

# Hidden errors in cryo-EM models of macromolecules

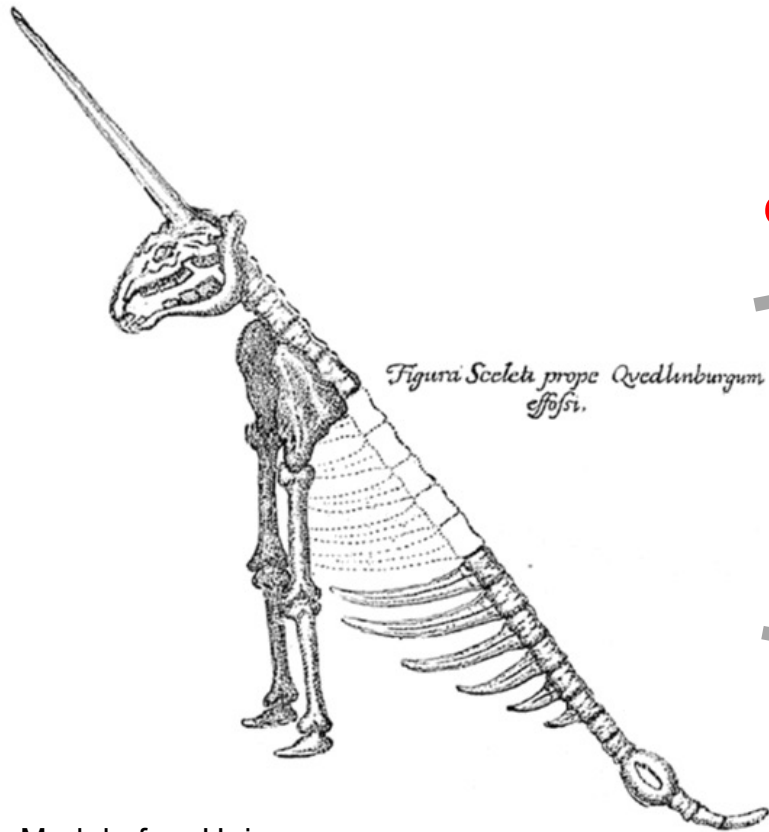
Grzegorz Chojnowski  
EMBL Hamburg

CCP-EM Icknield Model Building Workshop

09.10.2023

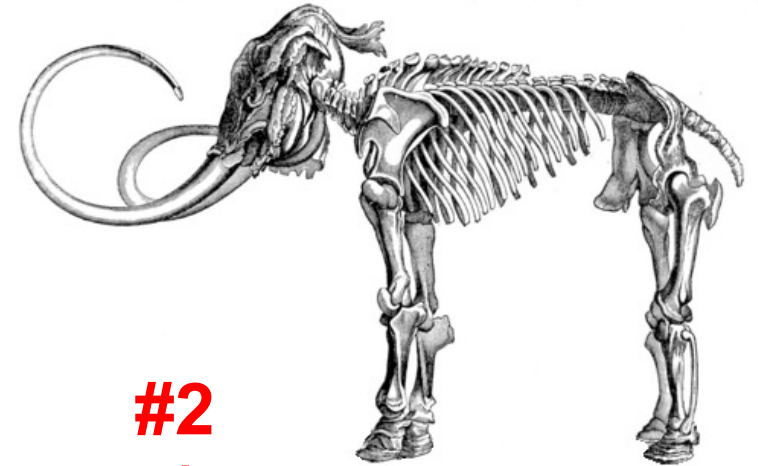


# Model building traps in EM and MX

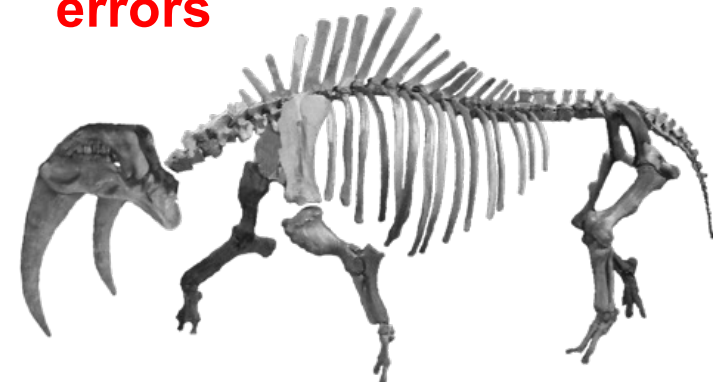


Model of an Unicorn  
Gottfried Leibniz after Otto von Guericke, *Protogaea* (1719)

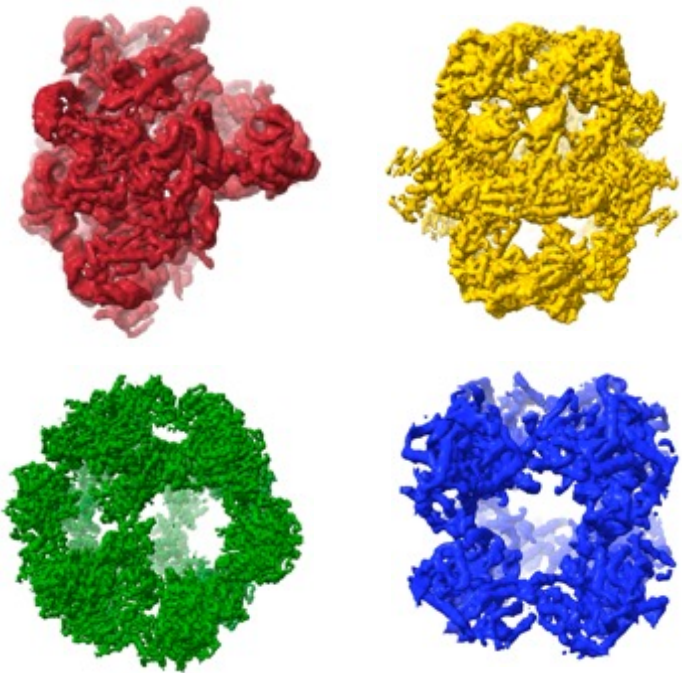
**#1**  
gross  
errors



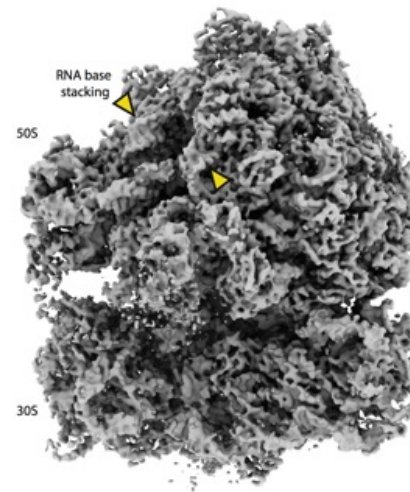
**#2**  
random  
errors



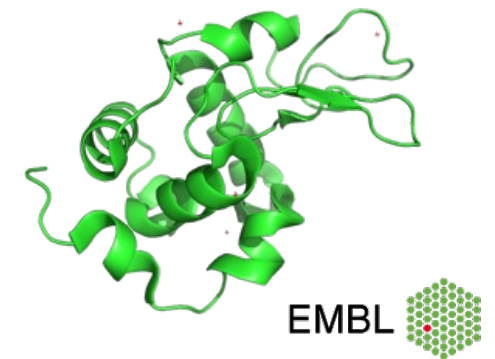
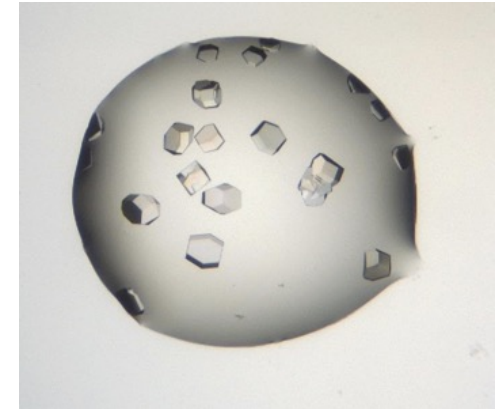
# Unknown proteins and gross errors in EM and MX



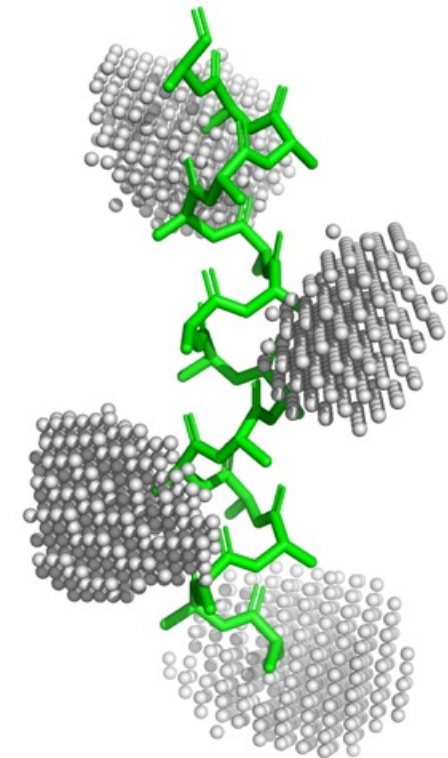
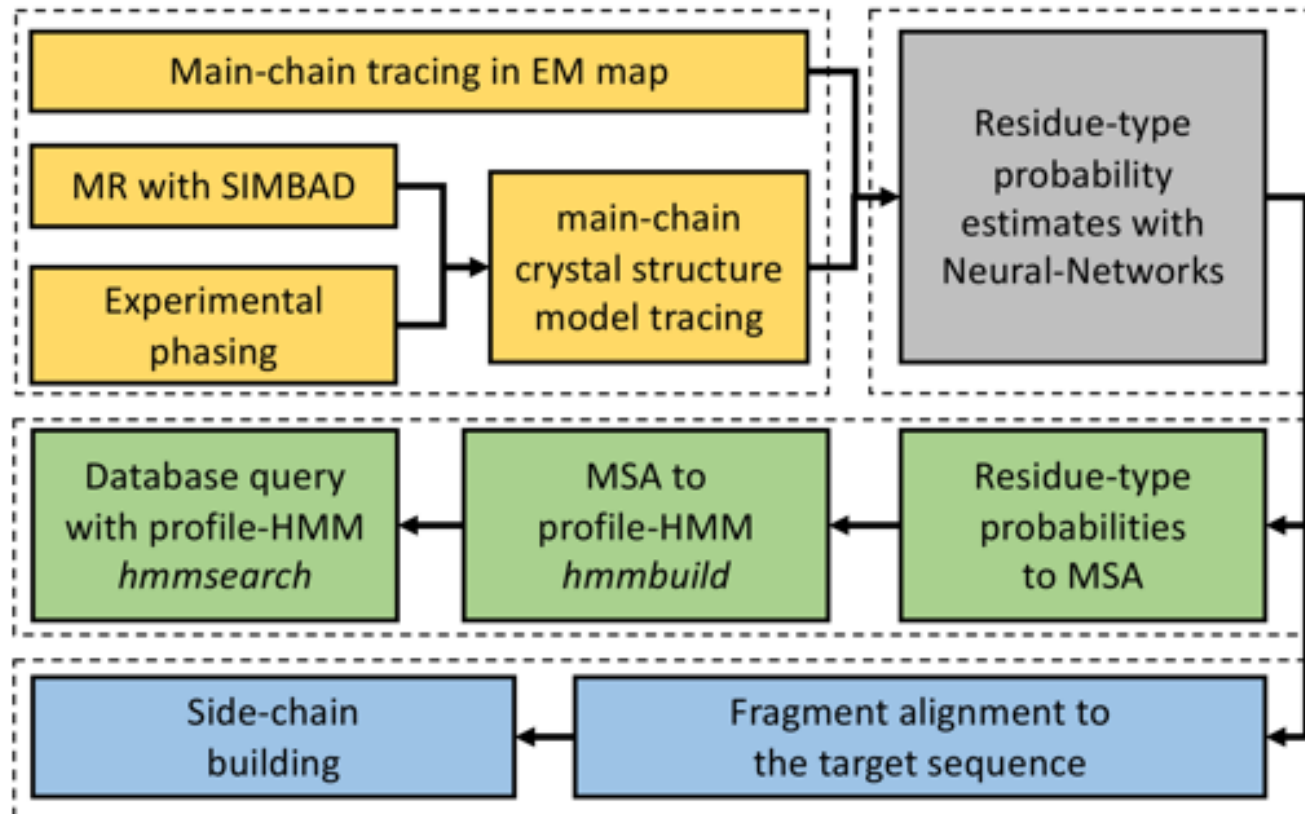
*Cryo-EM and artificial intelligence visualize endogenous protein community members*  
Skalidis et al. Structure 2022



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



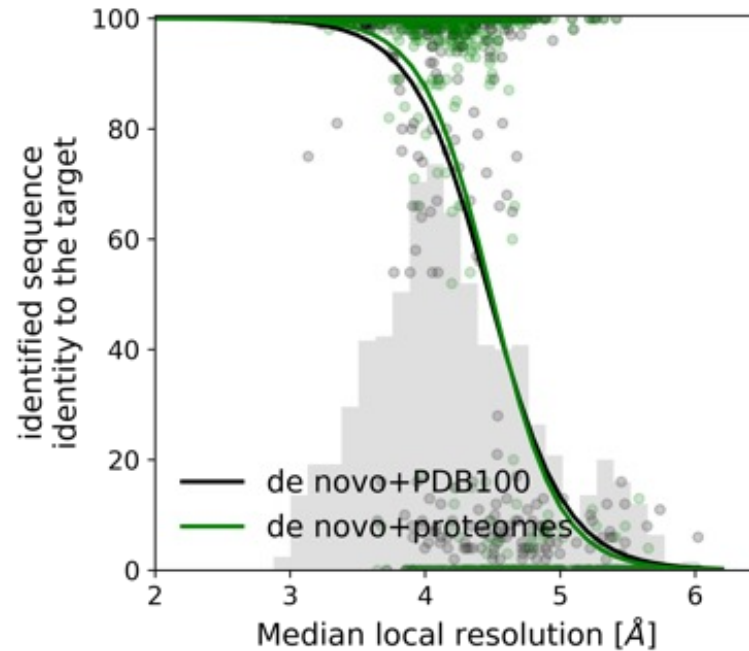
# Protein sequence identification with findMySequence



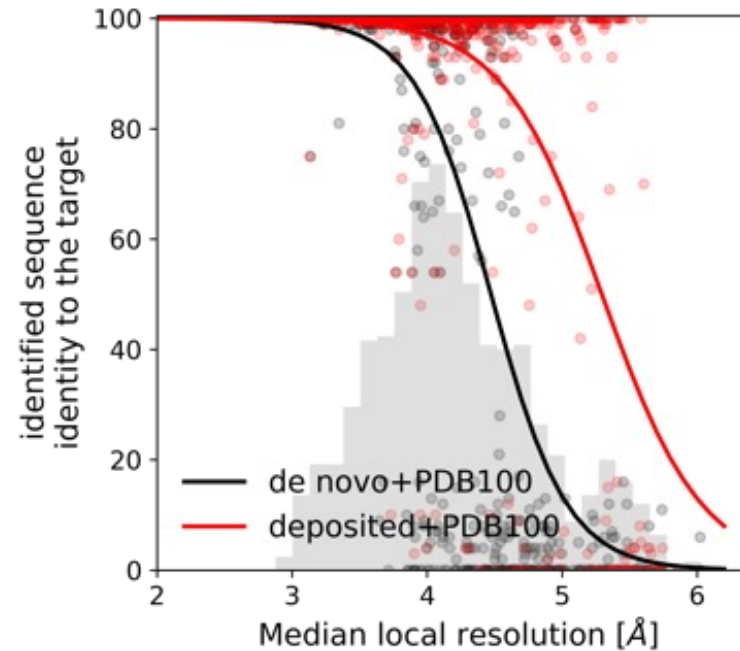
residue-type NN classifier  
328 input features  
111,800 parameters

