

Workshops series on enabling
AI approaches in biological research

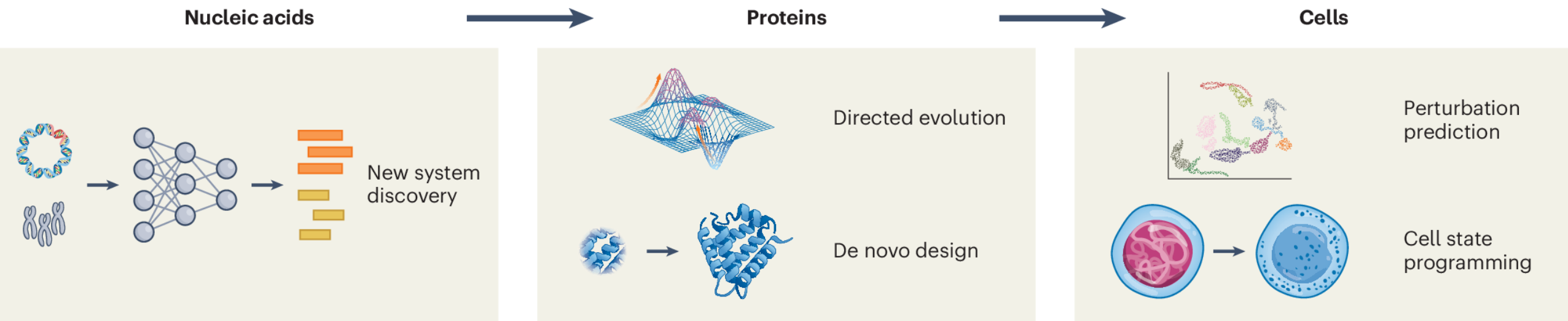
Workshop 1: Modern computational tools for molecular biosciences

18 February 2026 (Weds)

Darwin Building B05, UCL

Why?

- AI has revolutionised molecular biosciences
- Training on how to use these AI tools & how to interpret their results
- Grow a **community** of researchers in different fields but united in their interest in Computational & AI methods applied in Biology



What?

Supported by UCL Grassroot Research Culture Fund

Not just teaching – network with teas/coffees/snacks!!

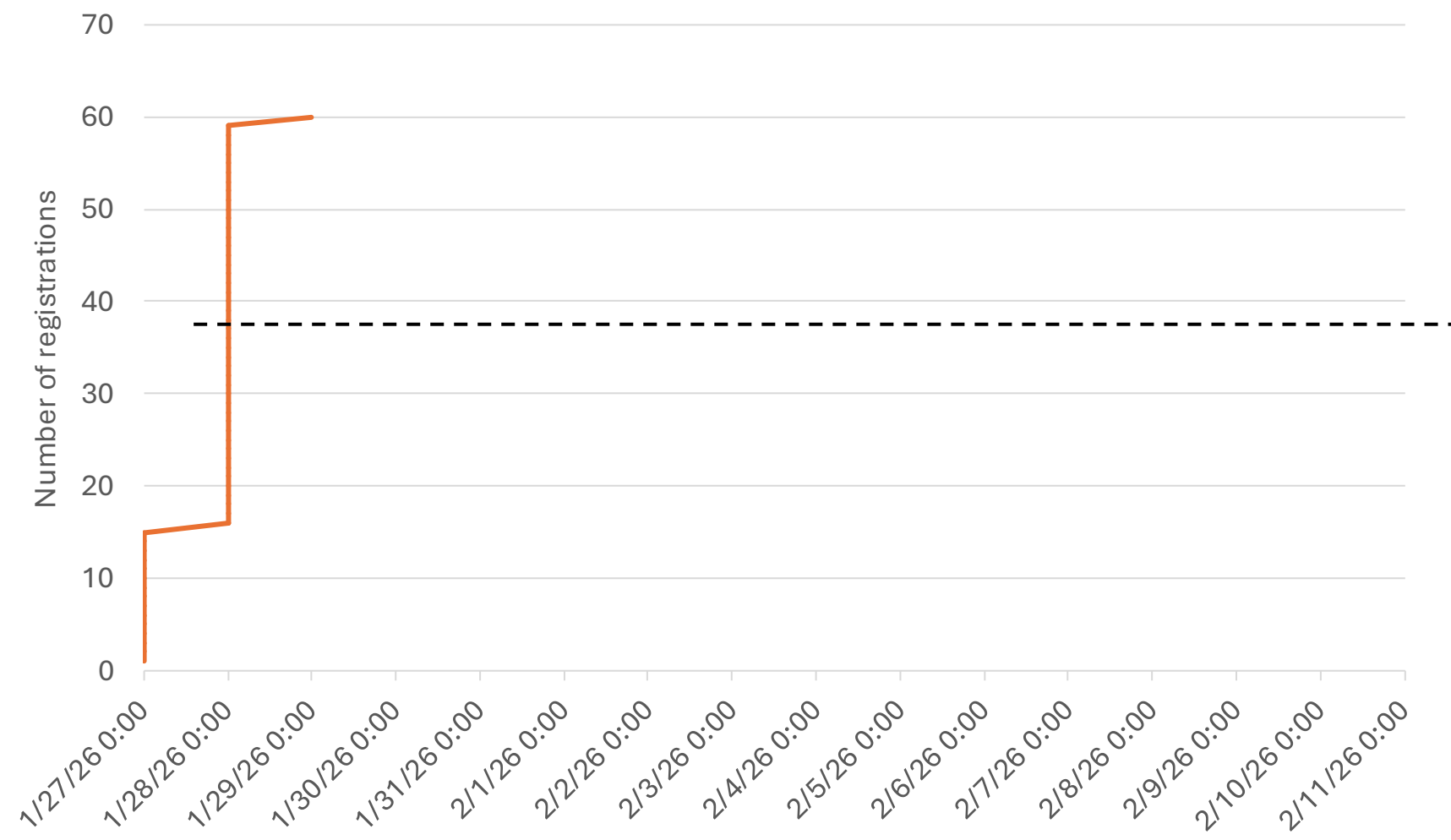
All workshops on Wednesday 2-5pm in Bloomsbury campus.

Workshops in May-July registration: To open in April!

Workshop	Date	Topic
1	18-Feb-26 (Weds)	Modern computational biology, databases and tools
2	25-Feb-26 (Weds)	Introduction to machine learning
3	18-Mar-26 (Weds)	Code Clinic I
4	11-Mar-26 (Weds)	Using AI tools for structural prediction
5	13-May-26 (Weds)	Using large language models (LLMs)
6	20-May-26 (Weds)	Code Clinic II
7	10-Jun-26 (Weds)	Using AI tools for protein design
8	17-Jun-26 (Weds)	Code Clinic III
9	01-Jul-26 (Weds)	Using AI tools for transcriptomic analysis
10	08-Jul-26 (Weds)	Code Clinic IV

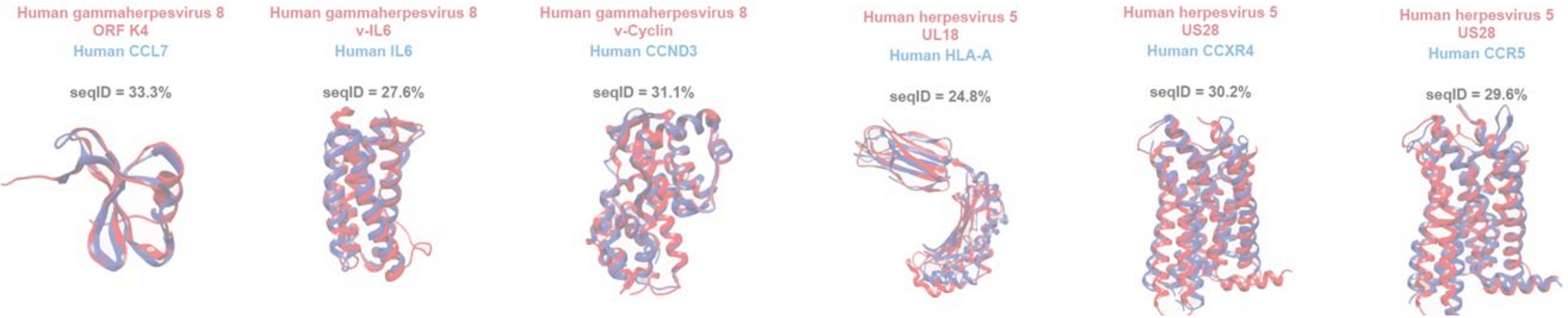
Code Clinics: Bring along your problems, we will try to help!

Thank you all for the enthusiasm!

