

Unknown sequences and hidden errors in macromolecular structure models

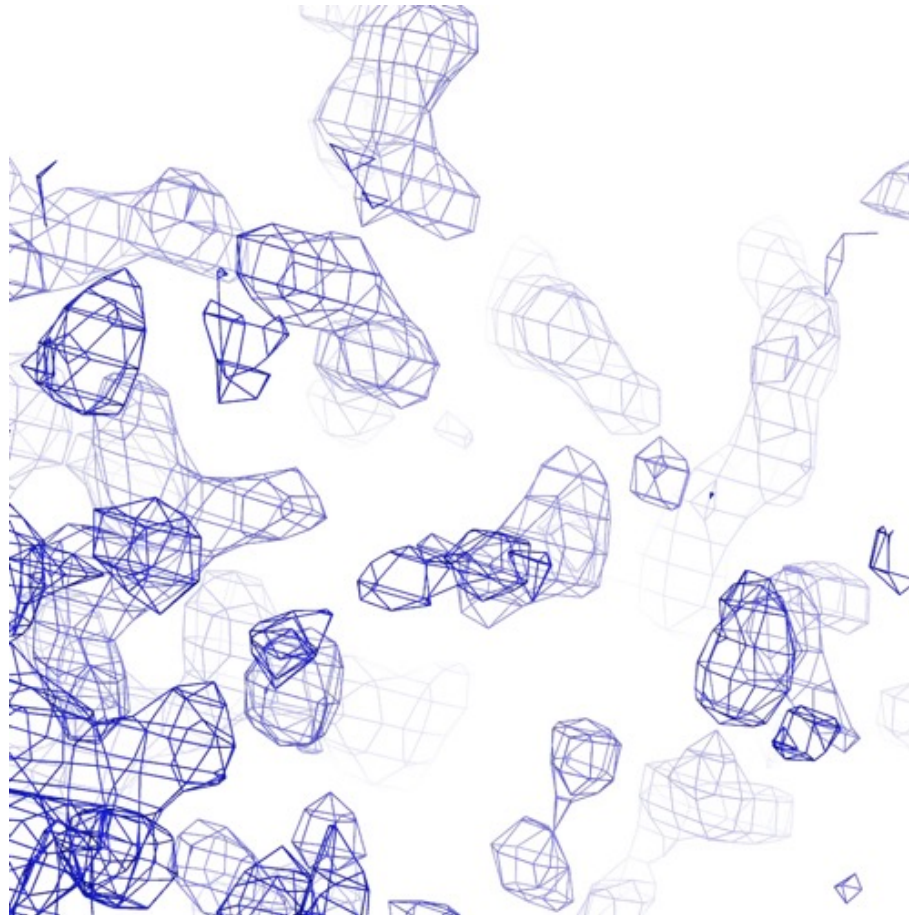
(New developments in model building and validation)

Grzegorz Chojnowski
EMBL Hamburg

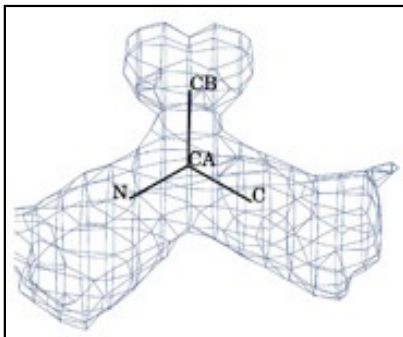
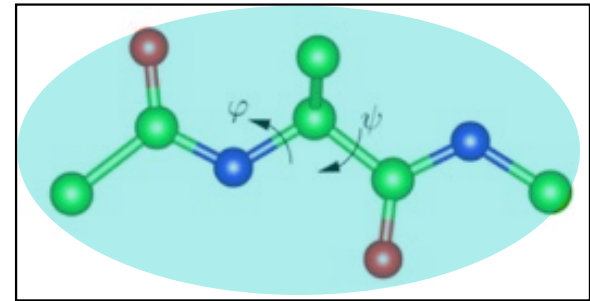
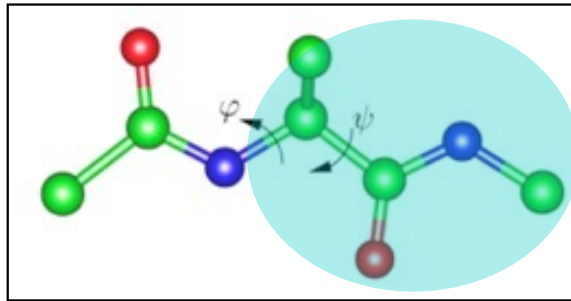
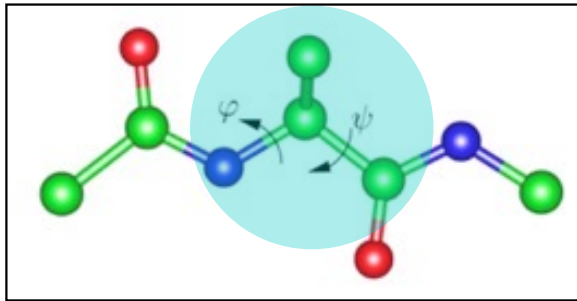
Diamond-CCP4 workshop 03.12.2022



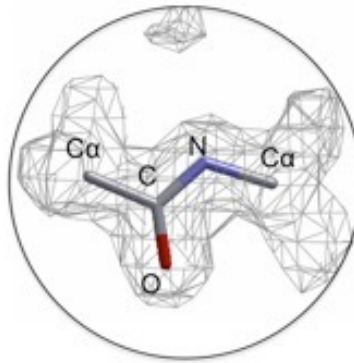
Electron density map interpretation



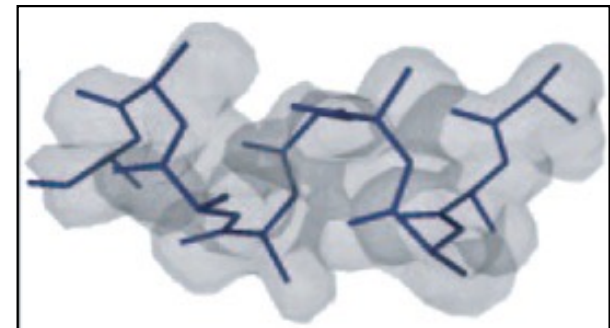
Crystal structure model building



Buccaneer
ModelCraft

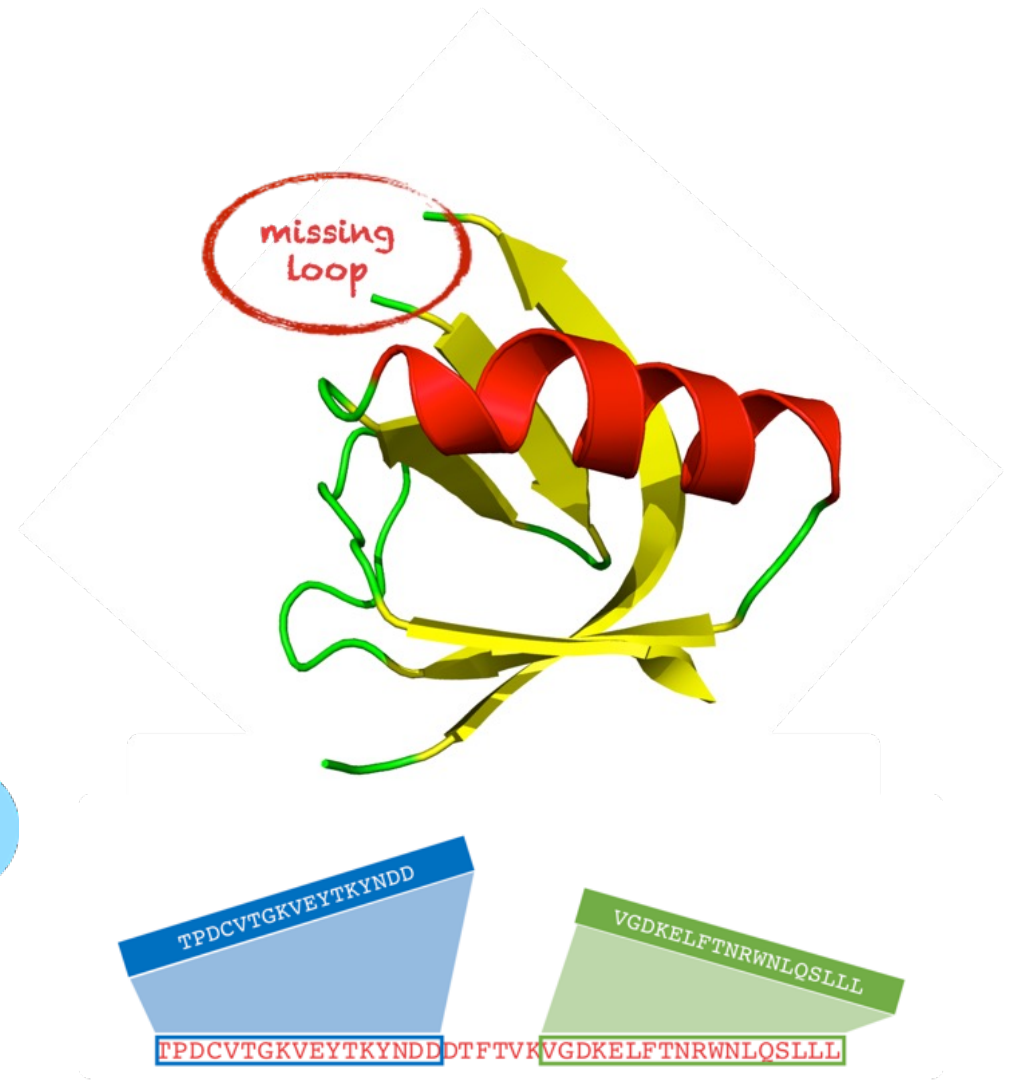
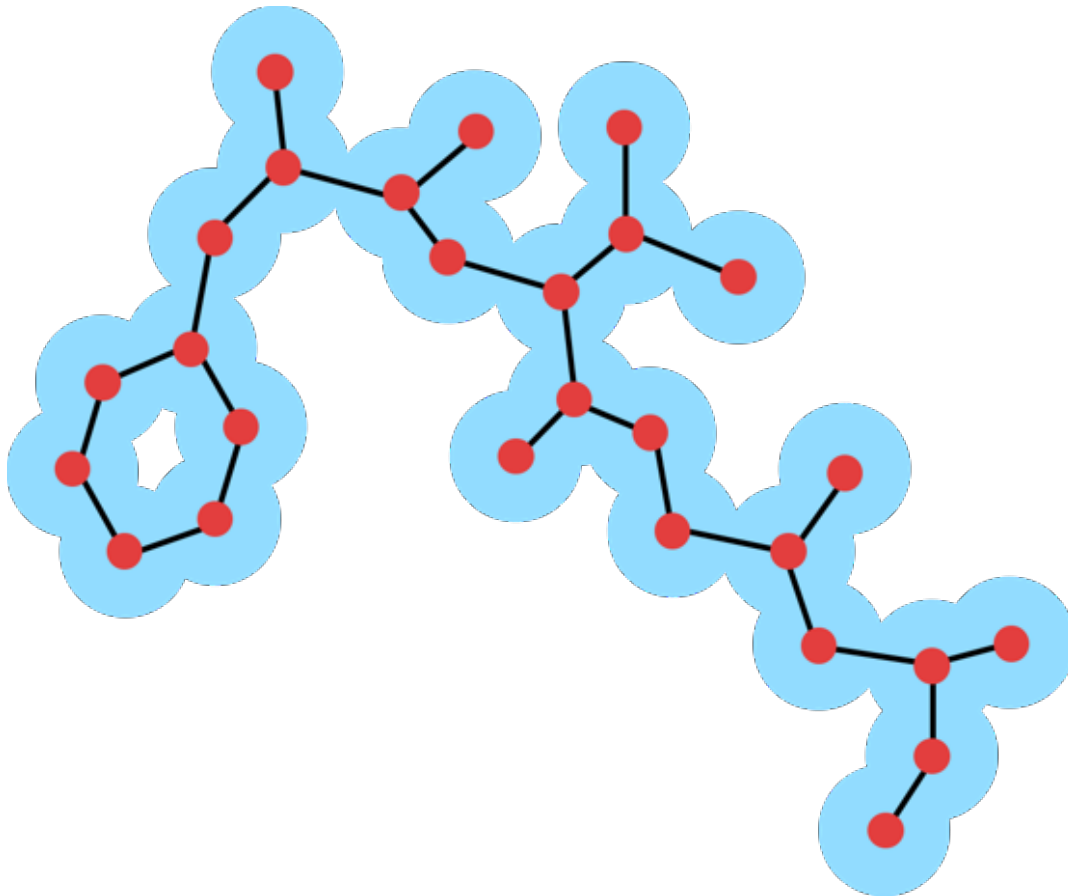


ARP/wARP

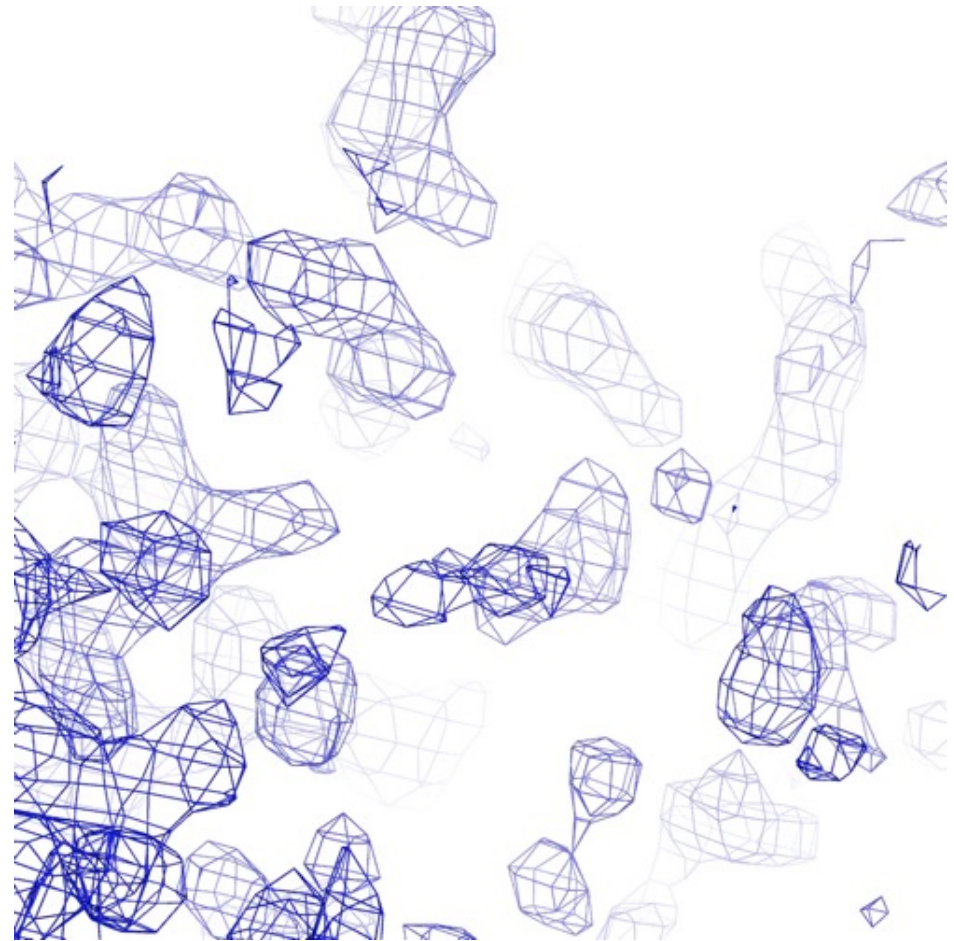
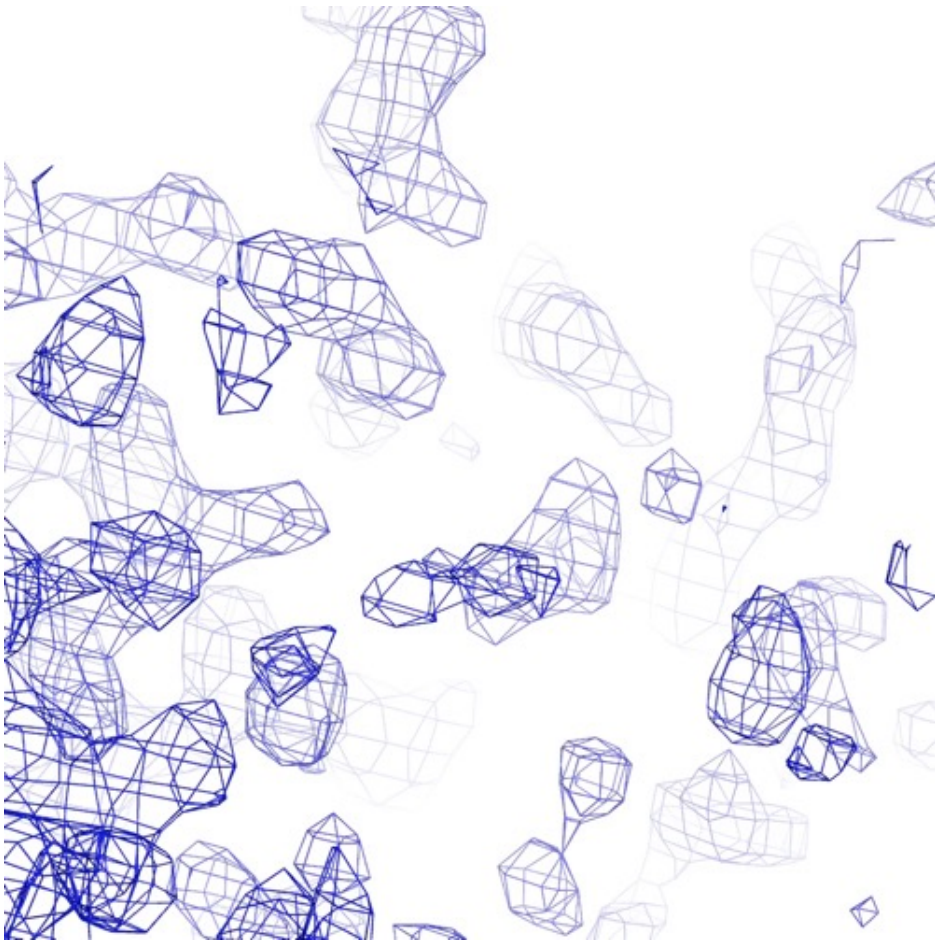