# Resolution, model quality, and sequence DB size

models built *de-novo* with ARP/wARP  
with different sequence DB sizes

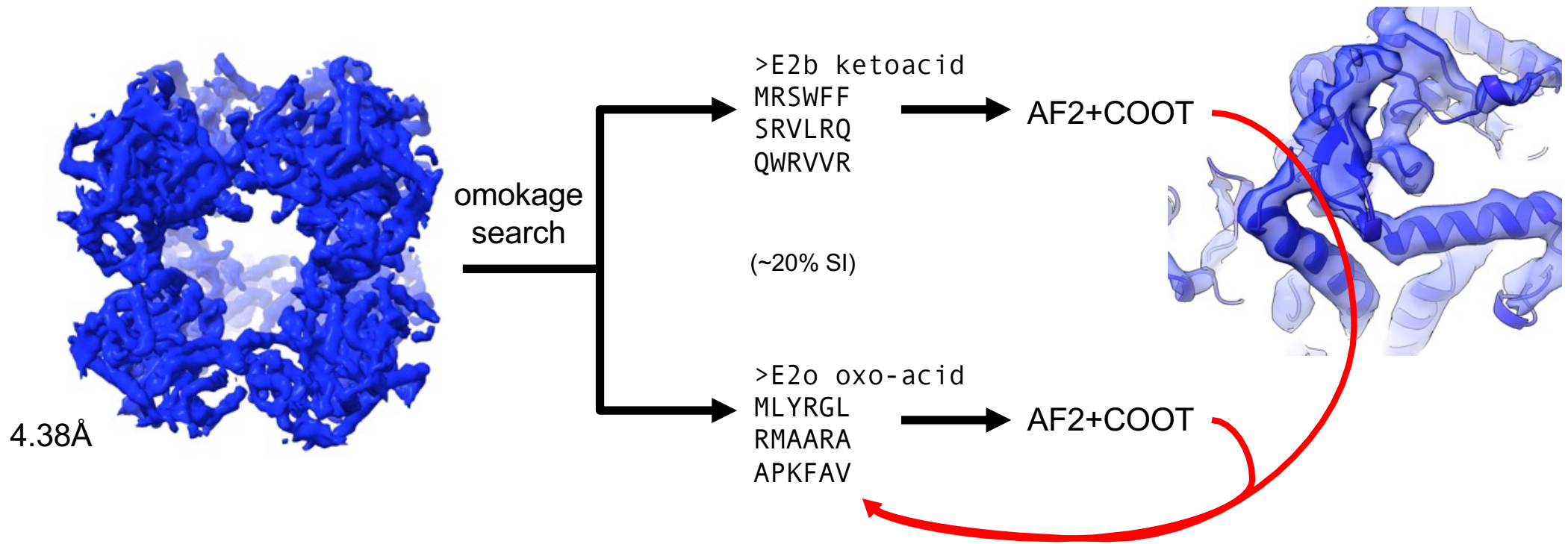
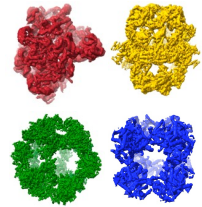


models built *de-novo* with ARP/wARP  
and deposited coordinates



# findMySequence: sequence validation in EM

A tale of two dehydrogenases from *Chaetomium thermophilum* native cell extracts



4.38Å

omokage search

>E2b ketoacid  
MRSWFF  
SRVLRQ  
QWRVVR

AF2+COOT

(~20% SI)

>E2o oxo-acid  
MLYRGL  
RMAARA  
APKFAV

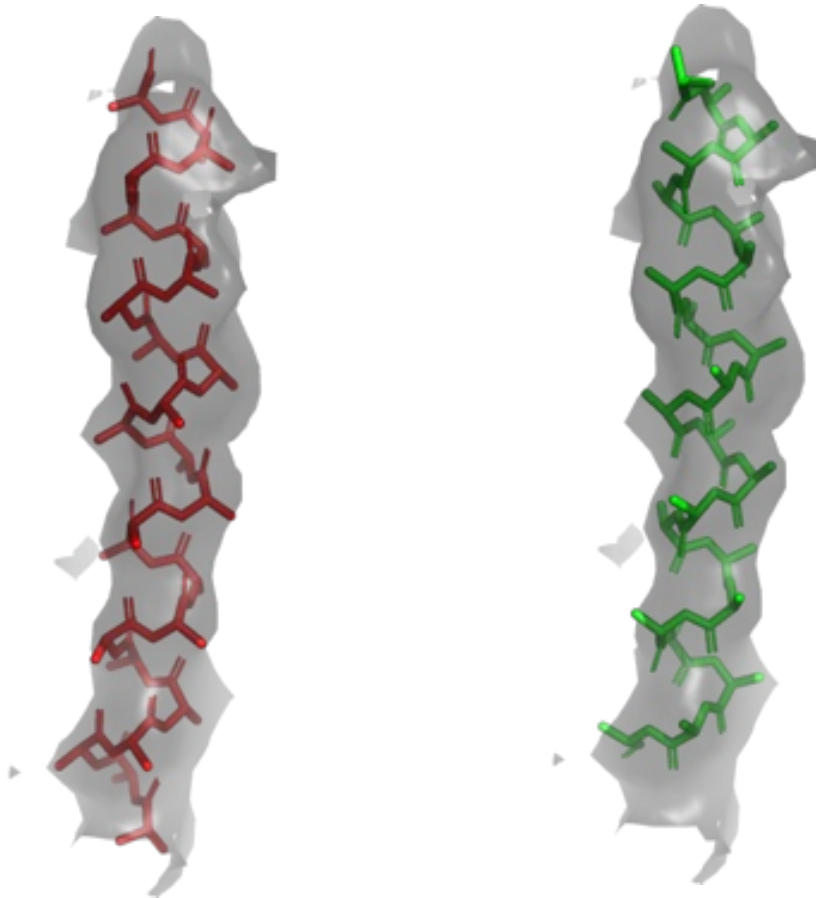
AF2+COOT

findMySequence

with Panos Kastiris  
Chojnowski et al IUCrJ 2022  
Skalidis et al Structure 2022



# findMySequence: model building



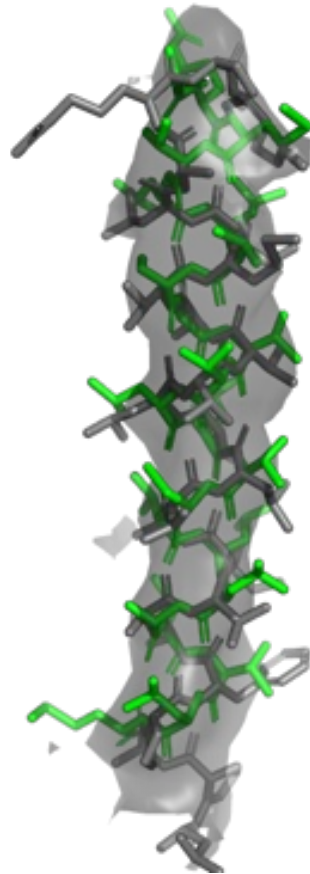
ESX-5 hexameric pore complex at 3.4 Å resolution

22aa  $\alpha$ -helix models  
built with COOT  
(two alternative directions)

# findMySequence: model building



p-value = 0.1



p-value = 0.99

Beckham, Ritter, Chojnowski et al *SciAdv* (2021)



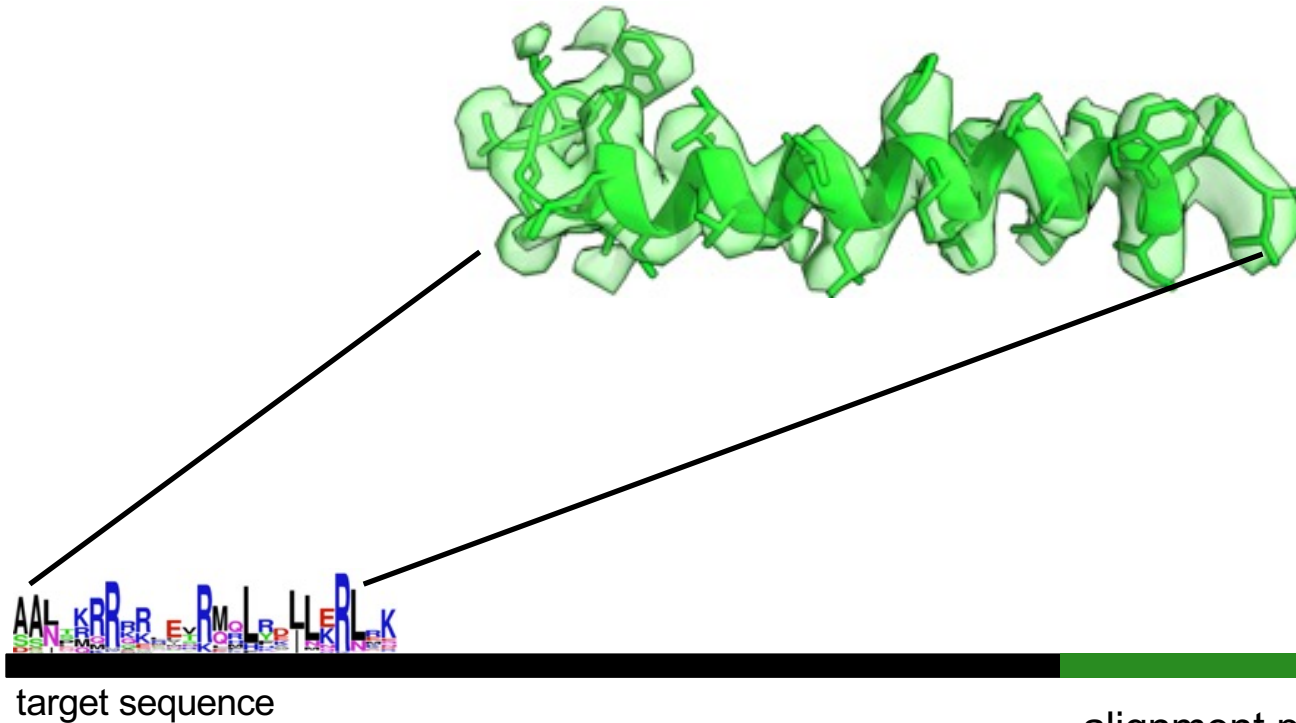
ESX-5 hexameric pore complex at 3.4 Å resolution

22aa  $\alpha$ -helix models  
built with COOT  
(two alternative directions)

... and assigned to the target  
sequence with findMySequence

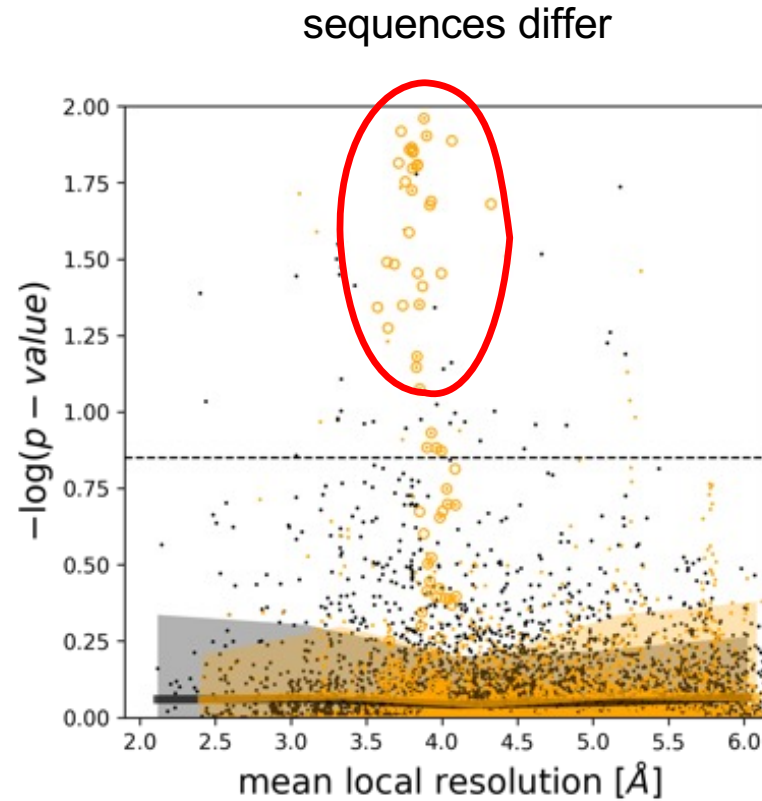
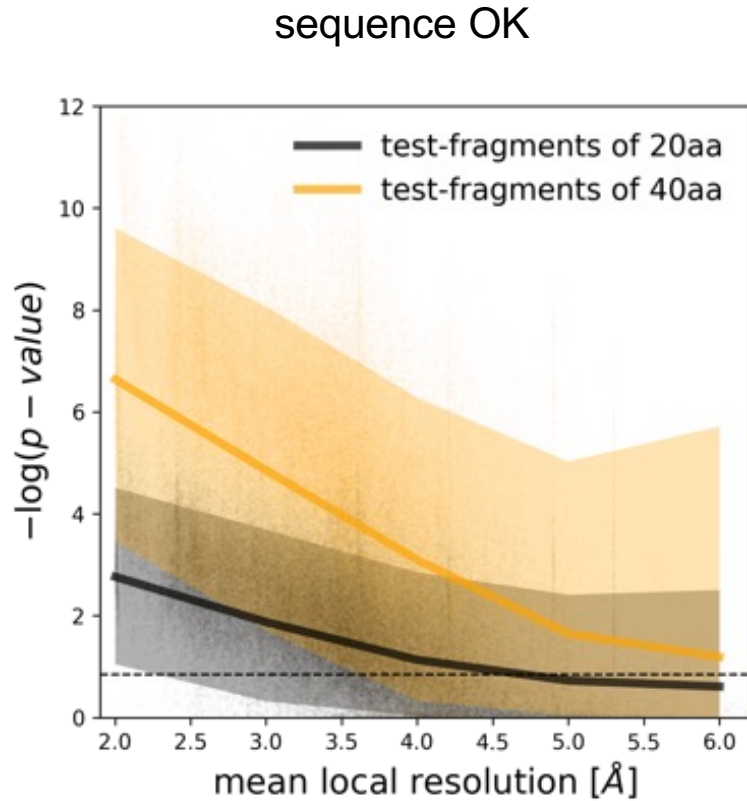


# Sequence re-assignment as a model validation



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

# Is the sequence assignment p-value a reliable score?

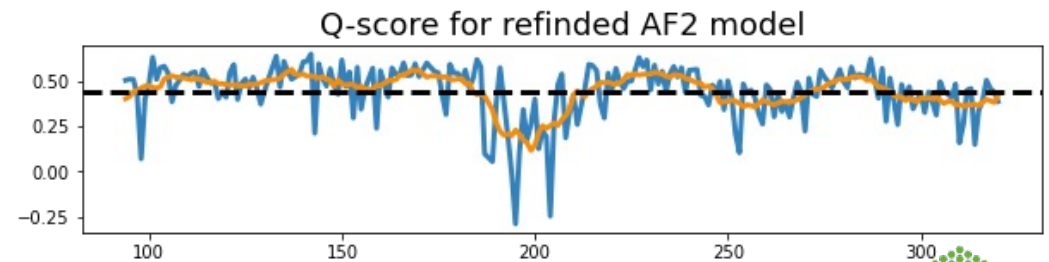
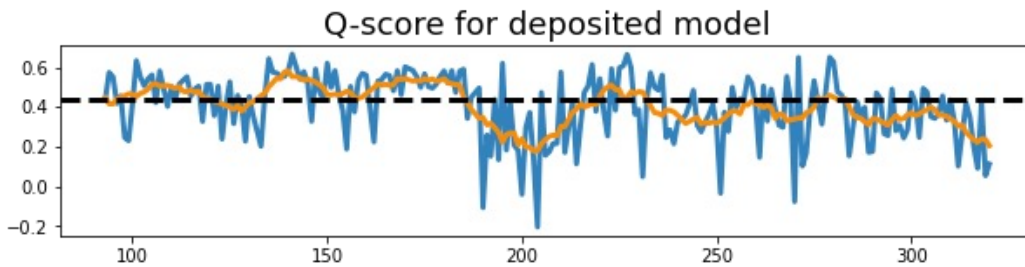
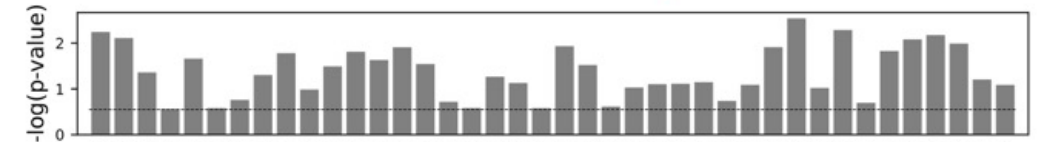
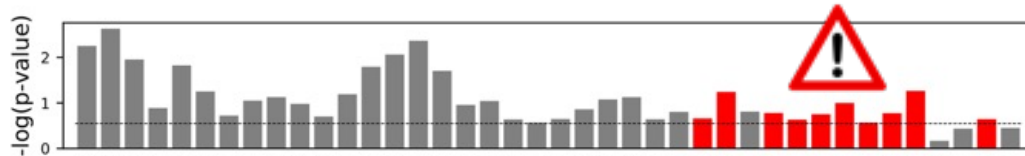
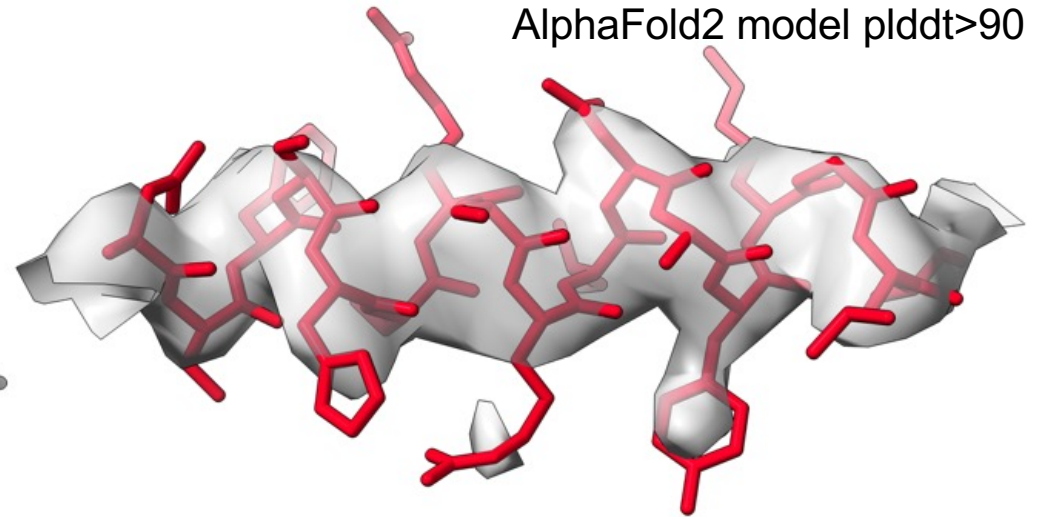
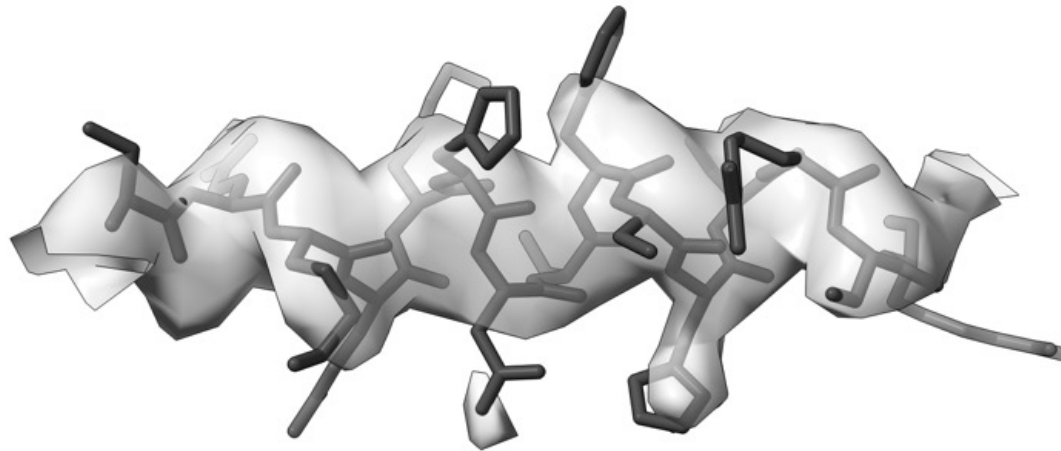