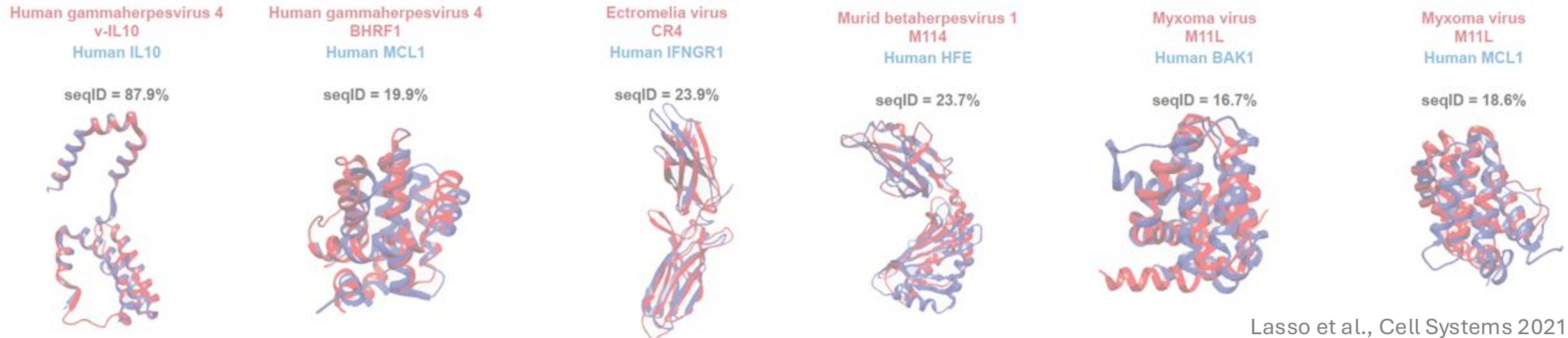


Agenda today

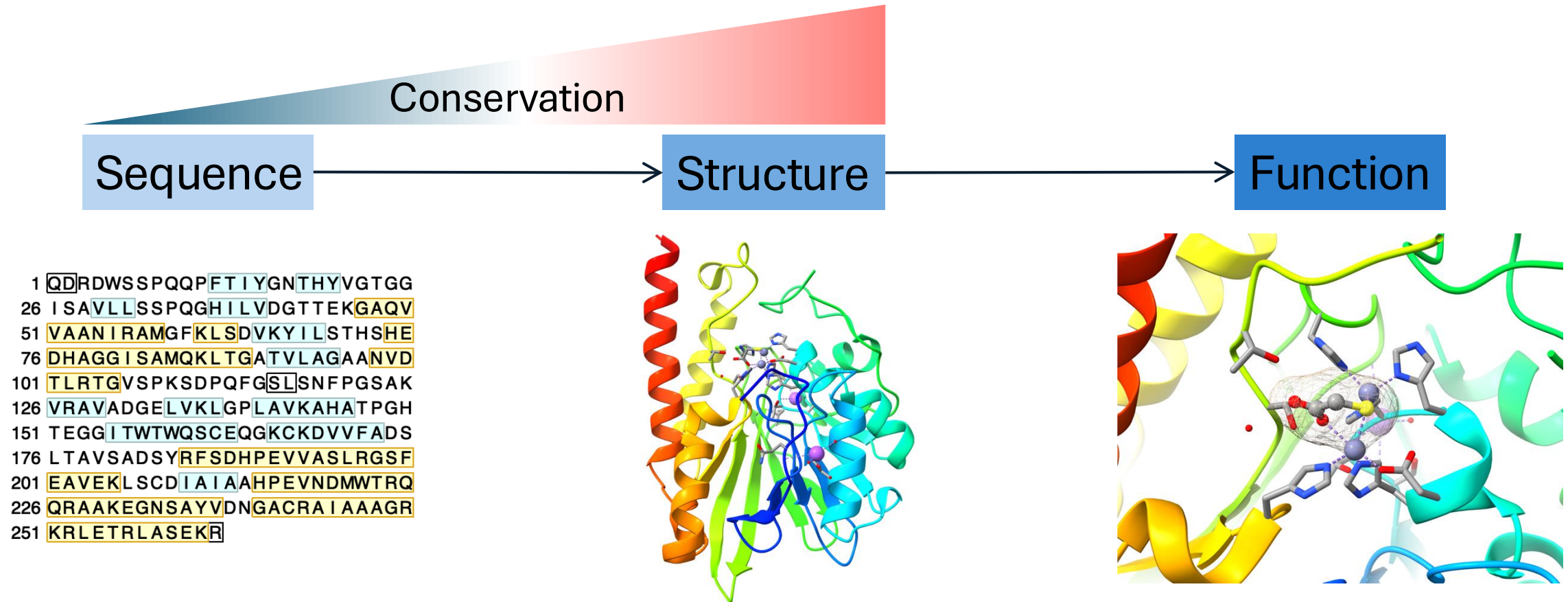
- 14:00-14:20 Background & Introduction to the workshops
- 14:20-15:00 FoldSeek
- 15:00-15:15 Break
- 15:15-16:15 TED/CATH
- 16:15-17:00 Networking (Tea, Coffee & Cookies)



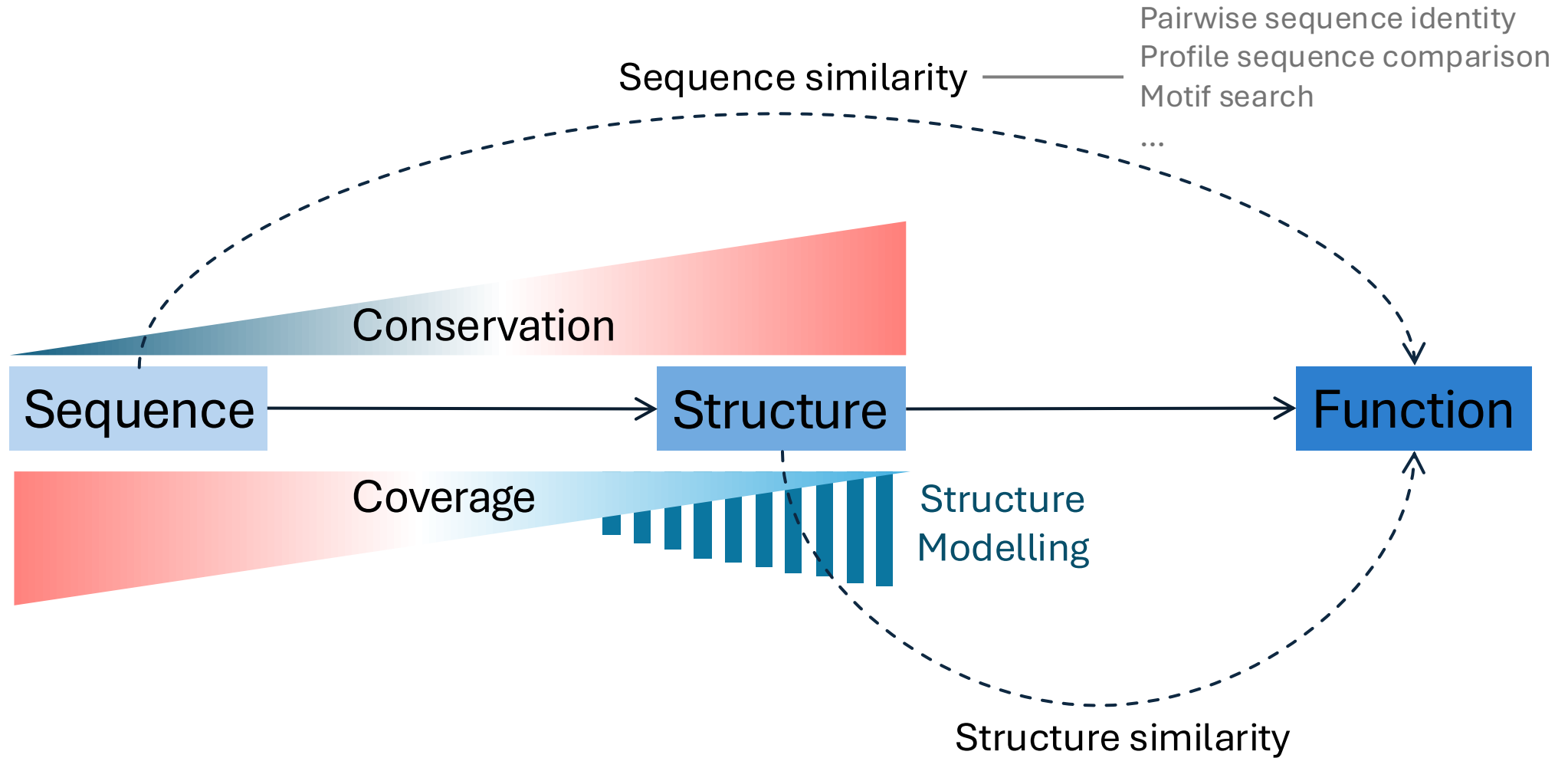
Structural Similarity Searches using Foldseek: What, Why, and How



Leveraging Structural Similarity to Predict Protein Function



Leveraging Structural Similarity to Predict Protein Function



Popular Structural Alignment Tools: DALI

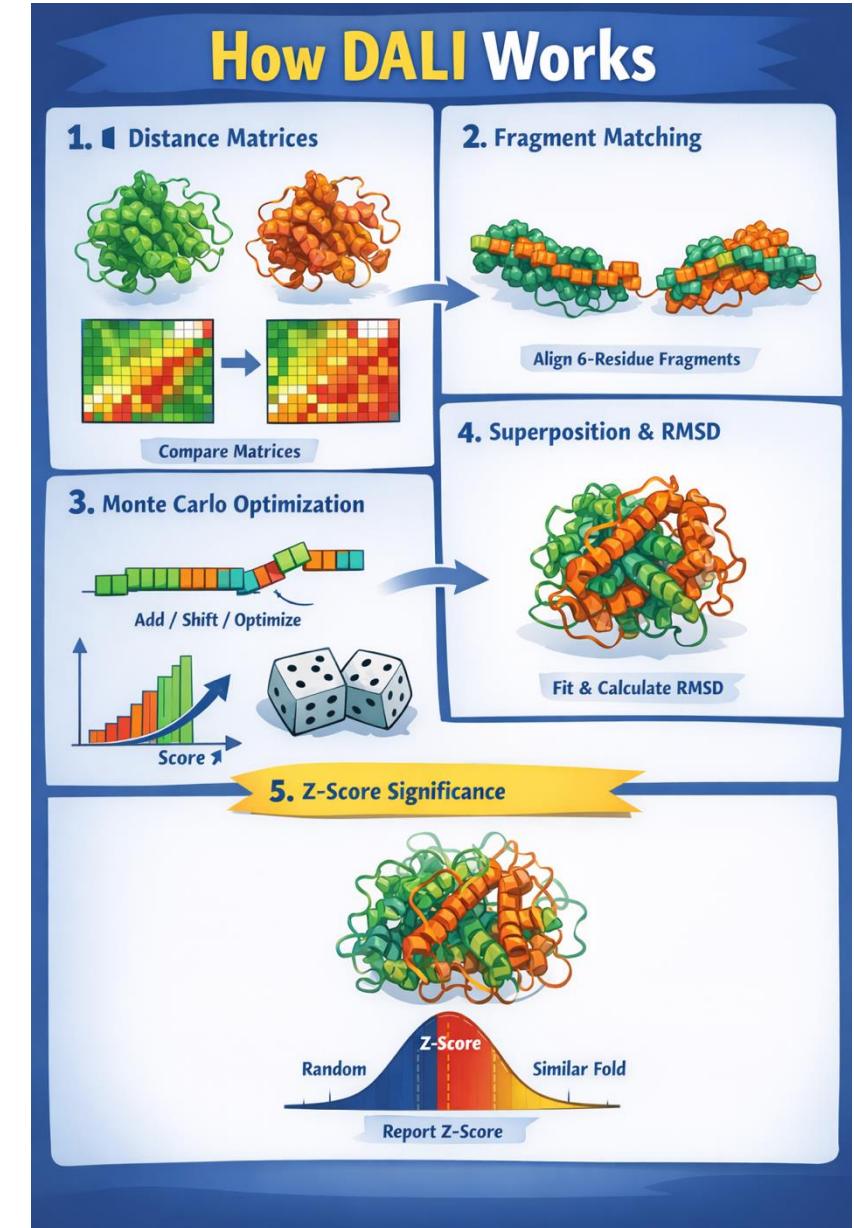
- Distance-matrix Alignment, a classic & highly sensitive method
- Among the most sensitive methods
- Structural similarity score: Z-score
 - $Z > 2$: possibly meaningful
 - $Z > 8$: strong structural similarity
 - $Z > 20$: essentially the same fold

<http://ekhidna2.biocenter.helsinki.fi/dali/>

DALI

PROTEIN STRUCTURE COMPARISON SERVER

About	PDB search	PDB25	AF-DB search	Pairwise	All against all	Tutorials	References	Statistics	Download
-------	------------	-------	--------------	----------	-----------------	-----------	------------	------------	----------



Popular Structural Alignment Tools: TM-align

- Template Modeling (TM) align structures to maximize the TM-score (global similarity)
- More sensitive to global topology than to local variations
- 20x faster than DALI
- TM-score [0-1]
 - Length independent
 - Give more importance to smaller distance errors
 - < 0.2: No structural similarity
 - 0.3 - 0.5: Weak structural similarity (questionable)
 - > 0.5: Same fold
 - > 0.7: Very similar structures

Major problems of RMSD metric
(Root Mean Square Deviation)



