

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

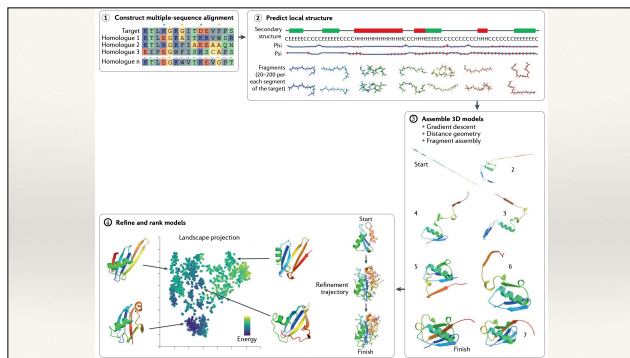
Ab initio prediction before AlphaFold

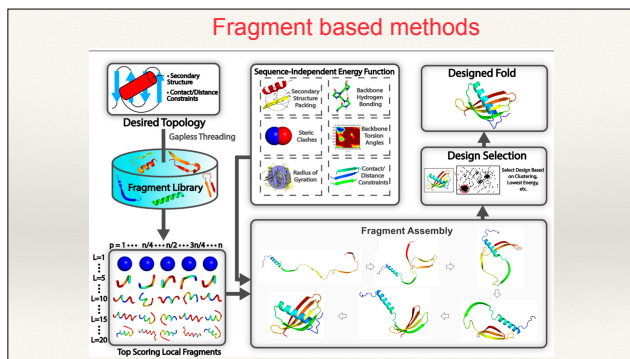
Ab initio prediction: Predicting Contacts