	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?

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

# hidden errors in EM models: finding a better hypothesis

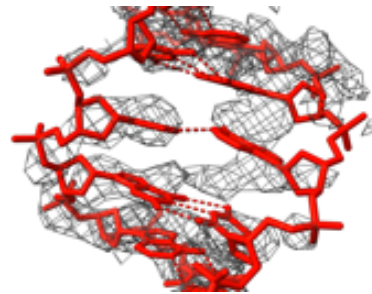
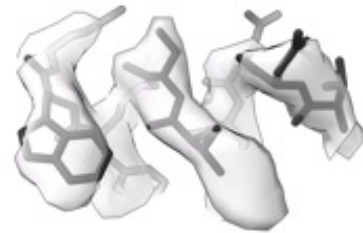
cytoplasmic domain of a cation channel at 3.8Å resolution



# checkMySequence: complete sequence assignment validation

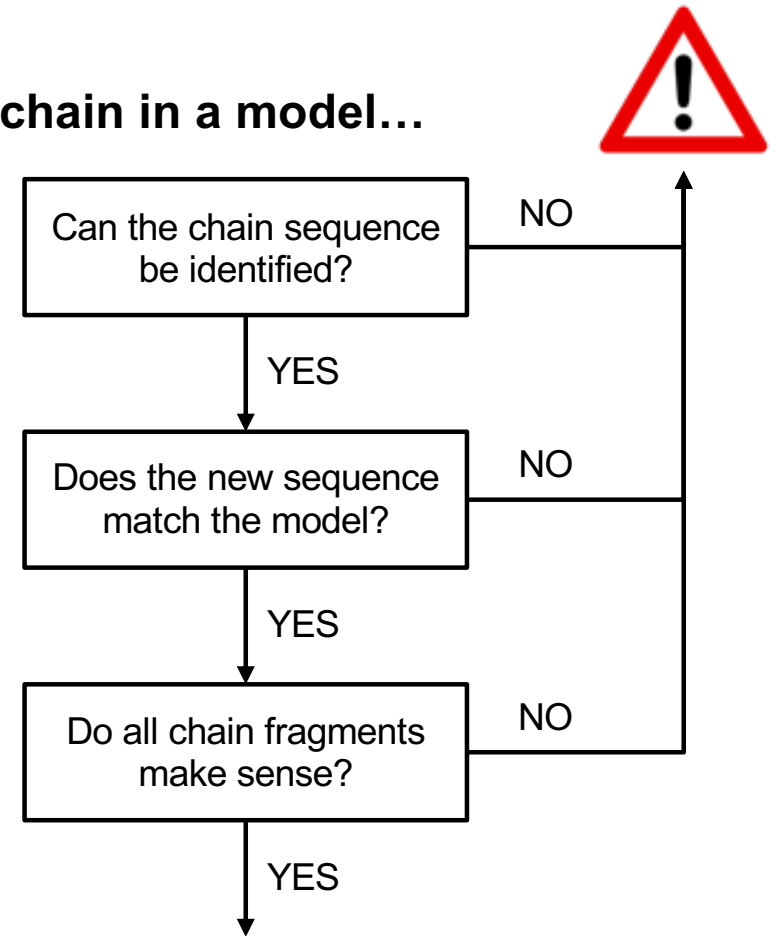
## checkMySequence dependencies

- findMySequence (aa, HMM)
- doubleHelix (na, CM/HMM)



**checkMySequence** (proteins, EM) Chojnowski ActaD 2022  
**checkMySequence** (proteins, MX) Chojnowski ActaD 2023  
**findMySequence** (proteins, EM and MX ) Chojnowski et al ActaD 2022  
**doubleHelix** (NA, EM and MX) Chojnowski NAR 2023

## ∀ chain in a model...

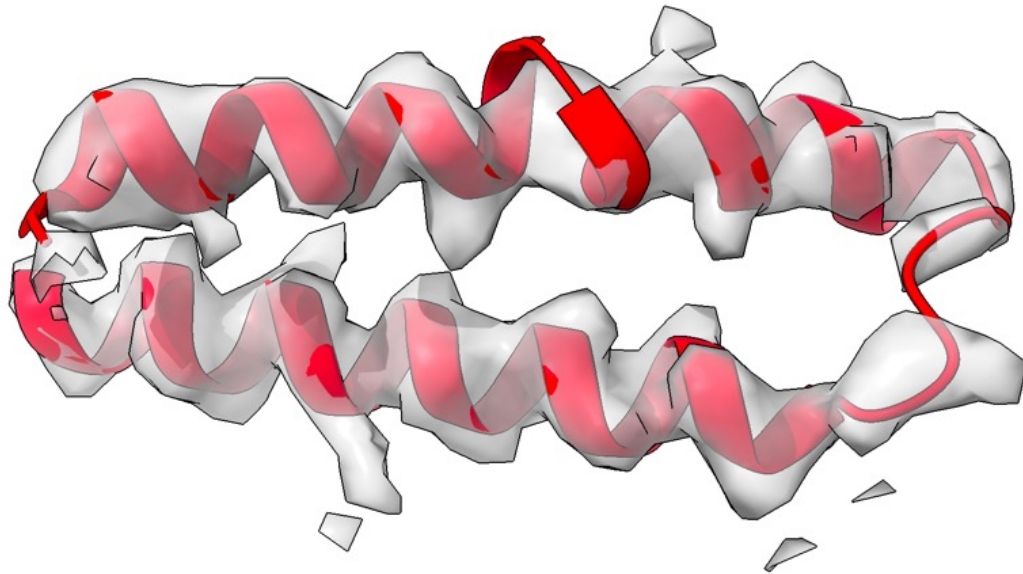


Your model **seems** fine!





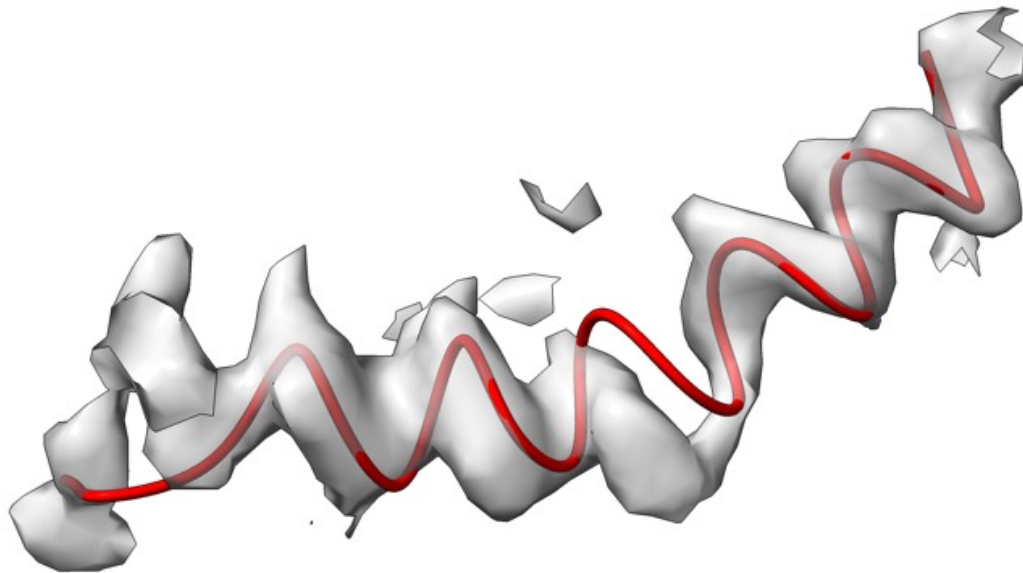
checkMySequence for 6,000 cryo-EM models at 5Å or better



- 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

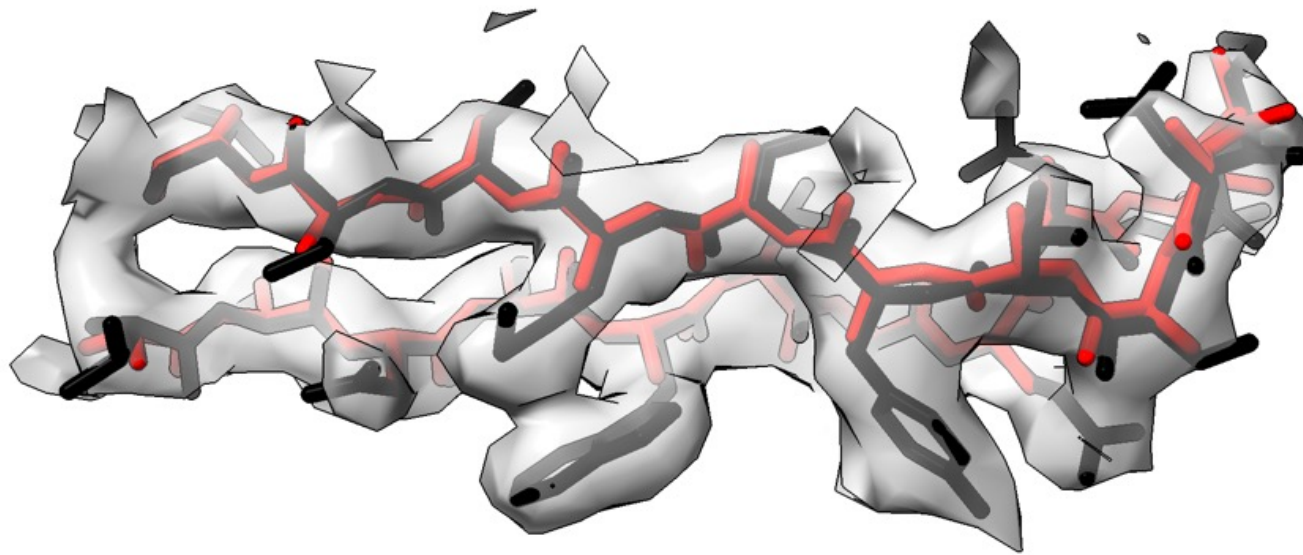
checkMySequence for 6,000 cryo-EM models at 5Å or better



- 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

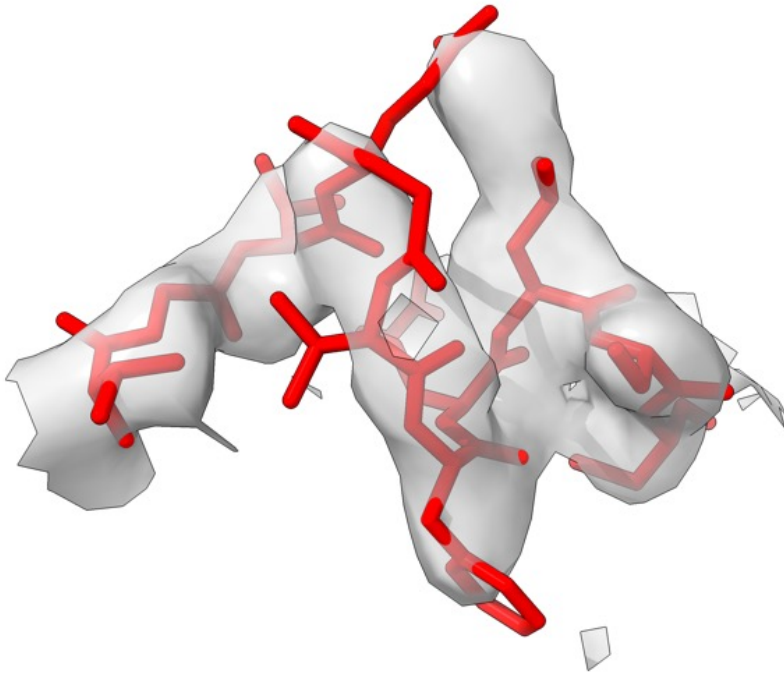


**deposited model**  
**rebuilt with findMySequence**

- Fragment C/161-200 is shifted by 2 residue [-log(p-value)=3.11]  
model seq 161-200  
lnnfypk**DINVKWKIDG**SERQNGVLNSWTDQDSK**DSTYSMSSTLTL**tkdeyer  
new seq 159-198  
lnnfy**PKDINVKWKIDG**SERQNGVLNSWTDQDSK**DSTYSMSSTLTL**tkdeyer

ferroportin at 3 Å

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



- Fragment A/2367-2426 is shifted by 1 residue [-log(p-value)=3.02]

model seq 2367-2426

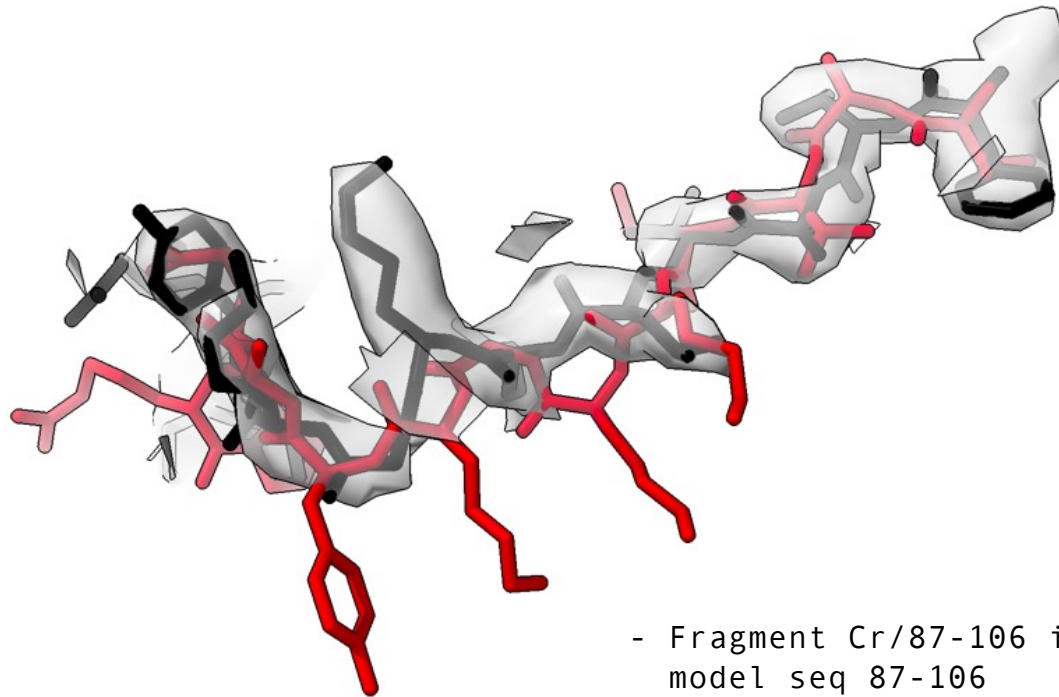
kgkslrVPEKVPFRMTQNIETALGVTGVEGVFRLSCEQVLHIMRRGRETLLTLLEAFVYDPLVDWtaggeag