<https://aideepmed.com/TM-align/>



TM-align

Quick & Accurate Structural Alignment

How TM-align Works

Protein Structure Alignment Process

1. Initial Alignment Seeds

Heuristic Starting Alignments



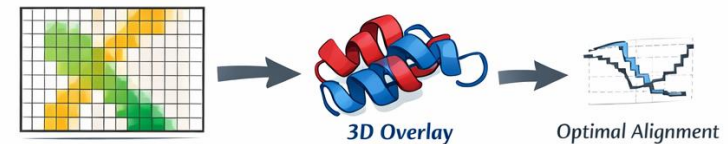
2. Superposition

Optimal Rotation & Translation



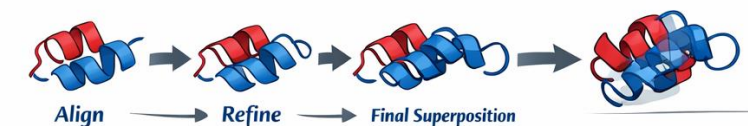
3. Dynamic Programming

TM-score Matrix & Alignment



4. Iterate to Converge

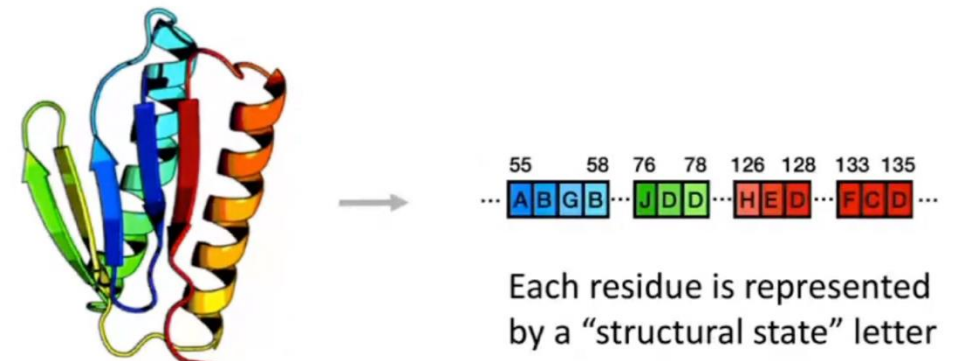
Repeat Until Best TM-score



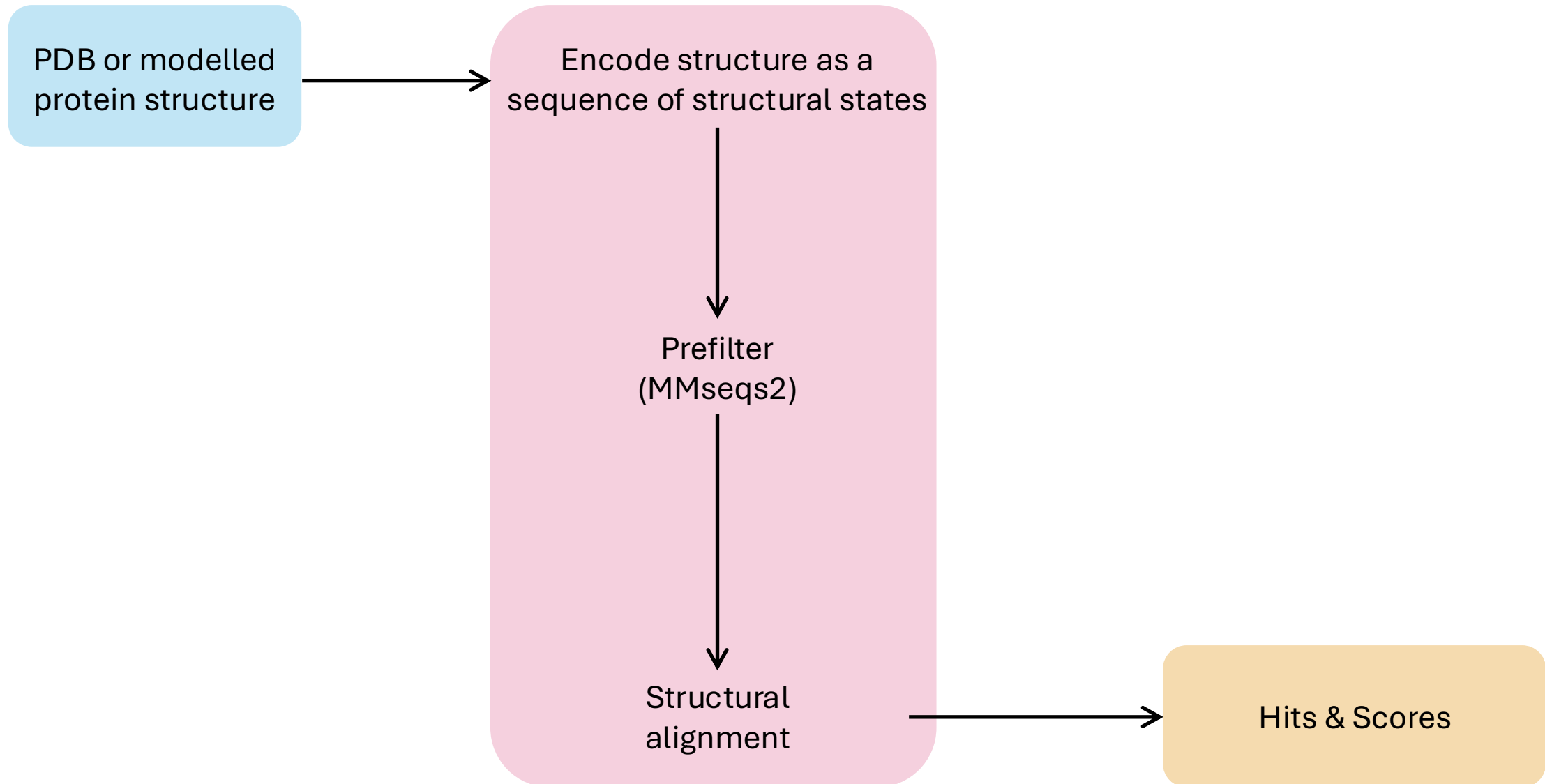
Maximize TM-Score for Best Structural Alignment

Foldseek: Structure Comparison at Scale

- Current protein structure prediction facilitates large-scale modelling (at near-experimental quality?)
 - The AlphaFold Database: 214M protein structures
 - The ESM Metagenomic atlas: 772M protein structures
- State-of-the-art structural similarity searches are not designed to cope for today's scale
 - TM-align search on a database with 100M entries: 1 month (1CPU core)
- To increase speed, one can describe residues in a protein structure using a structural alphabet, and compare structures using sequence alignments



The Strategy Behind Foldseek

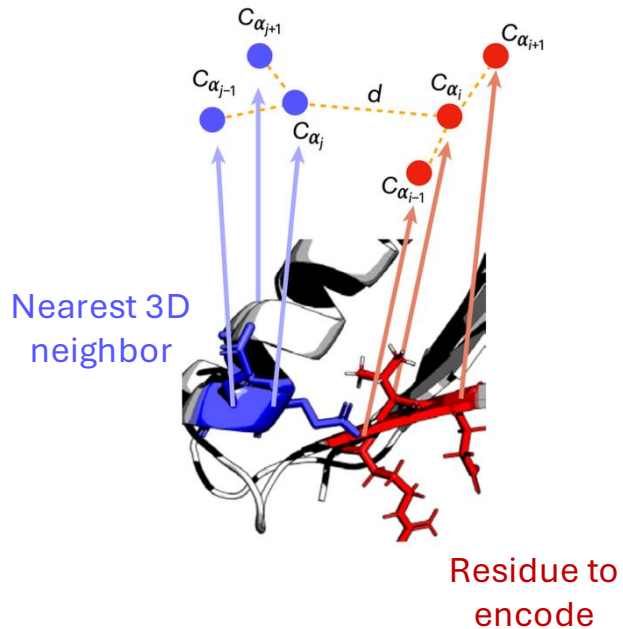


The Strategy Behind Foldseek

PDB or modelled
protein structure

Encode structure as a
sequence of structural states

Each residue is represented by 10
geometrical features representing
its spatial relationship with the
nearest neighbor

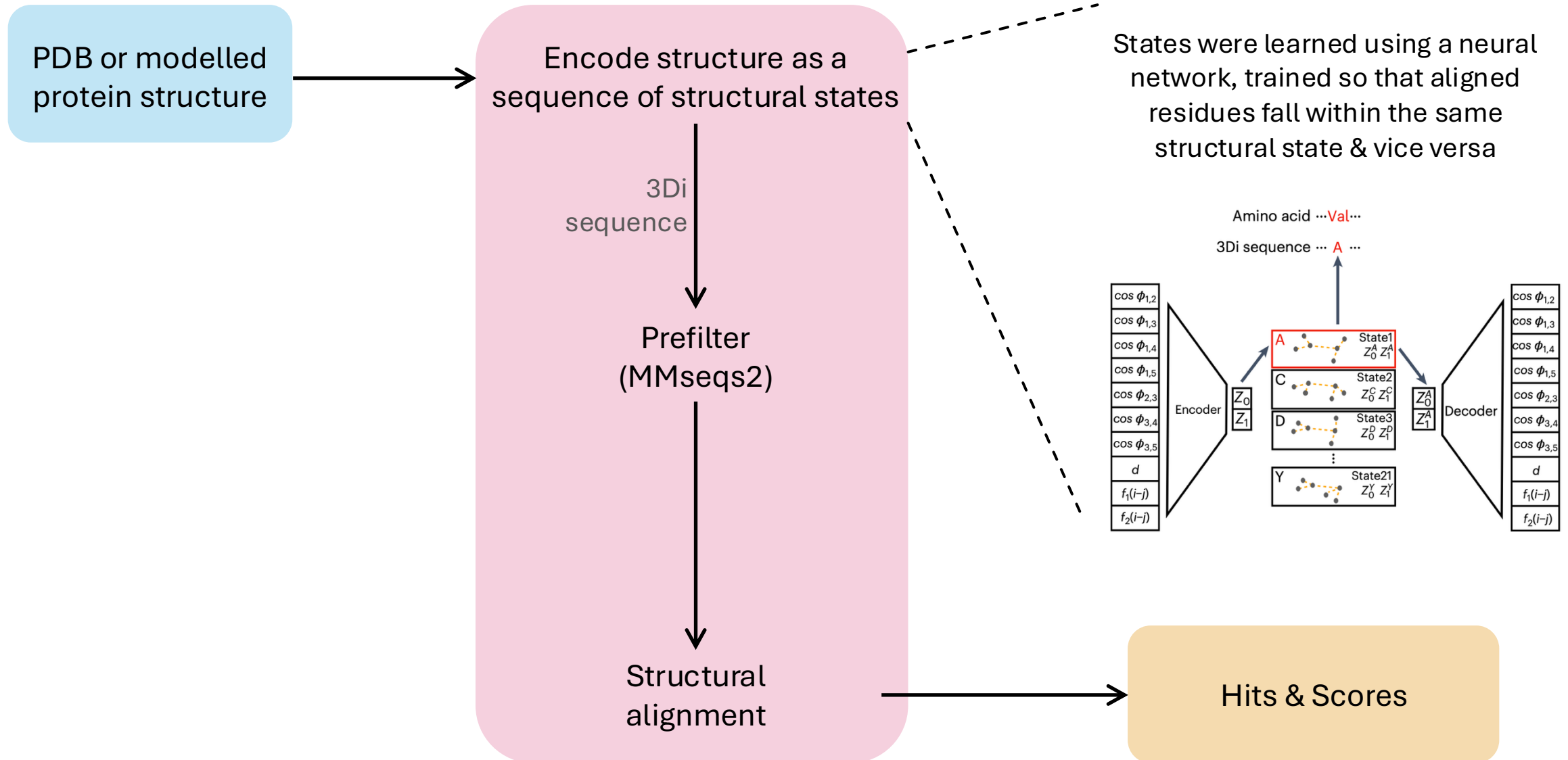


Prefilter
(MMseqs2)

