

Unknown sequences and hidden errors in macromolecular structure models

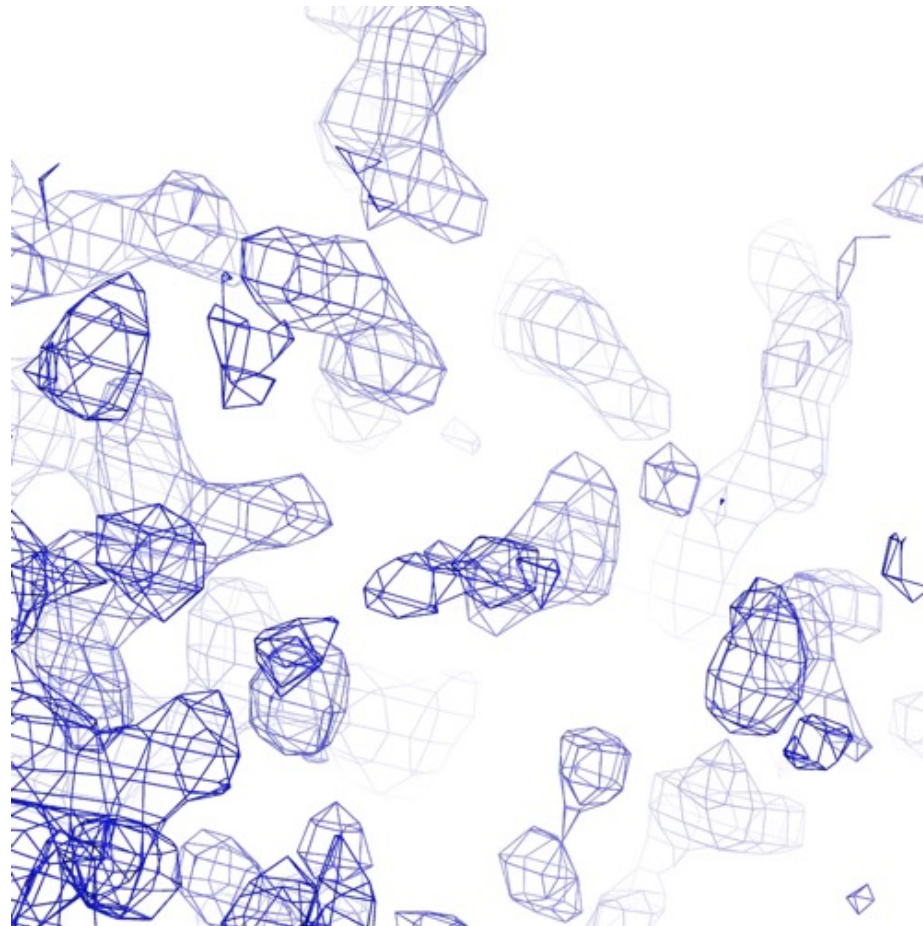
(New developments in model building and validation)

Grzegorz Chojnowski
EMBL Hamburg

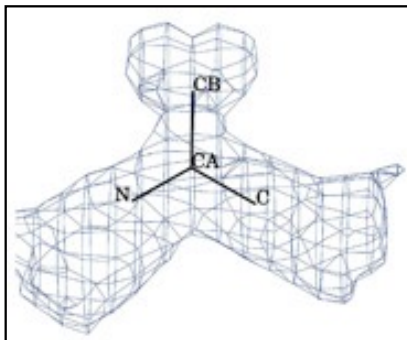
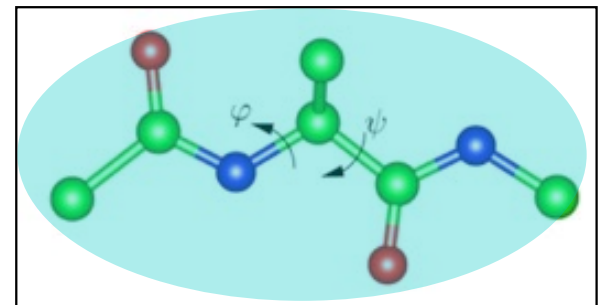
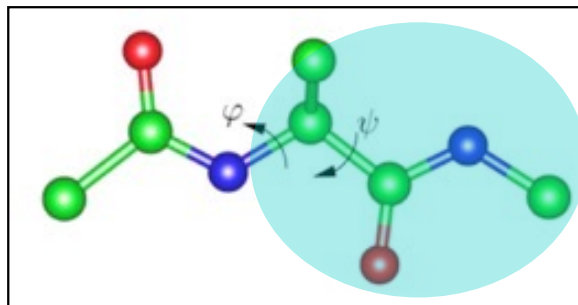
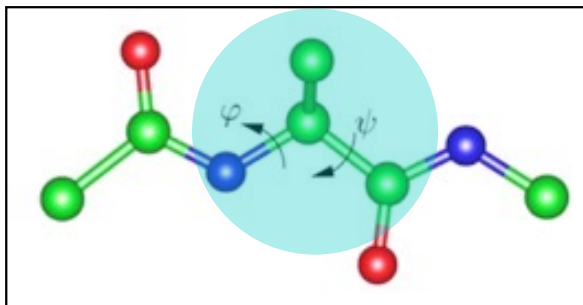
Diamond-CCP4 workshop 05.12.2022



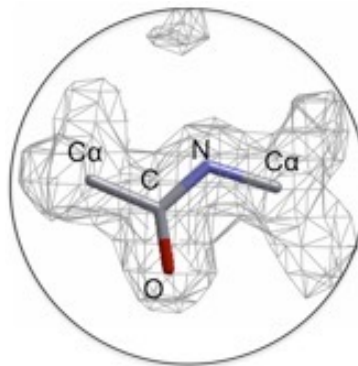
Electron density map interpretation



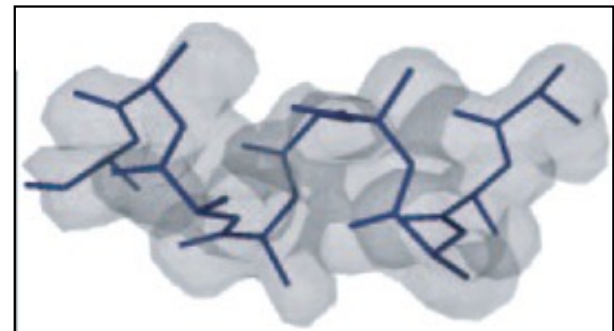
Crystal structure model building



Buccaneer
ModelCraft

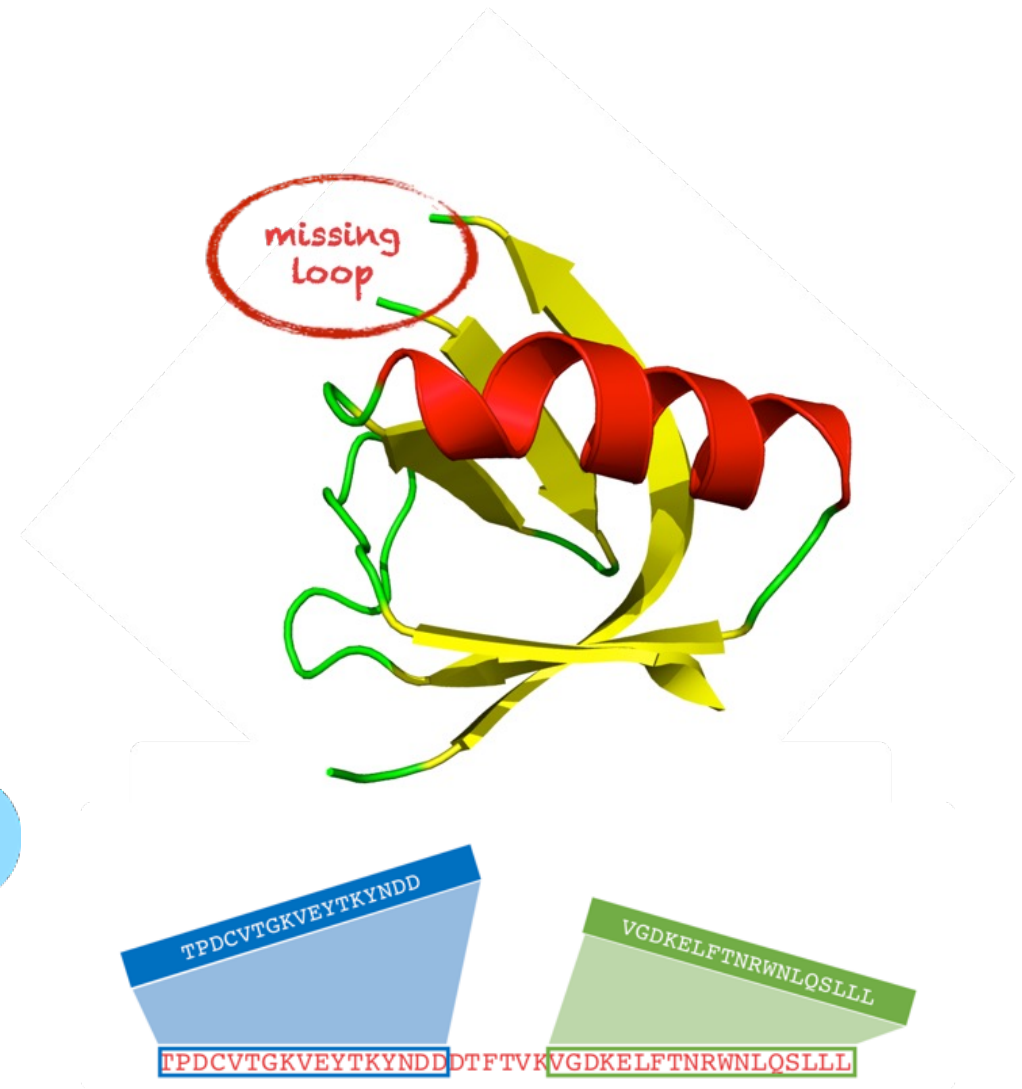
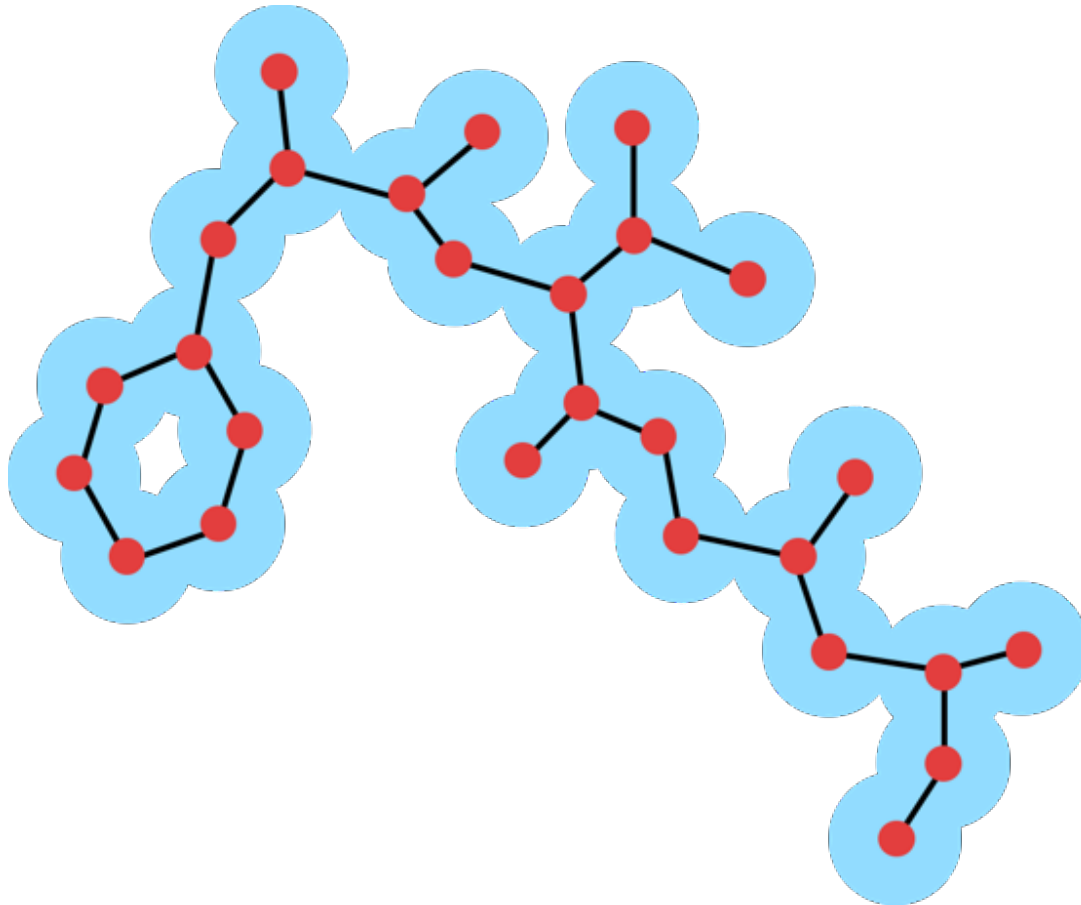


ARP/wARP

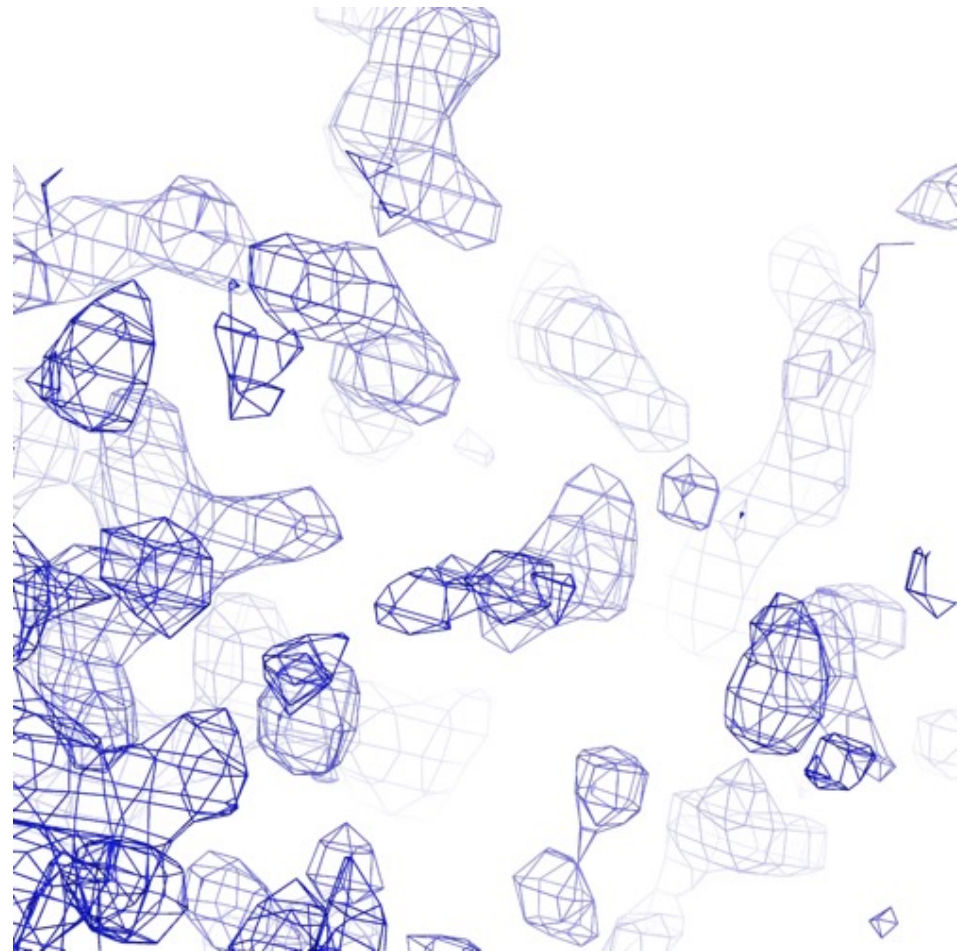
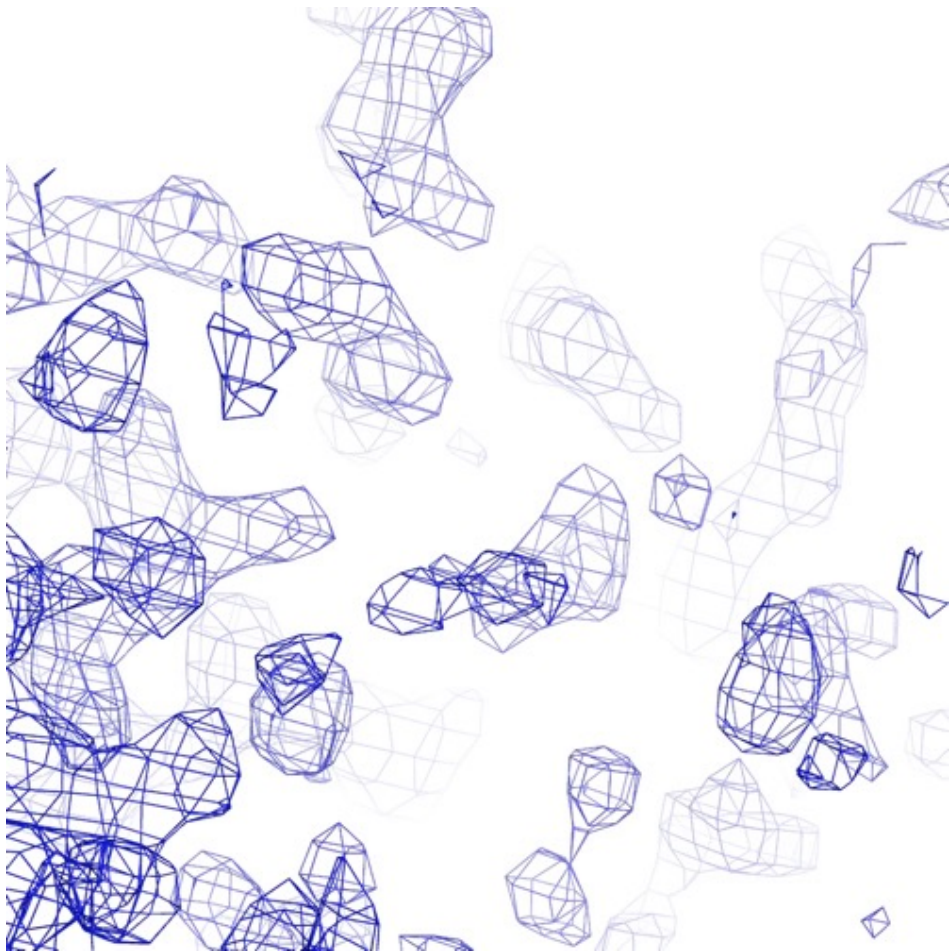