new seq 2366-2425

kgkslRVPEKVPFRMTQNIETALGVTGVEGVFRLSCEQVLHIMRRGRETLLTLLEAFVYDPLVDWtaggeag

a kinase at 3.6 Å

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

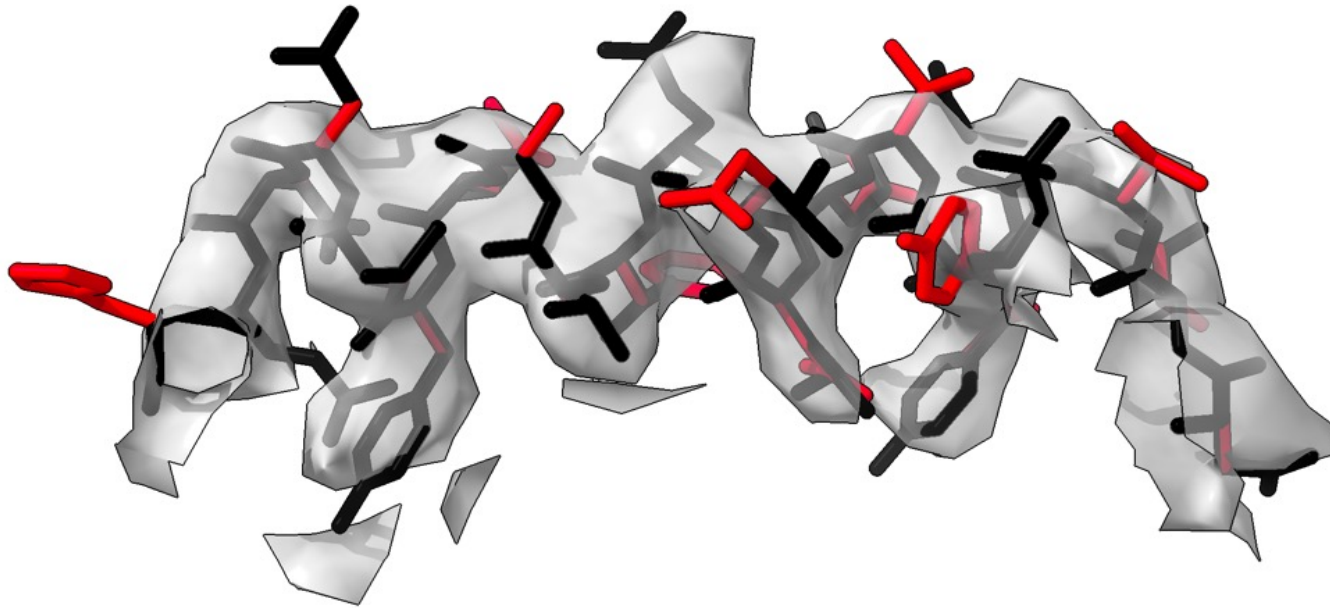


**deposited model**  
**rebuilt with findMySequence**

- Fragment Cr/87-106 is shifted by 4 residue [-log(p-value)=2.95]  
model seq 87-106  
akntvr~~vd~~fKAGPRRSLK~~KL~~KNLLIGSKYrkdltq  
new seq 83-102  
akntvRVDFKAGPRRSLK~~KL~~KNLLI~~g~~skyrkdltq

50S ribosomal protein 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

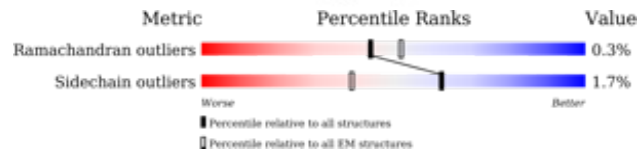
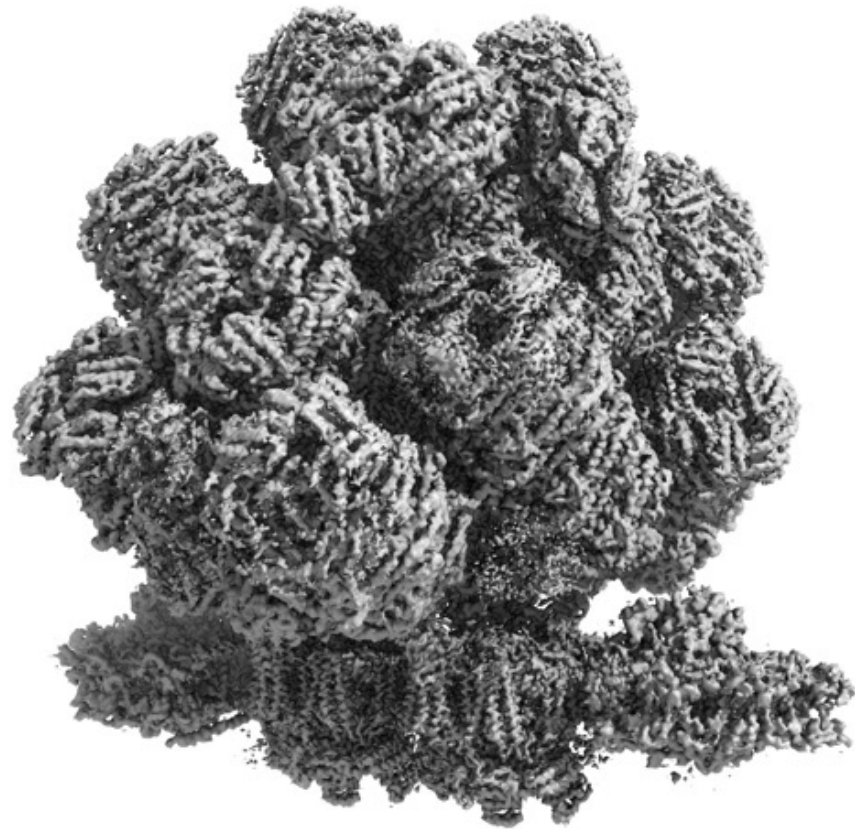
gfaelarryahnlan**ARFLWRNRVGAEAVEVRINHIRQGEVARAWRFDALAIGLRDFKADAELDALAELIASGLS**gsgshvl

new seq 135-194

gfael**ARRYAHNLANARFLWRNRVGAEAVEVRINHIRQGEVARAWRFDALAIGLRDFKADAELDA**laeliasglsghvl

CRISPR-associated protein at 3.8 Å

# recent EM use case: In situ photosystem megacomplex at 3.3Å



18MDa/153,000 residues/864 aa chains



Panagiotis Kastiris  
@3Dstructure

Replying to @emeKato

PDB validation report of \*only\* 4691 pages! Jaw-dropping work!



Tristan Croll  
@CrollTristan

Wow. At a rate of 1 second per residue, viewing every residue in this monster would take...

... checks notes...

... blinks...

... just under 43 \*hours\*.



Tristan Croll  
@CrollTristan



Tristan Croll  
@CrollTristan

Replying to @ktototak

OK, I'm happy to stand corrected - the example you showed me is clearly a bona fide register error.

7y5e/EMD-33618

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

# recent EM use case: In situ photosystem megacomplex

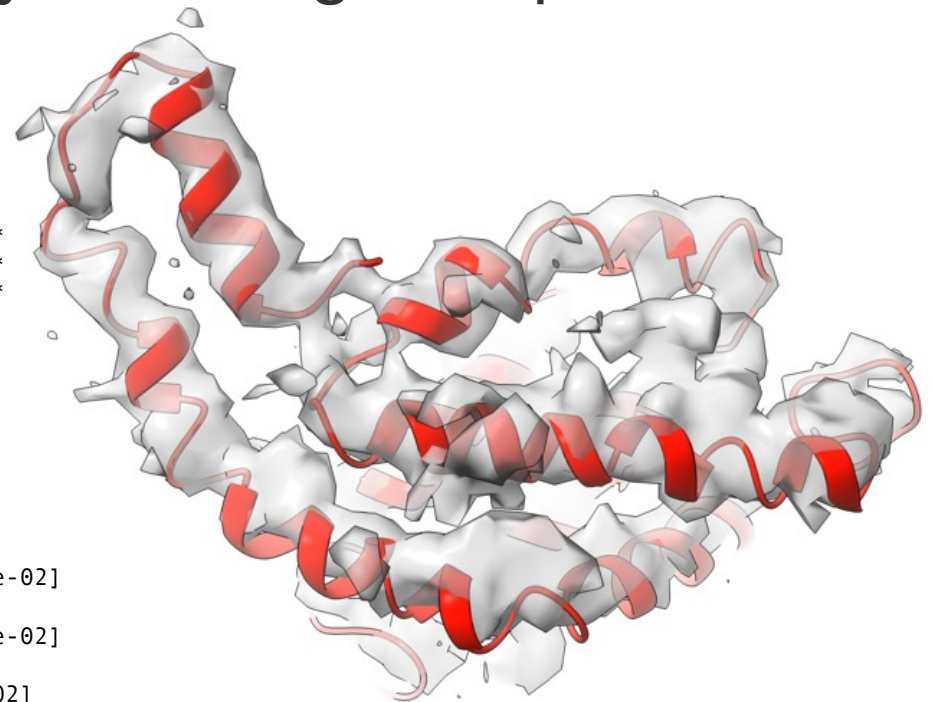
```
*****  
***** SUMMARY *****  
*****
```

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

AA/1:164  
... **[144 removed for the presentation clarity]**  
Z9/25:299

==> Sequence register shifts

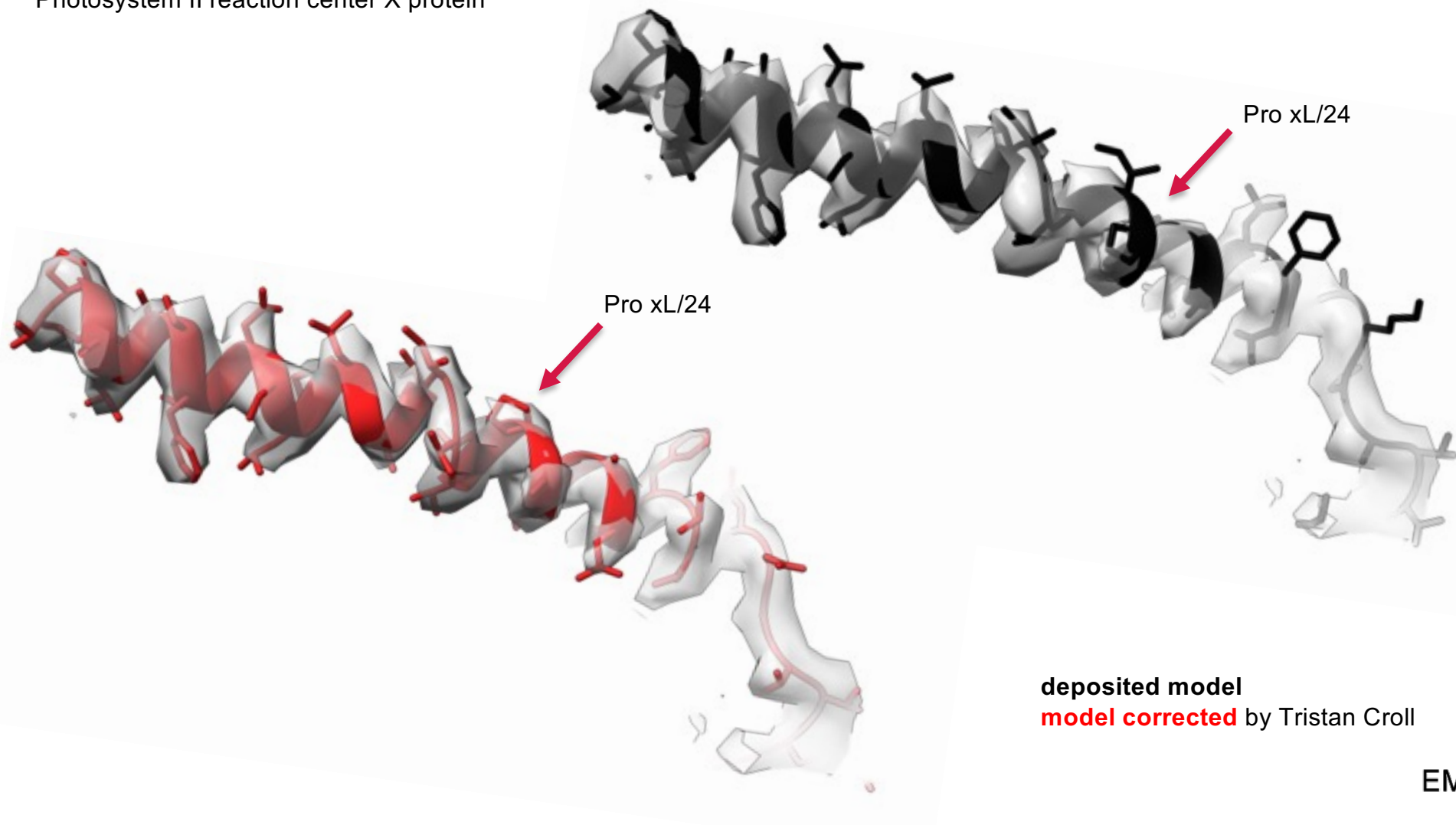
- protein chain fragment YP/227-266 may be shifted by 2 residues [p-value=4.14e-02]
- protein chain fragment bP/227-266 may be shifted by 2 residues [p-value=4.04e-02]
- protein chain fragment xL/18-41 may be shifted by -1 residue [p-value=1.19e-02]
- protein chain fragment XL/18-37 may be shifted by -1 residue [p-value=7.30e-03]
- protein chain fragment X6/18-37 may be shifted by -1 residue [p-value=5.30e-03]
- protein chain fragment x6/18-37 may be shifted by -1 residue [p-value=5.35e-02]



Time elapsed 0:43:23

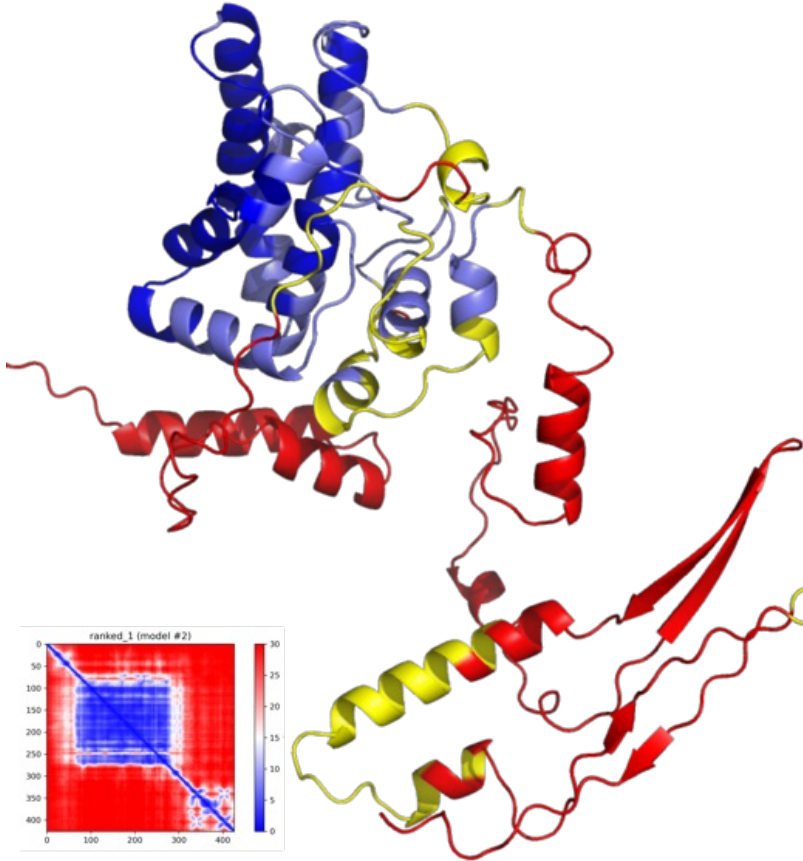
# recent EM use case: In situ photosystem megastructure

