# Deep learning for protein structure prediction and design

Tanja Kortemme

- Protein structure prediction – intro and significance
- Alphafold2 / concepts
- Applications: problem solved !?
- Design of new proteins

# AI & AlphaFold2 Revolution

- Breakthrough of 2021!

# Structural coverage of the proteome



https://www.deepmind.com/research/highlighted-research/alphafold

# Structural coverage of the proteome

**MetaAI ESMFold**

ESM Metagenomic Atlas:
The first view of the
'dark matter' of the
protein universe

blue: dark matter - no
similarity to previous
structures)



ESM Metagenomic Atlas

600,000,000 (!) Models available @https://esmatlas.com/

https://ai.facebook.com/blog/protein-folding-esmfold-metagenomics/

# Why predict protein structures?

(and what accuracy is needed?)

Protein structure prediction and structural genomics.
Baker D, Sali A.  Science. 2001 Oct 5;294(5540):93-6.

Which modeling accuracy is useful depends on the application

- Drug & protein design
- Docking

Design mutations for experimental tests

Hypotheses for function, effects of genetic variation

Protein structure prediction and structural genomics. Baker D, Sali A. Science. 2001 Oct 5;294(5540):93-6.

# CASP

-Blind structure prediction experiment

allows assessment of different approaches

• *every 2 years; summer 2022: CASP15*

Identification of major "winner strategies":
• CASP4: **fragments** (Rosetta)
• CASP11&12: **coevolution** and contact
prediction methods (contact-assisted
modeling)



**Starting CASP13: Deep learning** (Google alphafold)
CASP14 (2020): Google alphafold deepmind (AF2) "solved the problem"
CASP15 (2022): AF2-**based** methods lead; new, faster approaches using natural language
processing models (e.g., ESMFold) accelerate predictions

# Alphafold2: a game changer (CASP14 – 2020)

# What made the Alphafold2 breakthrough possible?

- **Basic research** - insights from > 70 years of protein research
- **Big data** - solved structures, large-scale sequencing *etc.*
- **Deep learning** - new architectures, optimization methods

# Sequence→ structure

**Anfinsen's dogma*:**

Native structure determined **only** by sequence

→ Native structure = global energy minimum

- unique
- stable
- kinetically accessible

Energy

GMEC

N

Conformation space

* true at least for a small globular protein, in its standard physiological environment

# Structure prediction

Sequence ➡ ➡

Basic Assumption:

Native structure = GMEC

Global Minimum Energy Conformation

➡A good energy function selects GMEC

➡A **good sampling technique** finds GMEC

Energy

N GMEC

Conformation space

# Why structure prediction is hard: Conformational space in "ab initio" structure prediction is enormous



- If only 3 states per residue, 100 residue protein: $3^{100} \sim 5 \times 10^{47}$

- Just considering 3 states isn't going to be detailed enough

- Clearly need methods to restrict degrees of freedom

# Breakthrough: contact maps

Sergey Ovchinnikov

**Rosetta GREMLIN (Generative REgularized ModeLs of proteINs)**

Long-standing idea: derive residue-residue contacts from sequence information



dMSA → co-evolution → Contact Map

**Learning: Apply techniques for object recognition on pictures… cats, street lights, faces, …**

# Image recognition using Deep NNs

Good at image recognition tasks:

Apply filters to image that highlight specific features

(for example: convoluted neural networks, CNN)



Input image

Edges highlighted

Residue covariance matrix

Contact probabilities

# Neural Networks

**Single** Neuron  -  **linear** separation

**Problem: not (linearly) separable**
**Solution: multiple** neurons, **multiple layers**



**a**  $x_{i...N}$ – inputs
$W_{i...N}$ – weights

$x_1$
$w_1$

$x_2$  $w_2$

$x_3$
$w_3$

$x_N$  $w_N$

$\sum_{i=1}^{N} w_i x_i$

Output T/F, 0,1

$g$ – transfer function
$t$ – threshold



Output

Hidden layer

$X_1$  $X_2$  $X_3$  $X_4$  $X_5$  $X_6$  $X_7$

# **Deep** Neural Networks

**Universal approximation theorem:**

A feed-forward neural network with a finite number of nodes can **approximate any continuous function**

**Deep** NN: many layers



$$output = \sigma\left(\sum_i w_i x_i + b\right)$$

# Breakthrough: ML

SQETRKKCTEMKKKFKNCEVRCDESNHCVEVRCSDTKYTL

**Protein Sequence**

**Neural Network**

**Databases**

*Collect sequence data*

**Distance Predictions**

**Angle Predictions**

*Contact / Distance prediction*

*Angle prediction*

**Score (Gradient Descent)**

*… build structure from restraints (as in Modeller)*

**Structure**

# CASP performance



**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

AlphaFold 2

AlphaFold

©nature

# Alphafold**2**: **End-to-End** architecture

## Unprecedented accuracy using novel representations

Trained on publicly available data consisting of ~170,000 protein structures (PDB) & large protein sequence databases.

Uses ~ 16 TPUv3s (= 128 TPUv3 cores ~100-200 GPUs) run over a few weeks

Provides confidence value & iterative improvement



New architectures
New approaches
New hardware

# Alphafold2 architecture in a nutshell

# Alphafold2 architecture in a nutshell



New architecture
New approaches
New hardware

Evoformer blocks

Structure modules

MSA & coevolution

Sequence

Structure & confidence

# Alphafold2 architecture in a nutshell



- **pLDDT:** residue confidence (predicted local distance difference test)
- **pAE:** residue *pair* confidence (predicted alignment error)

New architecture
New approaches
New hardware

Evoformer blocks

Structure modules

Sequence

Structure & confidence

# Alphafold2: a game changer (CASP14 – 2020)



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

**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.

AlphaFold 2

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

AlphaFold

Global distance test (GDT_TS; average)

Contest year

©nature

# The good, the bad and the ugly

The network also **models the uncertainty in its predictions** - when the s.d. of the predicted distribution is low, the predictions are more accurate:

Confidence measure for each residue:

**pLDDT**

(predicted local distance difference test)
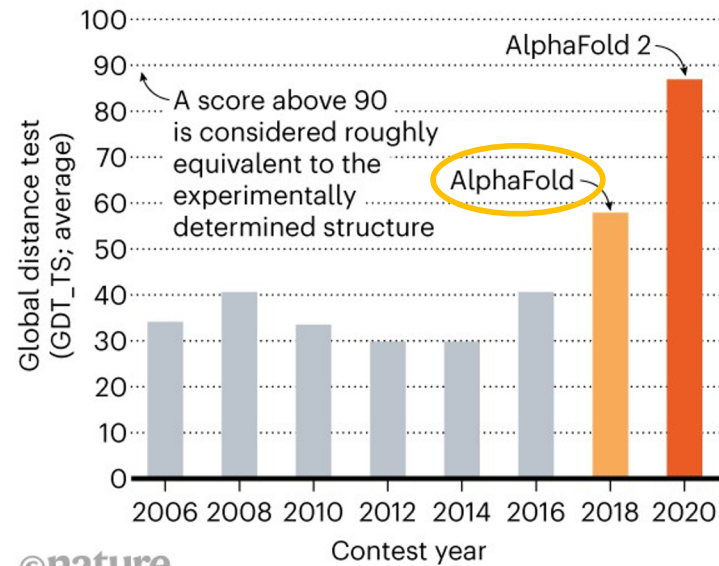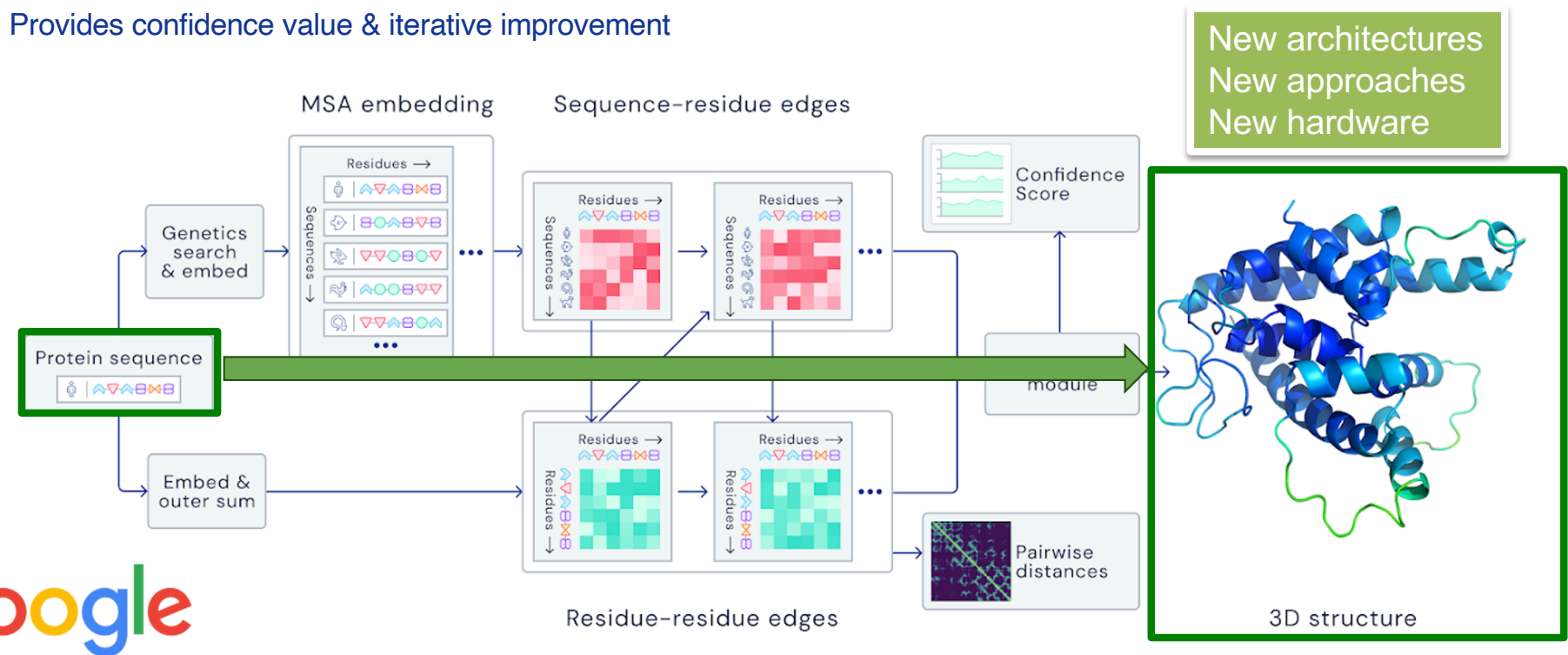


Protein Data Bank (PDB) structure

AlphaFold structure, with confidence estimates for each section.

Very high · High · Low · Very low

**Good**
AlphaFold model of phosphohistidine phosphatase overlaps closely with PDB structure.

**Bad**
AlphaFold model of human insulin bears no relation to the PDB structure.

**Ugly**
AlphaFold has little confidence across much of its prediction for this human ubiquitin-protein ligase. There is no PDB structure to compare it with.

# Accuracy on recent PDB structures

# Interpretation of Alphafold2 models

- pLDDT: predicted **local** distance from solved structure [0..1] > 0.7 precise



Identifying domains & possible disordered regions

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

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

Assessing confidence within a domain

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

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

*https://alphafold.ebi.ac.uk/*

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

# What can be done now? (and what is difficult)

- Combine AF2 predictions with experimental data to create models of complex proteins and assemblies
- Predict structures of complexes (limitation: MSA!)
- In some cases: use predictions for ligand docking

- Disorder?  Some indication from pLDDT
- "Orphan" sequences and de novo proteins – accuracy?
- Prediction of effect of mutations?  Difficult!

Which modeling accuracy is useful depends on the application

- Drug & protein design
- Docking

Design mutations for experimental tests

Hypotheses for function, effects of genetic variation

Protein structure prediction and structural genomics.
Baker D, Sali A. Science. 2001 Oct 5;294(5540):93-6.

# Why are AI models often insensitive to mutations?



AlphaFold   Experiment
r.m.s.d. = 0.59 Å within 8 Å of Zn

- In the example, the metal binding site is predicted accurately even though the metal was not included!

- Methods trained on metal-bound structures recognize the pattern of a metal binding site (even if a structure unfolds in the absence of the metal)

# Summary : Structure prediction

Enormous recent progress, enabled by:

**large databases of sequences and structures,** AI methods from other fields, new deep learning network architectures, hardware, computing power

- Informative and large sequence alignment is (typically) critical, but many sequences are available today (metagenomic data)
- ML & END-TO-END models (Alphafold2, ESMfold and more to come !)
- Language models to learn the Protein language (fast, perhaps more general?)

Accessible to all:

- Models available in Uniprot, EBI, MetaAI
- Modeling made easy on COLAB

Challenges: multiprotein assemblies, disordered proteins, mutations

# Outlook


An animation of the gradient descent method predicting a structure for CASP13 target T1008

**New applications**

Fast and accurate

- structures for research & medicine
- drug design



**Extend to protein design**

- inverse direction:



- protein "hallucinations": Dream new proteins with the NN and much much more

# Diffusion models for protein design



Fixed forward diffusion process

Data — Noise

Generative reverse denoising process

**Diffusion Model**

Forward (Noising) Process

$\mathcal{N}(0,1)$

**Gaussian Noise** ... single step ... **Protein Structure**

$X_T$  $X_t$  $X_{t-1}$  $X_0$

Reverse (Generative) Process

# Diffusion model for protein design

Generate a protein that binds to a helix:



https://www.ipd.uw.edu/2022/12/a-diffusion-model-for-protein-design/

# Diffusion model for protein design

Make assemblies



https://www.ipd.uw.edu/2022/12/a-diffusion-model-for-protein-design/

**A**

**Diffusion Model**

Forward (Noising) Process

$\mathcal{N}(0,1)$

single step

**Gaussian Noise**

**Protein Structure**

$X_T$  $X_t$  $X_{t-1}$  $X_0$

Reverse (Generative) Process

**RoseTTAFold**

Input Sequence  MADHTI?DTREE

Homologous Templates

Initial/Recycled Coordinates

**RF**

Predicted Structure

**RFdiffusion**

Masked Input Sequence  ?????????????

$\hat{X}_0^{t+1}$ (Self-conditioning)

$X_t$ **Diffused** Coordinates

**RF**

$\hat{X}_0$

**Single RFdiffusion step**

$X_t \rightarrow$ **RF** $\rightarrow \hat{X}_0 \rightarrow interp(X_t, \hat{X}_0) + \varepsilon \rightarrow X_{t-1}$

**Self-conditioning**

**B**

**Figure 1: RFdiffusion is a** denoising diffusion probabilistic model with RoseTTAFold **fined-tun**ed **as the denoising network. A)** Top panel: Diffusion models for proteins are trained to recover structur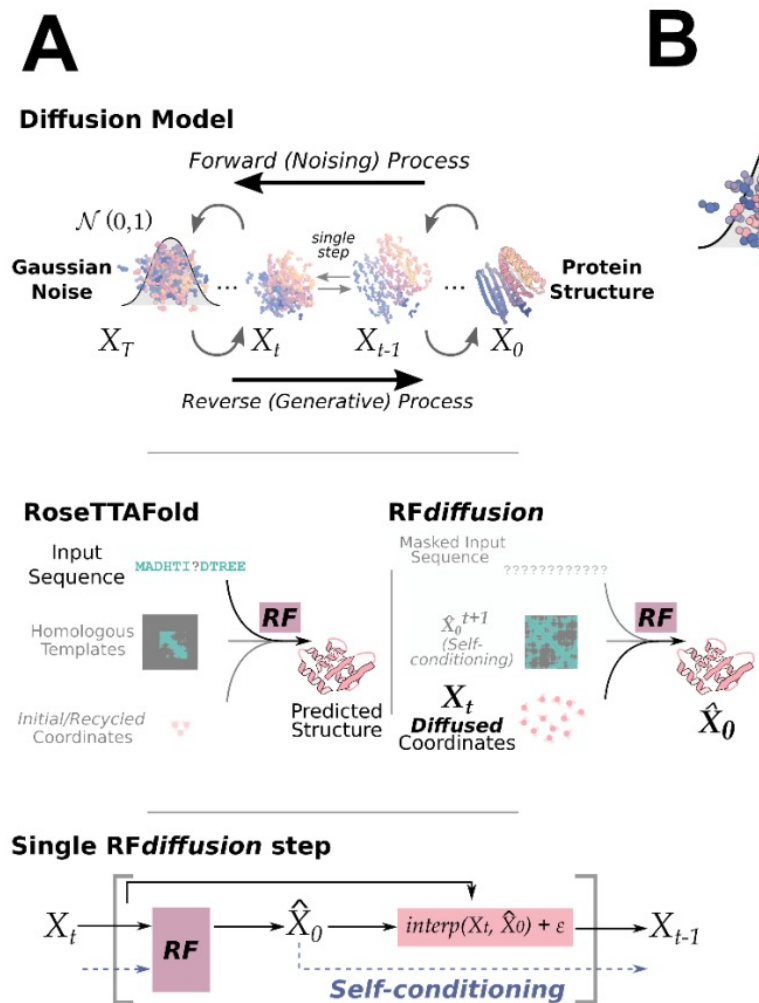es of proteins corrupted with noise, and generate new structures by reversing the corruption process through iterative denoising of initially random noise $X_T$ into a realistic structure $X_0$. Middle panel: RoseTTAFold (RF, left) can be fine-tuned as the denoising network in a DDPM. RFdiffusion (right) is trained from a *pre-trained* RF network with minimal architectural changes. While in RF, the primary input to the model is sequence, in RFdiffusion, the primary input is diffused residue frames. In both cases, the model predicts final 3D coordinates directly (denoted $\hat{X}_0$ in RFdiffusion). In RFdiffusion, the model receives its previous prediction as a template input ("self-conditioning", see Methods 2.4). Bottom panel: At each timestep "t" of a design trajectory (typically 200 steps), RFdiffusion takes $X_t$ and $\hat{X}_0^{t+1}$ from the previous step and then predicts an updated $X_0$ structure ($\hat{X}_0^t$). The coordinate input to the model at the next time step ($X_{t-1}$) is generated by a noisy interpolation toward $\hat{X}_0^t$. B) RFdiffusion is of broad