Phenix.autobuild

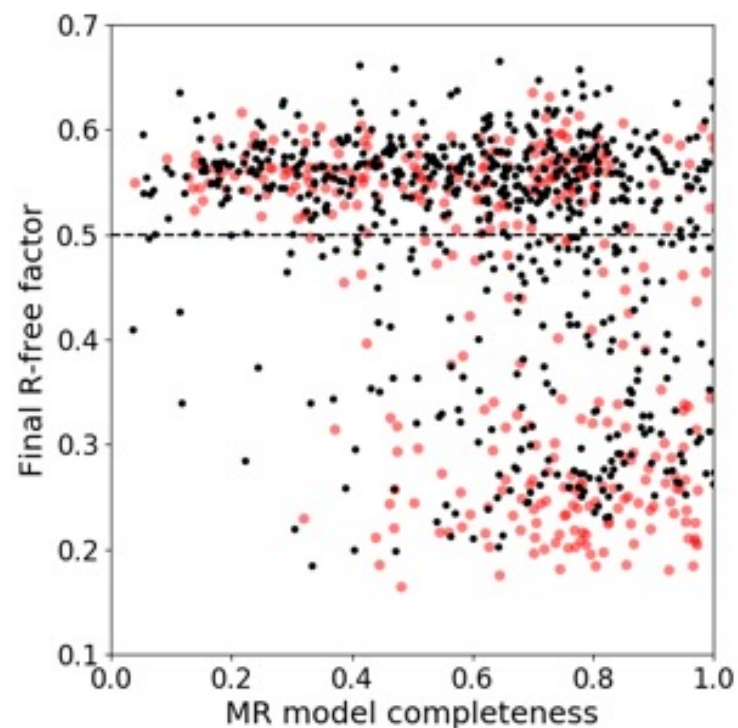
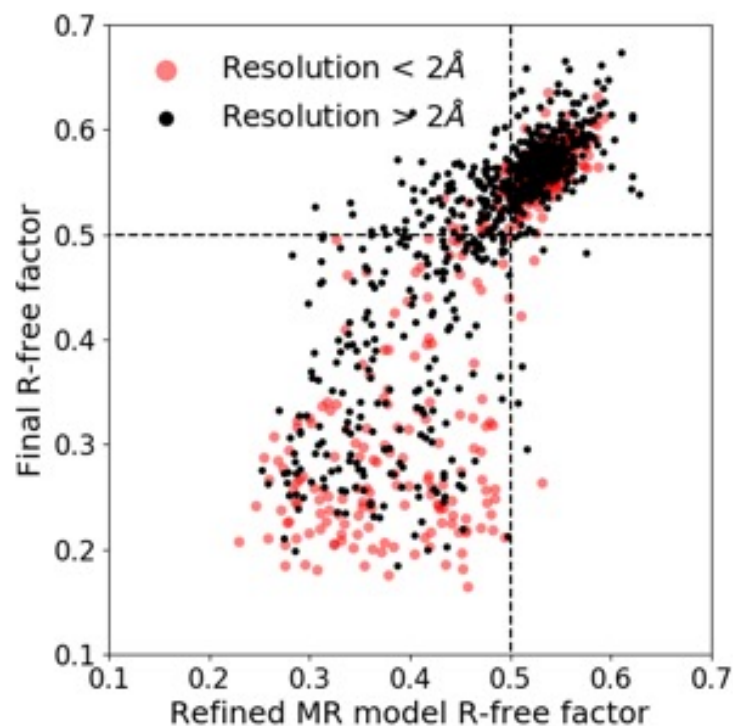
Crystal structure model building



Crystal structure model building

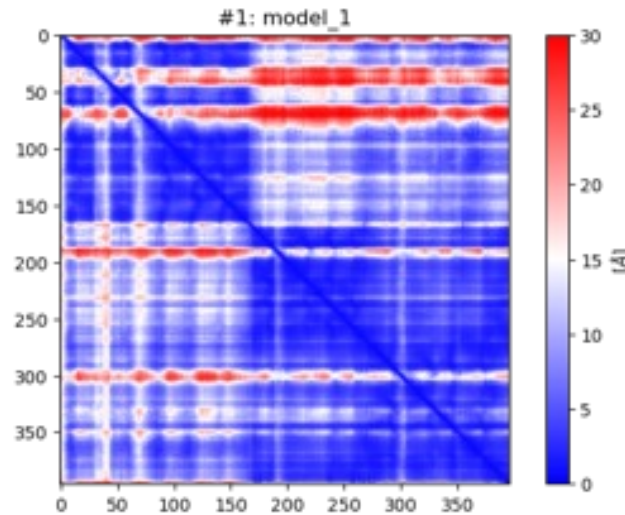
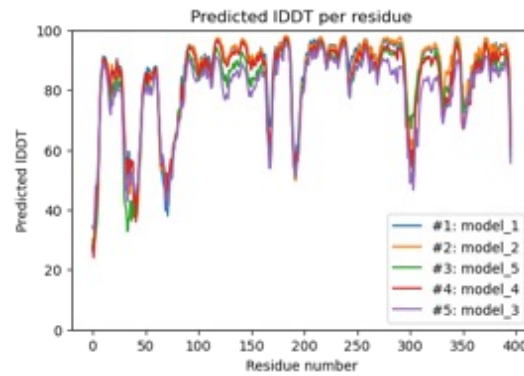


Crystal structure model building



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

Solving crystal structures with AlphaFold2 models

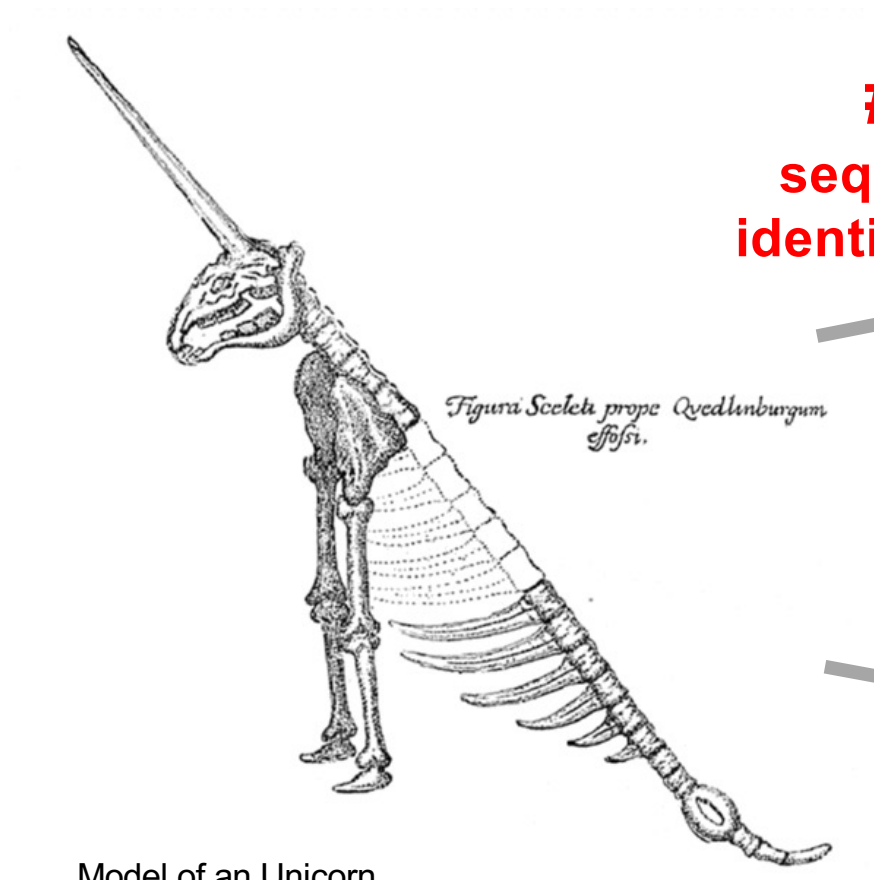


AlphaFold2
prediction



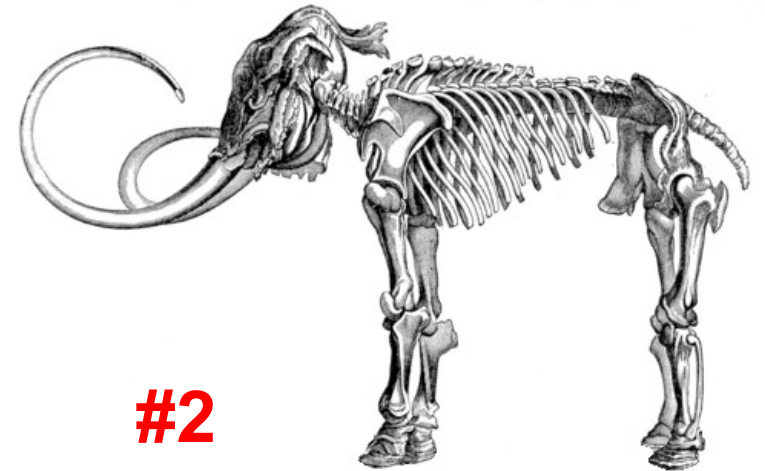
crystal structure
model

Model building traps in MX (and cryo-EM)

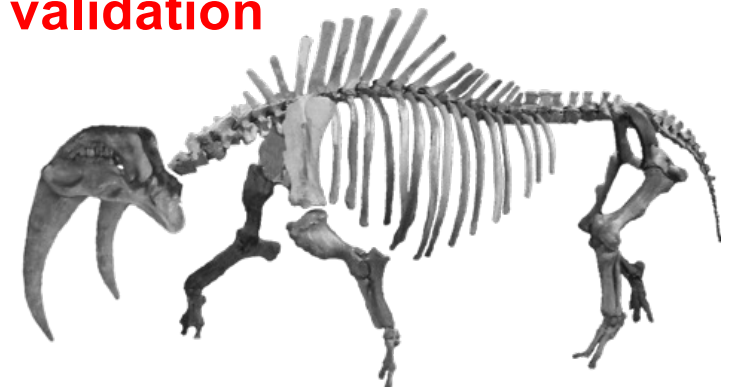


Model of an Unicorn
Gottfried Leibniz after Otto von
Guericke, *Protogaea* (1719)

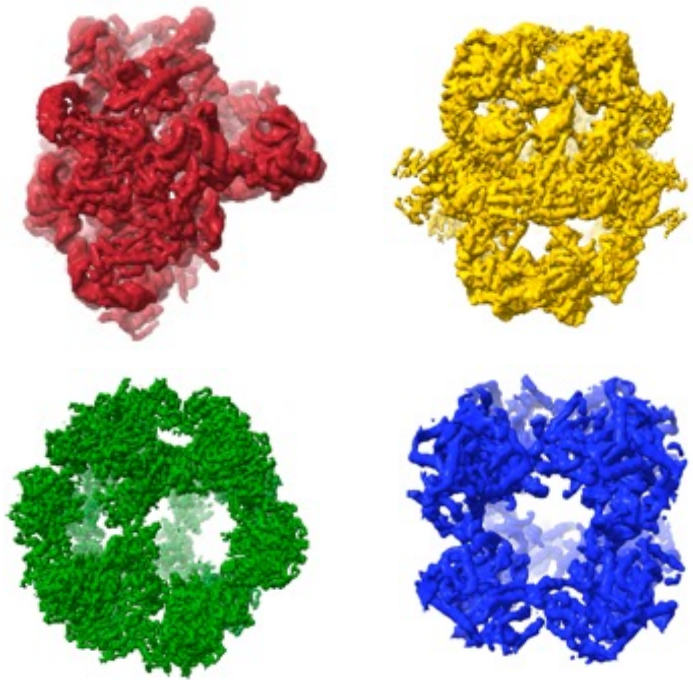
#1
sequence
identification



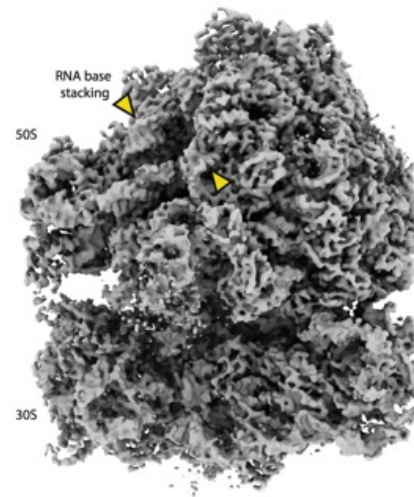
#2
sequence
validation



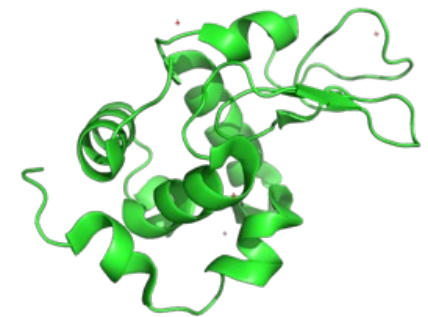
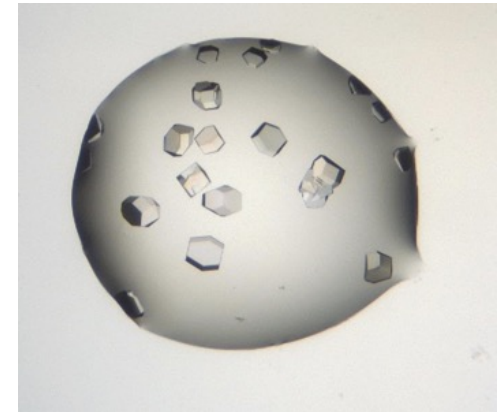
Unknown proteins in EM and MX



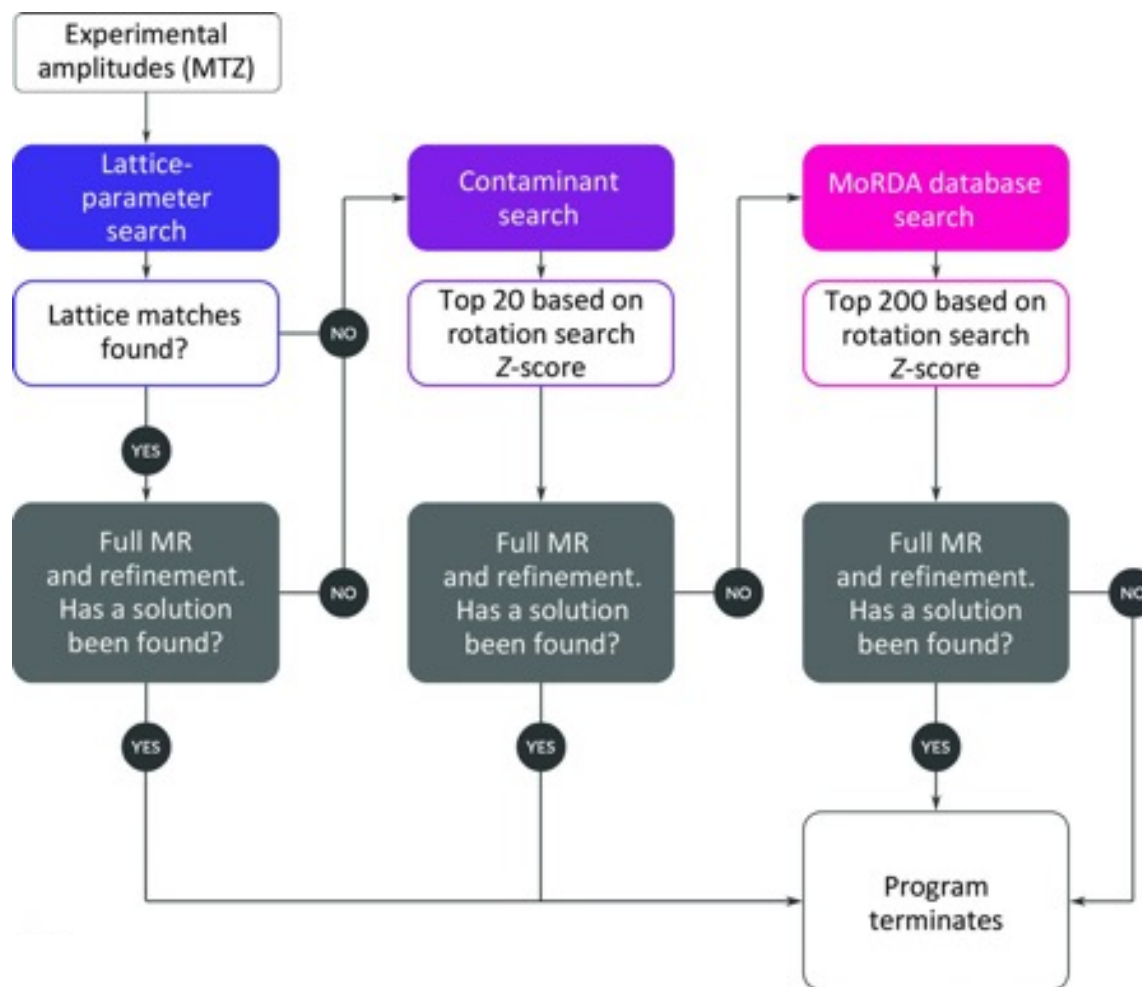
Cryo-EM and artificial intelligence visualize endogenous protein community members
Skalidis et al. Structure 2022