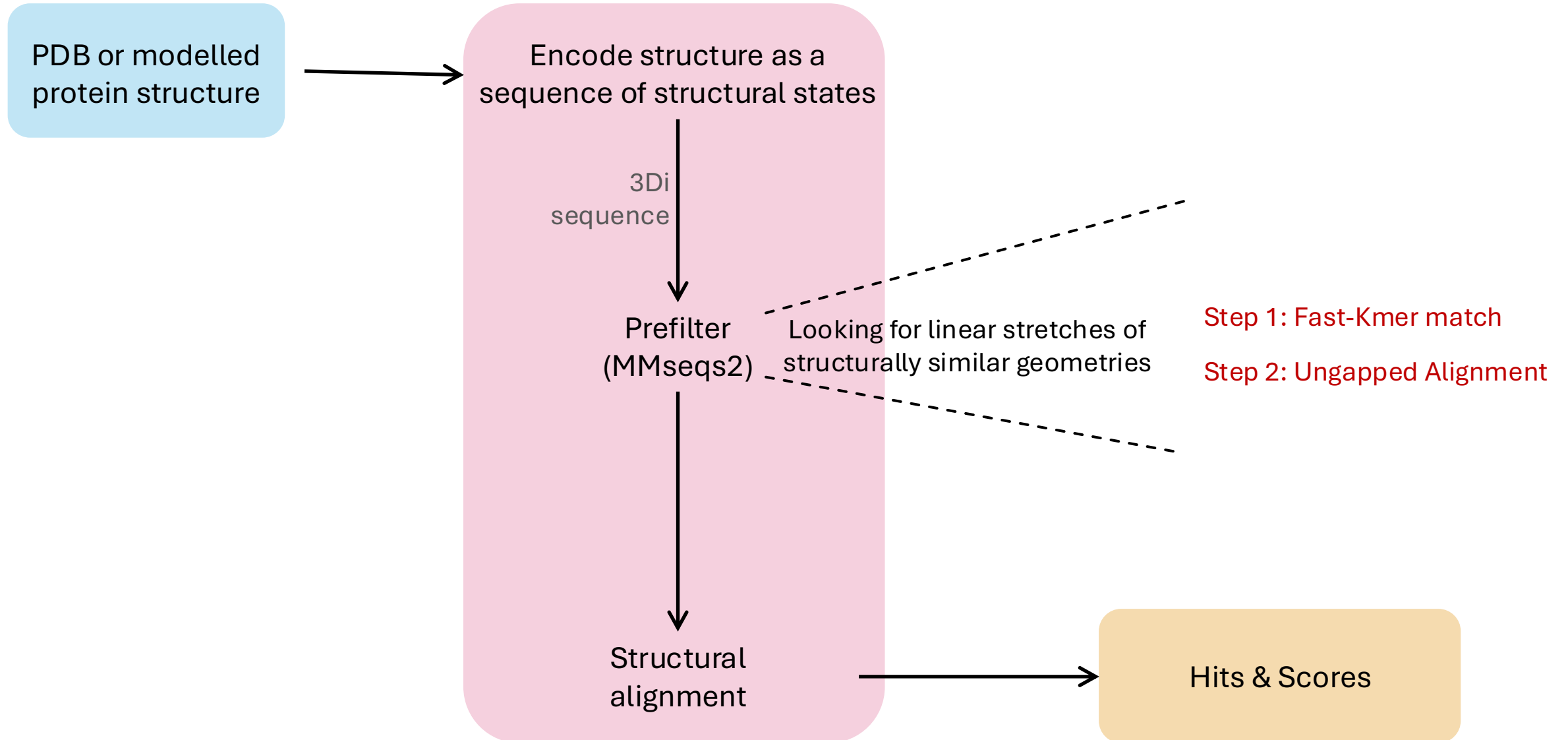
Structural
alignment

Hits & Scores

The Strategy Behind Foldseek



The Strategy Behind Foldseek



The Strategy Behind Foldseek

PDB or modelled
protein structure

Encode structure as a
sequence of structural states

3Di
sequence

Prefilter
(MMseqs2)

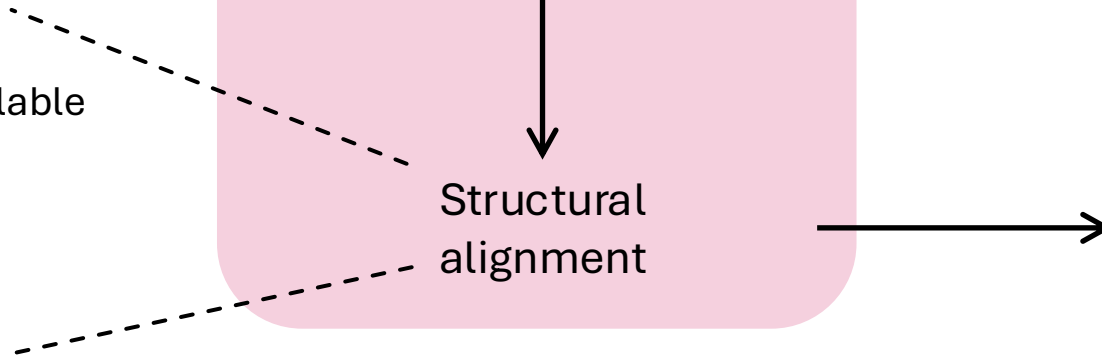
Hits

Structural
alignment

Hits & Scores

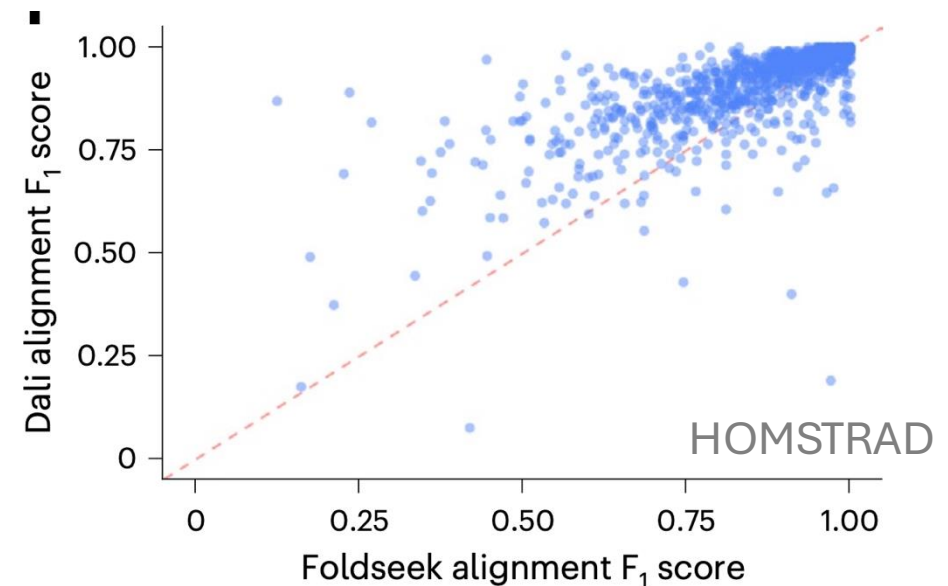
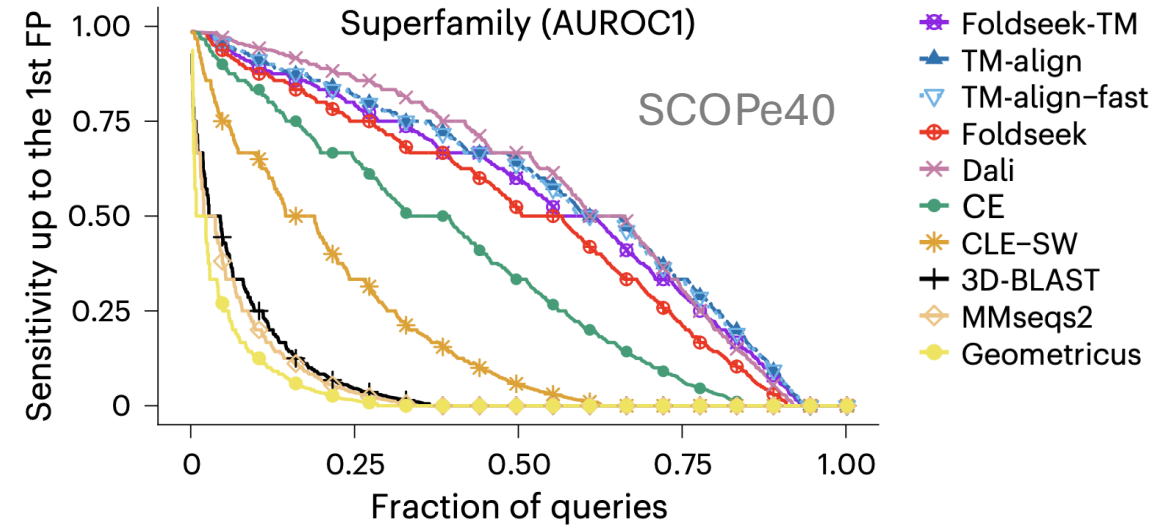
Gapped alignment, available
methods:

- 3Di/AA
- TM-align
- LoL-align



Foldseek Performance

- Accuracy varies from benchmark to benchmark
- It is not as accurate as TAlign & DALI but not falling far behind them
- It is remarkably faster. On AF dataset:
 - 184,600 faster than Dali
 - 23,000 faster than TM-align



But Then LoL-align Happened...