Photosystem II reaction center X protein



# recent EM use case: In situ photosystem megastructure

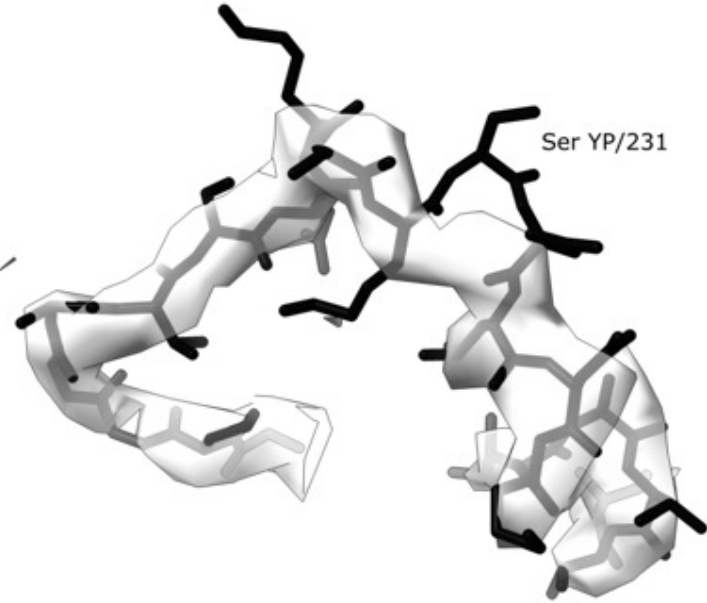
Phycobilisome linker polypeptide



AlphaFold2 model

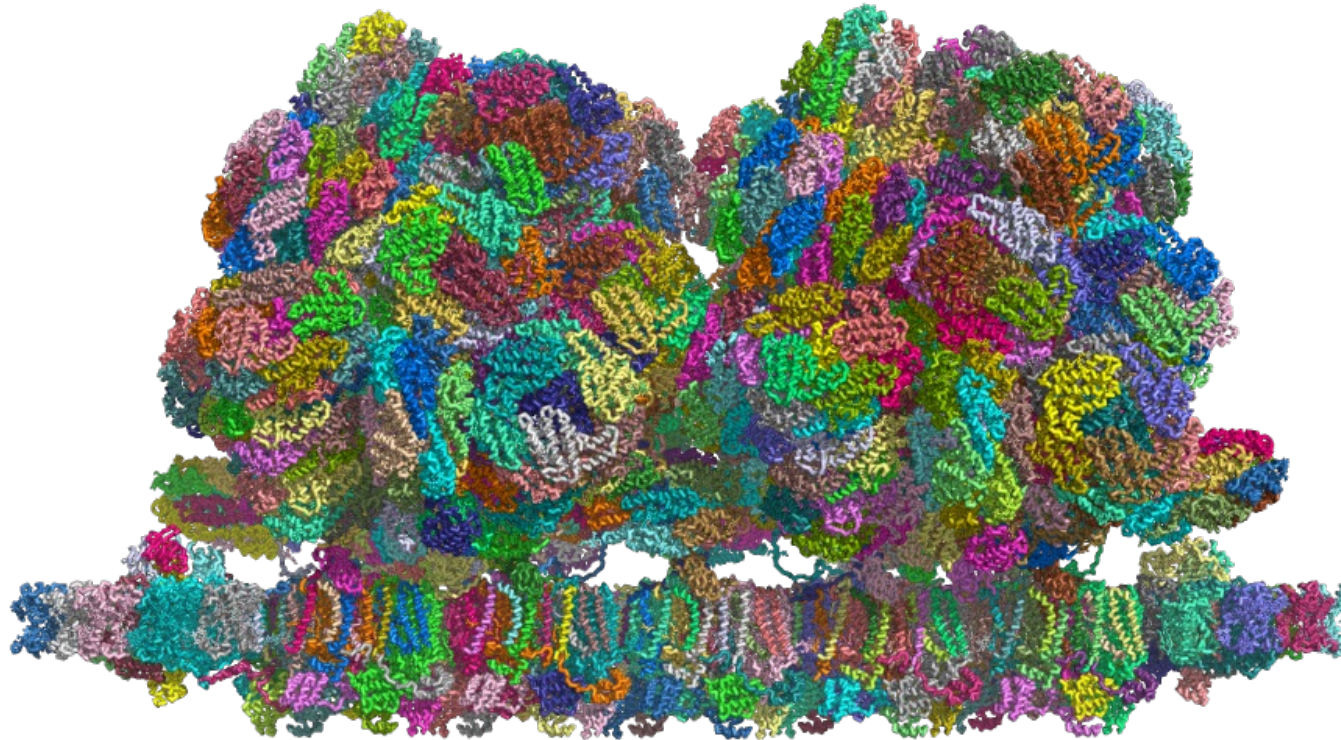


deposited model



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

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



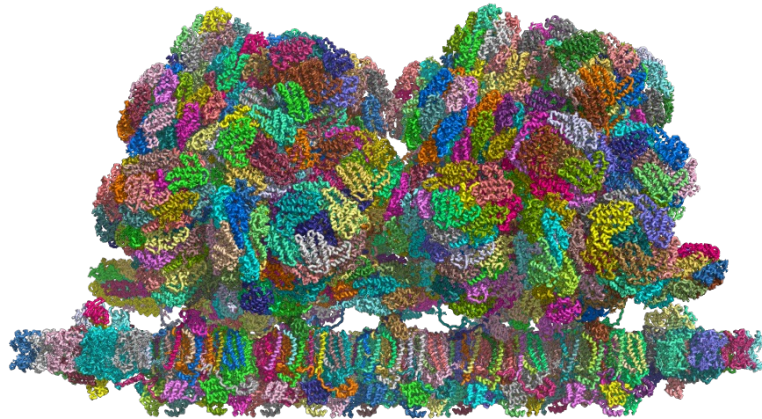
**37MDa/305,000 residues/1,792 aa chains**

7y5e/EMD-33618

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

# 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

```
- Protein chain Y5/62:298 - sequence identity to reference 65.22% [E-value=4.00e-10]
1234567890123456789012345678901234567890123456789012345678901234567890
model -----RRTVSFNARVA-----RN-KSQAKKILEKADEFFARSVTM
refseq          AAAAAAAAAAAMAFVGSAAASFTGASAVKANERKRSVCSLQMVAMPQTGLVNSKFSARMAKKTAKQTKNKVDEYMARSVQR

model          QYKAFACPNGVYDIQCTEGTVKGAAYEKRAMAVSAAFRAKQASPAAKARALFENRRHAIIASHECQHEEDLFVRFPKLSA
refseq          QYKQAAVATGVYGTQCTEGTVKGAAEASRSAALSRQFRIKQRSAFSKAHDLFEFRKHAIIAAGCSYEEKMVTRFPKLAA

model          AYMMGKTEAMRTCSRYVVPDSLEEEYMAASVDRQMKERACPGGVYASSSCVEGNAKGQAEQARVAALATAFRSAQKSASKT
refseq          AMVLGQTEMMRTCSRYVVPESVEEEYMAASVDKQMKRRGAPGGVYSLSCAEGVAKGQAEIARVSALGAAYRAASKSASAV

model          TAERYSSAAYGRDHFHAGCSYEEVSFNTYPATAAAMRSKSYNY-----
refseq          TAERYNSMAYGRVHFHAGCSYEEQQFNKYPAAAAAMRSDSYGYAAAAAAAAA
```

Time elapsed 1:46:54

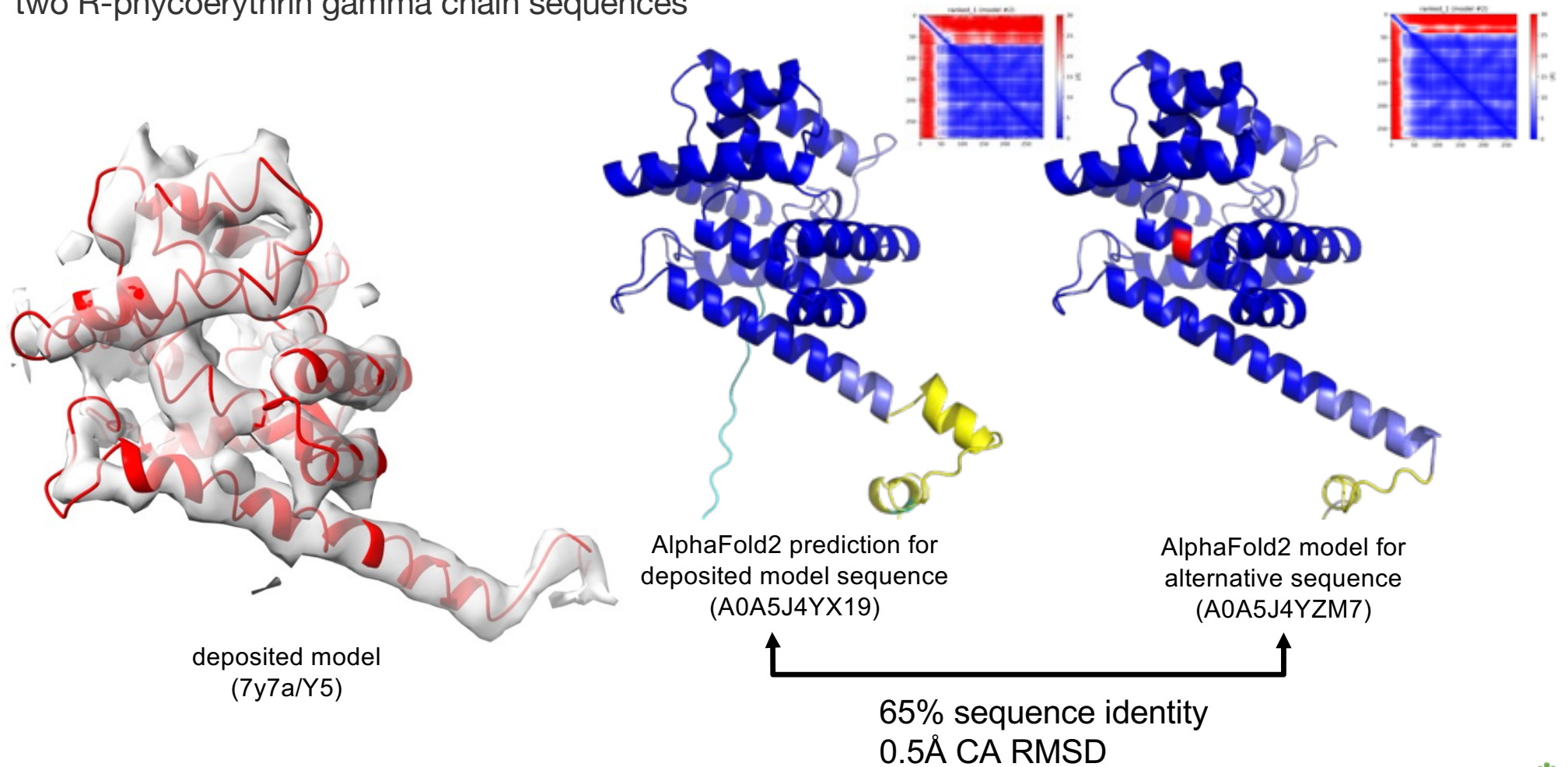
7y5e/EMD-33618

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

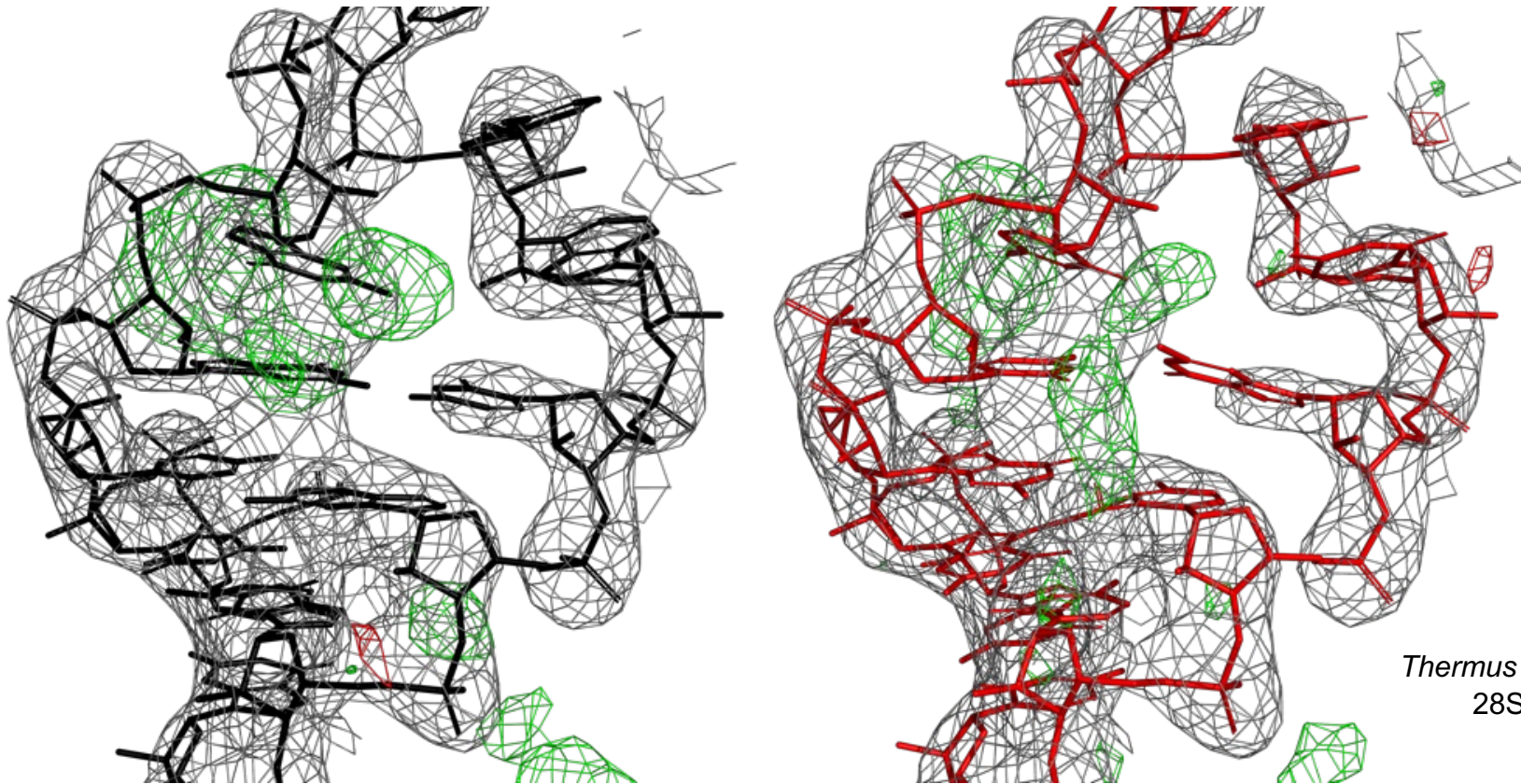


# When you think you've seen it all...

two R-phycoerythrin gamma chain sequences



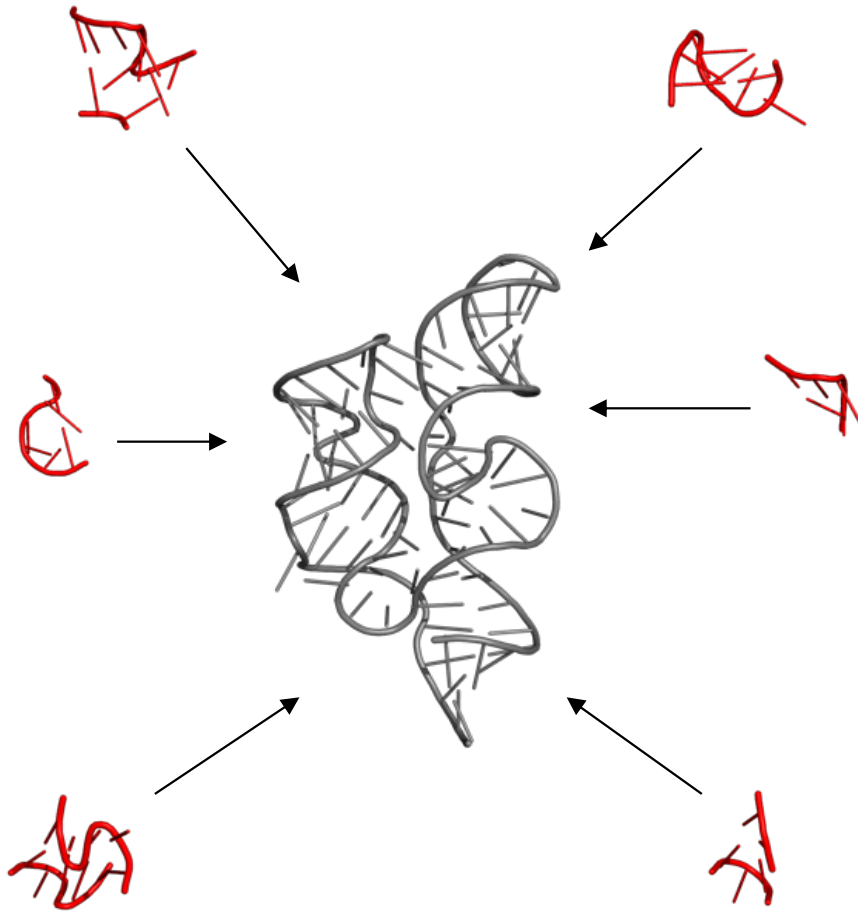
# Nucleic acid MX/EM model building is challenging



*Thermus thermophilus*  
28S MX @ 2.8Å

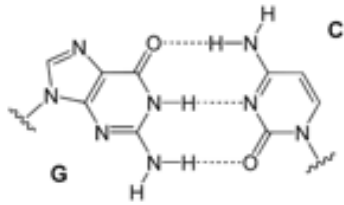
Different types of purines and pyrimidines are (usually) indistinguishable in MX or EM maps