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

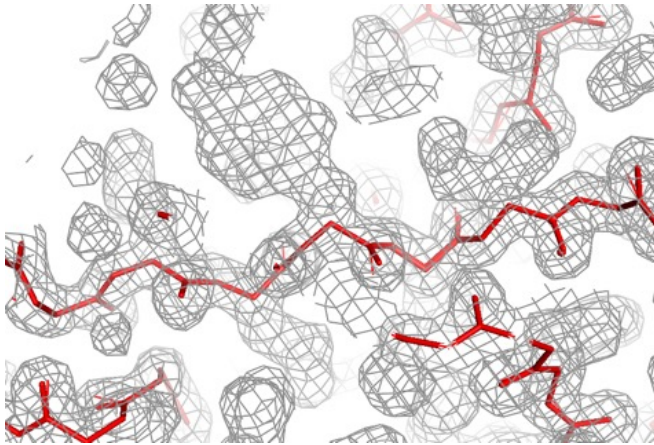


Sequence-free MR with SIMBAD



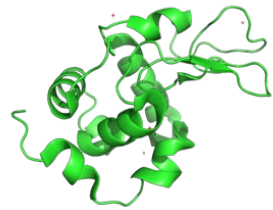
Available in


Protein sequence identification from a map



main-chain model and a map

GESAMT, FATCAT, DALI



protein structure database

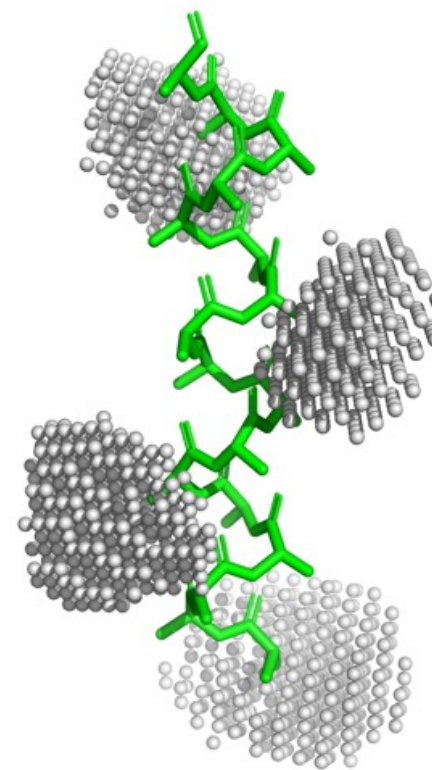
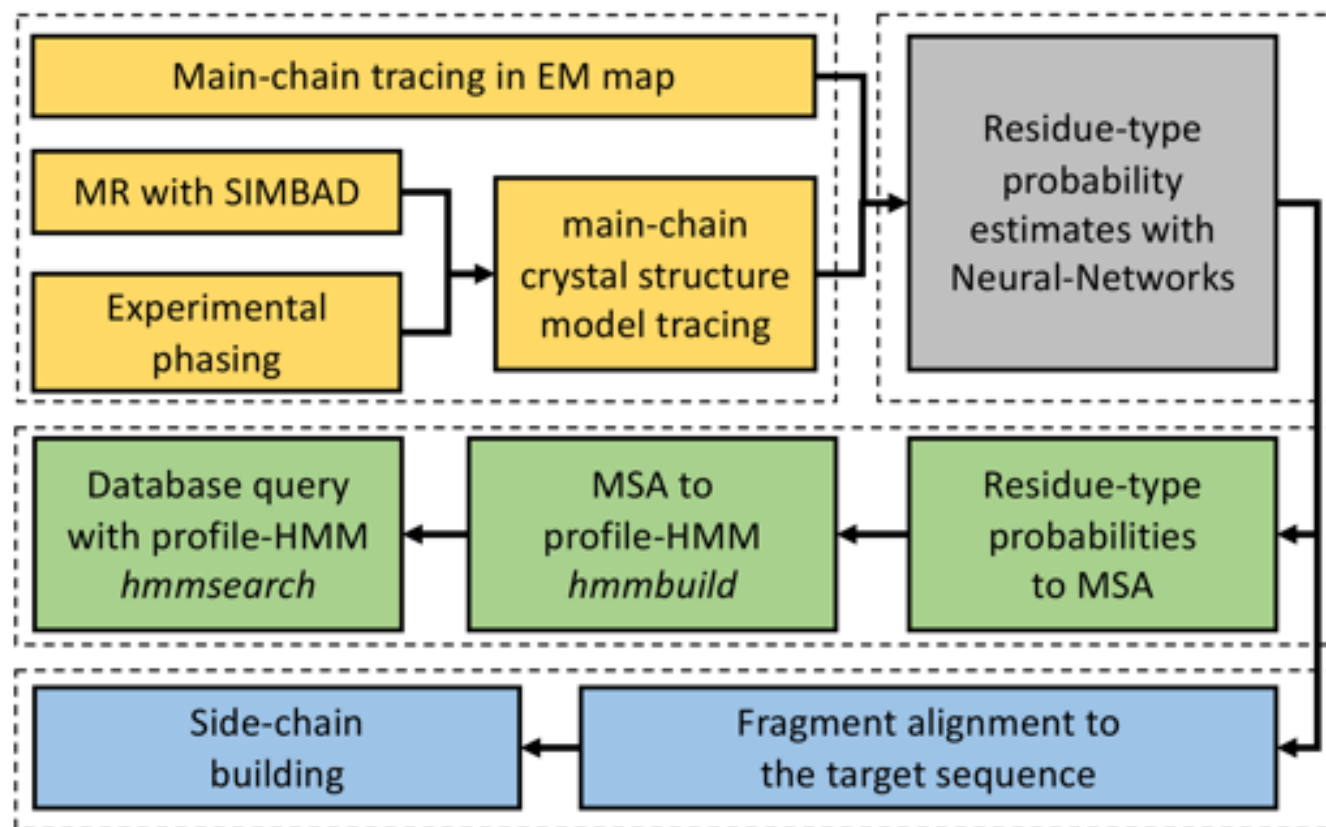


sequence profile



sequence database

Protein sequence identification with findMySequence

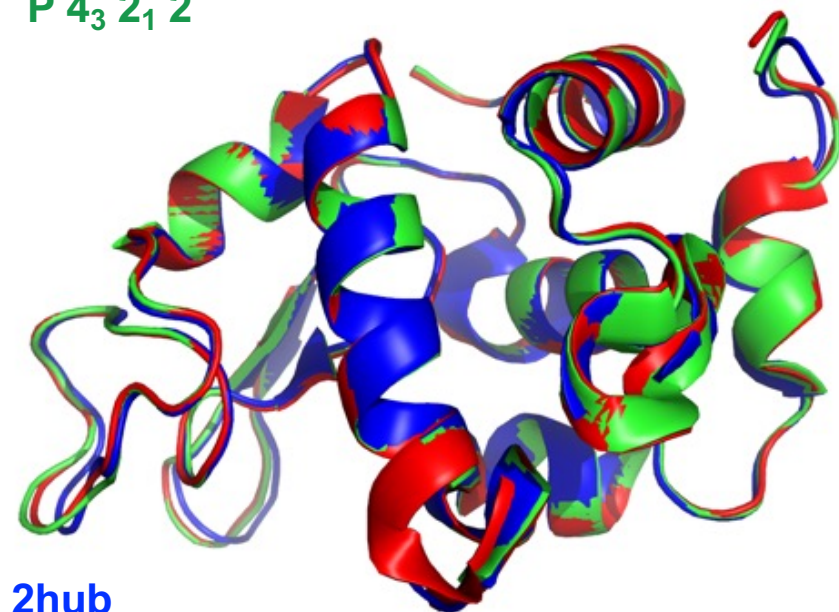


residue-type NN classifier
328 input features
111,800 parameters

MX benchmarks: three hen egg-white lysozyme targets

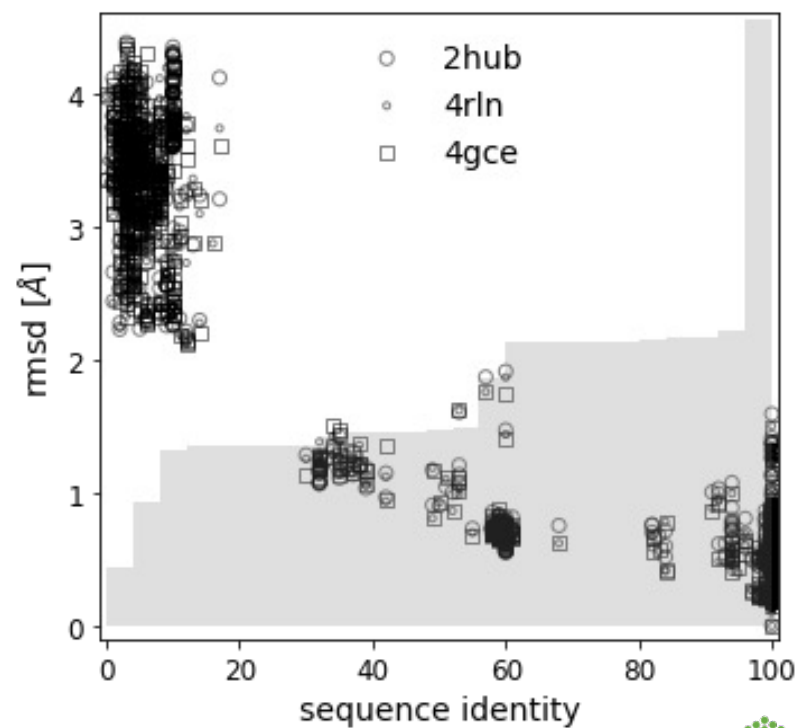
4gce
2.9Å
P 4₃ 2₁ 2

4rln
2.2Å
P 4₃ 2₁ 2

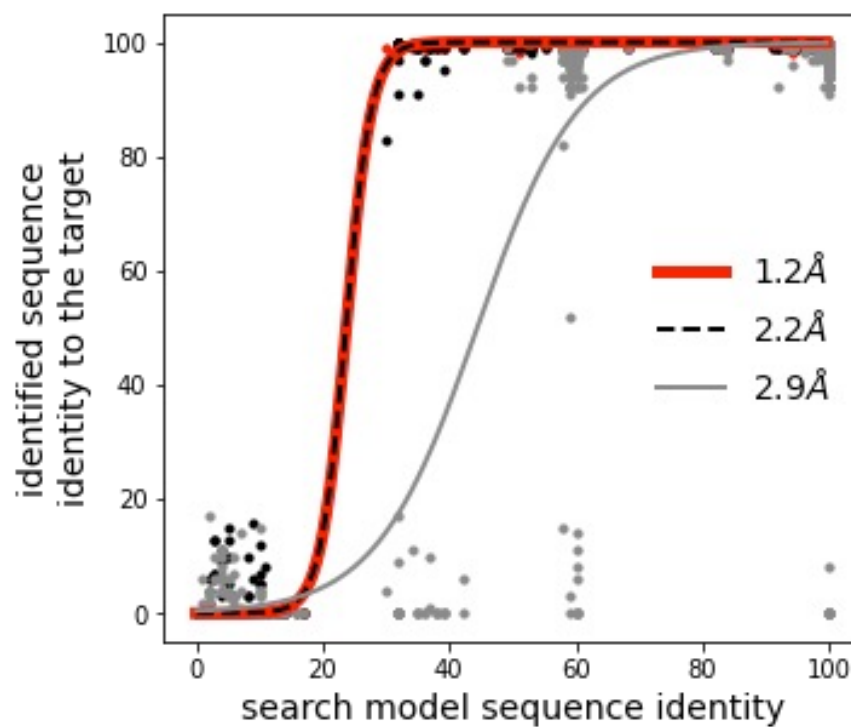
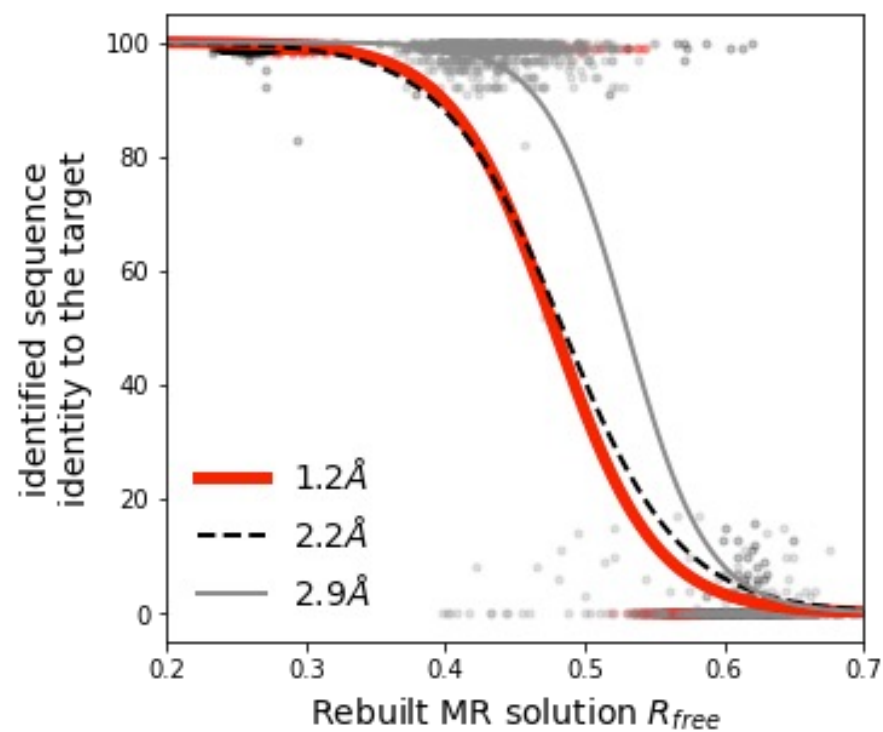


2hub
1.2Å
P 4₃ 2₁ 2

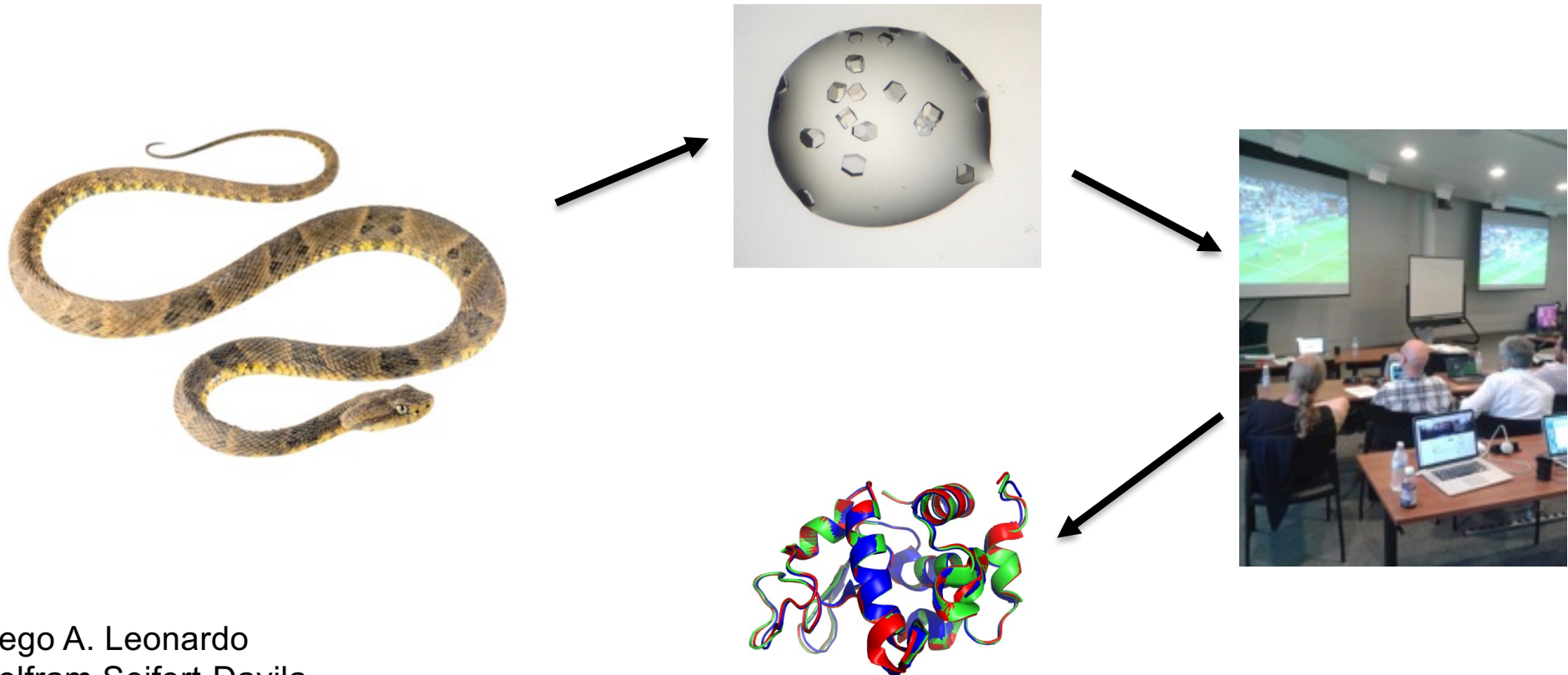
1,500 artificial
MR search models



MX benchmarks: three hen egg-white lysozyme targets



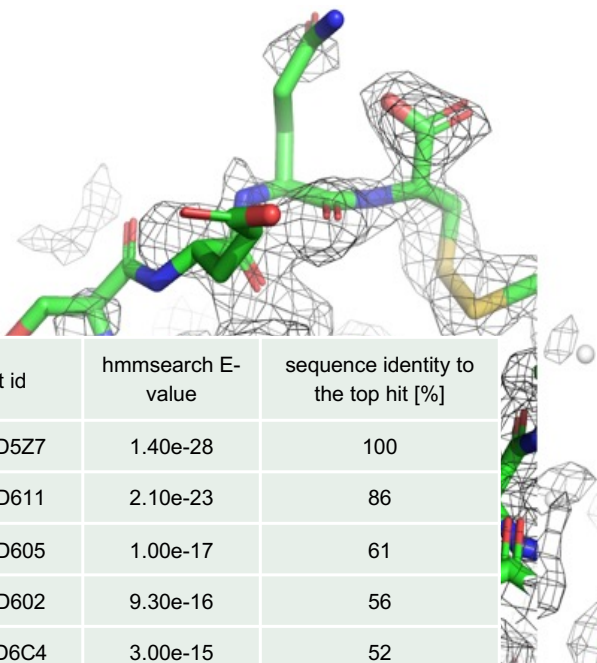
MX use-case: *Bothrops atrox* venom proteins



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

MX use-case: *Bothrops atrox* venom proteins

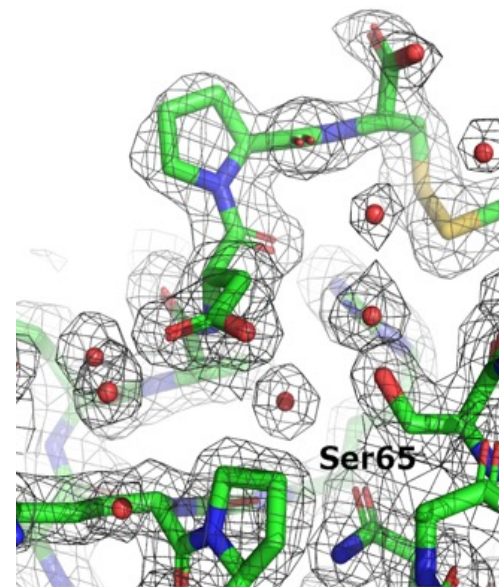
initial MR solution



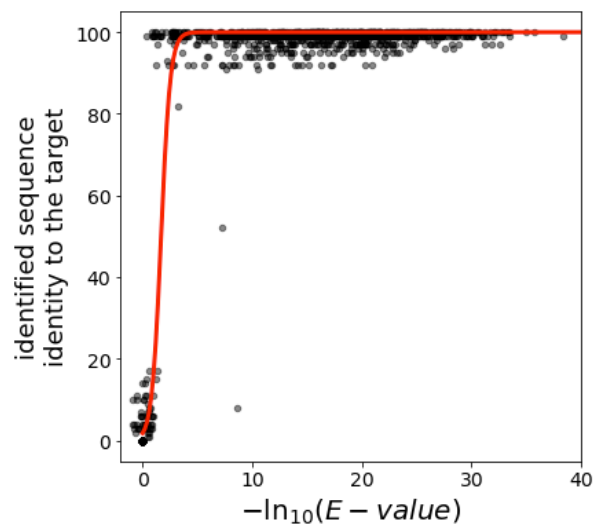
model rebuilt
wout sequence



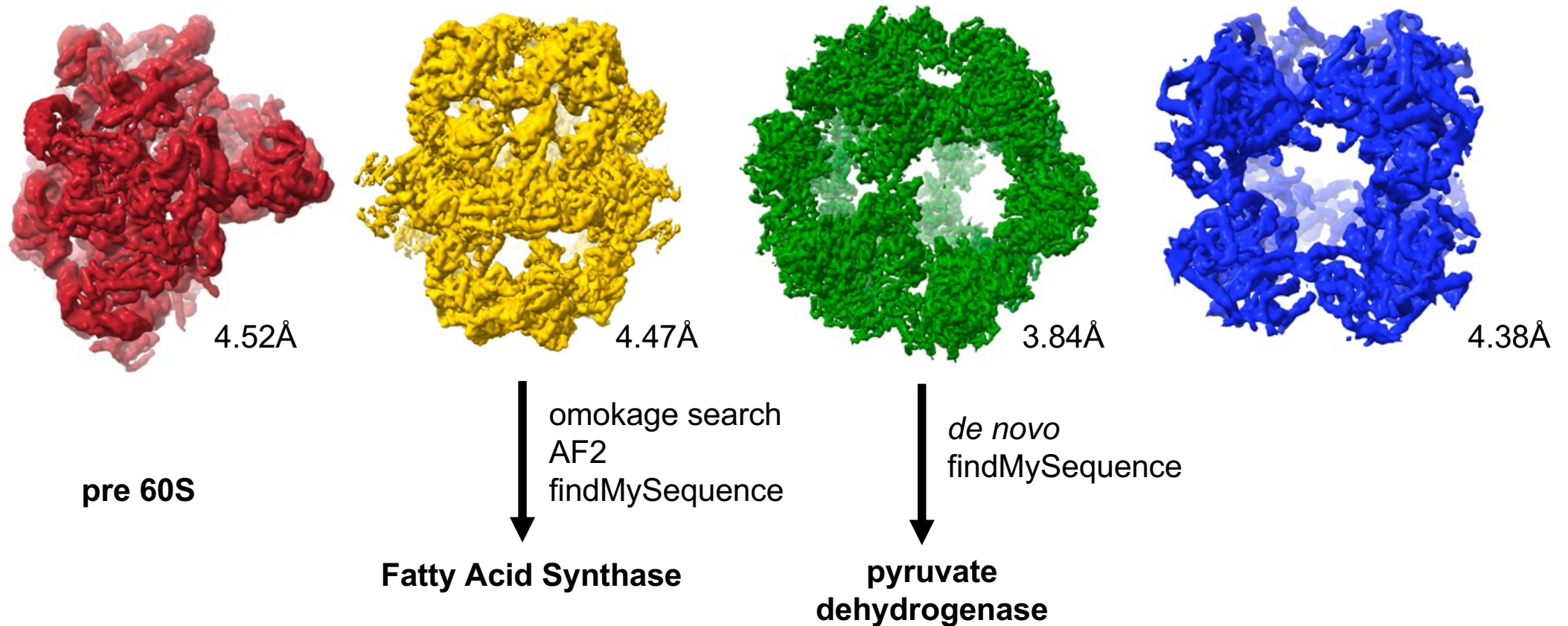
final model



uniprot id	hmmsearch E-value	sequence identity to the top hit [%]
A0A1L8D5Z7	1.40e-28	100
A0A1L8D611	2.10e-23	86
A0A1L8D605	1.00e-17	61
A0A1L8D602	9.30e-16	56
A0A1L8D6C4	3.00e-15	52
A0A1L8D5Z4	4.10e-12	50

