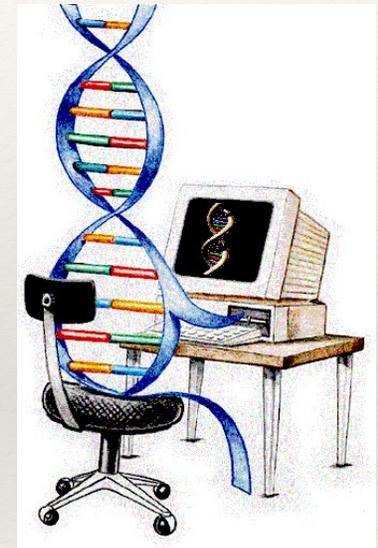# Ab Initio Protein Structure Prediction: AlphaFold

# Ab initio Protein Structure Prediction

Ab initio prediction before AlphaFold

Ab initio prediction: Predicting Contacts

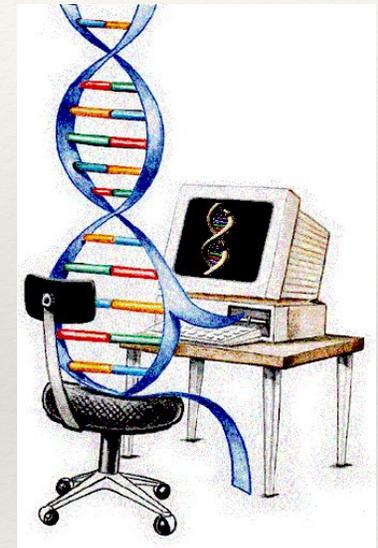AlphaFold 1

AlphaFold 2

# Ab initio Protein Structure Prediction

**Ab initio prediction before AlphaFold**

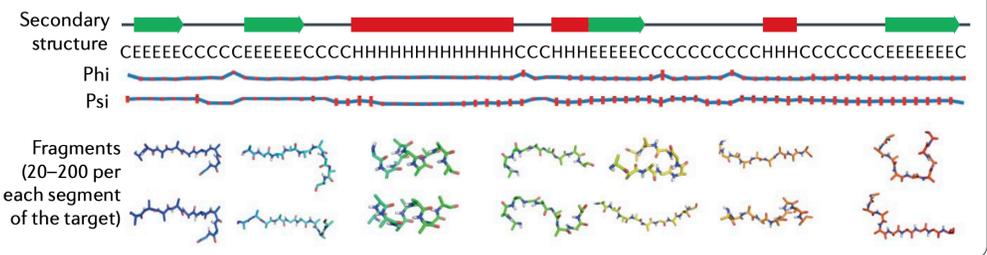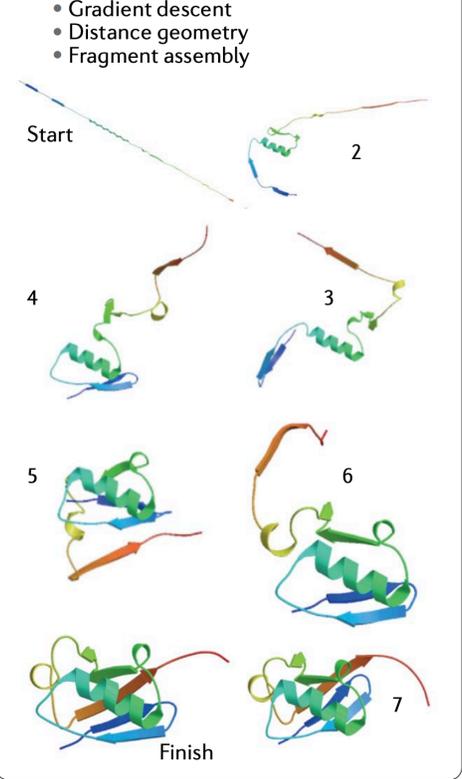Ab initio prediction: Predicting Contacts

AlphaFold 1

AlphaFold 2

① **Construct multiple-sequence alignment**

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | * | | * | | | * | | * | | |
| Target | K | T | L | R | G | K | G | I | T | D | E | V | F | P | S |
| Homologue 1 | K | T | L | E | G | K | A | I | T | K | K | V | W | S | R |
| Homologue 2 | K | T | L | R | G | K | F | I | A | E | E | A | A | Q | N |
| Homologue 3 | E | I | P | E | G | W | F | I | S | K | S | C | A | P | S |
| Homologue n | K | T | L | E | G | K | W | V | T | K | E | V | G | P | T |

② **Predict local structure**

Secondary structure

CEEEEECCCCCEEEEEECCCCHHHHHHHHHHHHHHCCHHHEEEEECCCCCCCCCCCHHHCCCCCCEEEEEEEEC

Phi

Psi

Fragments (20–200 per each segment of the target)

③ **Assemble 3D models**
- Gradient descent
- Distance geometry
- Fragment assembly

Start

2

4          3

5          6

Finish          7

④ **Refine and rank models**

Landscape projection

Energy

Start

Refinement trajectory

Finish

# Fragment based methods

**Exploring the energy landscape**

Gradient-based minimization

Energy

Start

Finish

Nearest energy minimum

Start

Finish

Conformation

**Metropolis Monte Carlo**

Energy (vertical axis) vs Conformation (horizontal axis)

$P = 0.03$
$\Delta E = 2.1$
Start
$P = 0.58$
$P = 1$
$\Delta E = 0.33$
$P = 1$
$P = 0.07$
$\Delta E = 1.6$
Global minimum
$P = 1$
Finish

**Monte Carlo moves**

Fragment replacement
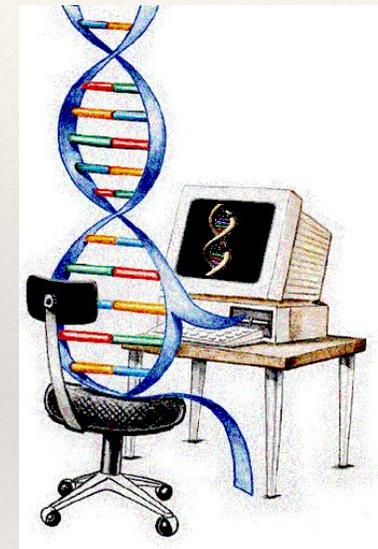
Rotamer substitution

# Ab initio Protein Structure Prediction

Ab initio prediction before AlphaFold

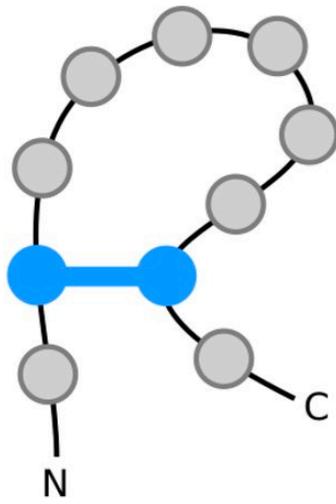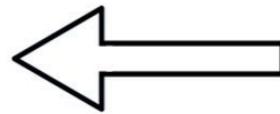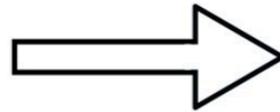**Ab initio prediction: Predicting Contacts**

AlphaFold 1

AlphaFold 2

# Predicting residue contacts



Interaction in structure → Constraint → Covariation in sequence alignment

Inference

# Predicting residue contacts

1. *Given a multiple sequence alignment (MSA):*

$X_1$

```
H A G D T A I L L M R W K D A
H L G D T A I L L M R W K D C
H L G D T S I L L M R W K D C
H A G E T T I L V M K W K D A
H I G E T A I L L M K W K D C
H A G E T T I L V M K W K D C
```

$X_N$

2. *Compute "mean" sequence and covariance matrix:*

$$\overline{\mathbf{X}} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{X}_n$$

$$\overline{C} = C(MSA, \overline{\mathbf{X}}) = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{X}_n - \overline{\mathbf{X}})^T (\mathbf{X}_n - \overline{\mathbf{X}})$$