Phenix.autobuild

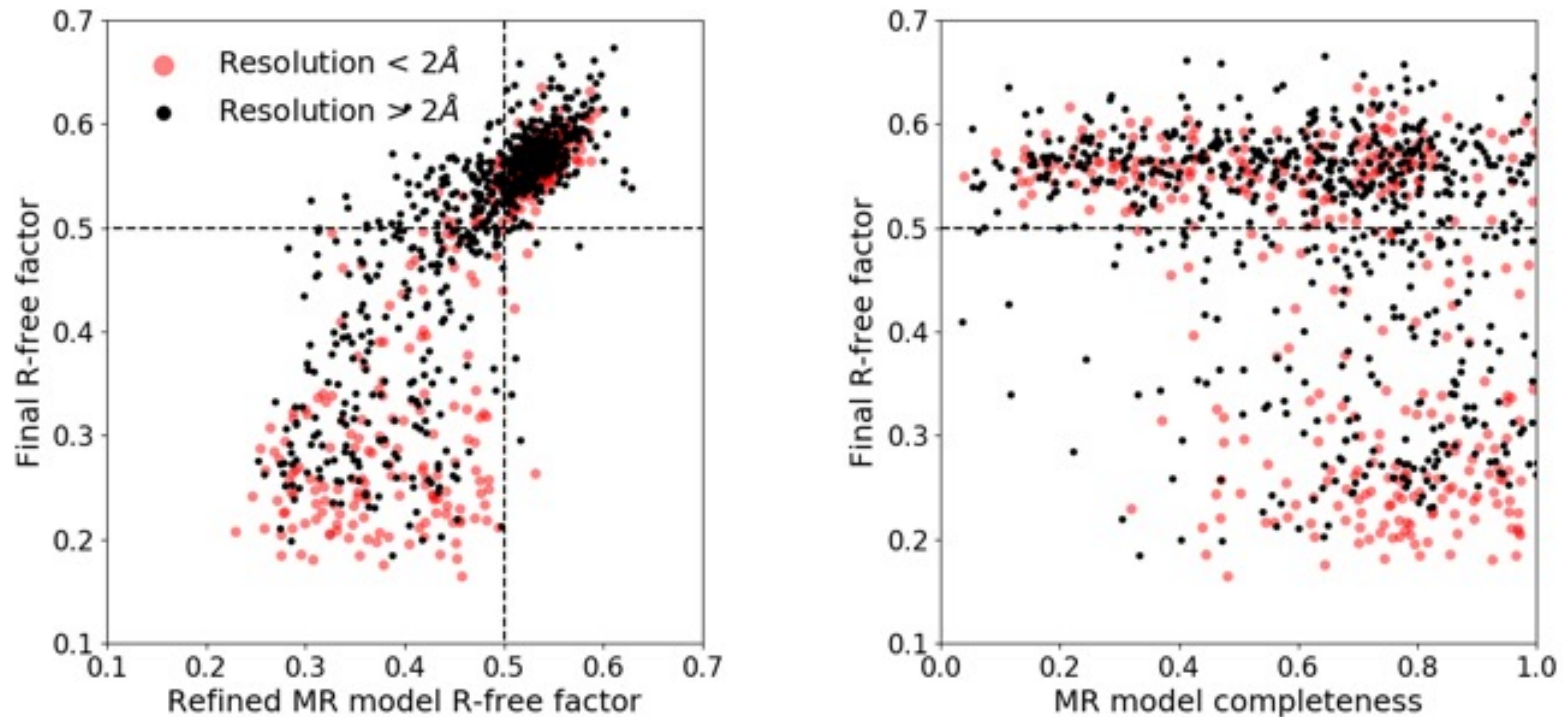
Crystal structure model building



Crystal structure model building

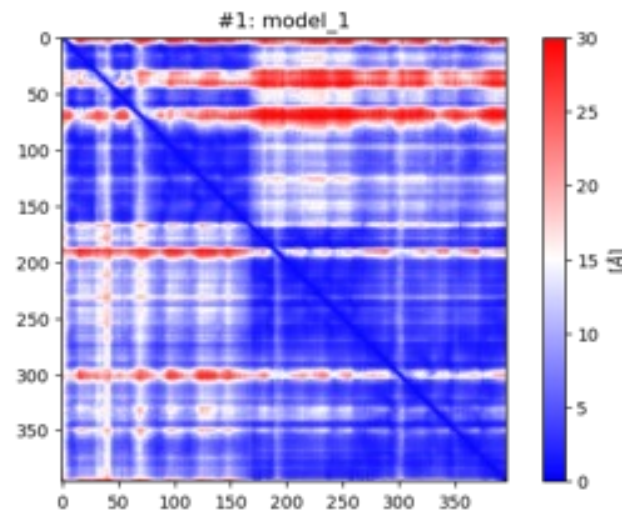
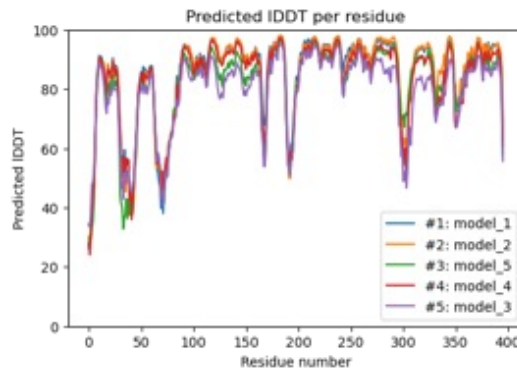


Crystal structure model building



1493 MR solutions submitted to the ARP/wARP web service from automated pipelines (BALBES, MrBUMP, MORDA)

Solving crystal structures with AlphaFold2 models



pTM=0.68 <pLDDT>=83

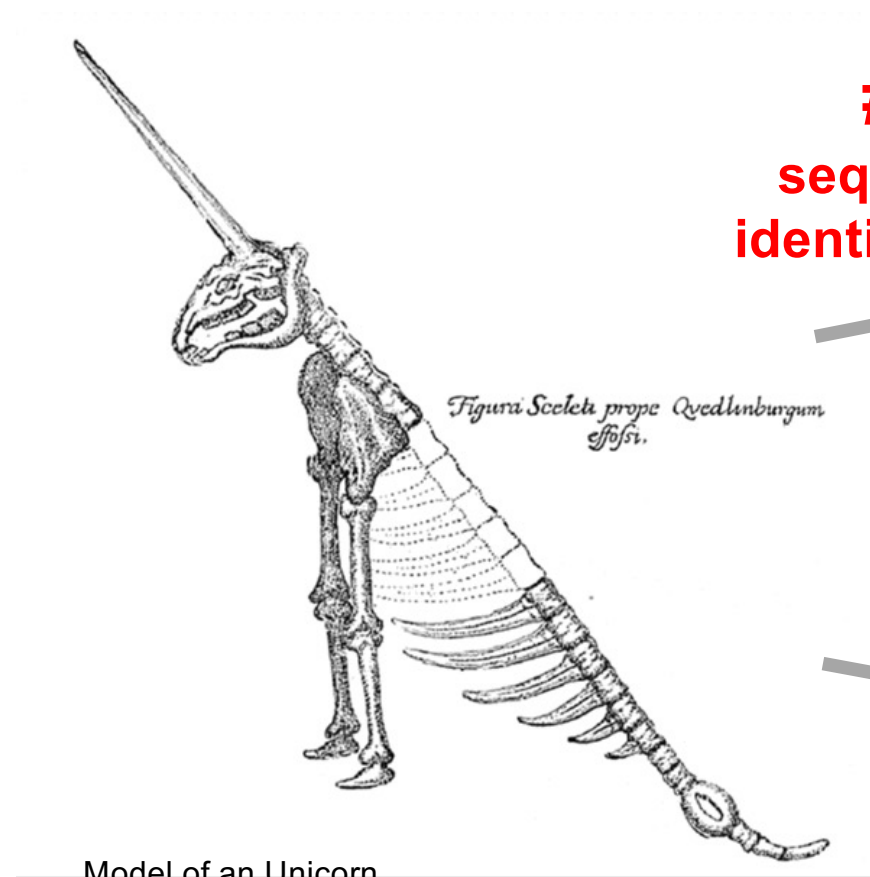


AlphaFold2
prediction



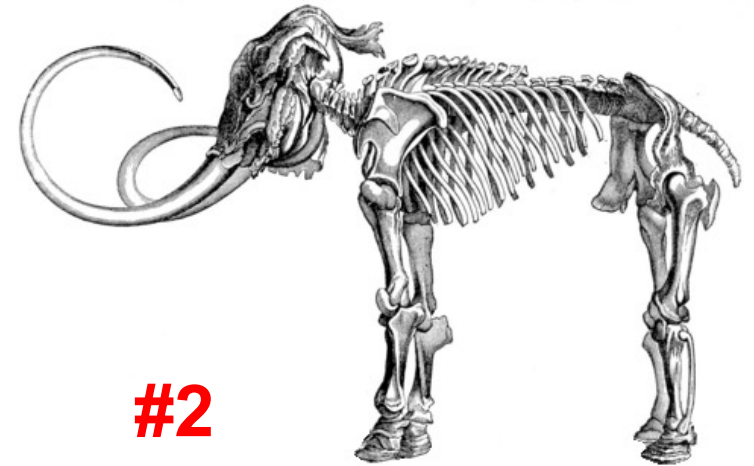
crystal structure
model

Model building traps in MX (and cryo-EM)

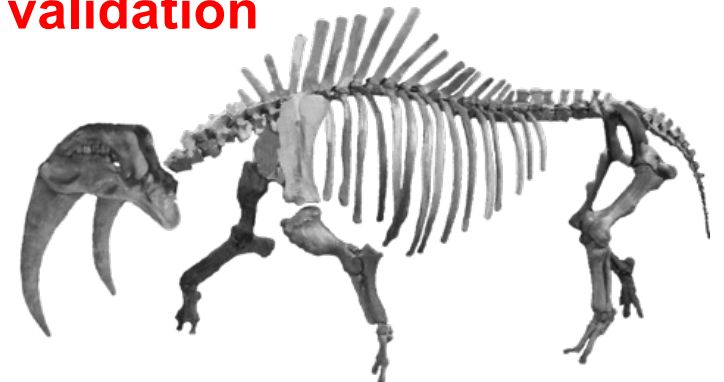


Model of an Unicorn
Gottfried Leibniz after Otto von
Guericke, *Protogaea* (1719)

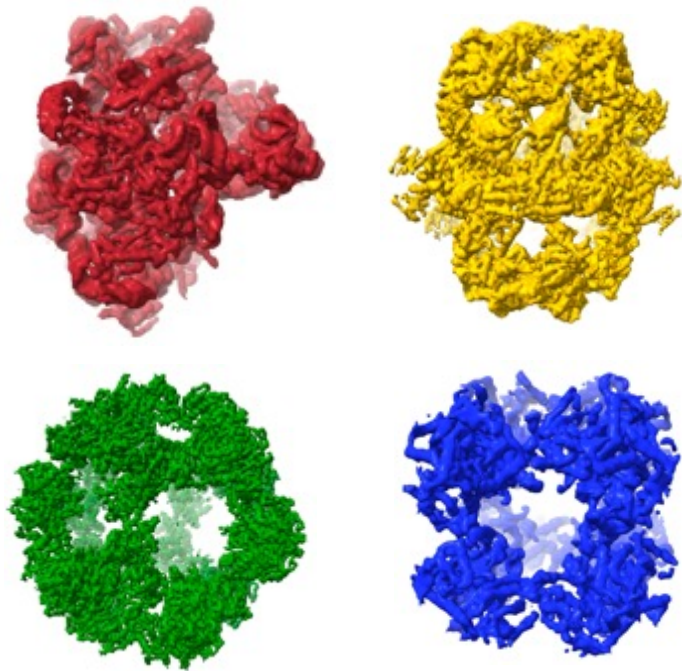
#1
sequence
identification



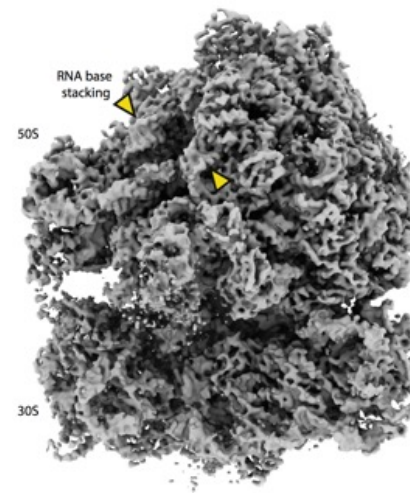
#2
sequence
validation



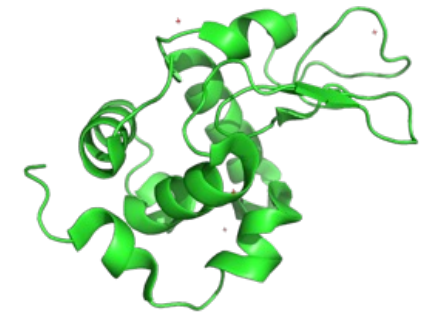
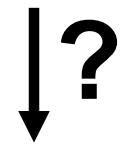
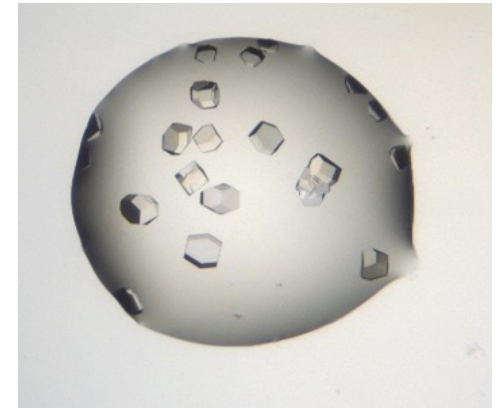
Unknown proteins in EM and MX



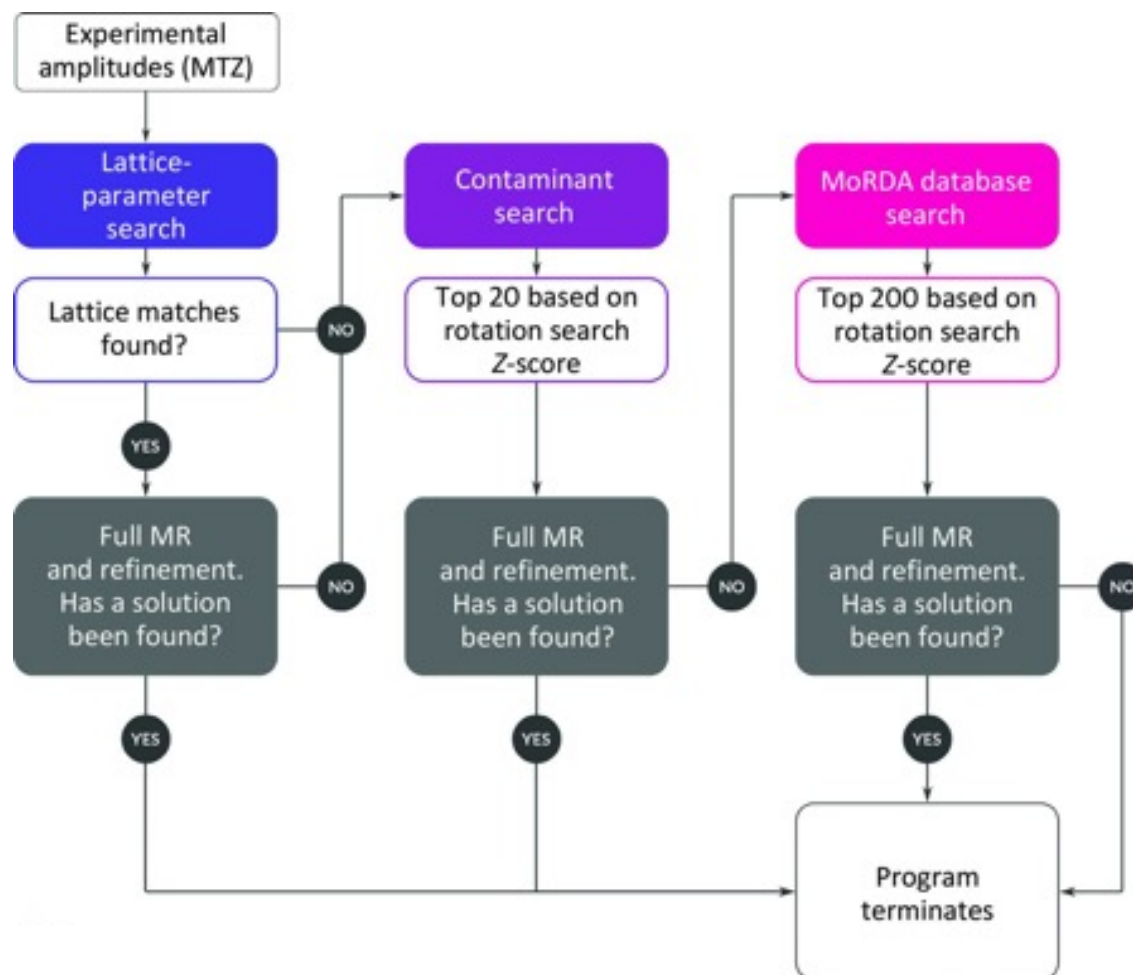
Cryo-EM and artificial intelligence visualize endogenous protein community members
Skalidis et al. Structure 2022



