

Generating/optimising predicted models for MR

Adam Simpkin



Talk overview

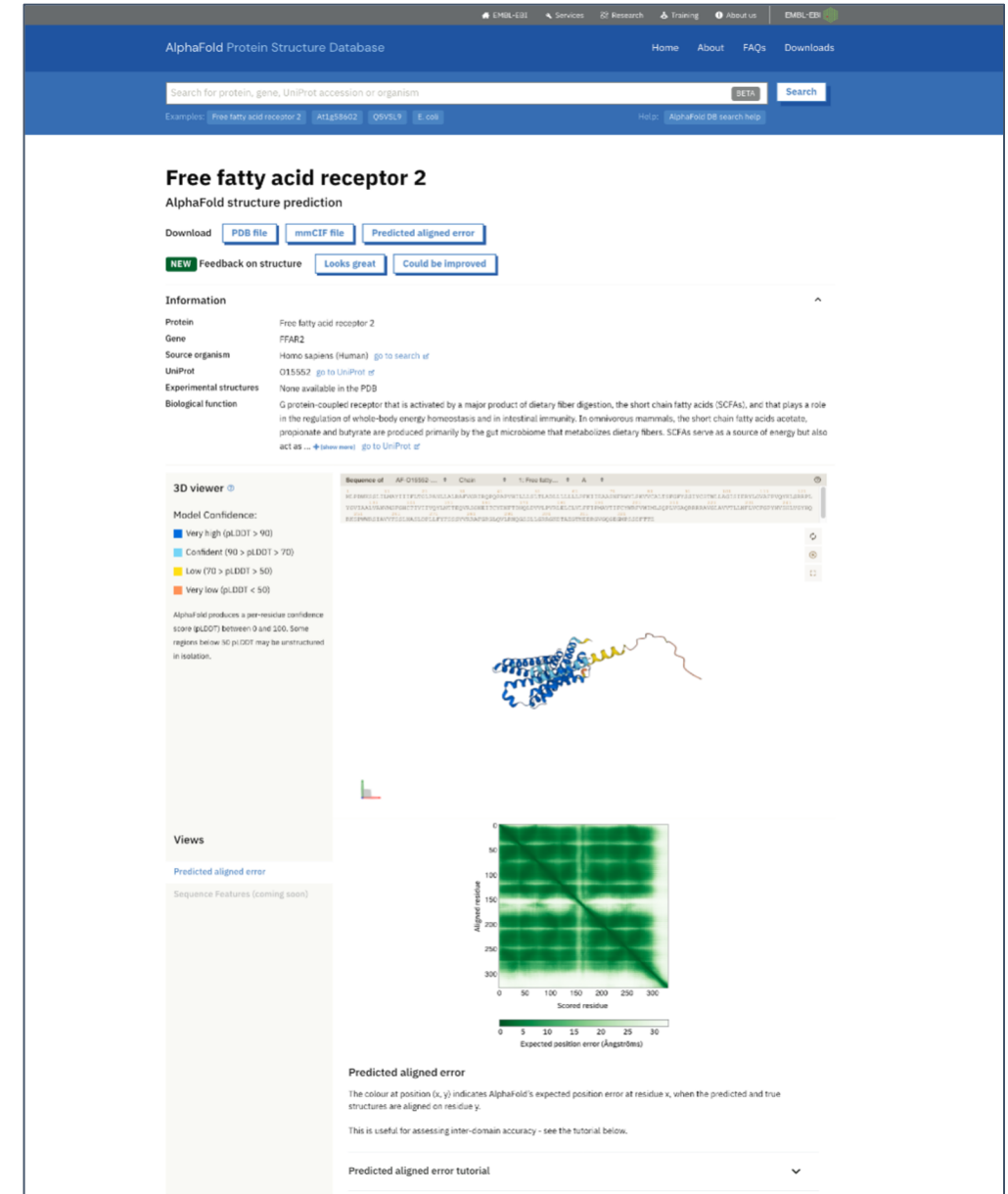
- The EBI AlphaFold database (AFDB)
 - How to query the AFDB for your target structure
 - How to assess the quality of hits in the AFDB
 - CCP4 tools that take advantage of the database
- Making your own models using Colab pages/online servers
 - ColabFold and AlphaFold's Colab page
 - Making complexes and other things to think about
 - Using RoseTTAFold's online server
 - Exploring conformational diversity with ColabFold
- ~~• Making your own models locally~~
 - ~~• What you need to run AlphaFold2 locally~~
- Optimising your models for MR
 - Problems with models
 - Slice'N'Dice

Obtaining models from AFDB

The EBI AlphaFold database

The AFDB should be your first port of call when looking for AlphaFold2 models for your target protein

