

AlphaFold: a new age for protein folding

Jasper Zuallaert





28/01/2022

PhD in Computer Science (promotors Wesley De Neve / Yvan Saeys / Nico Callewaert)

Jasper Zuallaert



(deep learning for sequence analysis)

GLOBAL CAMPUS

Post-doc at Nico Callewaert + Lennart Martens units at Flemish Institute for Biotechnology (VIB)







/UGent

Introduction

AlphaFold: protein folding software by Google DeepMind

One of the most frequent software requests ever (EasyBuild)

Table of contents:

- What are proteins?
- Why is protein folding important?
- What is AlphaFold?
- Using AlphaFold...



UNIVERSITEIT GENT



From DNA to protein



UNIVERSITEIT GENT





Central dogma of molecular biology



Central dogma of molecular biology



Protein structure prediction









Available sequences and structures



Main protein structure resource



~50 000 unique structures

Main experimental technique: X-ray crystallography

- Costly -
- Labour-intensive
- Slow



Image source: https://www.prweb.com/releases/2017/04/prweb14219956.htm





GENT

Why 3-D structure prediction? Low effort, time efficient

Examples:

- In silico drug candidate screening



- De novo protein design guidance







Structure prediction: homology models

Basic principle: 3-D structure more conserved than sequence

Query protein sequence: ELAIGILTVSYIPSAEKIRAPELTI

Sequence alignment:



Finetuning using statistical potentials and physics-based energy calculations





Critical Assessment of Structure Prediction (CASP14, 2020)



Median Free-Modelling Accuracy

GENT

UNIVERSITEIT

VIB-UGENT CENTER FOR MEDICAL

BIOTECHNOLOGY



 $Source: https://news.machinelearning.sg/posts/alphafold {\tt 2_10_things_you_want_to_know_about_biologys_imagenet_moment/sources/approx$





AlphaFold





DeepMind: a timeline



Critical Assessment of Structure Prediction (CASP14, 2020)

Median of 92.4 across all targets, 87.0 in free-modelling accuracy

Median Free-Modelling Accuracy



GENT

UNIVERSITEIT

VIB-UGENT CENTER FOR MEDICAL

BIOTECHNOLOGY

Reception of Alphafold 2

Mohammed AlQuraishi @MoAlQuraishi

CASP14 #s just came out and they're astounding —DeepMind looks to have solved protein structure prediction. Median GDT_TS went from 68.5 (CASP13) to 92.4!!!! Cf. their 2nd best CASP13 struct scored 92.8 (out of 100). Median RMSD is 2.1Å. I think it's over predictioncenter.org/casp14/zscores...

1:13 PM \cdot Nov 30, 2020 \cdot Twitter for iPhone

IB-UGENT CENTER

FOR MEDICAL

BIOTECHNOLOGY

619 Retweets 293 Quote Tweets 2,125 Likes

Dr. Mohammed AlQuraishi at Columbia University, who also participated in CASP, lauded the Al as transformational. "It's a breakthrough of the first order, certainly one of the most significant scientific results of my lifetime," he said to *Nature*.

A "gargantuan" leap today for #AI life science. @DeepMind @GoogleAI #AIphaFold2 prediction of #3D protein structure from amino acids nature.com/articles/d4158... by @ewencallaway @NatureNews twitter.com/demishassabis/... @demishassabis

17/33

'Once in a generation advance' as Google AIThe Telegraphresearchers crack 50-year-old biological
challenge

The development could 'significantly accelerate' drug development for cancer and other diseases

"This is a big deal," says John Moult, a computational biologist at the University of Maryland in College Park, who co-founded CASP in 1994 to improve computational methods for accurately predicting protein structures "In some sense the problem is solved."

All of the groups in this year's competition improved, Moult says. But with AlphaFold, Lupas says, "The game has changed. The organizers even worried DeepMind may have been cheating somehow. So Lupas set a special challenge: a membrane protein from a species of archaea, an ancient group of microbes. For 10 years, his research team tried every trick in the book to get an x-ray crystal structure of the protein. "We couldn't solve it."

But AlphaFold had no trouble. It returned a detailed image of a three-part protein with two long helical arms in the middle. The model enabled Lupas and his colleagues to make sense of their x-ray data; within half an hour, they had fit their experimental results to AlphaFold's predicted structure "It's almost perfect," Lupas says. "They could not possibly have cheated on this. I don't know how they do it."

SCIENCE MEETS LIFE

Reception of Alphafold 2

Garry Kasparov 🤄 😏	•••	
Congrats to @demishassabis and @DeepMind! Glad you're leaving the chessplayers alone for a little while!		
The New York Times @ nytimes The artificial intelligence lab DeepMind built a computer system that can	e 3) to	
identify the shape of a protein in mere hours. This long-sought breakthrough could accelerate the ability to understand diseases and help develop new medicines. nyti.ms/3fOIDyN	8	2
11:44 PM · Nov 30, 2020 (j)		D
♡ 1.2K ♀ 11 ♂ Copy link to Tweet		A

Dr. Mohammed AlQuraishi at Columbia University, who also participat CASP, lauded the Al as transformational. "It's a breakthrough of the fir certainly one of the most significant scientific results of my lifetime," h *Nature*.

A "gargantuan" leap today for #AI life s @DeepMind @GoogleAI #AIphaFoId2 protein structure from amino acids nature.com/articles/d4158... by @ewer @NatureNews twitter.com/demishassabis/... @demish Che Telegraph'Once in a generation advance' as Google AIThe Telegraphresearchers crack 50-year-old biological
challenge

Science

2021 BREAKTHROUGH OF THE YEAR

Protein structures for all

AI-powered predictions show proteins finding their shapes BY ROBERT SERVICE

In his 1972 Nobel Prize acceptance speech, American biochemist Christian Anfinsen laid out a vision: One day it would be possible, he said, to predict the 3D structure of any protein merely from its sequence of amino acid building blocks. With hundreds of thousands of proteins in the human body alone, such an advance would have vast applications, offering insights into basic biology and revealing promising new drug targets. Now, after nearly 50 years, researchers have shown that artificial intelligence (AI)-driven software can churn out accurate protein structures by the thousands—an advance that realizes Anfinsen's dream and is *Science*'s 2021 Breakthrough of the Year.

Multiple Sequence Alignment

AlphaFold input = query protein sequence

First step: calculate multiple sequence alignment

Source: https://www.blopig.com/blog/2021/07/alphafold-2-is-here-whats-behind-the-structure-prediction-miracle/

AlphaFold stages

AlphaFold training

Any neural network is trained using a loss/function

21/33

+ auxiliary losses:

- \rightarrow FAPE + torsion loss at intermediate structures
- \rightarrow distogram prediction loss
- \rightarrow Masked MSA prediction loss
- \rightarrow Structural violation loss
- → Predict local distance difference test (pLDDT)
- → Predicted Alignment Error (PAE)

GENT

UNIVERSITEIT

AlphaFold availability

Google COCD

AlphaFold Protein Structure Database

Developed by DeepMind and EMBL-EBI

AlphaFold on the HPC

IIII UNIVERSITEIT GENT

AlphaFold resource usage

- 1. Sequence database search
 - \rightarrow RAM-intensive
 - \rightarrow High disk reading speed desirable
 - → Datasets ~2.2 TB storage
- 2. Deep learning prediction models
 - \rightarrow GPU memory-intensive

Memory limits and time elapsed

Memory limits and time elapsed

5 4 time elapsed (hours) 3 ---MSA predict (5x) 2 --- relax (5x) 1 0 0 500 1000 1500 2000 2500 sequence length

joltik (V100) cluster, time elapsed for default AlphaFold implementation

UNIVERSITEIT GENT

AlphaFold on HPC

4. profit!

SCIENCE MEETS LIFE

Example comparison of time required

Target: SARS-CoV-2 spike RBD

>rbd ITNLCPFGEVFNATRFASVYAWNRKRISNCVADYSVLYNSASFSTFKCYGVSPTKLNDLC FTNVYADSFVIRGDEVRQIAPGQTGKIADYNYKLPDDFTGCVIAWNSNNLDSKVGGNYNY LYRLFRKSNLKPFERDISTEIYQAGSTPCNGVEGFNCYFPLQSYGFQPTNGVGYQPYRVV VLSFELLHAPATVCGPKK

(198 residues)

		Installation	MSA	Predict (5x)	Relax (1x/5x)	Total time elapsed
	Colab: Official AlphaFold (v2.1)	9 min	38 min	48 min (relax 1x)		1 hr 35 min
\int	HPC: Official AlphaFold (v2.0)	-	15 min	7 min	<1 min (5x)	23 min
	HPC: <i>adapted script</i> (ColabFold modified version)	-	14 min	4 min	<1 min (1x)	19 min

accelgor (A100) ^

What's next?

Protein complex prediction

1. Single chain with pseudolinker

UNIVERSITEIT GENT

2. Separate chains programmatically

3. AlphaFold-Multimer

DeepMind > Research > Protein complex prediction with AlphaFold-Multimer

31/33

Protein complex prediction with AlphaFold-Multimer

Abstract

While the vast majority of well-structured single protein chains can now be predicted to high accuracy due the recent AlphaFold [1] model, the prediction of multi-chain protein complexes remains a challenge in man cases. In this work, we demonstrate that an AlphaFold model trained specifically for multimeric inputs of kn stoichiometry, which we call AlphaFold-Multimer, significantly increases accuracy of predicted multimeric interfaces over input-adapted single-chain AlphaFold while maintaining high intra-chain accuracy. On a benchmark dataset of 17 heterodimer proteins without templates (introduced in [2]) we achieve at least

OpenFold

Mohammed AlQuraishi @MoAlQuraishi

An announcement I've been aching to make! After much sweat, we've built a trainable version of AlphaFold2, implemented in PyTorch, which we're calling OpenFold.

GitHub: github.com/aqlaboratory/o...

Colab: colab.research.google.com/github/aqlabor...

Why a trainable version of AlphaFold2 you ask? 🛃

4:57 PM \cdot Nov 12, 2021 \cdot Twitter Web App

VIB-UGENT CENTER FOR MEDICAL

BIOTECHNOLOGY

510 Retweets 46 Quote Tweets 2,052 Likes

GENT

UNIVERSITEIT

• • •

Mohammed AlQuraishi @MoAlQuraishi · Nov 12, 2021 Replying to @MoAlQuraishi

As we saw with the recent AlphaFold-Multimer, some applications can benefit from training new AF2 variants and possibly integrating AF2 within larger models. DeepMind's JAX version, while excellent, is missing training code. PyTorch is also more widely used, hence OpenFold.

...

Single sequence prediction

bioRxiv, August 2021

Single-sequence protein structure prediction using language models from deep

learning

Ratul Chowdhury^{a,*}, Nazim Bouatta^{a,*}, Surojit Biswas^{b,c,*}, Charlotte Rochereau^d, George M. Church^{a,b},

Peter K. Sorger^{a,†}, and Mohammed AlQuraishi^{a,e,†}

bioRxiv, January 2022

Single-sequence protein structure prediction using supervised transformer protein language models

Wenkai Wang¹, Zhenling Peng², Jianyi Yang^{1,*}

¹School of Mathematical Sciences, Nankai University, Tianjin 300071, China.

²Research Center for Mathematics and Interdisciplinary Sciences, Shandong University, Qingdao, 266237, China.

Molecular Dynamics

Source: https://2019.igem.org/Team:UANL/Simulations

Thank you for your attention.

INIVERSITEIT GENT

jasper.zuallaert@ugent.be