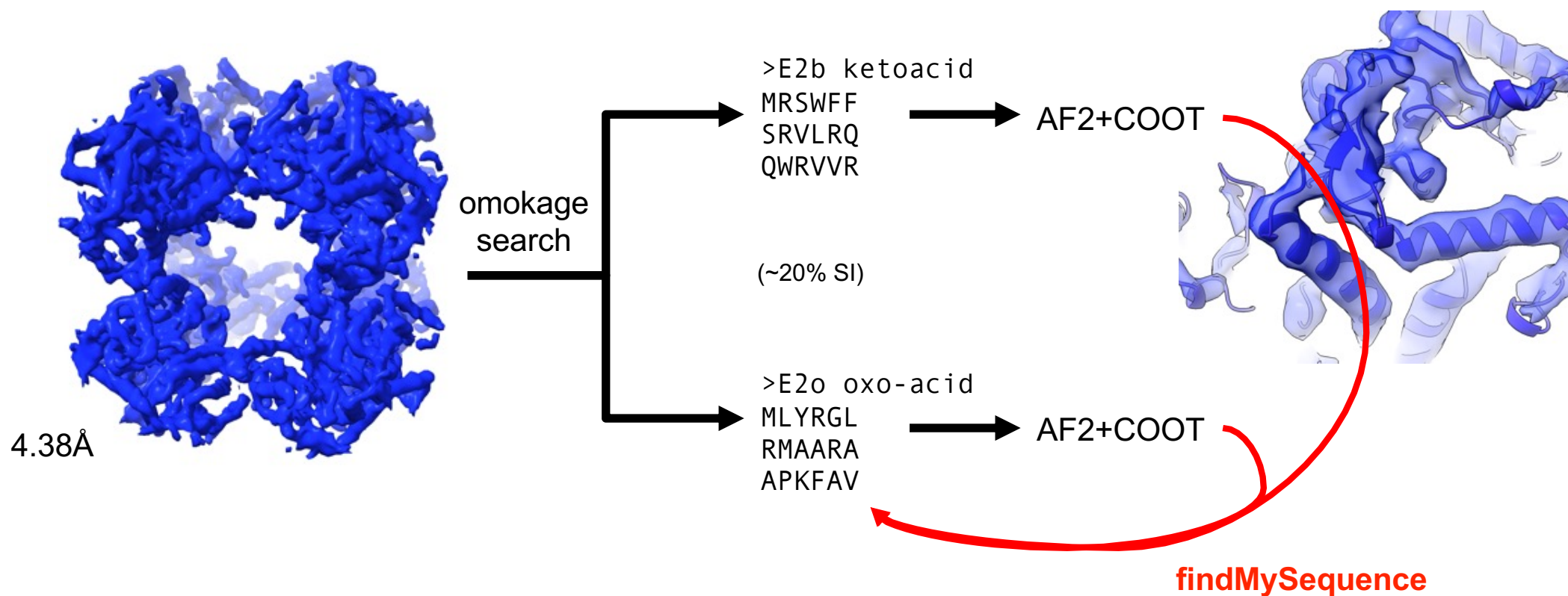


cryo-EM of *C. thermophilum* native cell extracts

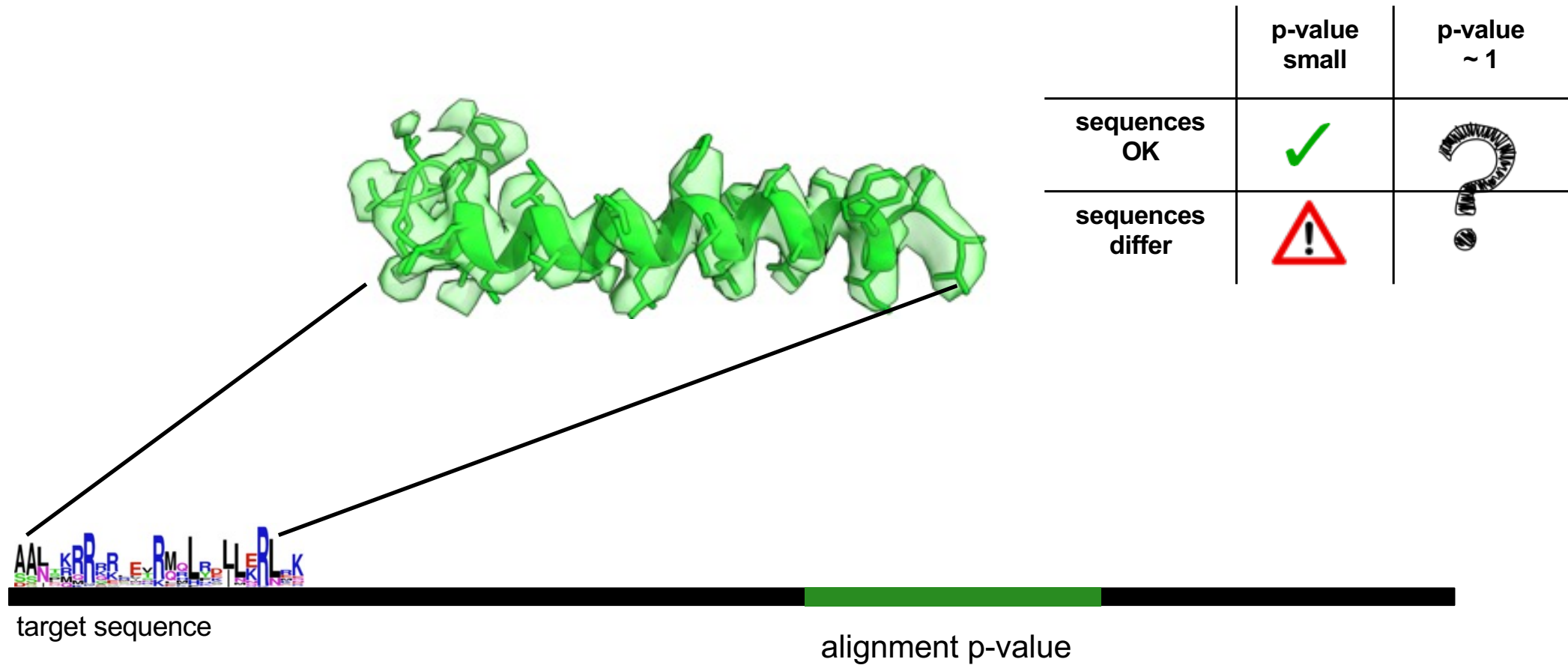


with Panos Kastitis and Ioannis Skolidis
Skolidis et al Structure 2022

two dehydrogenases in *C. thermophilum* native cell extracts

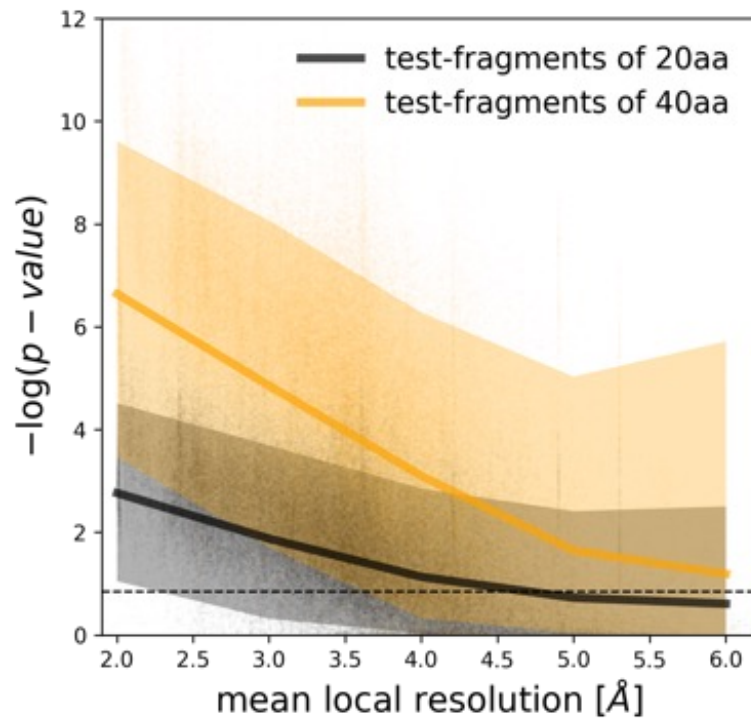


Sequence re-assignment as a model validation

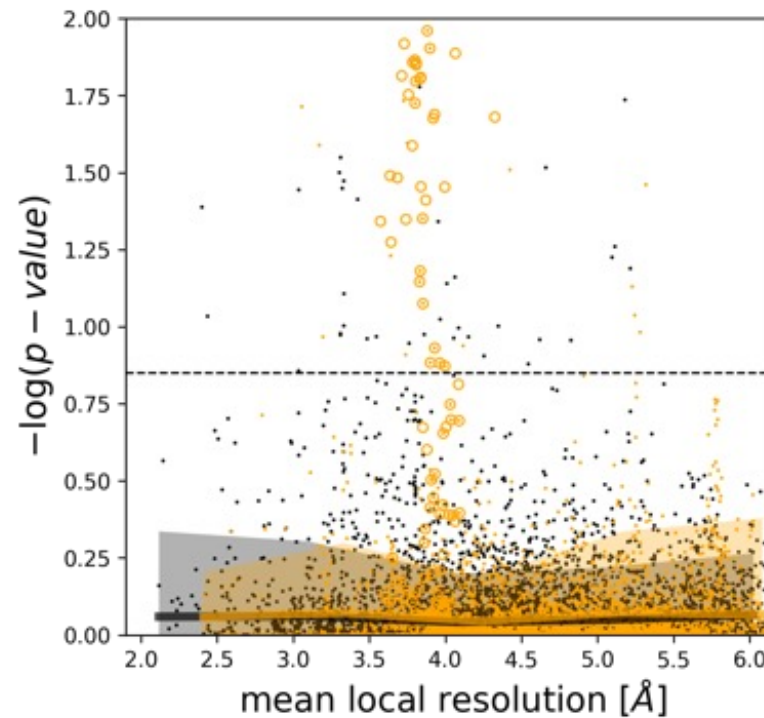


Is the sequence assignment p-value a reliable score?

sequence OK



sequences differ

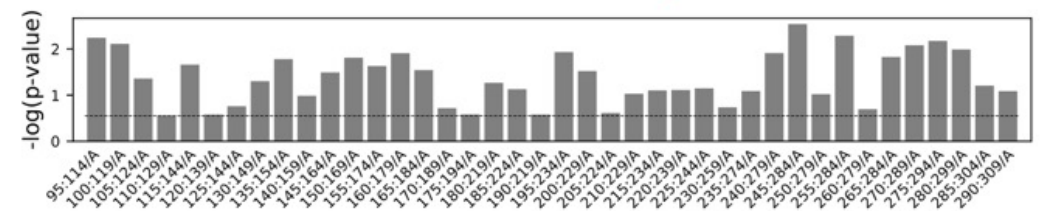
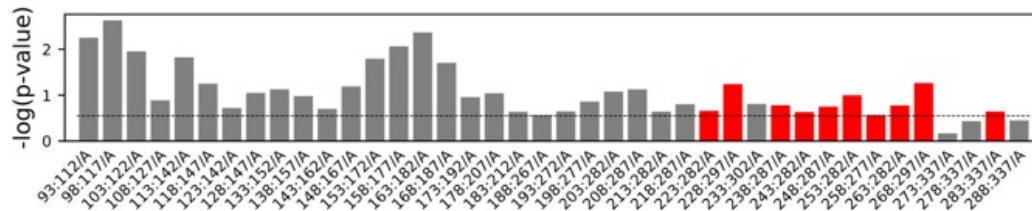
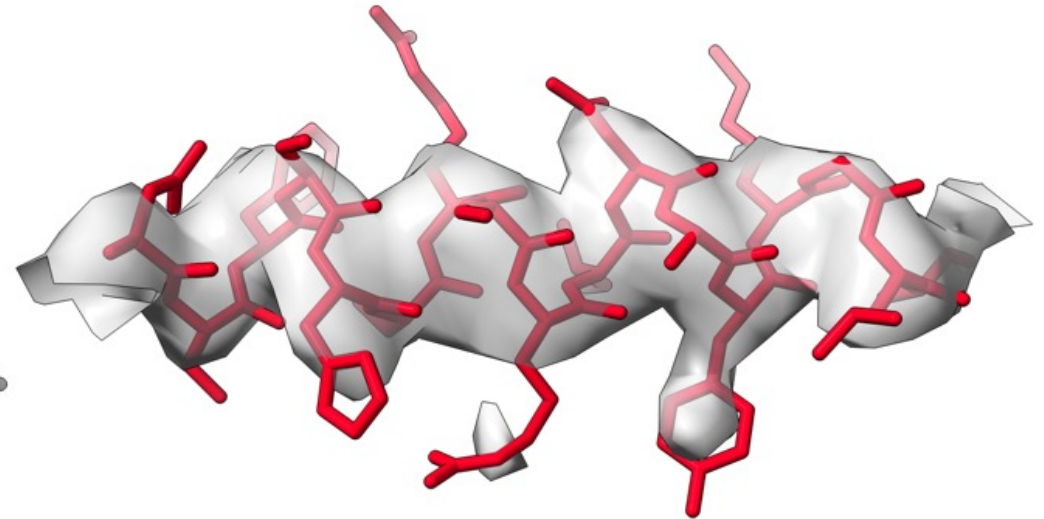
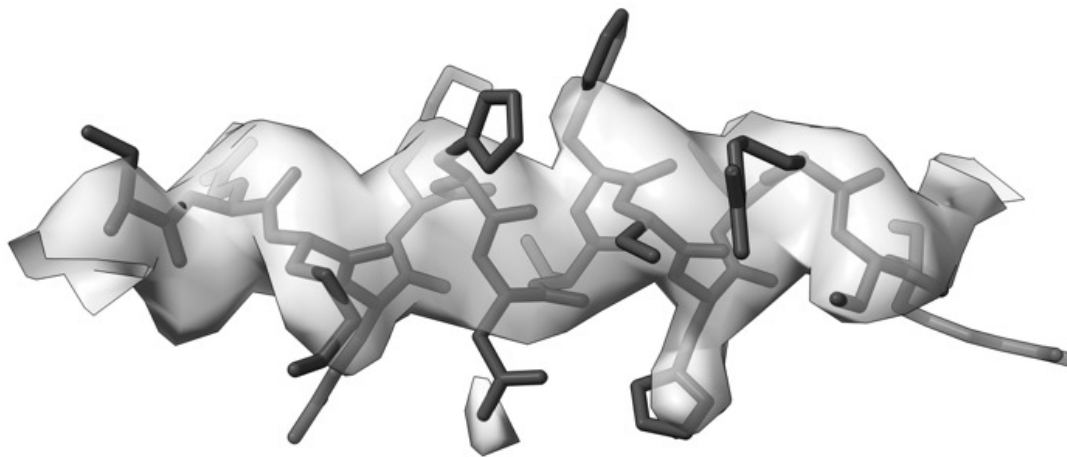


	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?