M. pneumoniae 70S ribosome at 3.5 Å
refined from in situ tilt-series data
Tegunov et al. 2021

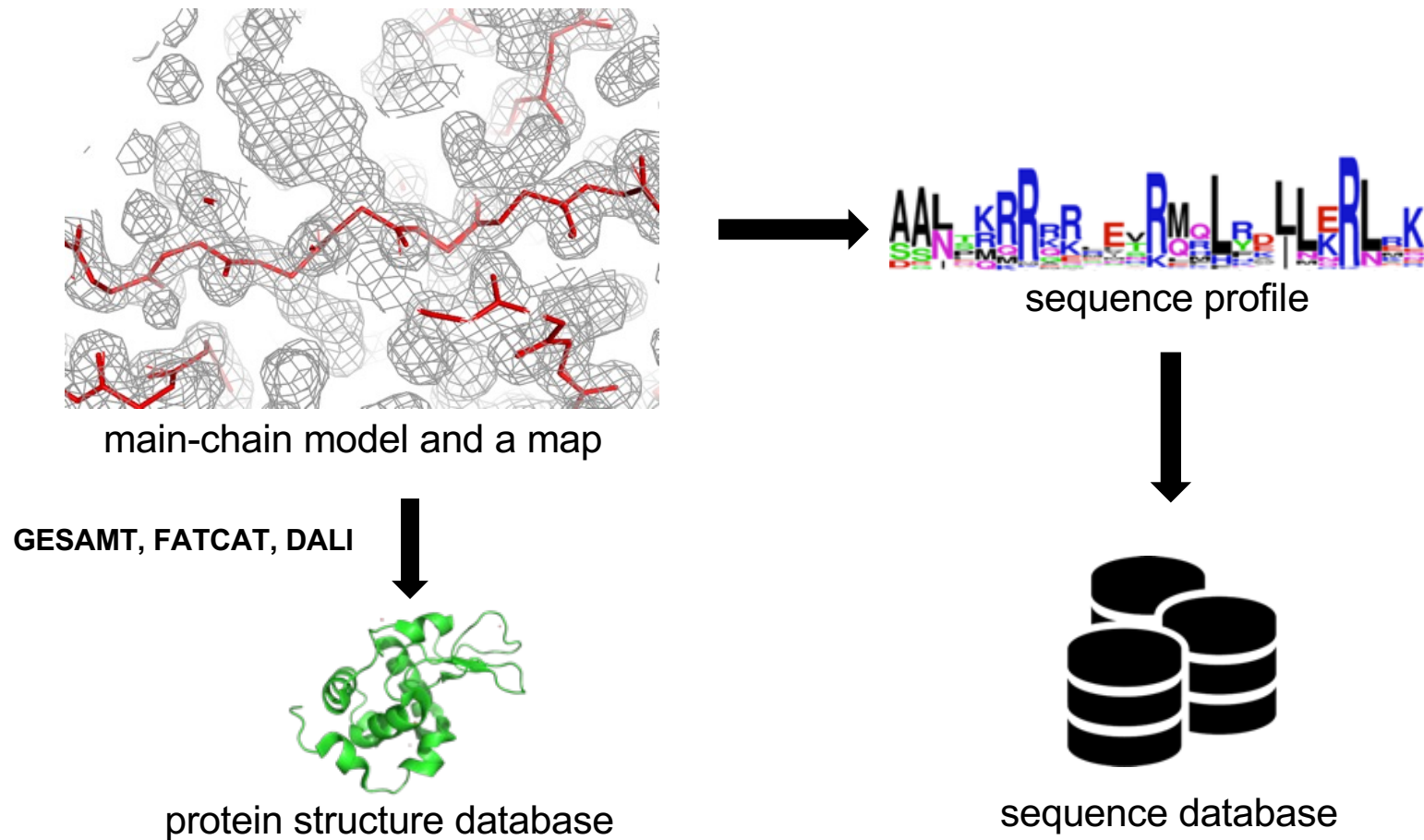


Sequence-free MR with SIMBAD

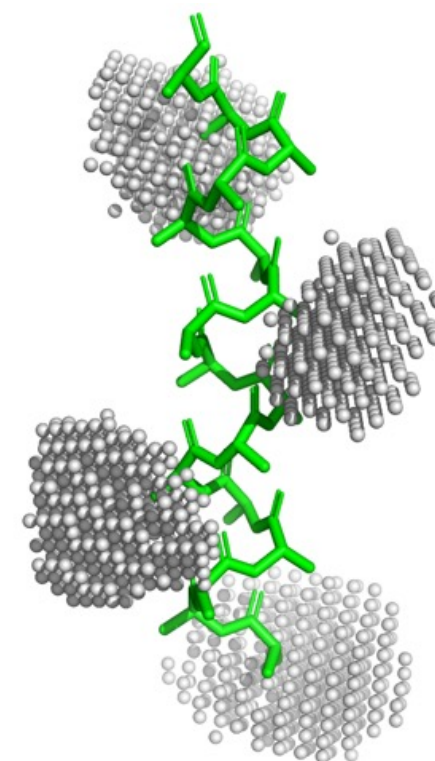
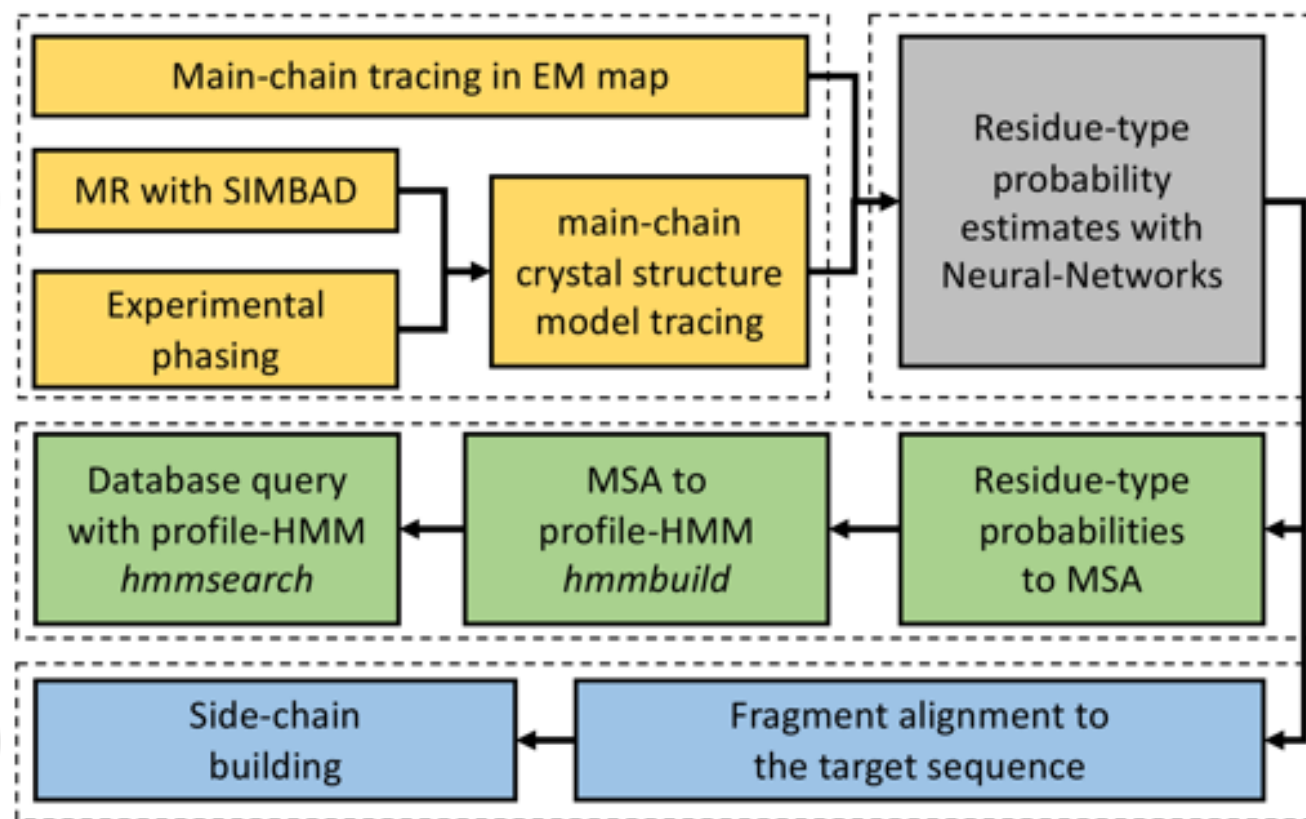


Available in


Protein sequence identification from a map



Protein sequence identification with findMySequence

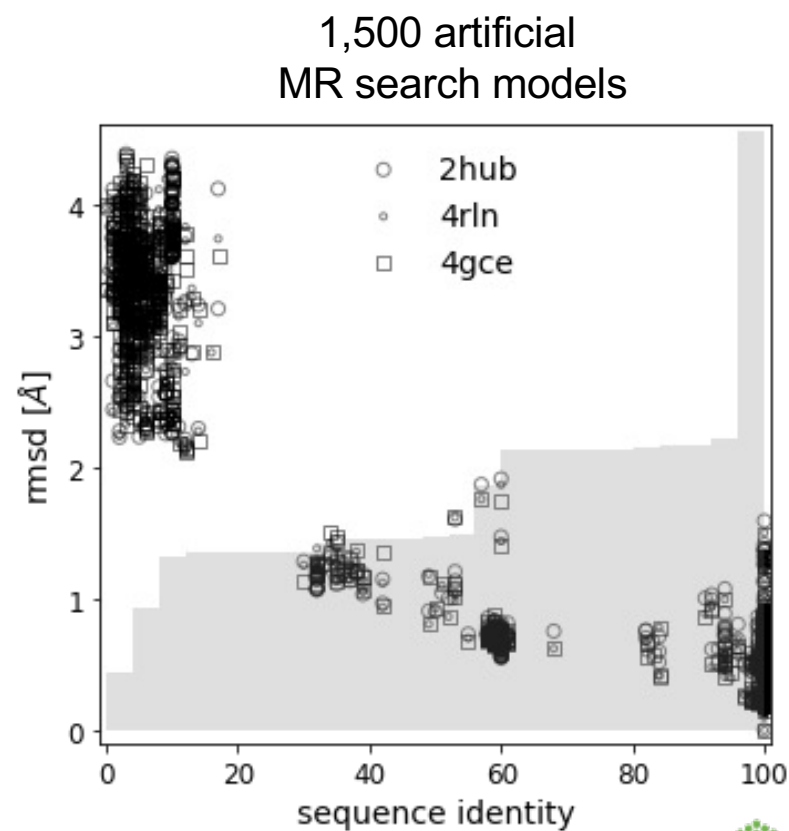
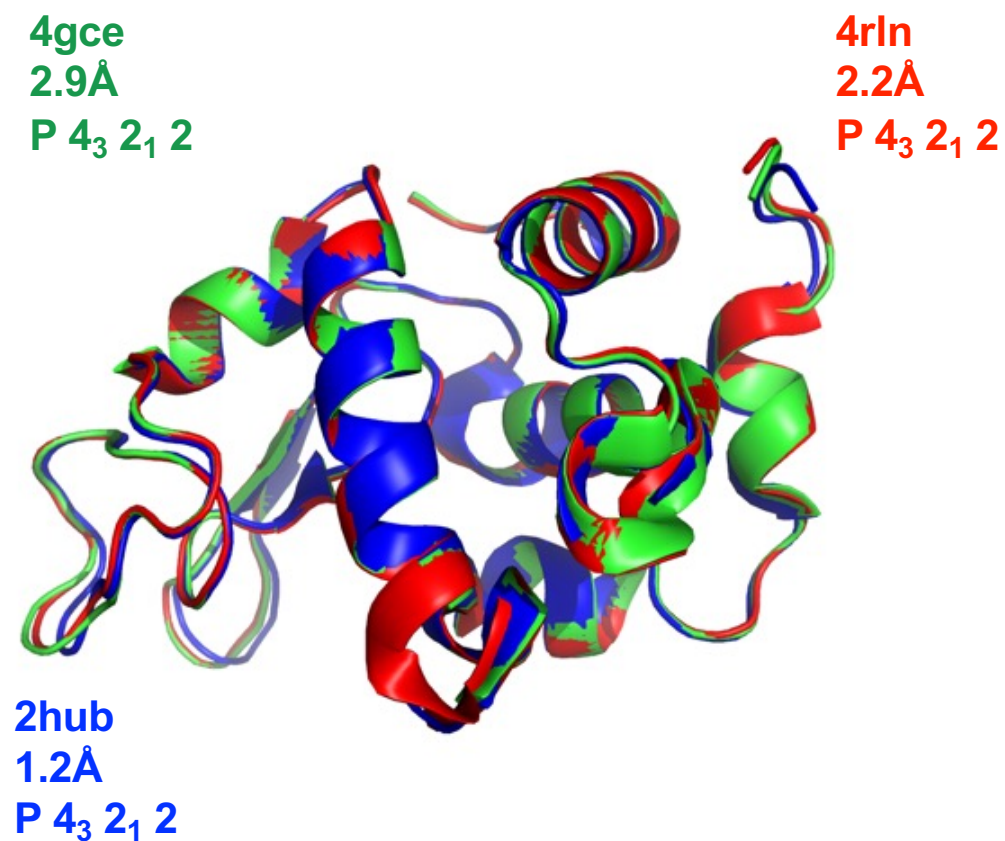


residue-type NN classifier

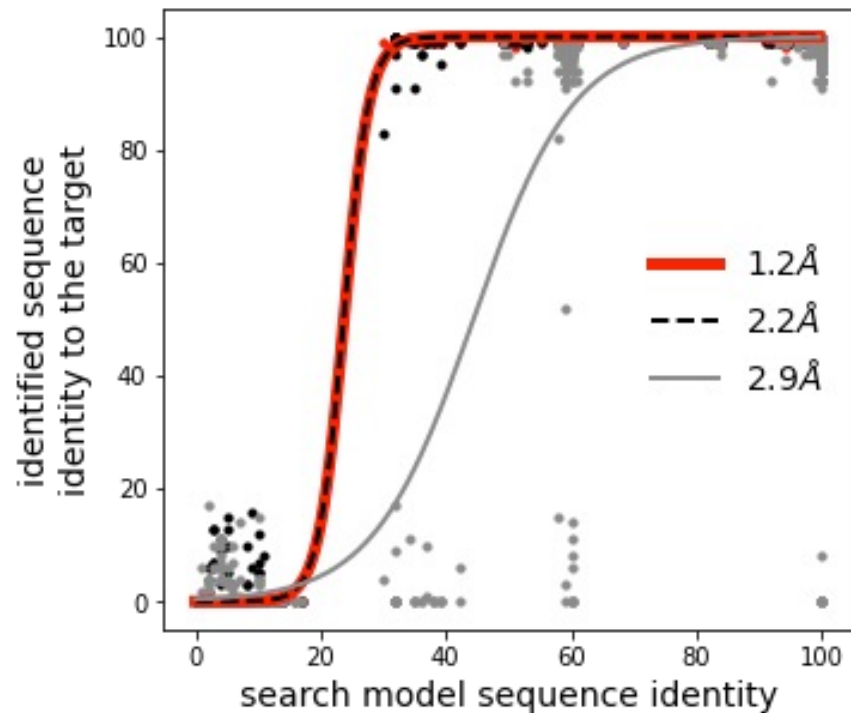
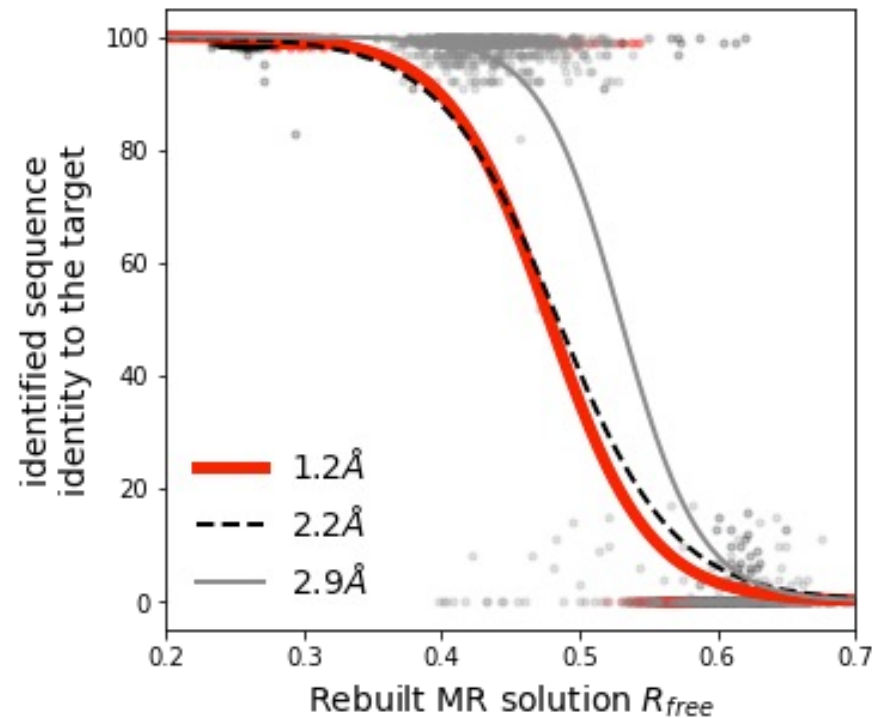
328 input features

111,800 parameters

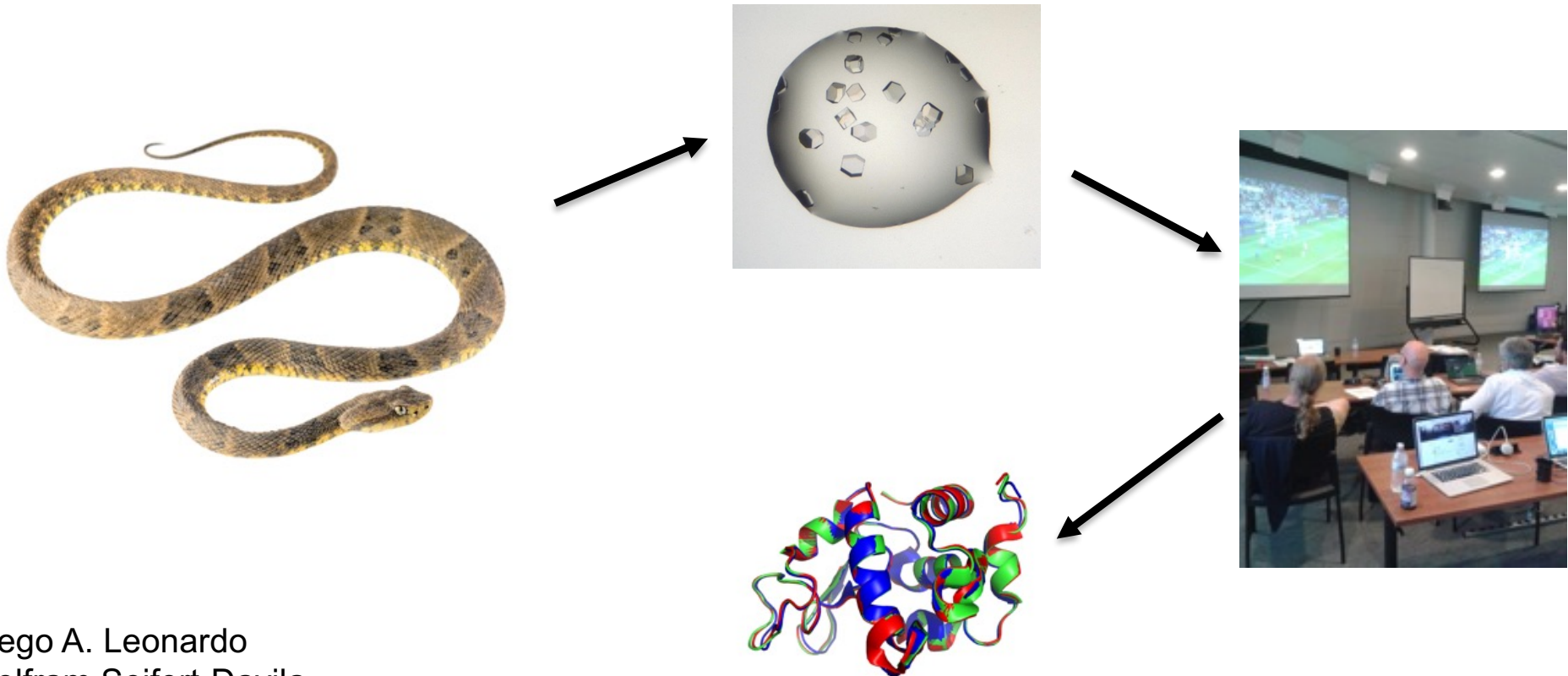
MX benchmarks: three hen egg-white lysozyme targets



MX benchmarks: three hen egg-white lysozyme targets



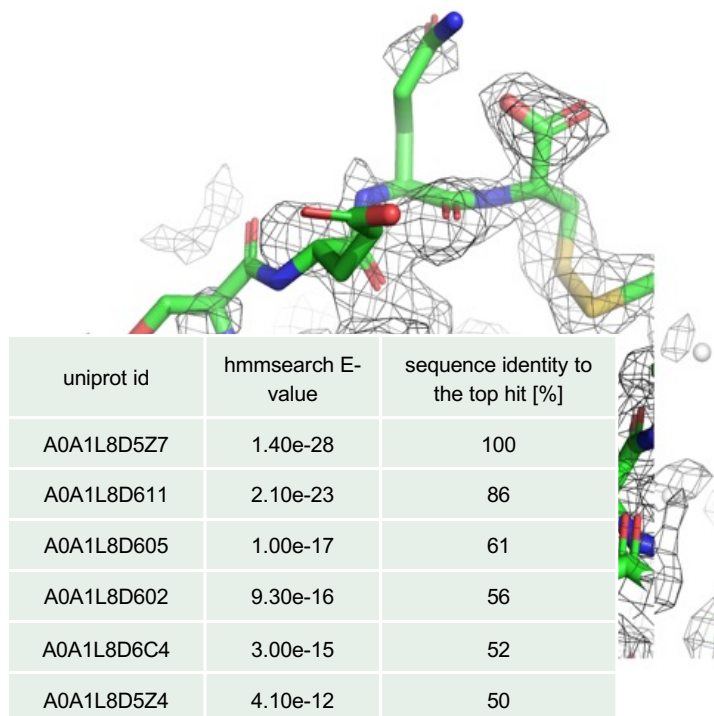
MX use-case: *Bothrops atrox* venom proteins



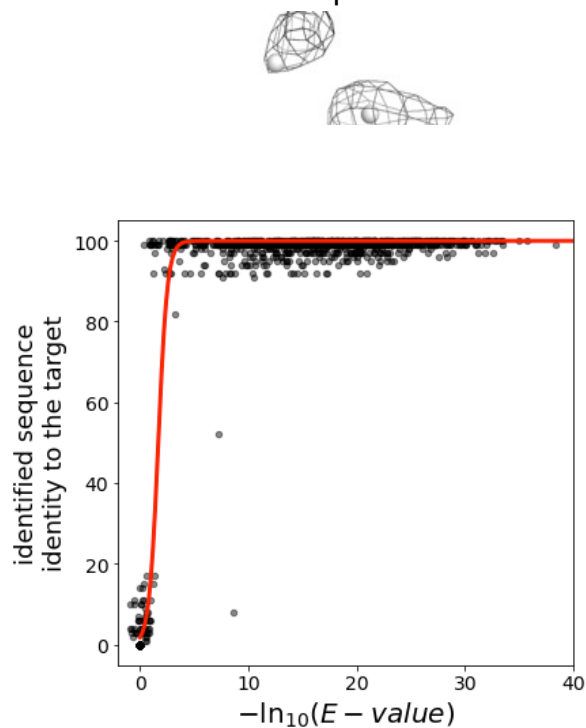
Diego A. Leonardo
Wolfram Seifert-Davila
Dan E. Vivas-Ruiz

