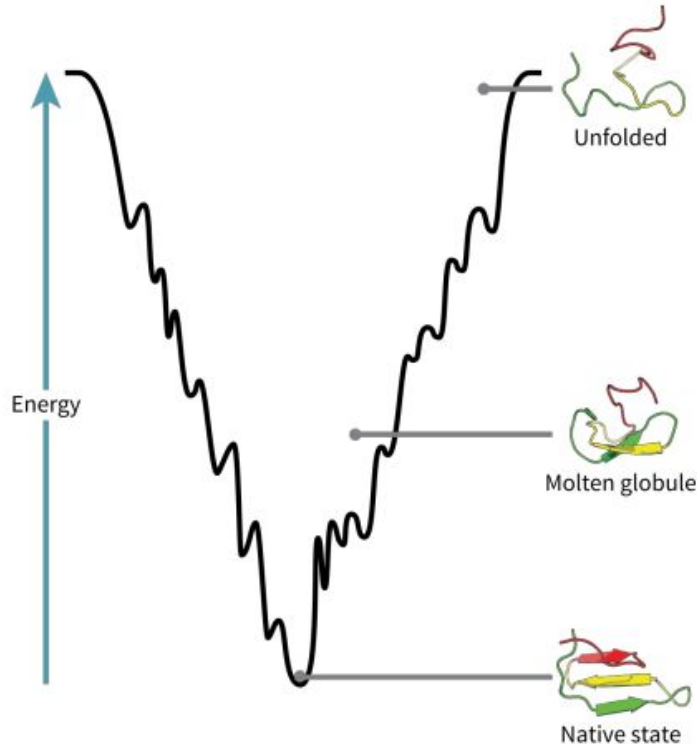


# Structure provides insight on function



Source: [Twitter](#)

# The protein folding problem



- The function of a protein is determined in large part by its 3D shape
- Can we predict the 3D structure of a protein given only its (1D) amino-acid sequence?

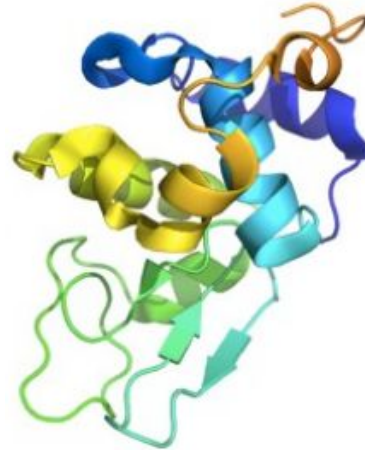
# Protein structure prediction

Amino acid sequence

MEKVNFLKNGVLRLLPPGFRFRPTDEELVVQYLKRKVFSPPLPASIPEVEVYKSDPWDLPGDMEQEKYFFSTK  
EVKYPNGNRSNRATNSGYWKATGIDKQIILRGRQQQQLIGLKKTLVFYRGKSPHGCRTNWIMHEYRLAN  
LESNYHPIQGNWVICRIFLKKRGNTKNKEENMTTHDEVNRNREIDKNPVSVMSSRDSEALASANSELKK



Algorithm / Model



Protein structure

# Classical approaches for protein structure prediction

- Homology modeling
- Fold recognition (threading)
- Fragment assembly
- Molecular dynamics



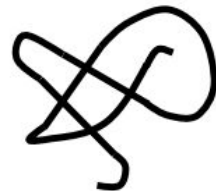
# Classical approaches for protein structure prediction

- **Homology modeling**

- given a query sequence **Q**, a database of protein structures, do:
  - find protein **P** such that
    - structure of **P** is known
    - **P** has high sequence similarity to **Q**
  - return **P**'s structure as an approximation to **Q**'s structure

- Fold recognition (threading)
- Fragment assembly
- Molecular dynamics

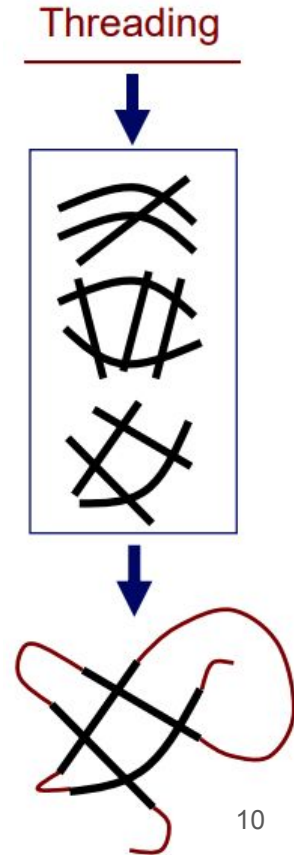
Homology  
modeling



# Classical approaches for protein structure prediction

- Homology modeling
- **Fold recognition (threading)**
  - given a query sequence **Q**, a database of known folds, do:
    - find fold **F** such that **Q** can be aligned with **F** in a highly compatible manner
    - return **F** as an approximation to **Q**'s structure
- Fragment assembly
- Molecular dynamics

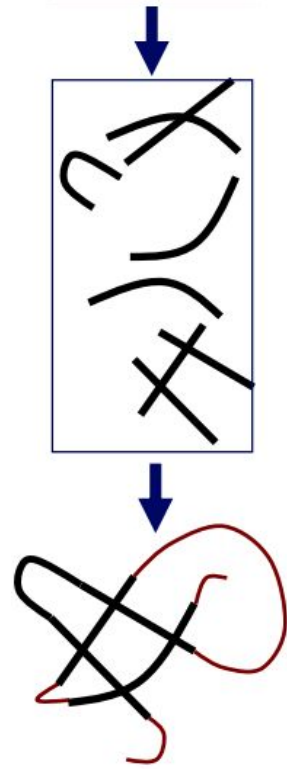
**Fold:** a description of the relative orientation of the secondary structure making up the tertiary structure.



# Classical approaches for protein structure prediction

- Homology modeling
- Fold recognition (threading)
- **Fragment assembly (e.g., Rosetta)**
  - given a query sequence **Q**, a database of structure fragments, do
    - find a set of **fragments** that **Q** can be aligned with in a highly compatible manner
    - return fragment assembly as an approximation to **Q**'s structure
- Molecular dynamics

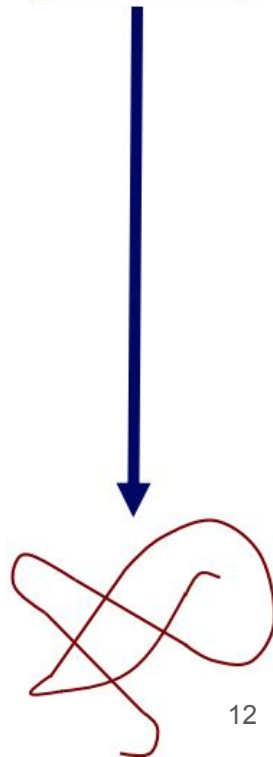
Fragment assembly  
(Rosetta)