3. *Compute contact J(i,j)*

$$J(i, j) = C(i, j)?$$

# Predicting residue contacts

*No! We need to pay attention to indirect effects:*

# Predicting residue contacts

*No! We need to pay attention to indirect effects:*

*Gaussian model:*

Each sequence $X_i$ in the MSA is drawn from a multivariate Gaussian distribution characterized by a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\Sigma$, with the probability:

$$P(\mathbf{X}_n | \boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{Ls}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\mathbf{X}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{X}_n - \boldsymbol{\mu})\right]$$
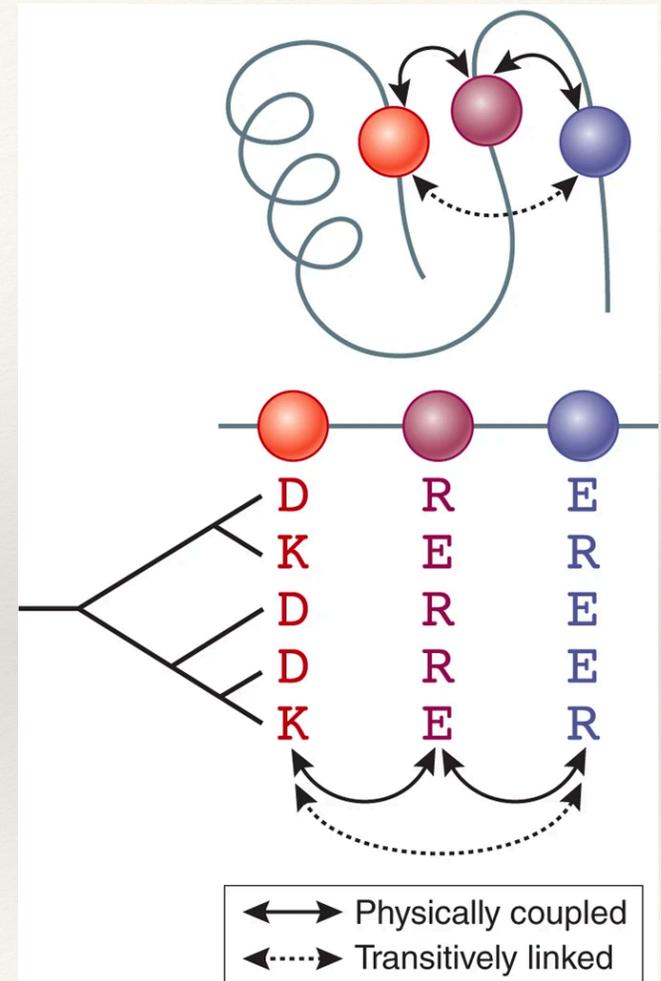
# Predicting residue contacts

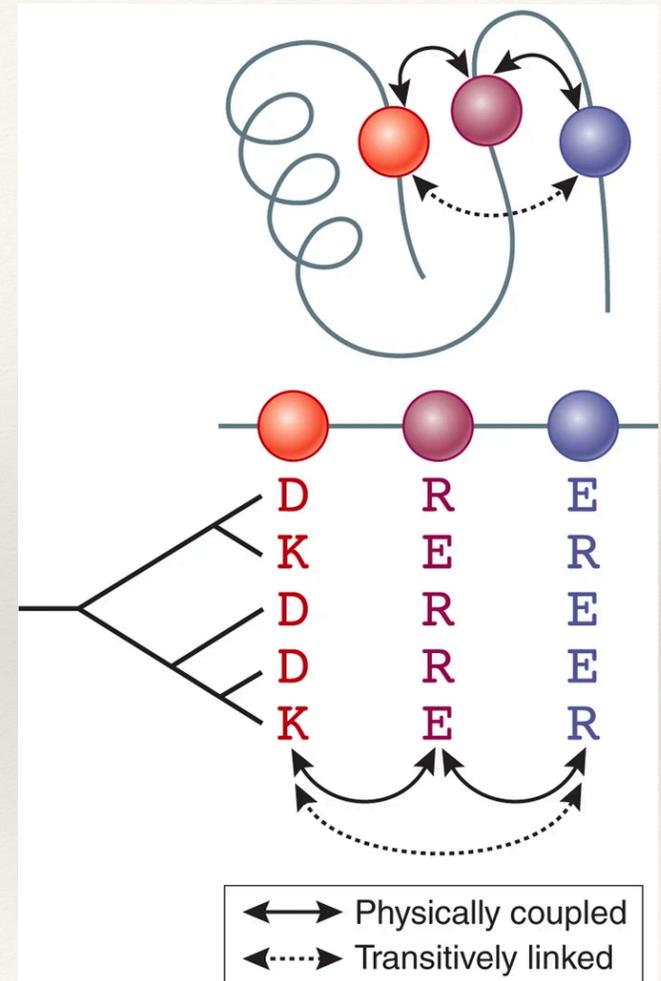*No! We need to pay attention to indirect effects:*

*Gaussian model:*

Each sequence $X_i$ in the MSA is drawn from a multivariate Gaussian distribution characterized by a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\Sigma$, with the probability:

$$P(\mathbf{X}_n | \boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{Ls}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\mathbf{X}_n - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{X}_n - \boldsymbol{\mu})\right]$$

Assuming that the $N$ sequences in the MSA are statistical independent, the probability, or likelihood of the data under this model is given by

$$P(MSA | \boldsymbol{\mu}, \Sigma) = \prod_{n=1}^{N} P(\mathbf{X}_n | \boldsymbol{\mu}, \Sigma)$$



| D | R | E |
| K | E | R |
| D | R | E |
| D | R | E |
| D | R | E |
| K | E | R |

↔ Physically coupled
◄┈┈► Transitively linked

# Predicting residue contacts

*No! We need to pay attention to indirect effects:*

*Gaussian model:*

Each sequence $X_i$ in the MSA is drawn from a multivariate Gaussian distribution characterized by a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\Sigma$, with the probability:

$$P(\mathbf{X}_n|\boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{Ls}{2}}|\Sigma|^{-\frac{1}{2}}\exp\left[-\frac{1}{2}(\mathbf{X}_n - \boldsymbol{\mu})^T\Sigma^{-1}(\mathbf{X}_n - \boldsymbol{\mu})\right]$$
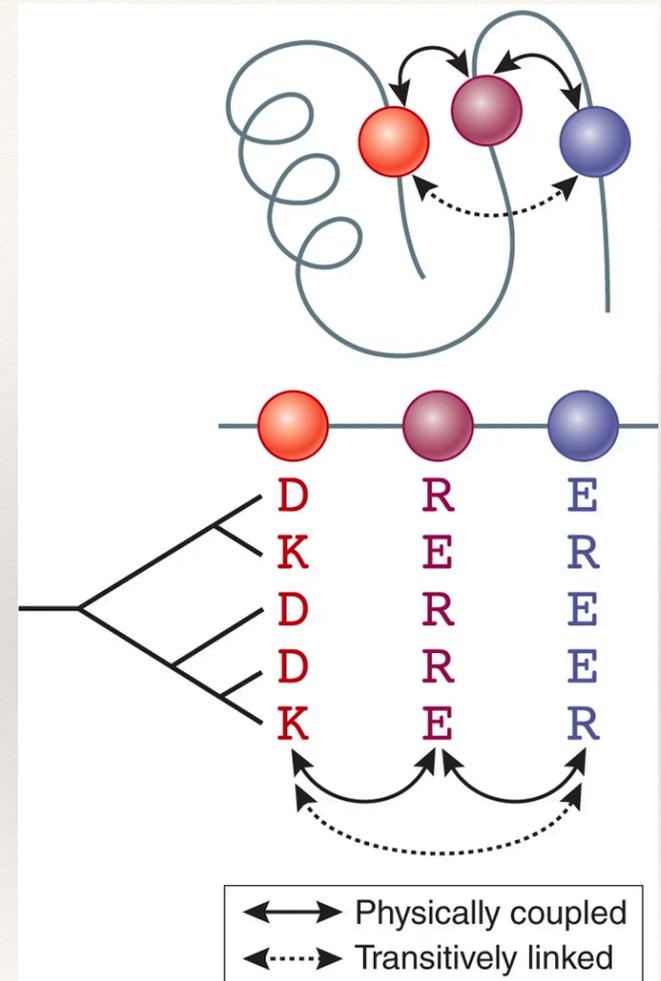
Assuming that the $N$ sequences in the MSA are statistical independent, the probability, or likelihood of the data under this model is given by

$$P(MSA|\boldsymbol{\mu}, \Sigma) = \prod_{n=1}^{N} P(\mathbf{X}_n|\boldsymbol{\mu}, \Sigma)$$