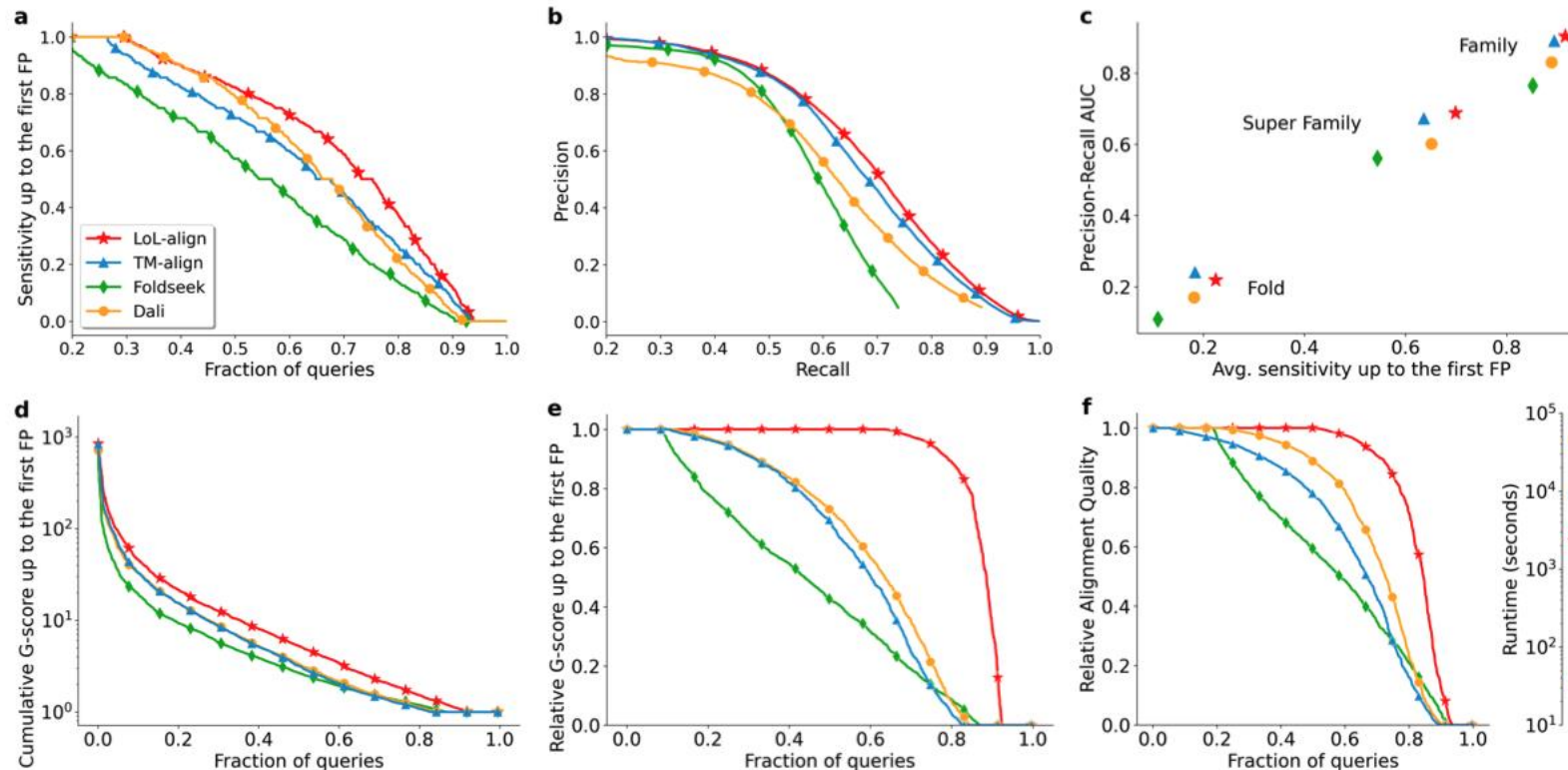
New Results

[Follow this preprint](#)

LoL-align: sensitive and fast probabilistic protein structure alignment

Lasse Reifenrath, Michel van Kampen, Gyuri Kim, Soo Hyun Kim, Mohammadreza Radnezhad, Milot Mirdita, Martin Steinegger, Johannes Söding

- A distance-based algorithm that maximizes a **Local-Log-odds** function, given their intra Ca-Ca distances
- More sensitive in detecting remote homologs than TM-align & DALI and 5-20 times faster
- Integrated in FoldSeek



From Foldseek to FoldDisco: a Cambrian Explosion of Structural Similarity Search



Foldseek

Protein structure

Foldseek-Multimer

Protein complexes

FoldMason MSA

Structure-based Multiple
Sequence Alignments

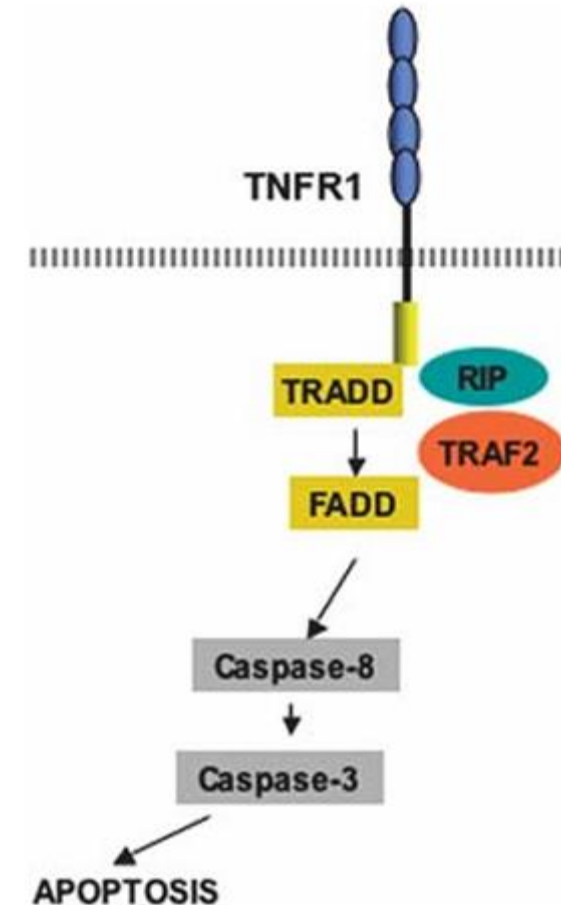
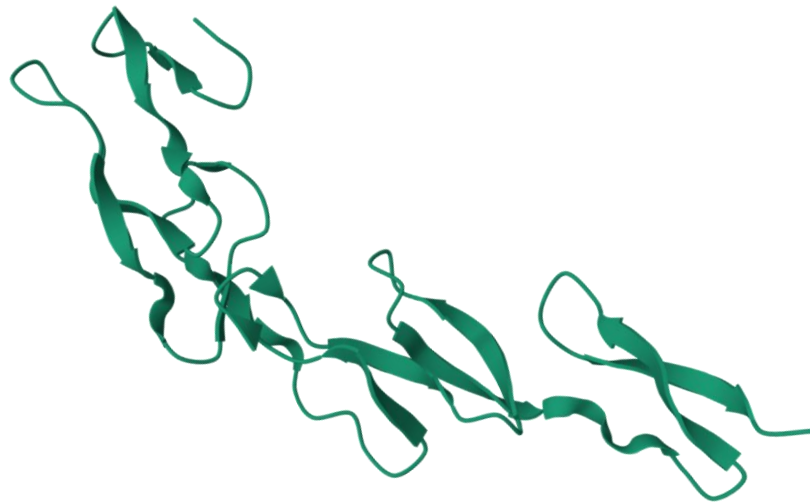
FoldDisco

Structural domains

FoldSeek examples

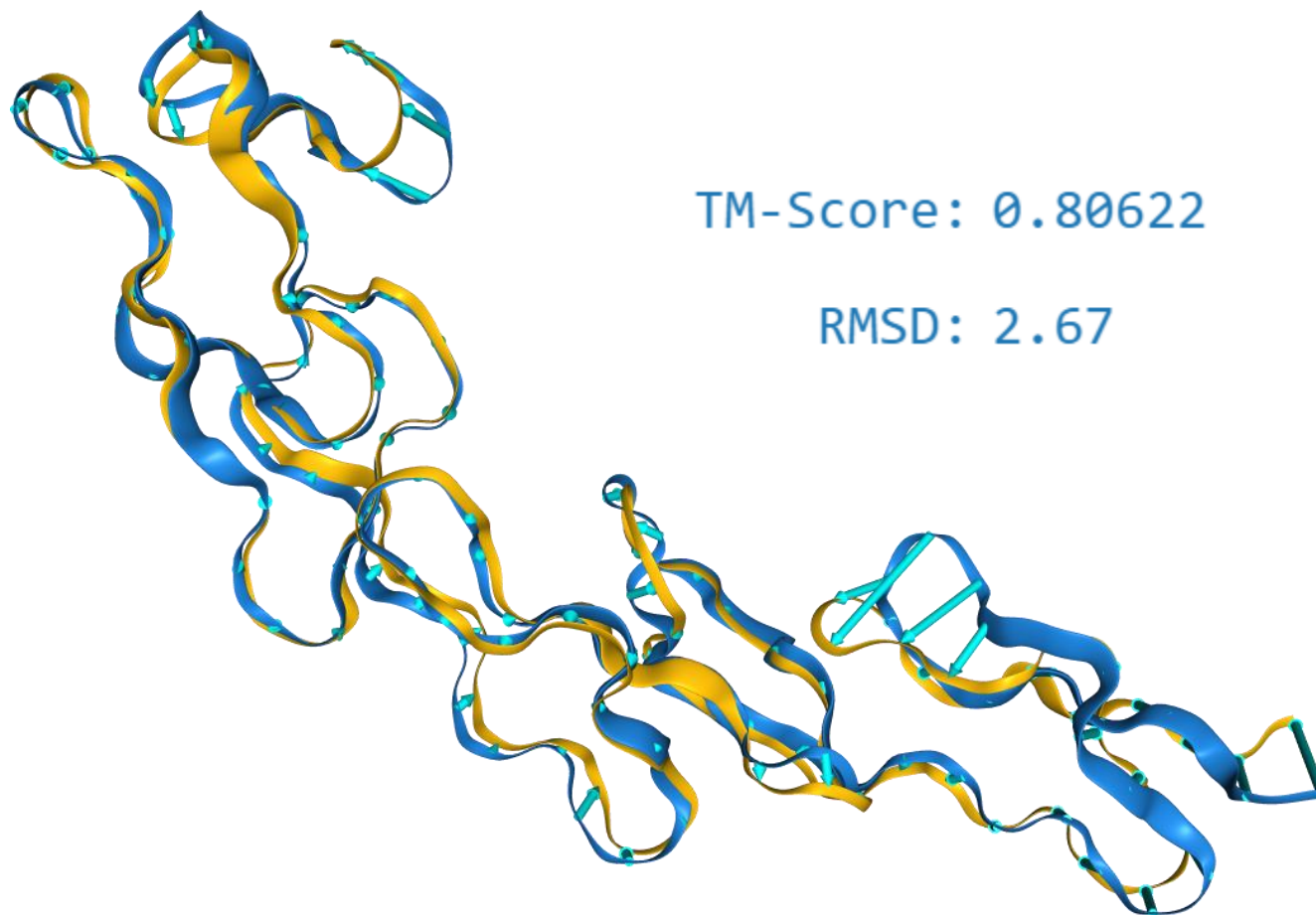
Cowpox virus TNF receptor mimic (PDB 2uwi)

- Binding of TNF α to TNF receptors on infected cells promote apoptosis \rightarrow protection
- Cowpox virus expresses a TNFR mimic CrmE (PDB 2uwi)



Structure databases available on the FoldSeek webserver

Database	Description
BFVD (Big Fantastic Virus Database)	Predicted structure database focussed on viral proteins
AFDB-PROTEOME	Proteome-level predicted structures from AlphaFold DB
AFDB-SWISSPROT	AlphaFold DB predicted structures but restricted to UniProt/SwissProt (i.e. manually reviewed protein entries)
AFDB50	AlphaFoldDB clustered at 50% sequence identity
BFMD (Big Fantastic Multimer Database)	Predicted structure database focussed on multimers (~300K multimer predictions from several community efforts)
CATH50	CATH structural domain classification – clustered at 50% Seq ID
GMGCL_ID	Gene catalog–based clustered structures from global metagenomic datasets
MGnify_ESM30	Protein sequences from the MGnify microbiome database. Structural models generated using ESMFold. Clustered at 30% Seq ID
PDB100	Non-redundant set of structures from the Protein Data Bank