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

hidden errors in EM models: finding a better hypothesis

AlphaFold2 model plddt>90



cytoplasmic domain of a cation channel at 3.8Å resolution

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

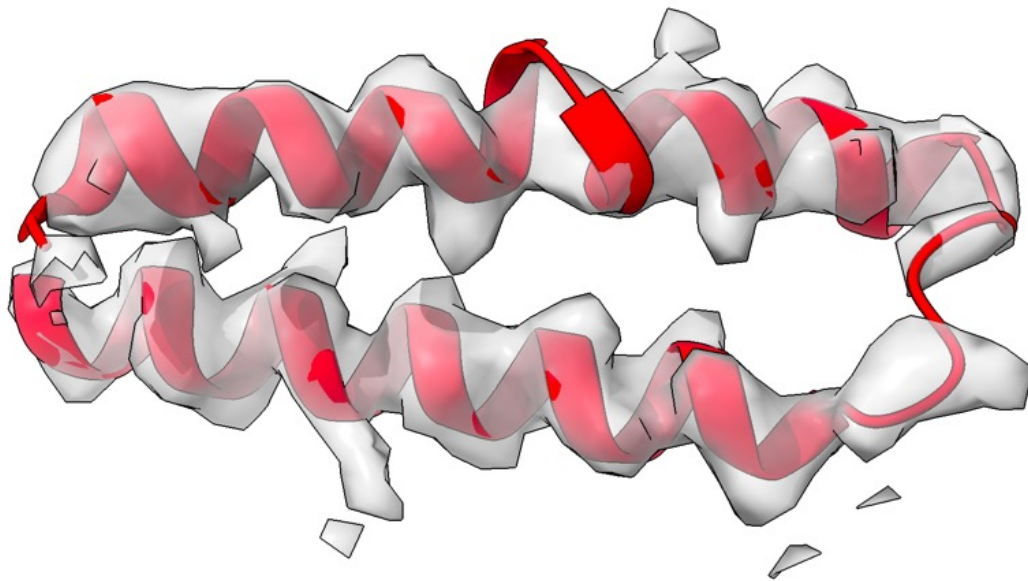
- Fragment N/5-24 has a **chain break at N/11 and no residue indexing gap**, check!

```
model      -----VELTEEE-LYISKKNLLFKRF-----
           ||||| |||||
refseq      AAAAAAMPVELTEEEKLYISKKNLLFKRFVEPGRLCLIE
```

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

```
model seq 356-395
  sknkkeKRVQKQIQKKELQKINHYYKGVAKAVKKKKKREEKKAKskkktanqavi
new seq 360-399
  sknkkekrvQKQIQKKELQKINHYYKGVAKAVKKKKKREEKKAKSKKTanqavi
```

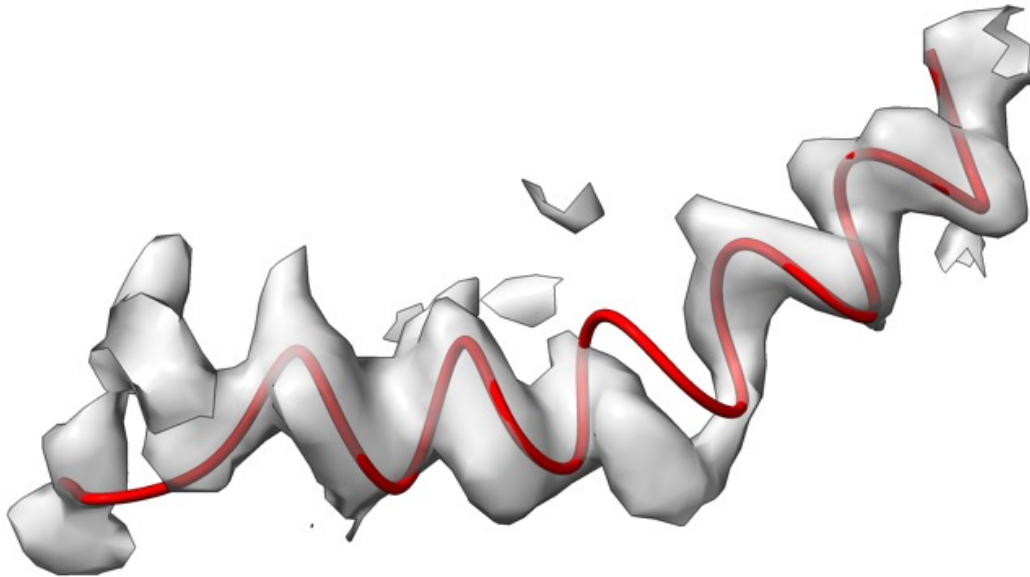
checkMySequence for 6,000 cryo-EM models at 5Å or better



- Fragment y/1-20 is shifted by -1 residue [-log(p-value)=2.95]
model seq 1-20
MKAKELREKSVEELNTELLNllreqfn
new seq 2-21
mKAKELREKSVEELNTELLNllreqfn

70S ribosomal protein 2.9 Å

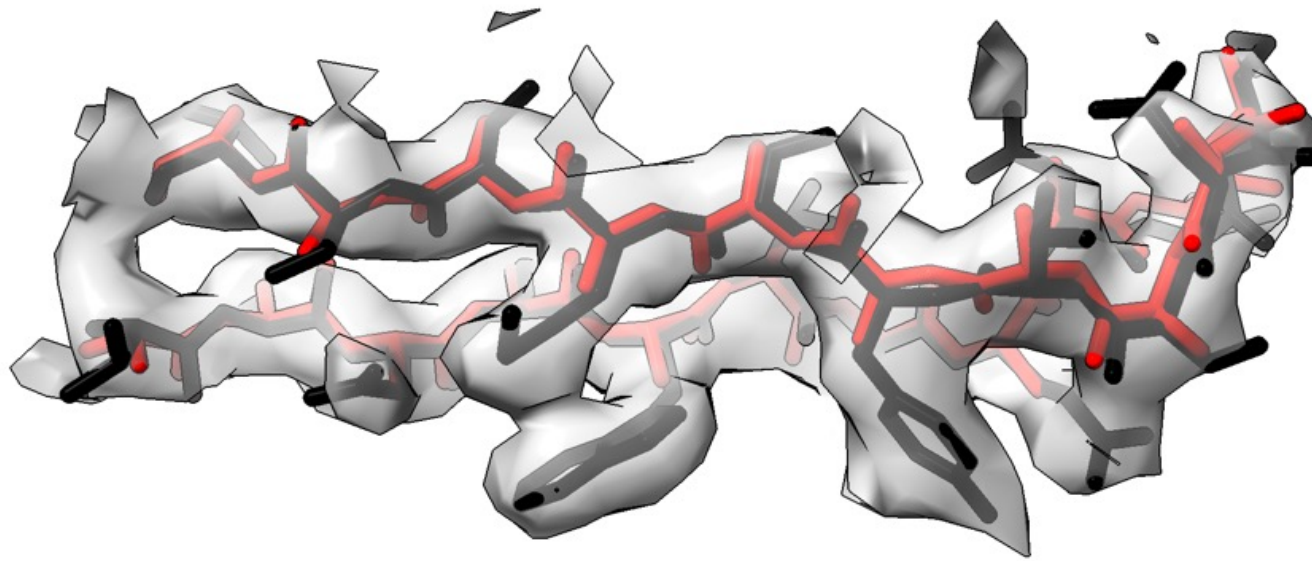
checkMySequence for 6,000 cryo-EM models at 5Å or better



- Fragment u/43-62 is shifted by -3 residues [-log(p-value)=2.40]
model seq 43-62
ekptt**ERKRAKASAVKRHAKKLARE**narrrt
new seq 46-65
ekptterk**RAKASAVKRHAKKLARENAR**rt

70S ribosomal protein at 3 Å

checkMySequence for 6,000 cryo-EM models at 5Å or better

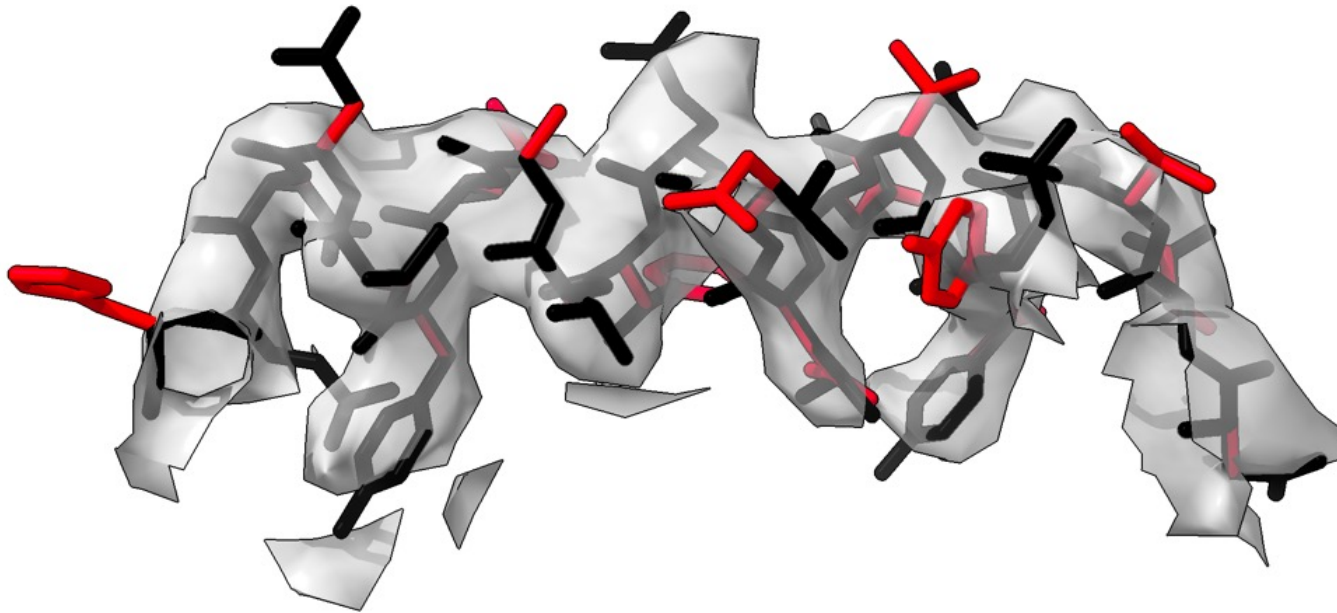


deposited model
rebuilt with findMySequence

- Fragment C/161-200 is shifted by 2 residue [-log(p-value)=3.11]
model seq 161-200
lnnfypk**DINV**KWKIDG**SER**QNGVLNSWTDQDSK**DSTYSMSSTLTL**kdeyer
new seq 159-198
lnnfy**PKDINV**KWKIDG**SER**QNGVLNSWTDQDSK**DSTYSMSSTLTL**kdeyer

ferroportin at 3 Å

checkMySequence for 6,000 cryo-EM models at 5Å or better

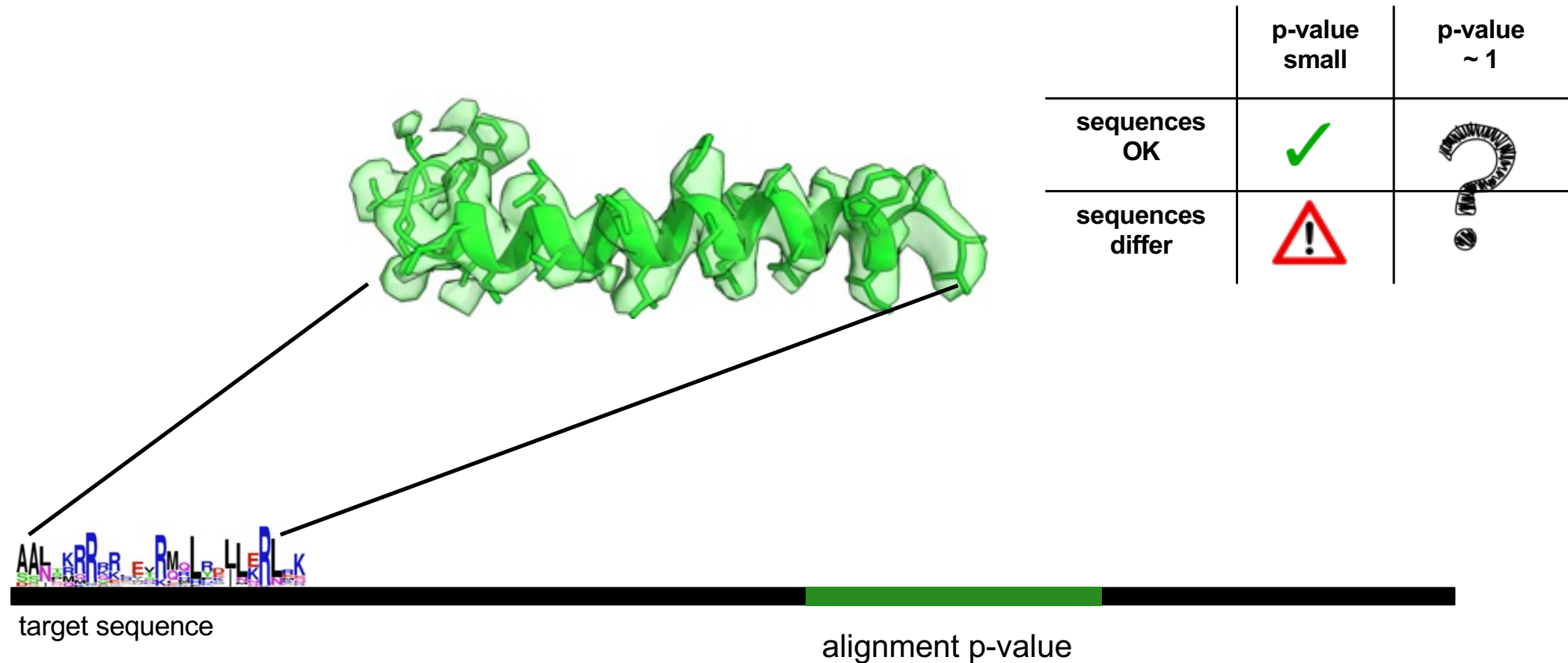


deposited model
rebuilt with findMySequence

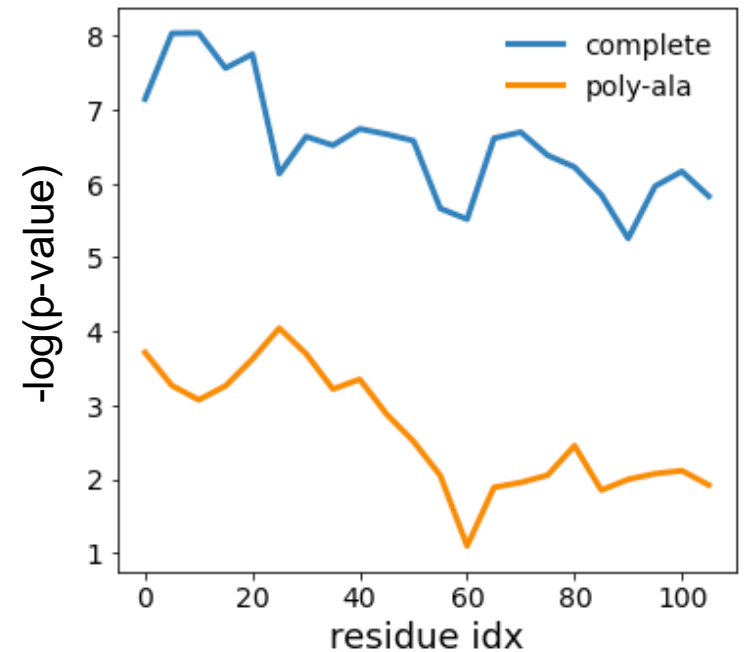
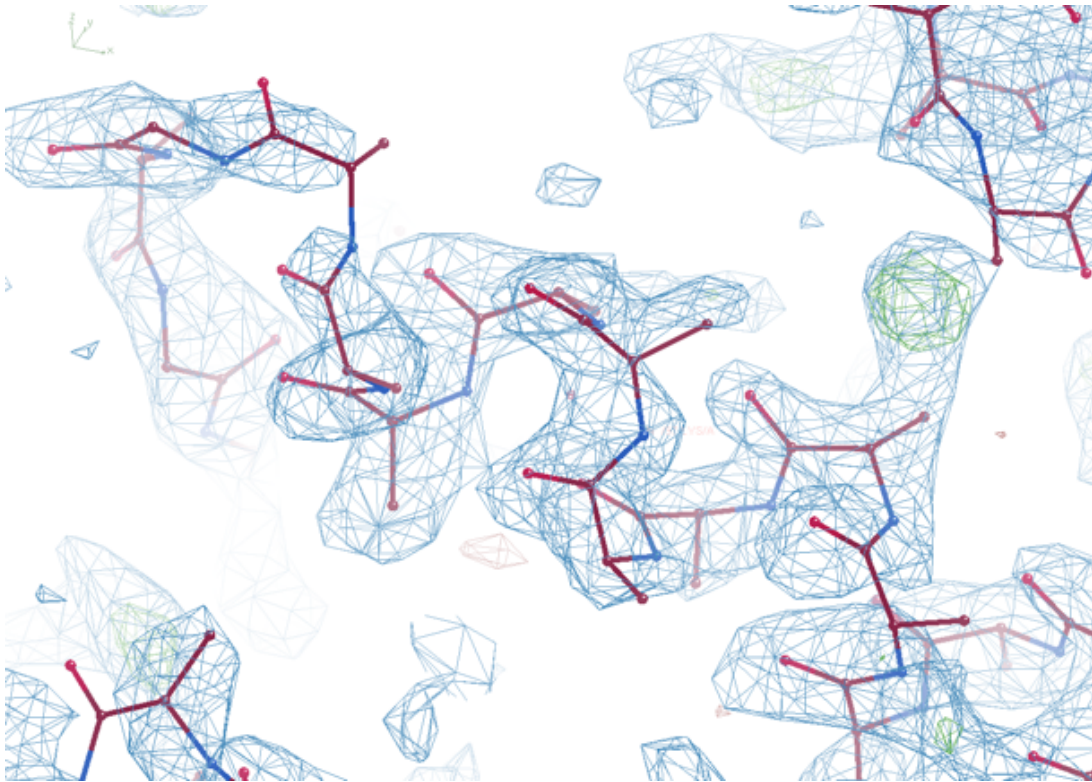
- Fragment C/145-204 is shifted by 10 residue [-log(p-value)=1.94]
model seq 145-204
gfaelarryahnlan**ARFLWRNRVGAEAVEVRINHIRQGEVARAWRFDALAIGLRDFKADAELDALAELIASGLS**gsgghvl
new seq 135-194
gfael**ARRYAHNLANARFLWRNRVGAEAVEVRINHIRQGEVARAWRFDALAIGLRDFKADAELDA**laeliasglsghvl

CRISPR-associated protein at 3.8 Å

checkMySequence and protein crystal structures – model bias

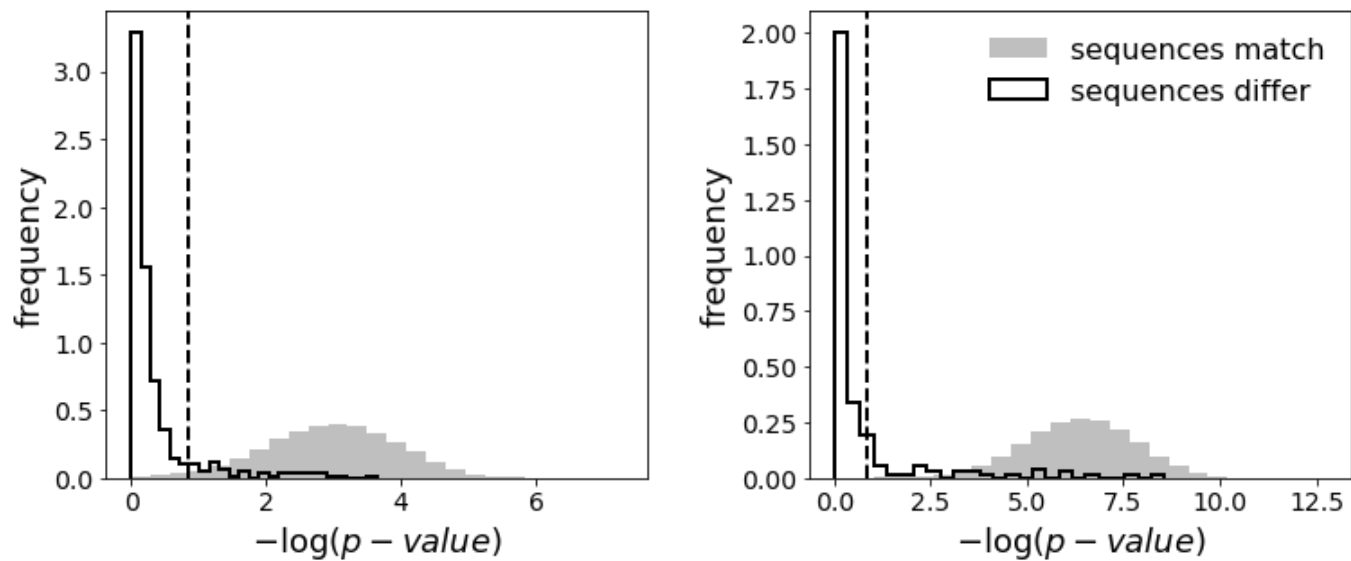


Sequence assignment in MX – model bias



hen egg-white lysozyme @2.9Å (poly-ala + Refmac5, 4gce)