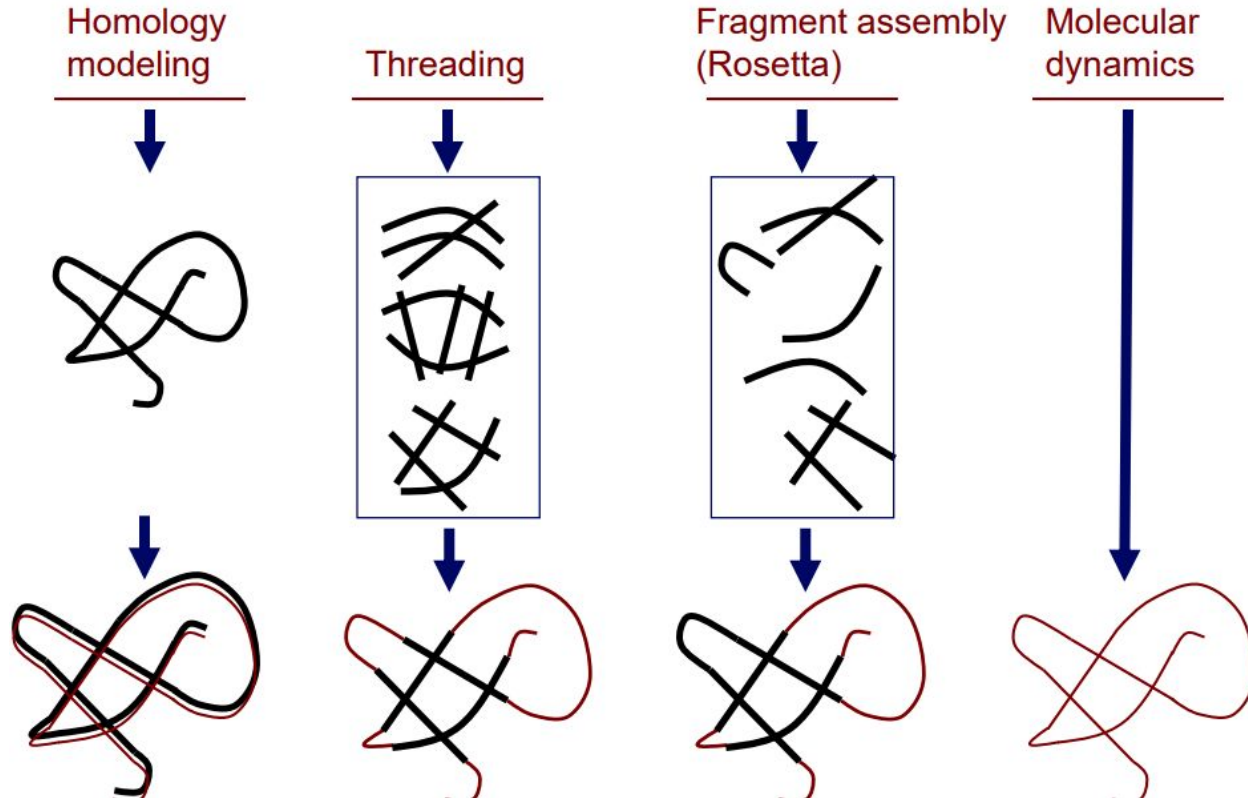
# Classical approaches for protein structure prediction

- Homology modeling
- Fold recognition (threading)
- Fragment assembly (e.g., Rosetta)
- **Molecular dynamics**
  - given a query sequence Q
  - do: use laws of physics to simulate folding of Q

Molecular  
dynamics

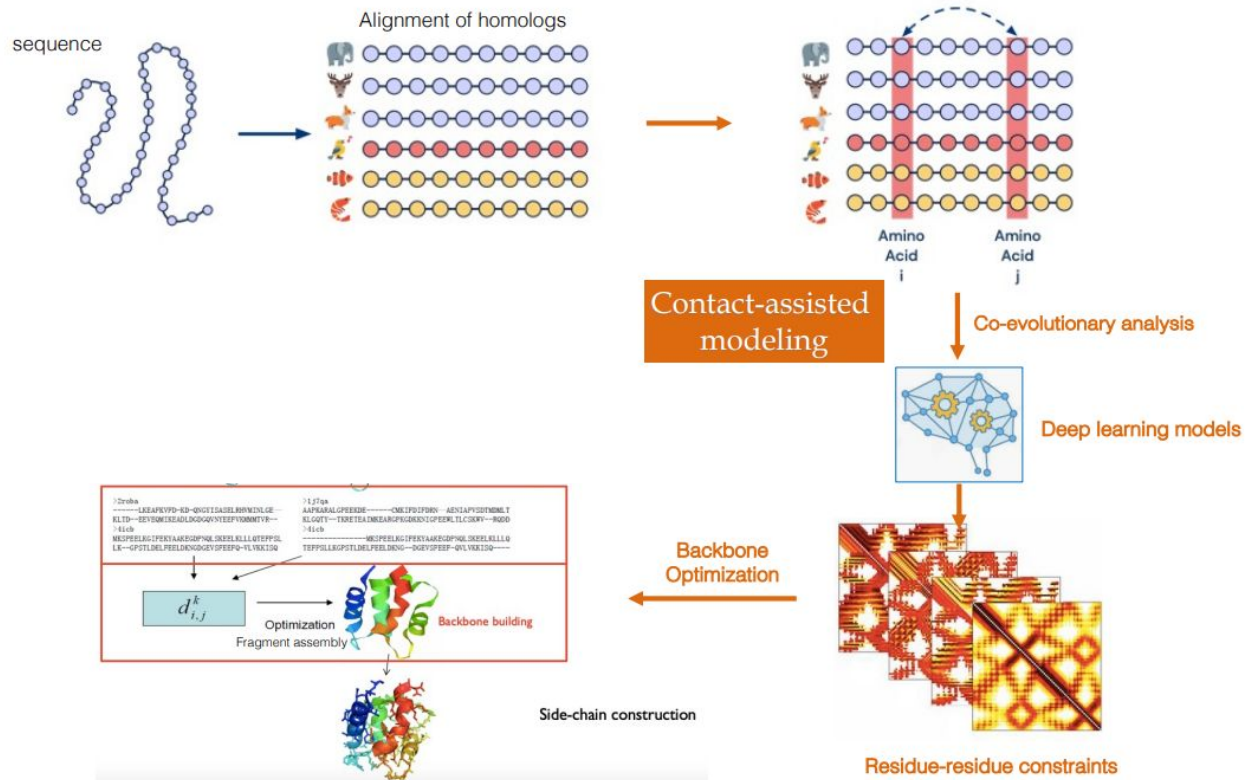


# Classical approaches for protein structure prediction

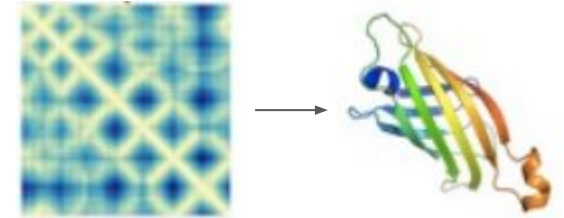
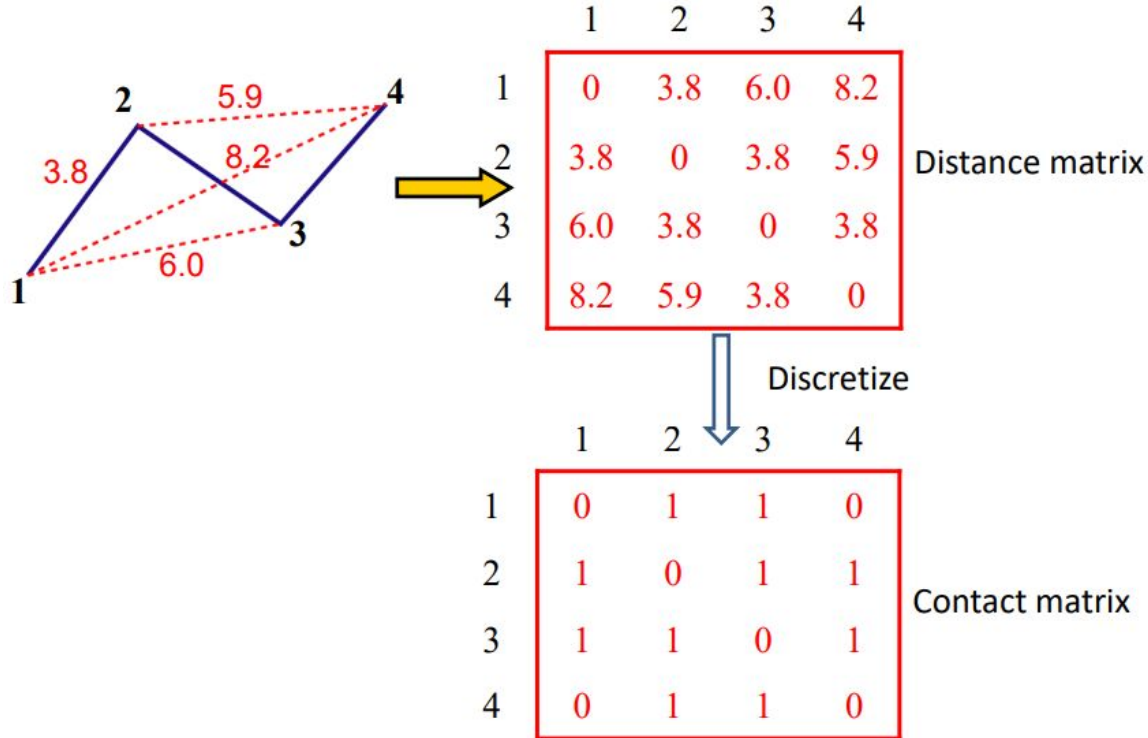




