



Science and
Technology
Facilities Council

Scientific Computing

Atomic Model Validation

Agnel Praveen Joseph
CCP-EM STFC

08/11/2024

Icknield workshop



Outline

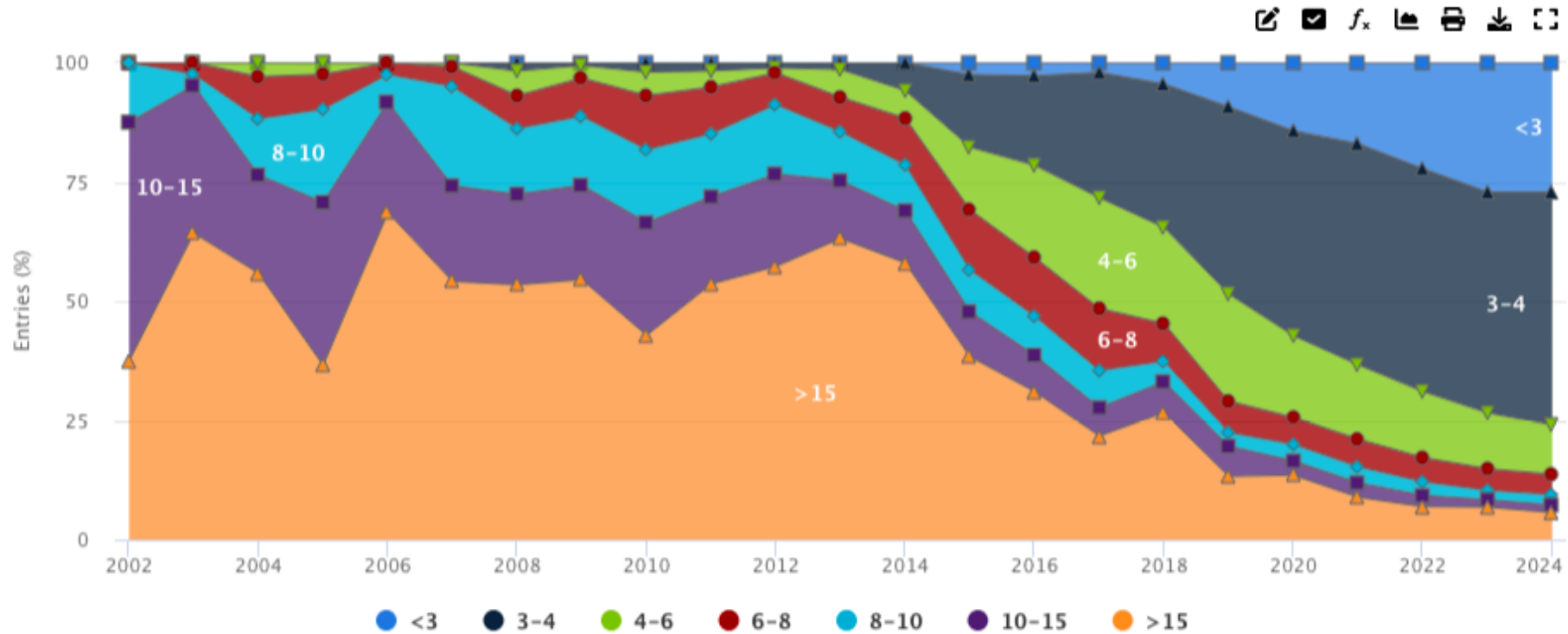
- 1 Background
- 2 Validation: Why is it important?
- 3 Tools for atomic model validation
- 4 Some examples
- 5 Model validation in Doppio



EM DataBank: map resolution

>39k cryo-EM reconstructions in EMDB

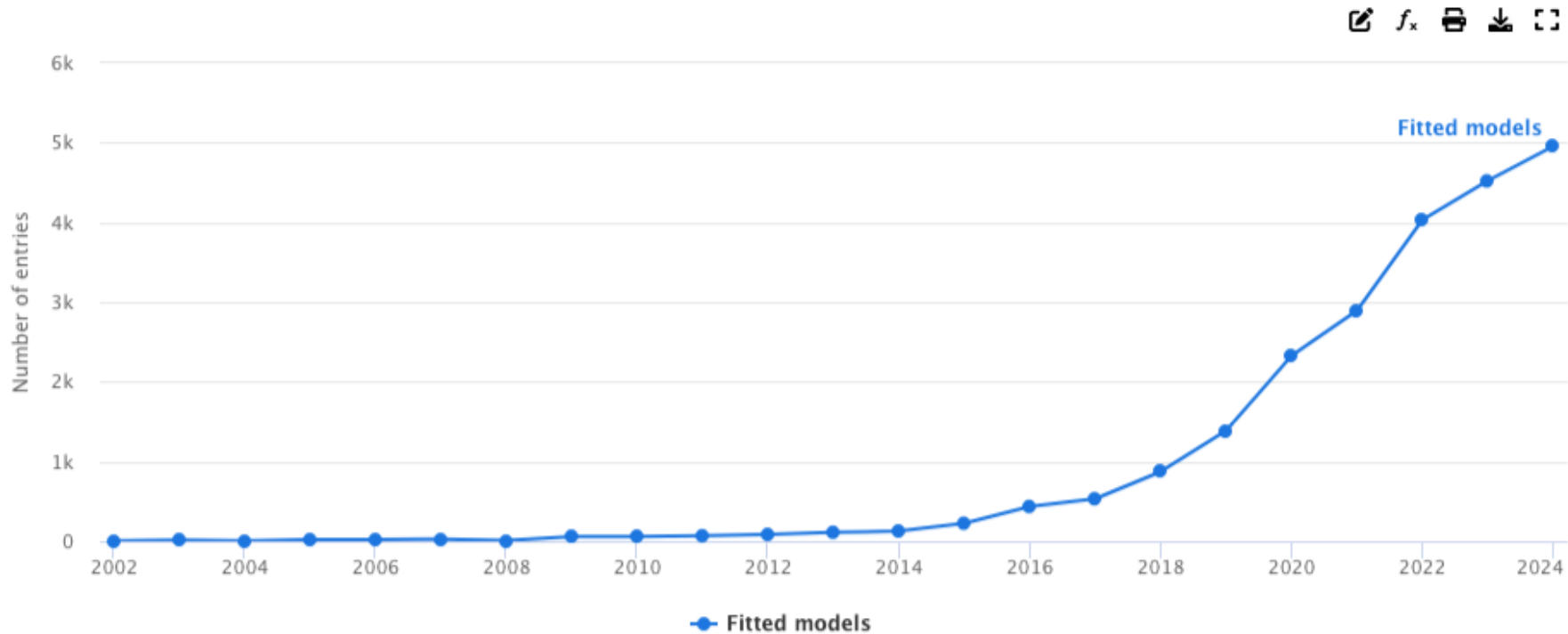
EMDB entry resolution in shells per year



Last 3 years
3-6 Å: ~60%

EMDB: maps with fitted models

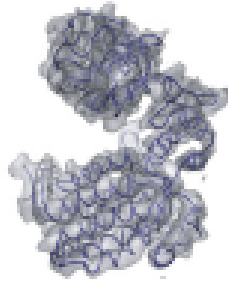
EMDB entries with fitted models in PDB per year



Resolution

- $< 2\text{\AA}$ [197]
- $\geq 2\text{\AA}$ and $< 3\text{\AA}$ [5360]
- $\geq 3\text{\AA}$ and $< 4\text{\AA}$ [12130]
- $\geq 4\text{\AA}$ and $< 5\text{\AA}$ [2728]
- $\geq 5\text{\AA}$ and $< 6\text{\AA}$ [448]
- $\geq 6\text{\AA}$ and $< 8\text{\AA}$ [865]
- $\geq 8\text{\AA}$ and $< 12\text{\AA}$ [611]
- $\geq 12\text{\AA}$ and $< 16\text{\AA}$ [171]

Features interpretable at different resolutions

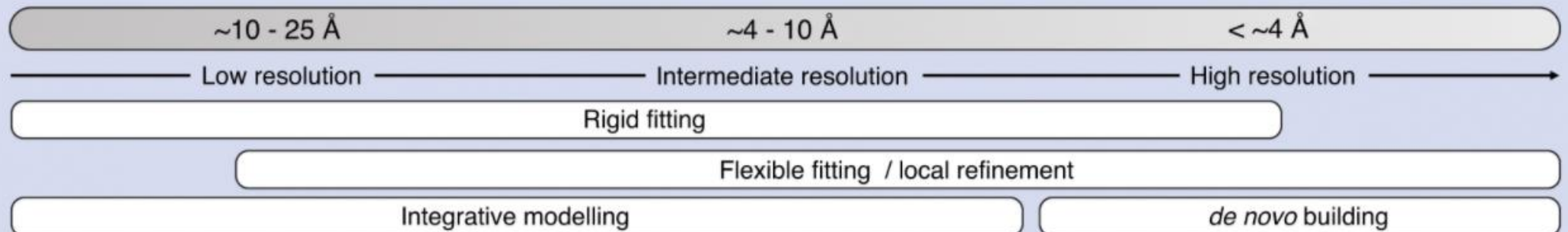


Villa et al. 2014

Up to 20 Å	Up to 9 Å	Up to 6 Å	Up to 4 Å
Conformational changes	Conformational changes	Conformational changes	Conformational changes
Domain boundaries	Domain boundaries	Domain boundaries	Domain boundaries
		Beta sheets	Beta sheets
	Alpha helices	Alpha helices	Individual beta strands
		Pitch of RNA helices	Alpha helices
			Pitch of alpha helices
			Pitch of RNA helices
			Phosphate "bumps"
			Side chains

Malhotra S. et al., COSB 2019

Resolution-based workflow





Science and
Technology
Facilities Council

Scientific Computing

Model validation



Science and
Technology
Facilities Council

Scientific Computing



Aspects of validation

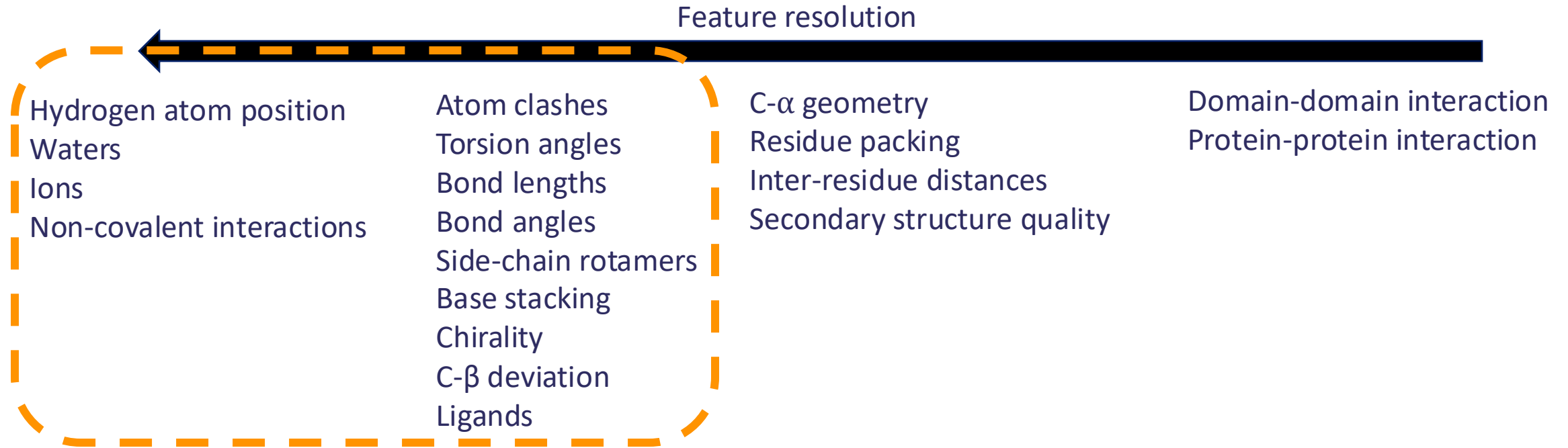
Model quality (geometry)

Fit to data (Global and Local)