Sequence assignment in MX – model bias



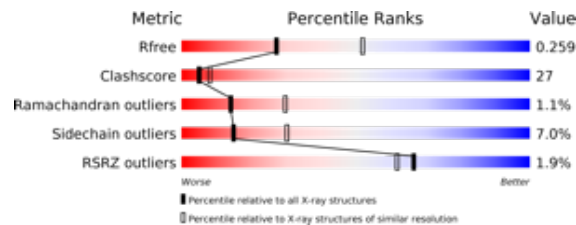
	p-value small	p-value ~ 1
sequences OK		
sequences differ		

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

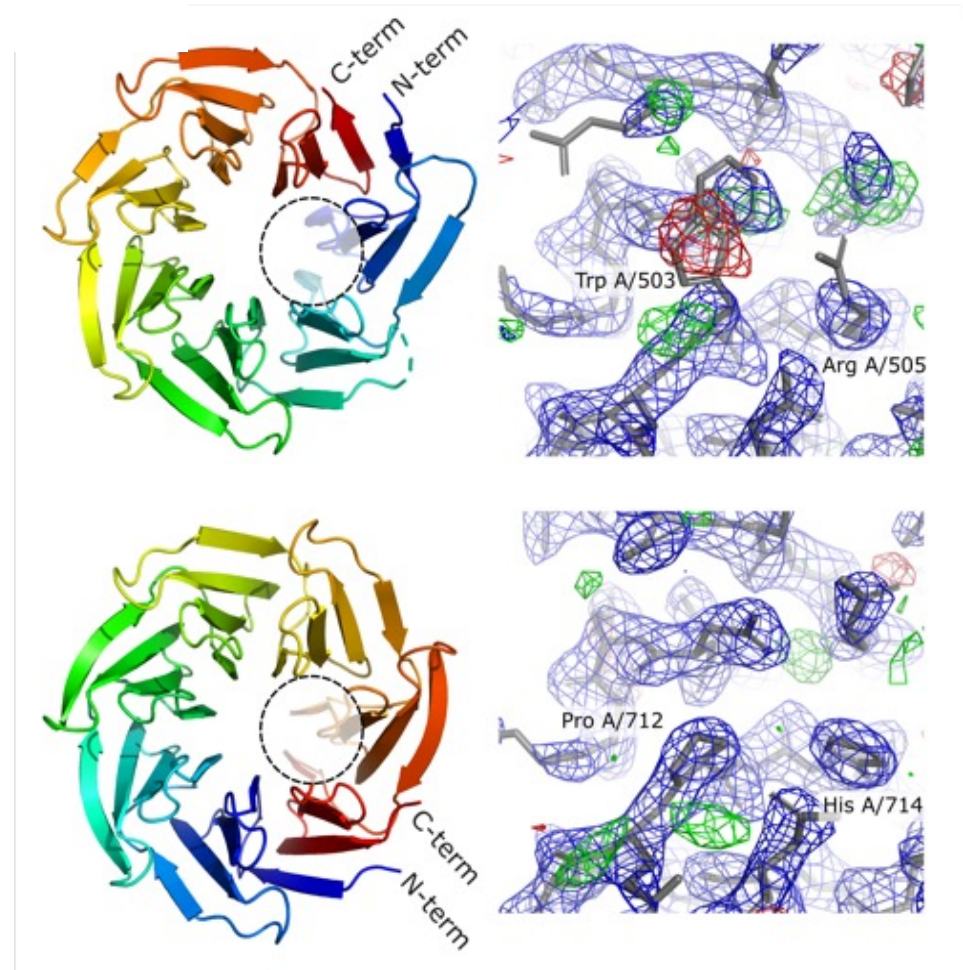
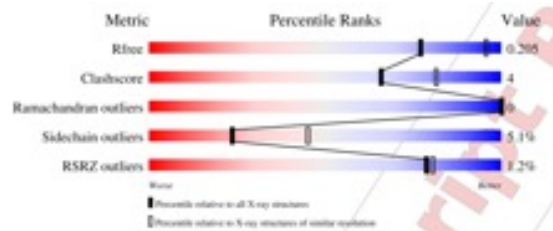
Register shifts in protein crystal structures – example 1

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

deposited
R/Rfree 22/26

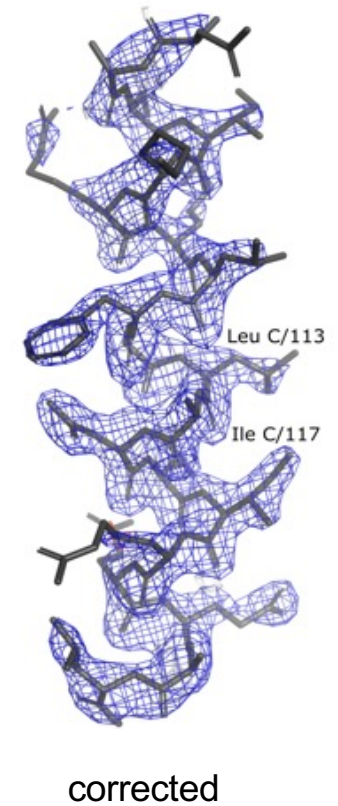
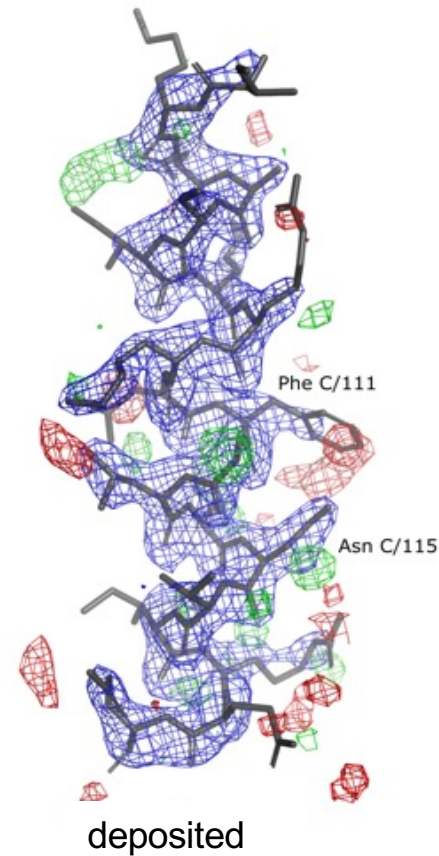
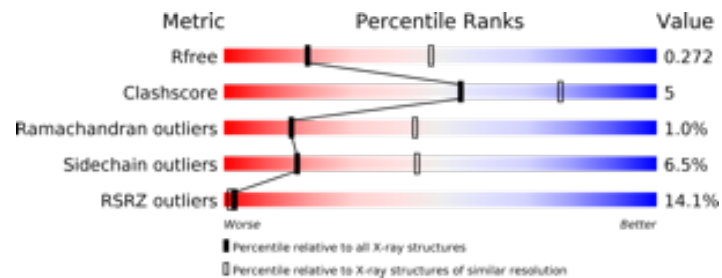


corrected
R/Rfree 17/21



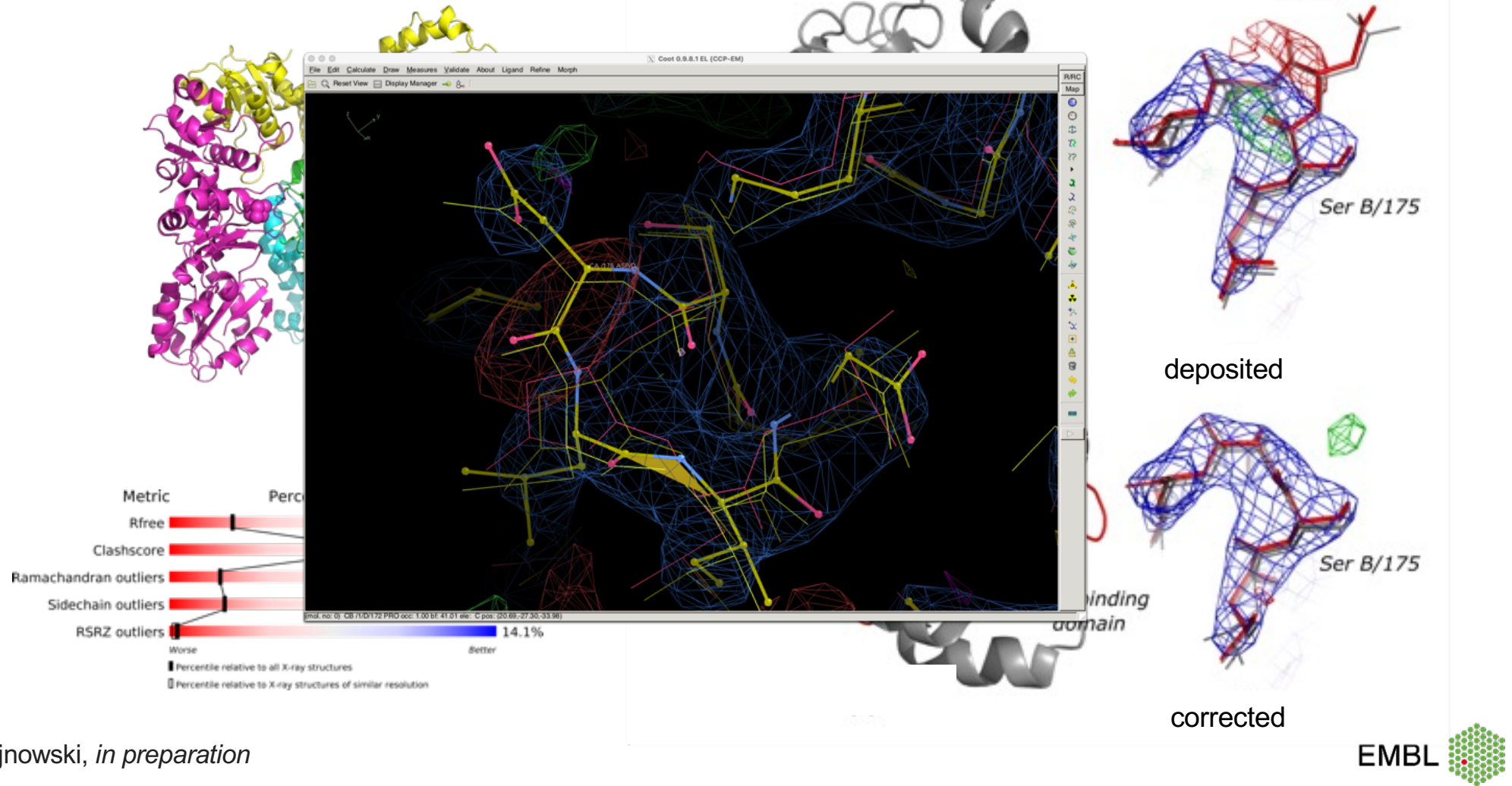
Register shifts in protein crystal structures – example 2

Helicase form *H. pylori* @2.5Å



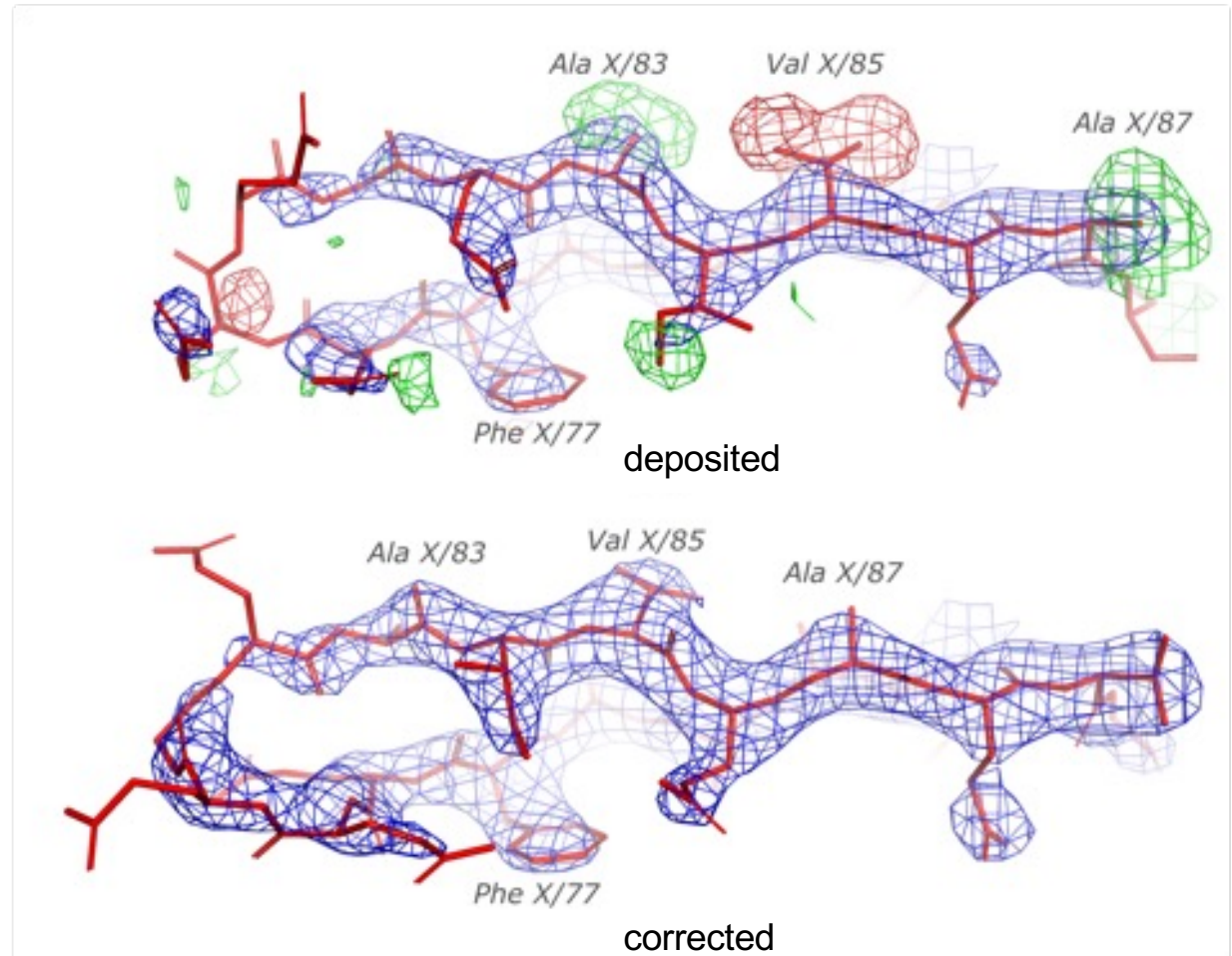
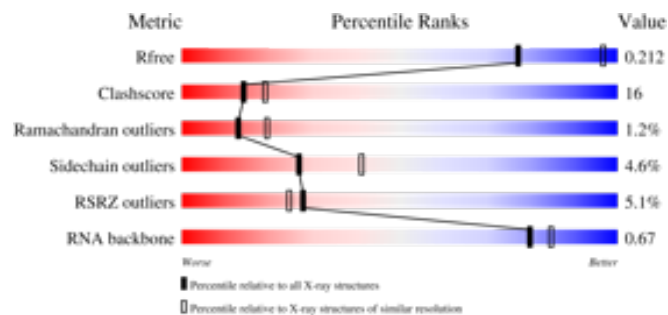
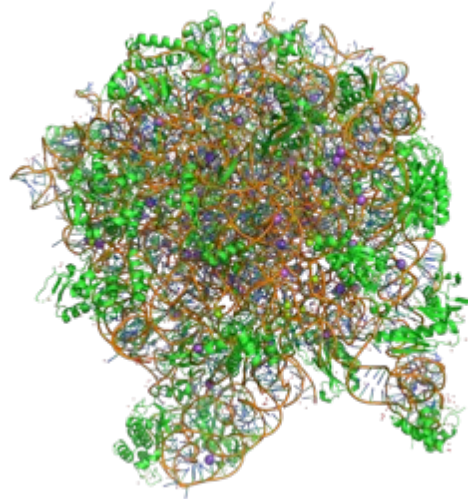
Register shifts in protein crystal structures – example 3

hydrogenase from *T. melanesiensis* @2.8Å



Register shifts in protein crystal structures – example 4

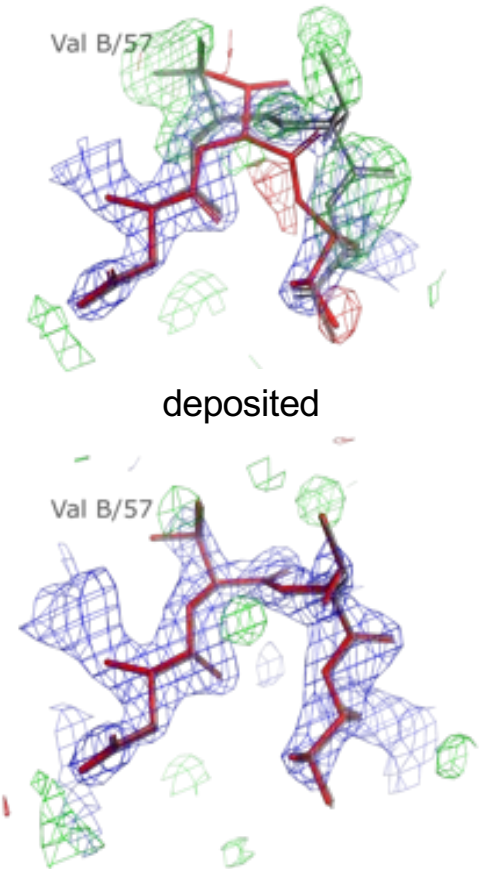
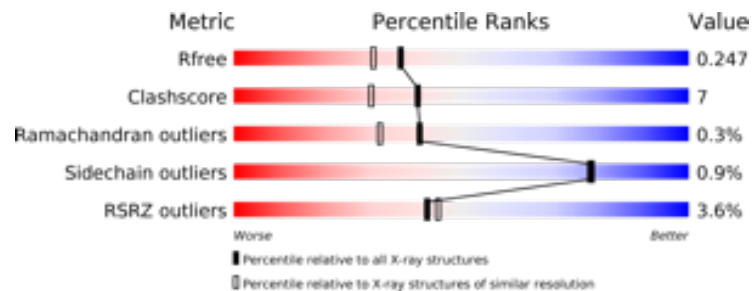
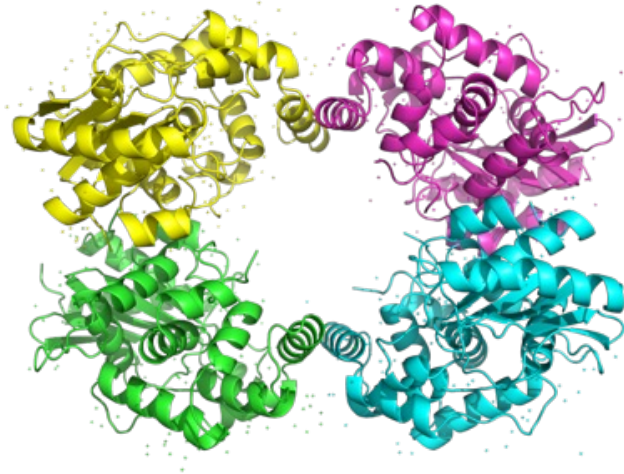
Ribosomal protein L31e form *H. Marismortui* @2.65Å