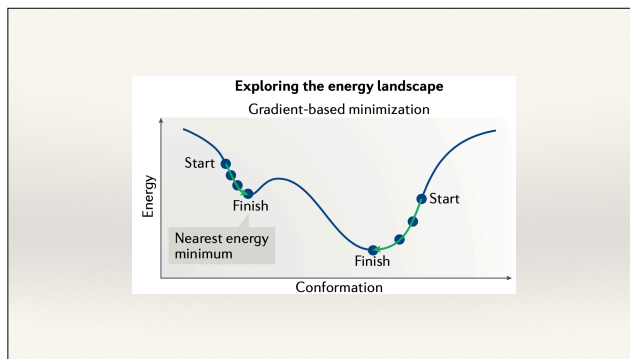
AlphaFold 1

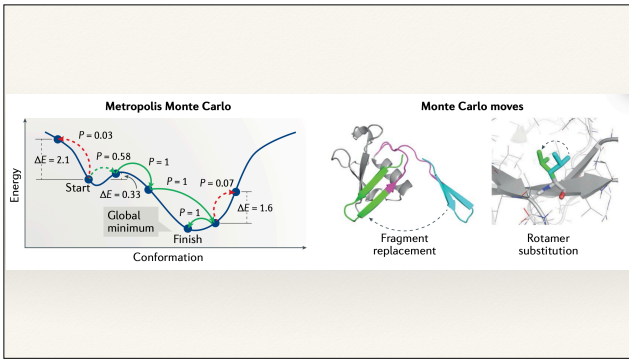
AlphaFold 2











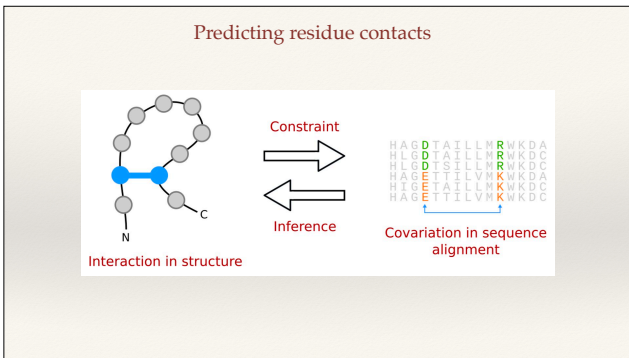
Ab initio Protein Structure Prediction

Ab initio prediction before AlphaFold

Ab initio prediction: Predicting Contacts

AlphaFold 1

AlphaFold 2



Predicting residue contacts

1. Given a multiple sequence alignment (MSA):

X_1	H	A	G	D	T	A	I	L	L	M	R	K	K	D	A
	H	L	G	D	T	A	I	L	L	M	R	K	K	D	C
	H	L	G	D	T	A	I	L	L	M	R	K	K	D	C
X_W	H	A	G	E	E	T	A	I	L	V	M	K	K	D	A
	H	A	G	E	T	A	I	L	V	M	K	K	D	C	

2. Compute "mean" sequence and covariance matrix:

$$\bar{X} = \frac{1}{N} \sum_{a=1}^N X_a$$

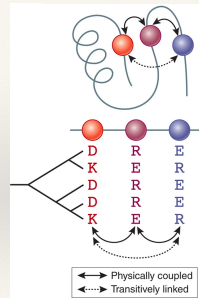
$$\bar{C} = C(MSA, \bar{X}) = \frac{1}{N} \sum_{a=1}^N (X_a - \bar{X})^T (X_a - \bar{X})$$

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

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

Predicting residue contacts

No! We need to pay attention to indirect effects:



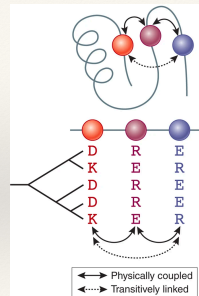
Predicting residue contacts

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 μ and a covariance matrix Σ , with the probability:

$$P(X_i | \mu, \Sigma) = (2\pi)^{-d} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right]$$



Predicting residue contacts

No! We need to pay attention to indirect effects:

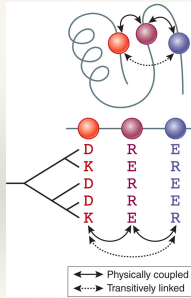
Gaussian model:

Each sequence X_i in the MSA is drawn from a multivariate Gaussian distribution characterized by a mean vector μ and a covariance matrix Σ , with the probability:

$$P(X_i | \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right]$$

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

$$P(\text{MSA} | \mu, \Sigma) = \prod_{i=1}^N P(X_i | \mu, \Sigma)$$



↔ 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 μ and a covariance matrix Σ , with the probability:

$$P(X_i | \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right]$$

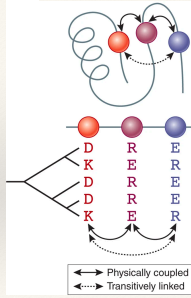
Assuming that the N sequences in the MSA are statistical independent, the probability, or likelihood of the data under this model is given by

$$P(\text{MSA} | \mu, \Sigma) = \prod_{i=1}^N P(X_i | \mu, \Sigma)$$

Using the maximum likelihood estimator for this probability

$$\mu = \bar{X}$$

$$\Sigma = \bar{C} = C(\text{MSA}, \bar{X})$$



↔ Physically coupled
 ⋈ Transitively linked

Predicting residue contacts

No! We need to pay attention to indirect effects:

Gaussian model:

$$P(X_i | \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (X_i - \mu)^T \Sigma^{-1} (X_i - \mu) \right]$$

$$\mu = \bar{X} \quad \Sigma = \bar{C} = C(\text{MSA}, \bar{X})$$

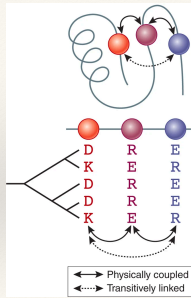
Note that:

$$(X_i - \mu)^T \Sigma^{-1} (X_i - \mu) = \sum_{k=1}^N \sum_{l=1}^N (X_i - \mu)_k (\Sigma^{-1})_{kl} (X_i - \mu)_l$$