# Base-pairs from backbone geometry

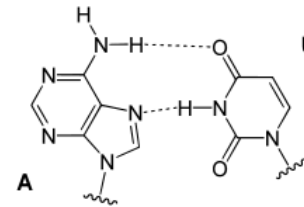
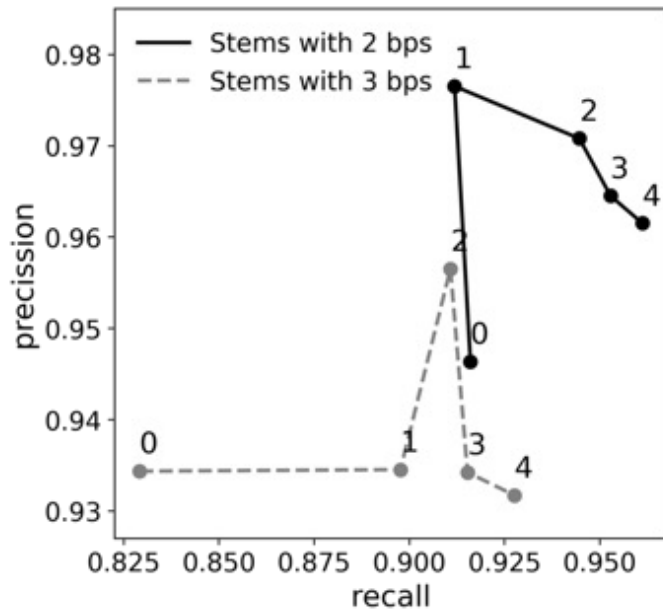


input RNA model backbone is “covered”  
with motifs with known secondary structure  
(from [genesilico.pl/rnabricks2](http://genesilico.pl/rnabricks2))

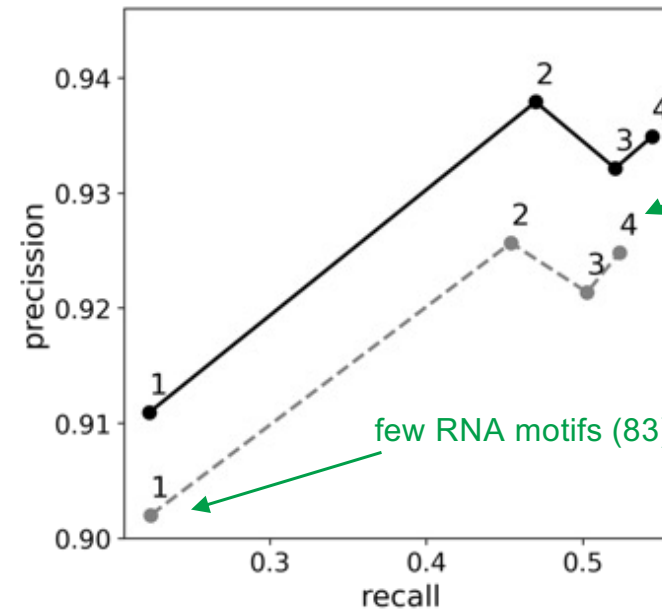
# Base-pairs from backbone geometry



Watson-Crick (canonical) base-pairs



non-canonical base-pairs

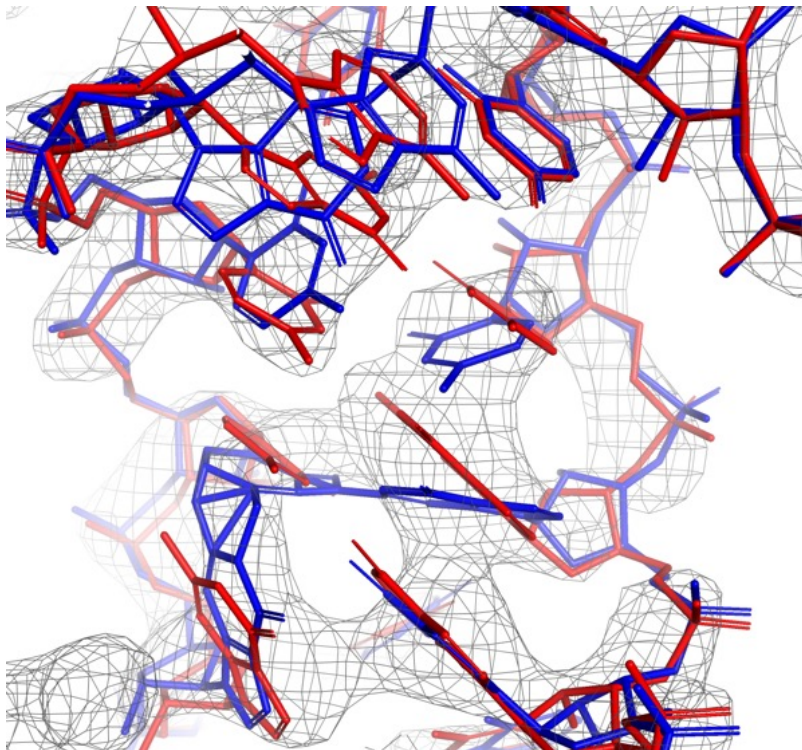


lots of RNA motifs (2,664)

few RNA motifs (83)

# Secondary structure restraints from doubleHelix

nucleic acid sequence identification and assignment tool



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

Standard tools

phenix.secondary\_structure\_restraints  
LibG  
ISOLDE

are very reliable if a model is good,  
but cannot restrain base-pairs if

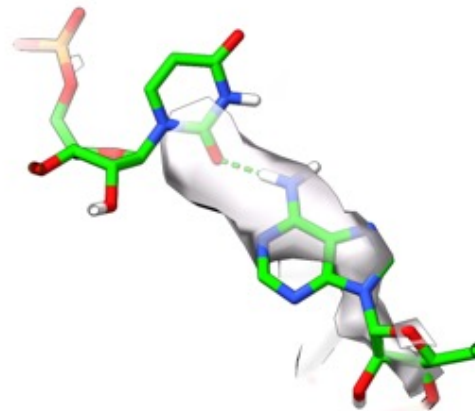
- geometry is very bad
- sequence is wrong

# Secondary structure restraints from doubleHelix

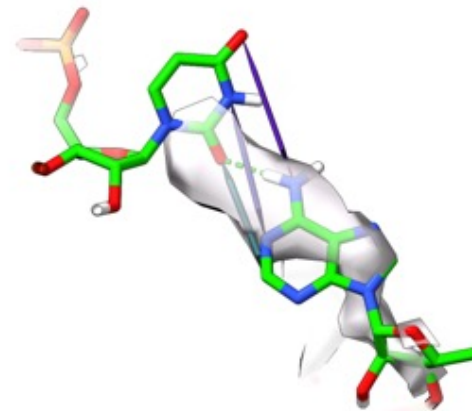
nucleic acid sequence identification and assignment tool



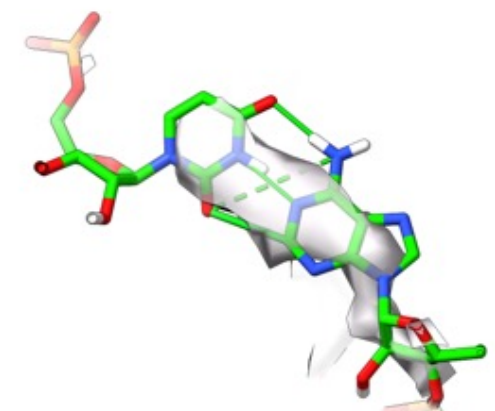
RNA polymerase  
@2.5Å (7bv2/30210)



deposited EM model  
(dashed line - auto-generated  
H-bond restraint)



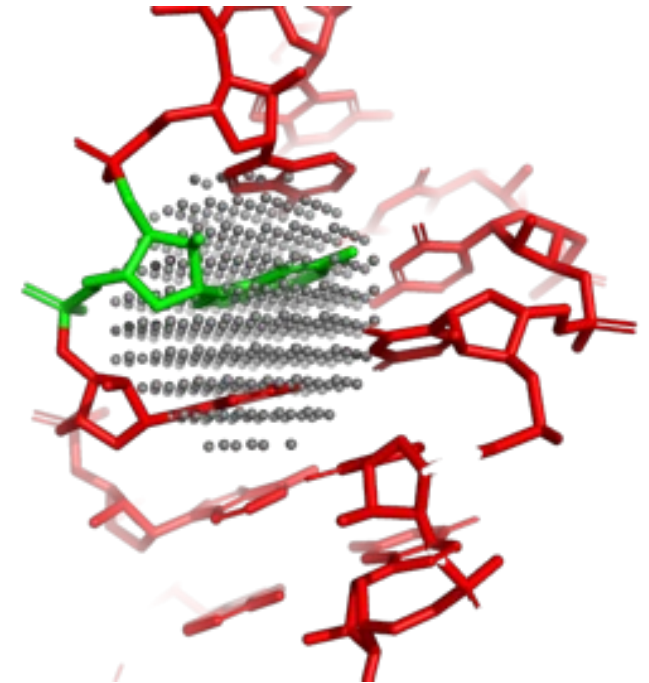
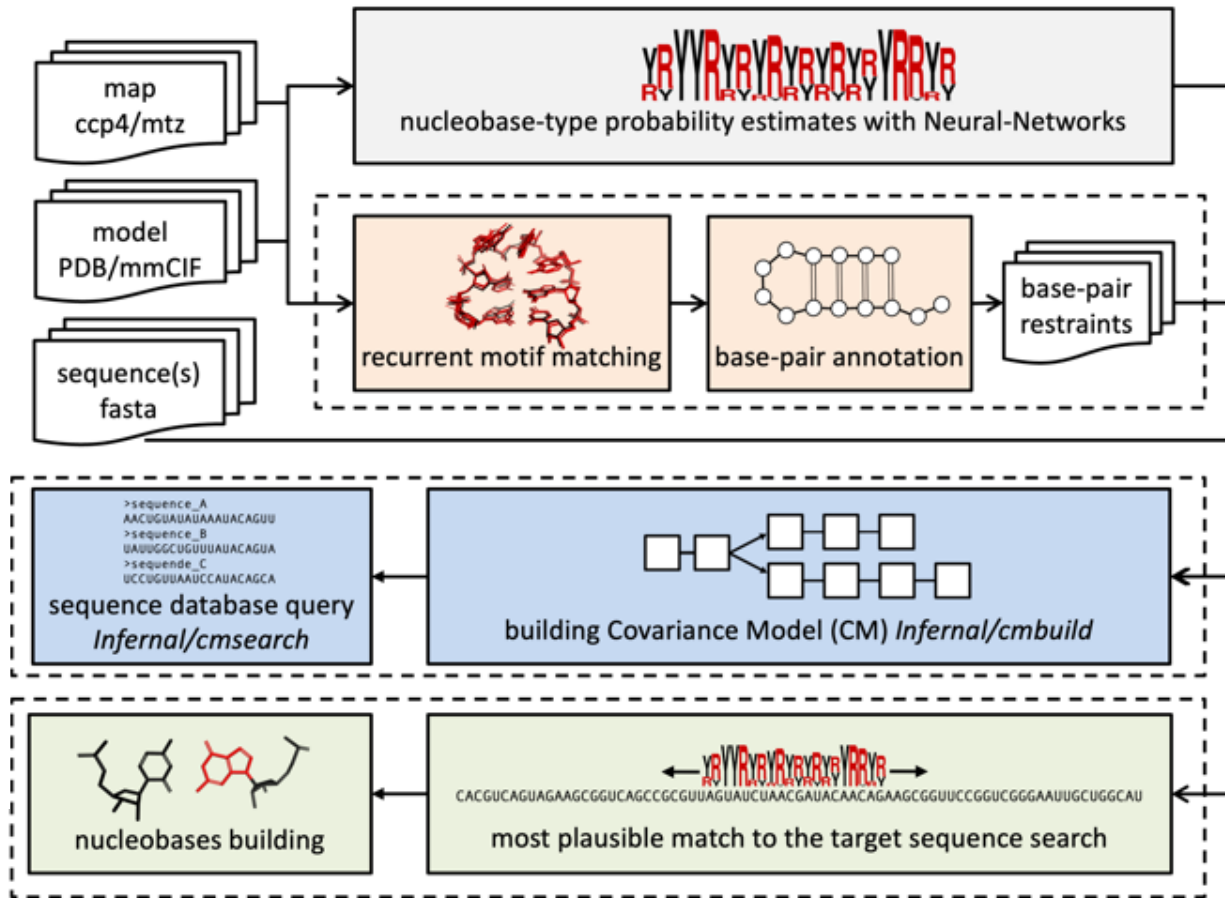
deposited EM model  
with doubleHelix restraints  
(violated - dark solid bars)



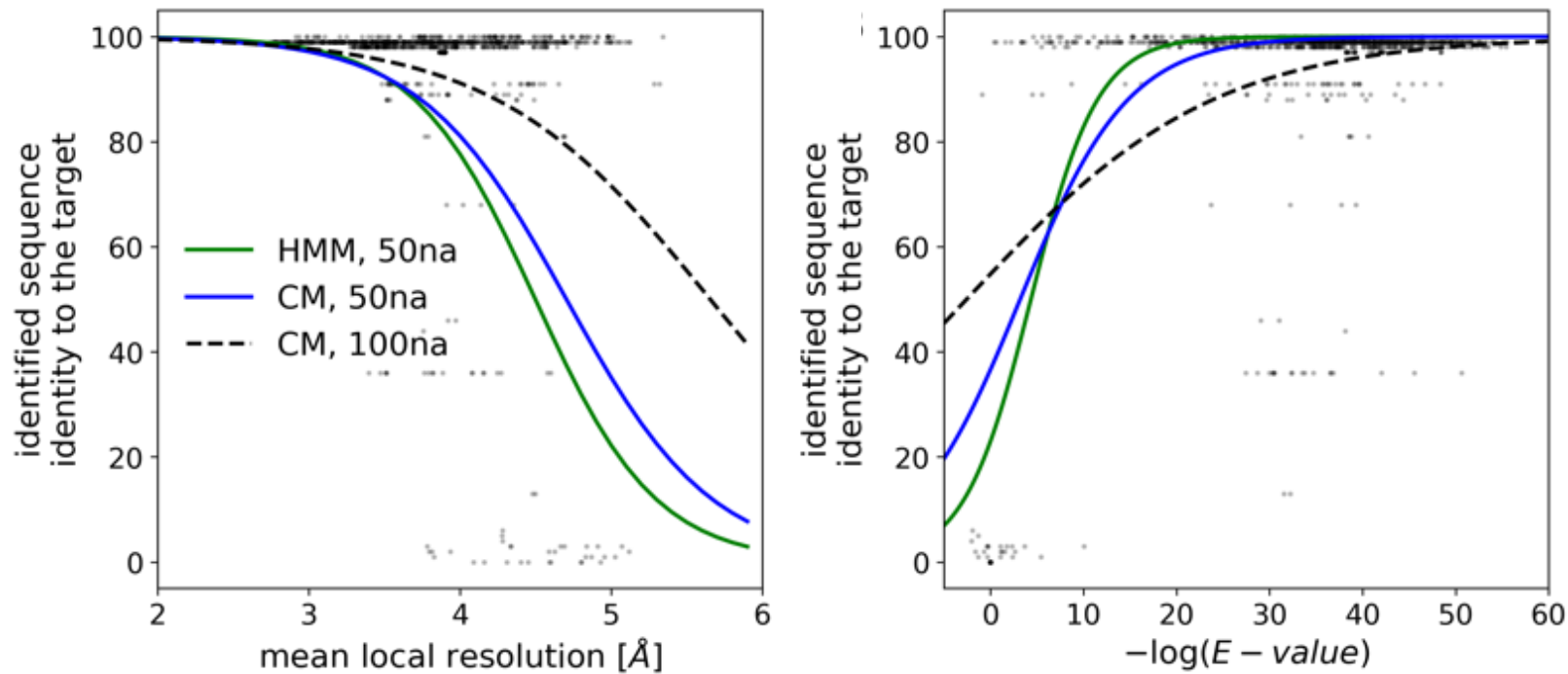
refined EM model  
with doubleHelix restraints  
(satisfied - green solid bars)