Using the maximum likelihood estimator for this probability

$$\boldsymbol{\mu} = \overline{\mathbf{X}}$$

$$\Sigma = \overline{C} = C(MSA, \overline{\mathbf{X}})$$

# Predicting residue contacts

*No! We need to pay attention to indirect effects:*

*Gaussian model:*

$$P(\mathbf{X}_n | \boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{Ls}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left[-\frac{1}{2}(\mathbf{X}_n - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{X}_n - \boldsymbol{\mu})\right]$$

$$\boldsymbol{\mu} = \overline{\mathbf{X}} \qquad\qquad \Sigma = \overline{C} = C(MSA, \overline{\mathbf{X}})$$
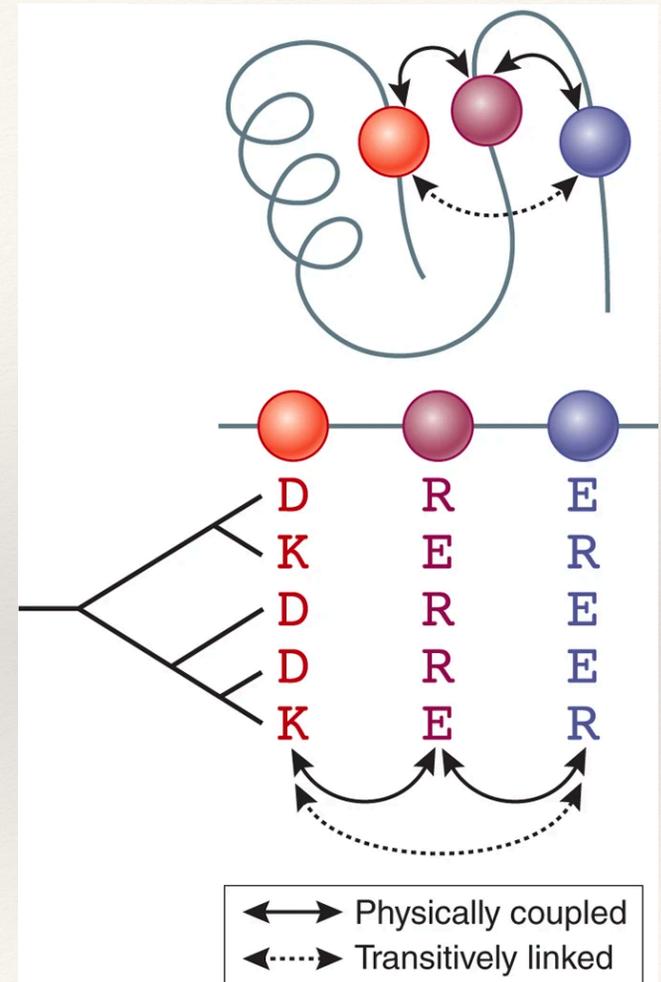
Note that:

$$(\mathbf{X}_n - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{X}_n - \boldsymbol{\mu}) = \sum_{k=1}^{N}\sum_{l=1}^{N}(X_k - \mu_k)(\Sigma^{-1})(k, l)(X_l - \mu_l)$$

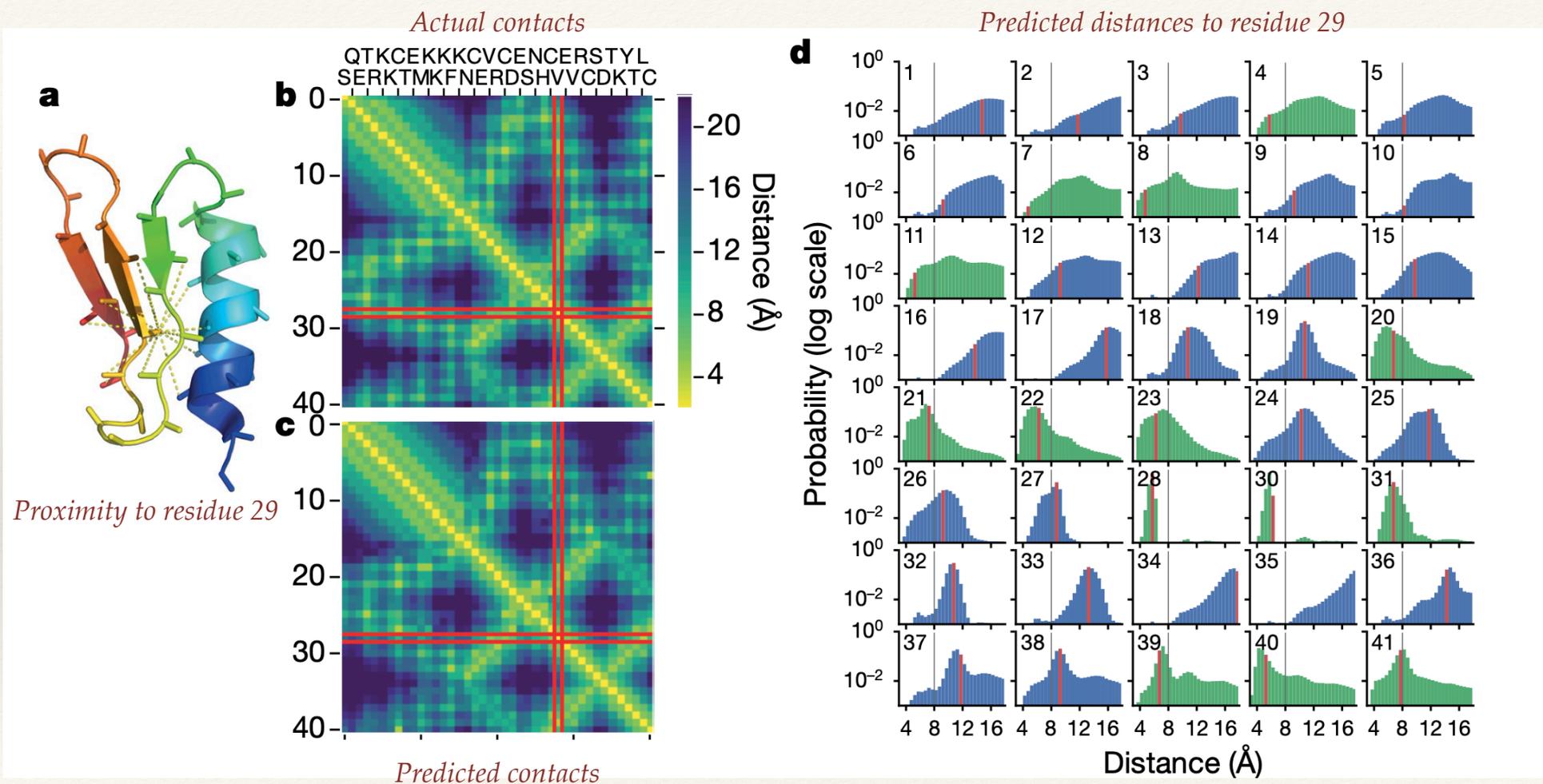This shows that $(\Sigma^{-1})(k, l)$ serves as a coupling between positions $k$ and $l$ in the MSA.

Therefore:

$$J = \Sigma = (C(MSA, \overline{X}))^{-1}$$

# Predicting residue contacts

# Predicting residue contacts: How well does it work?



*Actual contacts*

*Proximity to residue 29*

*Predicted contacts*

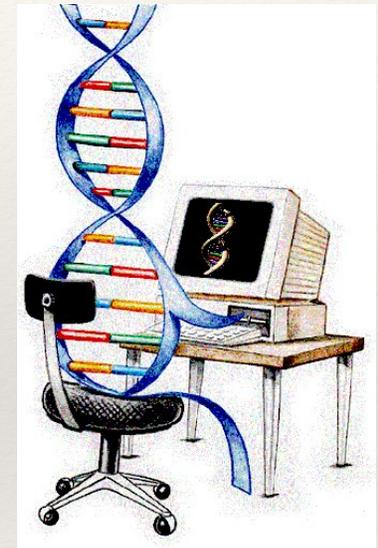*Predicted distances to residue 29*

*CASP target T0995*

# Ab initio Protein Structure Prediction
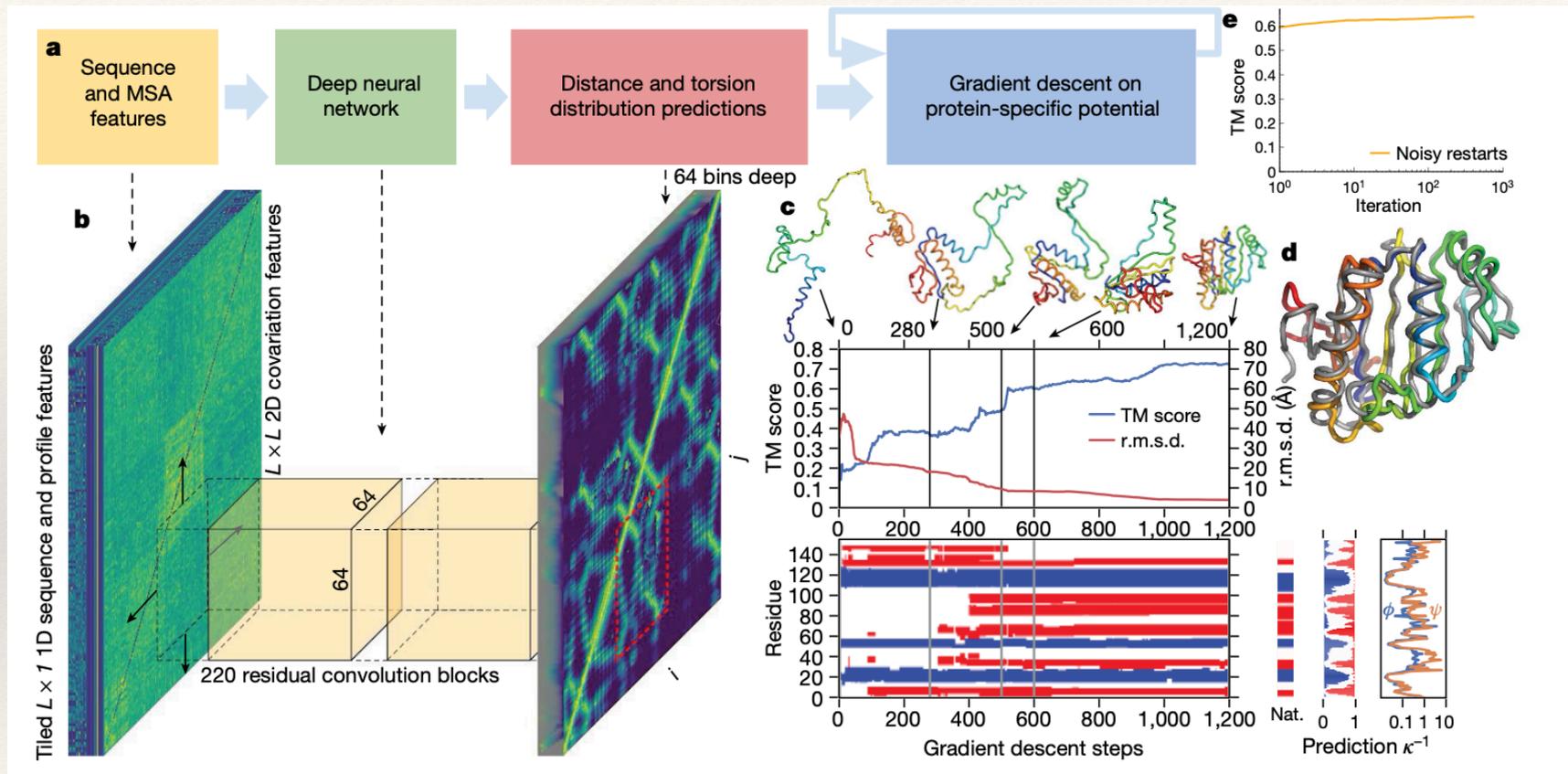
Ab initio prediction before AlphaFold

Ab initio prediction: Predicting Contacts

**AlphaFold 1**

AlphaFold 2

# AlphaFold 1



AlphaFold1

# AlphaFold 1

*Reminder:*

**To compare two sets of points (atoms) $A=\{a_1, a_2, \ldots a_N\}$ and $B=\{b_1, b_2, \ldots, b_N\}$:**

**-Define a 1-to-1 correspondence between A and B**

for example, $a_i$ corresponds to $b_i$, for all i in [1,N]

**-Compute RMS as:**                                                        **Compute TM score:**

$$RMS(A,B) = \sqrt{\frac{1}{N}\sum_{i=1}^{N} d(a_i, b_i)^2}$$        $$TM(A,B) = \frac{1}{N}\sum_{i=1}^{N} \frac{1}{1 + \left(\frac{d(a_i, b_i)}{d_0(N)}\right)^2}$$   with $d_0(N) = 1.24\sqrt[3]{N-15} - 1.8$
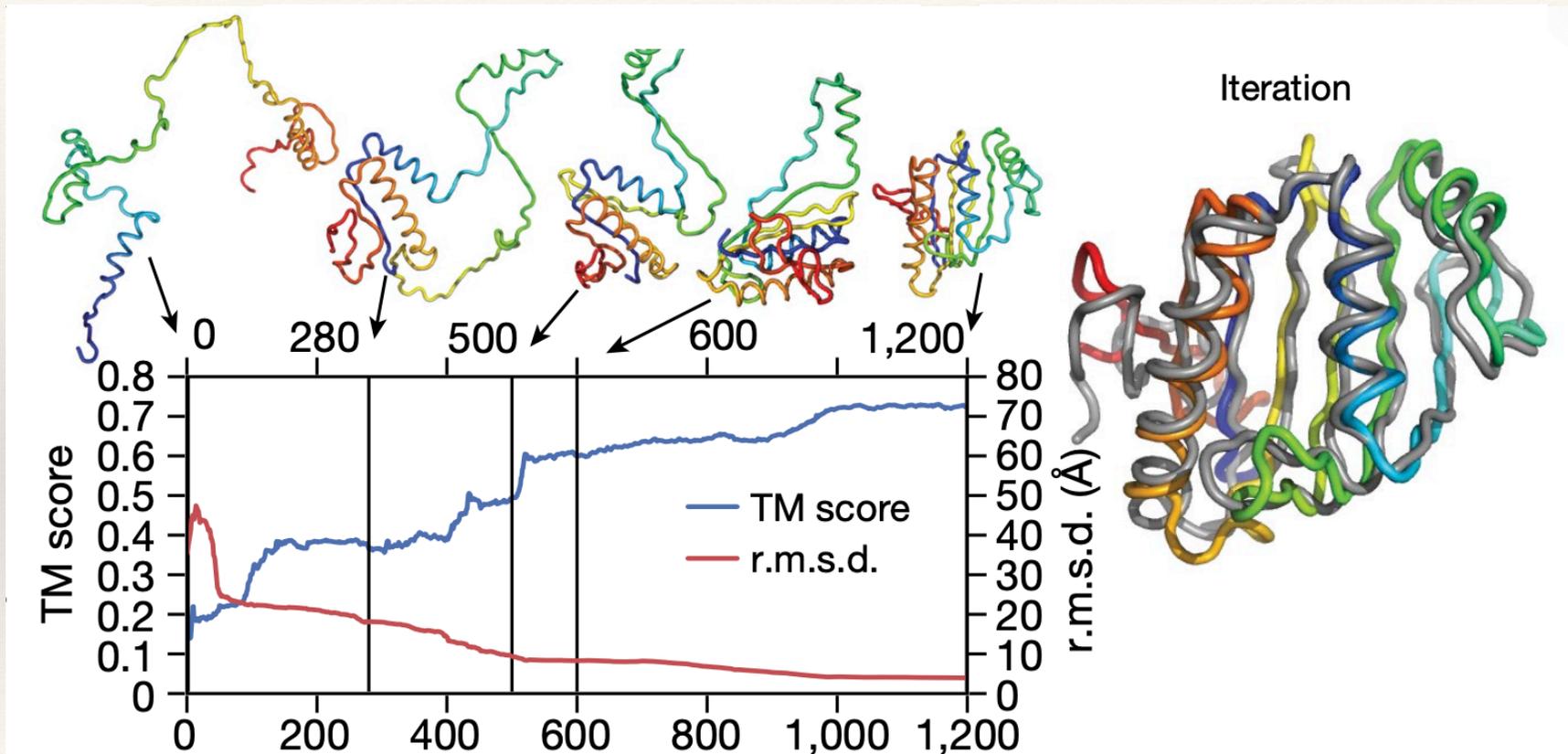
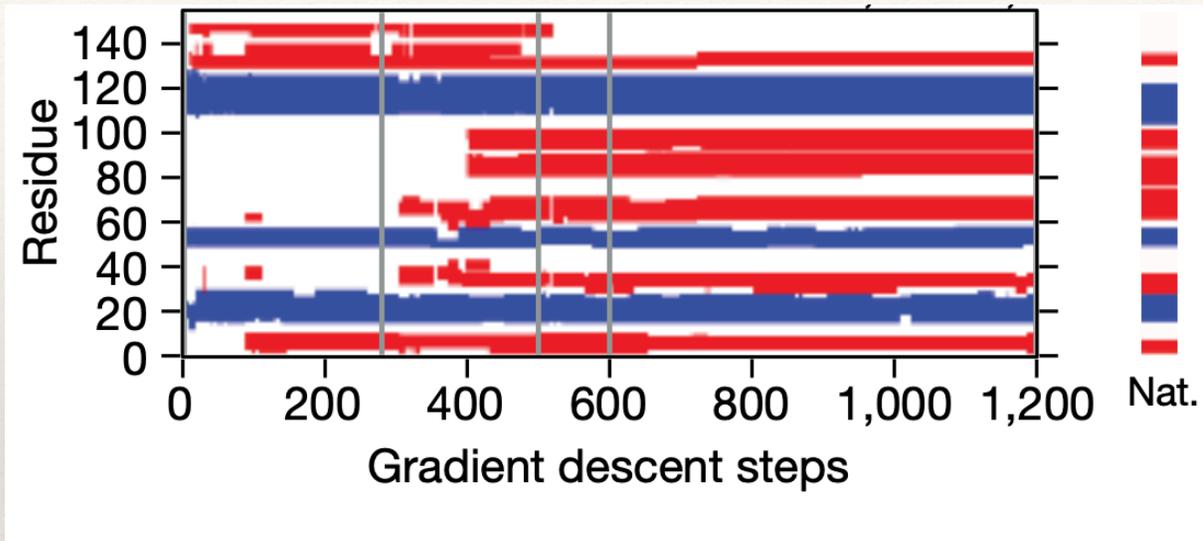$d(a_i,b_i)$ is the Euclidian distance between $a_i$ and $b_i$ after optimal alignment of B onto A