# New approach: Deep Learning

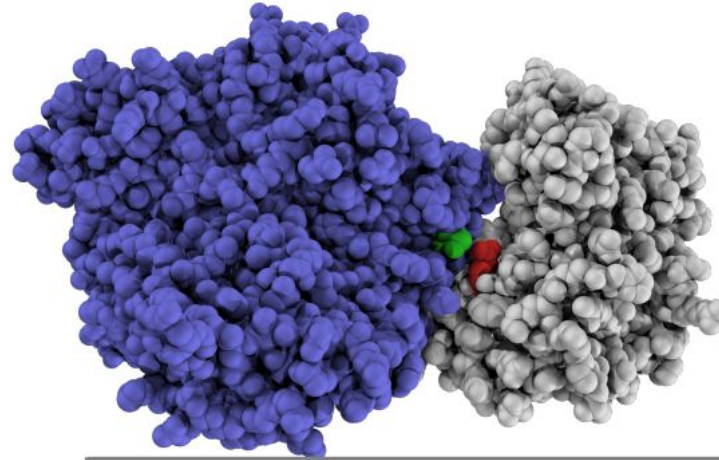


# Protein Distance & Contact Matrix

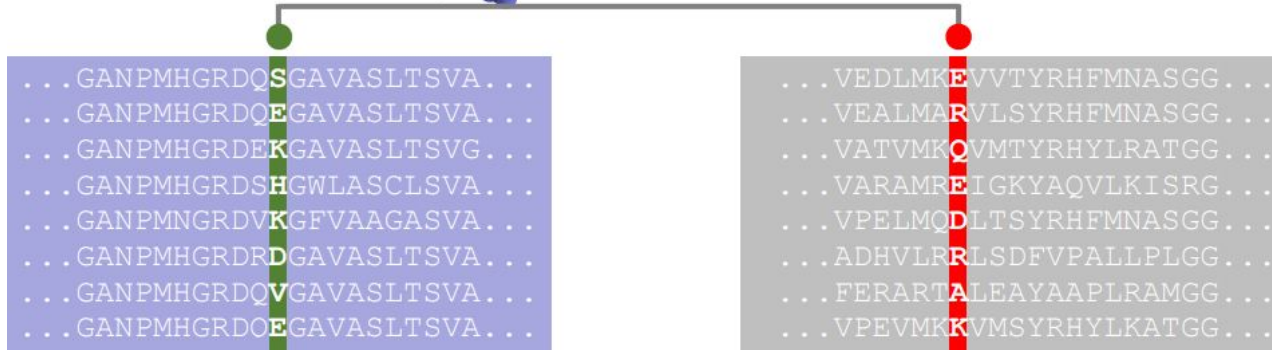


Liu, Palmedo, et al. Cell Systems 2018

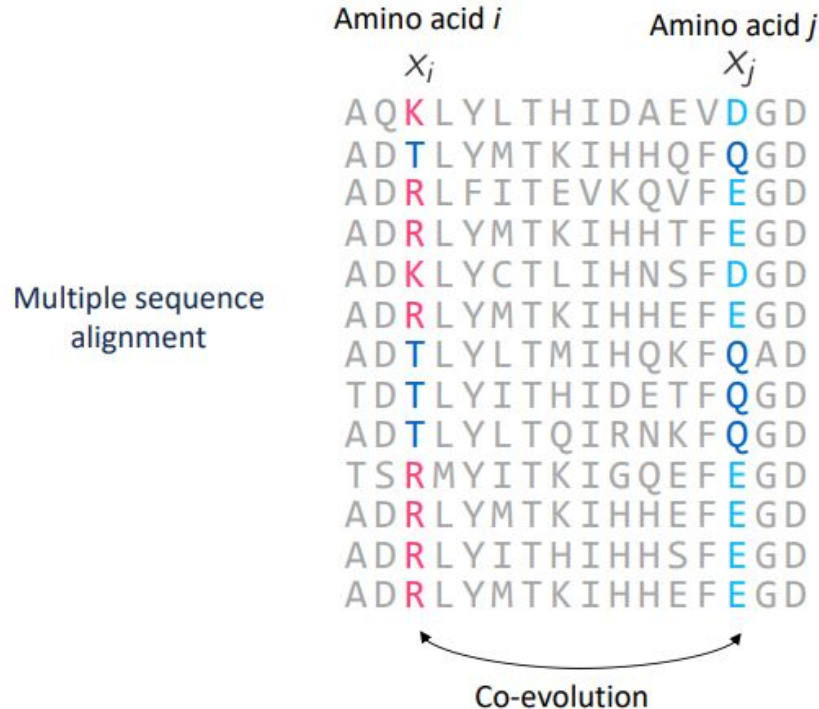
# Amino acids in direct physical contact tend to covary or “coevolve” across related proteins



For example, a mutation that causes one amino acid to get bigger is more likely to preserve protein structure and function (and thus survive) if another amino acid gets smaller to make space



# Learning co-evolution from multiple sequence alignment



$$P(\mathbf{x}) = \frac{1}{Z} \exp \left( \sum_i e_i(x_i) + \sum_{i \neq j} e_{ij}(x_i, x_j) \right)$$

Single  
potentials

Pairwise  
potentials

Local preference

Co-evolution strength

Markov random field  
Ising (Potts) model  
Undirected graphical model

# Learning with Markov Random Fields (MRF)

$$L(e) = \prod_{n=1}^N \frac{1}{Z_e^{(n)}} \prod_i^L \exp[e_i(x_i^n) + \sum_{j \neq i} e_{i,j}(x_i^n, x_j^n)]$$

Partition  
function

Singleton  
potentials

Pairwise  
potentials

Local AA preference

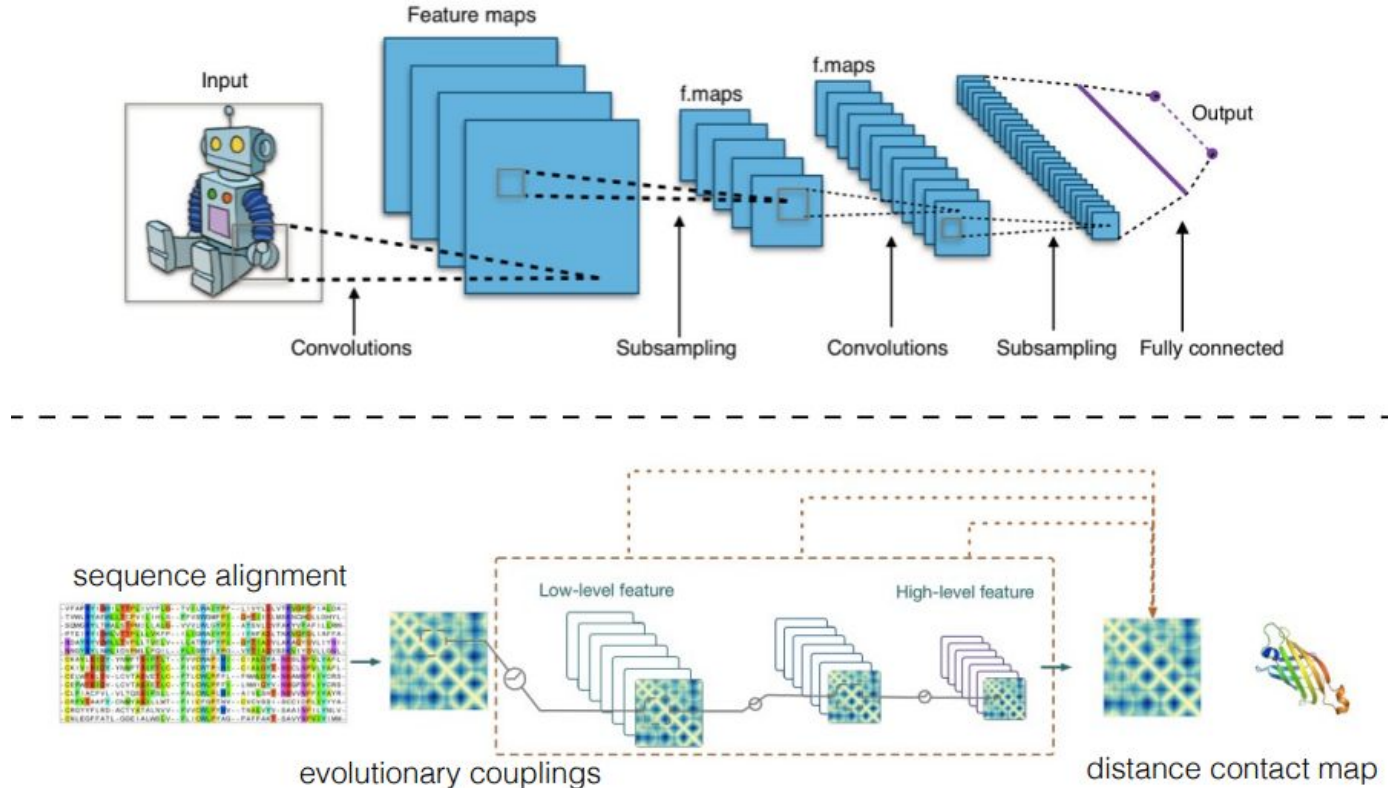
Pairwise AA couplings

Learning algorithms:

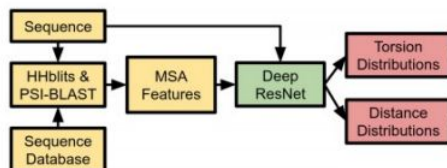
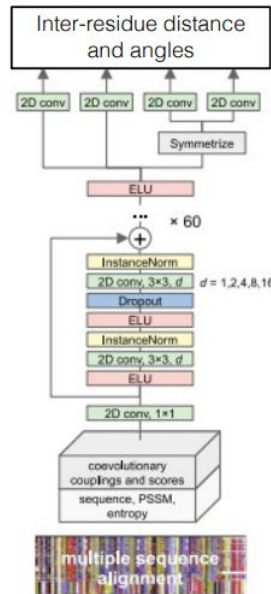
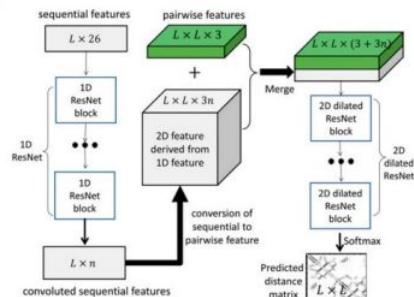
- Mean fields approximation: EVFold, DirectInfo
- Gaussian approximation: PSICOV
- Pseudolikelihood: GREMLIN, CCMpred



# Deep convolutional NNs recognize coevolutionary patterns



# Recent developments go beyond contact prediction



## Improved protein structure prediction using predicted interresidue orientations

PNAS, 2020

Jianyi Yang<sup>a,b</sup>, Ivan Anishchenko<sup>b,c,1</sup>, Hahnbeom Park<sup>b,c</sup>, Zhenling Peng<sup>d</sup>, Sergey Ovchinnikov<sup>a</sup>, and David Baker<sup>b,c,1,2</sup>

## Improved protein structure prediction using potentials from deep learning

<https://doi.org/10.1038/s41586-019-1923-7>

Received: 2 April 2019

Accepted: 10 December 2019

Andrew W. Senior<sup>1,2\*</sup>, Richard Evans<sup>1,2</sup>, John Jumper<sup>1,2</sup>, James Kirkpatrick<sup>1,2</sup>, Laurent Sifre<sup>1,2</sup>, Tim Green<sup>1</sup>, Chongli Qin<sup>1</sup>, Augustin Zidek<sup>1</sup>, Alexander W. R. Nelson<sup>1</sup>, Alex Bridgland<sup>1</sup>, Hugo Penedones<sup>1</sup>, Stig Petersen<sup>1</sup>, Karen Simonyan<sup>1</sup>, Steve Crossan<sup>1</sup>, Pushmeet Kohli<sup>1</sup>, David T. Jones<sup>1,2</sup>, David Silver<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup> & Demis Hassabis<sup>1</sup>

Nature, 2020  
(AlphaFold 1)

## Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

Check for updates

John Jumper<sup>1,2,3,4\*</sup>, Richard Evans<sup>1,2</sup>, Alexander Pritzel<sup>1,2</sup>, Tim Green<sup>1,2</sup>, Michael Figurnov<sup>1,2</sup>, Olaf Ronneberger<sup>1,2</sup>, Kathryn Tunyasuvunakool<sup>1,2</sup>, Russ Bates<sup>1,2</sup>, Augustin Zidek<sup>1,2</sup>, Anna Potapenko<sup>1,2</sup>, Alex Bridgland<sup>1,2</sup>, Clemens Meyer<sup>1,2</sup>, Simon A. A. Kohl<sup>1,2</sup>, Andrew J. Ballard<sup>1,2</sup>, Andrew Cowie<sup>1,2</sup>, Bernardino Romera-Paredes<sup>1,2</sup>, Stanislaw Nikolov<sup>1,2</sup>, Rishub Jain<sup>1,2</sup>, Jonas Adler<sup>1,2</sup>, Trevor Back<sup>1,2</sup>, Stig Petersen<sup>1,2</sup>, David Reiman<sup>1,2</sup>, Ellen Clancy<sup>1,2</sup>, Michal Zielinski<sup>1,2</sup>, Martin Steinegger<sup>1,2</sup>, Michalina Pacholska<sup>1,2</sup>, Tamas Berghammer<sup>1,2</sup>, Sebastian Bodenstein<sup>1,2</sup>, David Silver<sup>1,2</sup>, Oriol Vinyals<sup>1,2</sup>, Andrew W. Senior<sup>1,2</sup>, Koray Kavukcuoglu<sup>1,2</sup>, Pushmeet Kohli<sup>1,2</sup> & Demis Hassabis<sup>1,2,3,4</sup>

Nature, 2021  
(AlphaFold 2)

RESEARCH ARTICLE | PROTEIN FOLDING

f t in

## Accurate prediction of protein structures and interactions using a three-track neural network

Science, 2021

MINKYUNG BAEK<sup>1</sup>, FRANK DIMAGIO<sup>1</sup>, IVAN ANISHCHENKO<sup>1</sup>, JUSTAS DAUPARAS<sup>1</sup>, SERGEY OVCHINNIKOV<sup>1</sup>, SYU RIE LEE<sup>1</sup>, JUE WANG<sup>1</sup>, QIAN CONG<sup>1</sup>

LISA N. KINCH<sup>1</sup>, DAVID BAKER<sup>1</sup>

+23 authors

Authors Info & Affiliations

# Critical Assessment of protein Structure Prediction (CASP)

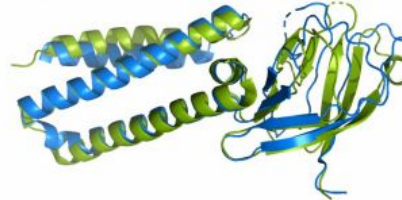
- Since 1994, every two years a contest is held to see who can best predict protein structures from peptide sequences
- Targets structures are held from publication until results are in