# doubleHelix - NA sequence assignment and validation in EM/MX

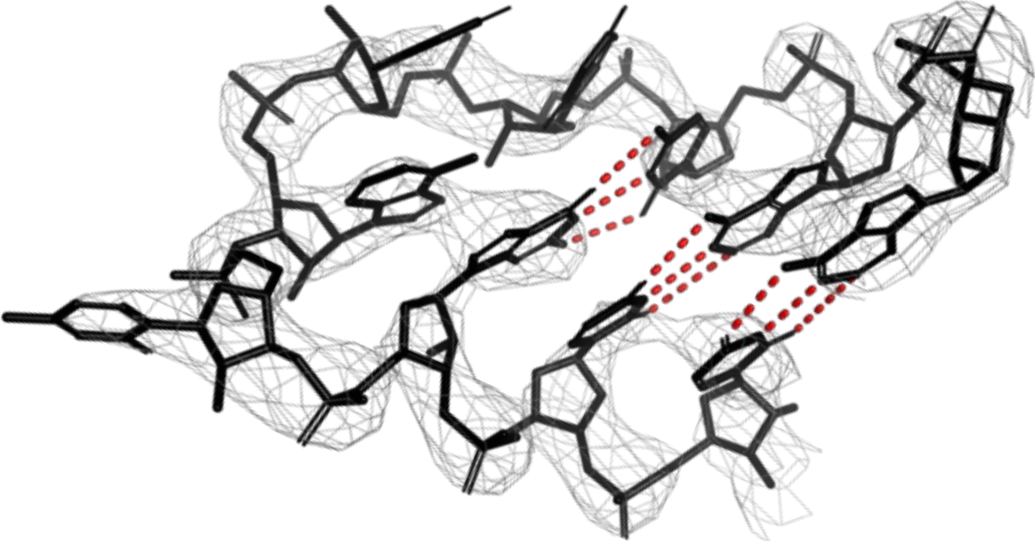


# Nucleic acid sequence identification in EM structures



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

# checkMySequence and RNA/DNA sequence validation



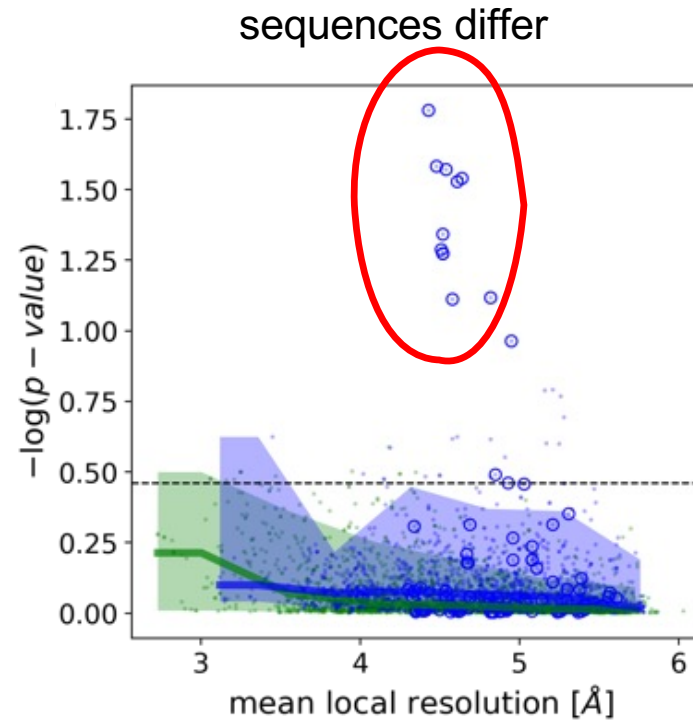
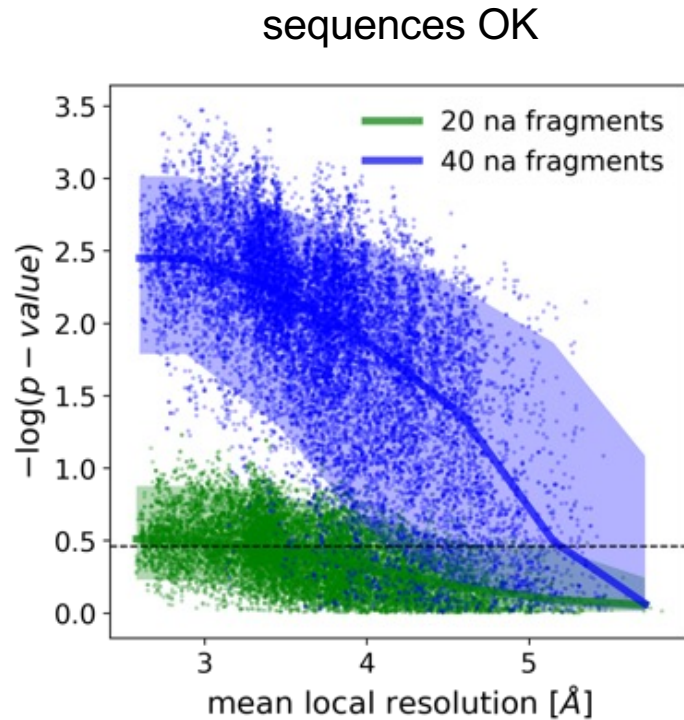
	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?



target sequence

alignment p-value

# checkMySequence and RNA/DNA sequence validation in EM

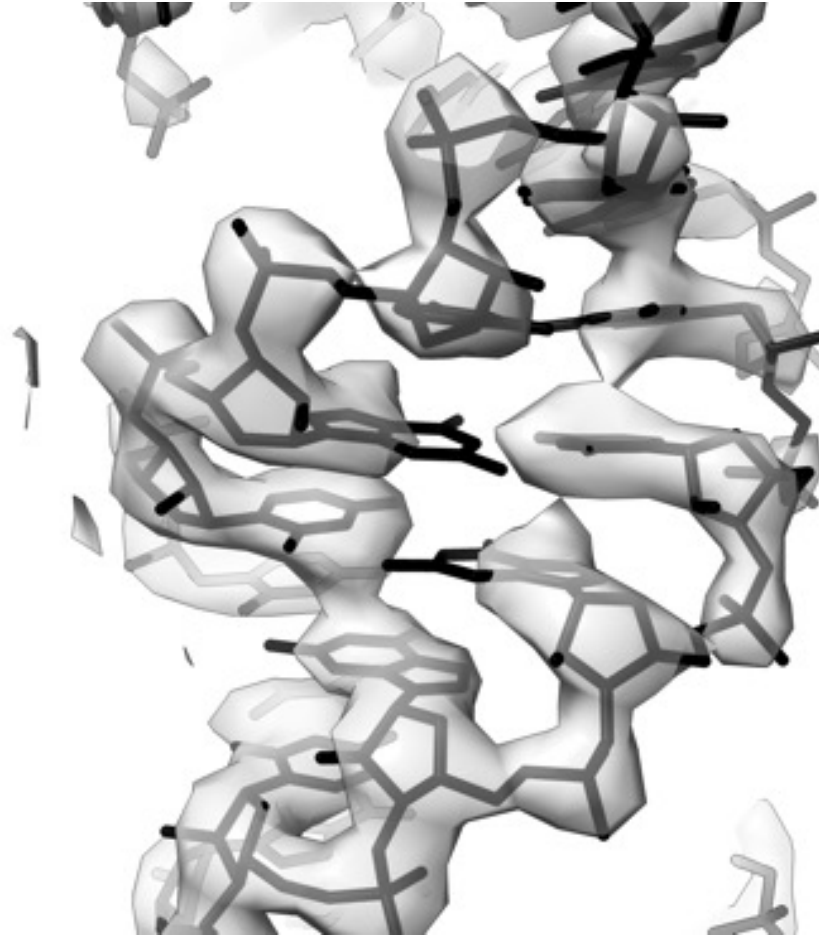
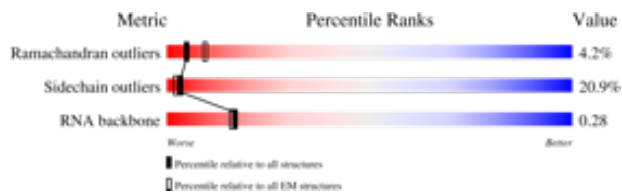
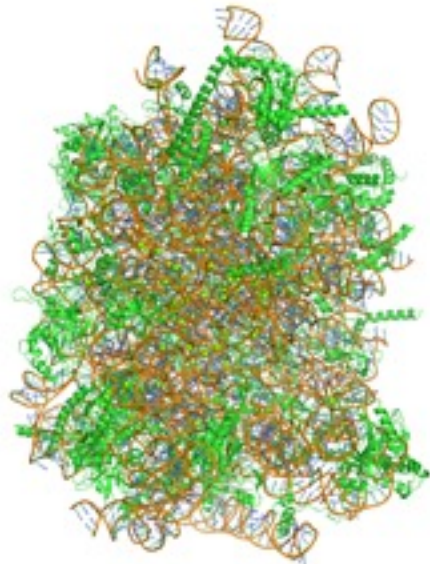


	p-value small	p-value ~ 1
sequences OK	✓	?
sequences differ	⚠	?

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

# #1 Sequence assignment issues in rRNA

expansion segment in an EM model of mammalian ribosome @3.4Å (deposited 2014)



# #1 Sequence assignment issues in rRNA

expansion segment in an EM model of mammalian ribosome @3.4Å (deposited 2014)



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

==> Unidentified chains; check input sequences and model-to-map fit  
 2/7:68  
 ==> Possible sequence assignment issues

- Nucleic-acid chain fragment 5/442-491 may be **shifted by 1 residue** [-log(p-value)=0.72]  
 model seq 442-491

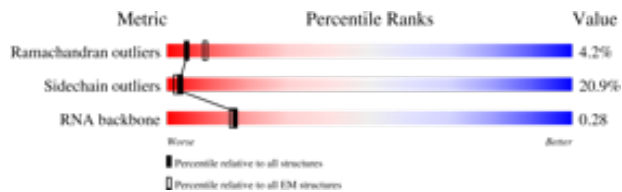
```
cgggugGGGUCCGCGCAGUCCGCCCGGAGGAUUCAACCCGGCGGGUCCGGCCgugucgg
new seq 441-490
cggguGGGUCCGCGCAGUCCGCCCGGAGGAUUCAACCCGGCGGGUCCGGCCgugucgg
```

- Nucleic-acid chain fragment 5/657-706 may be **shifted by -1 residue** [-log(p-value)=1.04]  
 model seq 657-706

```
gucCCCGACCGGCGACCGCCGCCCGGGCGCAUUUCCACCGGGCGGUGCGCcgcgacc
new seq 658-707
gucCCCGACCGGCGACCGCCGCCCGGGCGCAUUUCCACCGGGCGGUGCGCCcgcgacc
```

[...]

Time elapsed 0:05:53

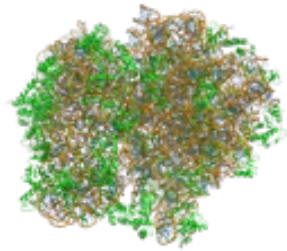




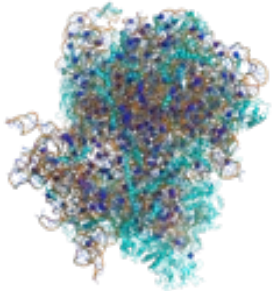
# #1 Sequence assignment issues in rRNA

expansion segment in an EM model of mammalian ribosome @3.4Å (deposited 2014)

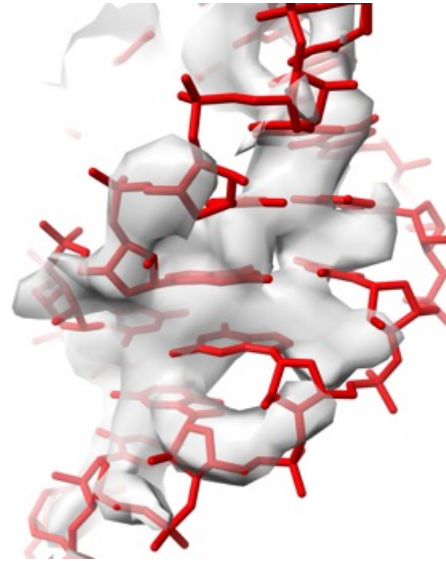
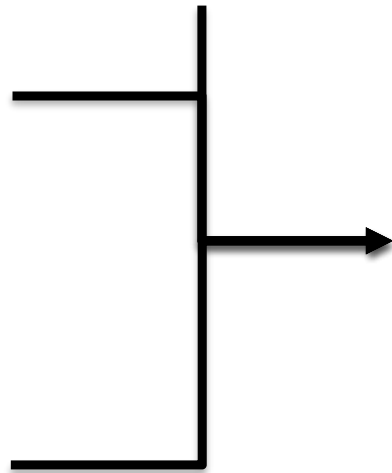
**Assemble  
RNAfold  
COOT**



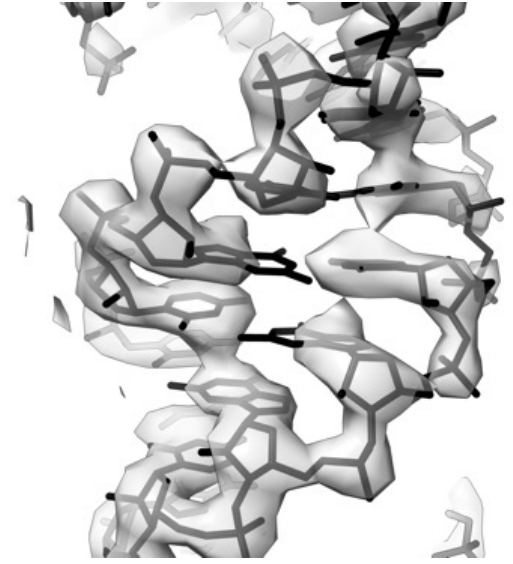
70S MX@3Å



80S MX@4Å



human 28S EM @ 5Å

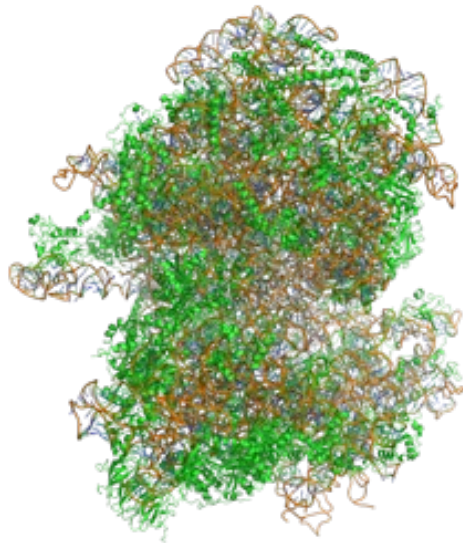


porcine 28S EM @ 3.4Å

# #2 Sequence assignment issues in rRNA

EM structure model of eukaryotic ribosome @3.5Å (deposited 2017)