MX use-case: *Bothrops atrox* venom proteins

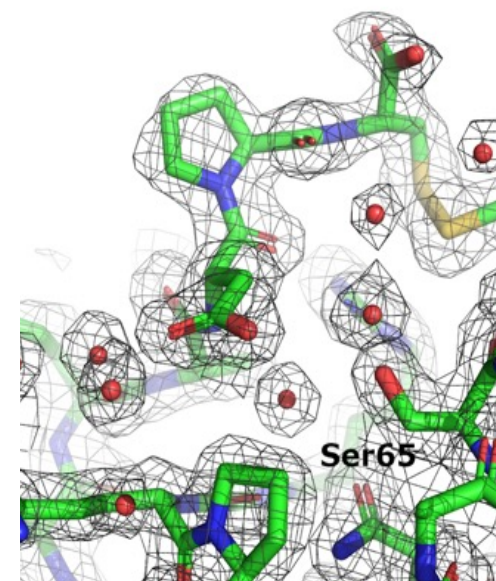
initial MR solution



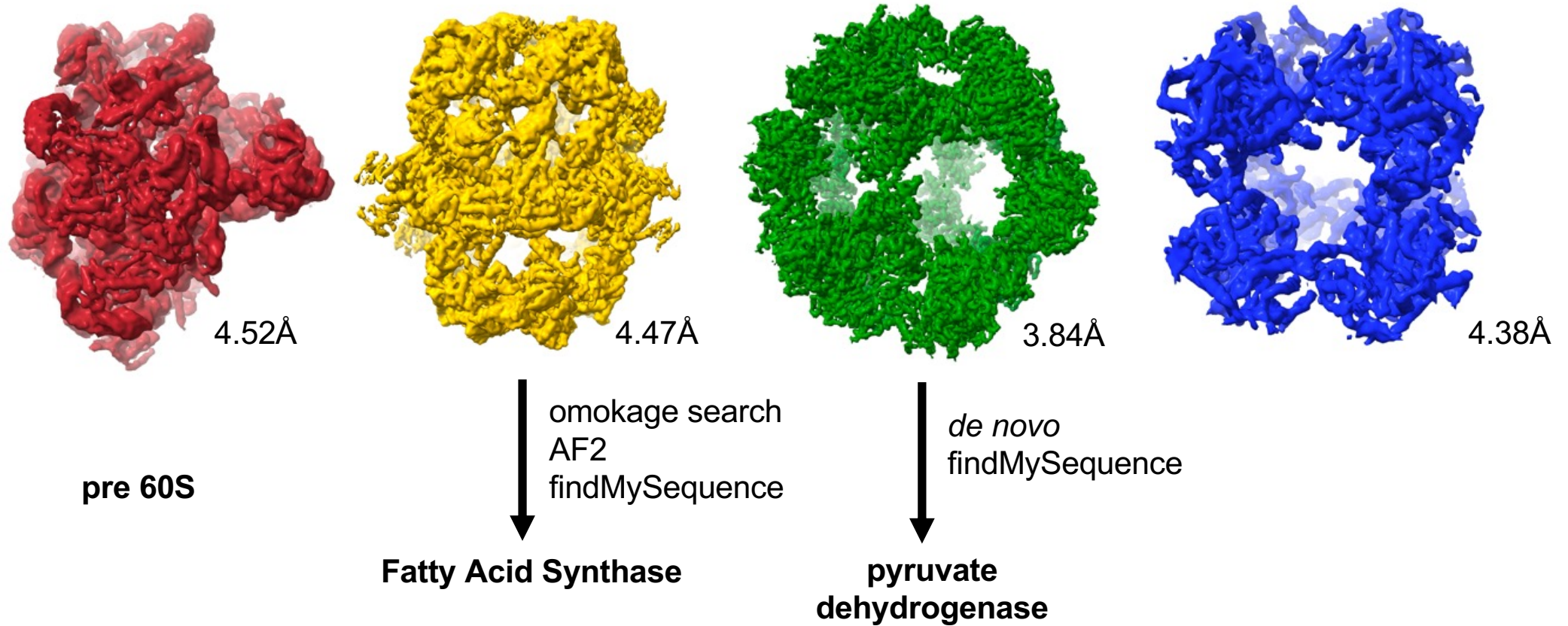
model rebuilt
wout sequence



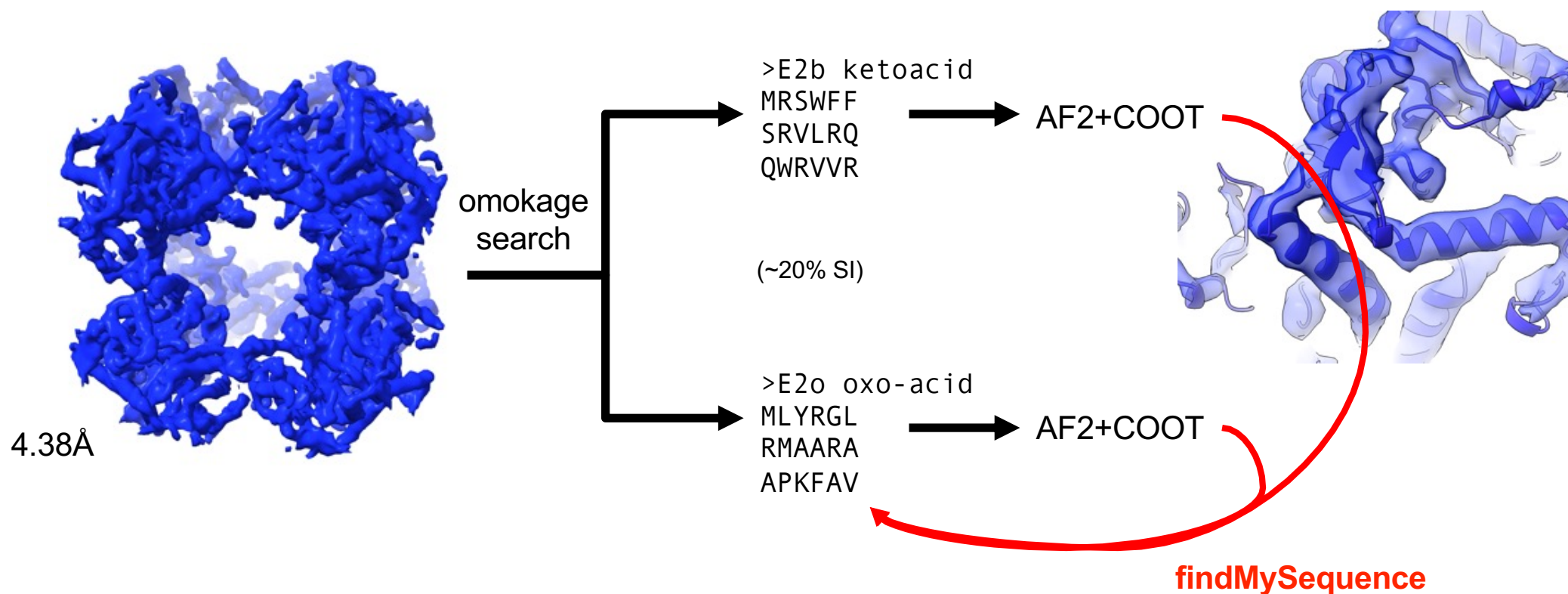
final model



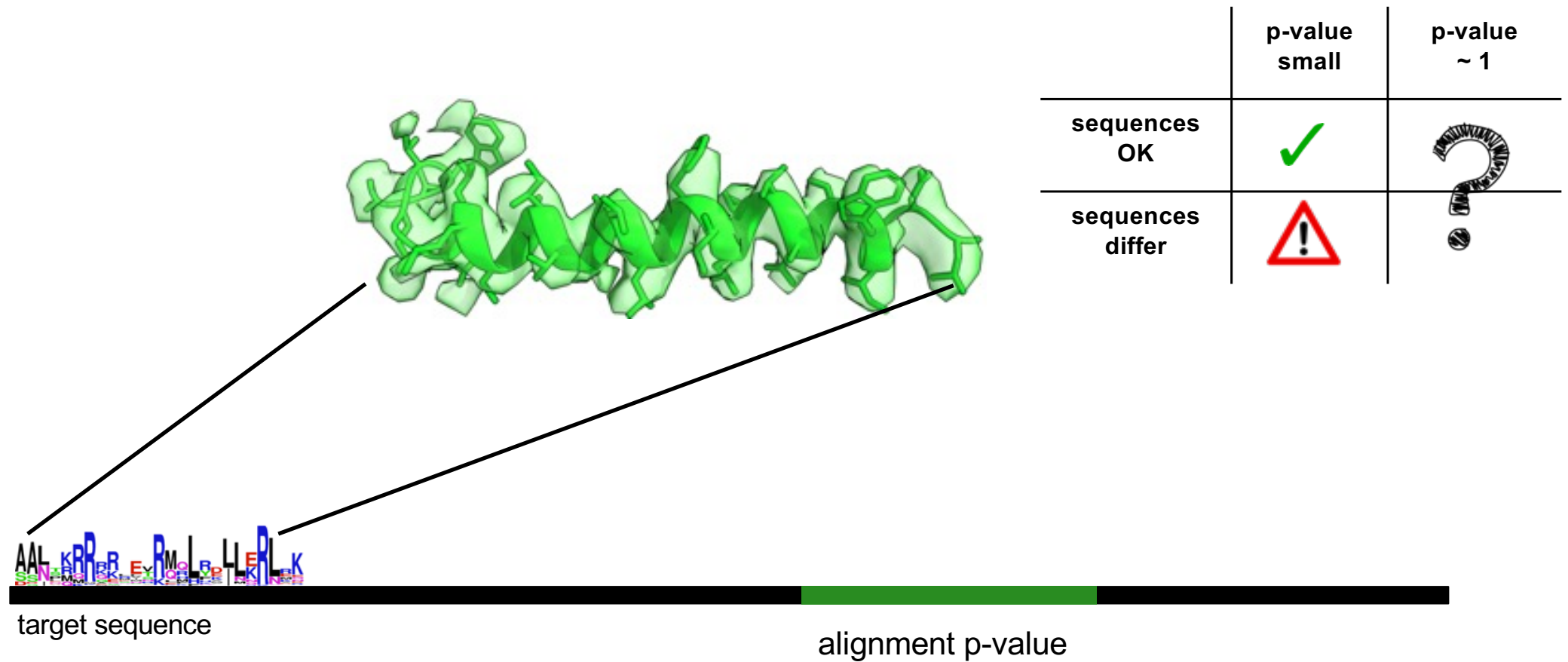
cryo-EM of *C. thermophilum* native cell extracts



two dehydrogenases in *C. thermophilum* native cell extracts

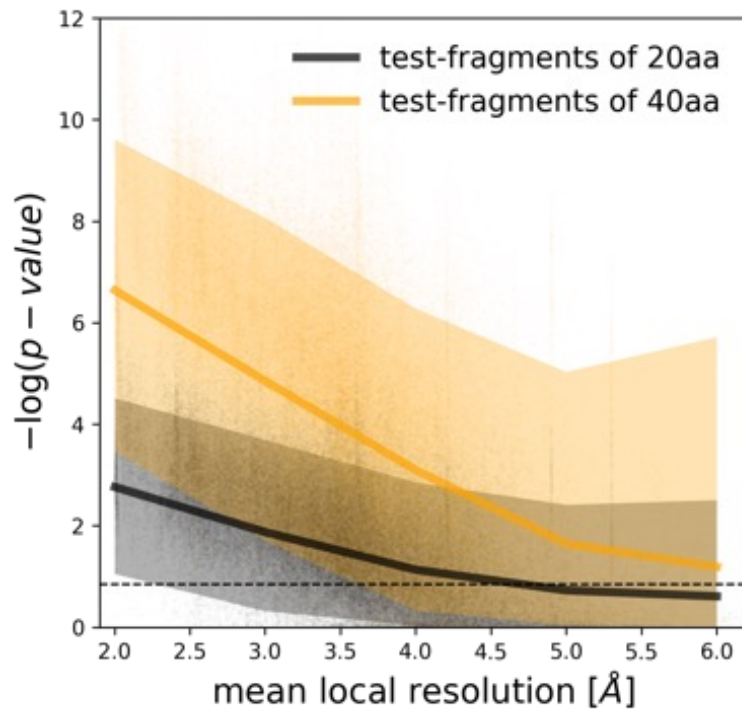


Sequence re-assignment as a model validation

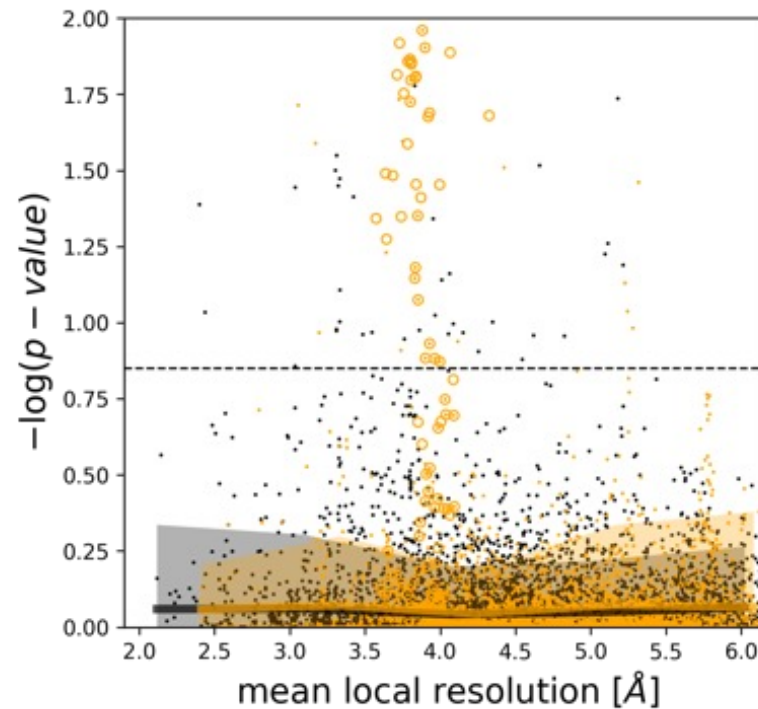


Is the sequence assignment p-value a reliable score?

sequence OK



sequences differ

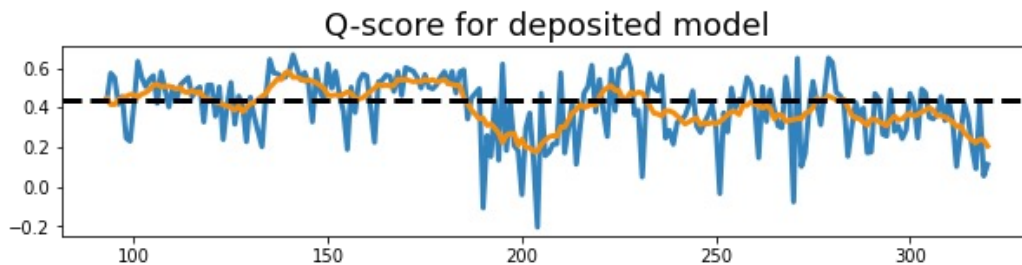
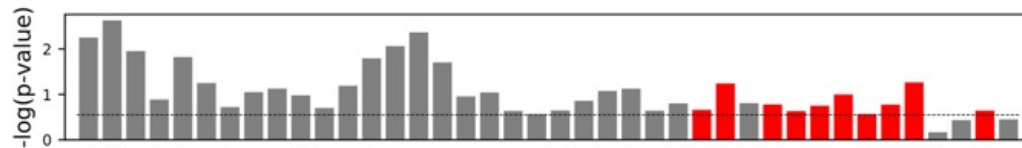
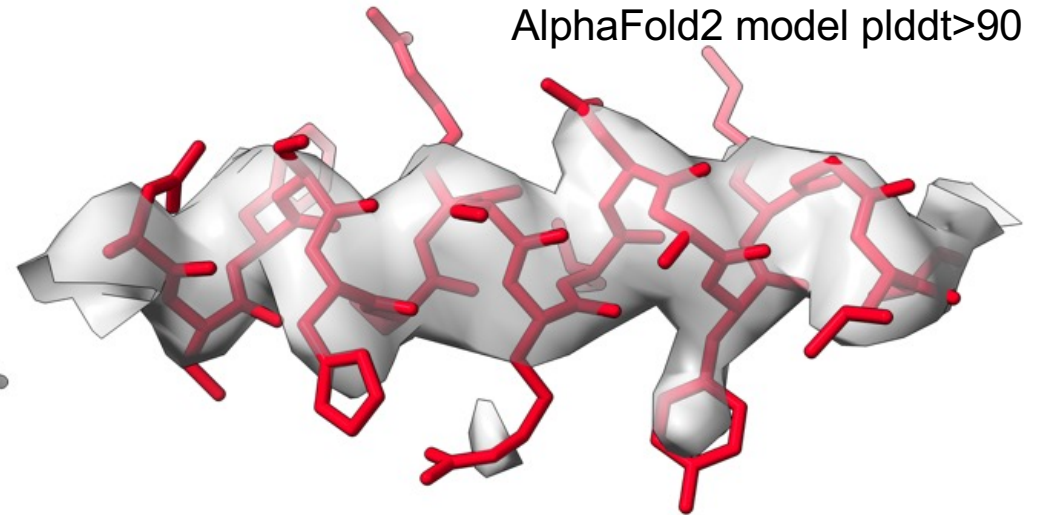
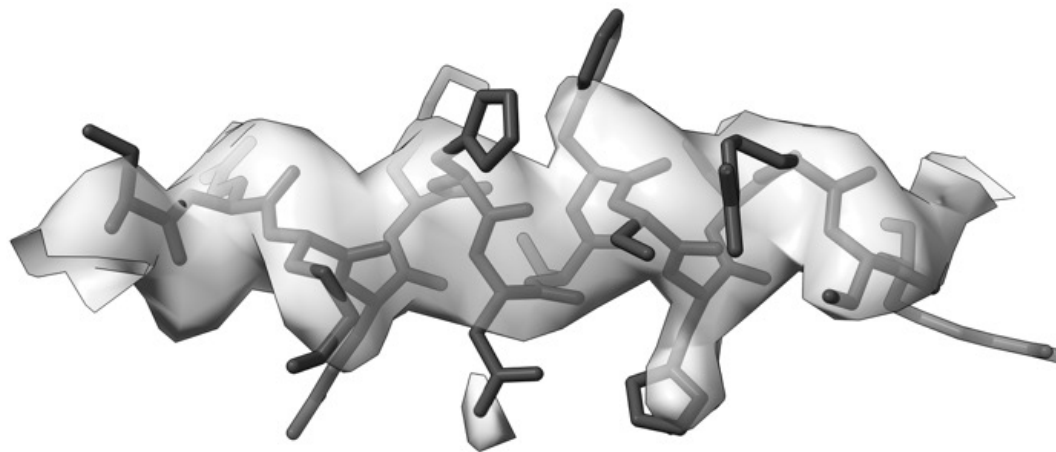


	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?

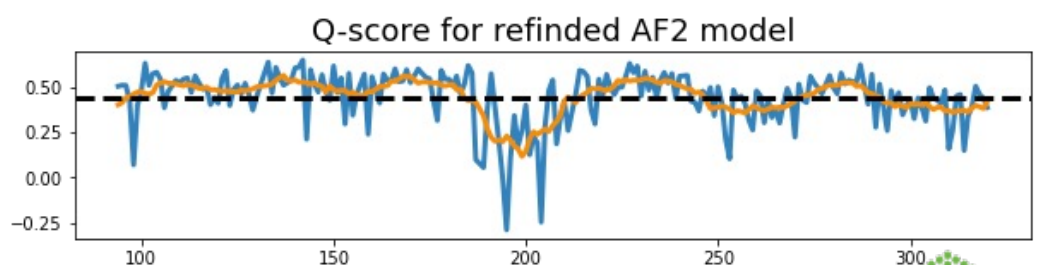
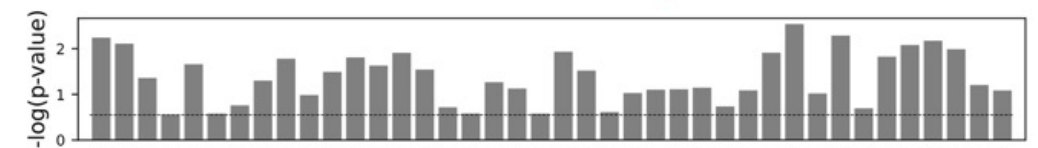
p-values for 30k protein chain fragments re-assigned to target sequence

checkMySequence: finding a better hypothesis

cytoplasmic domain of a cation channel at 3.8Å resolution



Chojnowski ActaD 2022



checkMySequence: complete sequence assignment validation

```
*****
***** SUMMARY *****
*****
```

==> Unidentified chains; check input sequences and model-to-map fit

```
e/2:51
g/3:39
```

==> Chains with sequence mismatches; you will have to fix them first!

```
model      KDNVVQMMNEKKSFDVSDFPKVYLTTAVEEDLDT--
           |||
refseq      KDNVVQMMNEKKSFDVSDFPKVYLTTTVEEDLDTRG
```

==> Possible sequence assignment issues

- Fragment N/5-24 has a **chain break at N/11 and no residue indexing gap**, check!