# CASP14: DeepMind's AlphaFold 2



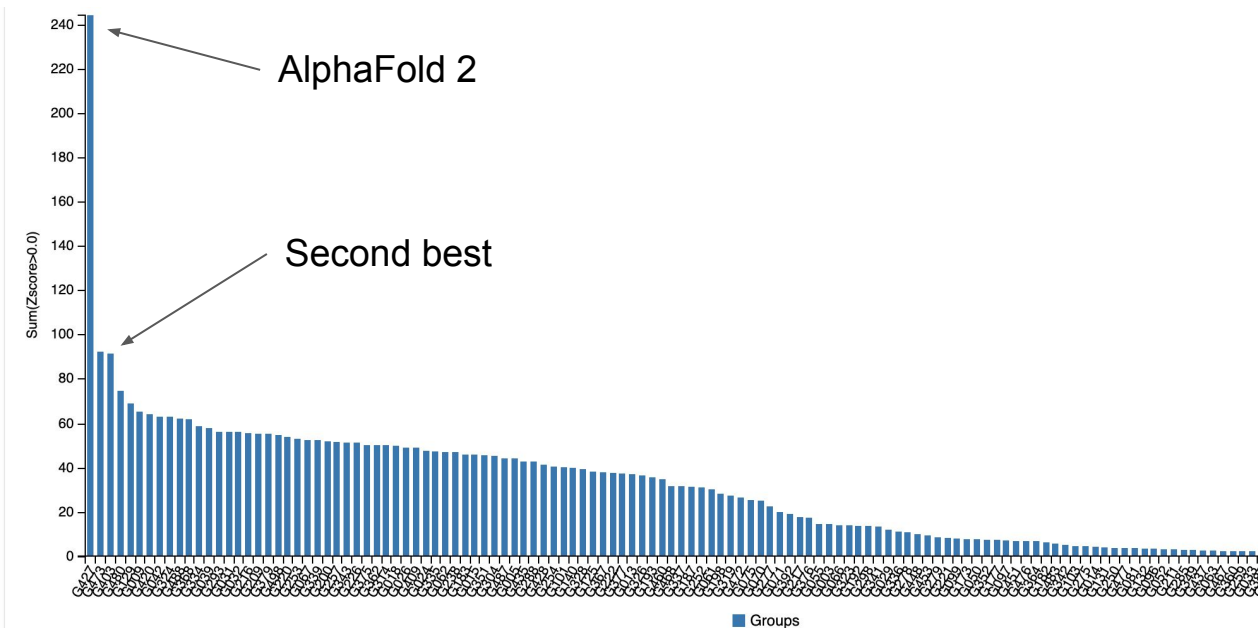
Blue: Predicted  
Green: Actual

ORF8



ORF3a

# CASP14: DeepMind's AlphaFold 2



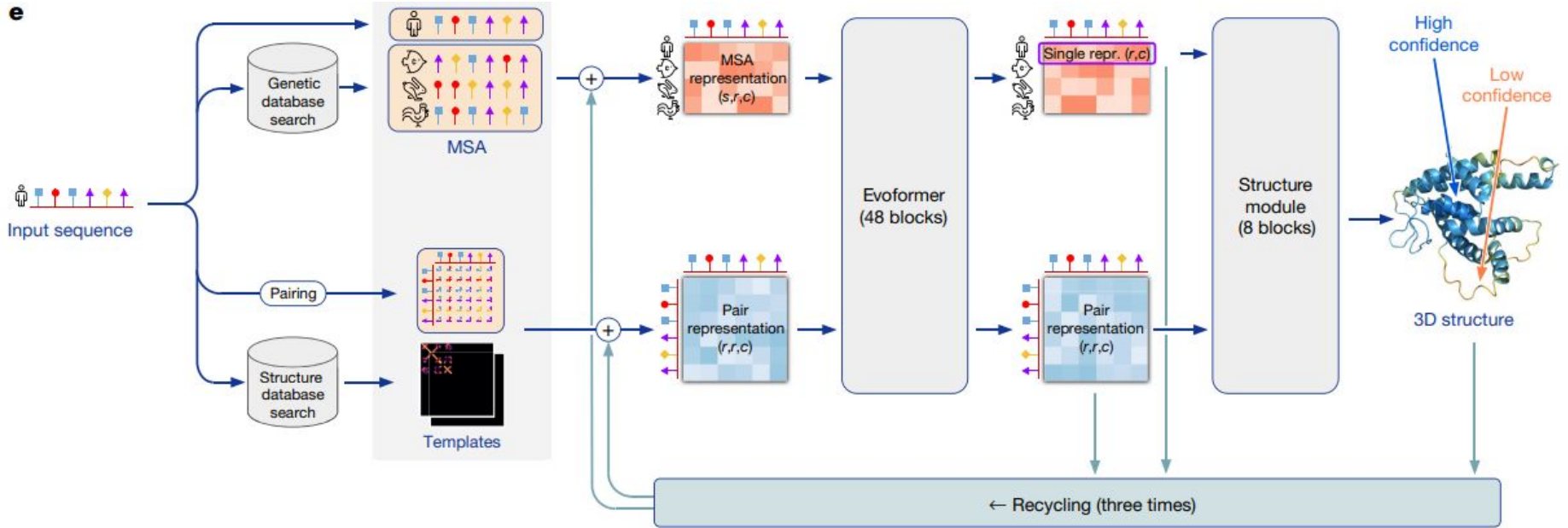
#	GR code	GR name	Domains Count	SUM Zscore (>-2.0)	Rank SUM Zscore (>-2.0)	AVG Zscore (>-2.0)	Rank AVG Zscore (>-2.0)	SUM Zscore (>0.0)	Rank SUM Zscore (>0.0)
1	427	AlphaFold2	92	244.0217	1	2.6524	1	244.0217	1
2	473	BAKER	92	90.8241	2	0.9872	2	92.1241	2
3	403	BAKER-experimental	92	88.9672	3	0.9670	3	91.4731	3
4	480	FEIG-R2	92	72.5351	4	0.7884	4	74.5627	4



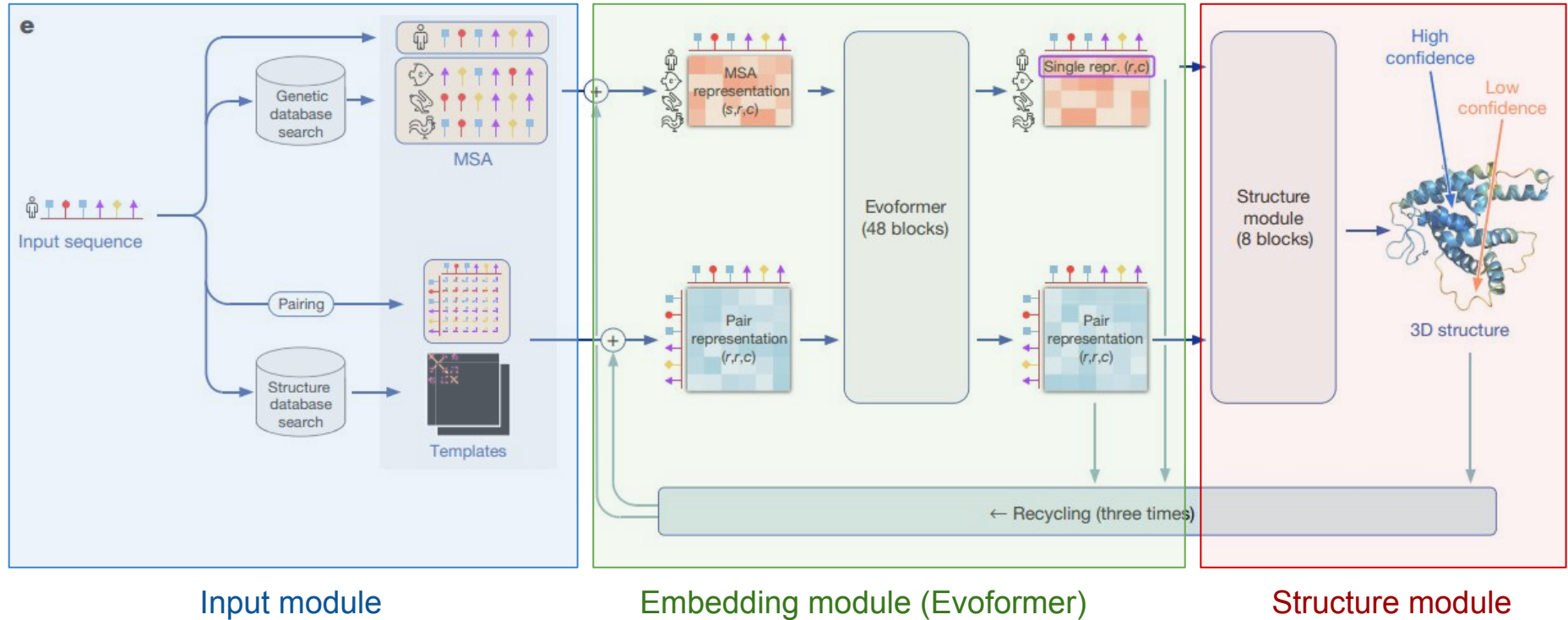
# AlphaFold

# AlphaFold model

e



# AlphaFold model

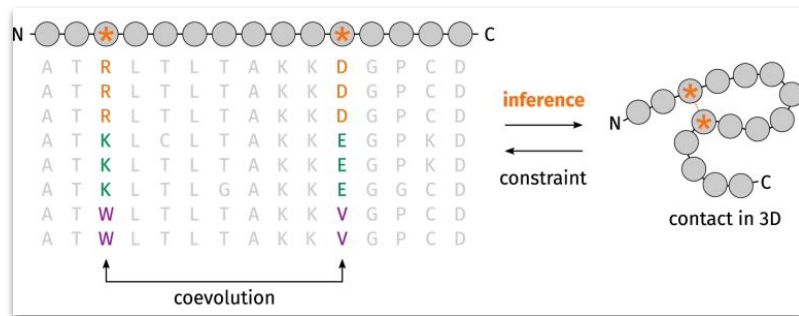


# Input module

Two types of input

- Multiple sequence alignment (MSA)