Overfitting

Model bias

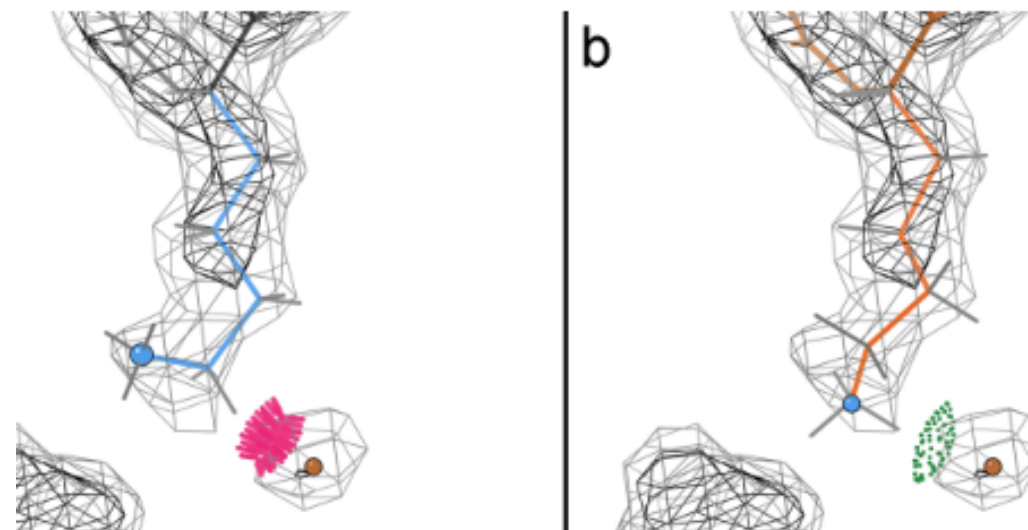
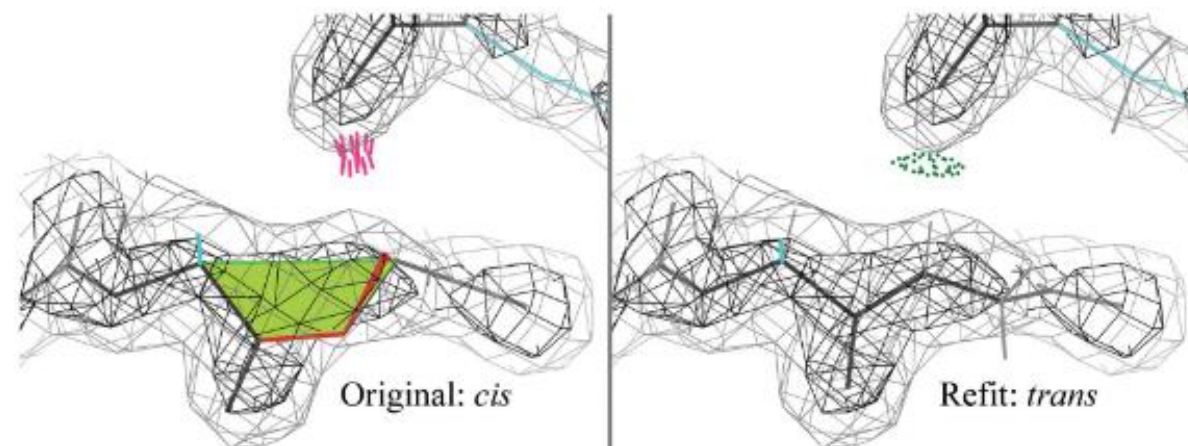
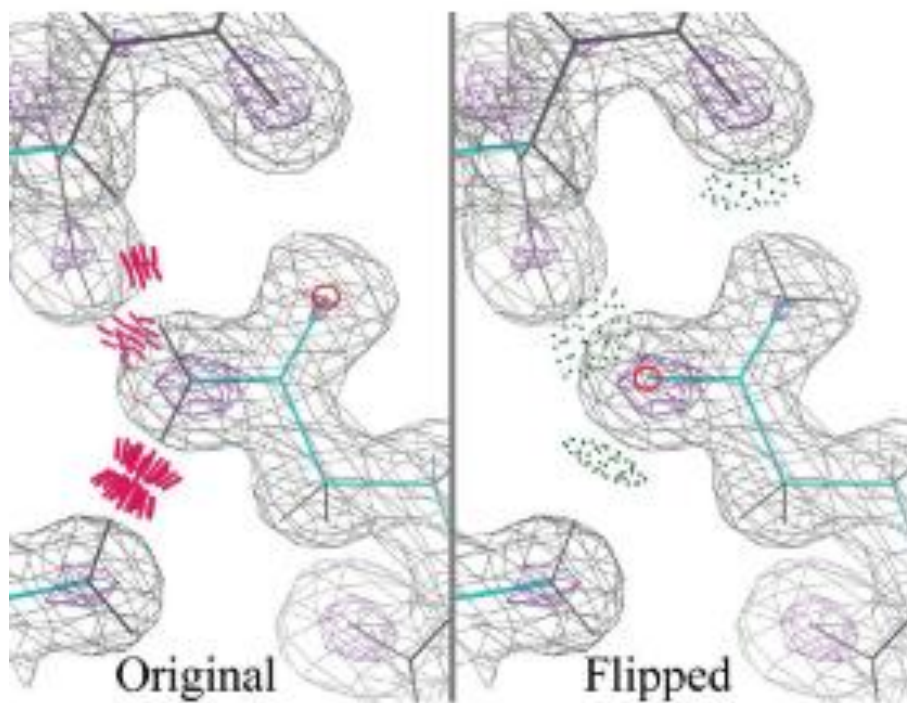
Model quality (geometry)



Molprobit, Coot/Moorhen, Isolde, PDB-REDO and others

Model quality (geometry)

Clashes



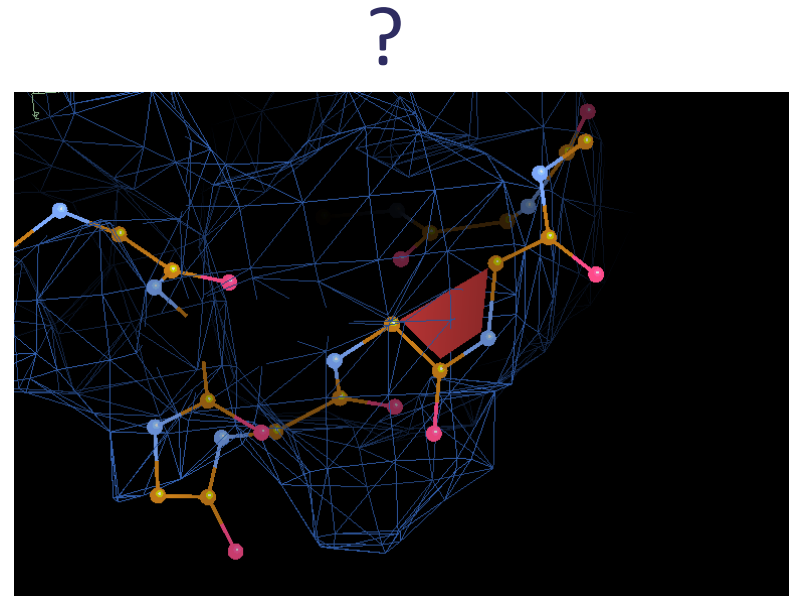
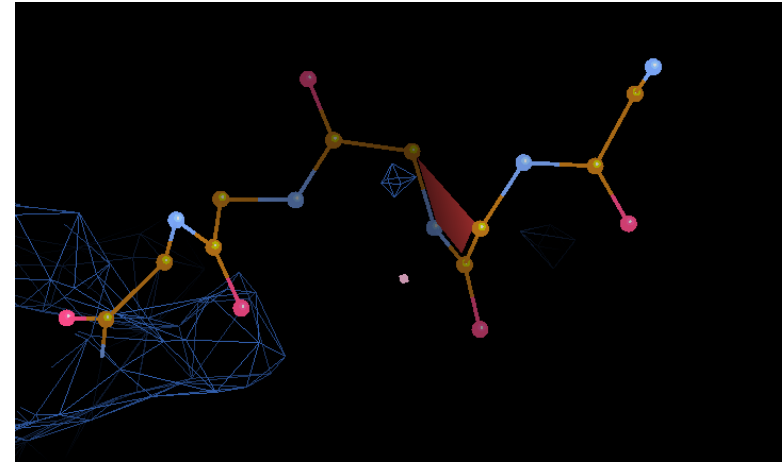
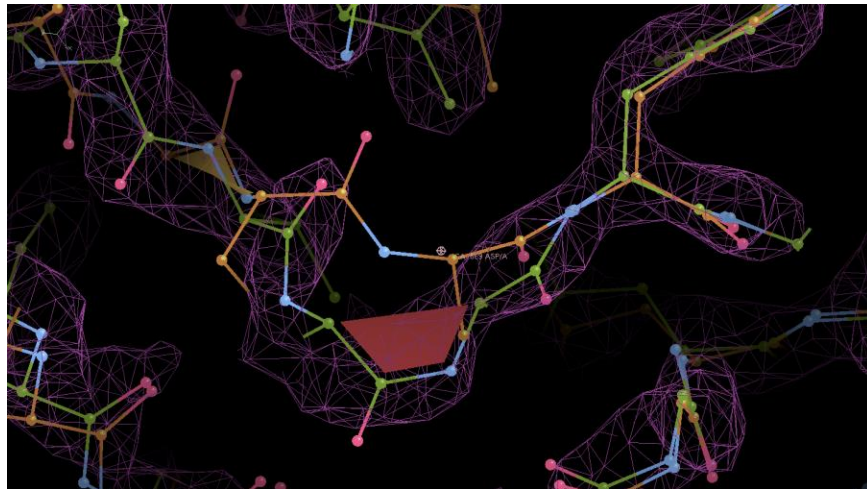
Model quality (geometry)

Cis-peptides

Cis-Pro : 5-6%

Cis-nonPro: ~0.05%

genuine

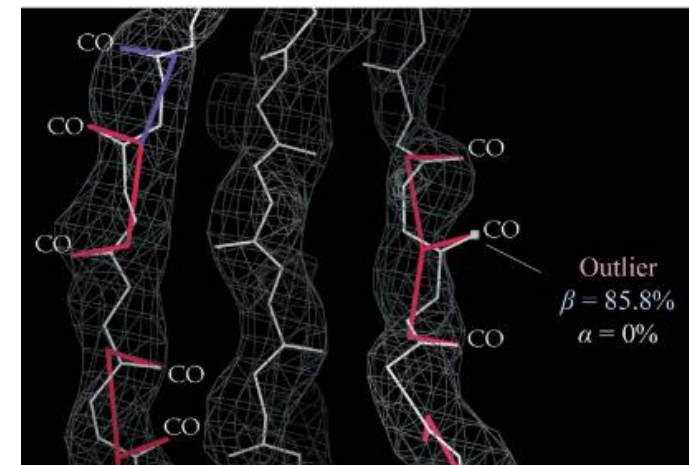
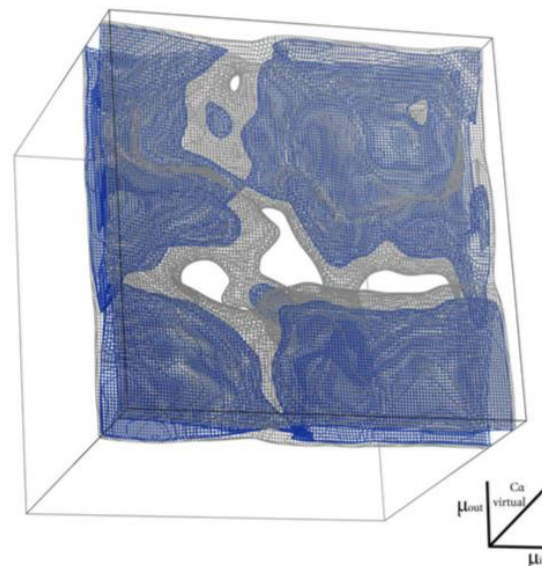
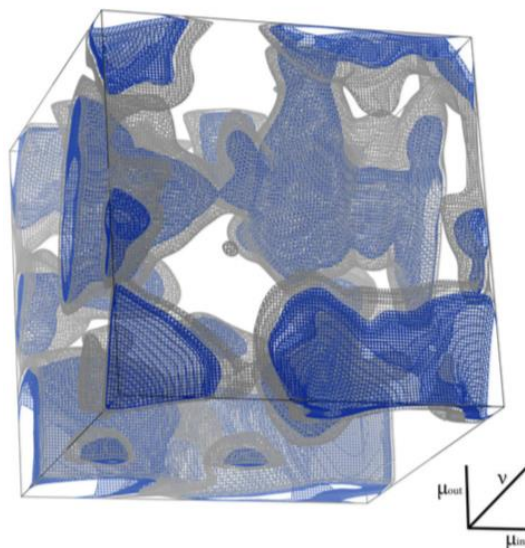
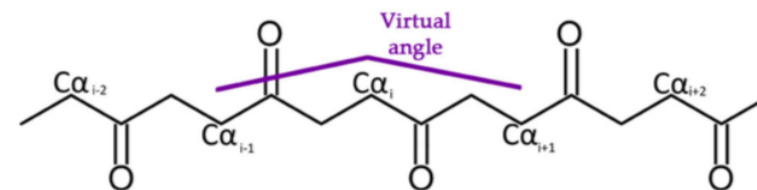
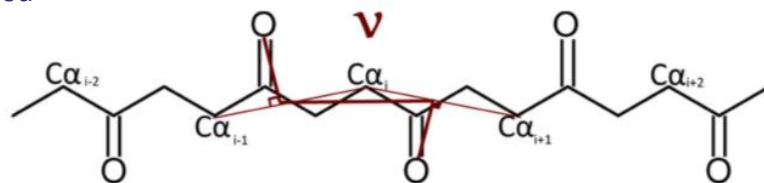
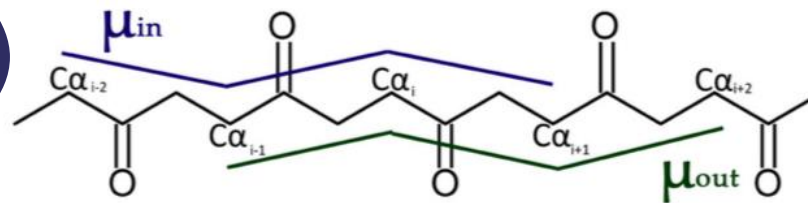


Model quality (geometry)

