# AlphaFold 2

Kovács Gábor research / engineering @ Turbine.Al

#### AlphaFold 'pushes science forward' by releasing structures of almost all human proteins

DeepMind's AI predicted over 365,000 protein structures, which are now freely available online by **Emily Harwitz** 

July 29, 2021 | A version of this story appeared in Volume 99, Issue 28

#### AlphaFold 2 is here: what's behind the structure prediction miracle

*Nature* has now released that <u>AlphaFold 2 paper</u>, after eight long months of waiting. The main text reports more or less what we have known for nearly a year, with some added tidbits, al-though it is accompanied by a painstaking description of the architecture in the <u>supplementary</u> <u>information</u>. Perhaps more importantly, the authors have released the entirety of the code, including all details to run the pipeline, <u>on Github</u>. And there is no small print this time: you can run inference on *any* protein (I've checked!).

Oct 3, 2021, 07:34pm EDT | 54,391 views

#### AlphaFold Is The Most Important Achievement In AI —Ever

BLOG POST RESEARCH



Rob Toews Contributor ①

I write about the big picture of artificial intelligence.

Follow

DeepMind

**AlphaFold: a solution** 

to a 50-year-old grand

challenge in biology

30 NOV 2020

DeepMind open-sources AlphaFold 2 for protein structure predictions



2

- An end-to-end neural network for the protein folding problem
- Proteins are everywhere (digestion, muscles, firing neurons, immune system....)
- Predict 3D structure from amino acid (AA) sequence

sequence: MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLSPLPSQAMDDLMLSPDDIEQWFTEDPGPDEAPRMPEAAPPVAPAPAAPTPAAPAPAPS







- Why is it important?
  - Shape -> implies function
  - Design new drugs
  - Outside medicine: design proteins that digest plastic...



Valine

- Why is it hard?
  - Measuring structure experimentally: slow & expensive (X-ray crystallography)
  - ~100K measured --> 100s of billions potential configurations
  - Levinthal's paradox
    - 10^300 conformations
    - Happens in msec scale in life
    - Nature doesn't brute force folding...

- Predictive model classes
  - Physical interactions-based methods
  - Evolutionary history-based methods
  - Since 2018: neural networks emerge on this field with dominating results
- CASP Critical Assessment of Techniques for Protein Structure Prediction
  - Biennial challenge for models to predict yet unpublished prot. structures
  - DeepMind's AlphaFold 2 model wins CASP 14 event... by a lot



## Why do they say: "problem solved" ?



## Why do they say: "problem solved" ?





## Why do they say: "problem solved" ? - Global distance test (GDT)



44% Match



#### Why do they say: "problem solved" ? - Global distance test (GDT)



#### Why do they say: "problem solved" ? - Global distance test (GDT)



## **Biology background – a layman's view**

- Inputs are proteins, defined by their amino acid (AA) sequences.
- But what is an amino acid?
- Cα alpha carbon
- NH2 amino group
- COOH carboxyl group
- R residue / R group



- R groups are very important. They determine each AA's identity
- Variations of R groups determine characteristics of molecules that form when AAs are joined together



## **Biology background – a layman's view**



#### Amino acids: "alphabet" --- Proteins: "very long words"



turbine

#### From AA sequence --> to 3D structure

- Backbone is essentially planar & rigid due to  $\pi$  bond (covalent subtype)
- $\sigma$  bond between C $\alpha$  and R group -> can freely rotate
- $\bullet$  Through these  $\sigma$  bonds' rotations local structures are formed
- R groups of 2 AAs can interact by hydrogen bond / electrostatically



# Secondary structure



alpha helix

turbine

## **Tertiary structure**



## AlphaFold 2 architecture



## AlphaFold 2 architecture



## AlphaFold 2 inputs – MSA – Multiple Sequence Alignment

Query sequence to fold: PAWKFIQLLYP....

Bioinfo databases

Result set. Each record is a similar matching pattern

Q5E940 BOVIN	MI	REDRATWK	SNY <mark>F</mark> LK <b>I</b> I	<mark>)LLDDYP</mark> K	CFIV <mark>G</mark> ADNV	<mark>G S</mark> K <mark>QMQ</mark> Q I RMS	LRGK-AVVLMC	KNTMMRKAIRGHI	ENNPALE
RLA0 HUMAN	MI	REDR <mark>A</mark> TWK	SNY <mark>F</mark> LK <b>I</b> I	<mark>)LLDDYP</mark> K	CFIV <mark>G</mark> ADNV	<mark>G S</mark> K <mark>QMQ</mark> Q I RMS	L <mark>RG</mark> K-AVVLM <mark>O</mark>	<mark>KNT</mark> MMRKAIRGHI	ENNPALE
RLA0 MOUSE	M	REDRATWK	SNY <mark>F</mark> LK <b>I</b> I	<mark>)LLDDYP</mark> K	CFIV <mark>G</mark> ADNV	G <mark>S</mark> K <mark>QMQ</mark> QIRMS	L <mark>RG</mark> K-AVVLM <mark>O</mark>	<mark>KNT</mark> MMRKAIRGHI	ENNPALE
RLA0_RAT	MI	REDR <mark>A</mark> TWK	SNY <mark>F</mark> LKII	<mark>)LL</mark> DD <mark>YP</mark> K	CFIV <mark>G</mark> ADNV	G S K <mark>Q M Q</mark> Q I R M S	LRGK-AVVLMC	<mark>KNT</mark> MMRKAIRGHI	ENN <mark>P</mark> ALE
RLA0 CHICK	MI	REDR <mark>A</mark> TWK	SNY <mark>F</mark> MK <b>I</b> I,	<mark>)LLDDYP</mark> K	CFVV <mark>G</mark> ADNV	<mark>gs</mark> k <mark>omo</mark> qirms	L <mark>RG</mark> K-AVVLM <mark>O</mark>	KNTMMRKAIRGHI	ENNPALE
RLAO RANSY	MI	REDR <mark>ATW</mark> K	SNY <mark>F</mark> LK <b>I</b> I	<mark>)LLDDYP</mark> K	CFIV <mark>G</mark> ADNV	G <mark>S</mark> K <mark>QMQ</mark> QIRMS	L <mark>RG</mark> K-AVVLM <mark>G</mark>	KNTMMRKAIRGHI	ENNSALE
Q7ZUG3_BRARE	MI	REDR <mark>A</mark> TWK	SNY <mark>F</mark> LKII	<mark>) L L</mark> D D <mark>Y P</mark> K	CFIV <mark>G</mark> ADNV	<mark>G S</mark> K <mark>QMQ</mark> T IRLS	L <mark>RG</mark> K-AVVLM <mark>O</mark>	<mark>KNT</mark> MMRKAIRGHI	ENN <mark>P</mark> ALE

Convert into a matrix to feed into a neural network



# Feature representation



#### Feature representation

#### Evoformer





#### Step 2 – Pairwise representations (distograms)



## Step 2 – Evoformer stack – inside 1 block

- Attention vs convolution
- 2 way of information flow for the 2 modalities
- Row-wise & column-wise multihead self-attention



## Step 2 – Evoformer stack – triangle update



(3

## **Step 3 – Structure module**

- Try to predict 3D positions of residues
- black hole initialization
- Starting point: all residues in the same spatial position (0,0,0)
- Iterative refinement
- Done by the structure module's 8 layers of recurrent attention blocks
- Refinement of predicted 3D structure step by step
  - Step 1: seq representation update
  - Step 2: invariant update on residue positions





## **Techniques in training**

- Recycling
- Self-distillation
- Intentionally wrong templates
- Multiple components of loss function
  - Locally good predicted atom coordinates (Frame Aligned Point Error)
  - cross entropy between true and predicted pairwise distance representations
  - Random masking applied to input MSA masks. 1 task during training: restore the masked items correctly. Something like in language models <-- local representations get stronger
  - tries to predict how confident 3D structure prediction will be



 $\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$ 





Cuccos

**Bold cuccos**