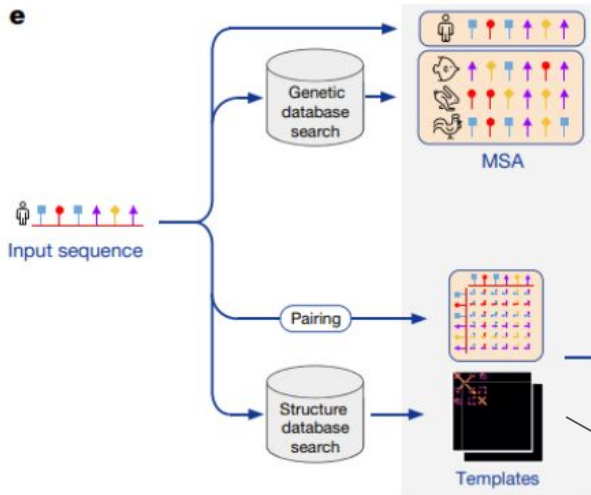
- Residues in contact tend to coevolve



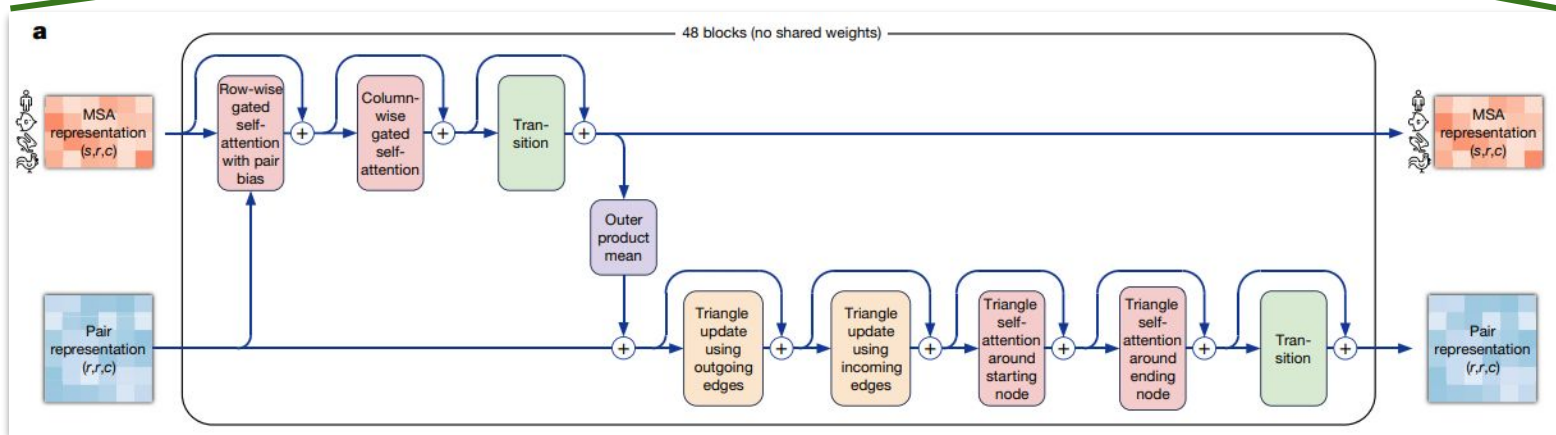
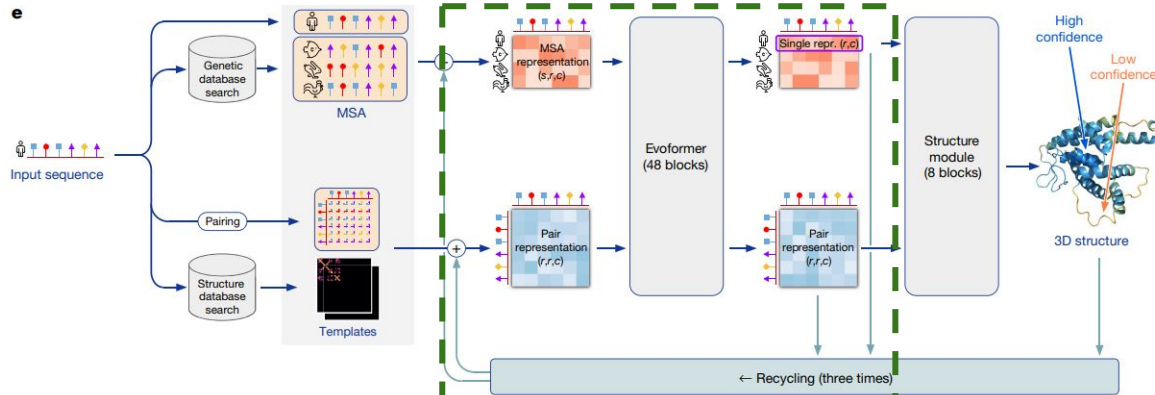
[\(Image source\)](#)

- Template structure

- Structures are more conserved than sequence
- Use conserved fragment to guide structure prediction



# Evoformer



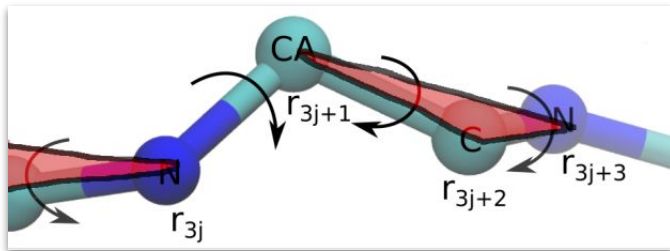
**Key idea:** learn **MSA Representations** and **Pair Representations** and iteratively exchange information between them



# Structure module

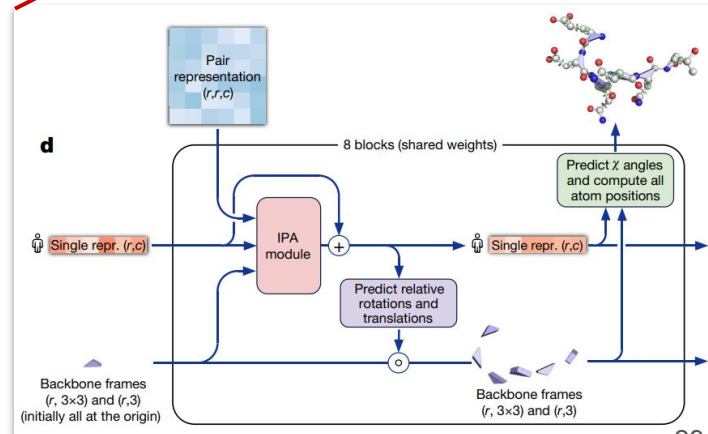
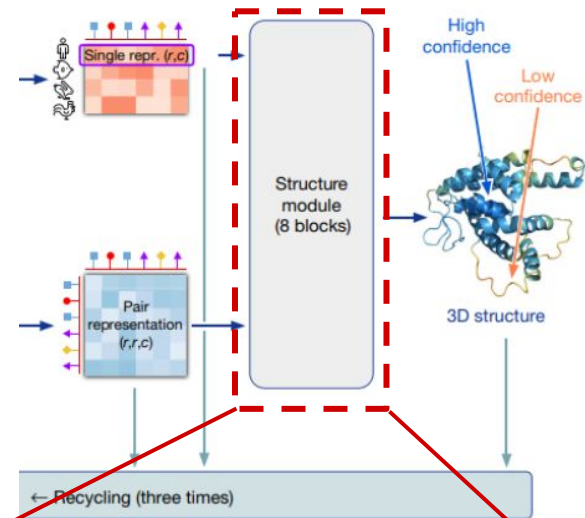
- Get structure from the MSA and pair representations
- Each residue is parameterized by a tuple representing rotation and translation

$$T_i := (R_i, \vec{t}_i)$$



(Image source)

- Iterative update to predict angle and atom positions



# Readings

## Article

### Highly accurate protein structure prediction with AlphaFold


<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

 Check for updates

John Jumper<sup>1,4,5</sup>, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4,5</sup>

AlphaFold (*Nature*, 2021)