TM-Score: 0.80622

RMSD: 2.67

Blue: query (i.e. Cowpox
TNFR mimic CrmE)

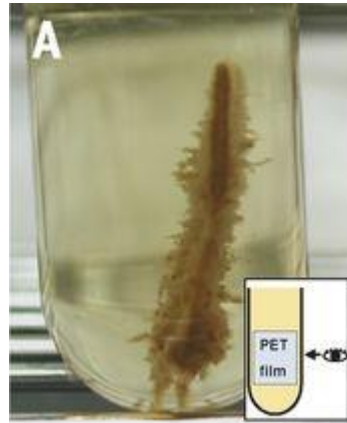
Gold: target (i.e. Human
TNFR1B [P20333])

Structural Similarity Metrics Reported by Foldseek

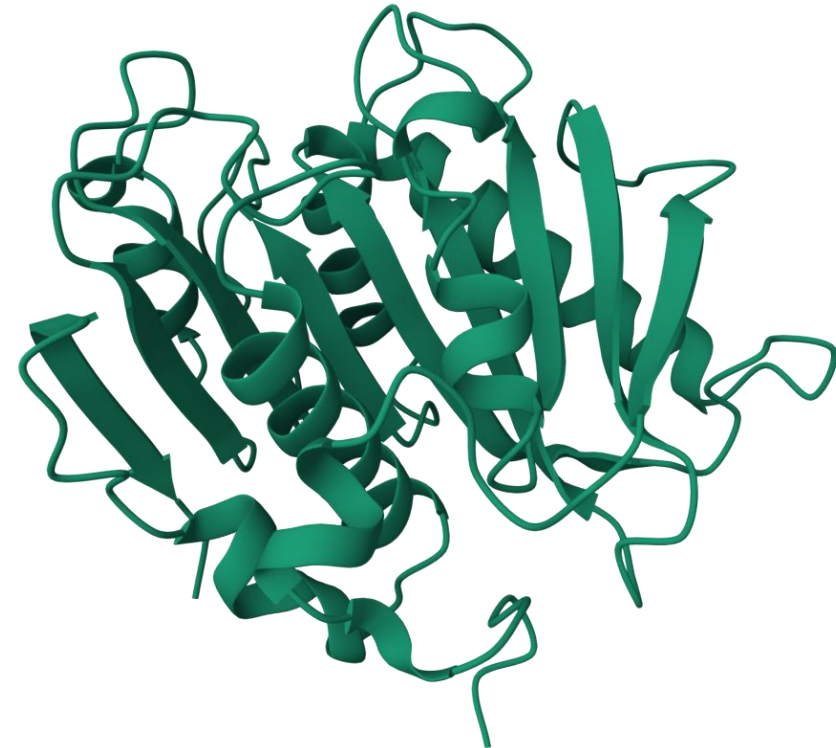
Metric	What it measures	Range / Units	Length-dependent	Alignment-dependent	Sensitive to local errors	Best used for
Probability	Confidence that the hit is <i>not random</i>	0–1	Yes (implicitly)	Yes	No	Ranking hits, filtering true positives
E-value	Expected number of random hits with equal or better score	≥ 0	Yes	Yes	No	Database searches, statistical significance
TM-score	Global fold similarity	0–1	Normalized	Yes	Low	Fold detection, global structure comparison
RMSD	Mean C α distance after superposition	Å	Yes (badly)	Yes	Very high	Local accuracy, refinement assessment
IDDT	Local distance agreement without superposition	0–1	No	Yes	High (locally)	Local model quality, per-residue accuracy

PET hydrolase (PDB 6eqg)

- Polyethylene terephthalate (PET) degrading hydrolase from *Ideonella sakaiensis* (PDB 6eqg)



[Yoshida et al Science \(2016\)](#)
[doi:10.1126/science.aad6359](https://doi.org/10.1126/science.aad6359)



☐ Target

☐ [GMGC10.019_473_839.BMUL_4291_trun_0](#)

Prob.

Seq. Id.

E-Value

Position in query

Alignment

1.00

17.6

8.23e-9

24

212

TM-Score: 0.56974

RMSD: 5.12

Colorscheme
Clustal2

CLEAR SELECTION

Select target residues to highlight their structure.
Click on highlighted sequences to dehighlight the corresponding chain.

A → GMGC10.019_473_839.BMUL_4291_trun_0

Q 24 VRSFTVS---RPSGYGAGTVY-YPTNAGGTVGAIIVPG---YTARQSSI---KWWGPRLASHGFVVITIDT-----

+ + S + + + +++G A+ ++ G +R ++ + WG LA+HG++ + +D+

T 19 PQHVEIPGAGLSSSPAPLNGFVFAPDGAGPHPAVMMHGC GGAYGRDGT LNPRHRMWGEFLAAHGYLALMLDSFGPRGVR

Q 85 -----NSTLDQPSSRSSQQAALRQVASLNGTSSSPIYGVDTARMGMVGMWSMGGGSLISAAN-NPS---LKAAAP

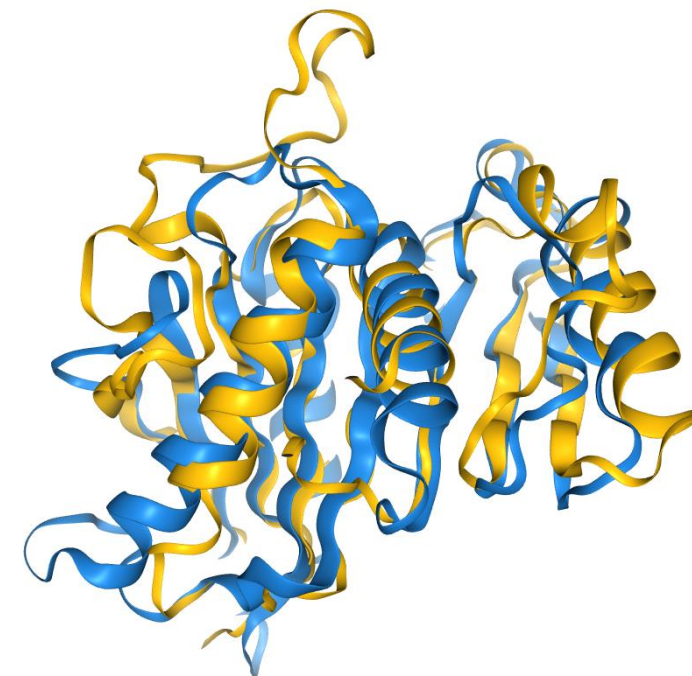
TL + R+ AAL + T +V +R++++GWS G+G L + P AA

T 99 ELCTQPMKERTL-KEHDRAGDADAALAYL----RT-----RPEVAAGRIALLGWSHGAGSVLATITGQRP GAPRYDAAIA

Q 154 QAPWDS----STNFSSVT--VPTLIFACENDSIAPVNSSALPIYDS-----MSRNAKQF-LEINGGSHSCA

P S VP L++ E D P + + S R + + A

T 169 FYPGCSARARHP--EDFHPAVPLLLLIGEADDWTPAEA-CRVLAASANARGDSVRLVTYPDTFHDF--DNPA



☐ [GMGC10.213_526_639.DPP_trun_3](#)

1.00

16

2.46e-7

22

223

☐ [GMGC10.305_916_099.UNKNOWN_trun_1](#)