RMS: the lower, the better                                    TM: between [0,1]; the higher the better
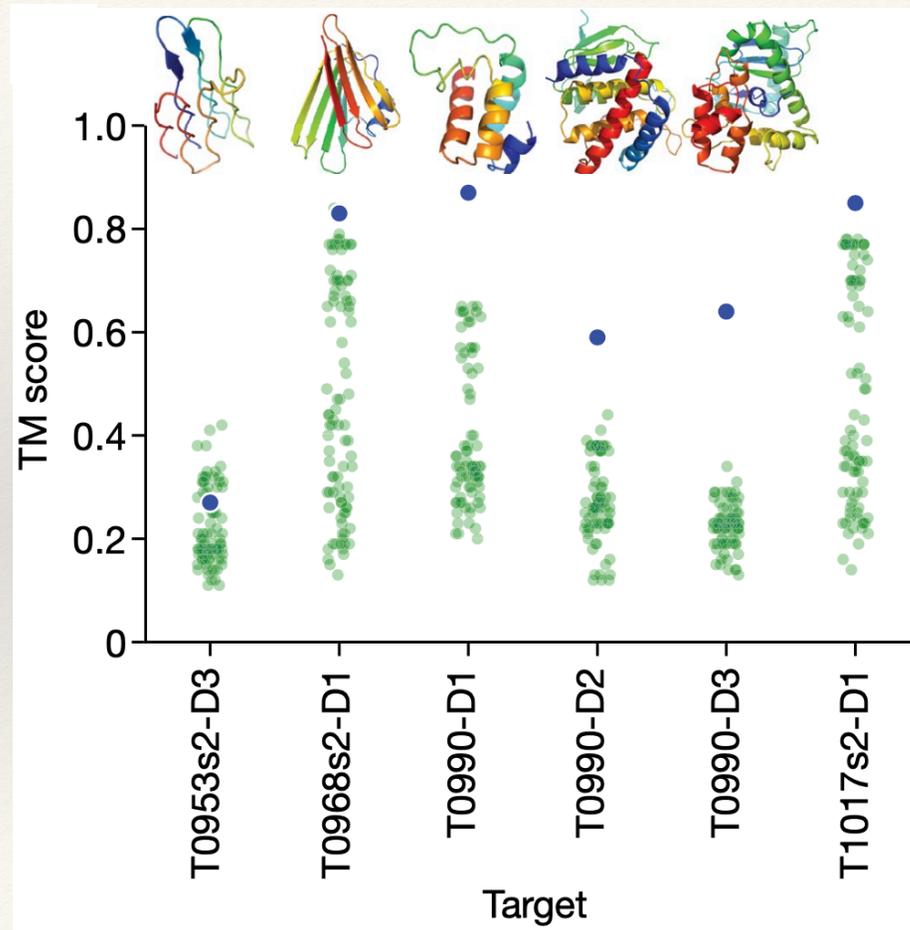
# AlphaFold 1