<https://www.nature.com/articles/s41586-021-03819-2>

## RESEARCH

### RESEARCH ARTICLE

#### PROTEIN FOLDING

### Accurate prediction of protein structures and interactions using a three-track neural network

Minkyung Baek<sup>1,2</sup>, Frank DiMaio<sup>1,2</sup>, Ivan Anishchenko<sup>1,2</sup>, Justas Dauparas<sup>1,2</sup>, Sergey Ovchinnikov<sup>3,4</sup>, Gyu Rie Lee<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Qian Cong<sup>5,6</sup>, Lisa N. Kinch<sup>7</sup>, R. Dustin Schaeffer<sup>6</sup>, Claudia Millán<sup>8</sup>, Hahnbeom Park<sup>1,2</sup>, Carson Adams<sup>1,2</sup>, Caleb R. Glassman<sup>9,10,11</sup>, Andy DeGiovanni<sup>12</sup>, Jose H. Pereira<sup>12</sup>, Andria V. Rodrigues<sup>12</sup>, Alberdina A. van Dijk<sup>13</sup>, Ana C. Ebrecht<sup>13</sup>, Diederik J. Opperman<sup>14</sup>, Theo Sagmeister<sup>15</sup>, Christoph Buhlheller<sup>15,16</sup>, Tea Pavkov-Keller<sup>15,17</sup>, Manoj K. Rathinaswamy<sup>18</sup>, Udit Dalwadi<sup>19</sup>, Calvin K. Yip<sup>19</sup>, John E. Burke<sup>18</sup>, K. Christopher Garcia<sup>9,10,11,20</sup>, Nick V. Grishin<sup>6,7,21</sup>, Paul D. Adams<sup>12,22</sup>, Randy J. Read<sup>8</sup>, David Baker<sup>1,2,23\*</sup>

RoseTTAFold (*Science*, 2021)

<https://www.science.org/doi/10.1126/science.abj8754>

- Two deep learning papers on protein structure prediction published on the same day
- Both in the presentation paper list of this course

# “Citizen science”: Volunteer/distributed computing

## Folding@Home

- <https://foldingathome.org/home/>
- Molecular dynamics (MD) simulations

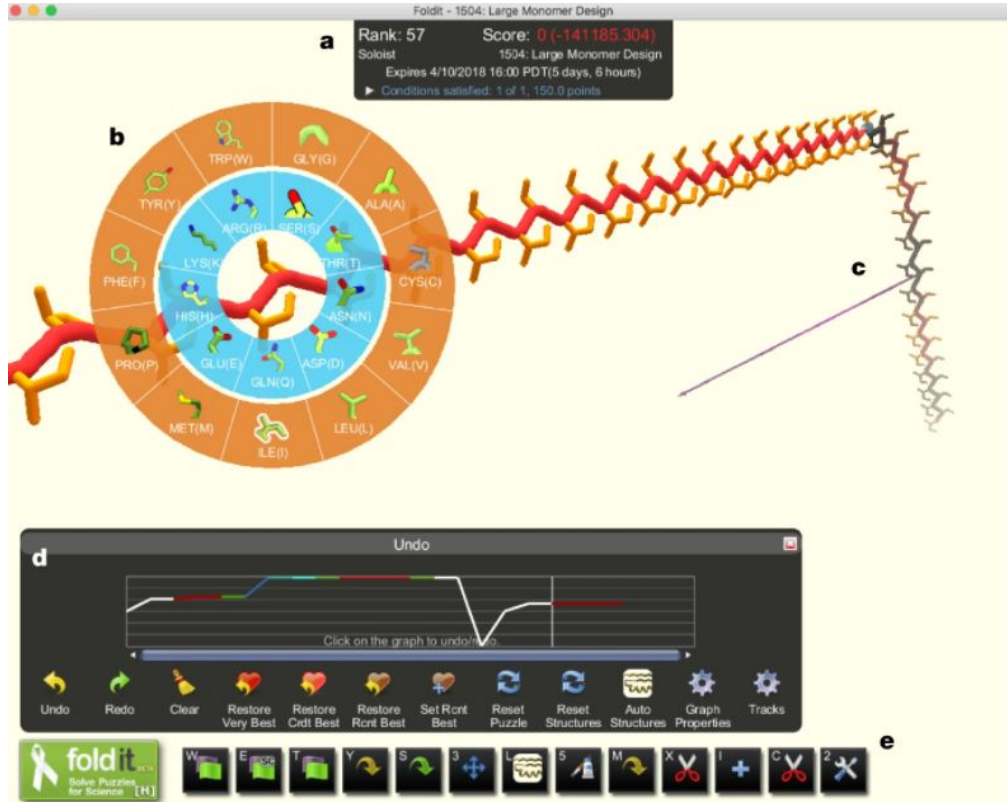


## Rosetta@home

- <http://boinc.bakerlab.org>
- Structure prediction



# Foldit



## nature

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [letters](#) > article

Letter | [Published: 05 June 2019](#)

### De novo protein design by citizen scientists

[Brian Koeppnick](#), [Jeff Flatten](#), [Tamir Husain](#), [Alex Ford](#), [Daniel-Adriano Silva](#), [Matthew J. Bick](#), [Aaron Bauer](#), [Gaohua Liu](#), [Yojiro Ishida](#), [Alexander Boykov](#), [Roger D. Estep](#), [Susan Kleinfelter](#), [Toke Nørgård-Solano](#), [Linda Wei](#), [Foldit Players](#), [Gaetano T. Montelione](#), [Frank DiMaio](#), [Zoran Popović](#), [Firas Khatib](#), [Seth Cooper](#) & [David Baker](#)

[Nature](#) **570**, 390–394 (2019) | [Cite this article](#)

**23k** Accesses | **56** Citations | **513** Altmetric | [Metrics](#)

# Protein Folding & Protein Design

- **Protein folding** (protein structure prediction)

- Sequence -> Structure

Amino acid sequence

```
MEKVFLKNGVLRLLPPGFRFRPTDEELVVQYLRKRVFSFPLPASIIPEVEVYKSDPWLPGDMEQEKYFFSTK  
EVKYPNGNRSNRATNSGYWKATGIDKQIILRGRQQQQLIGLKKTLVFYRGKSPHGCRTNWIMHEYRLAN  
LESNYHPIQGNWVICRIFLKKRGNTKNKEENMTTHDEVNRNREIDKNPVSVMSSRDSEALASANSELKK
```



Algorithm / Model



Protein structure



- **Protein design**

- Structure -> Sequence

Protein structure



Algorithm / Model



Amino acid sequence

```
MEKVFLKNGVLRLLPPGFRFRPTDEELVVQYLRKRVFSFPLPASIIPEVEVYKSDPWLPGDMEQEKYFFSTK  
EVKYPNGNRSNRATNSGYWKATGIDKQIILRGRQQQQLIGLKKTLVFYRGKSPHGCRTNWIMHEYRLAN  
LESNYHPIQGNWVICRIFLKKRGNTKNKEENMTTHDEVNRNREIDKNPVSVMSSRDSEALASANSELKK
```

Protein Design by Hand!

# Design sequences that folds into the following structures:

Make a helix!



Make two helices



Make beta sheets?



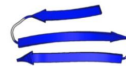
[Google Colab](#)

## Cheatsheet:

Secondary structure propensity

	C	H	E
A	-0.41	0.46	-0.40
E	-0.27	0.42	-0.53
Q	-0.24	0.36	-0.39
K	-0.05	0.21	-0.32
R	-0.22	0.24	-0.14
M	-0.37	0.26	0.03
L	-0.64	0.35	0.13
W	-0.37	0.06	0.36
Y	-0.36	-0.03	0.46
F	-0.39	-0.03	0.48
I	-0.77	-0.03	0.74
V	-0.69	-0.25	0.89
C	-0.06	-0.25	0.41
T	0.14	-0.39	0.30
H	0.14	-0.15	0.01
S	0.32	-0.25	-0.17
N	0.52	-0.30	-0.65
D	0.50	-0.21	-0.78
G	0.82	-1.05	-0.56
P	0.90	-0.94	-1.06

$\log(P(AA|SS) / P(AA))$

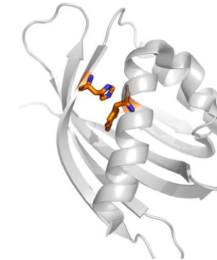
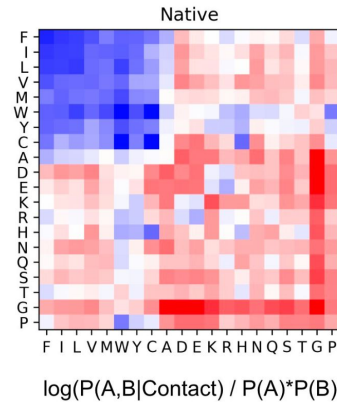