CaBLAM

C-Alpha Based Low-resolution Annotation Method

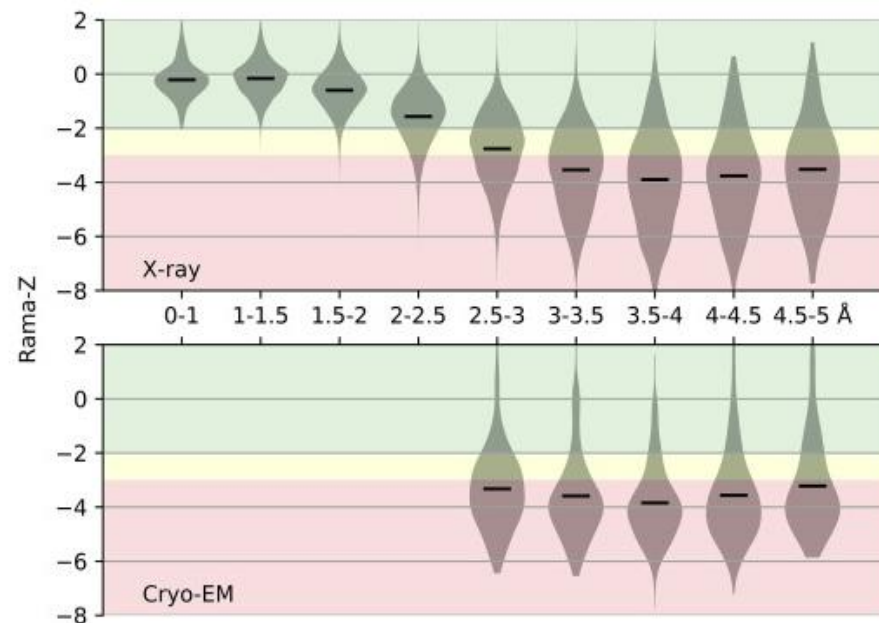
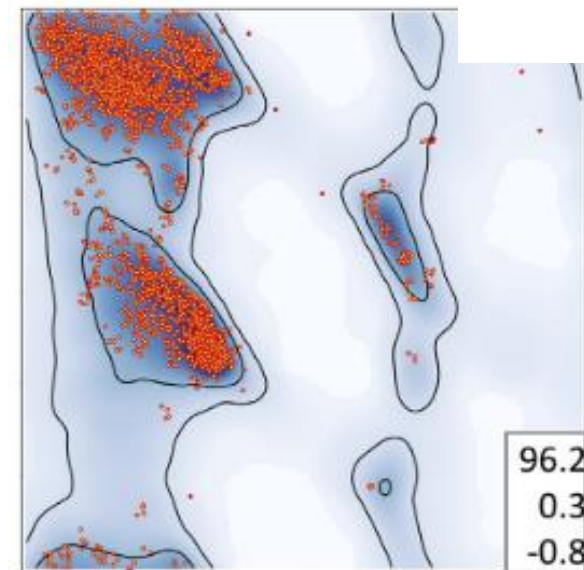
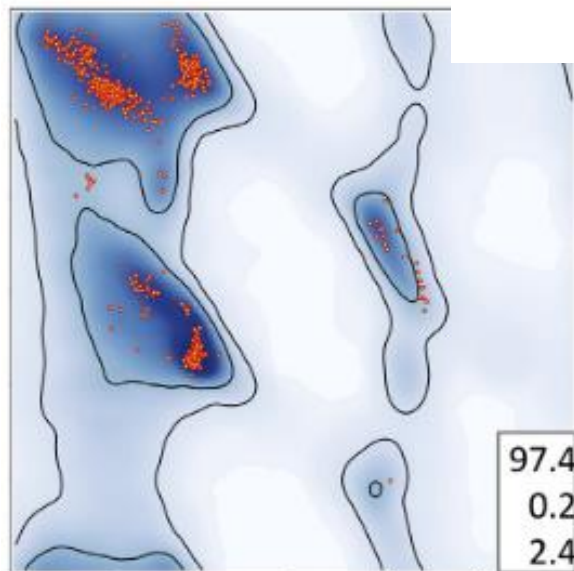
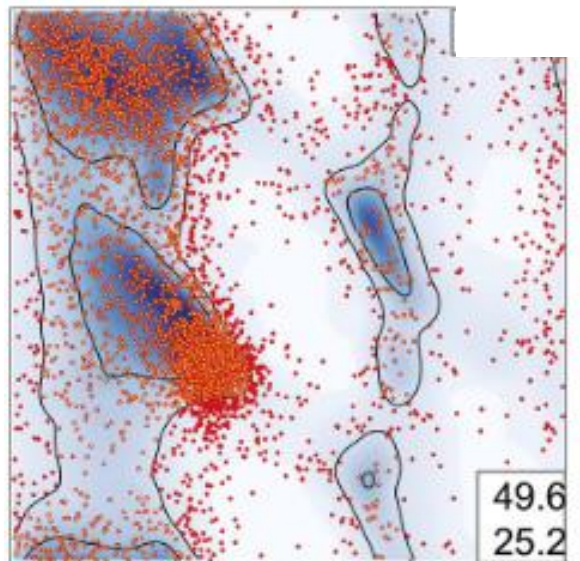
Backbone
geometry
check



Model quality (geometry)

Ramachandran z-scores

Ramachandran outliers



Fit to data

How well the model agrees with data both globally and locally ?

Is the model traced/fitted to the background rather than the molecular volume?

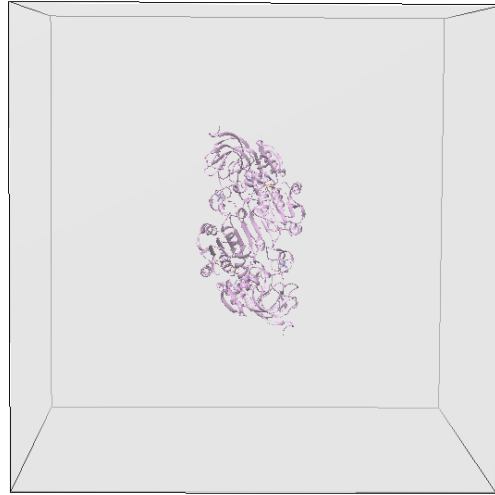
Is the model overfitted (to noise in the map) ?

