# Repliement des protéines
## avec la révolution d'AlphaFold2

Dirk Stratmann (dirk.stratmann@upmc.fr)

http://www.impmc.upmc.fr/~stratmann/

IMPMC, Sorbonne Université

octobre 2022

# Plan

# Processus du repliement

# Repliement



Disorder → Order

extended — transient secondary structure — compact globule — molten globule — disordered loop — folded protein
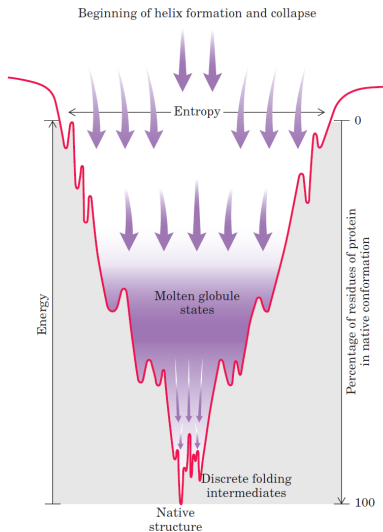
Compaction

- Molten globule: interactions hydrophobes
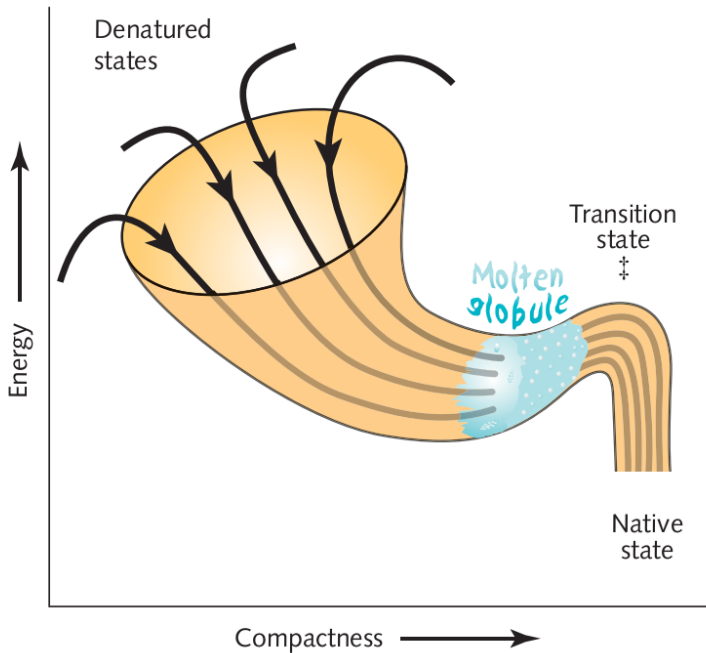- Structure 3D finale: liaisons hydrogènes

## Pourquoi est le repliement des protéines si rapide ?

Paradoxe de Levinthal (1968):

- chaque résidu n'a que deux possibilité de conformation
- une protéine de 100 résidu aurait $2^{100} \approx 10^{30}$ conformations possibles.
- Conclusion: une protéine ne peut pas se replier par une recherche au hasard de la conformation native
- Elle doit suivre un chemin de repliement (*folding pathway*) plus efficace.
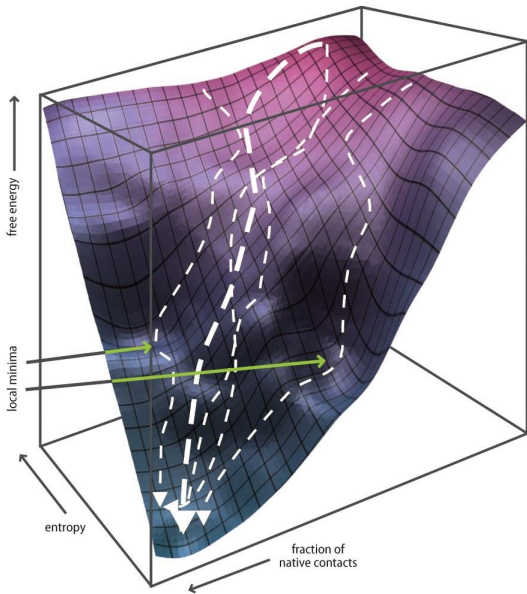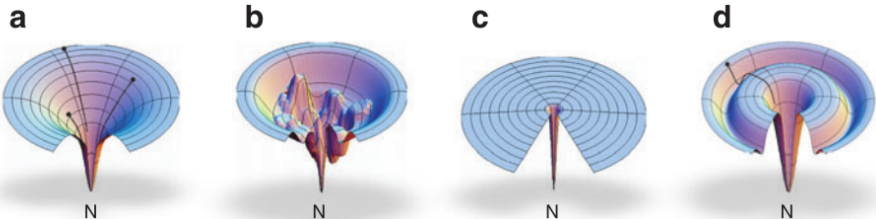
# Entonnoir du repliement

Denatured states

Transition state ‡

Molten globule

Native state

Energy

Compactness

8

Figure 6.15 How Proteins Work (©2012 Garland Science)

# Entonnoir du repliement
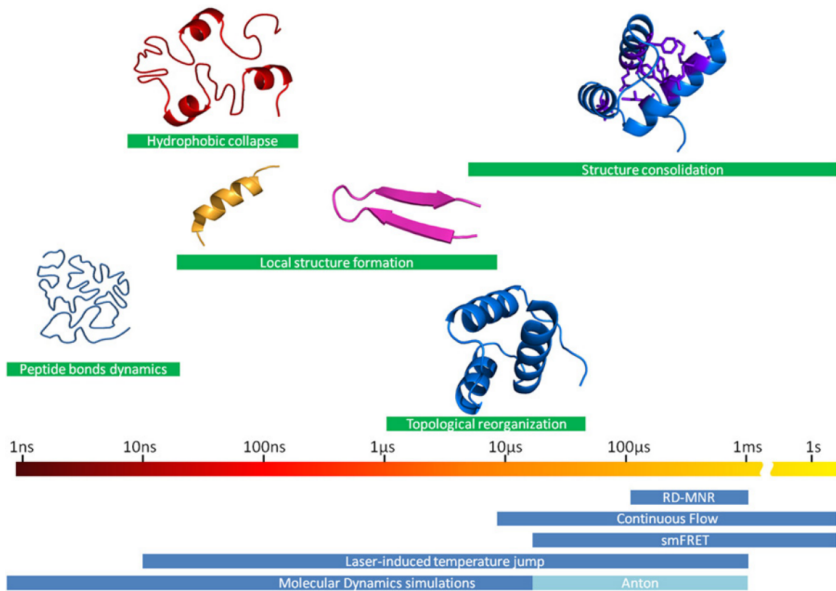


a: fast folder
b: kinetic traps
c: diffusional conformational search
d: obligatory intermediate

Ken A. Dill et al. (June 2008). en. In: *Annual Review of Biophysics* 37.1

# Vitesses du repliement

- Hélices se forment plus rapidement que des feuillet
- Hélices <=> contacts proches dans la séquence
- Feuillets <=> contacts distants dans la séquence
- Hélices: 0.1 - 1 $\mu s$, $\beta$-hairpins: 1 - 10 $\mu s$
- La vitesse 'limite' du repliement: N/50 $\mu s$, N: nombre de résidus

Hydrophobic collapse

Structure consolidation

Local structure formation

Peptide bonds dynamics

Topological reorganization

| 1ns | 10ns | 100ns | 1µs | 10µs | 100µs | 1ms | 1s |

RD-MNR

Continuous Flow

smFRET

Laser-induced temperature jump

Molecular Dynamics simulations    Anton

# État déplié

# Rayon de giration

| Protein | Number of residues | $R_G$ (native) | $R_G$ (denatured) | Ratio |
|---|---|---|---|---|
| PI3 kinase, SH3 domain | 90 | 18.6 | 27.5 | 0.7 |
| Horse heart cytochrome $c$ | 104 | 17.8 | 32.6 | 0.5 |
| Hen egg white lysozyme | 129 | 20.5 | 34.6 | 0.6 |
| Yeast triose phosphate isomerase | 247 | 29.7 | 49.7 | 0.6 |

- Si centre de gravité à l'origine, alors:

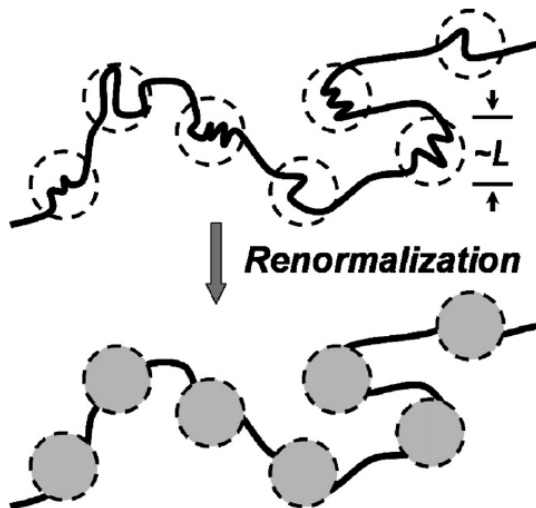$$R_G = \sqrt{\frac{\sum_i m_i(x_i^2 + y_i^2 + z_i^2)}{\sum_i m_i}}$$

- Random coil: $R_G \propto N^{0.6}$ avec N nombre de résidus
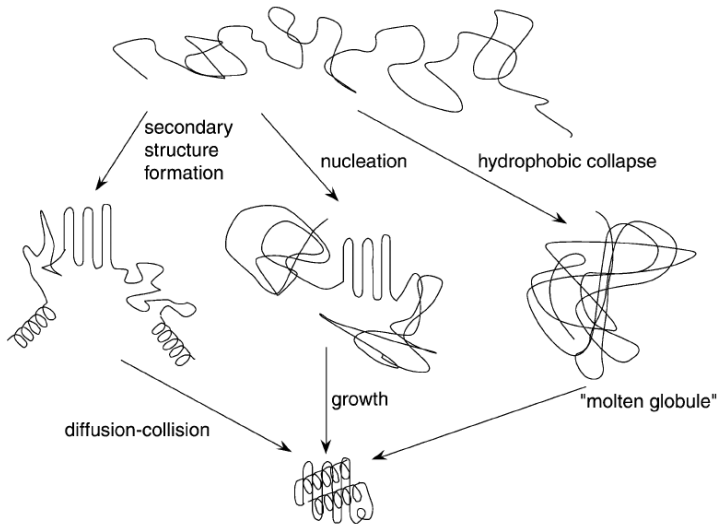- Globule compacte: $R_G \propto N^{0.33}$

# Rayon de giration

# Protéine dépliée = polymer sans interaction



Yiwen Chen et al. (Jan. 2008). In: *Archives of Biochemistry and Biophysics*. Highlight Issue: Protein Folding 469.1
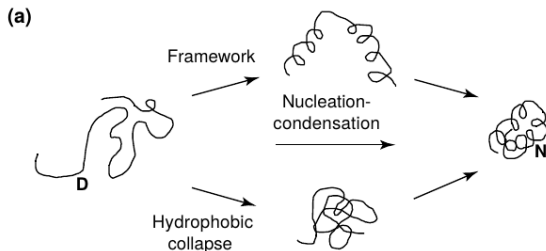
**16**

# Modèles classiques pour le repliement



à gauche: "framework" ou "hierarchical", Alan R. Fersht and Valerie Daggett (2002). In: *Cell*
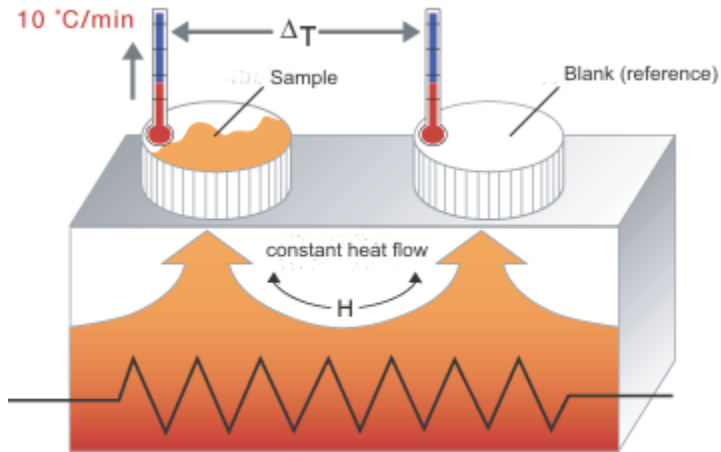108.4

# Modèles classiques pour le repliement



**(a)**

Framework

Nucleation-
condensation

Hydrophobic
collapse

D

N

**(b)**

Framework — engrailed homeodomain (α)
protein A (α)

Nucleation- — CI2 (α/β)
condensation — cMyb (a)
tenascin (β)
FKBP12 (α/β)

Barnase (α/β)
Hen egg white lysozyme (α/β)

Hydrophobic — ?
collapse

# Méthodes thérmodynamiques

# Differential scanning calorimetry (DSC)



Principe: référence et échantillion sont chauffés simultanément
Protéine absorbe de la chaleur lors du dépliement => $\Delta T$
www.itc.tu-bs.de/Abteilungen/Makro/Methods/dsc.htm
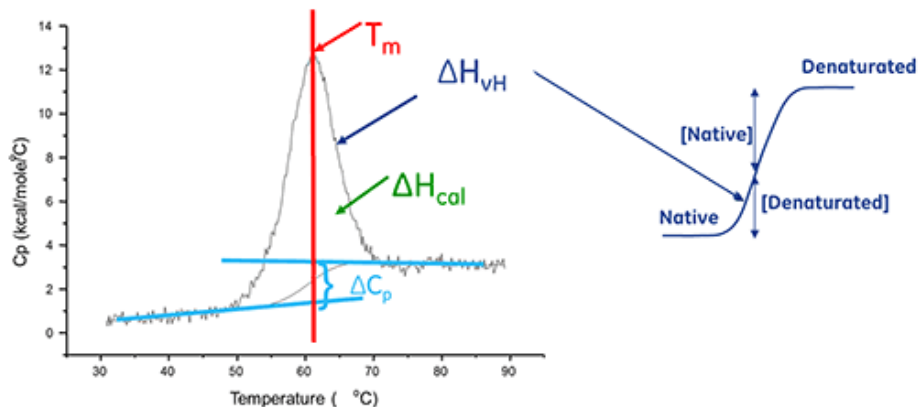
# Differential scanning calorimetry (DSC)



Tm: température médiane de transistion thermique (thermal transition midpoint)
Cp: capacité calorifique (heat capacity)
H: enthalpie
www.malvern.com/fr/products/technology/differential-scanning-calorimetry

# Differential scanning calorimetry (DSC)



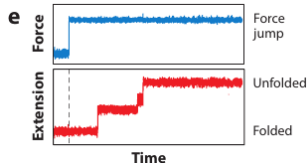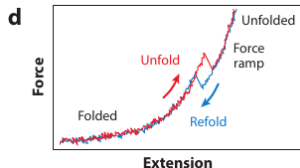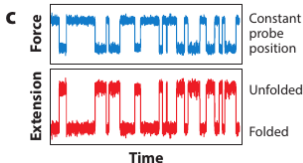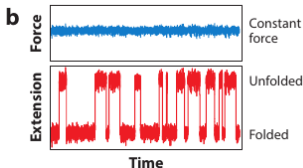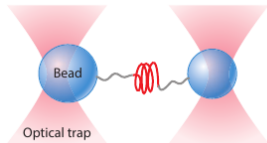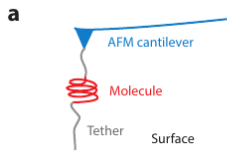Tm: température médiane de transistion thermique (thermal transition midpoint)
Cp: capacité calorifique (heat capacity)
H: enthalpie
www.malvern.com/fr/products/technology/differential-scanning-calorimetry