# AlphaFold 1



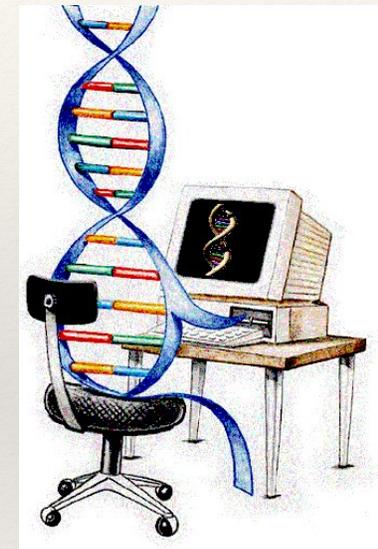*Helix in blue, strand in red*

# AlphaFold 1: Success

# Ab initio Protein Structure Prediction

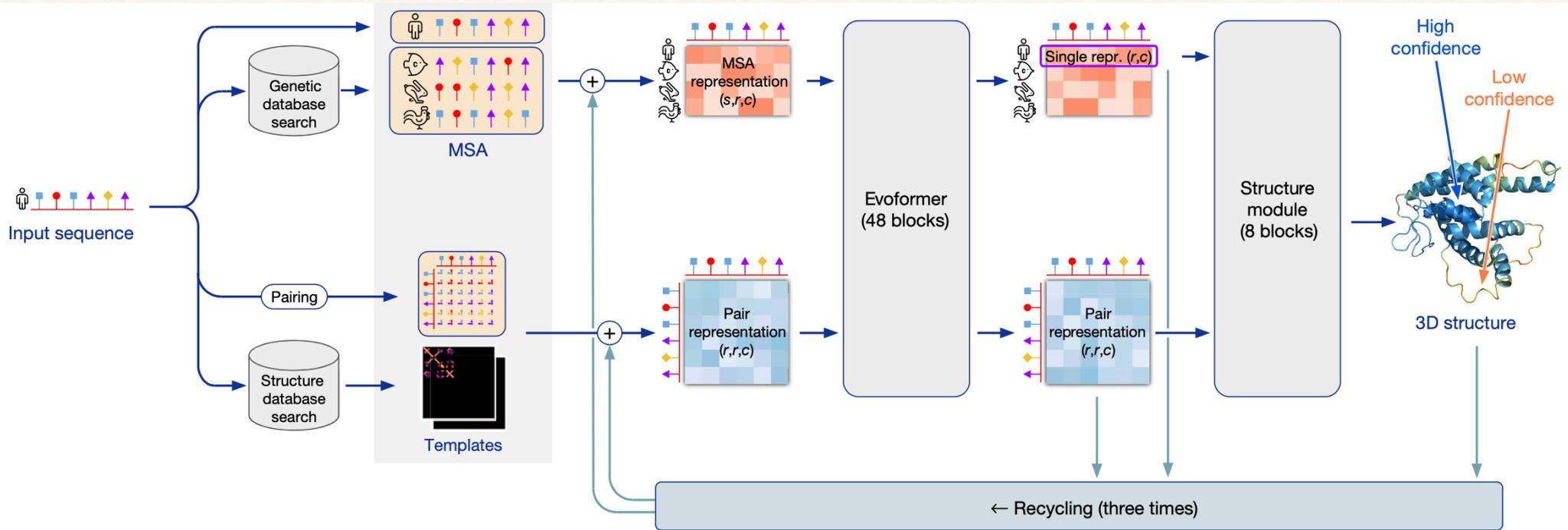Ab initio prediction before AlphaFold

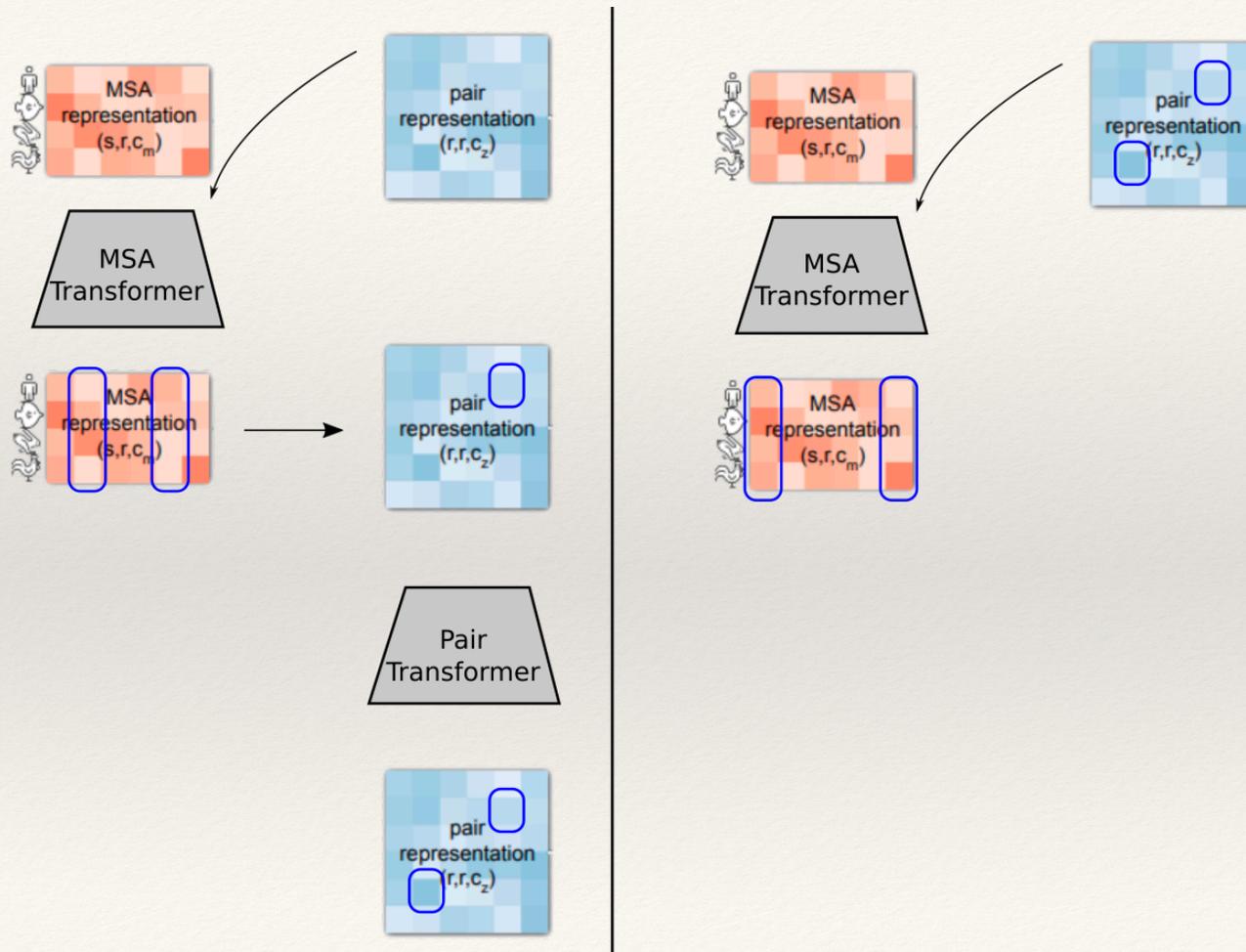Ab initio prediction: Predicting Contacts
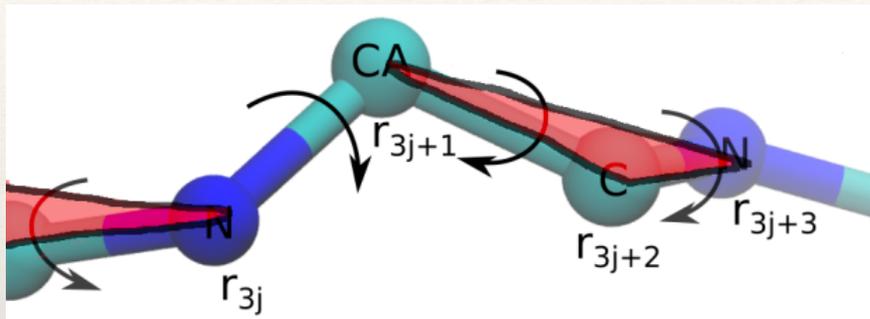
AlphaFold 1

AlphaFold 2

# AlphaFold 2

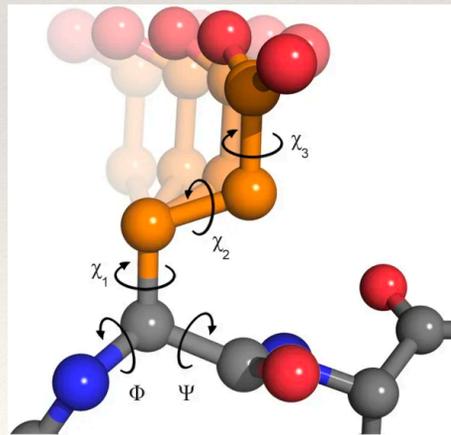# AlphaFold 2: some intuition

# AlphaFold 2: the structure module

*Predicting backbone:*
*the residues form a gas soup of triangles whose relative positions are characterized by affine transformation*
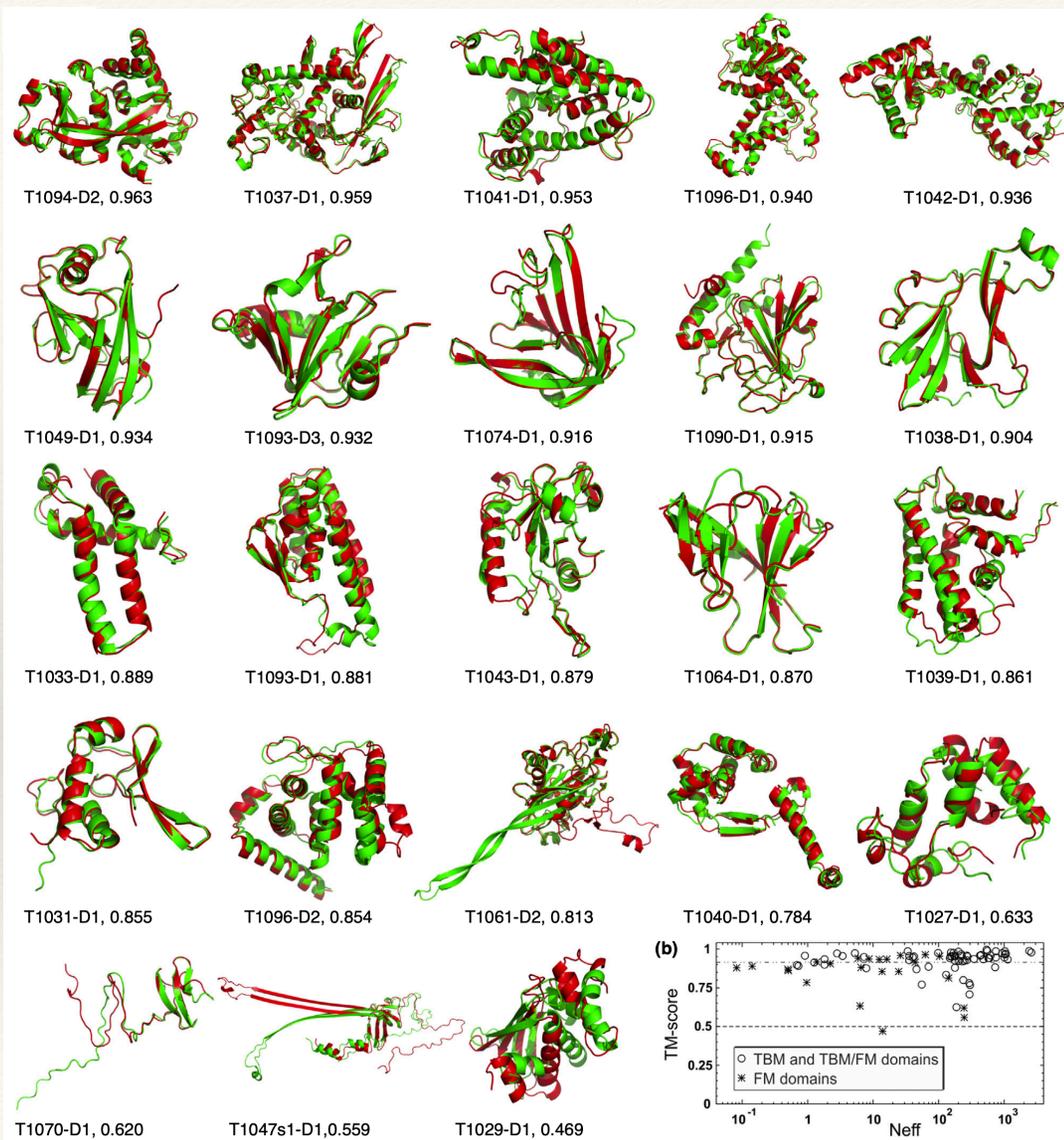