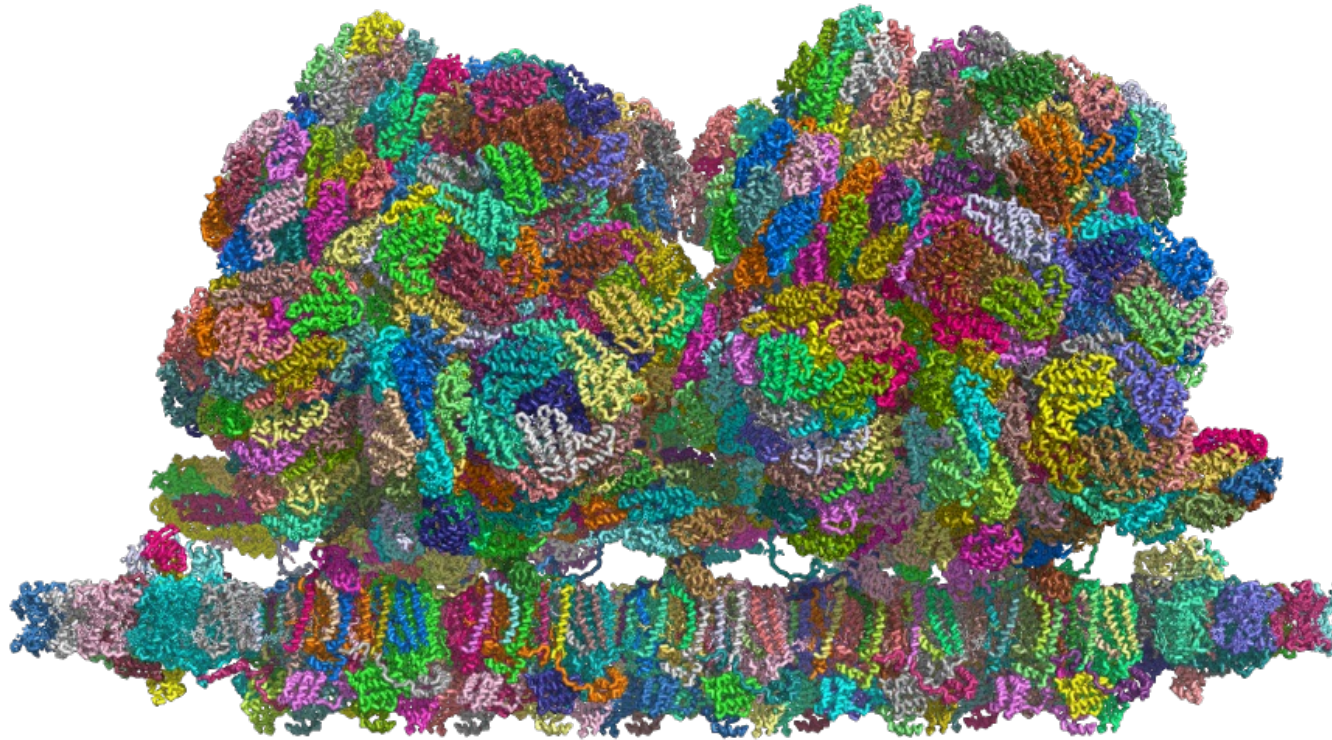
```
model      -----VELTEEE-LYISKKNLLFKRF-----
           |||
refseq      AAAAAAMPRVELTEEEKLYISKKNLLFKRFVEPGRLCLIE
```

- Fragment F/356-395 is shifted by -4 residues [-log(p-value)=1.99]

```
model seq 356-395
sknkkeKRVQKQIQKKELQKINHYYKGVAKAVKKKKKREEKKAkskktanqavi
new seq 360-399
sknkkekrvQKQIQKKELQKINHYYKGVAKAVKKKKKREEKKAkskktanqavi
```

When you think you've seen it all...

In situ **double**-PBS-PSII-PSI-LHCs megacomplex at 4.3Å



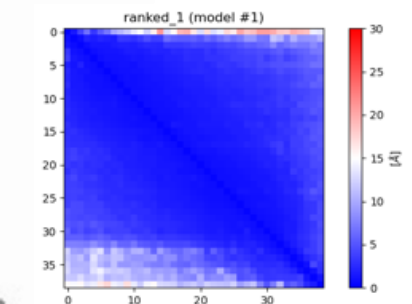
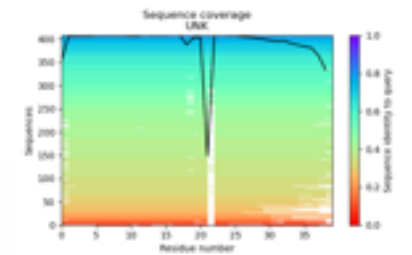
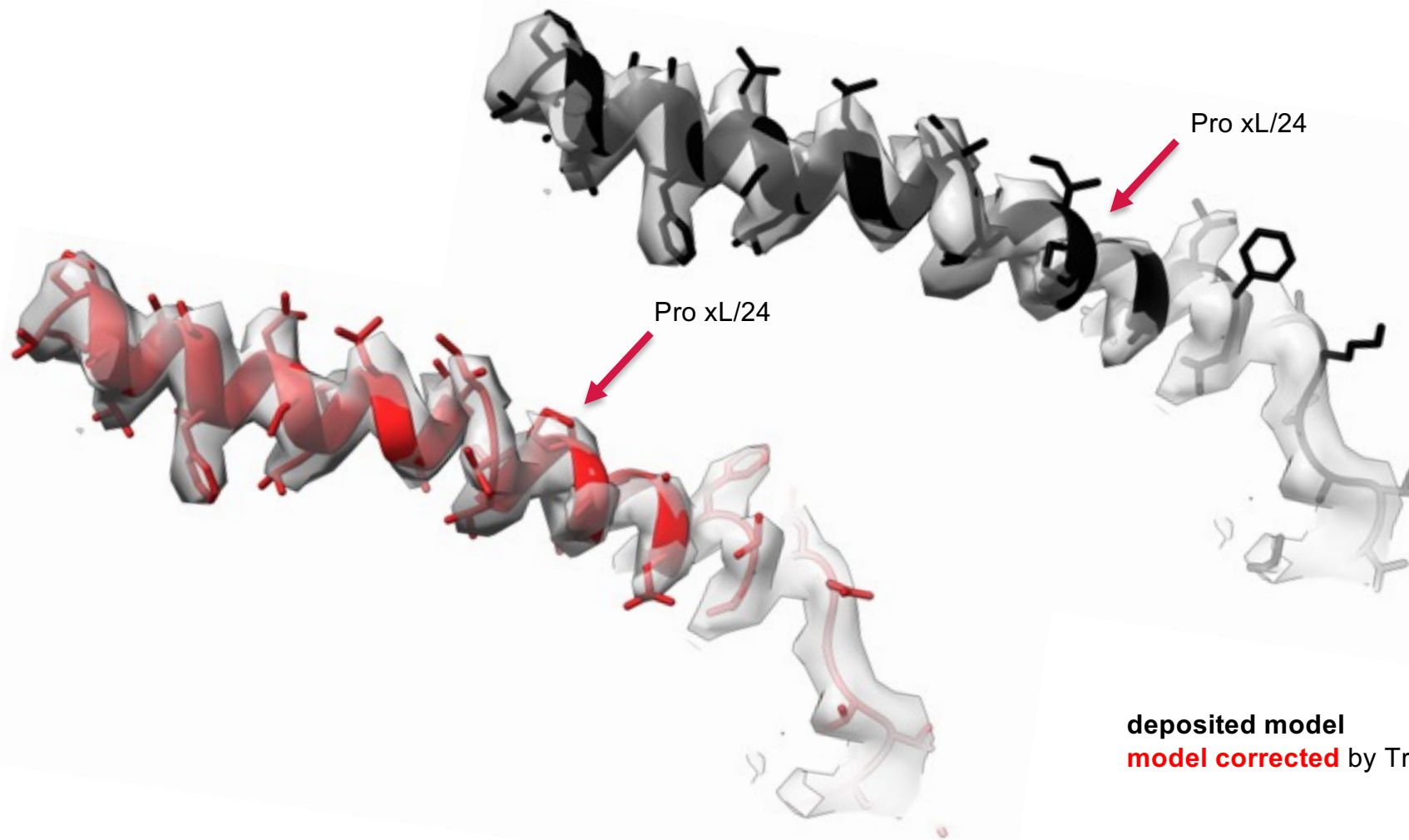
37MDa/305,000 residues/1,792 aa chains