**23**

# Techniques sur molécule unique

# Single molecule force spectroscopy (SMFS)



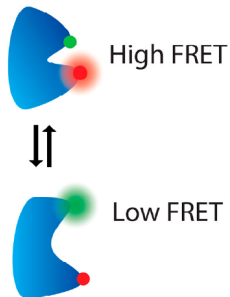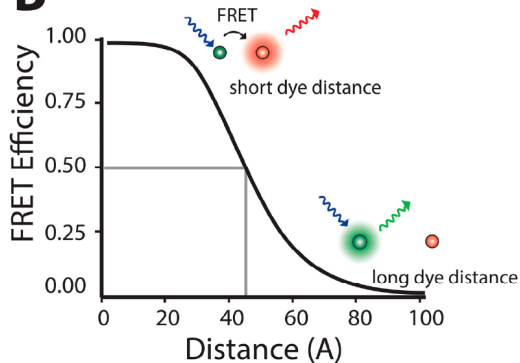a) Atomic force microscopy (AFM) and Optical tweezers
b) Constant force mode: extension fluctuates
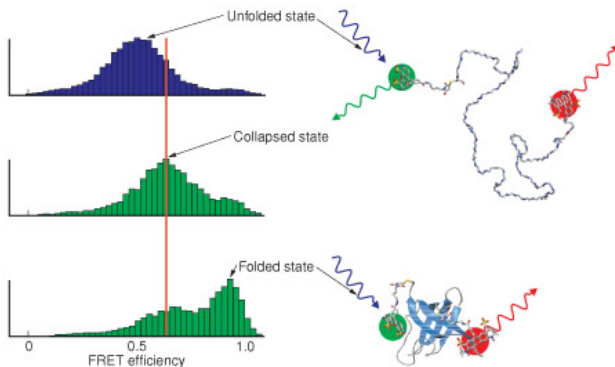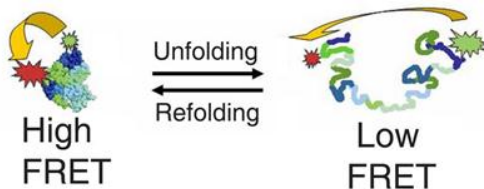c) Constant probe position: force and extension fluctuate
d) Force-ramp mode: elastic stretching is interrupted by a "rip", hysteresis indicates a nonequilibrium process
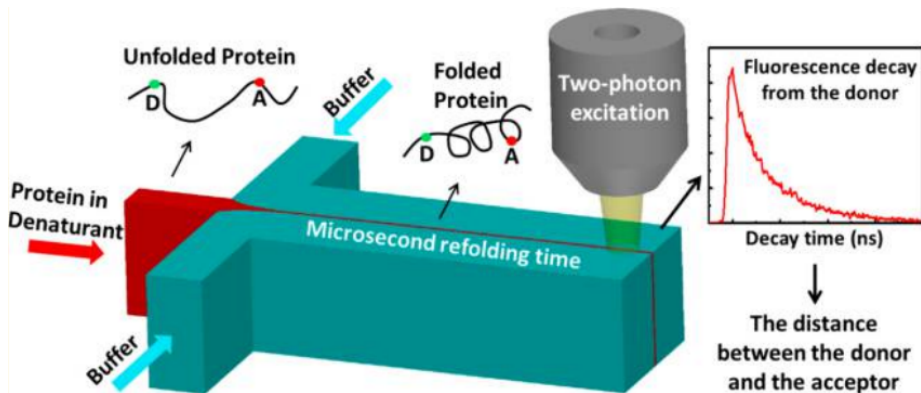e) Force-jump mode: extension changes in steps

# Fluorescence resonance energy transfer (FRET)
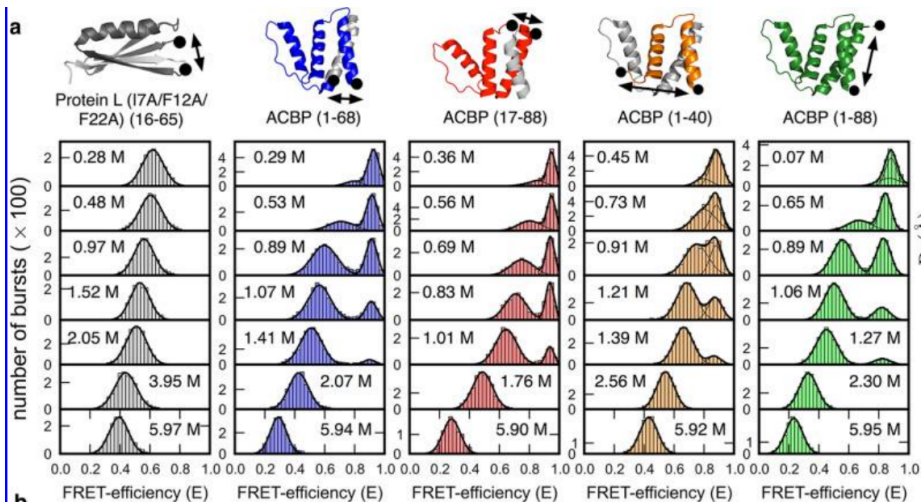
# Fluorescence resonance energy transfer (FRET)
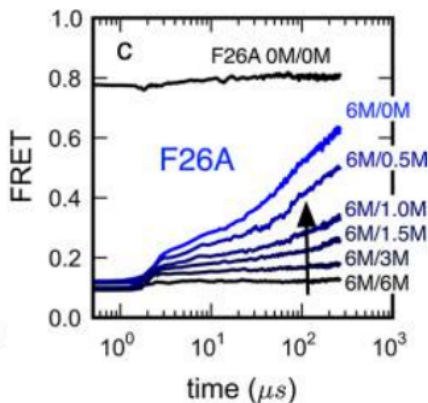
# Time-resolved FRET (TR-FRET)

# FRET sur ACBP

# FRET sur ACBP - Collapse très lent

Solution dénaturant (6 M GuHCl) vers solution de repliement (0 M GuHCl):



Indication pour: Formation d'une ensemble compact et hétérogènes de structures dépliées Vincent A. Voelz et al. (Aug. 2012). In: *J. Am. Chem. Soc.* 134.30

**30**

# RMN

# Pression et RMN



Kazuyuki Akasaka, Ryo Kitahara, and Yuji O. Kamatari (Mar. 2013). In: *Archives of Biochemistry and Biophysics.* Protein Folding and Stability 531.1–2

# Pression et RMN



Kazuyuki Akasaka, Ryo Kitahara, and Yuji O. Kamatari (Mar. 2013). In: *Archives of Biochemistry and Biophysics*. Protein Folding and Stability 531.1–2

# Pression et RMN



**B**

α-Domain intermediate
at 2000 bar & 5°C

Unfolded state

α/β Intermediates
N''' at 2000 bar & 10°C
N'' at 2000 bar & 15°C
N' at 2000 bar & 25°C

$Q_\beta$

$Q_\alpha$

F

Native state (N)
at 1 bar & 25°C
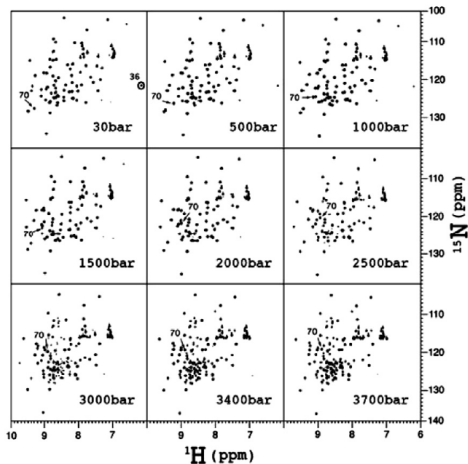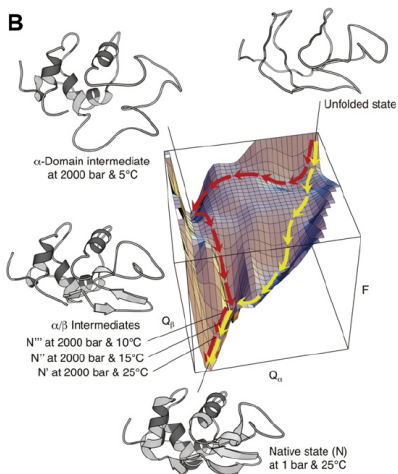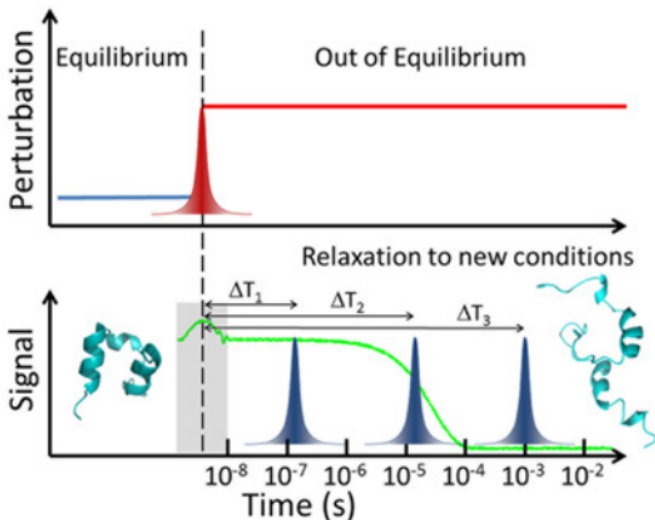
Kazuyuki Akasaka, Ryo Kitahara, and Yuji O. Kamatari (Mar. 2013). In: *Archives of Biochemistry and Biophysics*. Protein Folding and Stability 531.1–2
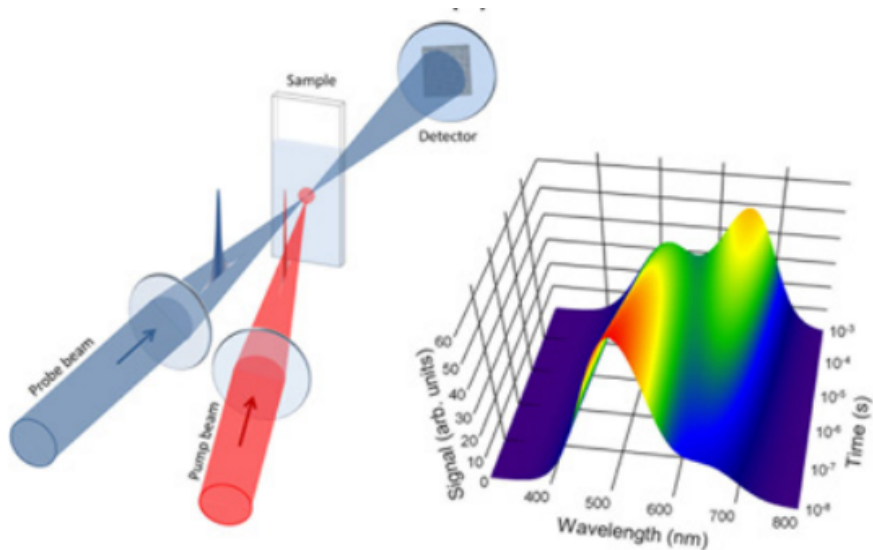
# Méthodes pour cinétique "ultra rapide"