- Initial release contained structure predictions for the human proteome and 20 other key organisms (~350,000 structures)
- December 2021, expanded to include manually curated UniProt (>400,000)
- Currently release contains >200 million models
 - Due to a compilation error, ~4% of these entries had low pLDDT entries
 - V4 of the database (releasing soon) will fix this

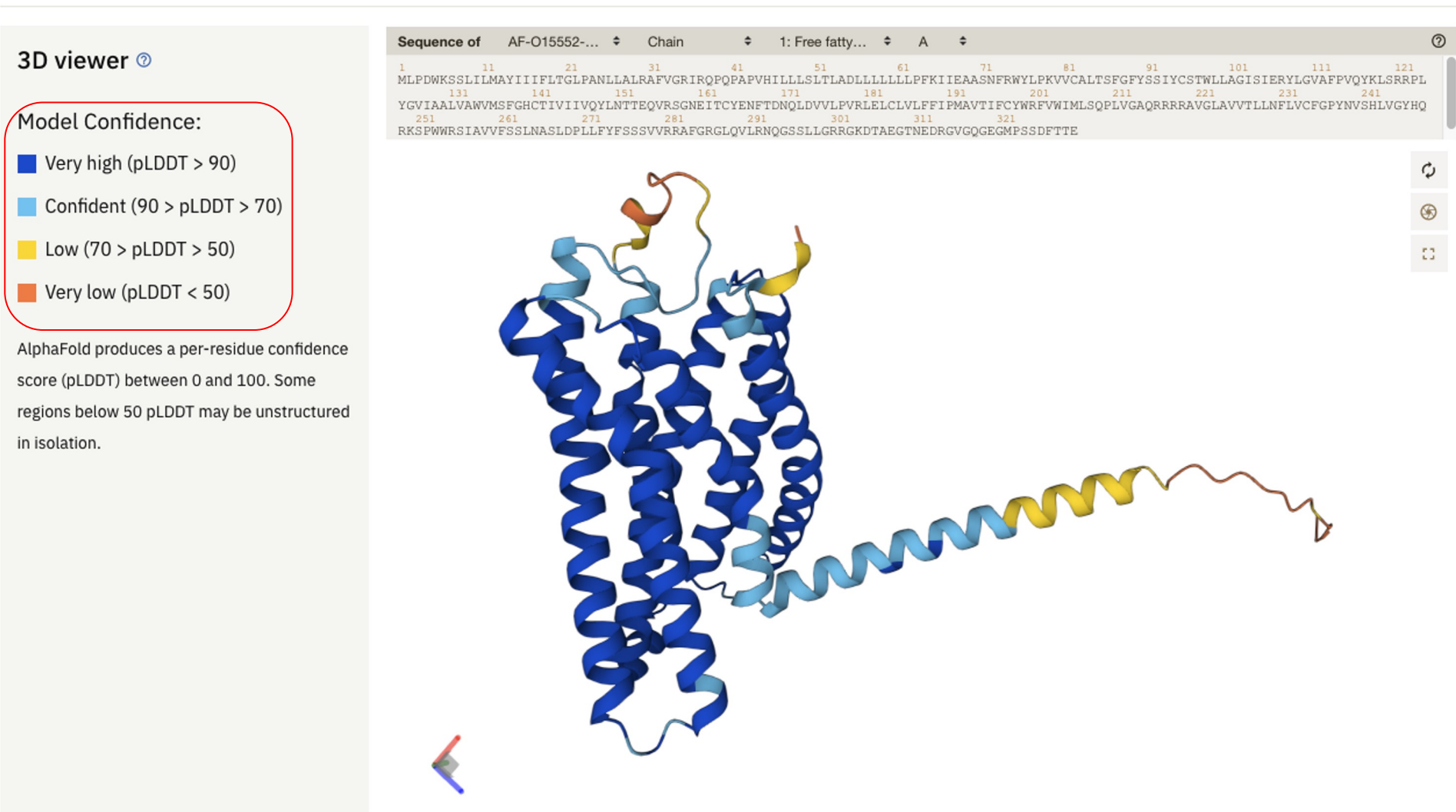


The screenshot displays the AlphaFold Protein Structure Database interface for the protein 'Free fatty acid receptor 2' (FFAR2). The page includes a search bar at the top, navigation links (Home, About, FAQs, Downloads), and a 'Search' button. Below the search bar, the protein name 'Free fatty acid receptor 2' is prominently displayed, followed by 'AlphaFold structure prediction'. A 'Download' section offers links for 'PDB file', 'mmCIF file', and 'Predicted aligned error'. A 'NEW' badge indicates 'Feedback on structure' with 'Looks great' and 'Could be improved' buttons. The 'Information' section provides details about the protein, including its name (FFAR2), source organism (Homo sapiens (Human)), UniProt accession (O15552), and biological function (G protein-coupled receptor). A '3D viewer' section shows a ribbon diagram of the protein structure, a 'Model Confidence' legend, and a 'Predicted aligned error' heatmap. The heatmap shows the expected position error at residue x when the predicted and true structures are aligned on residue y. A 'Predicted aligned error tutorial' link is also present.