7y5e/EMD-33618

You, X., Zhang, X., Cheng, J., Xiao, Y.N., Sui, S.F. *Nature* 2023

In situ photosystem megacomplex

Photosystem II reaction center X protein



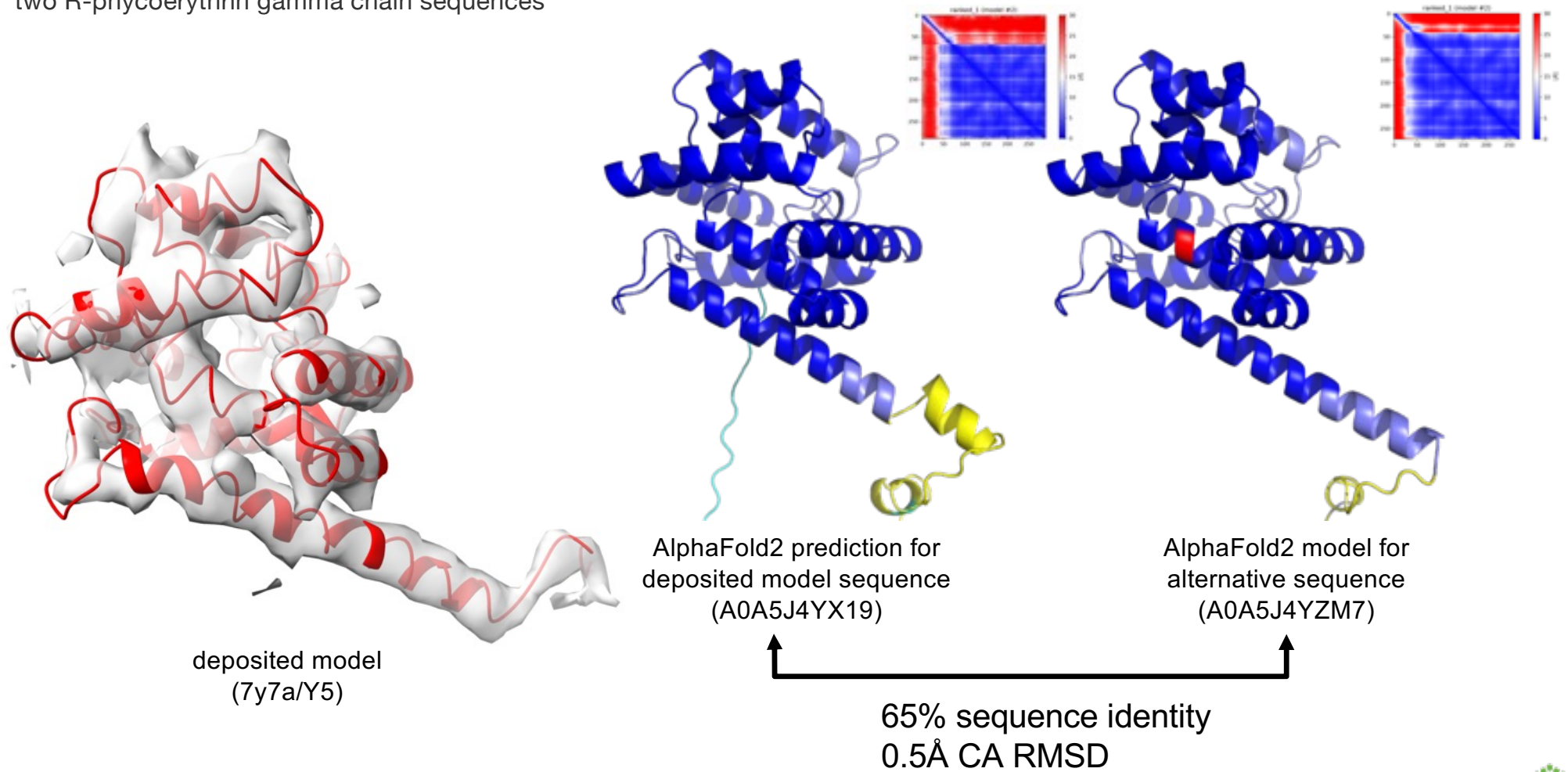
pTM=0.57
<pLDDT>=96

deposited model
model corrected by Tristan Croll



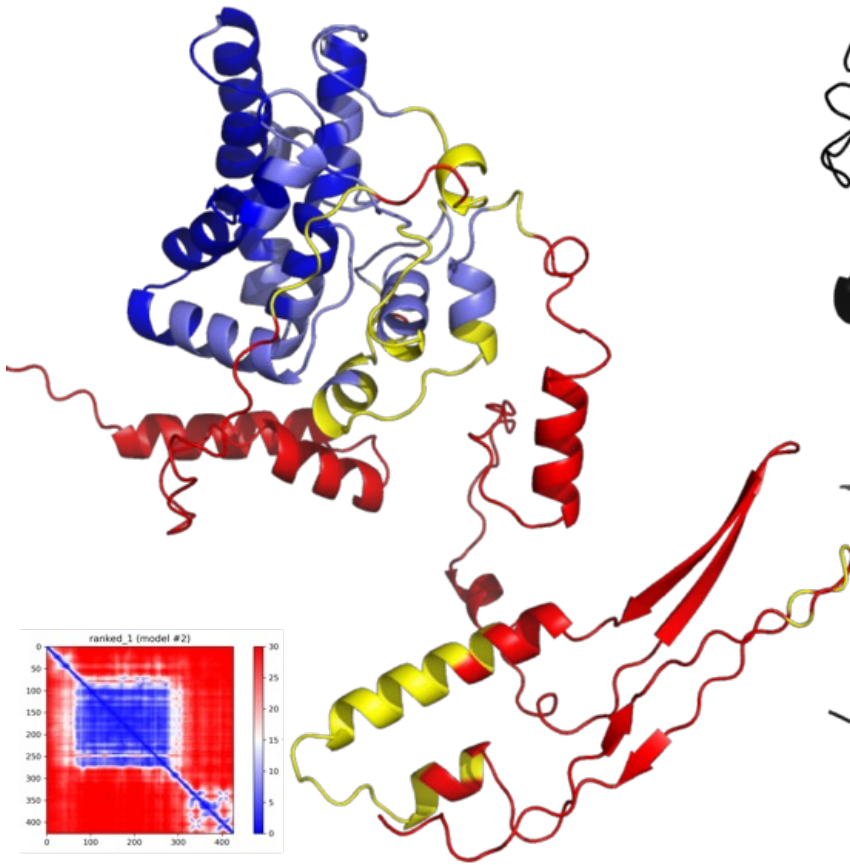
In situ photosystem megacomplex

two R-phycoerythrin gamma chain sequences

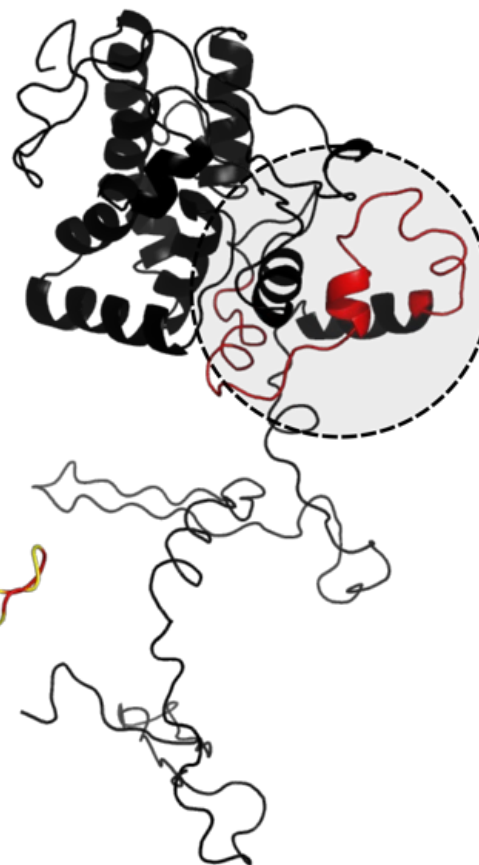


In situ photosystem megacomplex

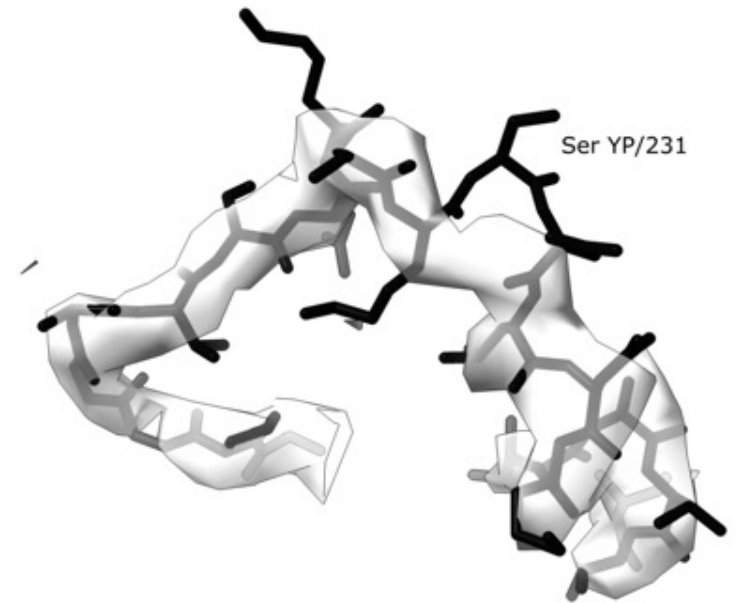
Phycobilisome linker polypeptide



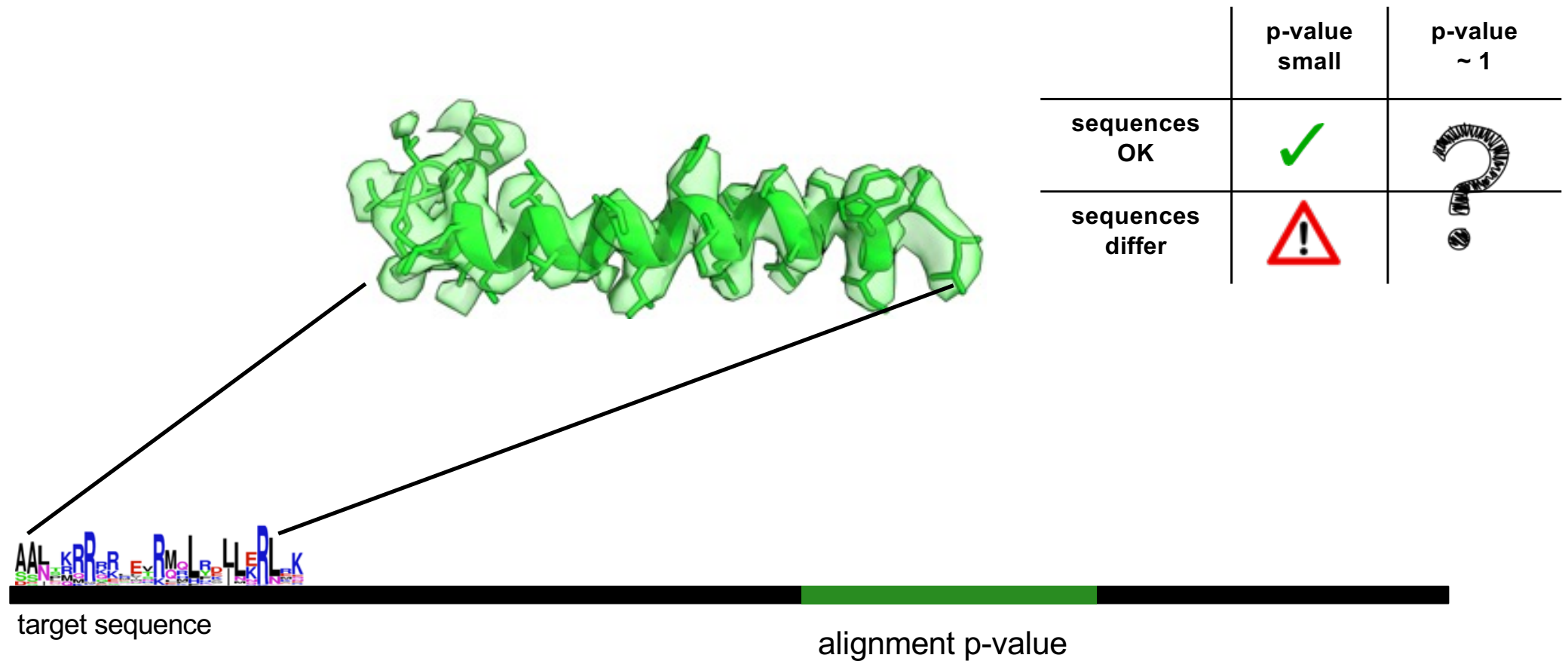
AlphaFold2 model



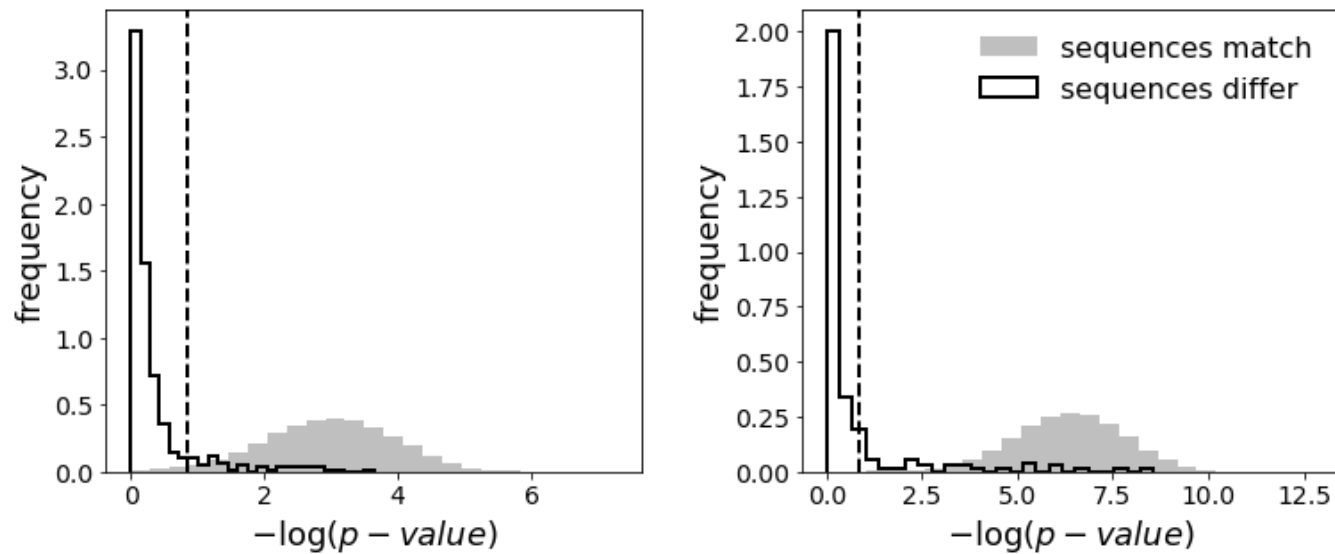
deposited model



checkMySequence and protein crystal structures – model bias



Sequence assignment in MX – model bias



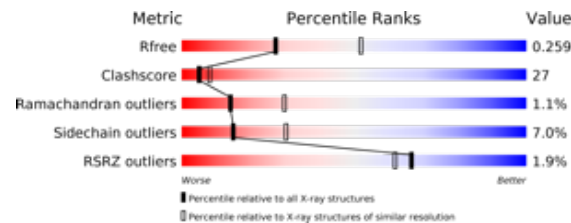
	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?

p-values for 30k protein chain fragments re-assigned to target sequence

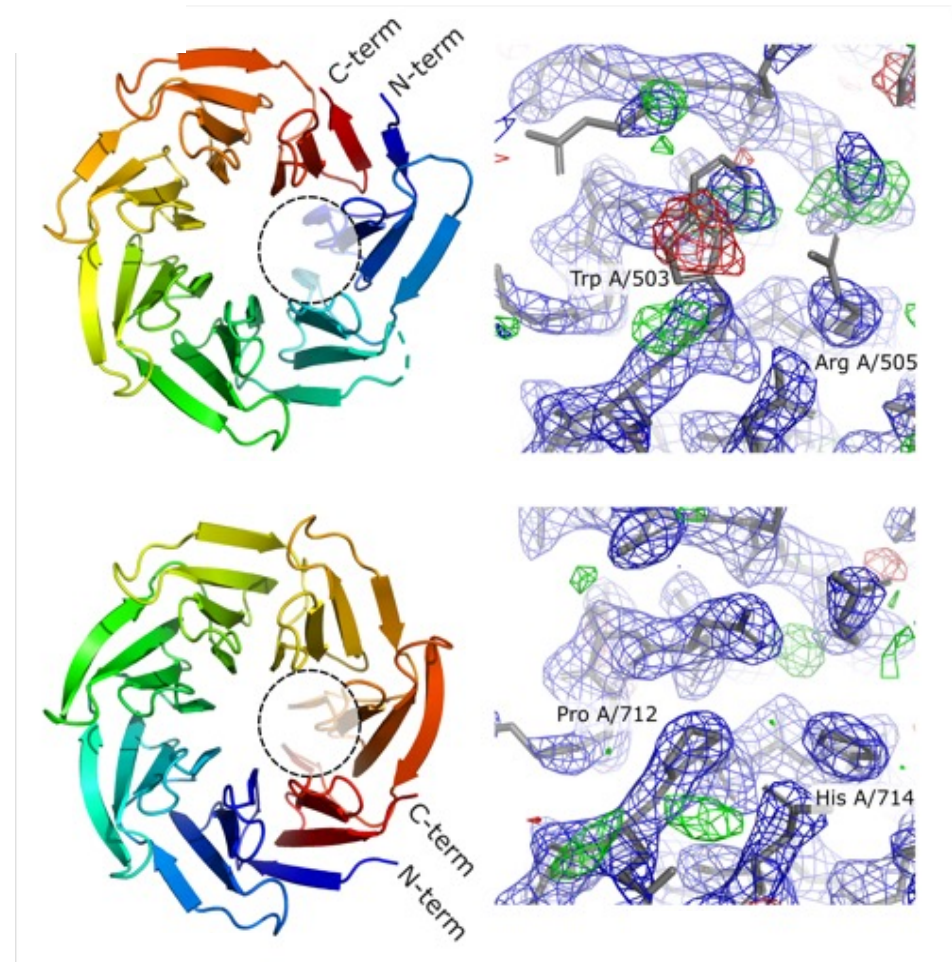
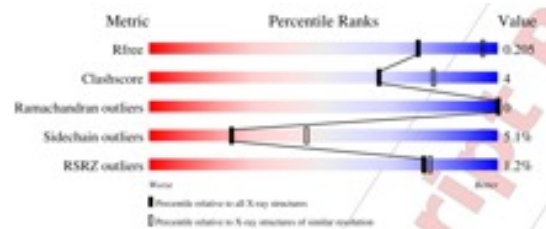
Register shifts in protein crystal structures – example 1

WD40-repeat domain from *T. curvata* @2.5Å

deposited
R/Rfree 22/26

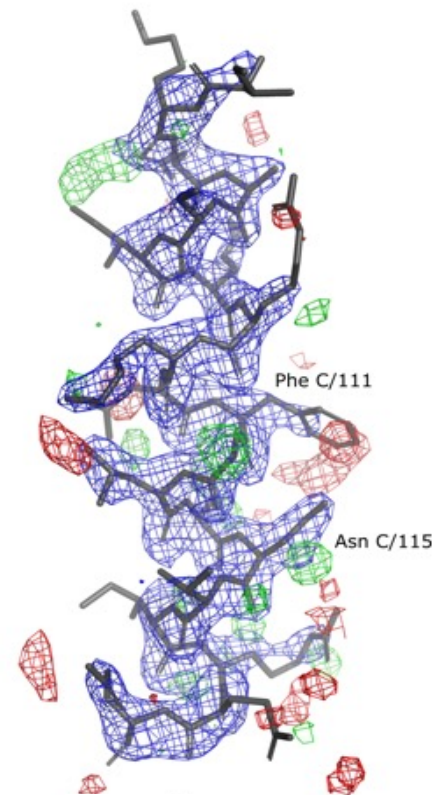
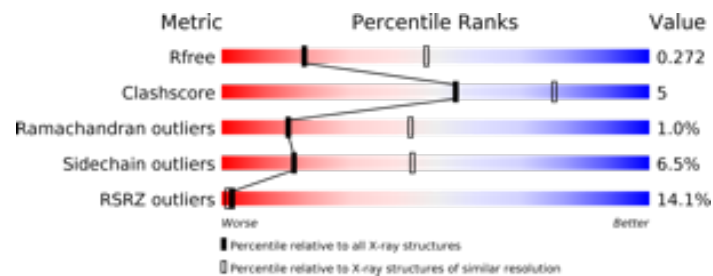


corrected
R/Rfree 17/21

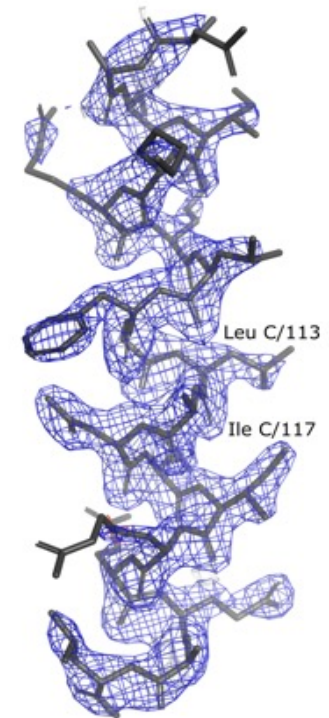


Register shifts in protein crystal structures – example 2

Helicase form *H. pylori* @2.5Å



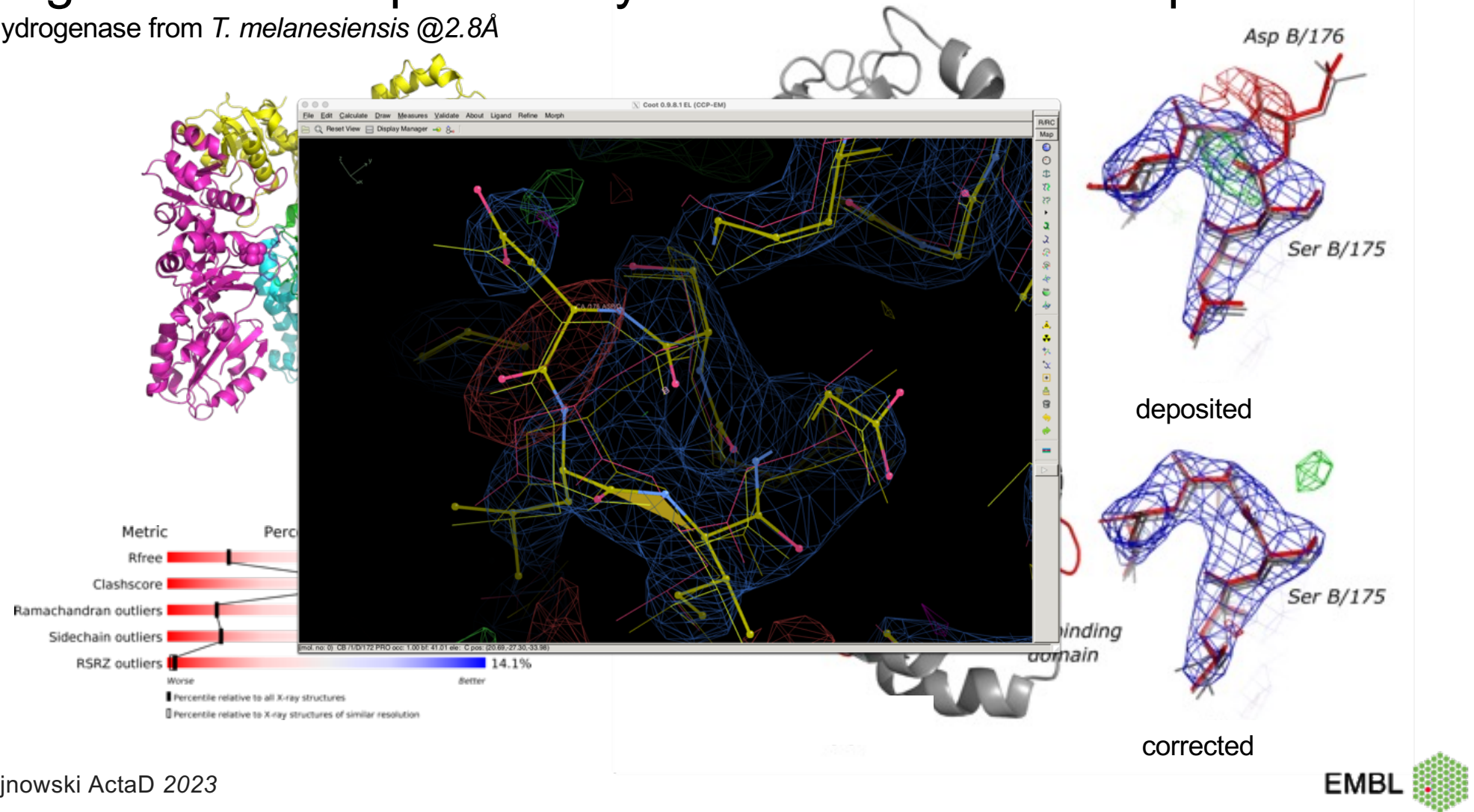
deposited



corrected

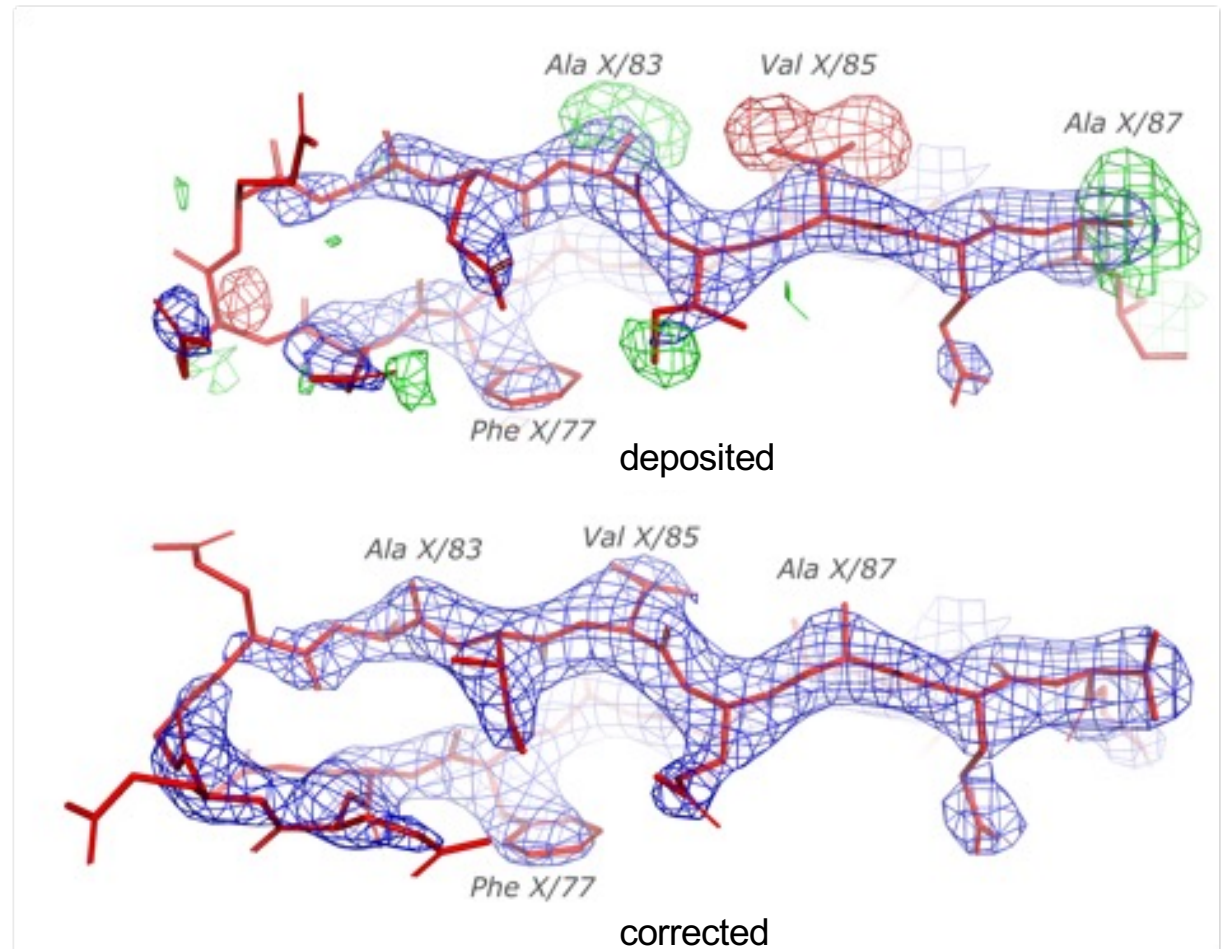
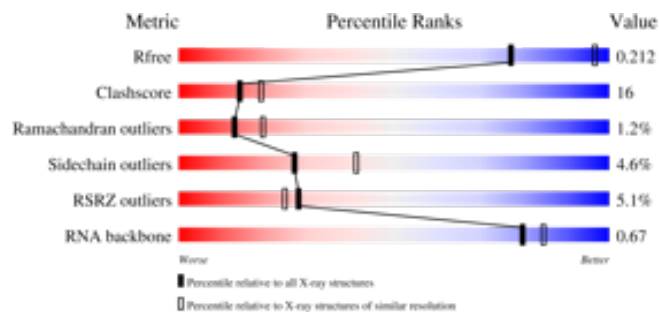
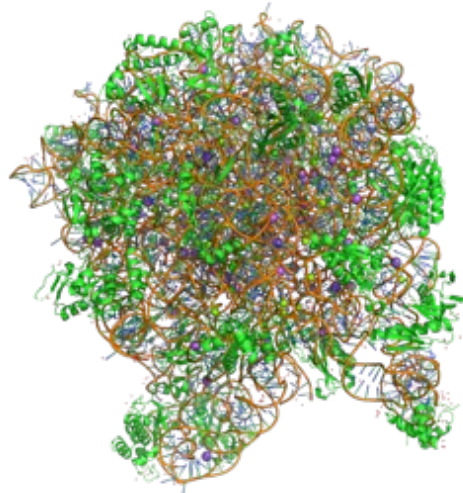
Register shifts in protein crystal structures – example 3

hydrogenase from *T. melanesiensis* @2.8Å



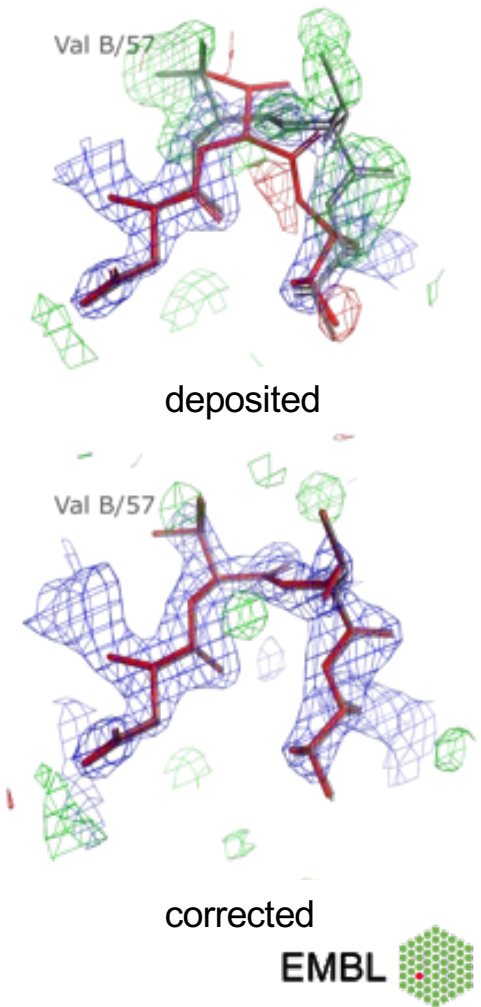
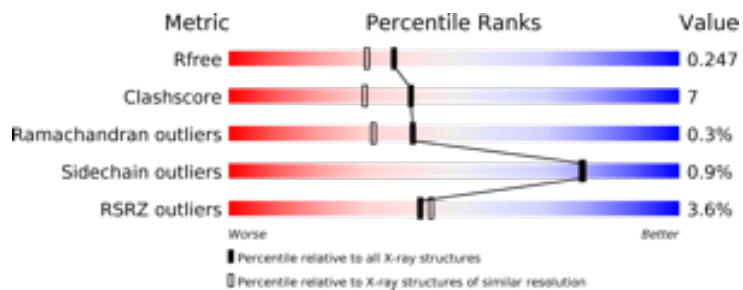
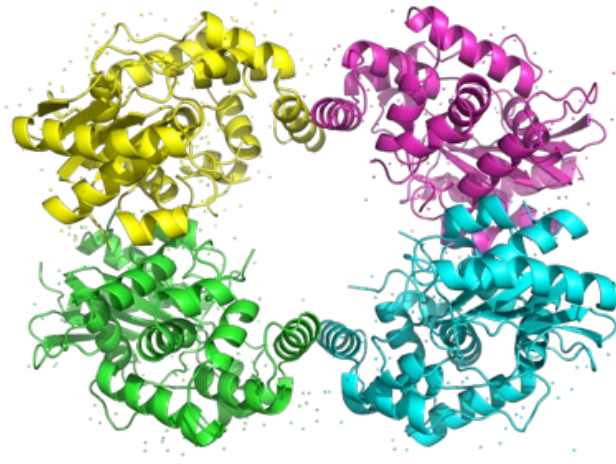
Register shifts in protein crystal structures – example 4