# Laser T-jump



Ultrafast Kinetic Perturbation Methods

# Laser T-jump



J7

# Introduction

# Anton



## Molecular Dynamics Simulations

### A. Force Field

$$U = \sum k_b (r - r_0)^2 + \sum k_\theta (\theta - \theta_0)^2 + \sum A[1 + \cos(nT - \varphi)] + \sum\sum q_i q_j / r_{ij} + \sum\sum B\left[\left(\frac{\sigma_{ij}}{r_{ij}}\right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}}\right)^{6}\right]$$

bond stretching    bending    torsional    electrostatics    van der Waals

### B. Simulation System

### C. Specialized High Performance Computing

ANTON2

### D. Molecular Trajectories

### E. Free Energy Surface calculation

# Temps de repliement

Current Opinion in Structural Biology

Thomas J Lane et al. (Feb. 2013). en. In: *Current Opinion in Structural Biology* 23.1

Thomas J Lane et al. (Feb. 2013). en. In: *Current Opinion in Structural Biology* 23.1

**43**

# Markov State Models (MSM)

protéine NTL9, Gregory R Bowman, Vincent A Voelz, and Vijay S Pande (Feb. 2011). en. In: *Current Opinion in Structural Biology* 21.1

protéine ACBP, Thomas J Lane et al. (Feb. 2013). en. In: *Current Opinion in Structural Biology* 23.1

# Etat déplié

# Etat déplié: difficle en MD



Current Opinion in Structural Biology

**48**

# Thermodynamique du repliement

# Prédiction de l'enthalpie du dépliement avec MD



Current Opinion in Structural Biology

**50**

# Conclusions de Piana et al.

Peuvent les champs de forces en dynamique moléculaire reproduire les données expérimentales du repliement ?
Oui pour :

- Structure native, repliée
- Taux de repliement

Non pour:

- Cinétique détaillée
- Structures dépliées
- Enthalpie du l'état repliée plus bas qu'en expérimental

# Introduction

A — Protein Synthesis

B — Unfolded Protein

C — Folded Protein

D — Functional Complex

E — Misfolded Protein

F — Aggregated Protein

G — Degraded protein

Proteasome

Chaperones

# Environnement cellulaire

# Molecular crowding

Molecular crowding

Localization

Cytoplasm

Organelles

Nucleus

Ligand binding

Osmolyte uptake

Posttranslational modifications

Volume increase (dilution)

External conditions (e.g., temperature, pressure, and osmolarity)

Complex formation

# David Baker - first work

# Contact order and folding kinetics



Low contact order

High contact order

$$CO = \frac{1}{L \cdot N} \sum^{N} \Delta S_{i,j} \qquad (1)$$

where $N$ is the total number of contacts, $\Delta S_{i,j}$ is the sequence separation, in residues, between contacting residues $i$ and $j$, and $L$ is the total number of residues in the protein.

Kevin W Plaxco, Kim T Simons, and David Baker (Apr. 1998). In: *Journal of Molecular Biology* 277.4

# Contact order and folding kinetics

# *De novo* protein design



**Structure prediction**
Sequence known, structure unknown

**Input:** Known amino-acid sequence

**Backbone sampling** Guided by local native sequence

**Side-chain sampling** Rotamers of native amino acids

**Output:** Predicted structure

**Fixed-backbone design**
Sequence unknown, structure known

Known backbone structure

**Backbone sampling** None

**Side-chain sampling** Rotamers of all amino acids

Designed sequence

***De novo* design**
Sequence unknown, structure unknown

Architecture definition

**Backbone sampling** Sequence independent

**Side-chain sampling** Rotamers of all amino acids

Designed backbone and designed sequence

Po-Ssu Huang, Scott E. Boyken, and David Baker (Sept. 2016). en. In: *Nature* 537.7620

# Protein design - first sucess story



**Fig. 1.** A two-dimensional schematic of the target fold (hexagon, strand; square, helix; circle, other). Hydrogen bond partners are shown as purple arrows. The amino acids shown are those in the final designed (Top7) sequence.

# Protein design - first sucess story



Brian Kuhlman et al. (Nov. 2003). en. In: *Science* 302.5649

**64**

# David Baker - 15 years of protein design

# Protein design - examples

# Protein design - examples



Po-Ssu Huang, Scott E. Boyken, and David Baker (Sept. 2016). en. In: *Nature* 537.7620

**67**

# Protein design - examples

# Protein design - examples

# Protein design - examples

# Protein design - macrocycles

# Protein design - macrocycles



AABBXBY    AABBX*ABY    AABBXXABY    AABBXXA*BBY

Parisa Hosseinzadeh et al. (Dec. 2017). en. In: *Science* 358.6369

# Protein design - macrocycles

# Protein design - macrocycles



Parisa Hosseinzadeh et al. (Dec. 2017). en. In: *Science* 358.6369

# Peptides cycliques ciblant des interactions protéine-protéine

# Drug design: cibler les interactions protéine-protéine

# Drug = peptide cyclique

# Drug = peptide cyclique

# Alpha Fold - La révolution

# Le press release de CASP du 30 nov 2020

"Artificial intelligence solution to a 50-year-old science challenge could 'revolutionise' medical research"

"Today (Monday) researchers at the 14th Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction (CASP14) will announce that an artificial intelligence (AI) solution to the challenge has been found."
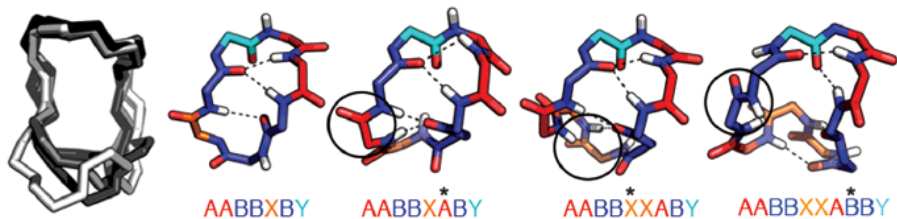
"Nearly 50 years ago, Christian Anfinsen was awarded a Nobel Prize for showing that it should be possible to determine the shape of proteins based on their sequence of amino acids – the individual building blocks that make up proteins. That's why our community of scientists have been working on the biennial CASP challenge."

```
https://predictioncenter.org/casp14/doc/CASP14_
press_release.html
```

# DeepMind de google basé à Londres

## Le relais immédiat dans Nature de l'annonce de CASP

30 nov 2020, Nature, "'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures"
https://www.nature.com/articles/d41586-020-03348-4

# La presse

30 nov 2020, The New York Times, London AI claims breakthrough that could accelerate drug discovery

1 déc 2020, France Culture, Le "problème de repliement des protéines" résolu par une intelligence artificielle

11 déc 2020, Les Échos, DeepMind met les chercheurs du monde entier au tapis

# CASP 14 et AlphaFold 2



**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.

# CASP 14 et 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 |
| 2 | 473 | BAKER | 92 | 90.8241 | 2 | 0.9872 | 2 | 92.1241 | 2 | 1 |
| 3 | 403 | BAKER-experimental | 92 | 88.9672 | 3 | 0.9670 | 3 | 91.4731 | 3 | 0 |
| 4 | 480 | FEIG-R2 | 92 | 72.5351 | 4 | 0.7884 | 4 | 74.5627 | 4 | 0 |
| 5 | 129 | Zhang | 92 | 67.9065 | 5 | 0.7381 | 5 | 68.8922 | 5 | 0 |