Chojnowski, *in preparation*

Register shifts in protein crystal structures – example 5

glutaminase from *G. kaustophilus* @2.1Å



Chojnowski, *in preparation*

corrected

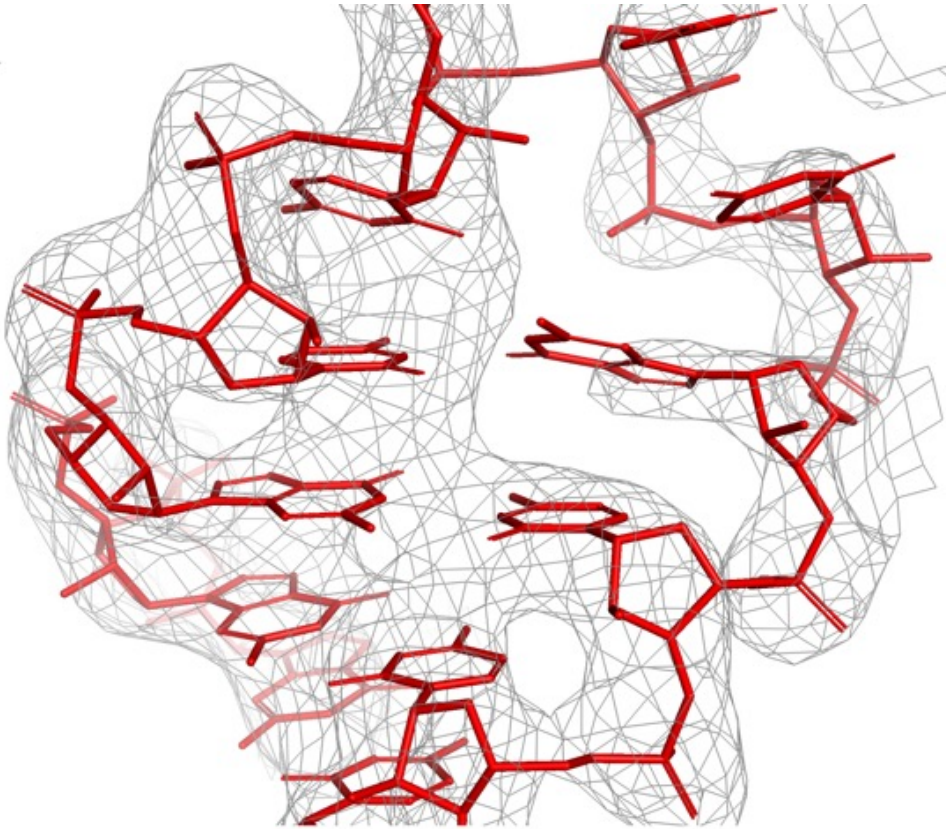


Register shifts in protein crystal structures

Statistics from a scan of 10,000 deposited protein crystal structures at 2-3Å resolution

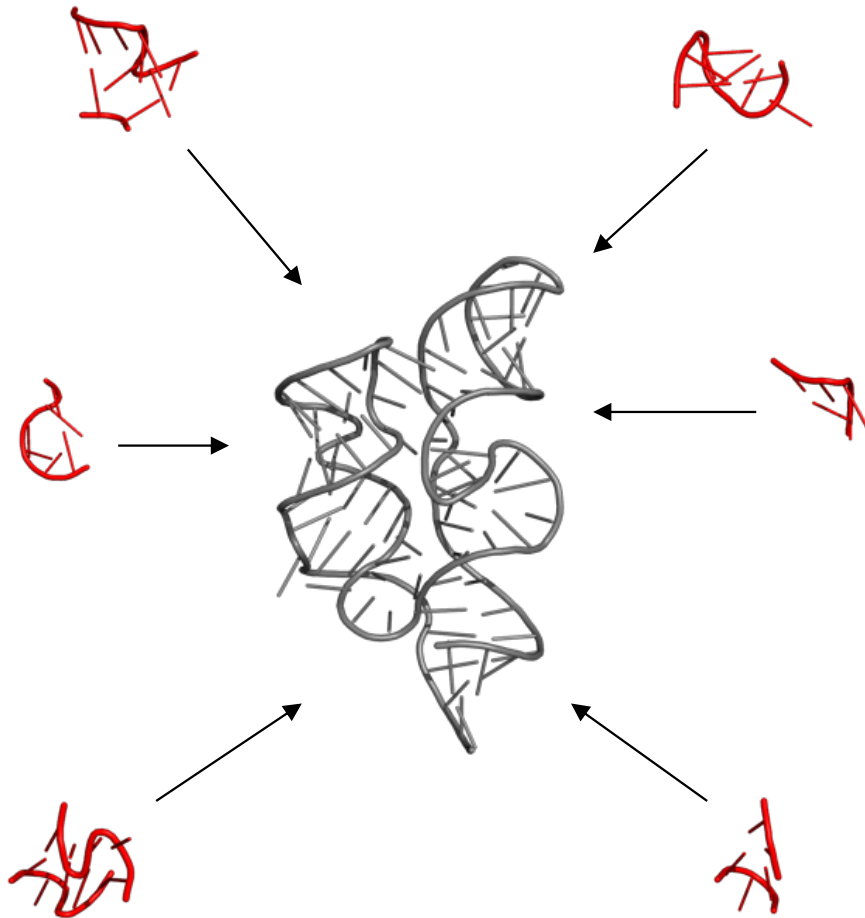
- 27 models with indexing issues (residue numbering ignores gaps)
- 89 structures with sequence mismatches
- 115 with poorly resolved chains (at least 1)
- 70 plausible register shifts

Nucleic acid sequence assignment in EM/MX



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

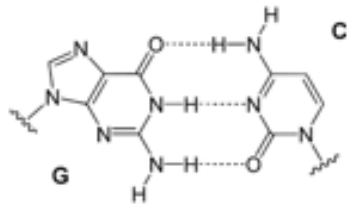
Base-pairs from backbone geometry



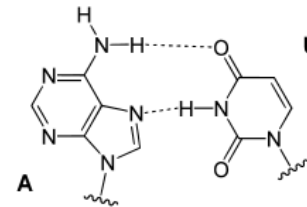
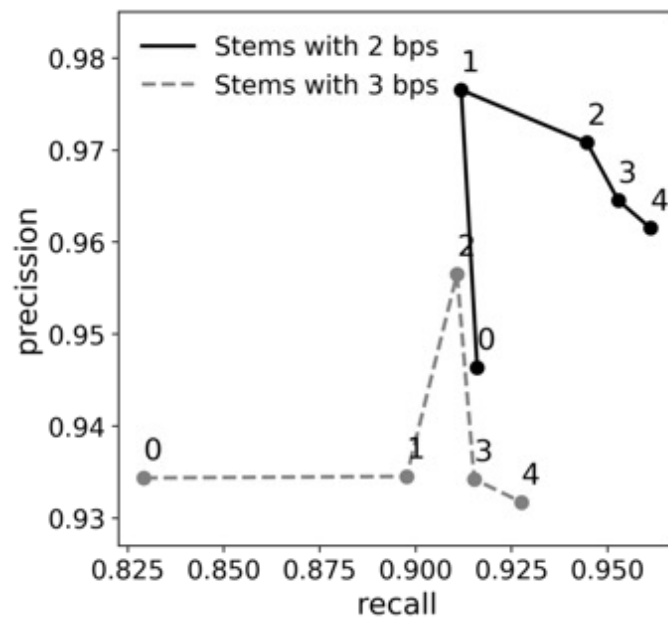
input RNA model backbone is “covered”
with motifs with known secondary structure

- base pairs assignment
- sequence validation
- refinement restraints

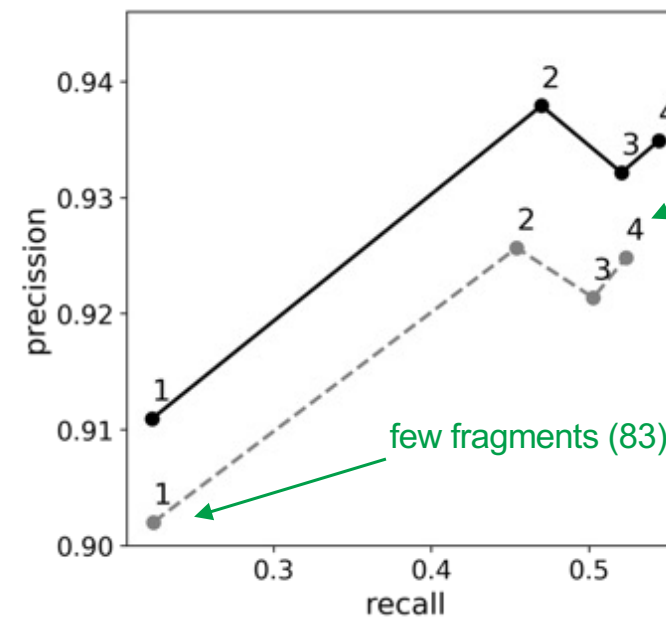
Base-pairs from backbone geometry



Watson-Crick (canonical) base-pairs



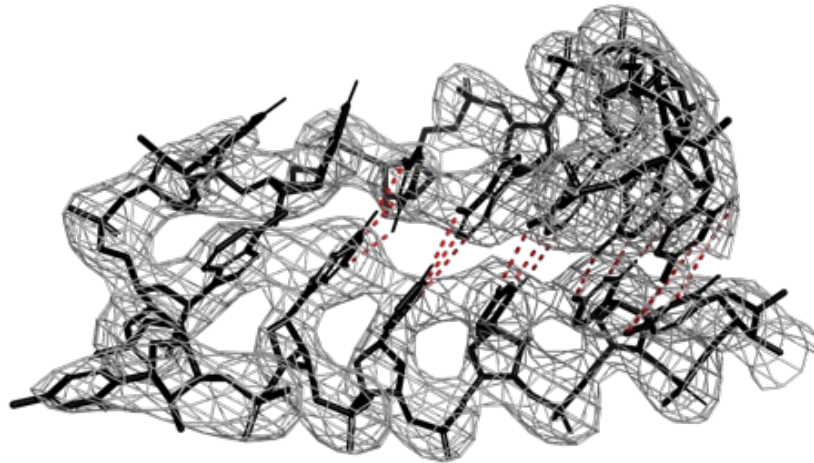
non-canonical base-pairs



lots of fragments (2,664)

few fragments (83)

checkMySequence and NA crystal structures



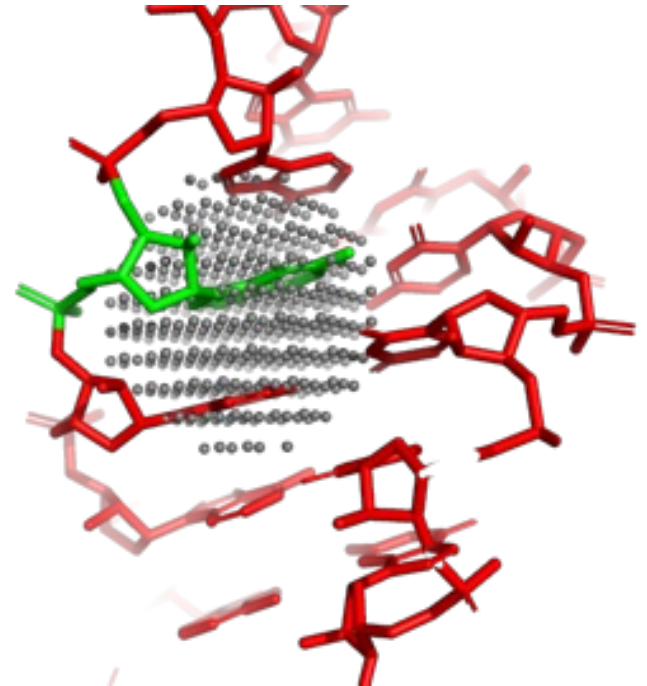
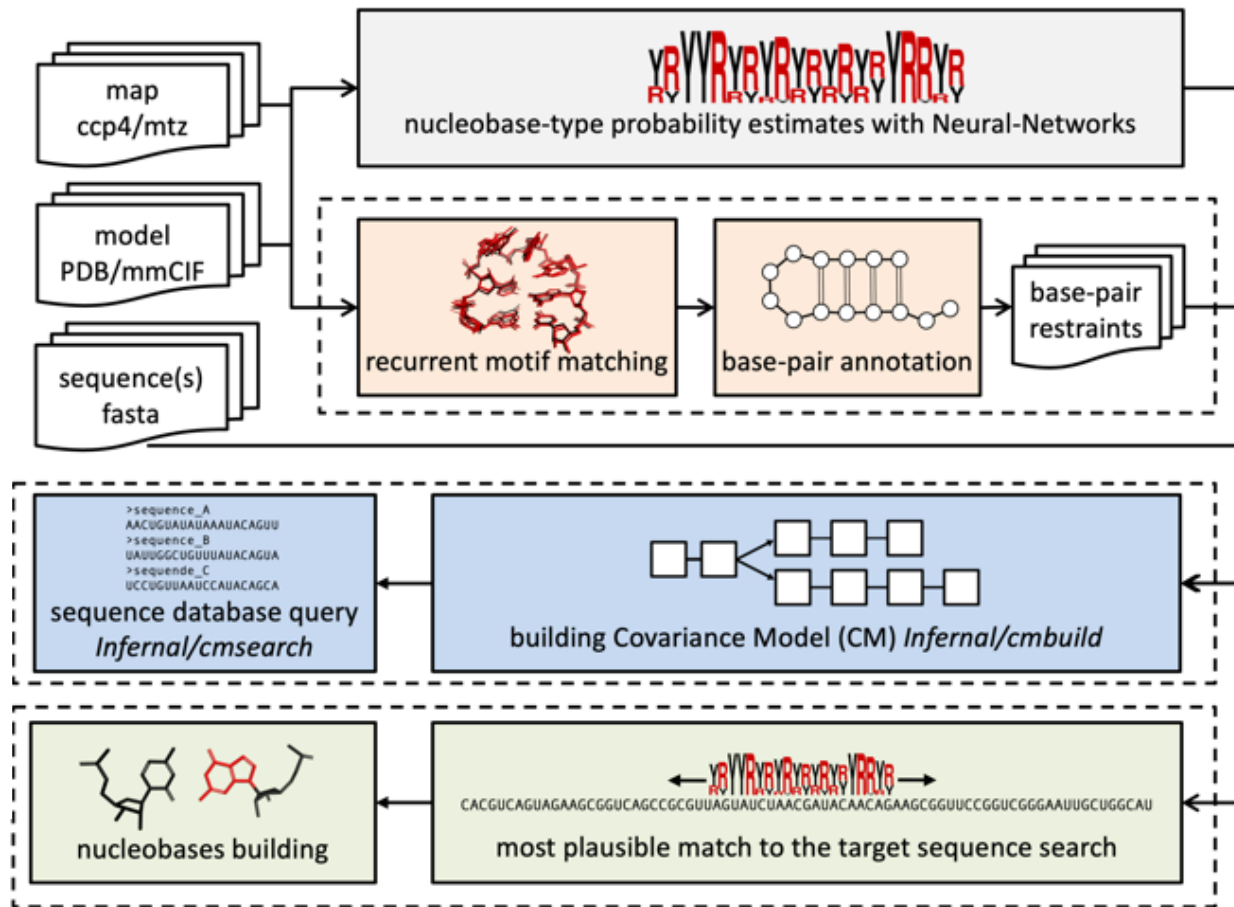
	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?