Ribosomal protein L31e from *H. Marismortui* @2.65Å



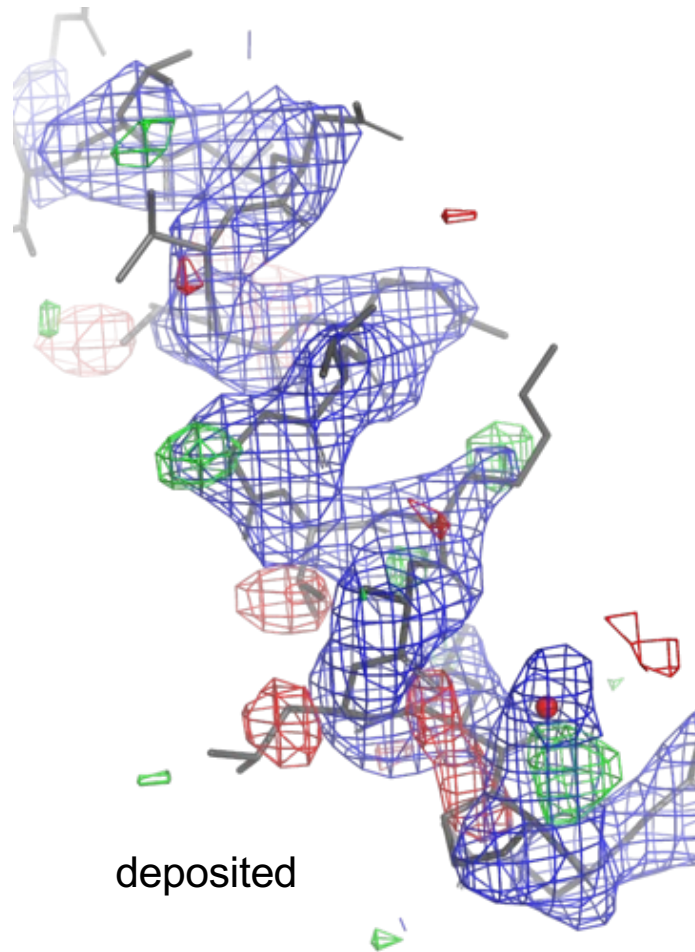
Register shifts in protein crystal structures – example 5

glutaminase from *G. kaustophilus* @2.1Å

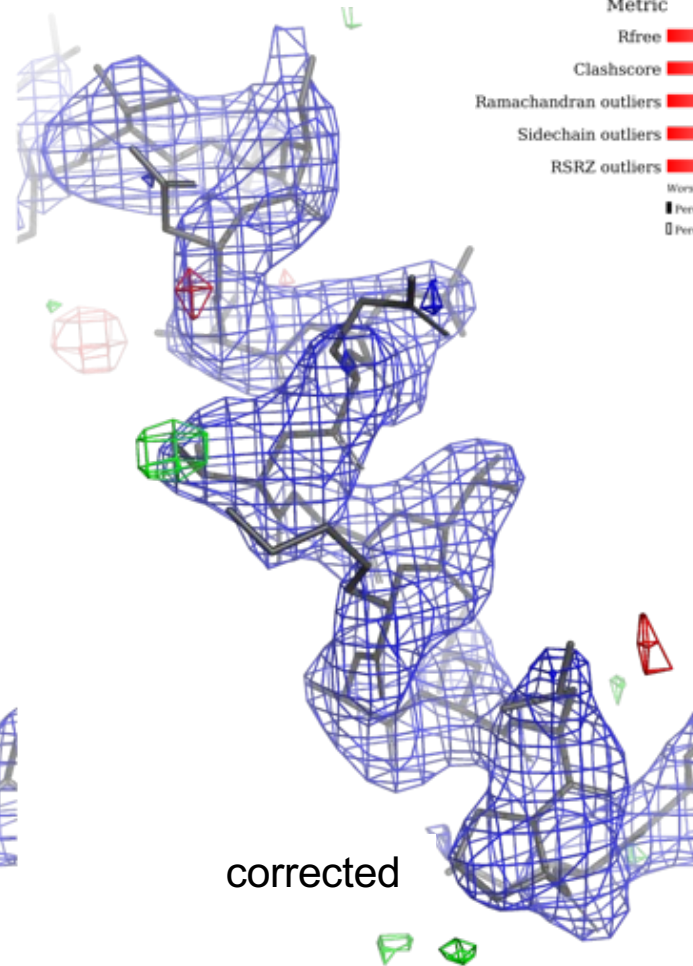


Example from last month's PDB release...

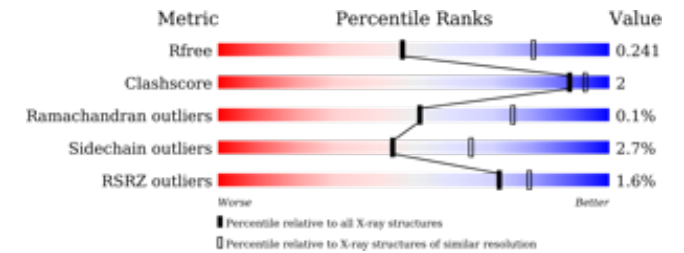
DNA binding protein at 2.75Å



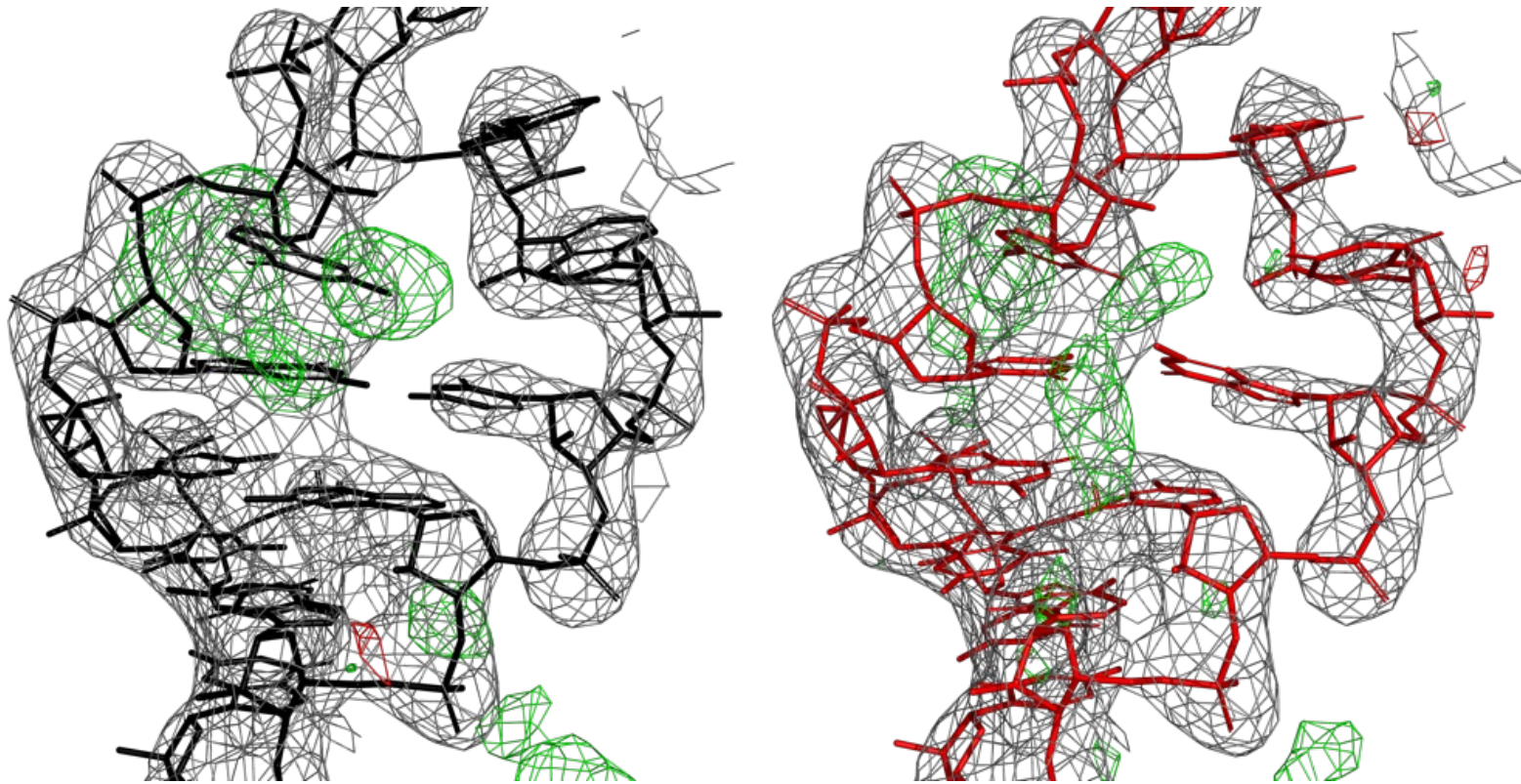
deposited



corrected

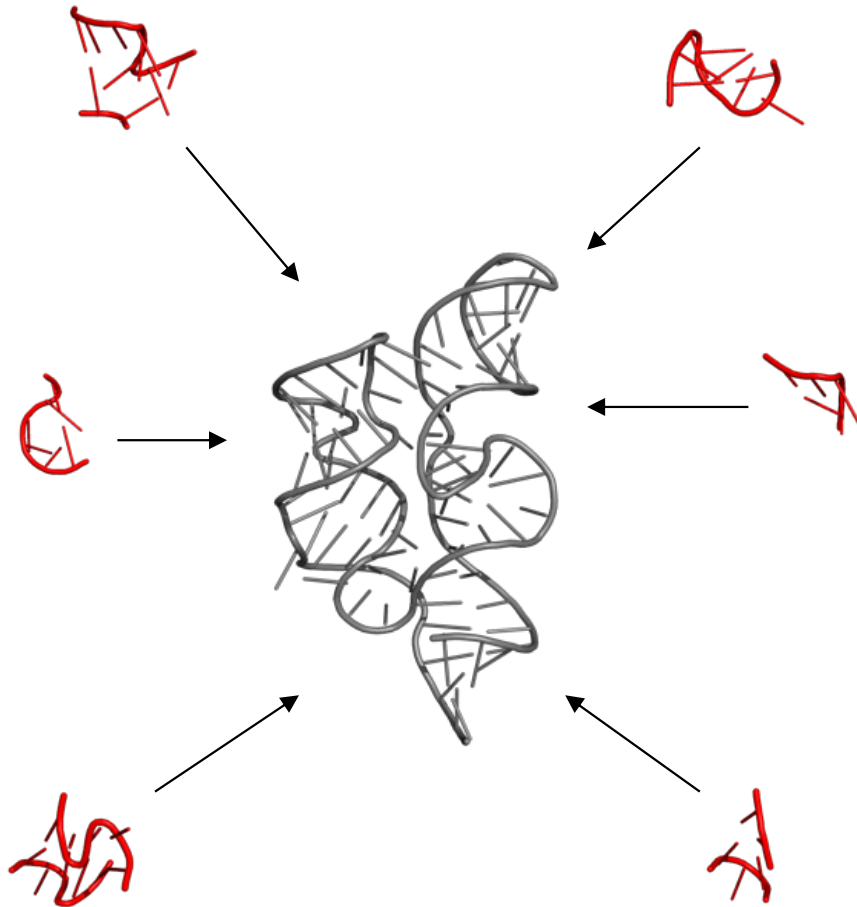


Nucleic acid model building can be challenging



Different types of purines and pyrimidines are (usually) indistinguishable in MX or EM maps

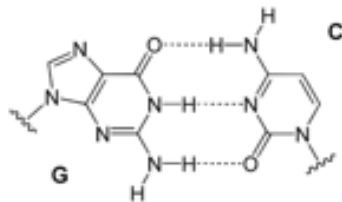
Base-pairs from backbone geometry



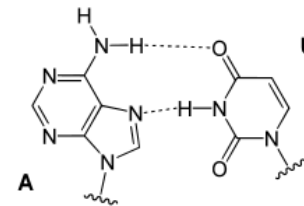
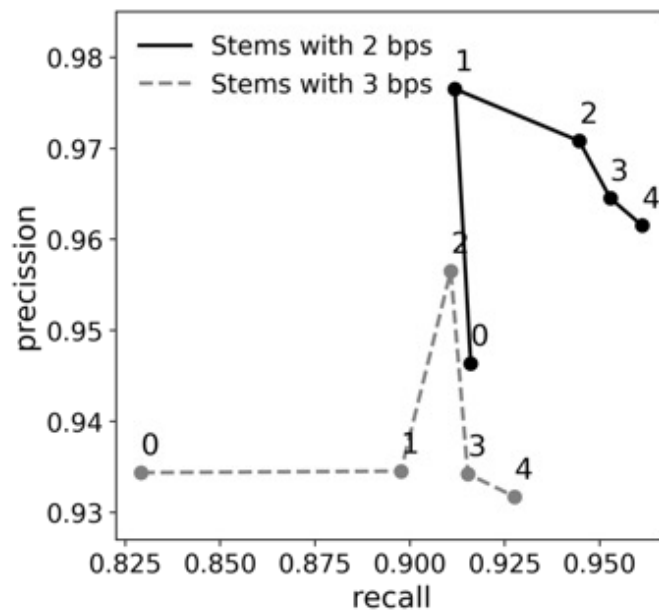
input RNA model backbone is “covered”
with motifs with known secondary structure

- base pairs assignment
- sequence validation
- refinement restraints

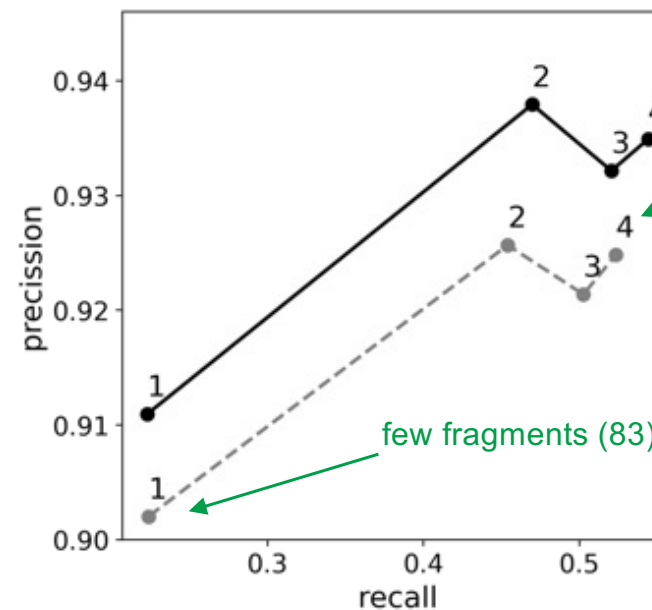
Base-pairs from backbone geometry



Watson-Crick (canonical) base-pairs



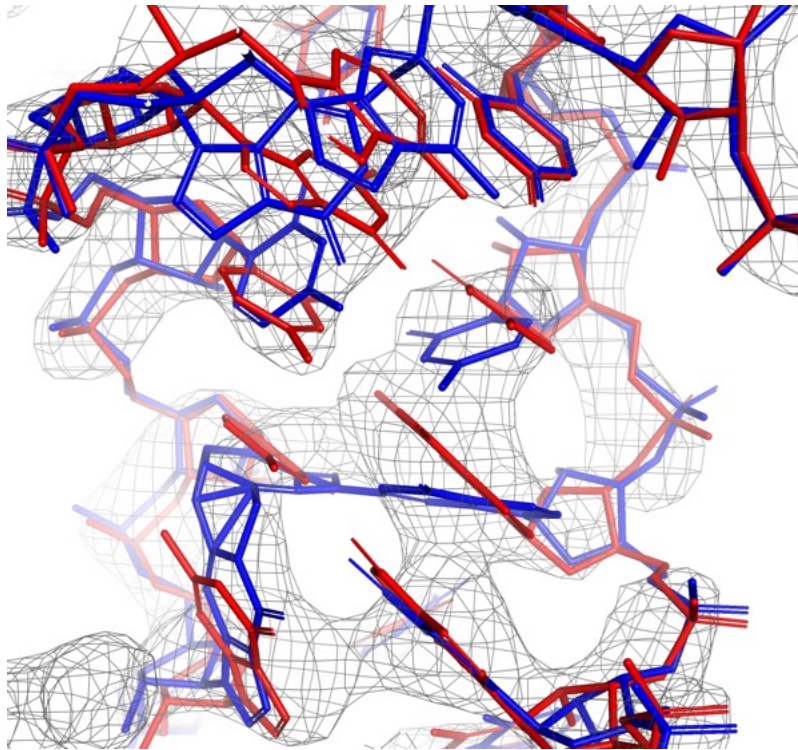
non-canonical base-pairs



lots of fragments (2,664)

few fragments (83)

doubleHelix - base-pairs from backbone geometry



Standard tools

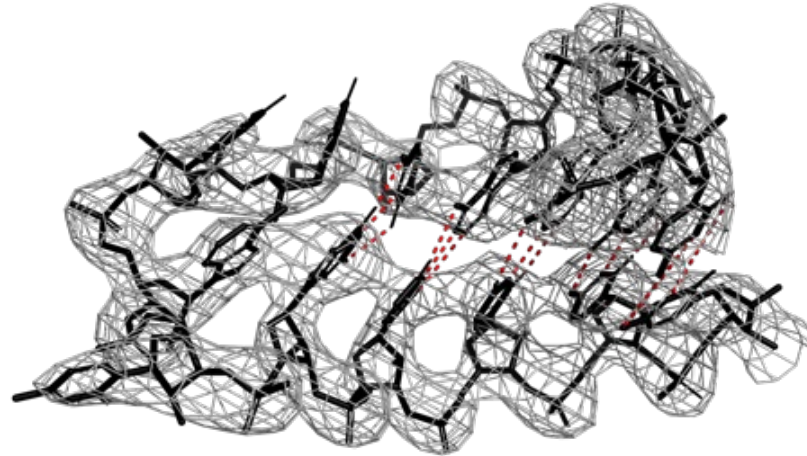
`phenix.secondary_structure_restraints`
LibG