Global agreement with map

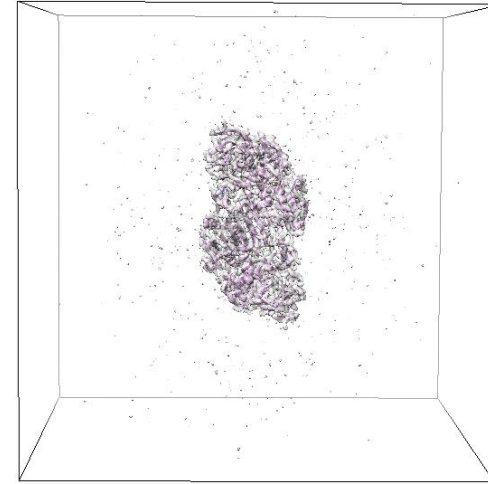
Cross-Correlation Coefficient

Values vary depending on map processing, resolution

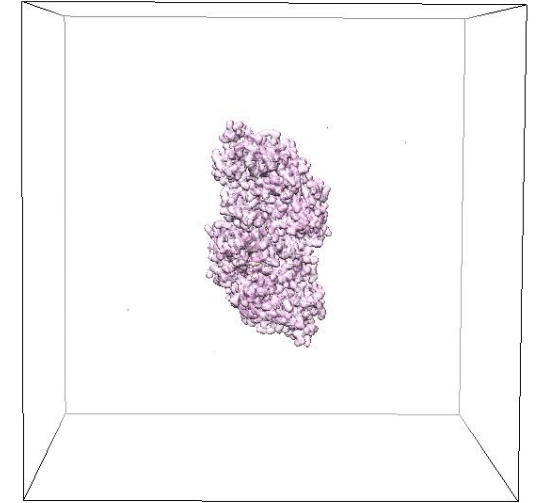
2.9Å full map: 0.32



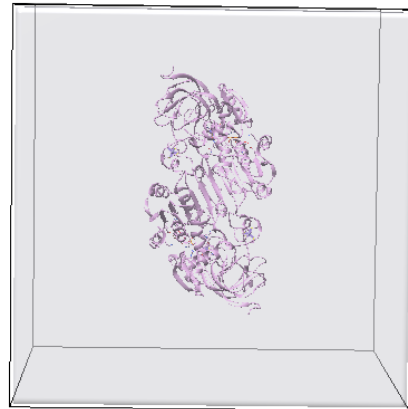
Contoured map: 0.59



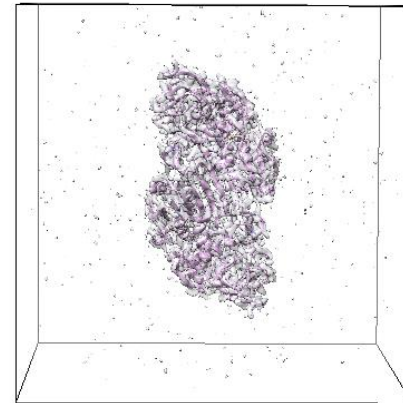
Map-model overlap mask: 0.42



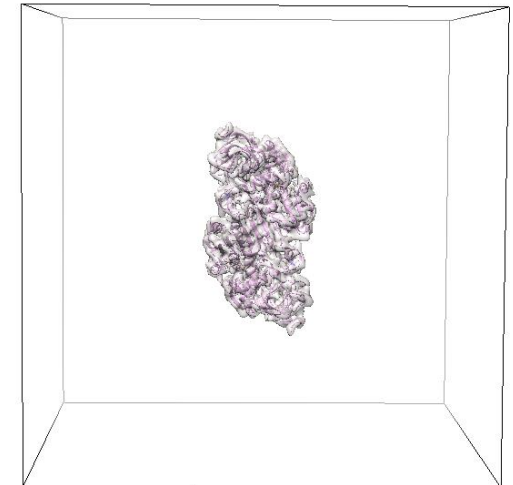
Cropped full map: 0.37



Cropped,contoured map: 0.55



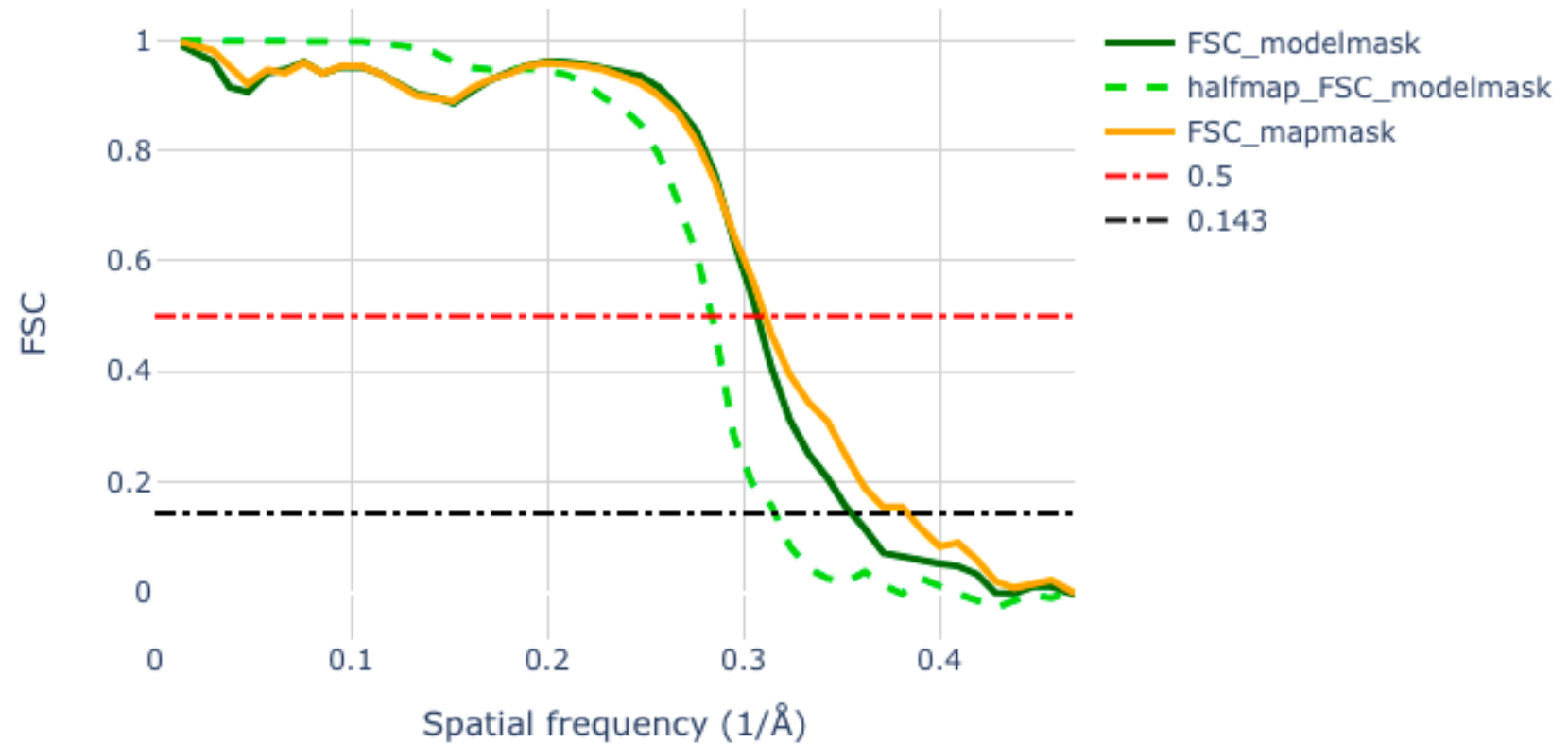
5Å contoured map: 0.76



Global agreement with map

Model-map FSC

model-map FSC plot

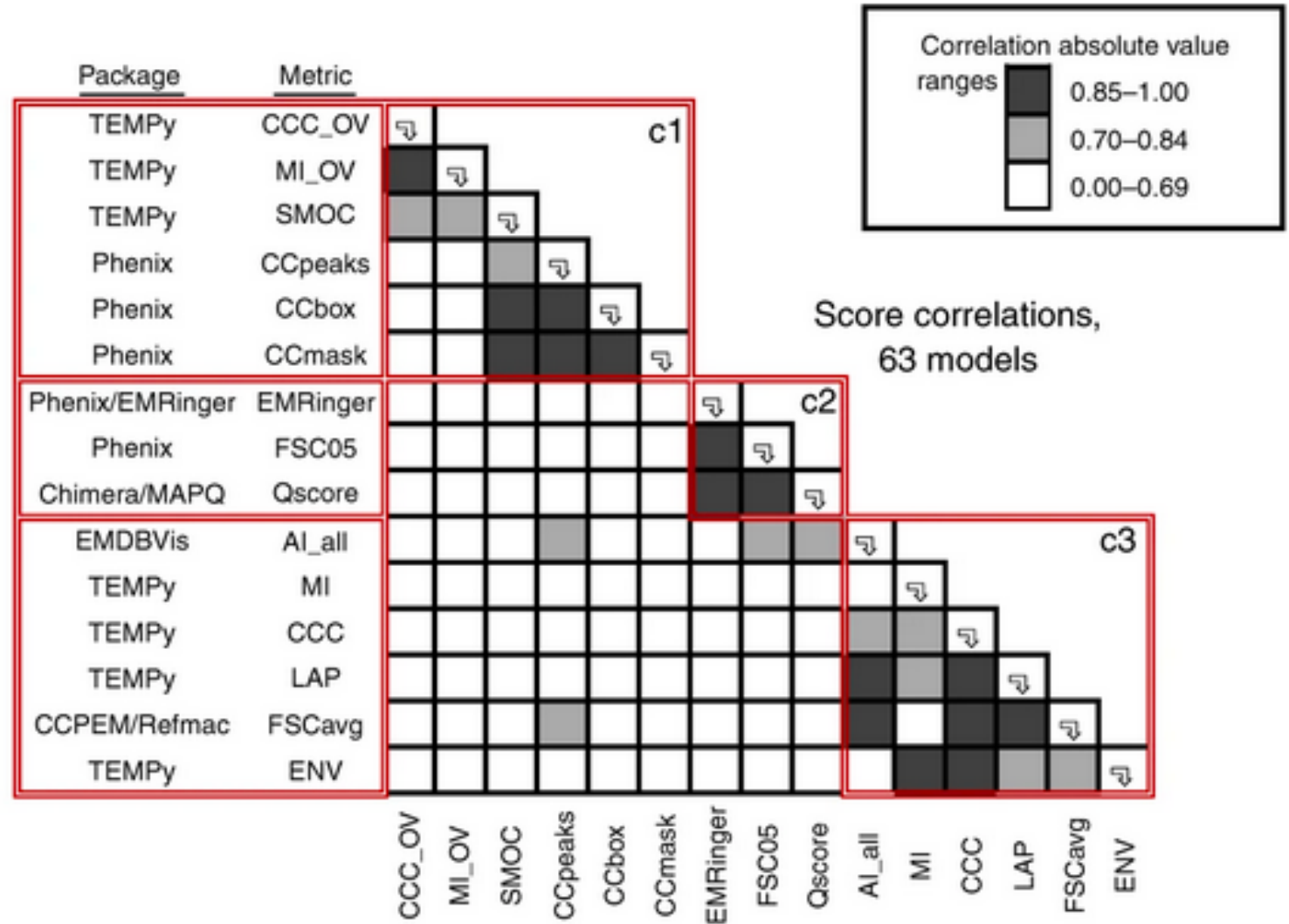


Brown et al. Acta D, 2015

Yamashita et al. Acta D, 2021

Rosenthal and Henderson JMB 2003

Global agreement with map



Lawson et al. *Nature Methods*, 2021

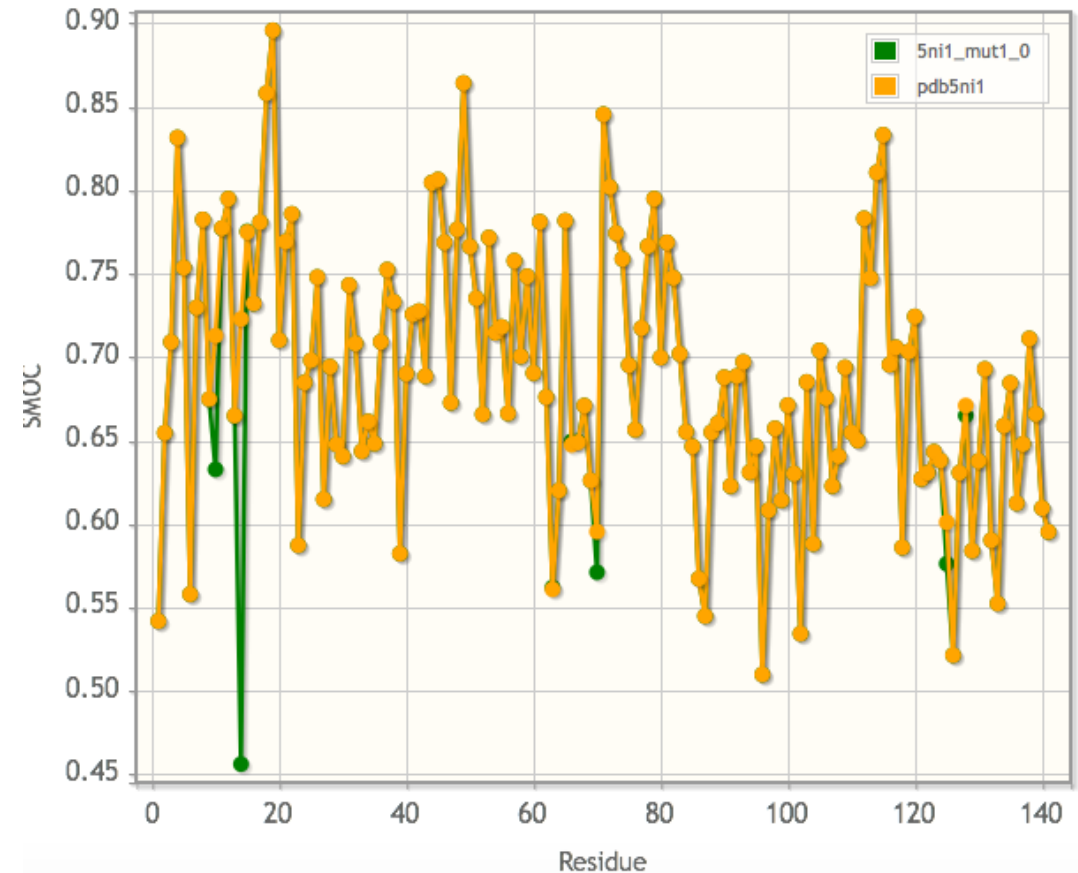
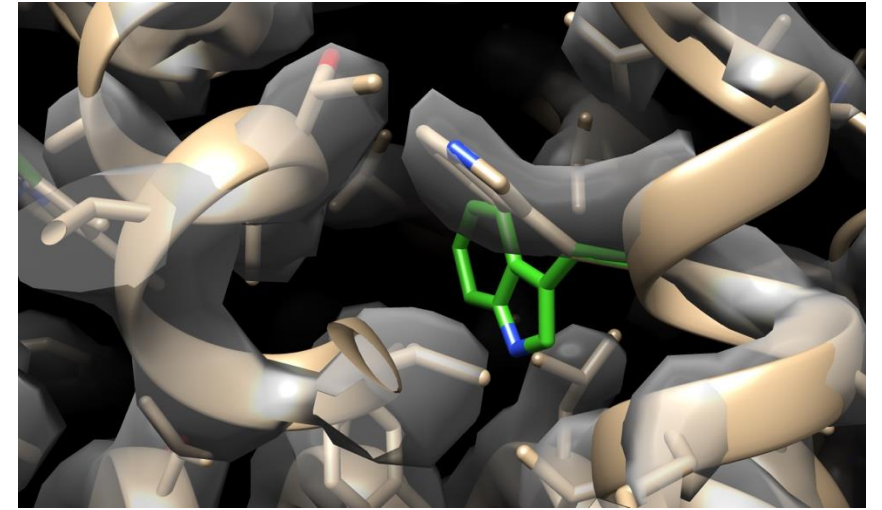
Local agreement with map

Segment based Manders' Overlap Coefficient (SMOC)

An overlap coefficient is calculated over voxels covered by each residue (and the local neighborhood)

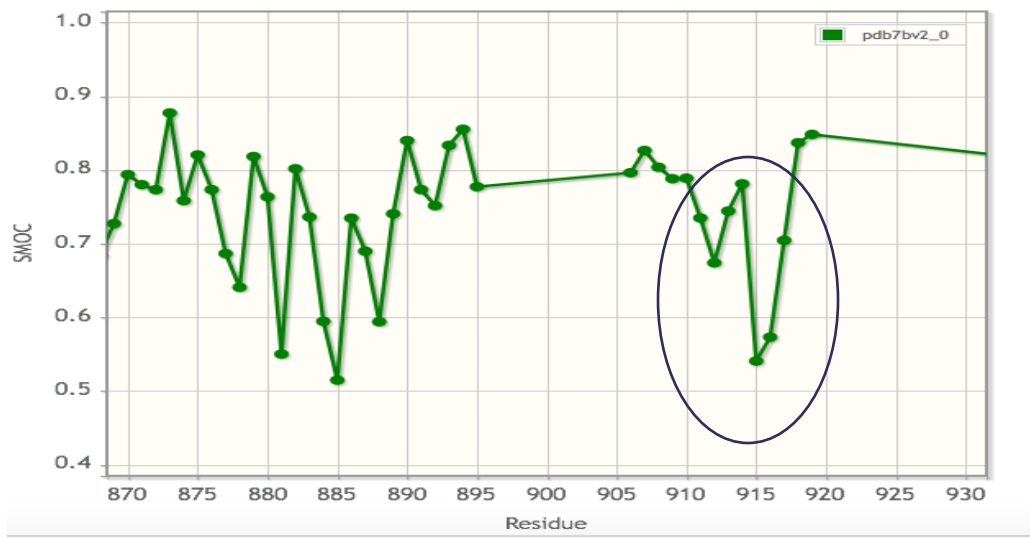
$$SMOC = \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n X_i^2} \sqrt{\sum_{i=1}^n Y_i^2}}$$

Joseph *et al. Methods* 2016 , Farabella *et al. Acta D* 2015



Local agreement with map

2.5Å SARS-CoV2 RNA pol



Chojnowski G. Acta Cryst 2022

CheckmySequence Chain: A

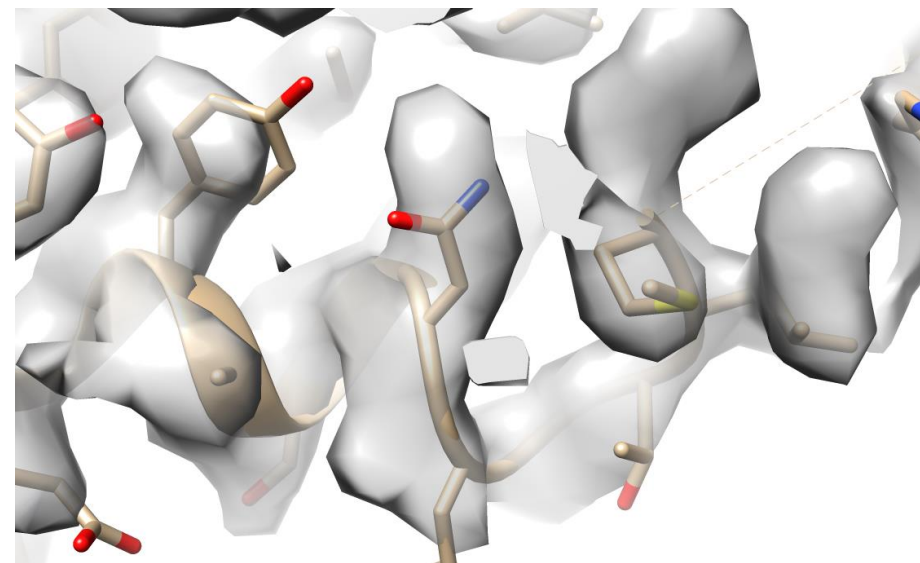
*** Register shift:

protein start: 906 end: 919 start_new: 915 end_new: 928 -log(pvalue): 1.0068789995424352 si: 100.0

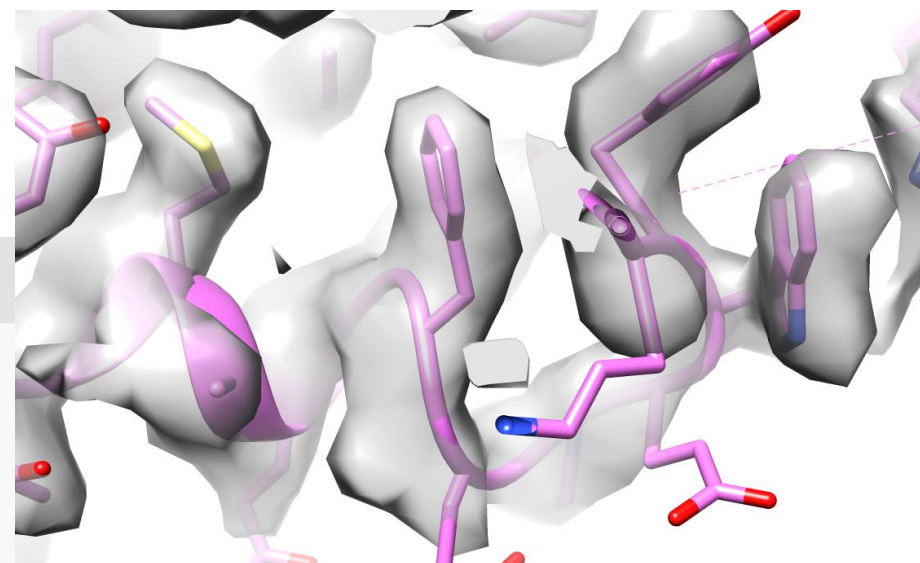
model seq: dmysvMLTNDNTRSRYWEPEfyeamytphtvlqgg