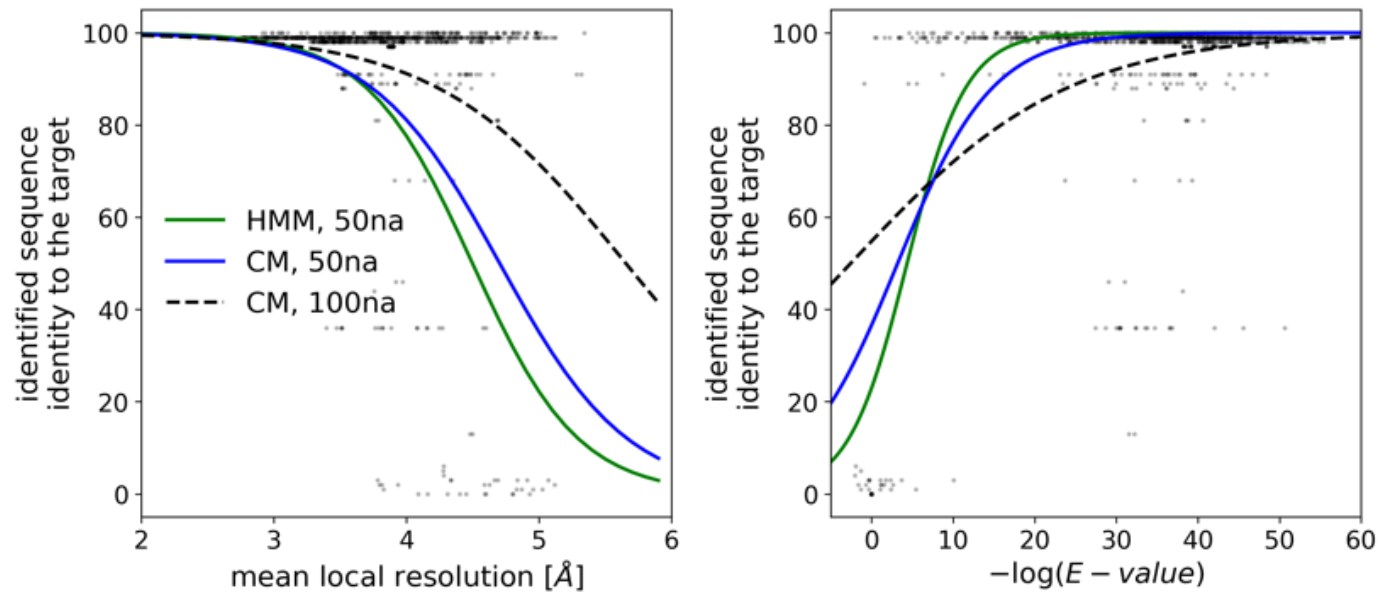
target sequence

alignment p-value

Nucleic acid sequence assignment and validation in EM/MX

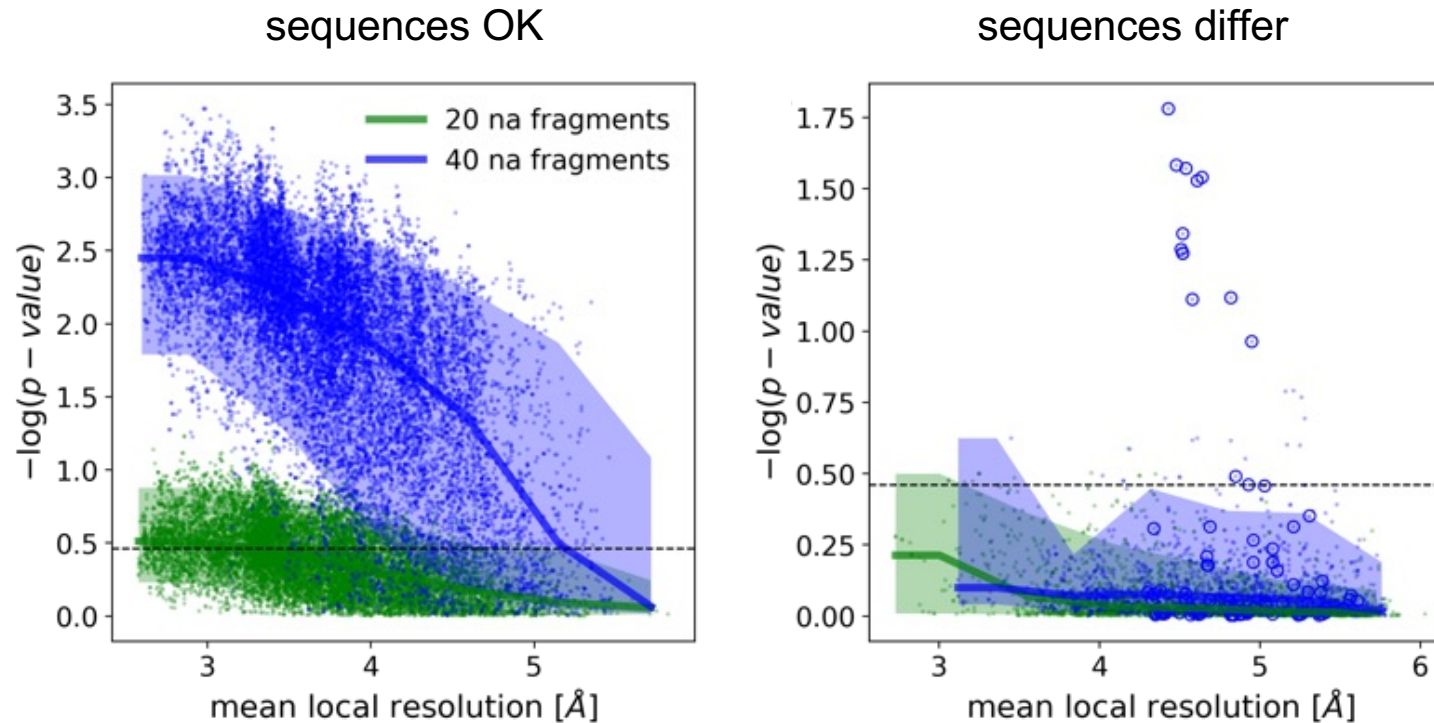


Nucleic acid sequence identification in EM structures



Sequence identification of short RNA fragments from
17 EM structures of ribosomes at 3.5Å or better

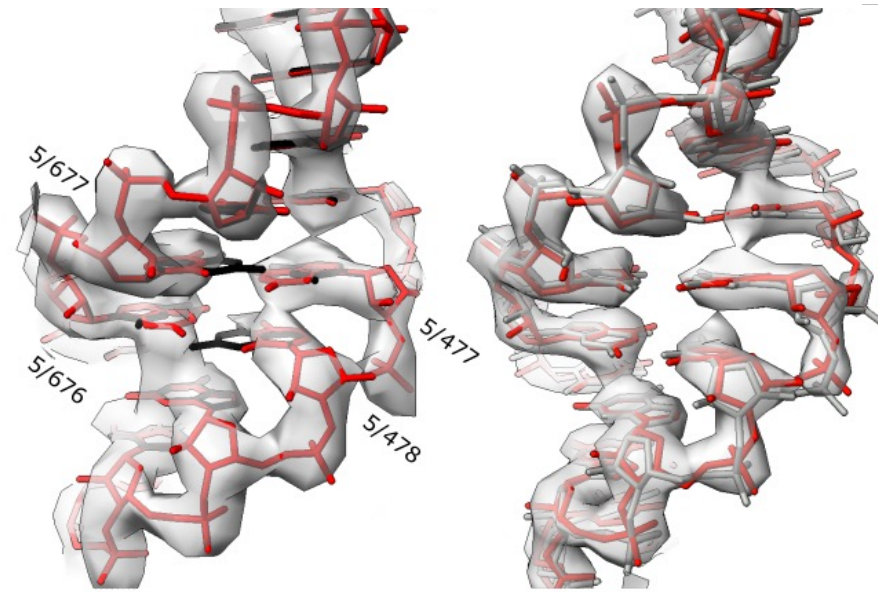
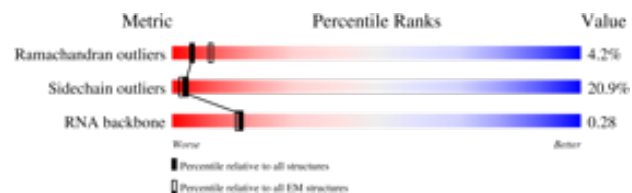
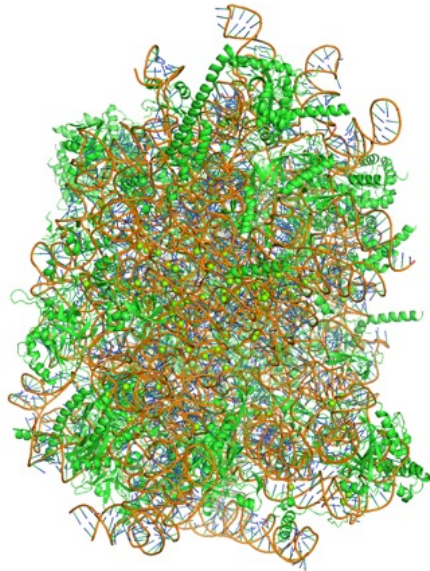
Nucleic acid sequence assignment in EM structures



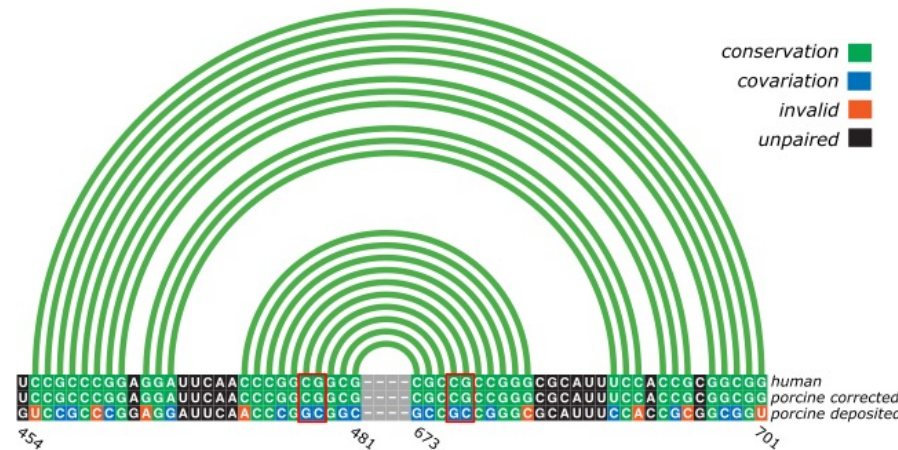
p-values for 20k RNA chain fragments re-assigned to target sequence

Sequence assignment issues in rRNA

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



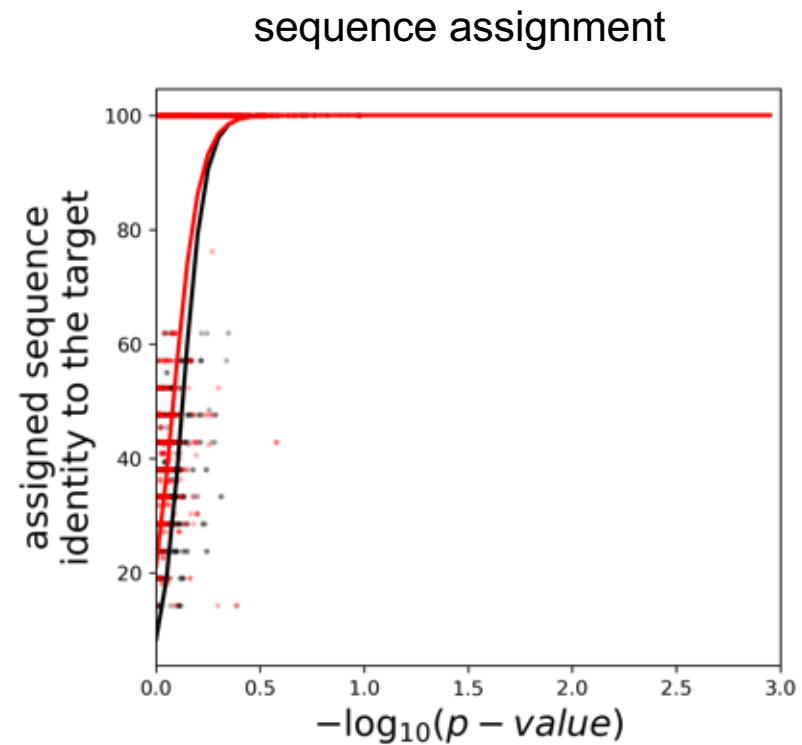
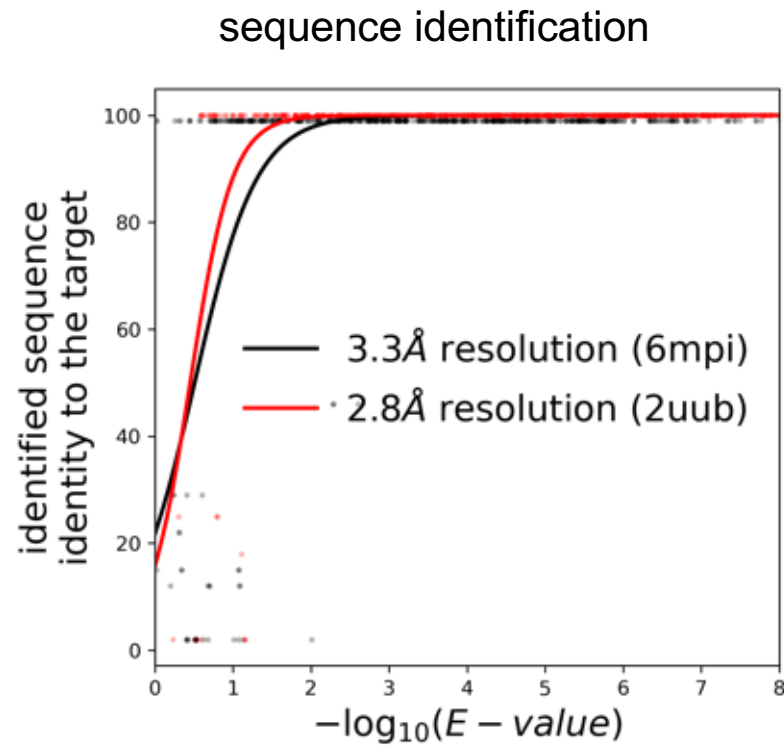
deposited
corrected
reference

Chojnowski, *in preparation*

EMBL

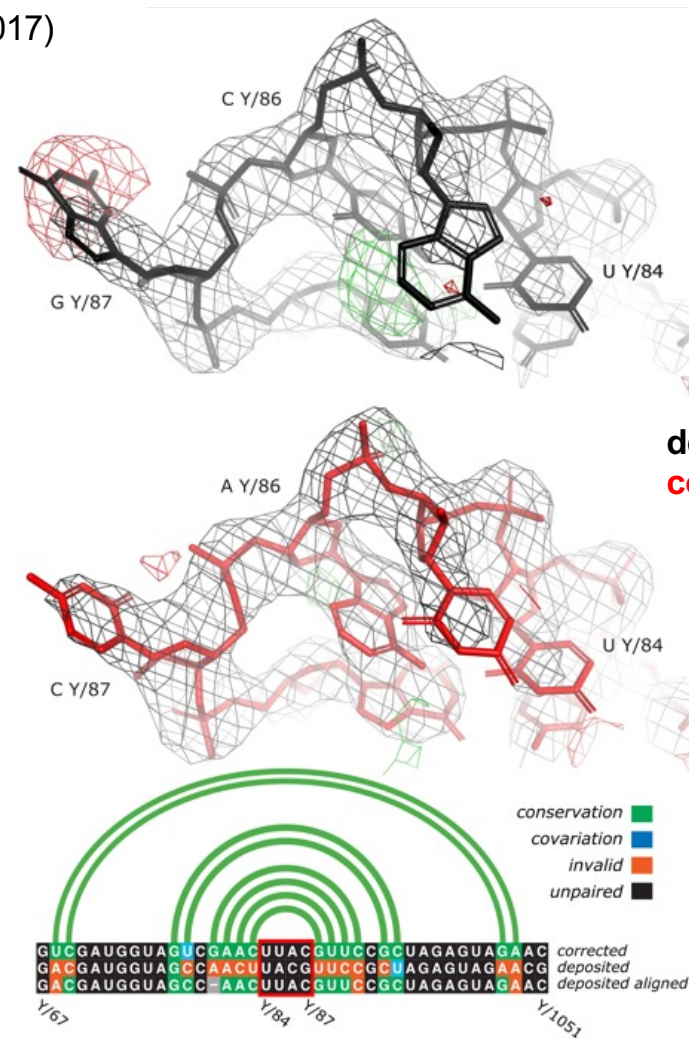
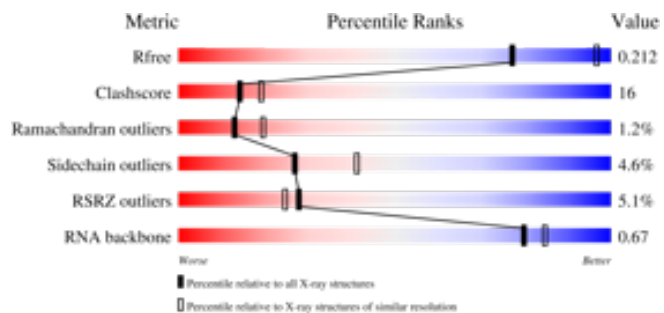
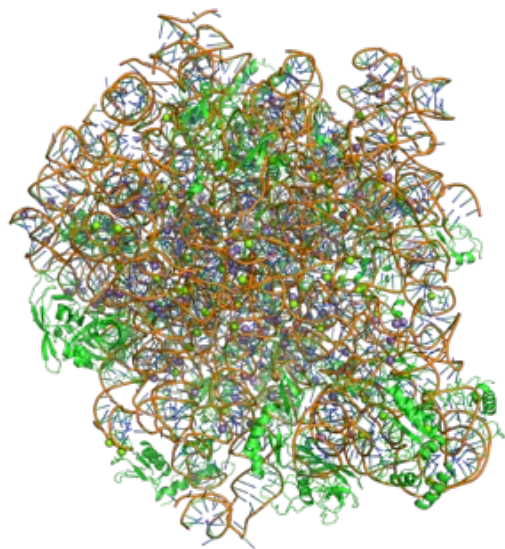


Nucleic acid sequence assignment and identification in MX



Sequence assignment issues in rRNA

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

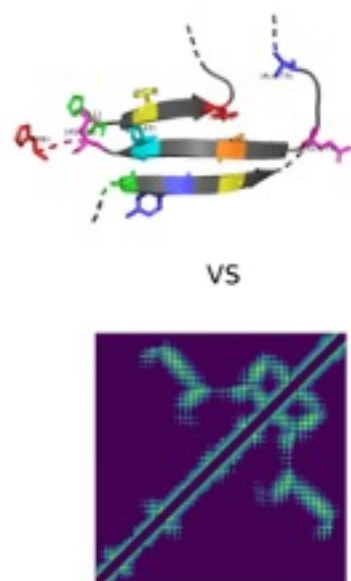


deposited
corrected

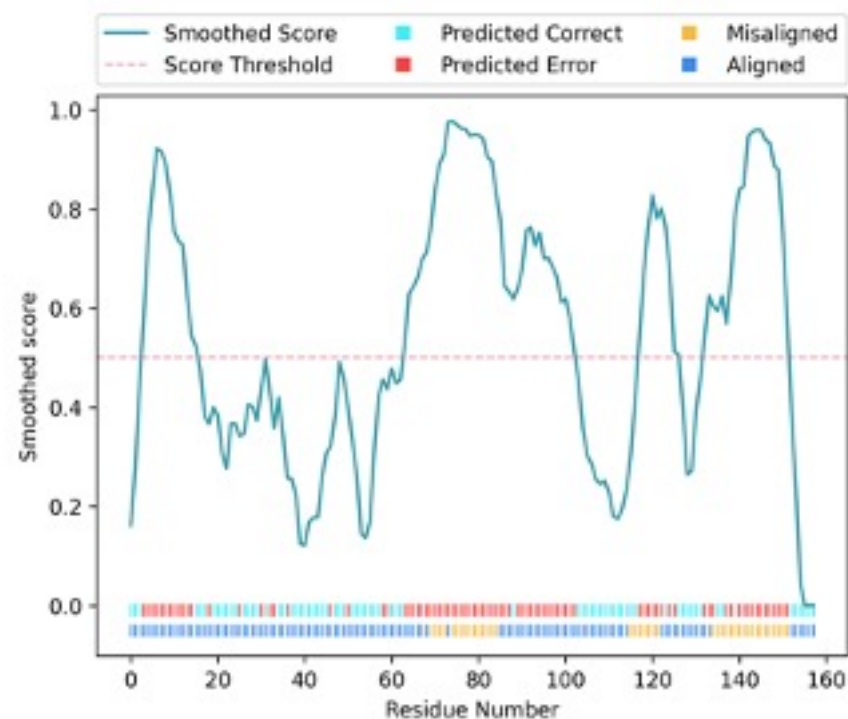
Chojnowski, *in preparation*

Predicted inter-residue distances for model validation

www.conkit.org



	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 Krissignel

São Carlos Institute of Physics

Diego A. Leonardo

EMBL Heidelberg

Wolfram Seifert-Davila

Laboratorio de Biología Molecular, Peru

Dan E. Vivas-Ruiz

Martin Luther University

Panos Kastitis

Ioannis Sklidis



Matthias Wilmanns

Kate Beckham

Jan Kosiński

Christina Ritter

Edukondalu Mullapudi

Isabel Bento

Alice Bochel

CCPEM and CCP4

workshop students!

findMySequence



checkMySequence



doubleHelix

available soon
at gitlab.com/gchojnowski