$$M = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ 0 & 0 & 0 & 1 \end{pmatrix}$$
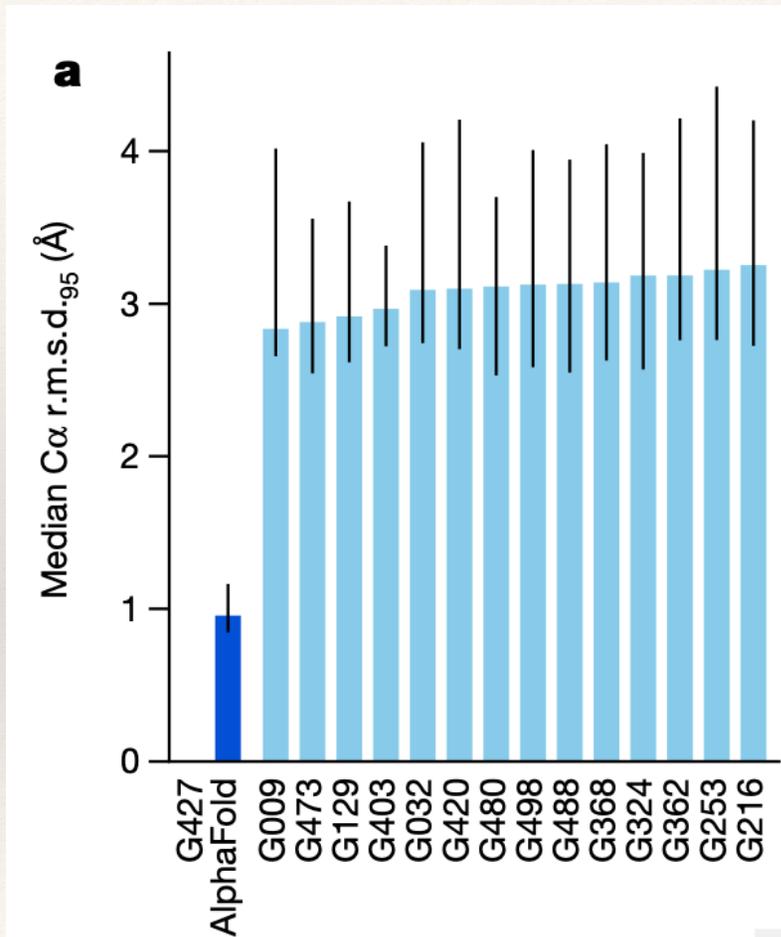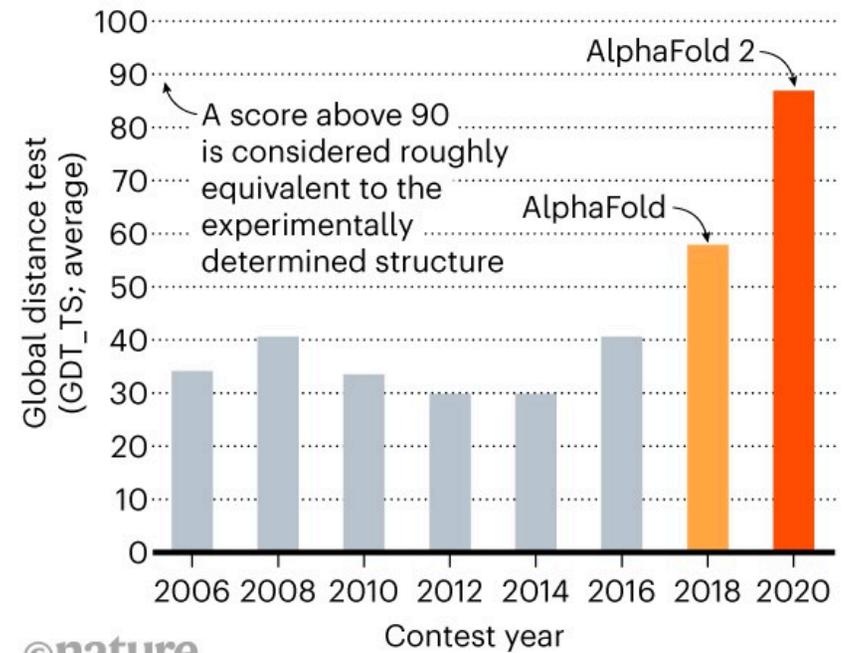
*Predicting side chains:*

**Successes at CASP14**



T1094-D2, 0.963   T1037-D1, 0.959   T1041-D1, 0.953   T1096-D1, 0.940   T1042-D1, 0.936

T1049-D1, 0.934   T1093-D3, 0.932   T1074-D1, 0.916   T1090-D1, 0.915   T1038-D1, 0.904

T1033-D1, 0.889   T1093-D1, 0.881   T1043-D1, 0.879   T1064-D1, 0.870   T1039-D1, 0.861

T1031-D1, 0.855   T1096-D2, 0.854   T1061-D2, 0.813   T1040-D1, 0.784   T1027-D1, 0.633

T1070-D1, 0.620   T1047s1-D1, 0.559   T1029-D1, 0.469

*TBM: template-based modeling*
*FM:    free modeling*
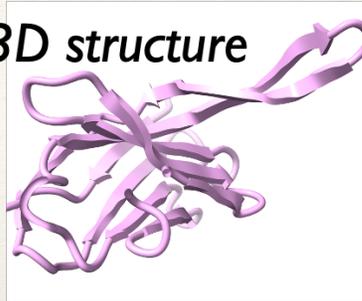
*Successes at CASP14*



**STRUCTURE SOLVER**

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.

A score above 90 is considered roughly equivalent to the experimentally determined structure

©nature

**Training**
- *Sequence*
- *Multiple sequence alignment*
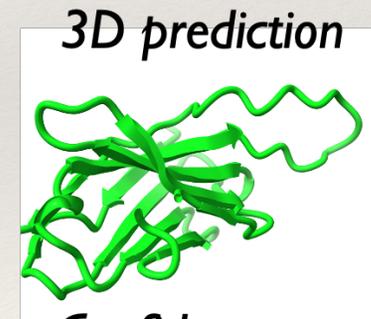  - *3D structure*

→ **21 million parameters**

**Prediction**

- *Sequence*
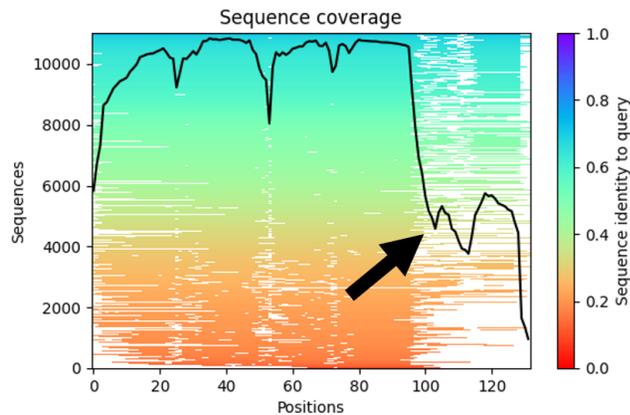- *Multiple sequence alignment*

```
EVQLVESGGGLVQPGGSLRLSCAASGFNIYSSSIHWVRQAPGKGLEWVAYI
...................................F...M.........Q........
...................K...........Y...L........A........
.................A.........................................V......
.................A.......................................
..............................L...V....E.............
.........A.....................................Q......
```
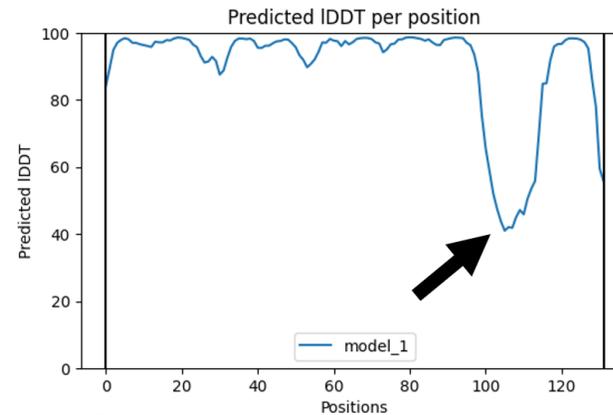
**21 million parameters**

*Focus attention on important relationships*

→ *3D prediction*

*Confidence estimates*

*Credit: Tom Terwilliger, Los Alamos NL*

# Multiple sequence alignment

EVQLVESGGGLVQPGGSLRLSCAASGFN**I**YSSS**I**HWVRQAPGKGLEWVAYI
.........................................**F**...**M**........Q........
.....................K.......**Y**...**L**......A.......
...........A....................................V......
.............................A..................................
.........................**L**...**V**...E...........
...........A..................................Q.....