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

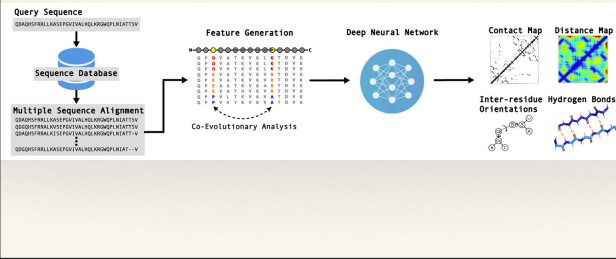
Therefore:

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

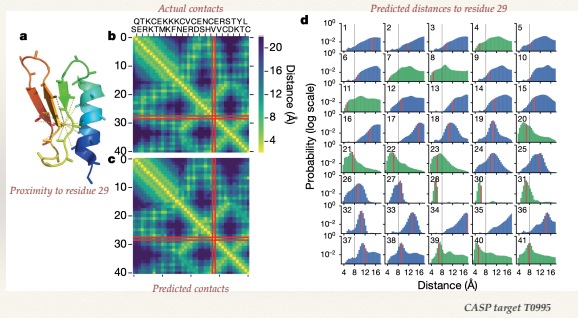


↔ Physically coupled
 ⋈ Transitively linked

Predicting residue contacts



Predicting residue contacts: How well does it work?



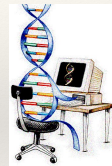
Ab initio Protein Structure Prediction

Ab initio prediction before AlphaFold

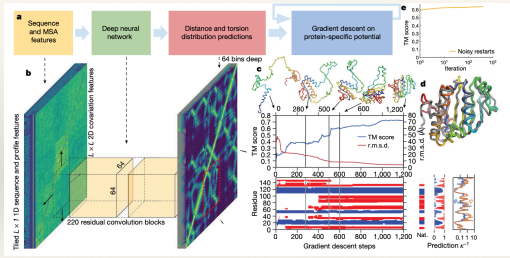
Ab initio prediction: Predicting Contacts

AlphaFold 1

AlphaFold 2



AlphaFold 1



AlphaFold1

AlphaFold 1

Reminder:

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

-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:

$$RMS(A, B) = \sqrt{\frac{1}{N} \sum_{i=1}^N d(a_i, b_i)^2}$$

Compute TM score:

$$TM(A, B) = \frac{1}{N} \sum_{i=1}^N \frac{1}{1 + \left(\frac{d(a_i, b_i)}{d_0(N)}\right)^2}$$

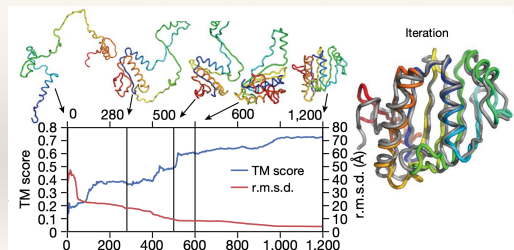
with $d_0(N) = 1.24\sqrt{N-15} - 1.8$

$d(a_i, b_i)$ is the Euclidian distance between a_i and b_i after optimal alignment of B onto A

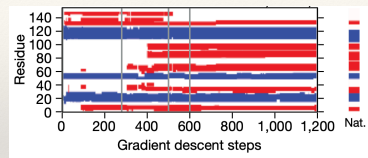
RMS: the lower, the better

TM: between [0,1]; the higher the better

AlphaFold 1

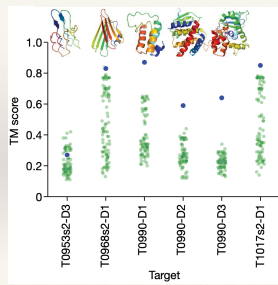


AlphaFold 1



Helix in blue, strand in red

AlphaFold 1: Success



Ab initio Protein Structure Prediction

Ab initio prediction before AlphaFold

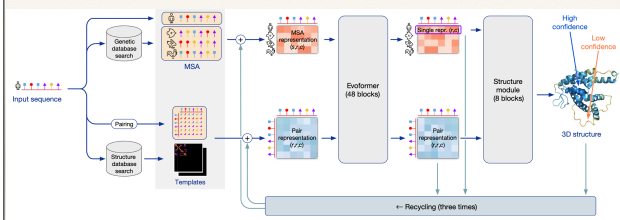
Ab initio prediction: Predicting Contacts