model-reference alignment of a 28S rRNA fragment



```

model/I-453      1 GGGGAUGAACCAAACGUAUUGUUA CG- UGCCCCAAAUUUACAACUCAUGCAGAUACCAUGAAA 60
reference/I-480  1 GGGGAUGAACCAAACGUAUUGUUA CGGUGCCCCAAAUUUACAACUCAUGCAGAUACCAUGAAA 61

model/I-453      61 GGCGUUGGUUGCUUAAAAACAGCAGGACGGUGAUCAUGGAAGUCGAAAUCCGCUAAGGAGUG 121
reference/I-480  62 GGCGUUGGUUGCUUAAAAACAGCAGGACGGUGAUCAUGGAAGUCGAAAUCCGCUAAGGAGUG 122

model/I-453      122 UGUAAACAACUCACCUGCCGAAGCAACUAGCCCUUAAAAUUGGAUGGCGCUUAAAGUUGUAUAC 182
reference/I-480  123 UGUAAACAACUCACCUGCCGAAGCAACUAGCCCUUAAAAUUGGAUGGCGCUUAAAGUUGUAUAC 183

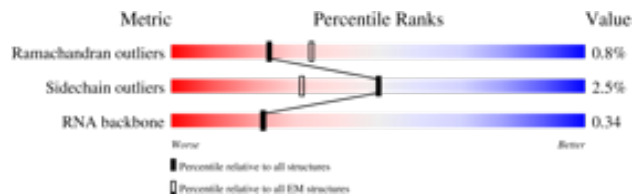
model/I-453      183 CUAUACAUUACCGCUAAAAGUAGAUGAUUUUUAUUUACUUGUGAUUAAAAUUUUGAAACUUUA 243
reference/I-480  184 CUAUACAUUACCGCUAAAAGUAGAUGAUUUUUAUUUACUUGUGAUUAAAAUUUUGAAACUUUA 244

model/I-453      244 GUGAGUAGGAAGGUACAAUGGUUUGCGUAGAAGUGUUUGGCGUAAAGCCUGCAUGGAGCUGC 304
reference/I-480  245 GUGAGUAGGAAGGUACAAUGGUUUGCGUAGAAGUGUUUGGCGUAAAGCCUGCAUGGAGCUGC 305

model/I-453      305 CAUUGGUACAGAUUCUUGGUGGAUAGUAGCAAAUUAUUCGAAUGAGAGCCUUGGAGGACUGAA 365
reference/I-480  306 CAUUGGUACAGAUUCUUGGUGGAUAGUAGCAAAUUAUUCGAAUGAGAGCCUUGGAGGACUGAA 366

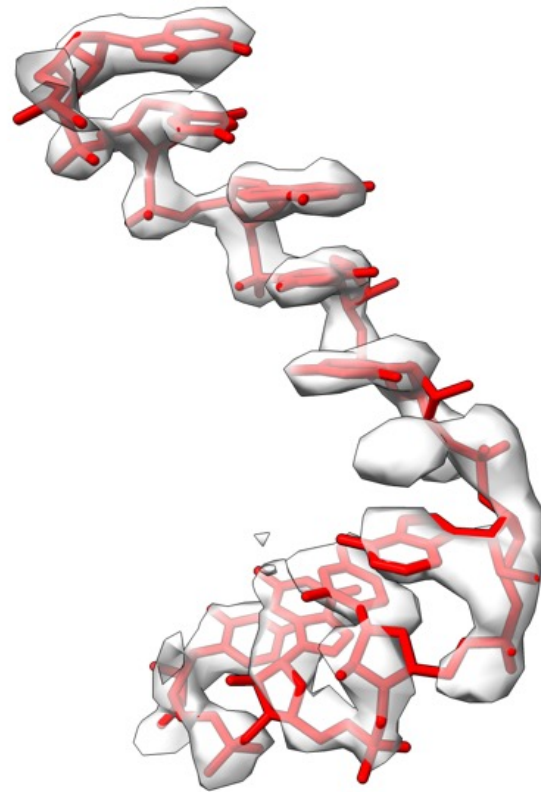
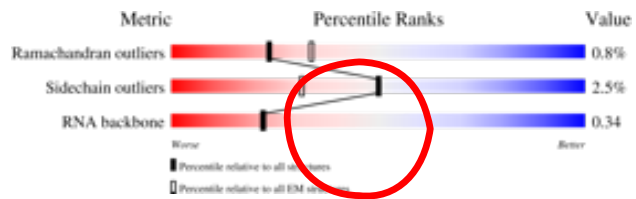
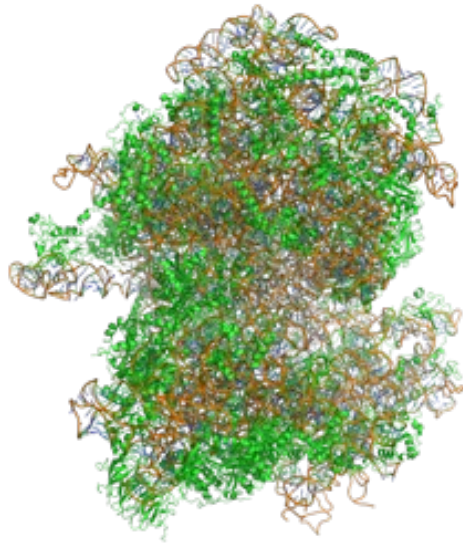
model/I-453      366 GUGGAGAAGGGUUUCGUGUGAACAGUGGUUGAUACAGAGUUAGUCGGUCCUAAAGUUCAAGG 426
reference/I-480  367 GUGGAGAAGGGUUUCGUGUGAACAGUGGUUGAUACAGAGUUAGUCGGUCCUAAAGUUCAAGG 427

model/I-453      427 CGAAA GC- GAAAAUUUUUCAAGUAAAAACA----- 453
reference/I-480  428 CGAAA GCCGAAAAAUUUUUCAAGUAAAAACAAAAAUGCCUAAACUAUUAUAAACAAAG 480
    
```

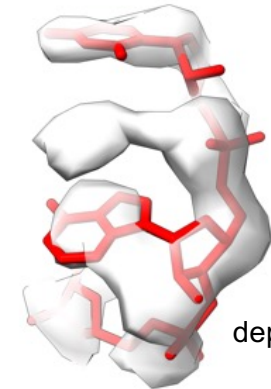


# #2 Sequence assignment issues in rRNA

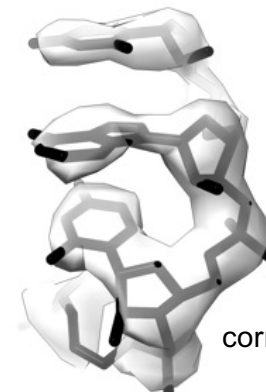
EM structure model of eukaryotic ribosome @3.5Å (deposited 2017)



fragment of 28S rRNA



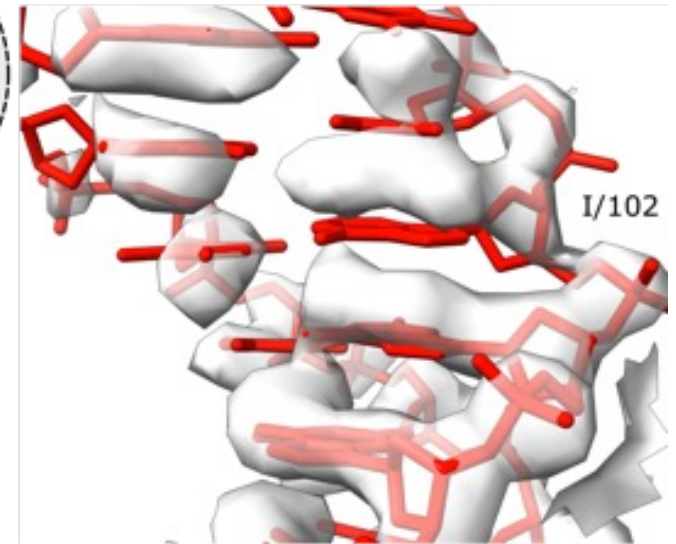
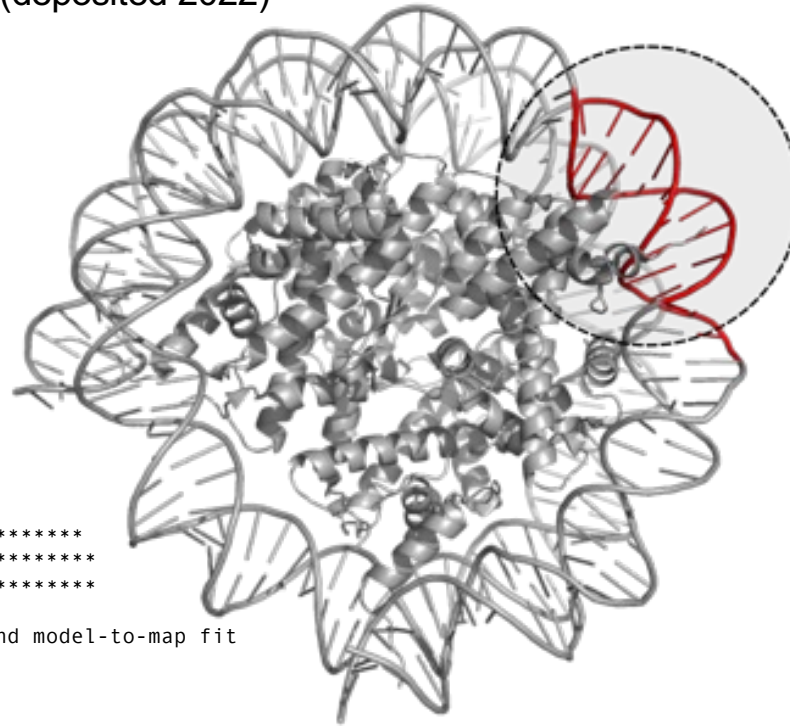
deposited



corrected

# #3 Sequence assignment issues in DNA

EM model of nucleosome at @3.4Å (deposited 2022)



```
*****  
***** 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

actagGGAGTAATCCCCTTGCGGTAAAACGCGGgggacag

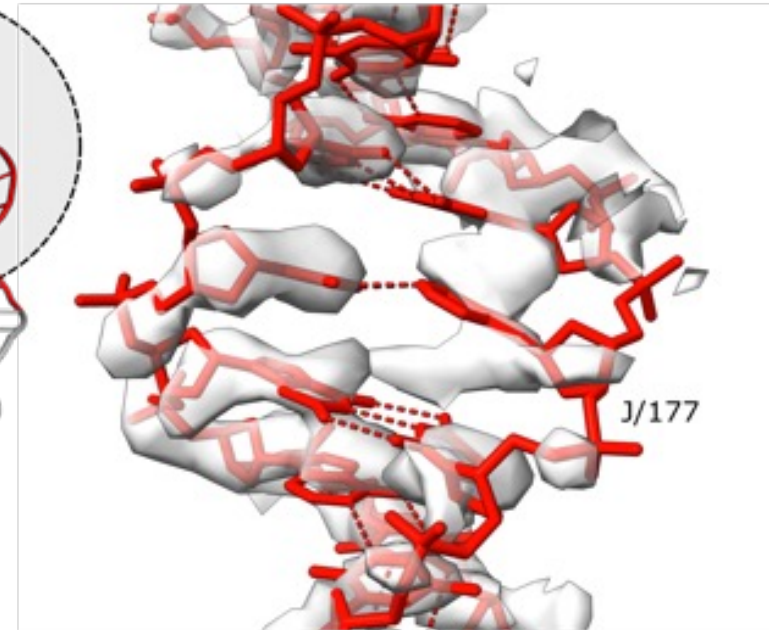
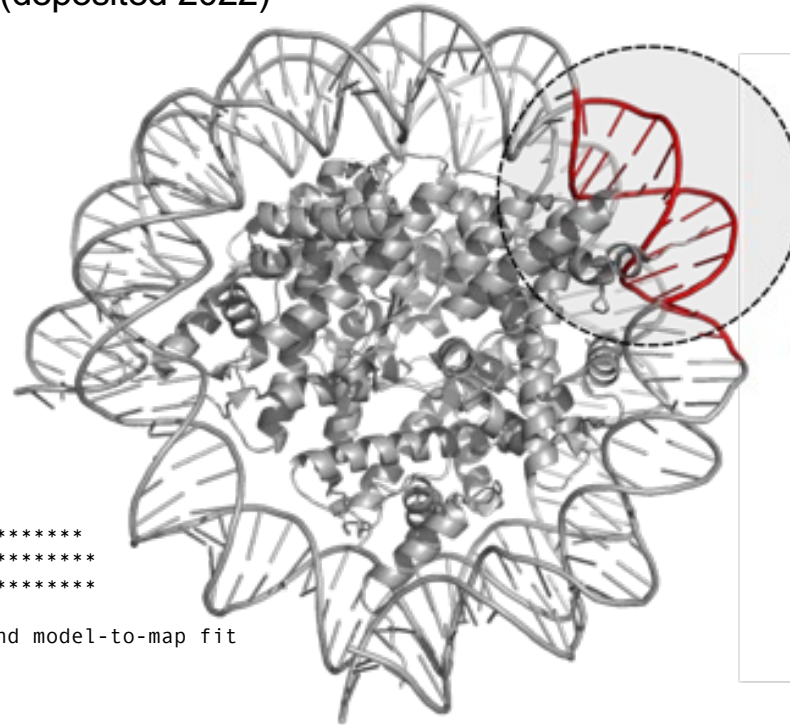
new seq 183-212

actaggGAGTAATCCCCTTGCGGTAAAACGCGGgggacag

Time elapsed 0:00:19

# #3 Sequence assignment issues in DNA

EM model of nucleosome at @3.4Å (deposited 2022)



```
*****  
***** 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

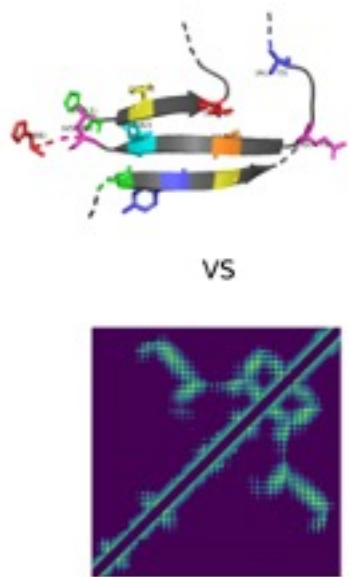
actagGGAGTAATCCCCTTGCGGTAAAACGCGGgggacag

new seq 183-212

actaggGAGTAATCCCCTTGCGGTAAAACGCGGgggacag

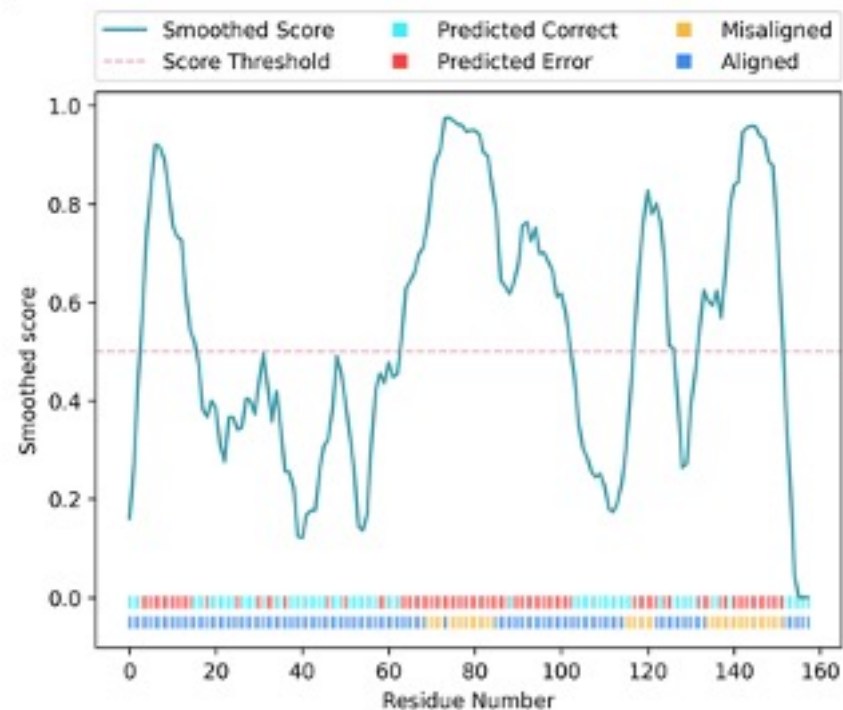
Time elapsed 0:00:19

# Predicted inter-residue distances for model validation



FEATURES	RESIDUE NUMBER				
	1	2	3	4	5
RMSD	0.3	2.5	0.5	1.2	0.6
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

SVM CLASSIFIER



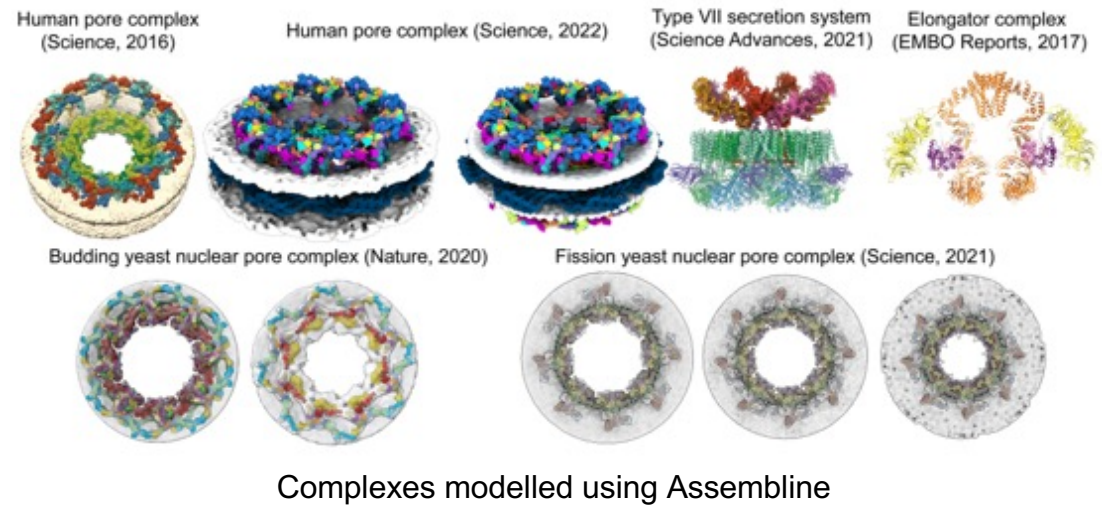
Filomeno Sanchez  
Ronan Keegan  
Daniel Rigden

[www.conkit.org](http://www.conkit.org)

*Filomeno Sánchez Rodríguez et al Acta Cryst. (2022). D78*

# Assembleline

## Assembly line of macromolecular complexes



Jan Kosinski group at EMBL Hamburg

<https://www.embl-hamburg.de/Assembleline/>

*Nature Protocols, 2022*



# Acknowledgements

## University of Liverpool

Daniel Rigden  
Adam Simpkin  
Filomeno Sánchez Rodríguez

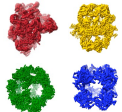
## CCP4 Core Team

Ronan Keegan  
Charles Ballard  
Eugene Krissinel  
Kyle Stevenson  
Maria Fando

## CCPEM team

Agnel Joseph  
Tom Burnley  
Colin Palmer  
Matt Iadanza


## Martin Luther University

 Panos Kastritis  
Ioannis Skalidis

EMBL  
Hamburg 

Matthias Wilmanns  
Kate Beckham  
Jan Kosiński  
Christina Ritter  
Edukondalu Mullapudi  
Isabel Bento  
Alice Bochel


## findMySequence

protein sequence identification  


## checkMySequence

sequence validation  


## doubleHelix

NA sequence identification  


EMBL 