# CASP 13 et AlphaFold 1



| # | 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 | 043 | A7D | 104 | 120.4307 | 1 | 1.1580 | 1 | 128.0693 | 1 |
| 2 | 322 | Zhang | 104 | 107.5948 | 2 | 1.0346 | 2 | 108.1948 | 2 |
| 3 | 089 | MULTICOM | 104 | 99.4661 | 3 | 0.9564 | 3 | 99.9886 | 3 |
| 4 | 145 | QUARK | 104 | 90.9915 | 4 | 0.8749 | 4 | 91.5625 | 4 |
| 5 | 261 | Zhang-Server | 104 | 88.9540 | 5 | 0.8553 | 5 | 89.7597 | 5 |

## La presse en 2021

3 oct 2021, Forbes: "AlphaFold is the most important achievement in AI - Ever"

https://www.forbes.com/sites/robtoews/2021/10/03/
alphafold-is-the-most-important-achievement-in-ai-ever
?sh=2857ff656e0a

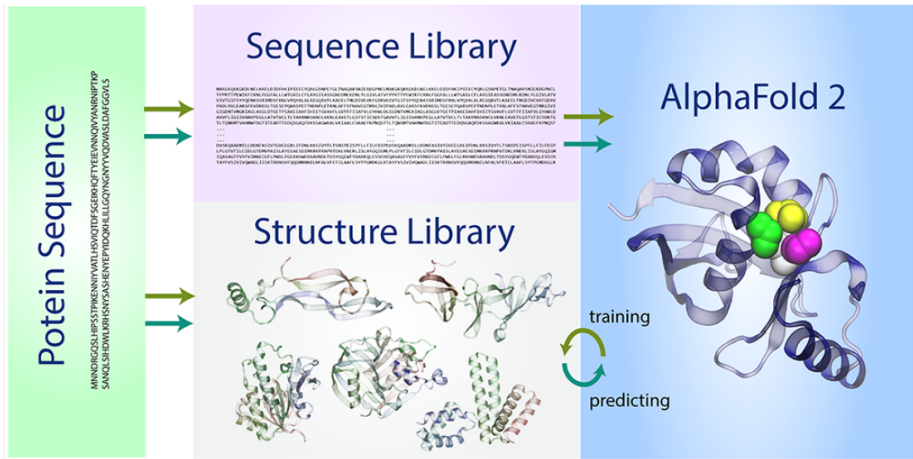22 juillet 2021, Fortune: "In giant leap for biology, DeepMind's A.I. reveals secret building blocks of human life"

https://fortune.com/2021/07/22/
deepmind-alphafold-human-proteome-database-proteins/

# Principe de Alpha Fold 2

# 15 juillet 2021: Détails de la méthode et code source

**Article**

# Highly accurate protein structure prediction with AlphaFold

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

+5000 citations sur google scholar à ce jour (6 oct 2022) !

# Principe de Alpha Fold 2

# Applications de Alpha Fold

# Alpha Fold DB

The international journal of science / 26 August 2021

# nature

**outlook**
Sickle-cell
disease

# PROTEIN POWER

AI network predicts highly
accurate 3D structures
for the human proteome

**Troubled waters**
The race to save the
Great Barrier Reef
from climate change

**Coronavirus**
Time is running out
to find the origins
of SARS-CoV-2

**Storage hunting**
Quantifying carbon
held in Africa's
montane forests

# Autres méthodes utilisant le deep learning

# Le pionnier: RaptorX de Jinbo Xu et al.

Serveur RaptorX:
http://raptorx.uchicago.edu/
CASP 14: score 38 vs 244 pour AlphaFold 2
"This server was officially ranked 1st in contact prediction in both
CASP12 and CASP13 and initiated the revolution of protein structure
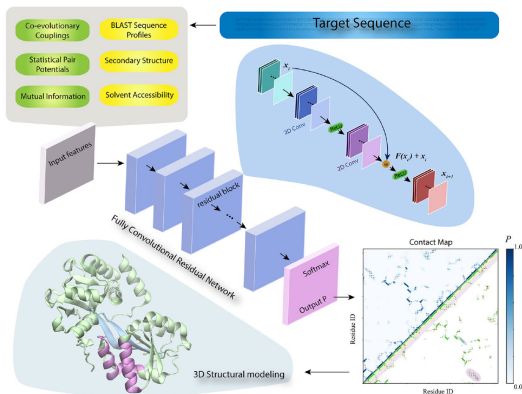prediction by deep learning."

# DESTINI de Jeffrey Skolnick et al.

Serveur DESTINI:

`https://sites.gatech.edu/cssb/destini/`

CASP 14: score 29 vs 244 pour AlphaFold 2

# Le rattrapage de David Baker



**RoseTTAFold: Accurate protein structure prediction accessible to all**

# Le rattrapage de David Baker

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

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

publié le même jour que le Nature sur AlphaFold2 (15 juillet 2021) !
mais "que" 1120 citations sur google scholar à ce jour (6 oct 2022)

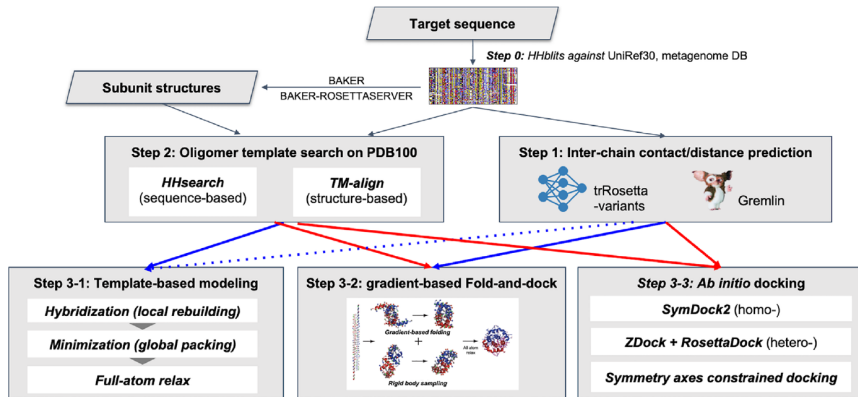**A** 2-track block [yellow box] 3-track block [blue box]

Cropped MSA
1D track
Attention

Outer product & aggregation
Attention

Cropped sequence & templates
2D track
Attention

Graph-Transformer

Masked Attention
Outer product & aggregation
Attention

SE(3)-Transformer
3D track

BB-only model
SE(3) iterative refinement
SE(3)-Transformers
Combine all crops
Gradient-based folding
Full atom model

**B** CASP14 targets

BAKER-ROSETTASERVER
Zhang-server
BAKER (human)
2-track (pyRosetta)
RoseTTAFold (end-to-end)
RoseTTAFold (pyRosetta)
AlphaFold2 (human)

TM-score (40 to 100)

**C** CAMEO targets

RoseTTAFold
Robetta
IntFOLD6-TS
BestSingleTemplate
SWISS-MODEL

TM-score (55 to 85)

**100**

# Next frontier: Oligomères



**FIGURE 1** The oligomer structure modeling procedure used by the BAKER-experimental group

# RosettaFold: github et serveur

```
https://github.com/RosettaCommons/RoseTTAFold
https://robetta.bakerlab.org/
```
voir aussi ici:
```
https://www.rosettacommons.org/
```

# Comprendre Alpha Fold

## Le blog de "Oxford Protein Informatics Group"

```
https://www.blopig.com/blog/2021/07/
alphafold-2-is-here-whats-behind-the-structure-predict
```

"... we have many new questions. What is the secret sauce before the news splash, and why is it so effective? Is it a piece of code that the average user can actually run? What are AlphaFold 2's shortcomings? And, most important of all, what will it mean for computational biology? And for all of us?"

# Le blog de Mohammed AlQuraishi

```
https://moalquraishi.wordpress.com/2020/12/08/
alphafold2-casp14-it-feels-like-ones-child-has-left-ho
amp/
```

# Mutations corrélées, MSA
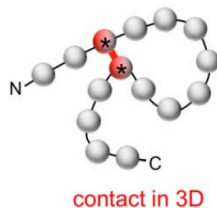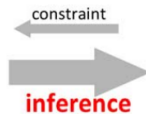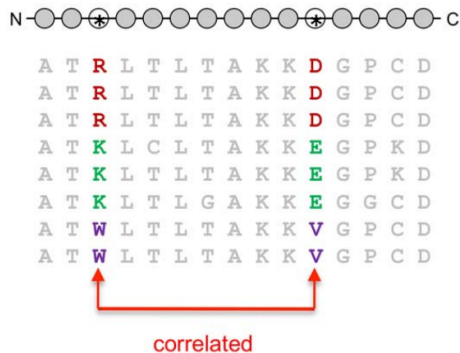
# Mutations corrélées, MSA

# Protein 3D Structure Computed from Evolutionary Sequence Variation

Debora S. Marks[1]*, Lucy J. Colwell[2], Robert Sheridan[3], Thomas A. Hopf[1], Andrea Pagnani[4], Riccardo Zecchina[4,5], Chris Sander[3]
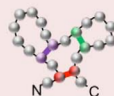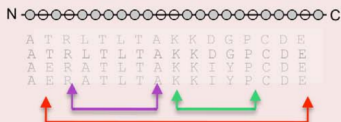
1 Department of Systems Biology, Harvard Medical School, Boston, Massachusetts, United States of America, 2 MRC Laboratory of Molecular Biology, Hills Road, Cambridge, United Kingdom, 3 Computational Biology Center, Memorial Sloan-Kettering Cancer Center, New York, New York, United States of America, 4 Human Genetics Foundation, Torino, Italy, 5 Politecnico di Torino, Torino, Italy

# Principe



correlated

# Protocole de 2011 sans deep-learning



Align evolutionary diverged sequences

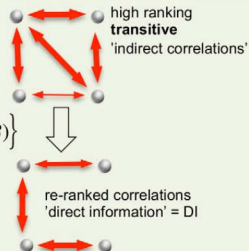Calculate covariance matrix for each pair of sequence positions for all pairs of amino acids (A,B)

$$C_{ij}(A,B) = f_{ij}(A,B) - f_i(A)P_j(B)$$

$$C_{ij}^{-1}(A,B) = -e_{ij}(A,B)_{i \neq j}$$

Identify maximally informative pair couplings using **statistical model** of entire protein to infer residue-residue co-evolution

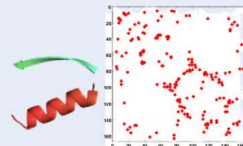$$P_{ij}^{Dir}(A,B) = \frac{1}{Z}\exp\left\{e_{ij}(A,B) + \tilde{h}_i(A) + \tilde{h}_j(B)\right\}$$

$$DI_{ij} = \sum_{A,B=1}^{q} P_{ij}^{Dir}(A,B)\ln\frac{P_{ij}^{Dir}(A,B)}{f_i(A)f_j(B)}$$

high ranking **transitive** 'indirect correlations'

re-ranked correlations 'direct information' = DI

# Protocole de 2011 sans deep-learning



Analyze the highest scoring pairs to produce ranked list of residue pairs which we predict to be close in 3D space. Use these pairs as predicted close "evolutionary inferred contacts" , EICs, in folding calculations

assign (resid 143 and name CA) (resid 123 and name CA)  4 4 3
assign (resid 16 and name CA) (resid 10 and name CA)  4 4 3
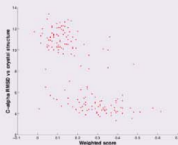assign (resid 141 and name CA) (resid 82 and name CA)  4 4 3
assign (resid 129 and name CA) (resid 87 and name CA)  4 4 3
assign (resid 92 and name CA) (resid 11 and name CA)  4 4 3
assign (resid 116 and name CA) (resid 81 and name CA)  4 4 3

predicted contacts (EICs)

Start with extended structure
use **distance geometry** and **simulated annealing** with predicted constraints, EICs, to fold the chain

Rank predicted structures using quality measure of backbone alpha torsion and beta sheet twist

good scores

bad scores

# Utiliser Alpha Fold

## google colab

```
https://colab.research.google.com/notebooks/intro.
ipynb?utm_source=scs-index
https://www.youtube.com/watch?v=inN8seMm7UI
```
Exemple pour faire de la MD ou aussi AlphaFold2 + MD:
```
https://github.com/pablo-arantes/Making-it-rain
```

## Alpha Fold avec google colab

une version amélioré d'alphaFold2:

https://colab.research.google.com/github/sokrypton/
ColabFold/blob/main/AlphaFold2.ipynb

ou l'original:

https://colab.research.google.com/github/deepmind/
alphafold/blob/main/notebooks/AlphaFold.ipynb

## Alpha Fold en local

On peux l'installer aussi sur une machine en local, mais il faut qu'elle soit très puissante:

https://github.com/deepmind/alphafold

" The simplest way to run AlphaFold is using the provided Docker script. This was tested on Google Cloud with a machine using the nvidia-gpu-cloud-image with 12 vCPUs, 85 GB of RAM, a 100 GB boot disk, the databases on an additional 3 TB disk, and an A100 GPU."

Rien que la carte graphique A100 coûte le prix d'une petite voiture:

https://fr.aliexpress.com/item/1005002408111365.html

# Pour aller plus loin: EMBL webinar

# EMBL webinar "How to interpret AlphaFold structures"
## 8 sept 2021

```
https://www.ebi.ac.uk/training/events/
how-interpret-alphafold-structures/
```

DeepMind

# Introduction to AlphaFold

**Presenter**: Kathryn Tunyasuvunakool
Research Scientist at DeepMind

# How does it work? (the short version)



See Jumper et al. 2021 (especially the SI) for details

118

# Inputs

Multiple Sequence Alignment

Input sequence

A key AlphaFold input is the MSA, containing sequences evolutionarily related to the target.
Related sequences are found using standard tools and public databases.

# Inputs

Multiple Sequence Alignment

Input sequence

Residue pairs

The input sequence is used to create an array of numbers representing all residue pairs.

# Inputs



AlphaFold can also use template structures from the PDB, found using standard tools.
However, it often produces accurate predictions without a template.

# Network



The Evoformer blocks extract information about the relationship between residues.
The MSA representation can update the pair representation and vice versa.

# Network

Input sequence

Genetic database search

Multiple Sequence Alignment

Residue pairs

Structure database search

Templates

MSA representation

Pair representation

Evoformer (48 blocks)

Structure Module (8 blocks)

The Structure Module predicts a rotation + translation to place each residue.
A small network predicts side chain chi angles. The final structure is run through a relaxation process.

# Network



Feeding certain outputs back through the network again improves performance

# Other outputs



As well as a predicted structure, AlphaFold produces two confidence estimates

# Interpreting predictions

The short version: use **both** confidence metrics

# Predicted LDDT: definition

AlphaFold's per-residue prediction of its lDDT-Cα score*

Roughly, lDDT measures the percentage of correctly predicted interatomic distances, not how well the predicted and true structures can be superimposed.

It rewards **locally correct** structures, and **getting individual domains right**. pLDDT behaves similarly, as a measure of **local confidence**



Alignment-based metric      lDDT

*Mariani et al. 2013

# Predicted LDDT: format

pLDDT ranges from 0 to 100 (100 is most confident)

We use a consistent "confidence bands" color scheme when displaying predictions

A pLDDT plot is also displayed by some of our tools

Prediction files always contain pLDDT in the B-factors
Therefore a **higher** B-factor is better!



- ■ Very high (>90)
- ■ Confident (70–90)
- ■ Low (50–70)
- ■ Very low (<50)

```
MODEL        0
ATOM     1  N   MET A   1      -9.212  -5.798  33.490  1.00 63.75           N
ATOM     2  CA  MET A   1     -10.075  -6.130  32.322  1.00 63.75           C
ATOM     3  C   MET A   1     -11.469  -6.615  32.714  1.00 63.75           C
ATOM     4  CB  MET A   1      -9.419  -7.112  31.341  1.00 63.75           C
ATOM     5  O   MET A   1     -12.429  -6.075  32.184  1.00 63.75           O
ATOM     6  CG  MET A   1      -8.311  -6.411  30.547  1.00 63.75           C
ATOM     7  SD  MET A   1      -7.766  -7.280  29.061  1.00 63.75           S
ATOM     8  CE  MET A   1      -7.045  -8.751  29.798  1.00 63.75           C
ATOM     9  N   ALA A   2     -11.624  -7.579  33.634  1.00 66.38           N
ATOM    10  CA  ALA A   2     -12.951  -8.096  34.007  1.00 66.38           C
```

# Predicted LDDT: usage

**Identifying domains & possible disordered regions**

**Assessing confidence within a domain**



**pLDDT > 70**
Residues 65–342
and 418–784 form
a confident domain

**pLDDT < 50**
A disorder
prediction not
a structure
prediction

**pLDDT > 90**
Reasonable to
investigate side
chains / active
site details

**pLDDT > 70**
Lower confidence on
these specific parts

# Predicted LDDT: pitfalls

High pLDDT on all domains does **not** imply AlphaFold is confident of their relative positions

# Predicted Aligned Error: definition

AlphaFold's prediction of its position error at residue $x$,
if the predicted and the true structures were aligned on residue $y$

PAE aims to measure confidence in the **relative positions** of **pairs of residues**

Mainly used to assess relative domain positions, but applicable whenever pairwise confidence is relevant

# Predicted Aligned Error: format

PAE is displayed as a 2D plot.

Suppose residue y were aligned to the true structure and we measured the position error at residue x. The color at (x, y) is AlphaFold's prediction of that error.

In this case the squares correspond to two domains.



Residues 400–722

Residues 1–375



Expected position error (Ångströms)

Aligned residue

Scored residue

# Predicted Aligned Error: format

PAE is displayed as a 2D plot.

Suppose residue y were aligned to the true structure.
And we measured the position error at residue x.
The color at (x, y) is AlphaFold's prediction of that error.

AlphaFold is confident in relative positions within each domain.



Expected position error (Ångströms)
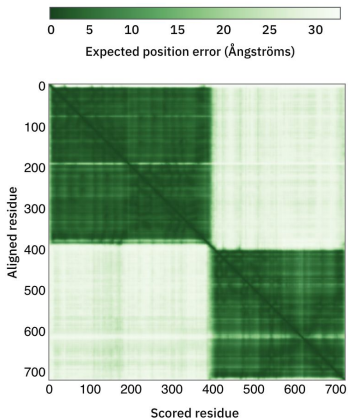
# Predicted Aligned Error: format



Expected position error (Ångströms)

# Predicted Aligned Error: format

PAE is displayed as a 2D plot.

Suppose residues y were aligned to the true structure
And we measured the position error at residue x
The color at (x, y) is AlphaFold's prediction of that error.

...but not between domains.

# Predicted Aligned Error: format

Expected position error (Ångströms)

# Predicted Aligned Error: usage

**No confidence in relative domain positions**



**Predicted domain packing**

# Things to be aware of

**Uncertain domain placement**

- If AlphaFold is uncertain, it won't necessarily place domains sensibly relative to each other
  - Membrane proteins won't leave space for the cell membrane
  - Clashes can occur

**Complexes**

- For proteins that exist in complex, AlphaFold is missing context about their binding partners
  - Heteromers more problematic than homomers
  - Worst case: the protein is flexible in isolation

- Some have had success predicting complexes by joining 2 sequences with a linker
  - We think it is possible to extend the ideas in AlphaFold to complexes
  - However, this linker setup remains to be benchmarked

# AlphaFold database

Sameer Velankar
EMBL-EBI

DeepMind

EMBL-EBI

## AlphaFold Database

- ~365K predicted structures for proteins from 21 model organisms

- For the organisms currently covered, predicted structures available for sequences in the UniProt reference proteome that are between 16 and 2700 amino acids long and contain only standard amino acids

- Only one prediction (out of 5 independent predictions) is made available in the release

- Accession – AF-P12345-F[1-N]

- Files
  - AF-P12345-F[1-N]-model_V[1-N].[pdb, mmcif]
  - AF-P12345-F[1-N]-predicted-aligned-error_V[1-N].json

DeepMind

EMBL-EBI

# AlphaFold web pages

- Basic search system

- Allows search using UniProt accession, UniProt id, protein name, gene name and organism



- Clear indication that the structure shown is a prediction

- Allow easy download of structure data

- Basic information about protein

- Clearly indicates if there are experimental structures available

- Display residue-quality information in 3D viewer (pLDDT – predicted Local Distance Difference Test)

- Predicted Aligned Error (PAE viewer)



DeepMind

# Models available across EBI resources



UniProt

Pfam

PDBe-KB

InterPro

DeepMind

DeepMind   EMBL-EBI

# AlphaFold Database – limitations

- Information on complexes with other proteins, nucleic acids (DNA or RNA) or ligands. In some cases, the single-chain prediction may correspond to the structure adopted in a complex. The missing context from surrounding molecules may lead to an uninformative prediction

  - AlphaFold does not make any predictions about any of the non-protein components such as cofactors, metals, ligands including drug-like molecules, ions, carbohydrates and other post-translational modifications

- Protein dynamics - AlphaFold will usually only produce one of multiple conformations

- AlphaFold has not been trained or validated for predicting the effect of mutations

- May (or may not) lead to hypotheses about protein function – any hypotheses have to be tested by further experimentation

DeepMind

EMBL-EBI

## What's next – under discussion

- Remove signal peptides from predictions

- Making 5 independent predictions available for each protein

- Additional metadata
  - MSA – need to consider data size
  - information on templates
  - quality criteria e.g. predicted TM score

- Updating database to UniRef90 dataset (~130 million structures)

DeepMind

EMBL-EBI

Impact of AlphaFold database on life science research

# Structural bioinformatics – (structure/function)

- Predicting complexes between macromolecules
  - Homo- and Hetero- Protein-protein; Protein-nucleic acid complexes
  - Intrinsically disordered proteins

- Provide information on protein dynamics
  - Relevant conformational states

- Functionally important residues
  - Impact of mutation; Binding sites; Conformationally important residues
  - Interfaces

- Ligand prediction – What binds?
  - What might bind in a pocket



DeepMind

EMBL-EBI

# Structural biology

- Accelerating structure studies
  - Improved construct design
  - Starting model for structure determination
  - Fitting models in low resolution EM maps
  - Time resolved studies to understand mechanism

# Structural biology

- Integrative/hybrid methods
  - Models for individual components

- Combination of sparse experimental data and predicted model may lead to actionable data to test hypothesis
  - Chemical foot printing
  - Hydrogen-Deuterium exchange
  - smFRET - Single molecule fluorescence resonance energy transfer



I/H Methods Structures
552-protein yeast Nuclear Pore Complex
Kim et al. (2018) *Nature* 555, 475-82
PDBDEV_00000010; PDBDEV_00000011; PDBDEV_00000012



Jochem Smit @Jhsmit · 23 Jul
Replying to @eitan_lerner
If we had a curated **smFRET** / structure database it could probably serve as an input to a modified **Alphafold** which might gives us structures of transient species or ratios of subpopulations

(and/or HDX-MS etc)

♡ 2          ⟲          ♡ 3          ⇧

Dina Grohmann @DinaGrohmann · 23 Jul
I'm amazed by the accuracy of the predicted structure in the Mid/PIWI lobe. Apart from that, I think that #AlphaFold could help us **smFRET** folks to find suitable positions for dye engineering.

♡ 1          ⟲          ♡ 3          ⇧

Show this thread

DeepMind

EMBL-EBI

# ColabFold

Making Protein folding accessible to all via Google Colab
**(and the unintended uses of AlphaFold)**



**github.com/sokrypton/ColabFold**

# ColabFold - Advanced options

- Modify MSA input
  - Custom or MMseqs2 (much faster)
  - Trim
- **Complexes**
  - **Homo-oligomers**
  - **Hetero-oligomers**
- Fine control
  - Number of recycles
- Sample (Output more than 5 models)
  - Generate ensembles by iterating through random seeds, enabling dropout

# Can predict protein-protein/peptide interactions



Yoshitaka Moriwaki @Ag_smith · Jul 19
AlphaFold2 can also predict heterocomplexes. All you have to do is input the two sequences you want to predict and connect them with a long linker.

**G-linker!**

**Don't actually need a G-linker!**

Minkyung Baek @minkbaek
Adding a big enough number for "residue_index" feature is enough to model hetero-complex using AlphaFold (green&cyan: crystal structure / magenta: predicted model w/ residue_index modification).
#AlphaFold #alphafold2

```
# add big enough number to residue index to indicate chain breaks
idx_res = feature_dict['residue_index']
L_prev = 0
# Ls: number of residues in each chain
for L_i in Ls[:-1]:
    idx_res[L_prev+L_i:] += 200
    L_prev += L_i
feature_dict['residue_index'] = idx_res
```

Hiroki Onoda @onoda_hiroki

Unknown linker may be useful for multimer prediction on the local Alphafold2!!

**UNK-linker!**

大上雅史 | Ohue M 2.5G @tonets
あ、AlphaFold2でペプチドドッキングでき…
Translated from Japanese by Google
Oh, I was able to dock the peptide with AlphaFold2

**Protein-peptide interaction**

# Can predict protein-protein/peptide interactions



**dimer-swap**

**intertwined dimer**

**consistent w/ biochem data**

**Cross-species**

**Preprints rolling in...**

Can AlphaFold2 predict protein-peptide complex structures accurately?

Junsu Ko, Juyong Lee

bioRxiv 2021.07.27.453972; doi: https://doi.org/10.1101/2021.07.27.453972

Harnessing protein folding neural networks for peptide-protein docking

Tomer Tsaban, Julia Varga, Orly Avraham, Ziv Ben-Aharon, Alisa Khramushin, Ora Schueler-Furman

bioRxiv 2021.08.01.454656; doi: https://doi.org/10.1101/2021.08.01.454656

Improved Docking of Protein Models by a Combination of AlphaFold2 and ClusPro

Usman Ghani, Israel Desta, Akhil Jindal, Omeir Khan, George Jones, Sergey Kotelnikov, Dzmitry Padhorny, Sandor Vajda, Dima Kozakov

bioRxiv 2021.09.07.459290; doi: https://doi.org/10.1101/2021.09.07.459290

# ColabFold se diversifie

**Making Protein folding accessible to all via Google Colab!**

| Notebooks | monomers | complexes | mmseqs2 | jackhmmer | templates |
|---|---|---|---|---|---|
| AlphaFold2_mmseqs2 | Yes | Yes | Yes | No | Yes |
| AlphaFold2_batch | Yes | Yes | Yes | No | Yes |
| RoseTTAFold | Yes | No | Yes | No | No |
| AlphaFold2 (from Deepmind) | Yes | Yes | No | Yes | No |
| **BETA (in development) notebooks** | | | | | |
| OmegaFold | Yes | No | No | No | No |
| AlphaFold2_advanced | Yes | Yes | Yes | Yes | No |
| **OLD retired notebooks** | | | | | |
| AlphaFold2_complexes | No | Yes | No | No | No |
| AlphaFold2_jackhmmer | Yes | No | Yes | Yes | No |
| AlphaFold2_noTemplates_noMD | | | | | |
| AlphaFold2_noTemplates_yesMD | | | | | |

# ColabDesign

`https://github.com/sokrypton/ColabDesign`

## ColabDesign

**Making Protein Design accessible to all via Google Colab!**

```
pip install git+https://github.com/sokrypton/ColabDesign.git
```

- TrDesign - using TrRosetta for design (support for TrMRF coming soon)
- AfDesign - using AlphaFold for design

(WIP) Not yet fully integrated into ColabDesign

- MSA_transformer
- Potts models (GREMLIN, mfDCA, arDCA, plmDCA, bmDCA, etc)
- ProteinMPNN
- RfDesign - using RoseTTAFold for design

## Presentations

Slides Talk

## Contributors:

- Sergey Ovchinnikov @sokrypton
- Justas Dauparas @dauparas
- Weikun.Wu @guyujun Levinthal.bio
- Shihao Feng

# AlphaFold-Multimer de DeepMind

Richard Evans et al. (Oct. 2021). en. preprint. Bioinformatics

# Some AlphaFold use cases

Alex Bateman

EMBL-EBI

# Just because AlphaFold can fold it doesn't mean nature can

- Pfam use case 2: CPB_BcsS family (PF17036)

- AlphaFold prediction of region matched by Pfam identifies incomplete domain

- Structurally similar to MBB clan

**There are 129 sequences with the following architecture: CBP_BcsS**

A0A3S2XL24_9RHIZ [Methylobacterium sp. TER-1] Cellulose biosynthesis protein BcsS {ECO:0000313|EMBL:RVU17441.1} (241 residues)

**CBP_BcsS**

Show all sequences with this architecture.

- This will not be stable in vitro!

Natural linker

+GG linker

+GGGG linker

Application of AF2 structures for variant effect prediction and pocket detection

Pedro Beltrao, Group Leader

www.ebi.ac.uk/beltrao
@pedrobeltrao

EMBL-EBI

# Comparing experimental vs structure based prediction of missense mutations



- Experimental impact of mutations from deep mutational scanning experiments (30 proteins)
- Predicted ddG of mutation using FoldX on alphafold structures or experimental structures
- Alphafold structures give equal or better predictions and it holds for regions with no templates

# Comparing experimental vs structure based prediction of missense mutations

# Pocket detection and how to filter the models



We retained 230 of 304 proteins from a dataset by Clark et al., 2020. Pocket detection was performed using ghecom (Kawabata, 2010), as done previously in (Clark et al., 2020).

EMBL-EBI

# Pocket detection and how to filter the models



Overlap of top predicted pocket to UBS (F-score)

Filtering the pocket residues by confidence (likely also predicted aligned residue) improves pocket detection

**AlphaFold2**
for detecting intrinsically
disordered protein regions

Bálint Mészáros
EMBL Heidelberg
08/09/2021

EMBL

# AlphaFold2 indicates the presence of IDRs

AF2 generates coordinates for every residue, even ones that have no fixed structure



human p53

Model Confidence:
- ■ Very high (pLDDT > 90)
- ■ Confident (90 > pLDDT > 70)
- ■ Low (70 > pLDDT > 50)
- ■ Very low (pLDDT < 50)

**Two interpretations for low confidence:**
- AF2 isn't good enough to predict the structure
- There is no structure to predict



pLDDT scores along the human p53 sequence

Order
Disorder

pLDDT is a good indicator of disordered regions (in this case)
Let's test the generic case – binary disordered prediction

EMBL

# AlphaFold2 as a disorder prediction method



IUPred2 vs AlphaFold pLDDT

IUPred2 (AUC=0.870)
AF2 pLDDT₁₅ (AUC=0.902)
random

false positive

You're pregnant

false negative

You're not pregnant

pLDDT score distribution on the human proteome

EMBL

# L'avenir selon DeepMind

## Une petite vidéo pour terminer...

```
https://www.youtube.com/watch?time_continue=8&v=
KpedmJdrTpY&feature=emb_title
```
à regarder chez vous (regarder sur les écrans, vous reconnaissez le
logiciel de visualisation que les chercheurs de DeepMind utilisent ?):
```
https://www.youtube.com/watch?v=gg7WjuFs8F4
```

# The end

- MERCI pour votre attention!

Akasaka, Kazuyuki, Ryo Kitahara, and Yuji O. Kamatari (Mar. 2013). "Exploring the folding energy landscape with pressure". In: *Archives of Biochemistry and Biophysics*. Protein Folding and Stability 531.1–2, pp. 110–115.

Bowman, Gregory R, Vincent A Voelz, and Vijay S Pande (Feb. 2011). "Taming the complexity of protein folding". en. In: *Current Opinion in Structural Biology* 21.1, pp. 4–11.

Chen, Yiwen et al. (Jan. 2008). "Protein folding: Then and now". In: *Archives of Biochemistry and Biophysics*. Highlight Issue: Protein Folding 469.1, pp. 4–19.

Daggett, Valerie and Alan R. Fersht (2003). "Is there a unifying mechanism for protein folding?" In: *Trends in biochemical sciences* 28.1, pp. 18–25.

Dill, Ken A. et al. (June 2008). "The Protein Folding Problem". en. In: *Annual Review of Biophysics* 37.1, pp. 289–316.

Evans, Richard et al. (Oct. 2021). *Protein complex prediction with AlphaFold-Multimer*. en. preprint. Bioinformatics.

Fersht, Alan R. and Valerie Daggett (2002). "Protein folding and unfolding at atomic resolution". In: *Cell* 108.4, pp. 573–582.

Hosseinzadeh, Parisa et al. (Dec. 2017). "Comprehensive computational design of ordered peptide macrocycles". en. In: *Science* 358.6369, pp. 1461–1466.

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

Kuhlman, Brian et al. (Nov. 2003). "Design of a Novel Globular Protein Fold with Atomic-Level Accuracy". en. In: *Science* 302.5649, pp. 1364–1368.

Lane, Thomas J et al. (Feb. 2013). "To milliseconds and beyond: challenges in the simulation of protein folding". en. In: *Current Opinion in Structural Biology* 23.1, pp. 58–65.

Muñoz, Victor and Michele Cerminara (Sept. 2016). "When fast is better: protein folding fundamentals and mechanisms from ultrafast approaches". In: *Biochem J* 473.17, pp. 2545–2559.

Piana, Stefano, John L Klepeis, and David E Shaw (Feb. 2014). "Assessing the accuracy of physical models used in protein-folding simulations: quantitative evidence from long molecular dynamics simulations". en. In: *Current Opinion in Structural Biology* 24, pp. 98–105.

Plaxco, Kevin W, Kim T Simons, and David Baker (Apr. 1998). "Contact order, transition state placement and the refolding rates of single domain proteins11Edited by P. E. Wright". In: *Journal of Molecular Biology* 277.4, pp. 985–994.

Thirumalai, D. et al. (2010). "Theoretical Perspectives on Protein Folding". In: *Annual Review of Biophysics* 39.1, pp. 159–183.

Voelz, Vincent A. et al. (Aug. 2012). "Slow Unfolded-State Structuring in Acyl-CoA Binding Protein Folding Revealed by Simulation and Experiment". In: *J. Am. Chem. Soc.* 134.30, pp. 12565–12577.

Woodside, Michael T. and Steven M. Block (2014). "Reconstructing Folding Energy Landscapes by Single-Molecule Force Spectroscopy". In: *Annual Review of Biophysics* 43.1, pp. 19–39.