AlphaFold 1

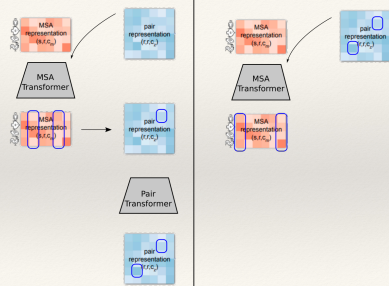
AlphaFold 2



AlphaFold 2

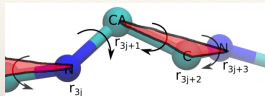


AlphaFold 2: some intuition



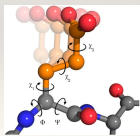
AlphaFold 2: the structure module

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

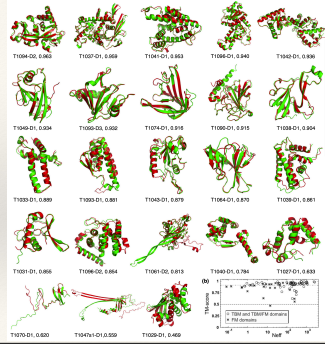


$$M = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Predicting side chains:

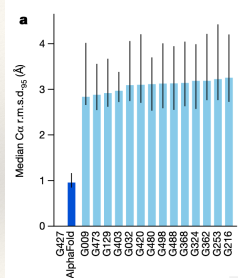


Successes at CASP14



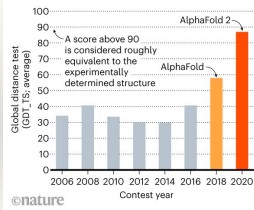
TBM: template-based modeling
FM: free modeling

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.



©nature

Training

- Sequence
- Multiple sequence alignment
- 3D structure



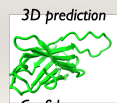
21 million parameters

Prediction

- Sequence
- Multiple sequence alignment

21 million parameters

Focus attention on important relationships



Confidence estimates

Credit: Tom Terwilliger, Los Alamos NL

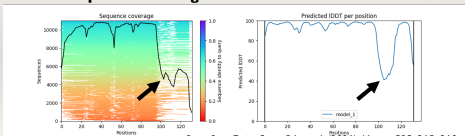
Multiple sequence alignment

```
EPGQYFSSGSSGSLVPSGSLGSLACAGGDFRFSSEVFHFFNAPAGLQGVNVT  
.....  
.....  
.....
```

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

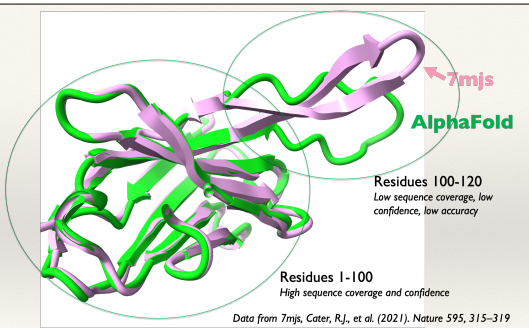
All sequences in alignment should be compatible with the right structure

Sequence coverage → **Confidence**



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

Credit: Tom Terwilliger, Los Alamos NL



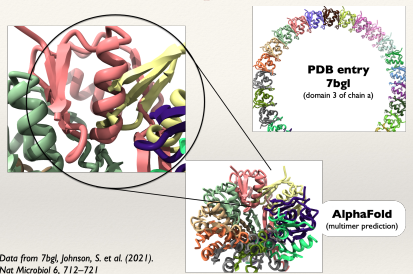
Residues 100-120
Low sequence coverage, low confidence, low accuracy

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



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

Credit: Tom Terwilliger, Los Alamos NL