Searching the database

- Models can be obtained directly from <https://alphafold.ebi.ac.uk/> using a UniProt ID or keywords relating to the protein
- For sequence search can use:
 - HMMER www.ebi.ac.uk/Tools/hmmer currently only searches older AFDB (2022-01-15)
 - MrParse in the CCP4 suite, searches UniprotKB and then finds which entries are in the AFDB (also does some model pre-processing to make them more effective in molecular replacement)
- For structural search can use:
 - DALI (<http://ekhidna2.biocenter.helsinki.fi/dali/>) is a well-established method, but only looks at one species at a time. Gives a Z-score. >2 non-random, ~>4 fold similarity
 - Foldseek (<https://search.foldseek.com/search>) uses Tmalign method to produce TM-score. Range is 0-1. Scans AFDB
 - RUPEE (<https://ayoubresearch.com/>) is a fast shape-based method. Maybe not the most sensitive. Scans all AFDB.

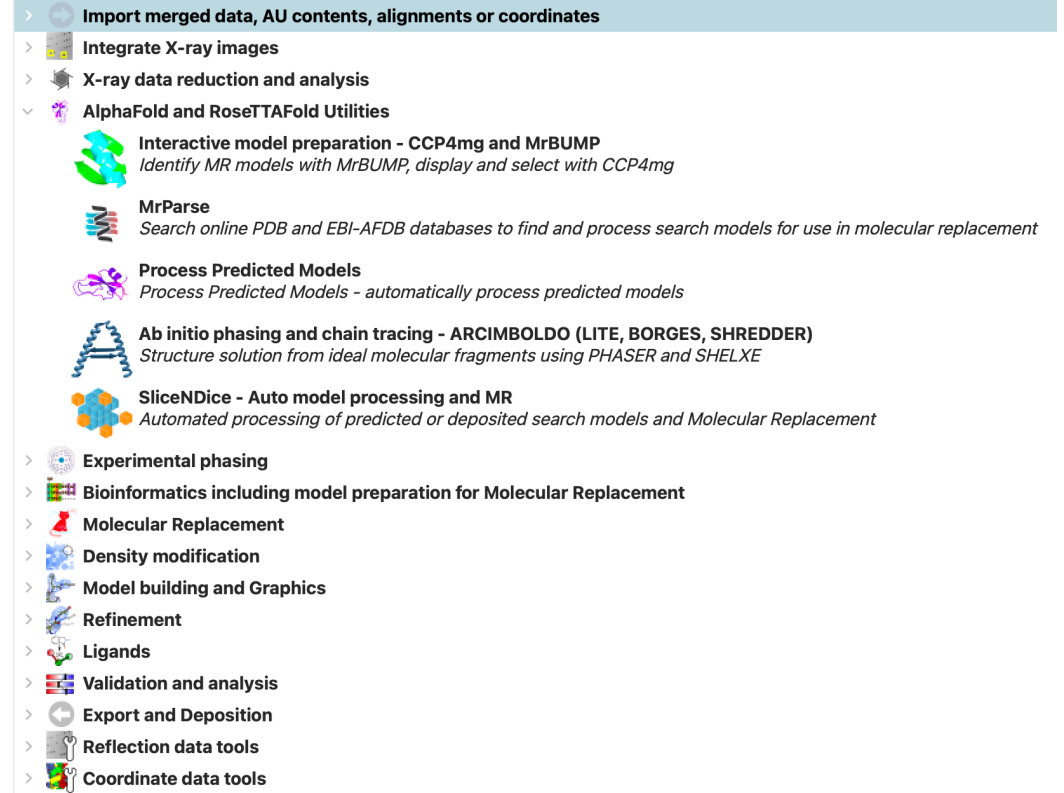
Assessing quality of models in the AFDB



CCP4 programs that use the AFDB

MrParse

- Available through CCP4i2 and CCP4 Cloud
- Located under the AlphaFold and RoseTTAFold utilities tab
- Sequence must be provided, reflections are optional (used to calculate eLLG and provide a summary of the reflection data)



A screenshot of the MrParse job configuration form. The 'Job title' field is set to 'mrparse'. The 'Use data from job' dropdown is set to 'No'. The 'Select input data' section has a 'Sequence' field with the text '..must be selected' and a 'Reflections' field with the text '..is not used'. Below this, a note states: 'Observed intensities (or amplitudes) are optional but recommended for assessing crystal pathology and calculating eLLG scores'. The 'Do sequence classification?' dropdown is set to 'Yes'.