new seq : dmysvmltndntsrYWEPEFYEAMYPHTvlqgg

Deposited

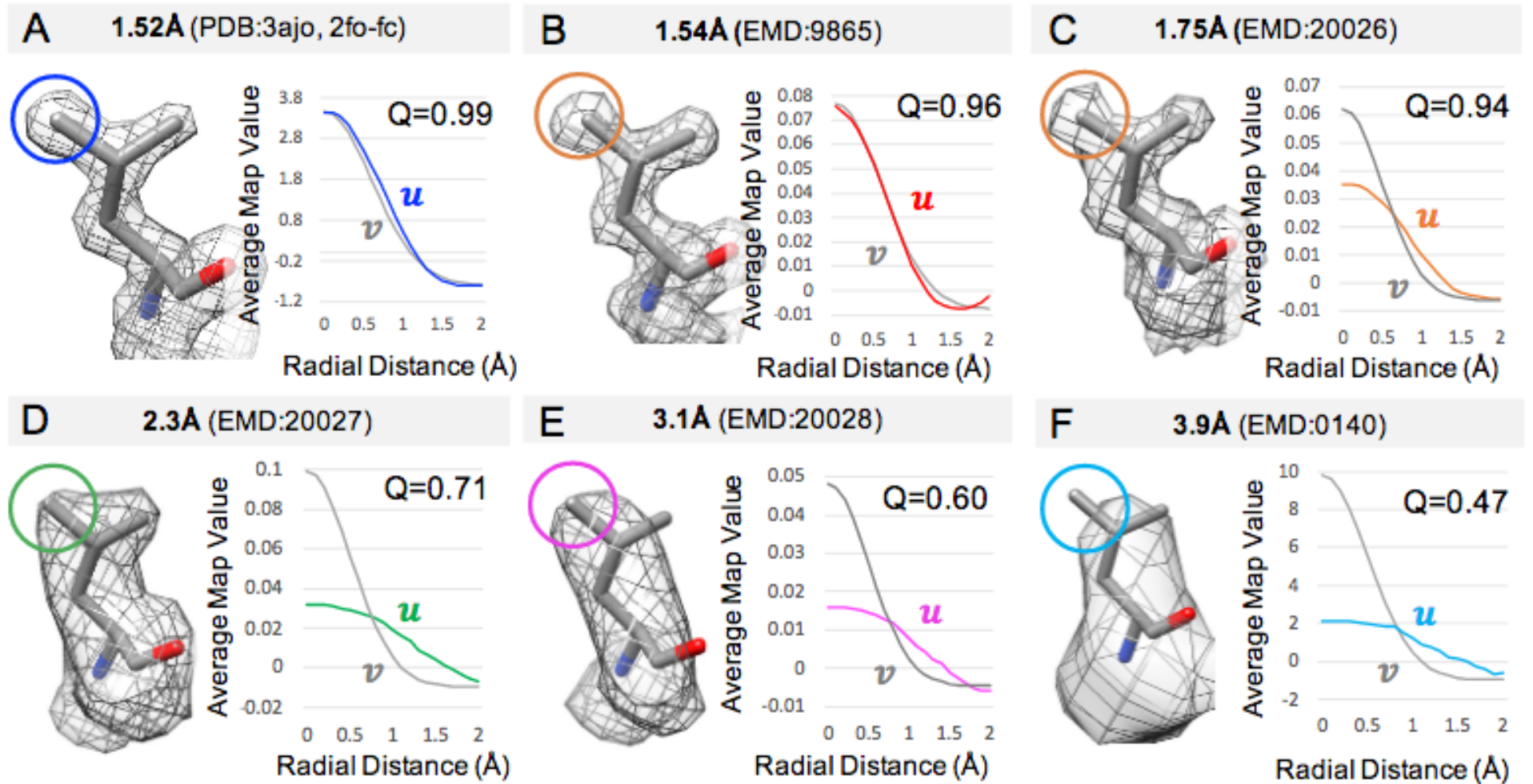


Remodelled



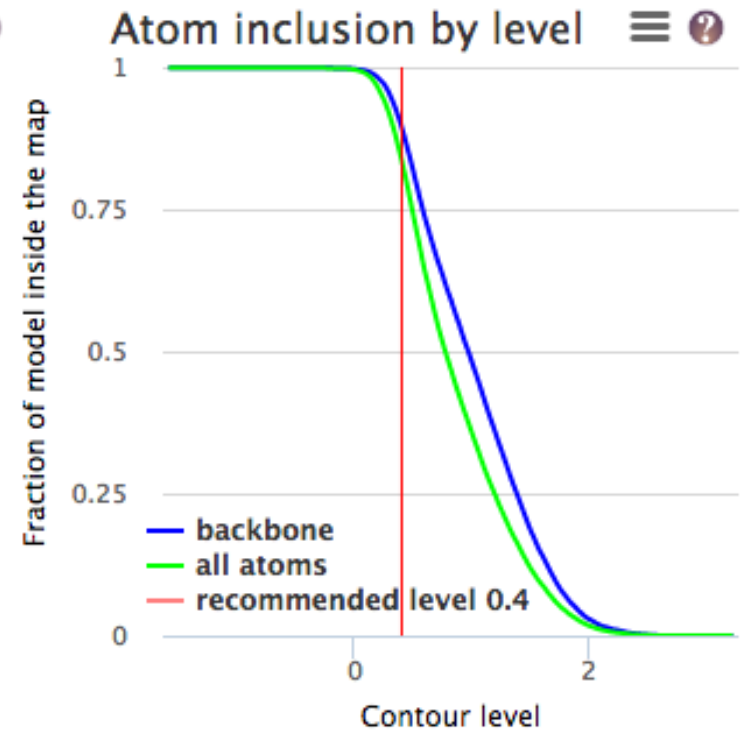
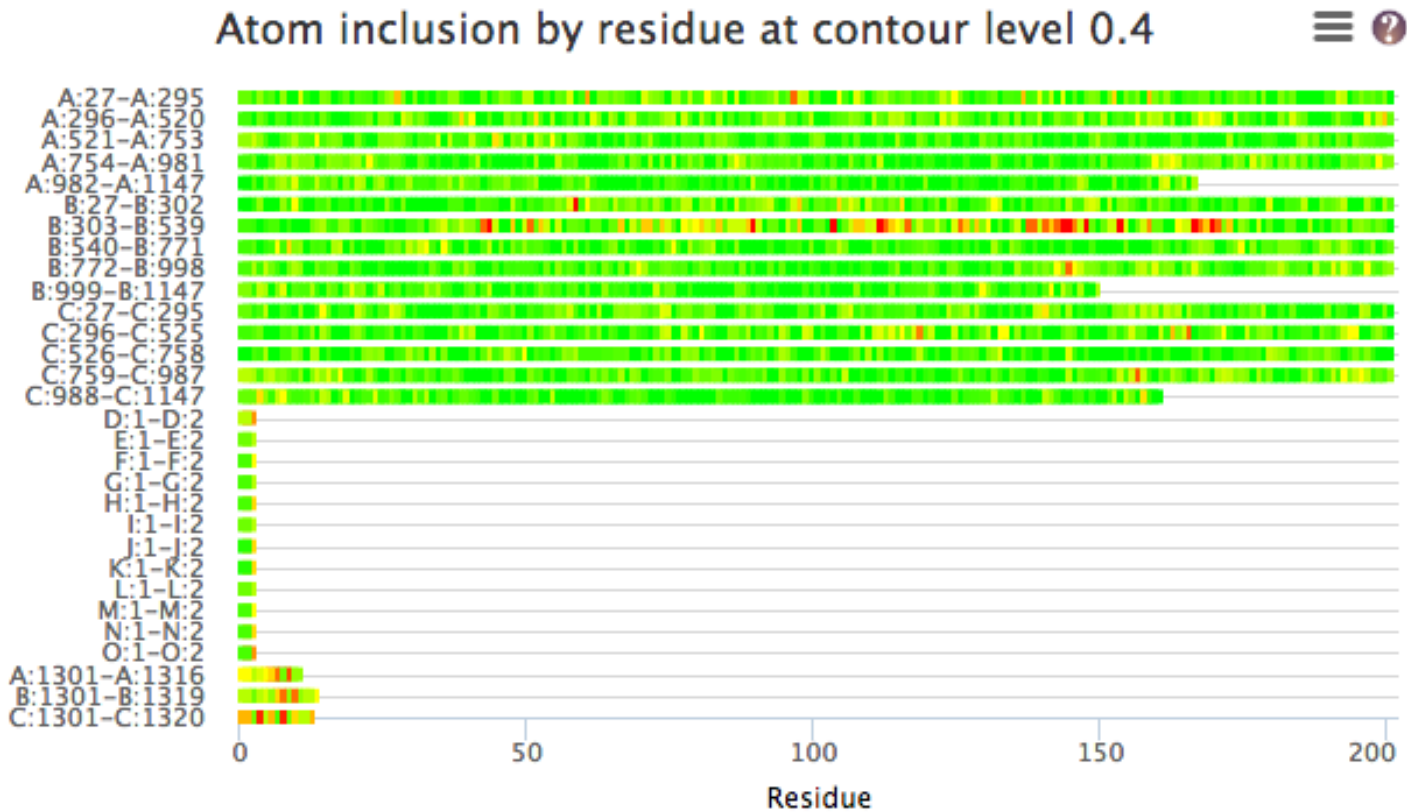
Local agreement with map: resolvability

Q score



Residue position/trace

Atom inclusion score



Backbone trace?

Lawson et al. *Nature Methods*, 2021

FDR validation score

MapQ
SMOC
SCCC
PHENIX
FSC_Q
FDR_score

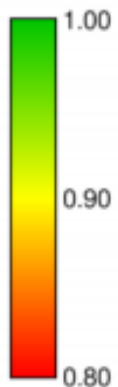
Value	Z-score
0.20	-3.34
0.78	-1.40
0.60	-1.25
0.73	-0.76
1.37	4.89
0.49	-4.85

TRP188

T0007EM192_2



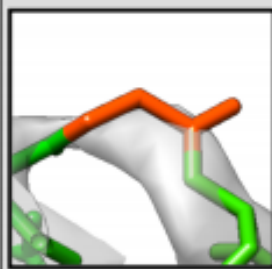
5a63



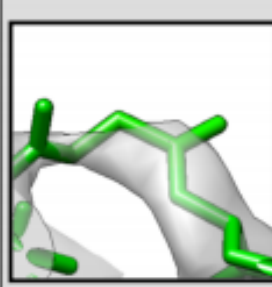
Value	Z-score
0.53	-0.66
0.70	-0.29
0.55	-0.64
0.73	-0.34
0.25	0.11
0.83	-1.23

GLY86

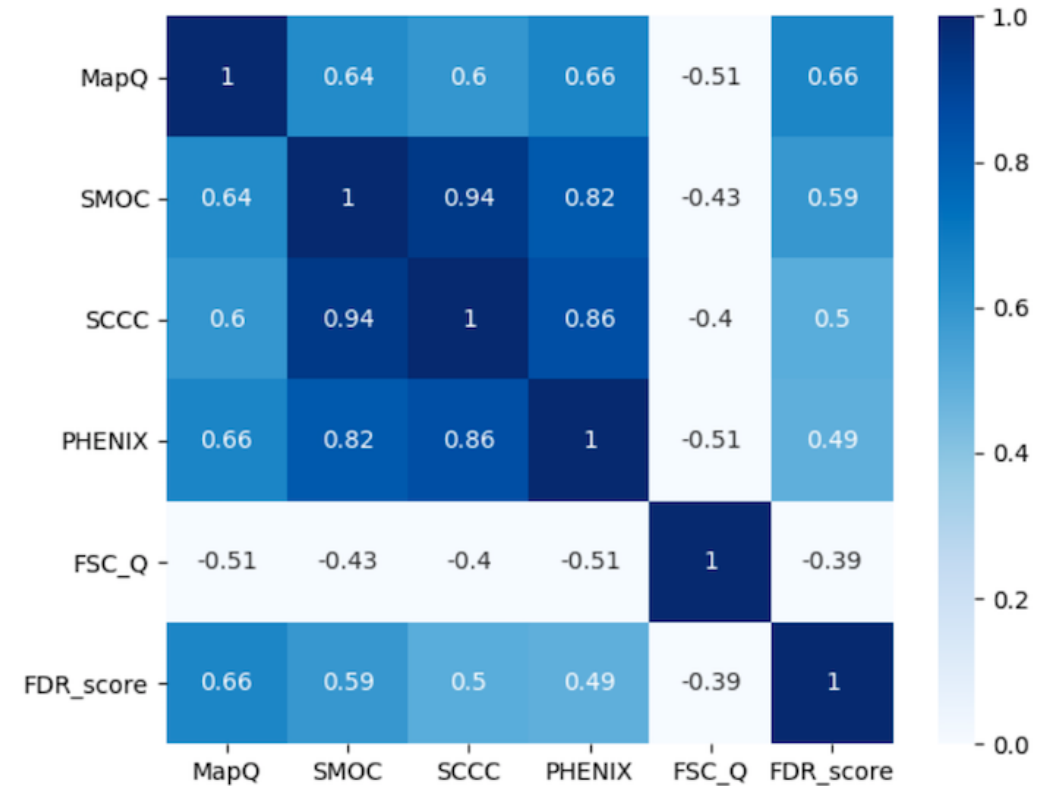
T0104EM060_2



6nbb.2



T0104EM060_2



Model geometry vs fit-to-data

Geometry

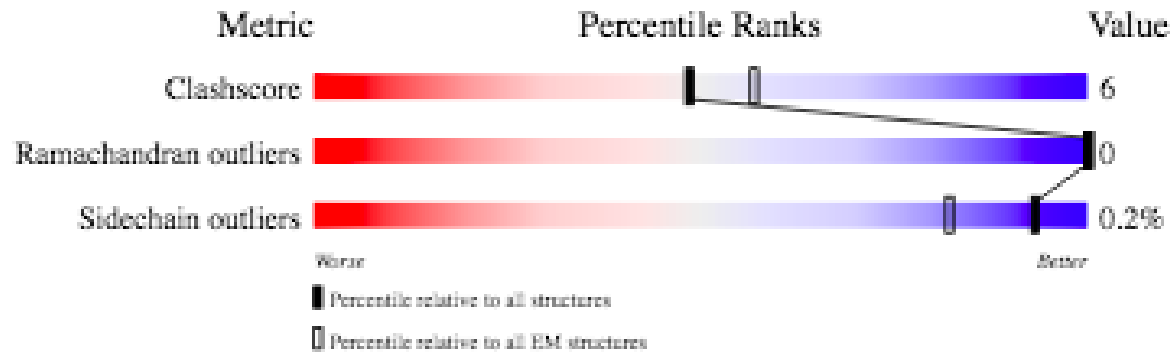
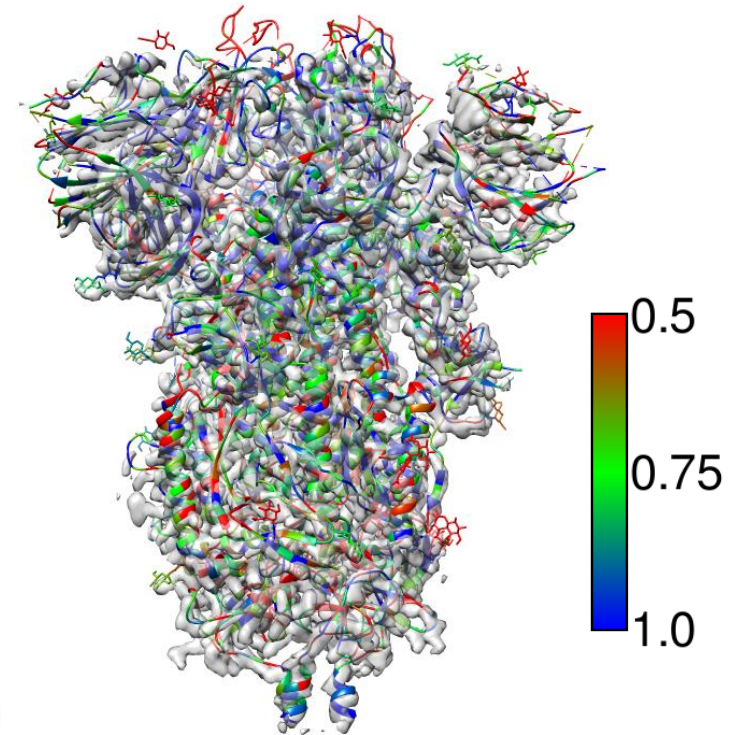


Image: PDBe/wwPDB

Map agreement ?