1.00

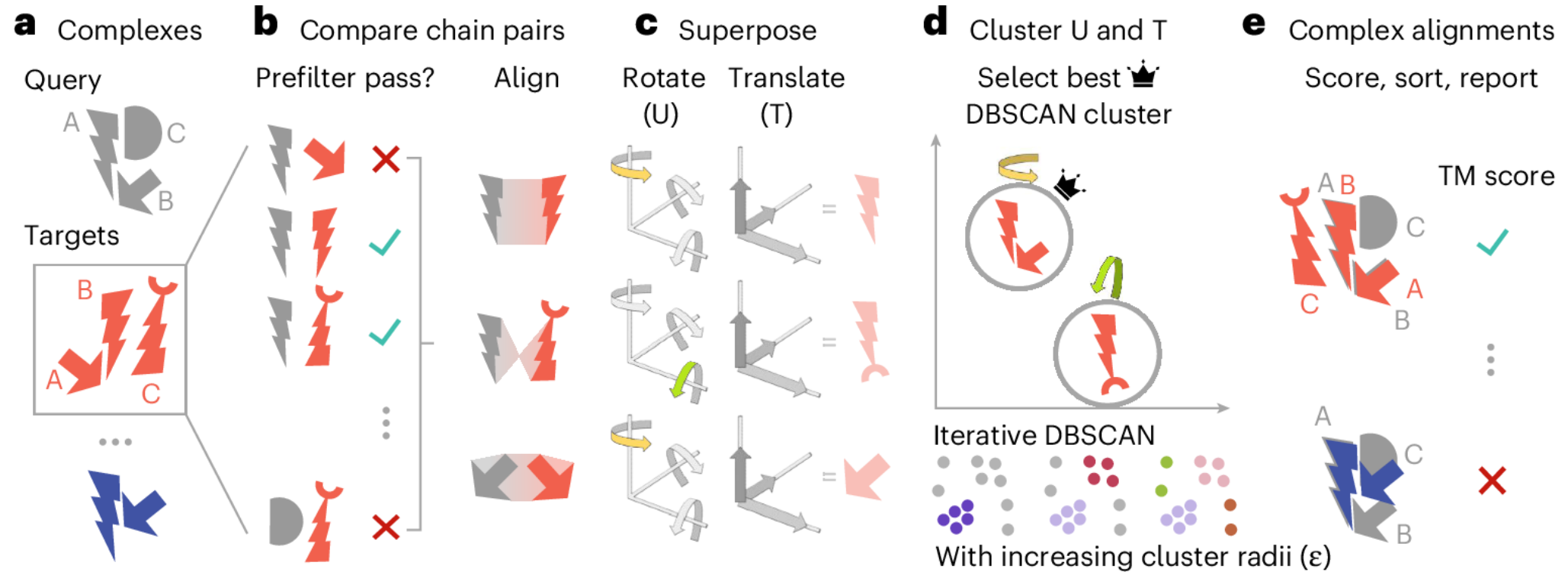
13.5

1.92e-7

91

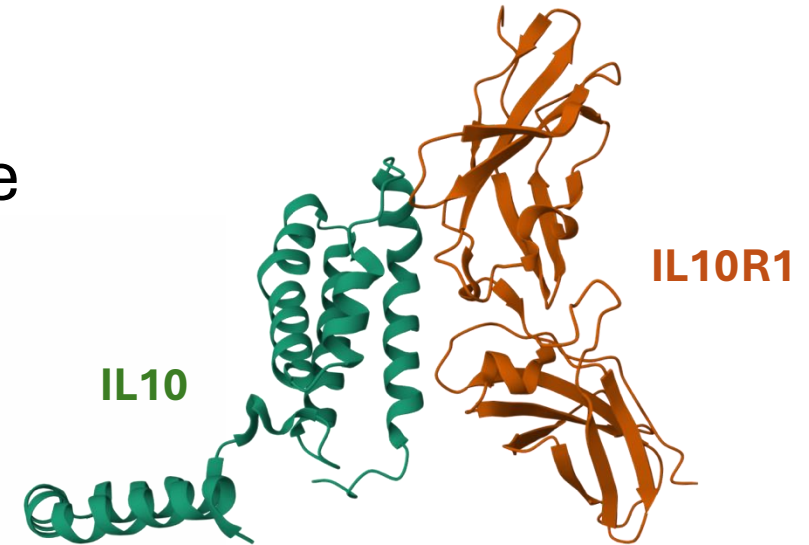
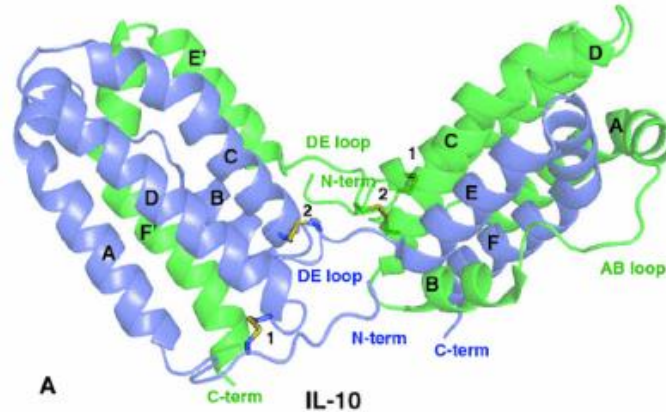
225

FoldSeek-Multimer

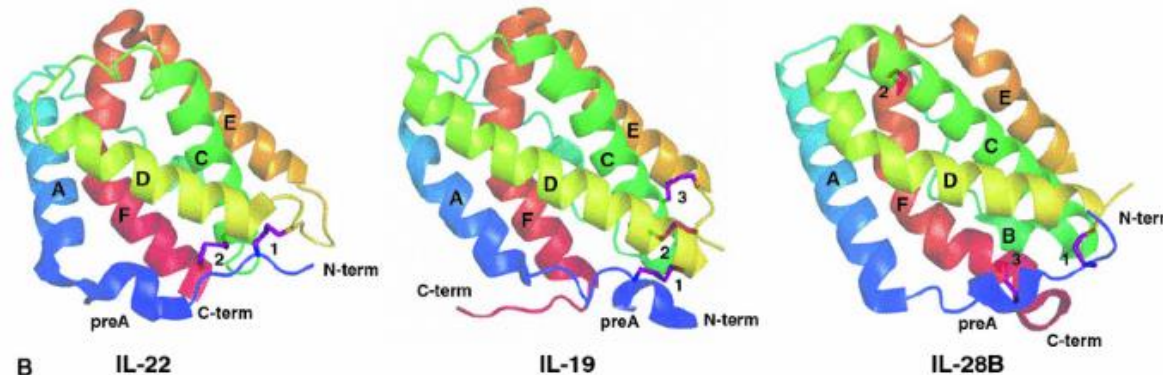


Human IL-10/IL-10R1 complex (PDB 1j7v)

- IL-10 is a well-known immunoregulatory cytokine that ‘dampens’ the immune response
- The IL-10 superfamily contains other cytokines sharing similar structural folds











Josephson et al. Immunity. 2001 Jul;15(1):35-46.



Trivella et al. Cell Mol Life Sci. 2010 Sep;67(17):2909-35.

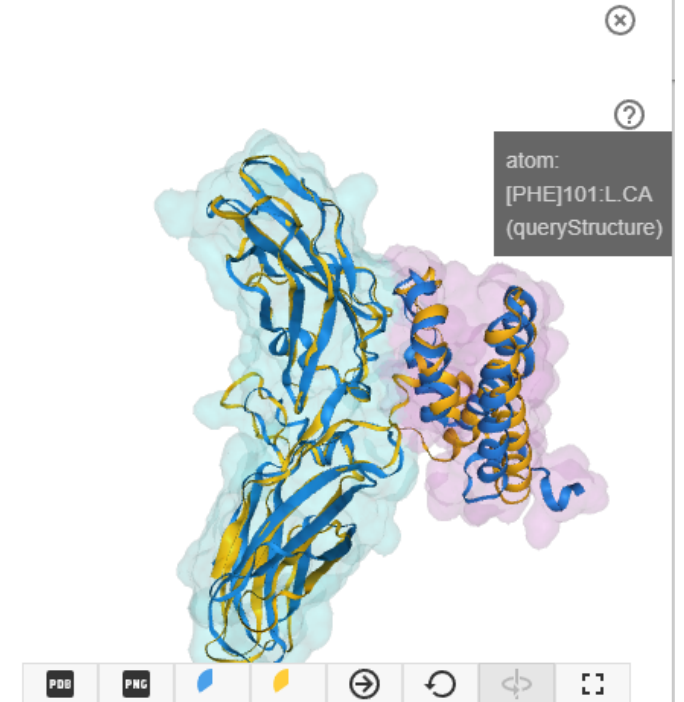
Left Clicking these links would direct you to UniProt and AFDB entries of the two proteins

	Complex		Chain								
	<input type="checkbox"/>	qTM	tTM	Chain pairing	Scientific Name	Prob.	Seq. Id.	E-Value	Position in query		Alignment
IL10	<input type="checkbox"/>	0.84	0.39	L → ProtVar_P22301_Q1365...	Homo sapiens	1.00	100	2.83e-13			
				R → ProtVar_P22301_Q1365...	Homo sapiens	1.00	97	7.19e-36			
IL20	<input type="checkbox"/>	0.71	0.35	L → ProtVar_Q9NYY1_Q9UH...	Homo sapiens	0.94	28	1.93e+0			
				R → ProtVar_Q9NYY1_Q9UH...	Homo sapiens	1.00	22.3	1.93e-13			
IL19	<input type="checkbox"/>	0.71	0.35	L → ProtVar_Q9UHD0_Q9UH...	Homo sapiens	0.44	21	5.75e+0			
				R → ProtVar_Q9UHD0_Q9U...	Homo sapiens	1.00	22.3	3.18e-13			

e.
he corresponding chain.

LKESLLEDFKG---YLGQCALSEMIQFYLEEVMPPAENQDPDIKAHVNSLGENLKTLR
S LE ++ C ++ FY++ V+ + +P I ++S++ ++ ++
:-LSTLETLQIIKPLDVCCTKNLLAFYVDRVFKDHQEPNPKILRKISSIANSFLYMQ

PQ-QSESTCYEALLRYGIESWNSI--SQCSQTLSDYDLTAVTLDLYHSNGYRVRVA
Q Y V + YG + W + + + DL+A T D Y + Y A+V+A
EGLQGKVTYTVQYFIYGQKKWLNKSECRNINRTYCDLSAETSD-YEHQ-YYAKVKA
TLTVGSVNLEIHNGFILGKIQLPRPKMAPAQDT---YESIFSHFREYEIAIRKVPQGQF
V L I + P +D + I+S +Y +++
QIGPPEVALTTDEKSISVVLTAPEKWKRNPELDPVSMQIQIYS-NLKYNVSVLNTK-SN
FCVQVKPSVASRSNKGMSKEECI
V+V+ V ++ S+++C
YCVHVESFVPGPPRRAQPSEKQCA

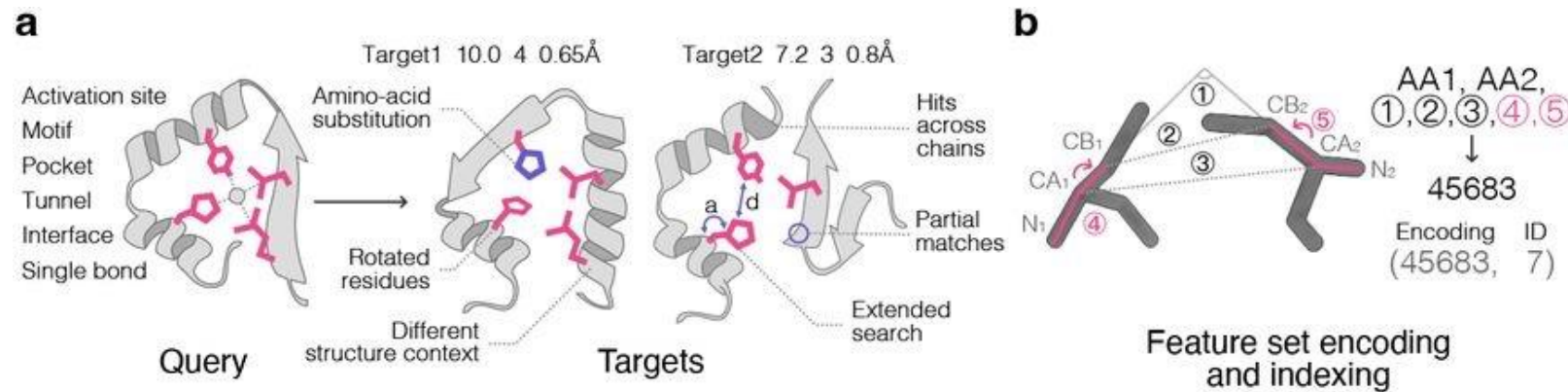


ProtVar is a protein-variant structure database and annotation resource that links **genetic variants (mutations)** to **protein structures and functional features**.