Residues that **co-vary** are probably close in 3D structure

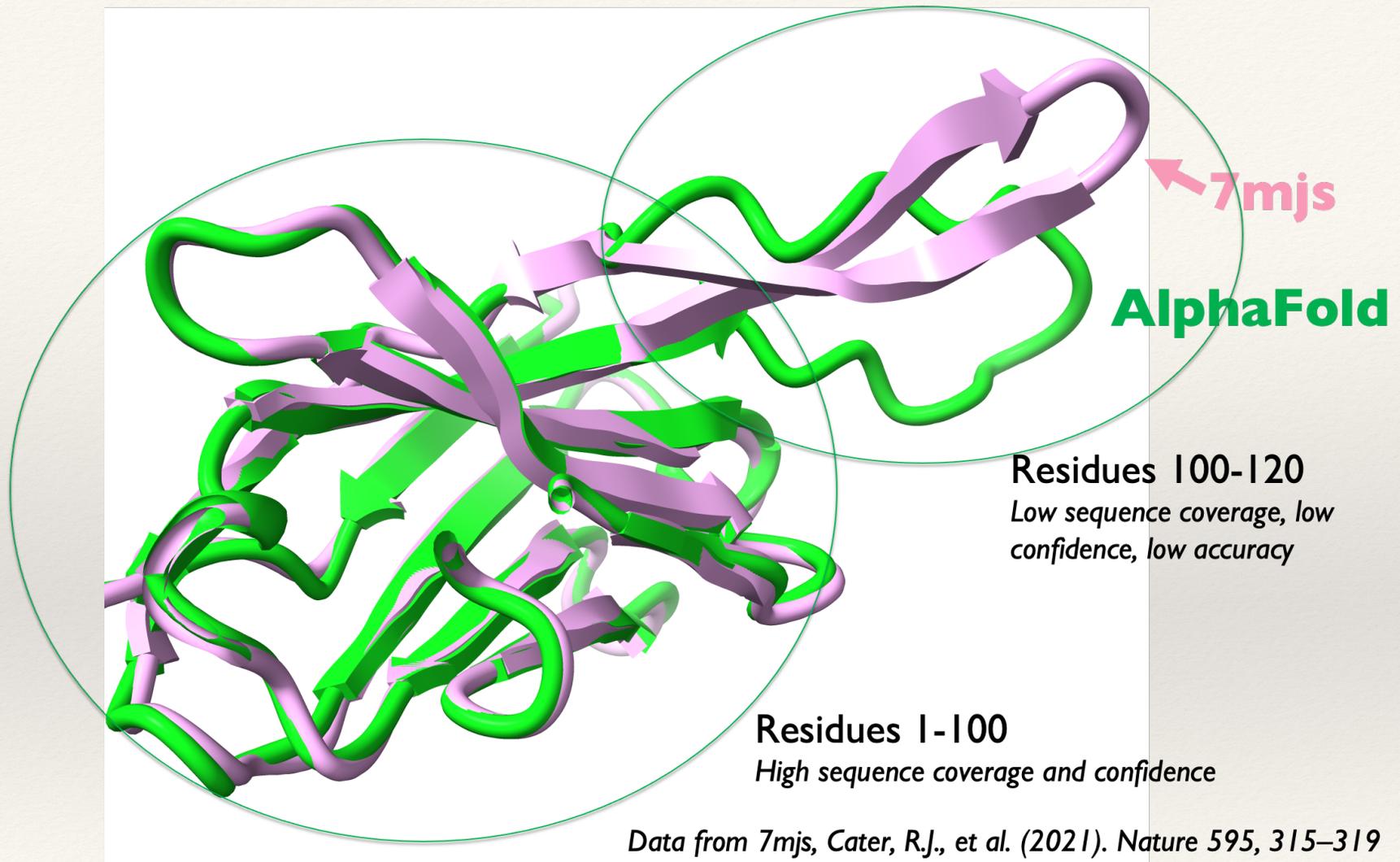All sequences in alignment should be compatible with the right structure
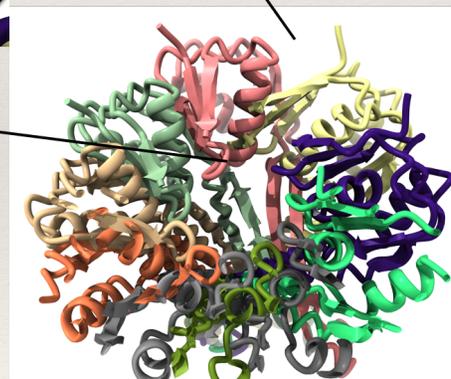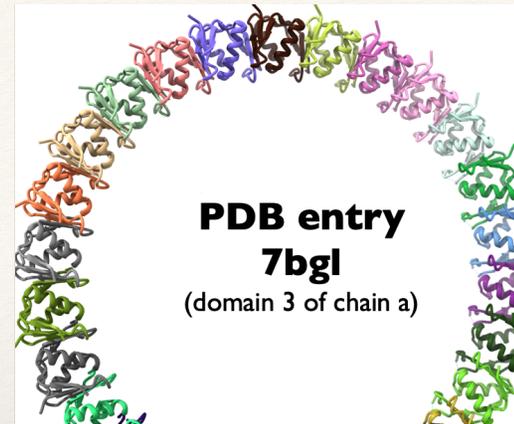
Sequence coverage → Confidence

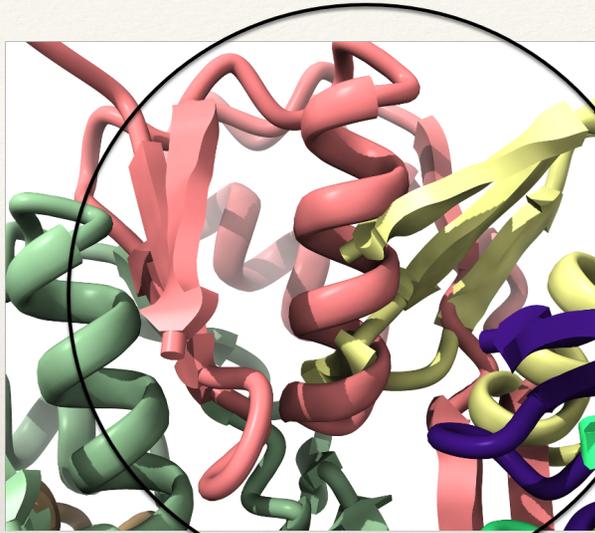

Data from 7mjs, Cater, R.J., et al. (2021). Nature 595, 315–319

7mjs

AlphaFold

Residues 100-120
*Low sequence coverage, low confidence, low accuracy*

Residues 1-100
*High sequence coverage and confidence*

*Data from 7mjs, Cater, R.J., et al. (2021). Nature 595, 315–319*

*Credit: Tom Terwilliger, Los Alamos NL*

# Multimeric proteins



**PDB entry
7bgl**
(domain 3 of chain a)

**AlphaFold**
(multimer prediction)

*Data from 7bgl, Johnson, S. et al. (2021).
Nat Microbiol 6, 712–721*

*Credit: Tom Terwilliger, Los Alamos NL*
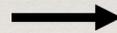
Credit: Tom Terwilliger, Los Alamos NL