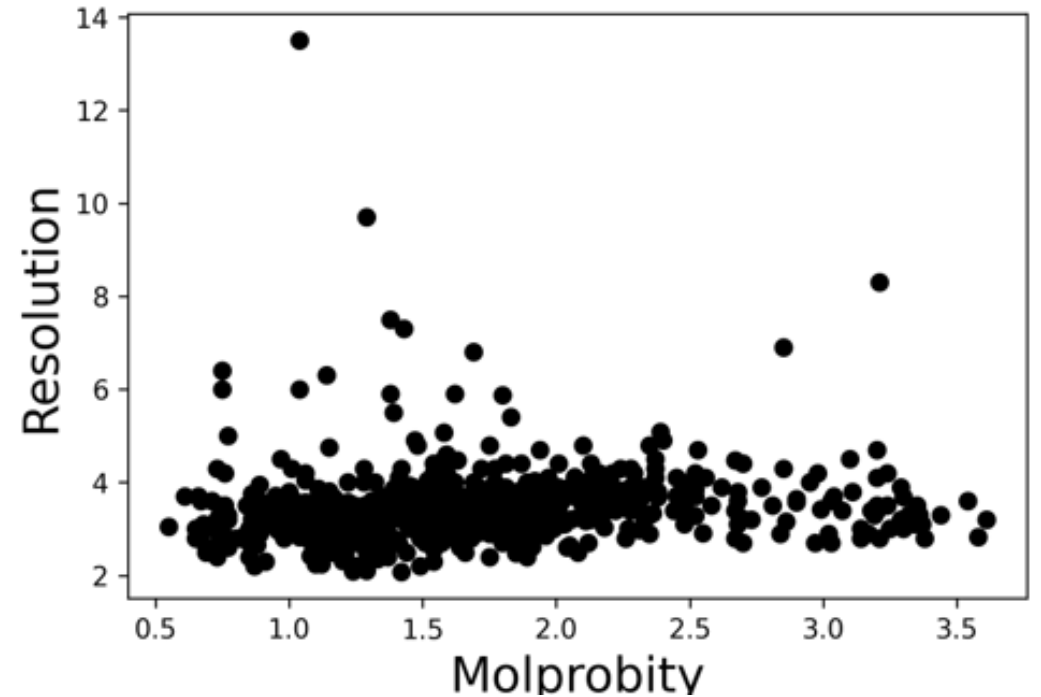
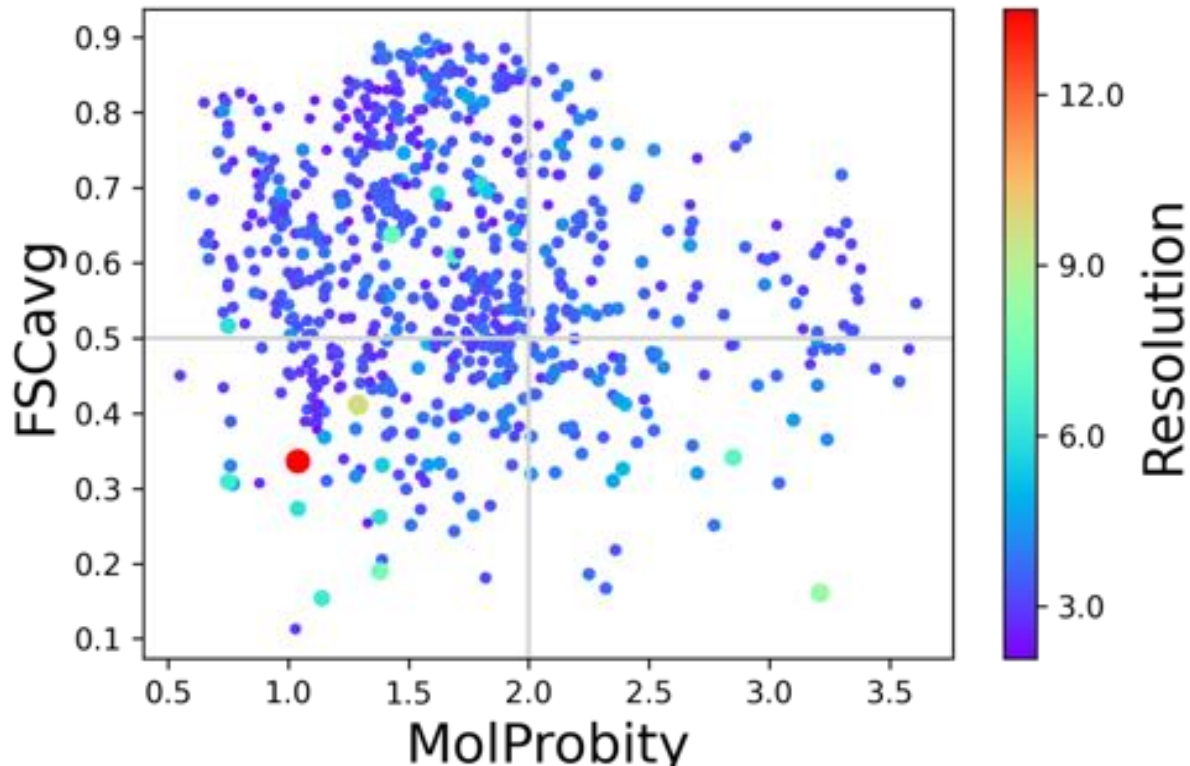
Overfitting ?



Model geometry vs fit-to-data

Do the models represent the map data well?

Do we need more validation metrics for publication?



Mean score of structures worse than 3.5Å resolution is 1.8 (< 3.5Å is 1.6).

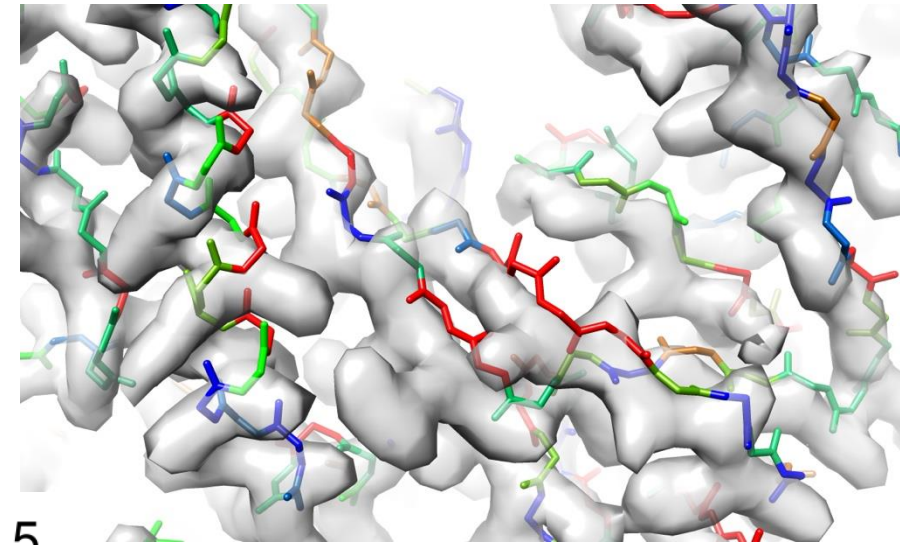
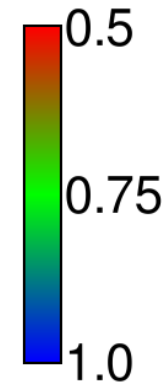
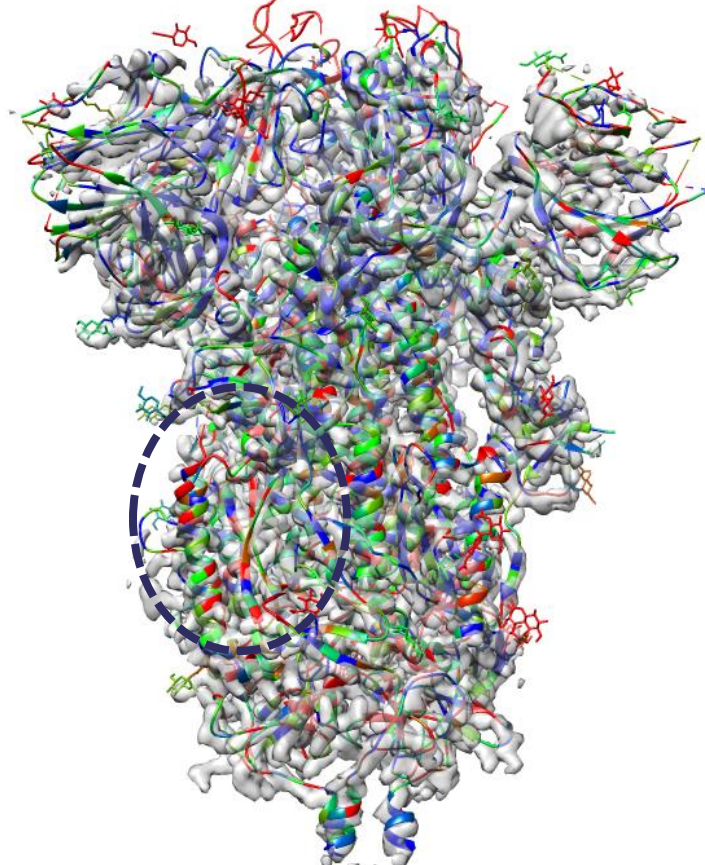
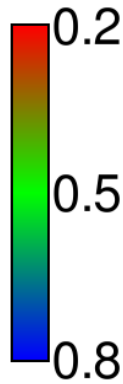
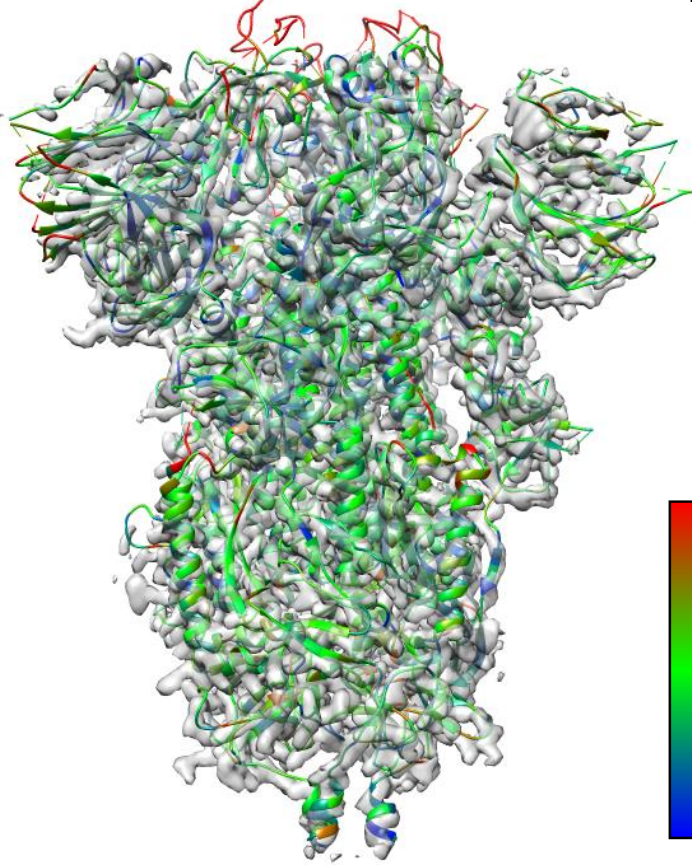
31.2% of the structures had FSCavg scores worse than 0.5.

Poor fit

~28% of structures have potential issues with backbone trace (FDR-backbone score)

Molprobtity score : 1.39

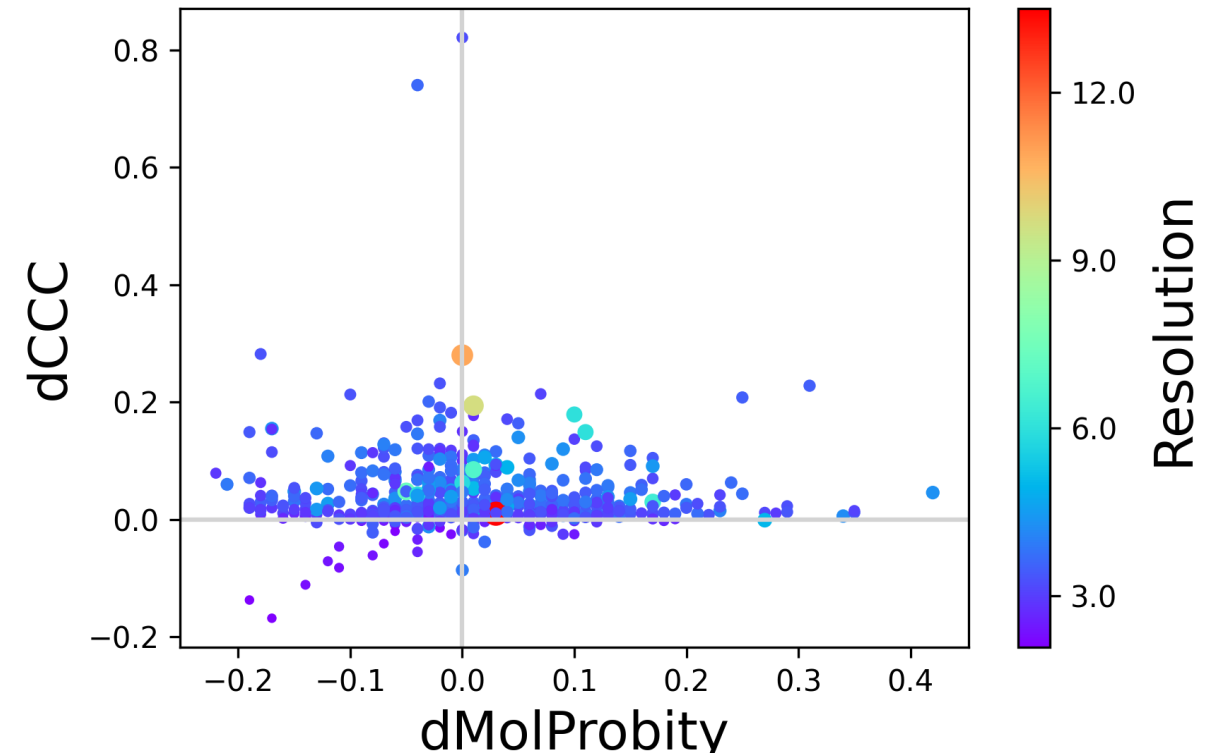
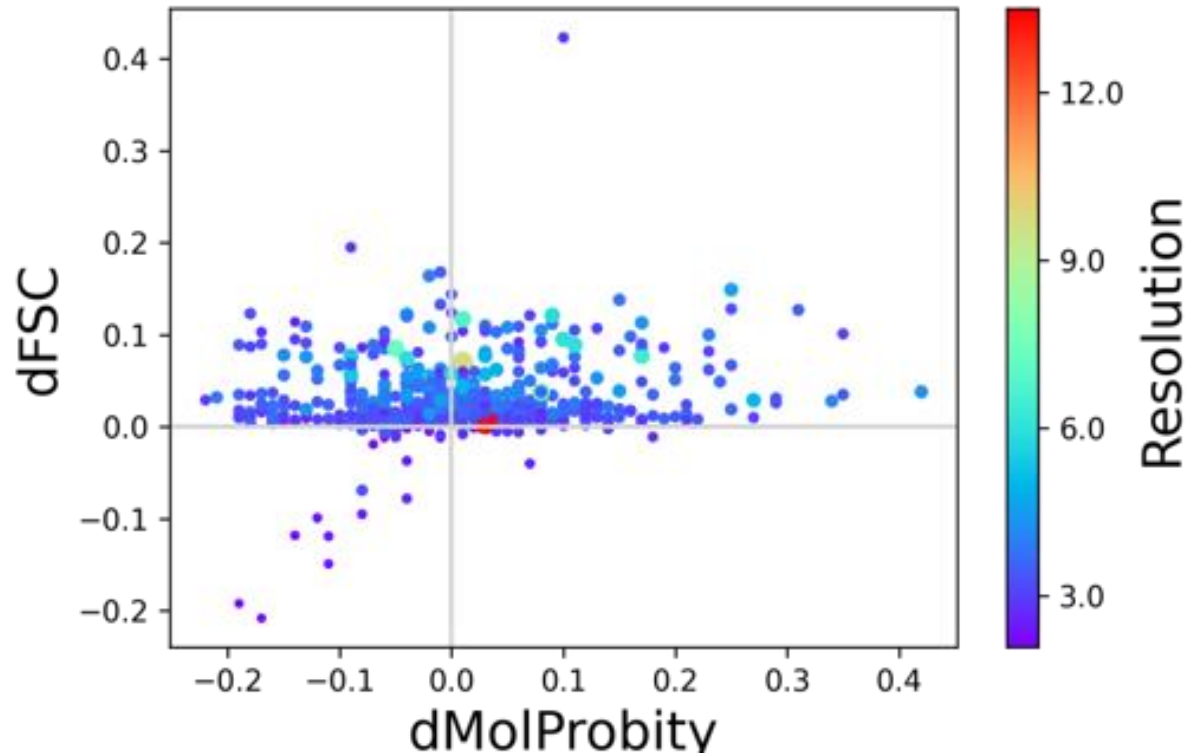
0 Ramachandran outlier, 1 poor rotamer, clashscore: 6.3



3.4Å

Model geometry vs fit-to-data

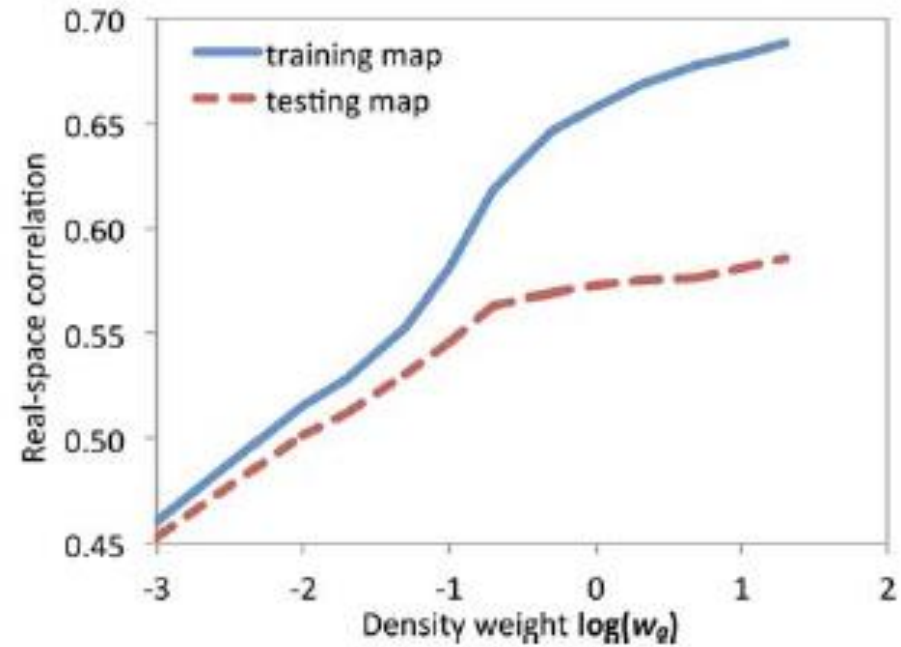
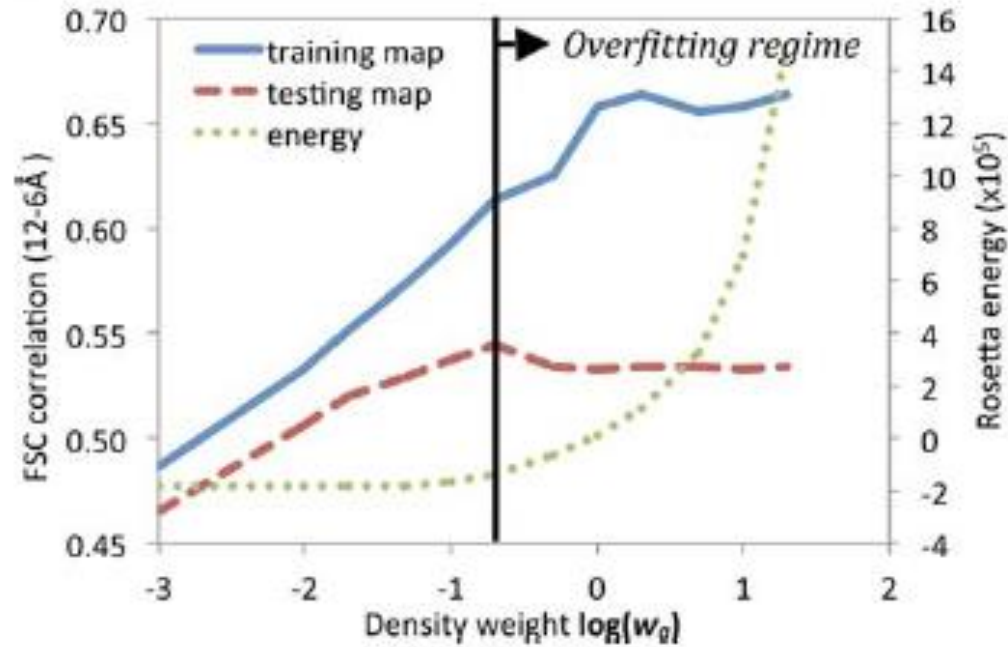
Can the map fit / representation be improved ?



FSCavg of 94% of structures in the dataset improved with further refinement

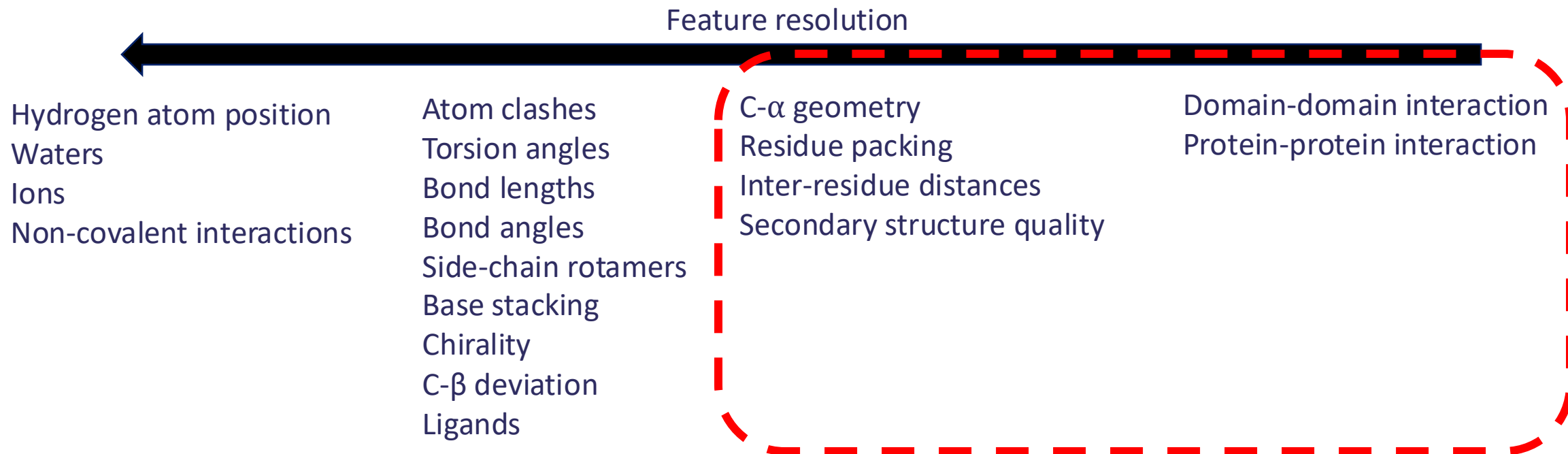
~45% of the dataset had improved MolProbity scores as well

Test against equivalent but independent data



Mm-cpn in the ATP/AlFx induced closed state
4.3 Å resolution

Model quality



Many of the deposited models are optimized more towards model geometry than fit to data.

Evaluation of low-resolution features is not common!