(E) Sheets

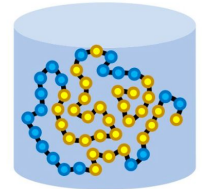
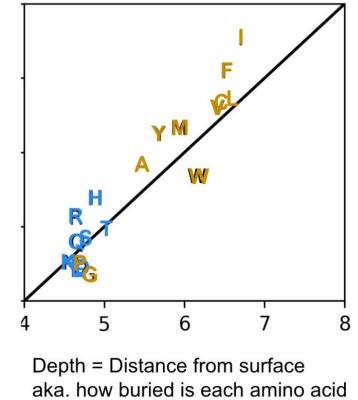


(H) Helix

Pairwise potential



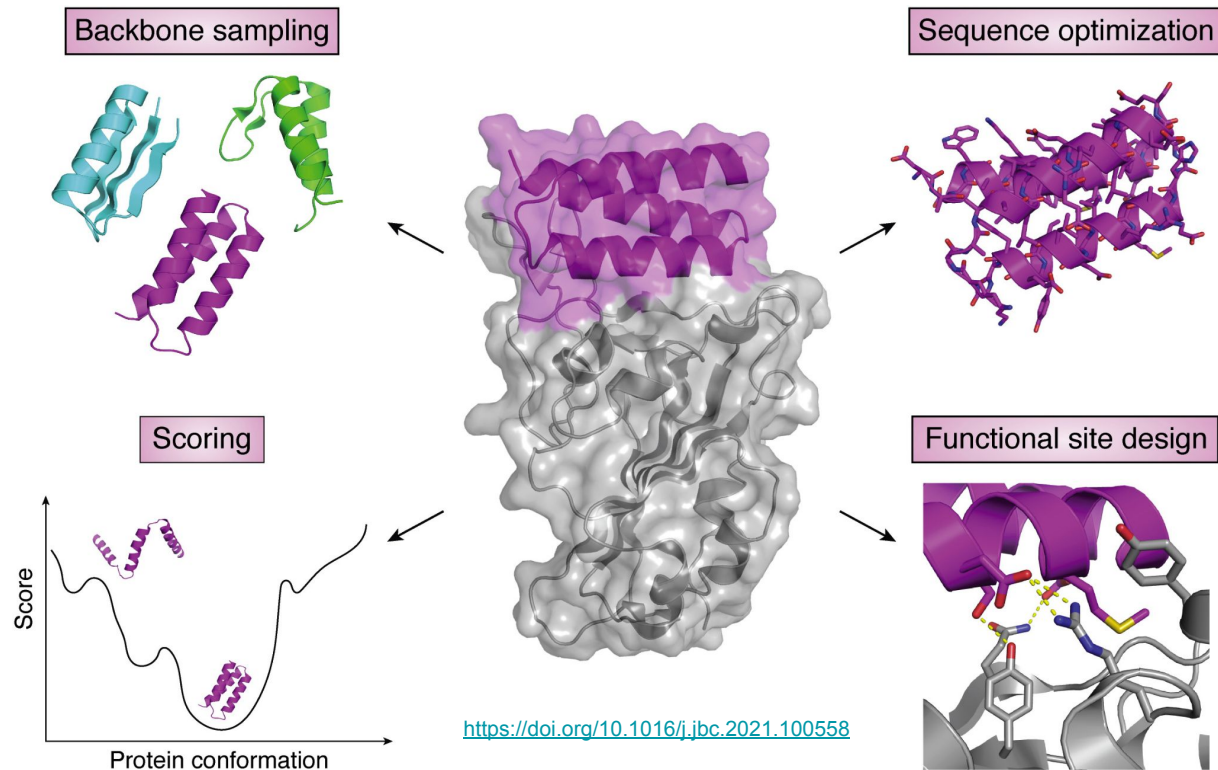
Average depth per amino acid type



(C) Random Coil (residues which are not in any of the above conformations)



# De novo Protein design



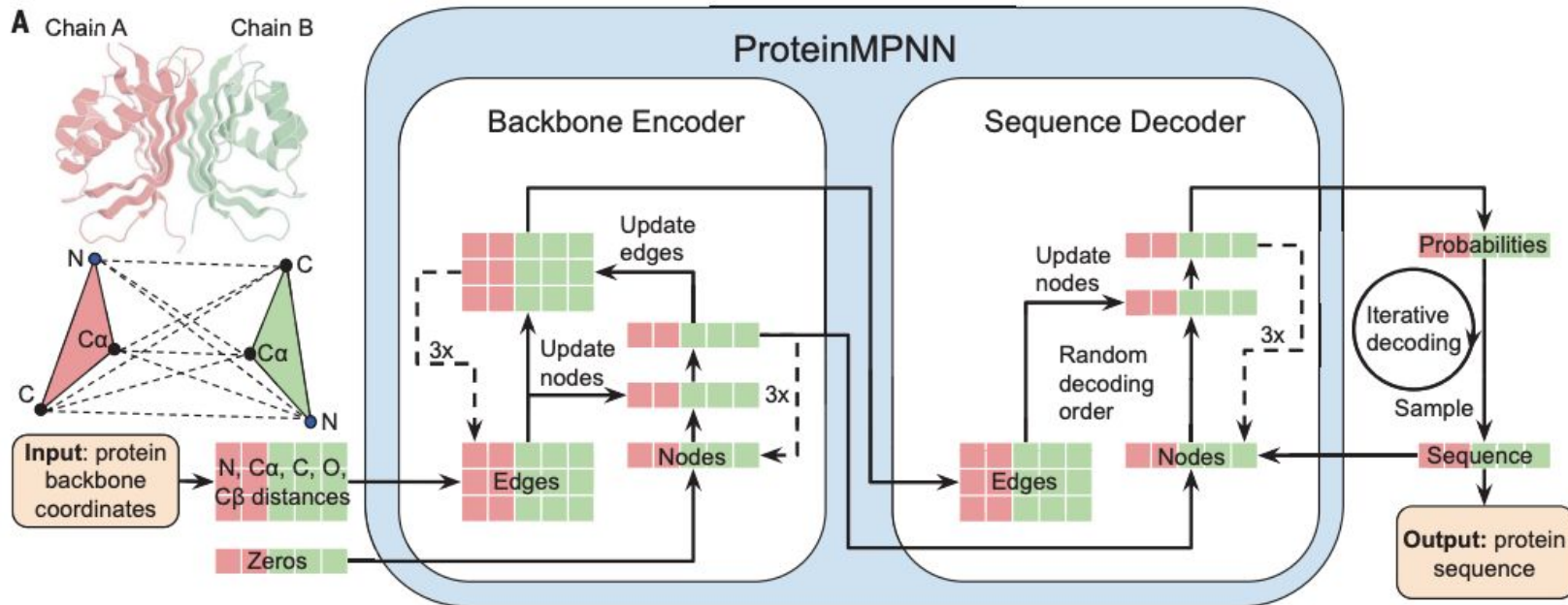
The design of a functional *de novo* protein, for example, a binder (*middle, magenta*) to a target protein (*middle, gray*), requires sampling of the backbone structure space to find a backbone compatible with the function, sequence optimization to stabilize the backbone, and designing the functional site interactions. A scoring function is necessary to select designs with desired properties, typically by identifying low-energy sequence–structure combinations.

## Robust deep learning-based protein sequence design using ProteinMPNN

J. DAUPARAS, I. ANISHCHENKO, N. BENNETT, H. BAI, R. J. RAGOTTE, L. F. MILLES, B. J. M. WICKY, A. COURBET, B. J. DE HAAS, I. J. AND D. BAKER

+12 authors Authors Info & Affiliations

SCIENCE • 15 Sep 2022 • Vol 378, Issue 6615 • pp. 49-56 • DOI: 10.1126/science.adg2187



# Conclusion

- Protein folding
  - Sequence -> Structure
- Protein design
  - Structure -> Sequence

# Reminder:

- You are encouraged to ask a question during the Q&A time of a presentation
- **Log your question in a survey form to receive the points (by 11:59pm, same day).** Once we verified your asked question, you will receive 2 points.
- The question is expected to be in-depth and preferably can prompt discussion/debate
- clarification questions will not count, e.g.:
  - Did the paper compare their model with Random forest?
  - How did they process the input data?

Spring 2024

[Home](#)

[Syllabus](#)

[Ed Discussion](#)




[Assignments](#)

[Gradescope](#)

[Grades](#)

## ML in CompBio - CSE7850/CX4803-MLB

### Resources

- Syllabus: [Link](#) 
- Schedule: [Link](#) 
- In-class discussion check-in form: [Link](#) 

### Office hours

- Instructor: