Clustering predicted structures at the scale of the known Protein Universe

CASP 2023, Worldwide

Martin Steinegger Steinegger Lab @ SNU



Mission: Accelerate biological science through fast and easy to use methods



Sung-eun Jang Graduate Student



Woosub Kim Graduate Student



We Want You

Preprint out today! Petascale Homology Search for Structure Prediction



https://www.biorxiv.org/content/10.1101/2023.07.10.548308v1

AlphaFold2 is revolutionizing protein bioinformatics

"Everything that relies on a protein sequence, we can now do with protein structure" Mohammed AlQuraishi, Columbia U.

EMBL EBI released structural models of over **214 million proteins (25TB)**

A.I. Predicts the Shape of Nearly Every 'THE ENTIRE PROTEIN UNIVERSE': **Protein Known to Science**

DeepMind has expanded its database of microscopic biological mechanisms, hoping to accelerate research into all living things.



AI PREDICTS SHAPE OF NEARLY EVERY KNOWN PROTEIN

DeepMind's AlphaFold tool has determined around 200 million protein structures, which are now available to scientists in a database.





Slide idea by Sameer Velankar



Foldseek

Ultra fast searching and clustering of protein structures









Stephanie S. Charlotte Kim* Tumescheit

MPI Göttingen SNU, Korea SNU, Korea SNU, Korea SNU, Korea MPI Göttingen MPI Göttingen * Contributed equally

Structure alignments reveal evolutionarily distant homologs

To speed up search, reduce structures to sequences and use fast sequence searches

	55	58	76	78	126	128	133 135
•••	A B G	в	JD	D ···	HE	D	FCD…

Each residue is represented by a **structural state** letter

Foldseek describes tertiary interactions (not backbone)

Foldseek – virtual centers are useful to find structural relevant nearest neighbors

Foldseek – encoding structures into 3Di sequences and training of 3Di alphabet

Variational autoencoder trained with structural alignments

Foldseek 3Di sequences are highly conserved

Foldseek 3Di sequences are highly conserved

Foldseek search sensitivity is comparable to TMalign

Benchmark on single domains (SCOP Superfamily)

search.foldseek.com

Foldseek summary

Encodes structure as sequence by discretizing tertiary interactions

Searches billions of structures in minutes, $10^4 - 10^5 \times \text{faster than Dali/TMalign}$...

- Accurate E-values
 Accurate E-values<
 - Local alignments, robust against domain movements
 - TP hits up to 1st FP >100 downloads/week, 1000 repo views/week, used in CATH, NCBI iCn3D,...

Foldseek will help to organize our structural space

10

15

20

github.com/steineggerlab/foldseek

Fast and accurate protein structure search with Foldseek, van Kempen & Kim et al. (2023), Nat. Biotechnol.

Can we cluster the AlphaFold TrEMBL database using structures?

Clustering can reduce redundancy and generate biological insights

Linclust

Linear-time clustering of huge sequence sets

(or now structures)

https://foldseek.com for protein structures

https://mmseqs.com for protein sequences

1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer (by numerically sorting seq-*k*mer pairs)

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence ____ per group (the longest)

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence ____ per group
- 3. It aligns each sequence in the group *only with the centre sequence*

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence e per group
- 3. It aligns each sequence in the group *only with the centre sequence,* **instead** of with all sequences in the group

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence e per group
- 3. It aligns each sequence in the group **only with the centre sequence**, instead of with all sequences in the group
- 4. Sequences similar enough to a centre sequence will form a cluster

After Linclust we perform an all vs all cascaded clustering

Existing sequence clustering algorithms scale almost quadratically in number of sequences

Linclust scale linearly in number of sequences

Analyze the protein universe using MMseqs2 and Foldseek cluster

MMseqs2 cluster 90% sequence overlap 50% sequence identity

214M proteins AFDB

Inigo Barrio Hernandez*

Jingi Yeo* Jürgen Jänes * Contributed equally

highest pLDDT

52.3M clusters

Representative

Foldseek cluster

90% structure overlap

E-value < 0.01

Tanita Wein Mihaly Varadi Sameer Velankar Pedro Beltrao

AN'S

18.8M cluster

Foldseek clusters

🕂 fragment

Clustering predicted structures at the scale of the known protein universe, Barrio&Yeo et al, biorxiv

Singletons

Remove

Fragments

• • •

2.30M cluster AFDB clusters

We looked at the data from different angles

Proteins that we don't know

37.7k C. Bathyarchaeota Archaea 914 archaeon 11k 12.6k Euryarchaeota 2.27M 1.38k C. bacteriu 885k Cellular 1.85k A. bacteriu Bacteria 40.2k Bacteroidetes Organisms 79.7k Firmicutes 11.1k Cyanobacteria 1.30k A. bacteriu 366k 87.2k Actinobacteria 1.10k P. bacteriur 10.7k Planctomycetes 1.25k D. bacteriu 208k 1.66k E. coli Proteobacteria 17k 769 H. pylori 817k Fungi 720 S. enterica Eukaryota 56.4k Ascomycota 1.06k S. irregulari 25.2k Basidiomycota 182k 1.40k S. pharaoni 29.4k Arthropoda 1.76k A. ventricos 307k Metazoa 16.9k Nematoda 28.4k Chordata 102k 1.34k T. cinerariif 532k 94 3k Streptophyta Viridiplantae 0 Domain Species Life Kingdom Phylum **Evolutionary analysis** New evolutionary insight on proteins

35% of clusters do not have have a member with annotation. Cluster w/o annotation tend to have less members.

Known proteins tend to cluster together The larger cluster the more likely it is to be annotated.

Predict the function of the dark clusters

39,180 proteins (7%) score > 0.8

Top predicted molecular functions

Examples of structures with predicted pockets and functional annotations

Predicted GO:0005215 - transporter activity

altering nucleic acid conformation

Explore the dark cluster at https://af-protein-universe.streamlit.app/

Putative novel enzymes and small molecule binding proteins

Browse examples Global statistics

All structures/pockets

Examples: Fig 2B <u>A0A849TG76</u> and <u>A0A2D8BRH7</u>; Fig 2C <u>A0A849ZK06</u>; Fig 2D <u>S0EUL8</u>.

Hide structures with a general lack of compactness (struct_resid_in_pockets > 0.4)

UniProtKB_ac	n_resid	mean_pLDDT	pocket_id	pocket_score
A0A7K0F2F4	194	96.34	1	65.5
A0A178P4I2	187	93.67	1	65.2
A0A2D8BRH7	225	96.58	1	69.1
A0A1G2SFE8	209	94.93	1	66.3
A0A5J4E7J2	151	95.59	1	77.6
A0A7K0F2F4	194	96.34	2	64.6

DeepFRI GO/EC terms for A0A7K0F2F4

Hide non-significant terms (Score < 0.5)</p>

Hide terms without saliency data

Protein	GO_term/EC	Score	GO_term/EC	DeepFRI_ont	
A0A7K0F2F4	GO:0016757	0.96	transferase acti	MF	
A0A7K0F2F4	GO:0140096	0.66	catalytic activity	MF	

Structure/visualisation for A0A7K0F2F4

Residues colored by saliency for GO:0016757

Lowest common ancestor (LCA) analysis highlights the taxonomical

distribution of the clusters

Rediscovering **known** and **novel** links between bacterial and human proteins

Many more cross-kingdom clusters need to be explored

Use Foldseek's local alignments to detect and connect domains

How to explore the data

cluster.foldseek.com

Check also out the protein universe by Durairaj et al. (2023) at <u>https://uniprot3d.org/atlas/AFDB90v4</u>

https://www.biorxiv.org/content/10.1101/2023.03.14.532539v3

Conclusion

Sequence analysis is at the basis of most protein bioinformatics

We can now super charge the analysis using structures

Acknowledgments

Stephanie S. Charlotte Michel van Kempen Kim Tumescheit

Jürgen Jänes

Jingi

Yeo

Inigo Barrio Hernandez

Cameron

Jeongjae Mihaly Varadi

Johannes Söding

Lee

Milot Mirdita

Sameer Tanita Wein

Pedro

Beltrao

GEFÖRDERT VOM

Bundesministerium für Bildung und Forschung

+ Foldseek logo by Doyoon Kim

Use Foldseek's local alignments to detect and connect domains

Conclusion

Sequence analysis is at the basis of most protein bioinformatics

We can now super charge the analysis using structures

Foldseek – virtual centers are useful to find structural relevant nearest neighbors

1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer (by numerically sorting seq-*k*mer pairs)

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence ____ per group (the longest)

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence ____ per group
- 3. It aligns each sequence in the group *only with the centre sequence*

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence e per group
- 3. It aligns each sequence in the group *only with the centre sequence,* **instead** of with all sequences in the group

- 1. Linclust selects **20** *k*-mers per sequence and finds groups of sequences sharing a *k*-mer.
- 2. It selects one centre sequence e per group
- 3. It aligns each sequence in the group **only with the centre sequence**, instead of with all sequences in the group
- 4. Sequences similar enough to a centre sequence will form a cluster