Fit to data

We need scores normalised for resolutions or reference distributions for different resolution bins

Better validation approaches especially to assess overfitting

EMDB statistics (geometry)

Molprobrity score

Resolution	Q1	Q2	Q3
1.0-2.5	1.2	1.39	1.67
2.5-3.5	1.27	1.53	1.8025
3.5-4.5	1.4025	1.695	1.96
4.5-6.5	1.32	1.77	2.1175
> 6.5	1.08	1.68	2.05

Clashscore

Resolution	Q1	Q2	Q3
1.0-2.5	3.1575	4.485	6.79
2.5-3.5	3.6175	5.765	8.17
3.5-4.5	4.23	7.165	11.075
4.5-6.5	3.33	7.82	13.803
> 6.5	2.26	6.57	12.65

Cis-proline (%)

Resolution	Q1	Q2	Q3
1.0-2.5	0	1.02	5.41
2.5-3.5	0	0	2.8
3.5-4.5	0	0	2.4675
4.5-6.5	0	0	2.94
> 6.5	0	1.22	2.46

Ramachandran Z-score

Resolution	Q1	Q2	Q3
1.0-2.5	-0.955	0.35	1.44
2.5-3.5	-1.32	-0.12	0.91
3.5-4.5	-1.705	-0.445	0.71
4.5-6.5	-2.12	-0.615	0.11
> 6.5	-1.73	-0.68	0.31

Ramachandran outlier (%)

Resolution	Q1	Q2	Q3
1.0-2.5	0	0	0.04
2.5-3.5	0	0	0
3.5-4.5	0	0	0.1
4.5-6.5	0	0.045	0.405
> 6.5	0	0.09	0.86

EMDB statistics (fit to map)

FSCavg (upto FSC0.5)

Resolution	Q1	Q2	Q3
1.0-2.5	0.746	0.773	0.8
2.5-3.5	0.738	0.7725	0.797
3.5-4.5	0.7	0.747	0.778
4.5-6.5	0.706	0.738	0.767
> 6.5	0.7055	0.754	0.8015

MI (mask)

Resolution	Q1	Q2	Q3
1.0-2.5	0.004	0.012	0.029
2.5-3.5	0.015	0.025	0.041
3.5-4.5	0.019	0.032	0.055
4.5-6.5	0.0245	0.05	0.07525
> 6.5	0.02675	0.0635	0.10825

CCC (mask)

Resolution	Q1	Q2	Q3
1.0-2.5	0.57225	0.615	0.6585
2.5-3.5	0.565	0.616	0.658
3.5-4.5	0.55925	0.615	0.657
4.5-6.5	0.55	0.626	0.693
> 6.5	0.5365	0.635	0.7105

SMOC Z-score < -2 (%)

Resolution	Q1	Q2	Q3
1.0-2.5	2.49675	2.845	3.27925
2.5-3.5	2.353	2.8	3.216
3.5-4.5	2.64725	3.1345	3.55925
4.5-6.5	3.116	3.552	3.995
> 6.5	2.6775	3.2875	3.60225

SMOC average

Resolution	Q1	Q2	Q3
1.0-2.5	0.6	0.6555	0.69325
2.5-3.5	0.6855	0.737	0.779
3.5-4.5	0.775	0.817	0.847
4.5-6.5	0.85975	0.8905	0.906
> 6.5	0.89	0.914	0.935

Model validation

Multiple complementary validation metrics are often useful than a single score

Faster validation scores

Better link between model building, refinement and validation

Doppio model validation

Metrics

Molprobrity (Geometry statistics)	<input checked="" type="radio"/> Yes	<input type="radio"/> No
Servalcat FSC (Model-map FSC)	<input type="radio"/> Yes	<input checked="" type="radio"/> No
TEMPy global (Real space map fit)	<input checked="" type="radio"/> Yes	<input type="radio"/> No
TEMPy SMOC (Per-residue map fit)	<input checked="" type="radio"/> Yes	<input type="radio"/> No
FDR backbone (Per-residue backbone trace)	<input type="radio"/> Yes	<input checked="" type="radio"/> No
CheckMySequence (Sequence agreement)	<input type="radio"/> Yes	<input checked="" type="radio"/> No

PROJECT JOBS NODES NEW JOB

Filter jobs by name or description ✕

Expand all

Atomic Model Validation

- CA** checkMySequence
checkmysequence.atomic_model_validation - Check modeled sequence against expected sequence and EM map
- MA** Molprobrity
molprobrity.atomic_model_validation.geometry - Atomic model geometry assessment
- PA** Atomic model validation
pipeliner.atomic_model_validation - Validate atomic models with multiple methods
- PA** Privateer
privateer.atomic_model_validation.glycans - Validate Glycan conformation
- TA** TEMPy Scores
tempy.atomic_model_validation.ma

Atomic model validation

RUN **JOB INFO** **RESET OPTIONS**

Job alias:

Main

Input model *

Quick evaluation? Yes No

Metrics

Molprobrity (Geometry statistics)	<input checked="" type="radio"/> Yes	<input type="radio"/> No	<input type="button" value="i"/>
Servalcat FSC (Model-map FSC)	<input type="radio"/> Yes	<input checked="" type="radio"/> No	<input type="button" value="i"/>
TEMPy global (Real space map fit)	<input checked="" type="radio"/> Yes	<input type="radio"/> No	<input type="button" value="i"/>
TEMPy SMOC (Per-residue map fit)	<input checked="" type="radio"/> Yes	<input type="radio"/> No	<input type="button" value="i"/>
FDR backbone (Per-residue backbone trace)	<input type="radio"/> Yes	<input checked="" type="radio"/> No	<input type="button" value="i"/>
CheckMySequence (Sequence agreement)	<input type="radio"/> Yes	<input checked="" type="radio"/> No	<input type="button" value="i"/>

Metric specific inputs

Willams et al. 2018, Chen et al. 2015,

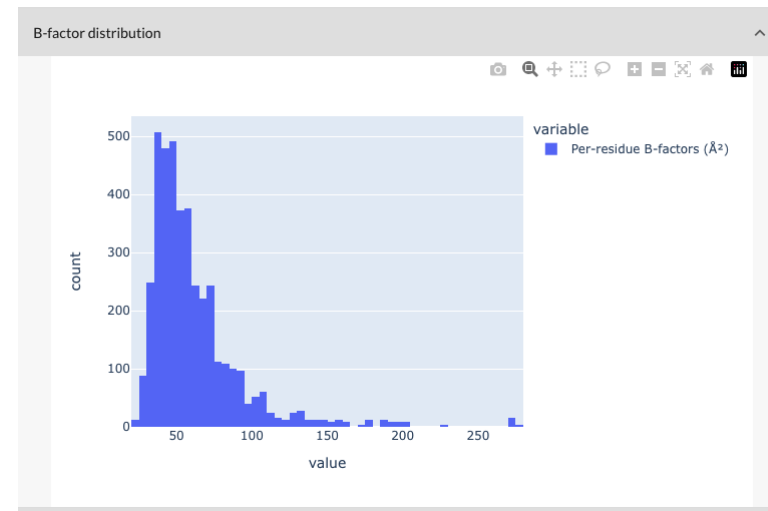
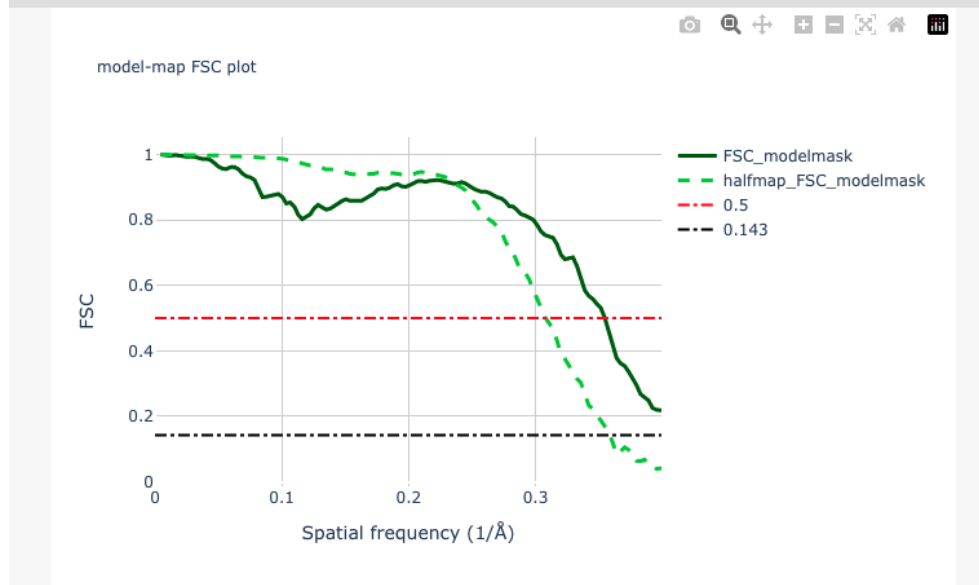
Yamashita et al. 2021,

Joseph et al. 2016; Grzegorz C Acta D 2022

Doppio model validation

Global map-fit scores	
Metric	Score
FSCavg_modelmask_FSC0.5	0.784
FSCavg_modelmask	0.639
FSCavg_h1_modelmask	0.568
FSCavg_h2_modelmask	0.568
CCC overlap mask	0.657
Model overlap fraction	0.733
CCC contoured map	0.573
MI overlap mask	0.049

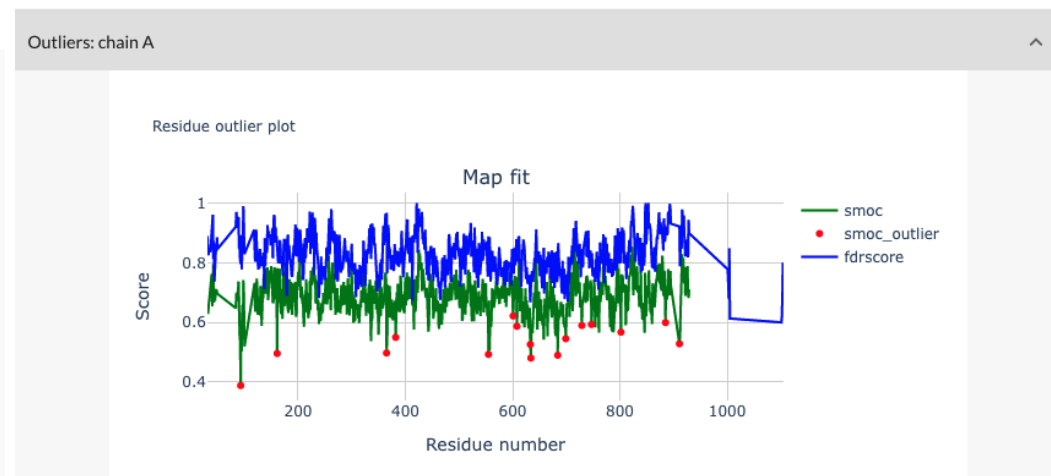
Model-map FSC	
---------------	--



Global geometry scores		
Metric	Score	Expected/Percentile
Ramachandran_outliers	0.10%	< 0.05%
Ramachandran_favored	97.26%	> 98%
RamachandranZ_whole	-0.72	bad Rama-Z > 3, suspicious 3 > Rama-Z ...
RamachandranZ_helix	-1.49	bad Rama-Z > 3, suspicious 3 > Rama-Z ...
RamachandranZ_sheet	0.12	bad Rama-Z > 3, suspicious 3 > Rama-Z ...
RamachandranZ_loop	-0.48	bad Rama-Z > 3, suspicious 3 > Rama-Z ...
Rotamer_outliers	0.34%	< 0.3%
CBeta_deviations	0	0
Clashscore	3.07	76.7 (percentile)
Molprobrity_score	1.24	83.0 (percentile)
Cis_proline	8.06	0%
Cis_general	0.31	0%
Rms_bonds	0.0085	< 0.02
Rms_angles	1.42	< 2.0

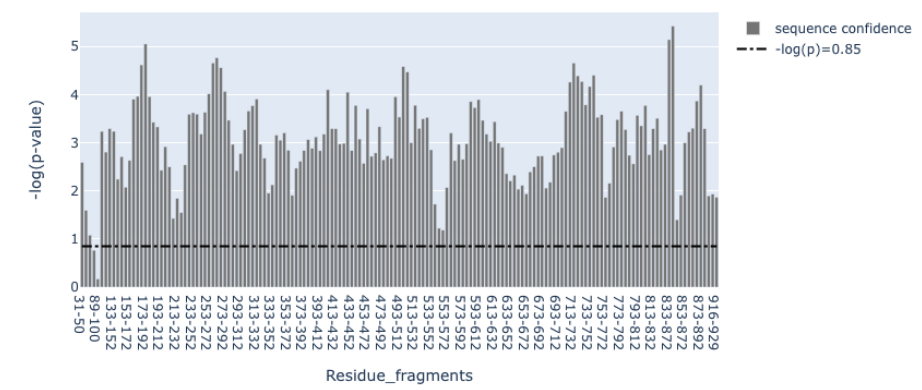
Doppio model validation

Cluster	Chain_Residue	Outlier Type(s)
1	A_462	ramachandran_outlier
1	A_423	omega_outlier bad_clash (atom:CG - D_280_H...
1	A_424	bad_clash (atom:N - D_281_OD2) bad_clash (a...
1	A_463	bad_clash (atom:HB3 - A_428_HG21)
1	D_280	bad_clash (atom:HD11 - A_423_CG) bad_clas...
1	D_281	bad_clash (atom:OD2 - A_424_N) bad_clash (a...
1	D_283	bad_clash (atom:NH1 - A_104_HG11) bad_cla...
2	A_280	bad_clash (atom:HD11 - D_423_CG) bad_clas...
2	A_281	bad_clash (atom:OD2 - D_424_N) bad_clash (a...
2	A_283	bad_clash (atom:NH1 - D_104_HG11)
2	D_462	ramachandran_outlier
2	D_423	omega_outlier bad_clash (atom:CG - A_280_H...
2	D_424	bad_clash (atom:N - A_281_OD2) bad_clash (a...
2	D_463	bad_clash (atom:HB3 - D_428_HG21)
3	B_462	ramachandran_outlier
3	B_423	omega_outlier bad_clash (atom:CG - C_280_H...
3	B_424	bad_clash (atom:N - C_281_OD2) bad_clash (a...
3	B_463	bad_clash (atom:HB3 - B_428_HG21)
3	C_280	bad_clash (atom:HD11 - B_423_CG) bad_clas...
3	C_281	bad_clash (atom:OD2 - B_424_N) bad_clash (a...
3	C_283	bad_clash (atom:NH1 - B_104_HG11)

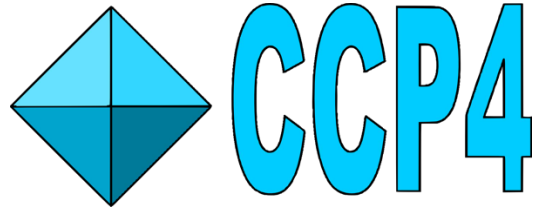


CheckMySequence sequence p-values, protein chain A

The bars are proportional to the $-\log(p\text{-value})$ and the higher they are, the more reliable the sequence of a fragment. Grey bars: Expected 'r' fragment sequence and model sequence match.



Acknowledgements



Scientific Computing





Science and
Technology
Facilities Council

Scientific Computing

Questions?





Science and
Technology
Facilities Council

Scientific Computing

Thank you

scd.stfc.ac.uk

 [@SciComp_STFC](https://twitter.com/SciComp_STFC)