cannot restrain RNA base-pairs if

- geometry is very bad
- sequence is wrong

ARP/wARP model of 23S fragment
in cryo-EM map @ 3.2Å resolution

checkMySequence and NA crystal structures



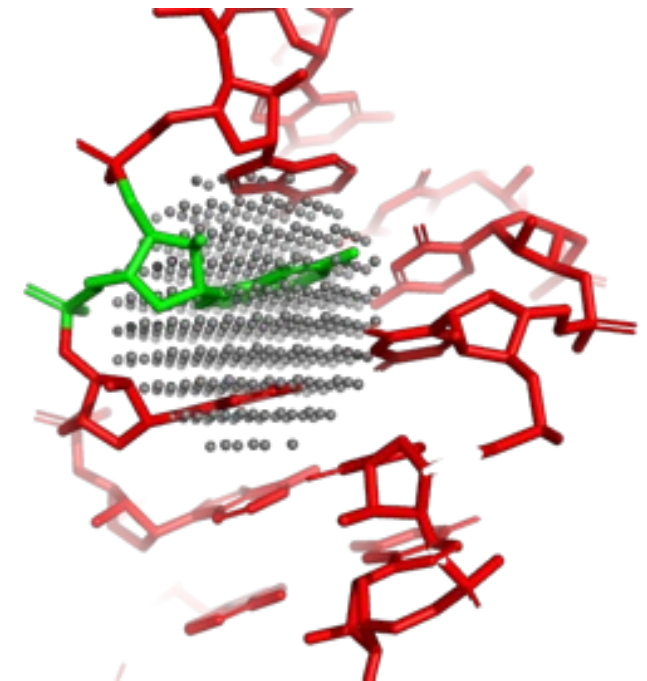
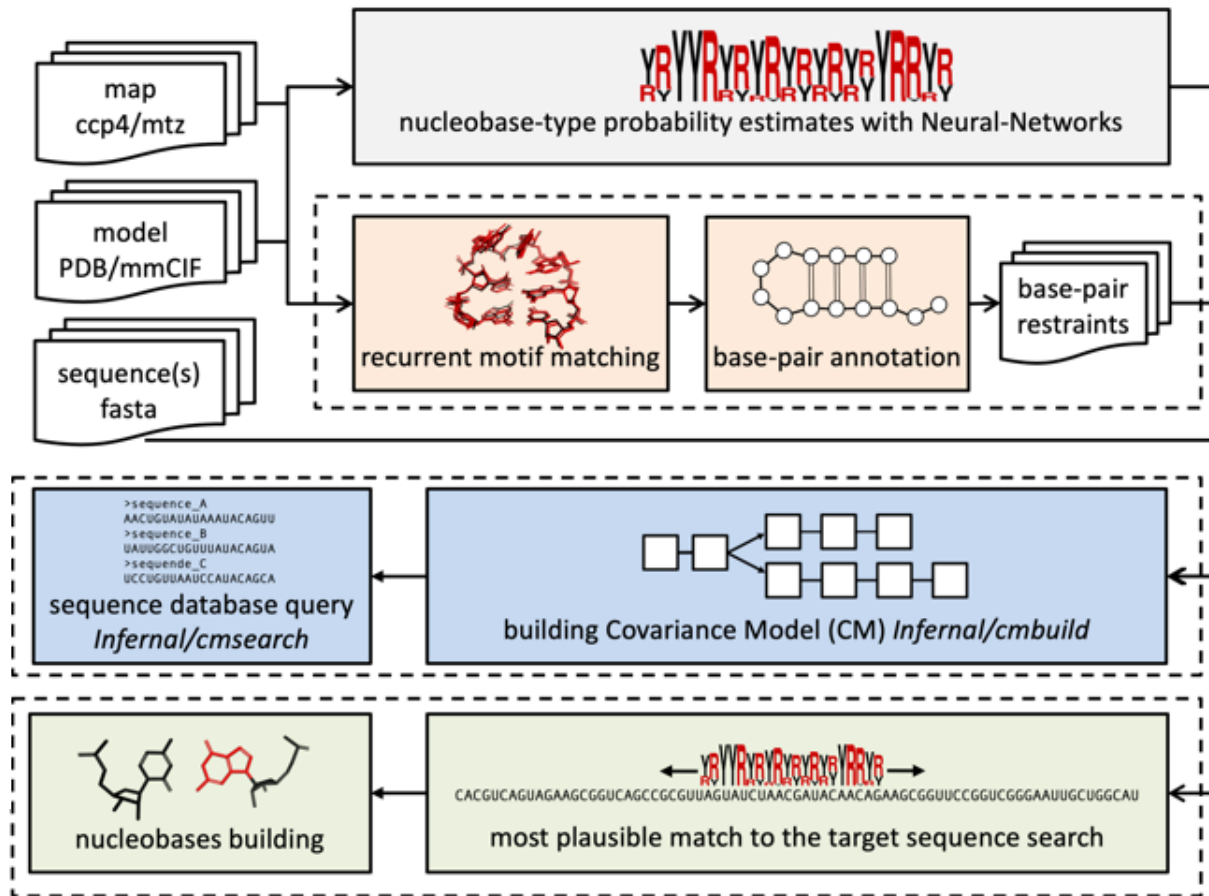
	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?



target sequence

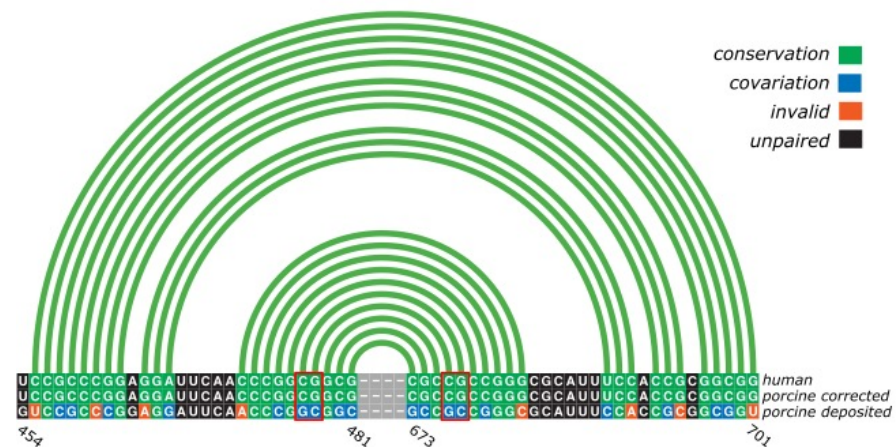
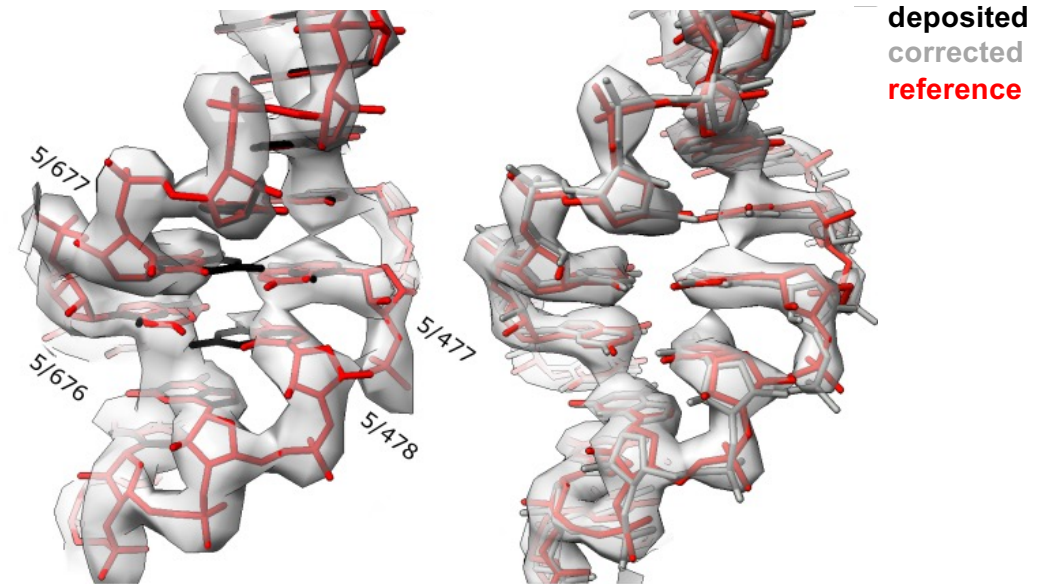
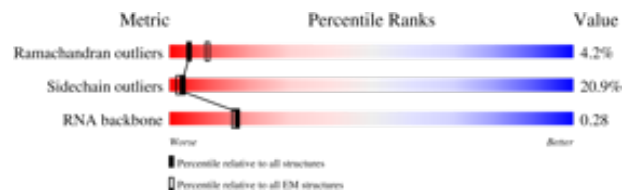
alignment p-value

Nucleic acid sequence assignment and validation in EM/MX



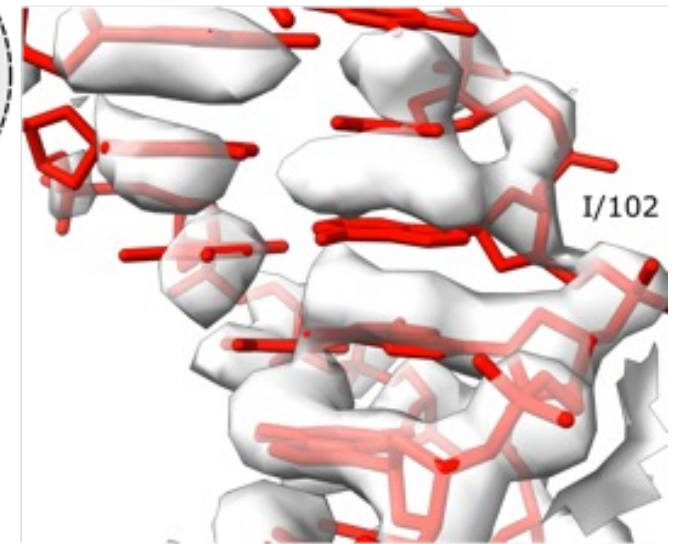
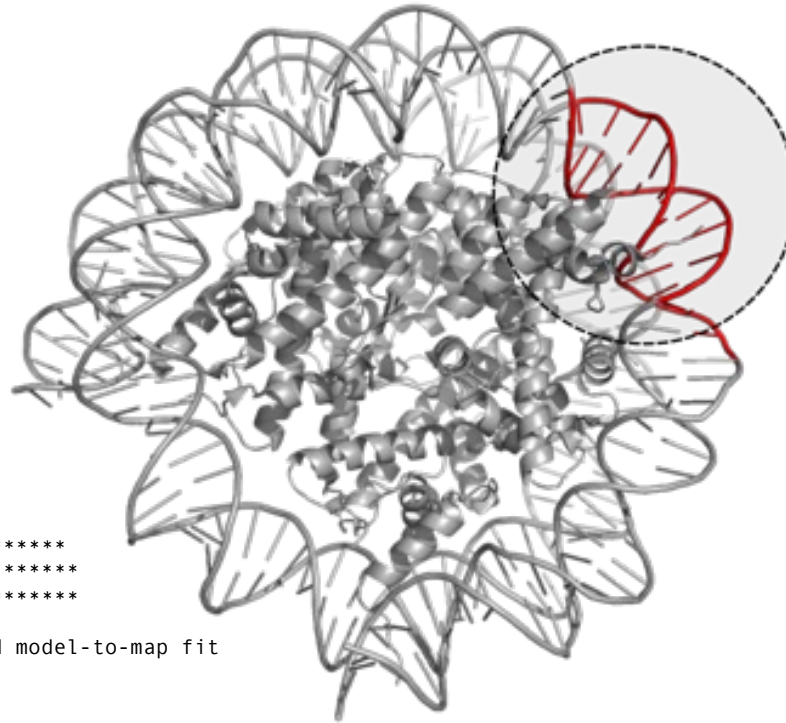
Sequence assignment issues in rRNA

EM model of mammalian ribosome @3.4Å (deposited 2014)



Chojnowski NAR, 2023

recent EM use case: nucleosome



```
*****  
***** SUMMARY *****  
*****
```

==> Unidentified chains; check input sequences and model-to-map fit

I/2:144

==> Sequence register shifts

- nucleic-acid chain fragment J/182-211 may be shifted by -1 residue [p-value=9.28e-02]

model seq 182-211

actagGGAGTAATCCCCTTGGCGGTAAAACGCGGgggacag

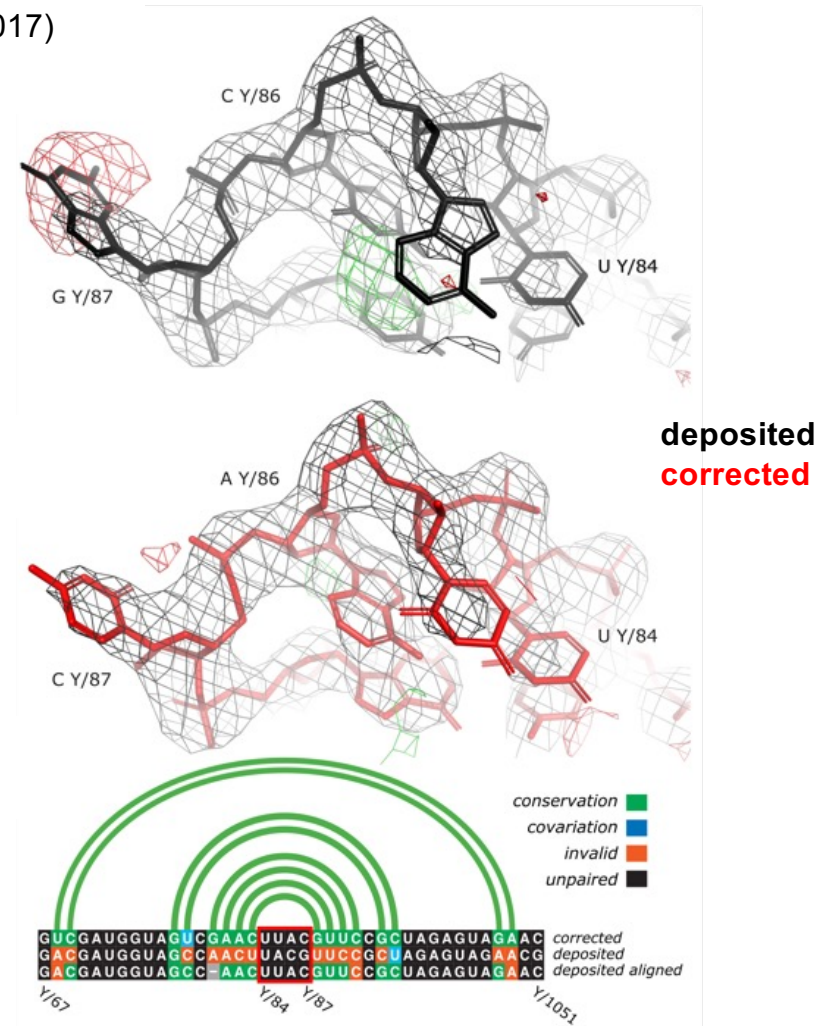
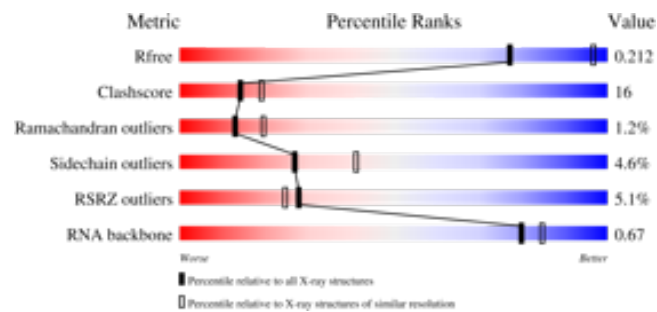
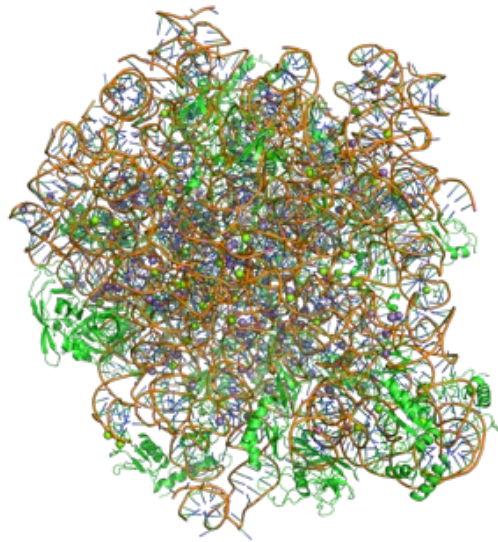
new seq 183-212

actaggGAGTAATCCCCTTGGCGGTAAAACGCGGgggacag

Time elapsed 0:00:19

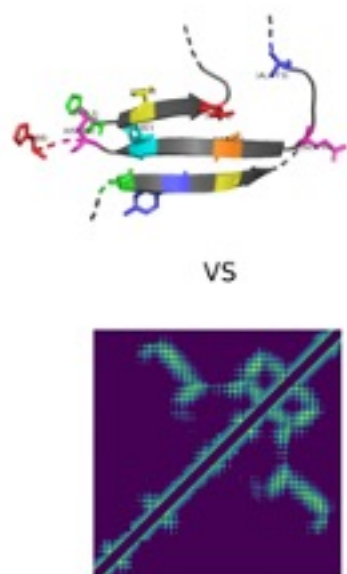
Sequence assignment issues in rRNA

Crystal structure model of bacterial ribosome @3.5Å (deposited 2017)

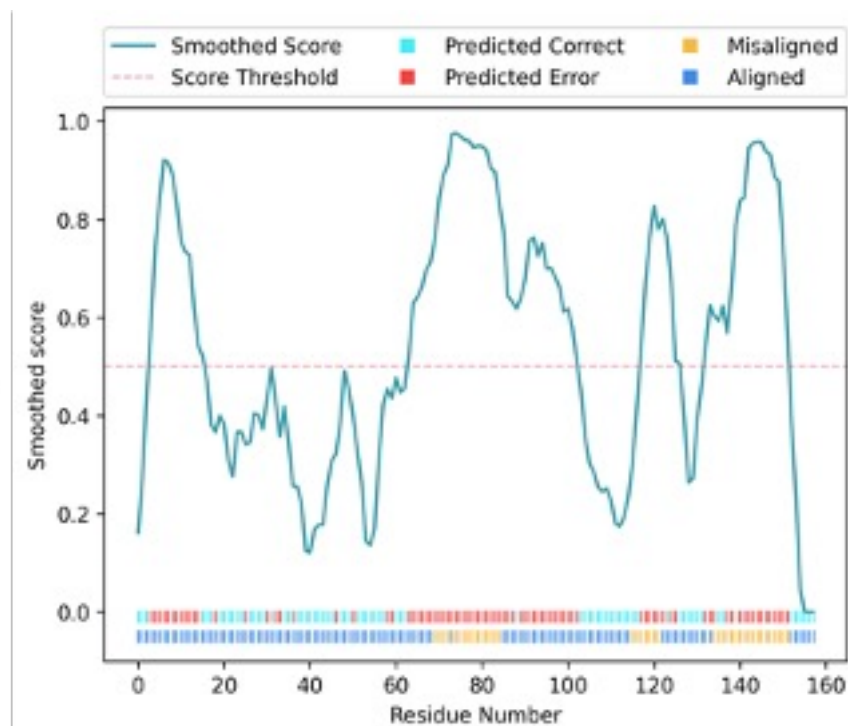


Predicted inter-residue distances for model validation

www.conkit.org



FEATURES	RESIDUE NUMBER					
	1	2	3	4	5	
RMSD	0.3	2.5	0.5	1.2	0.6	→ SVM CLASSIFIER →
FN RATE	3.5	0.2	1.5	2.8	1.1	
FP RATE	0.3	2.5	0.1	0.2	0.2	
SENSITIVITY	0.9	0.8	0.8	0.3	0.9	
ACCURACY	0.3	0.1	0.8	0.8	0.1	



Sanchez Rodriguez, F., Chojnowski, G., Keegan, R. M. & Rigden, D. J. (2022). Acta Cryst. D78, 1412-1427.

Acknowledgements

University of Liverpool

Daniel Rigden
Adam Simpkin
Filomeno Sánchez Rodríguez

CCP4 Core Team

Ronan Keegan
Charles Ballard
Eugene Krissinel
Ville Uski
Kyle Stevenson
Maria Fando

CCPEM team

Agnel Joseph
Tom Burnley
Colin Palmer
Matt Iadanza

Martin Luther University

Panos Kastitis 
Ioannis Sklidis 

EMBL Hamburg





Matthias Wilmanns
Kate Beckham
Jan Kosiński
Christina Ritter
Edukondalu Mullapudi
Isabel Bento
Alice Bochel
Wolfram Seifert-Davila

São Carlos Institute of Physics

Diego A. Leonardo
Laboratorio de Biologia
Molecular, Peru
Dan E. Vivas-Ruiz








findMySequence

protein sequence identification
  
 

checkMySequence

sequence validation
   


doubleHelix

NA sequence identification
  
 
EMBL 