ProtVar reference:
Stephenson et al NAR 2024
52(W1):W140-W147

FoldDisco

[Structural motif search across the protein-universe with Folddisco | bioRxiv](#)

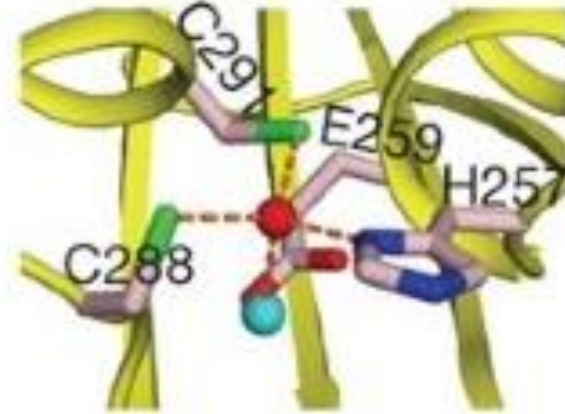


For each pair of residues involved:

- the amino acid (AA1 and AA2)
- distance between their $C\alpha$ atoms
- distance between their $C\beta$ atoms
- intersecting angle between the $C\alpha$ - $C\beta$ vectors.
- two dihedral angles ($N1-C\alpha1-C\beta1-C\beta2$ and $N2-C\alpha2-C\beta2-C\beta1$)

APOBEC3G bound to zinc (PDB 3e1u)

- APOBEC3G is an antiviral protein known for restricting HIV-1 replication
- Cytidine deaminase with zinc at its active site (similar to other enzymes binding and catalysing deamination of DNA/RNA bases)



Search

Multimer search

FoldMason MSA

Folddisco search

History

Input protein motif structure (PDB/CIF)

data_3E1U

#

_entry.id 3E1U

#

_audit_conform.dict_name mmcif_pdbx.dic

_audit_conform.dict_version 5.387

_audit_conform.dict_location http://mmcif.pdb.org/dictionaries/ascii/mmcif_pdbx.dic

#

loop_

_database_2.database_id

_database_2.database_code

Selected Motif
A257,A288,A291

Format: <chain ID><position>

LOAD ACCESSION

UPLOAD PDB

Click here to select which binding site to consider (alternatively manually enter the motif)

Filter

Query residues



i.e. remove redundancy

Cluster

☐ Cluster

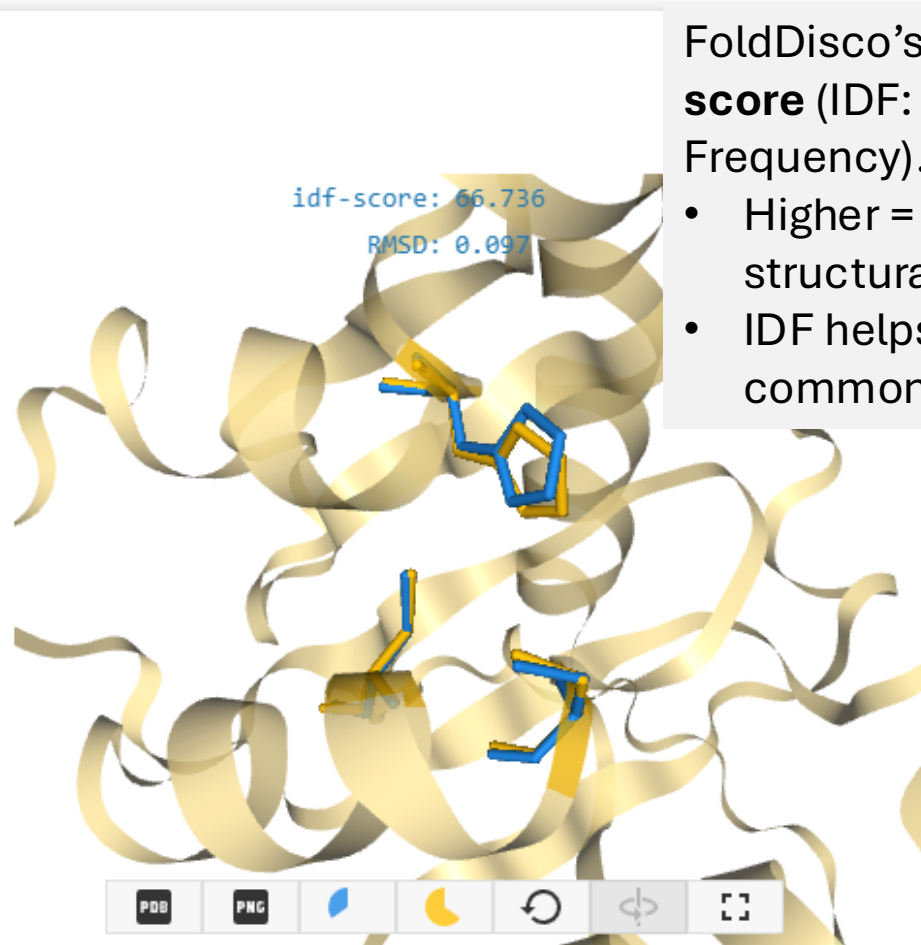
Clustering Min Points

2 minimum how many hits in a cluster?

Clustering Epsilon

8 Smaller value = stricter clustering

<input type="checkbox"/> Target	Description	Scientific Name	IDF-score	RMSD	Nodes	Residues	Structure
<input type="checkbox"/> AF-A0A524MHR0-F1-MODEL...	CMP/dCMP-type deaminase domain-containin...	Candidatus Thorarchaeot...	66.736	0.097	3	A91,A118,A121	
<input type="checkbox"/> AF-A0A535A005-F1-MODEL...						A80,A108,A111	
<input type="checkbox"/> AF-A0A524L9V8-F1-MODEL...						A72,A99,A102	
<input type="checkbox"/> AF-A0A350PD33-F1-MODEL...							
<input type="checkbox"/> AF-A0A2R5GWQ0-F1-MODEL...							
<input type="checkbox"/> AF-A0A3C0IJQ0-F1-MODEL...						A49,A75,A78	
<input type="checkbox"/> AF-A0A1M3E0W4-F1-MODEL...						A55,A83,A86	
<input type="checkbox"/> AF-A0A3N5HE05-F1-MODEL...			66.736	0.108	3	A30,A58,A61	
<input type="checkbox"/> AF-A0A3S0CR58-F1-MODEL...			66.736	0.109	3	A54,A88,A91	
<input type="checkbox"/> AF-A0A5C7Q2I5-F1-MODEL...			66.736	0.109	3	A73,A96,A99	
<input type="checkbox"/> AF-A0A381THS7-F1-MODEL...			66.736	0.110	3	A21,A49,A52	
<input type="checkbox"/> AF-A0A6C0BBK0-F1-MODEL...			66.736	0.111	3	A71,A109,A112	
<input type="checkbox"/> AF-A0A853KLB5-F1-MODEL...			66.736	0.111	3	A76,A104,A107	
<input type="checkbox"/> AF-A0A7C7N9E7-F1-MODEL...			66.736	0.114	3	A71,A99,A102	
<input type="checkbox"/> AF-A0A496MAH3-F1-MODEL...			66.736	0.114	3	A32,A60,A63	

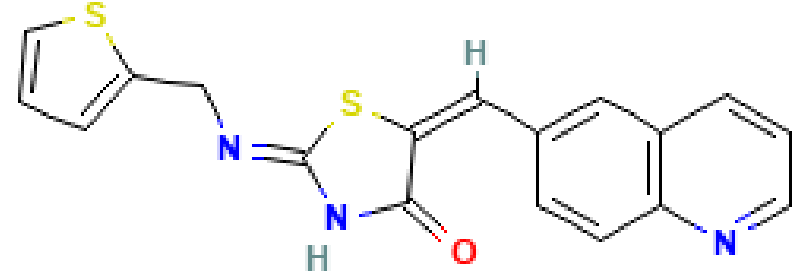


FoldDisco's **motif match score** (IDF: Inverse Document Frequency).

- Higher = more distinctive structural motif match
- IDF helps downweigh common/simple motifs

Number of residue matches

Kinase inhibitor Ro-3306

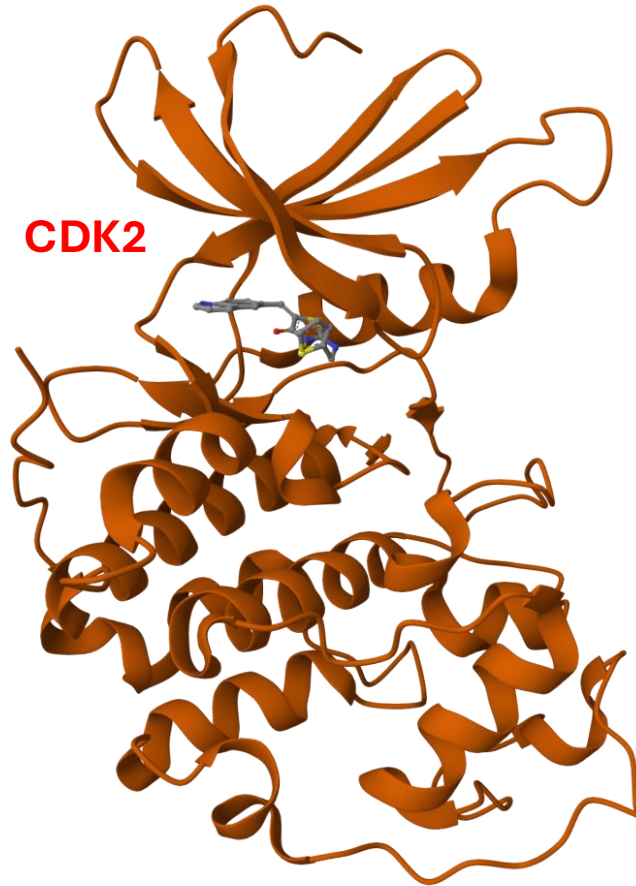


- Ro-3306 is a Cyclin dependent kinase (CDK) inhibitor
- Reported to have off-target effects



PIM1

PDB 5o12



CDK2

PDB 4eon

Additional background slides

Root Mean Square Deviation (RMSD)

- A quantitative measure of the average distance between corresponding atoms of two superimposed structures.
- In structural biology, it is most commonly used to compare protein 3D conformations.

$$\text{RMSD} = \sqrt{\frac{1}{N} \sum_{i=1}^N [(x_i^A - x_i^B)^2 + (y_i^A - y_i^B)^2 + (z_i^A - z_i^B)^2]}$$

Structure A: (x_i^A, y_i^A, z_i^A)

Structure B: (x_i^B, y_i^B, z_i^B)

N = number of aligned atom pairs

x_i, y_i, z_i = coordinates of atom i in the two structures

Critical Limitations

- Length dependence
 - RMSD increases with protein size, even for similar folds.
- Outlier sensitivity
 - Because distances are squared, a few large deviations dominate.
- Global metric
 - Local similarity can be masked by flexible regions.

When to Use RMSD

- ✓ Comparing very similar structures
- ✓ Molecular dynamics trajectory analysis
- ✓ Local region comparison

When to avoid relying solely on RMSD

- ✗ Comparing proteins of different lengths
- ✗ Assessing fold similarity
- ✗ Flexible regions are present

Template Modeling (TM) score

- A length-normalized metric for comparing protein 3D structures
- Designed to answer a specific question: *Do these two protein folds belong to the same overall topology?*

$$\text{TM-score} = \max \left[\frac{1}{L_{\text{ref}}} \sum_{i=1}^{L_{\text{ali}}} \frac{1}{1 + \left(\frac{d_i}{d_0(L_{\text{ref}})} \right)^2} \right]$$

$$d_0 = 1.24 \sqrt[3]{L_{\text{ref}} - 15} - 1.8$$

L_{ref} = length of the **reference protein**

L_{ali} = number of aligned residues

d_i = distance between the i -th aligned residue pair (usually Ca–Ca)

d_0 = **length-dependent scaling factor**

“max” means the score is optimized over all possible superpositions

Feature	RMSD	TM-score
Length dependence	Bad	Normalized
Loop sensitivity	Very high	Low
Fold detection	Poor	Excellent
Comparability	Weak	Strong
Optimization	Single superposition	Global max

IDDT score

- A superposition-free metric that evaluates how well the **local atomic environments** of a model match a reference structure.
- If local distances are preserved, the score is high, even if domains are shifted relative to each other.
- Think of IDDT as checking whether the **local distance map** around each residue is preserved, rather than whether the whole structure overlaps globally.

$$IDDT = \frac{1}{N} \sum_i \frac{1}{M_i} \sum_{j \in neighbours} \text{fraction of thresholds satisfied}$$

N = number of residues

M_i = number of valid neighbours for residue i