Avoiding the miscompiled models



MrParse Analysis

Version: 0.3.5

MrParse: a program to find and analyse search models for crystallographic Molecular Replacement. The program is being developed by [Dan Rigden's group](#) at the University of Liverpool.

MrParse is currently under development and we are keen to make it as useful to the community as possible. If you have any suggestions for it's development, or ideas on how we could improve it, please [get in touch](#).

Experimental structures from the PDB

Name	PDB	Resolution	Region	Range	Length	eLLG	Mol. Wt.	eRMSD	Seq. Ident.
4mo1_B	4mo1	2.10	1	2-144	141		15231		0.99

Structure predictions from the EBI AlphaFold database

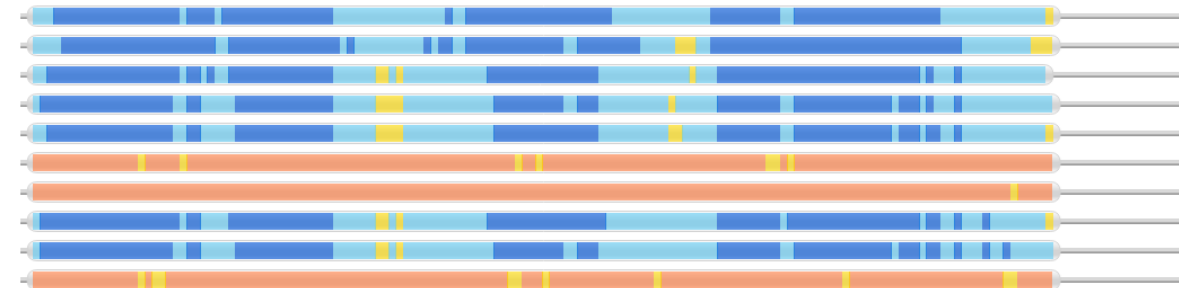
Name	model	Date Made	Region	Range	Length	Avg. pLDDT	H-score	Seq. Ident.
AF-A0A6B5D4Y7-F1	AF-A0A6B5D4Y7-F1	01-JUN-22	1	1-147	145	90.31	84	0.99
AF-A0A7R8VQK2-F1	AF-A0A7R8VQK2-F1	01-JUN-22	1	1-147	145	90.11	85	0.99
AF-A0A644RJ15-F1	AF-A0A644RJ15-F1	01-JUN-22	1	1-146	144	88.03	81	0.99
AF-A0A5U9QSB5-F1	AF-A0A5U9QSB5-F1	01-JUN-22	1	1-147	145	87.68	81	0.99
AF-A0A3U6ZW02-F1	AF-A0A3U6ZW02-F1	01-JUN-22	1	1-147	145	87.63	81	0.99
AF-A0A4Q8JJP8-F1	AF-A0A4Q8JJP8-F1	01-JUN-22	1	1-147	145	40.54	41	0.99
AF-A0A7U9FWD5-F1	AF-A0A7U9FWD5-F1	01-JUN-22	1	1-147	145	34.25	35	0.99
AF-A0A860YL56-F1	AF-A0A860YL56-F1	01-JUN-22	1	1-147	145	88.65	82	0.98
AF-A0A0K3YHM2-F1	AF-A0A0K3YHM2-F1	01-JUN-22	1	1-147	145	87.93	82	0.98
AF-A0A0A0FDZ1-F1	AF-A0A0A0FDZ1-F1	01-JUN-22	1	1-147	145	36.53	37	0.98

Visualisation of Regions



** Sequence Based Prediction step was skipped: append --do_classify argument to run **

Visualisation of Regions



Problematic entries are very obvious

Using AlphaFold database models in MrBUMP

- AFDB models can now also be used in MrBUMP via CCP4mg
- Can select threshold score so that models are retrieved without low quality regions
- Ensembles can be made from PDB and AFDB entries combined
- Slider allows interactive truncation
- Full instructions available here:
https://ccp4i2.gitlab.io/rstdocs/tutorials/alphafold/alphafold_tutorial.html

Job 8: Interactive model preparation - CCP4mg and MrBUMP The job is Pending

Input Results Comments

Input data

Job title: CCP4mg MrBUMP

Use data from job: 3 Define AU contents as input below..

Sequences from AU content:

AU contents: 3 Asu content file from Define AU contents

If a suitable ASU is not available above, you can press the cross & then button to quickly create one.

Select one sequence

☒ gamma gamma

Model databases

☒ Search PDB for possible MR search models

Non-redundancy level for homologue search: 100

☒ Search EBI-AFDB for possible MR search models

EBI-AFDB pLDDT residue score cut-off: 50

Maximum no. of search models to create: 10

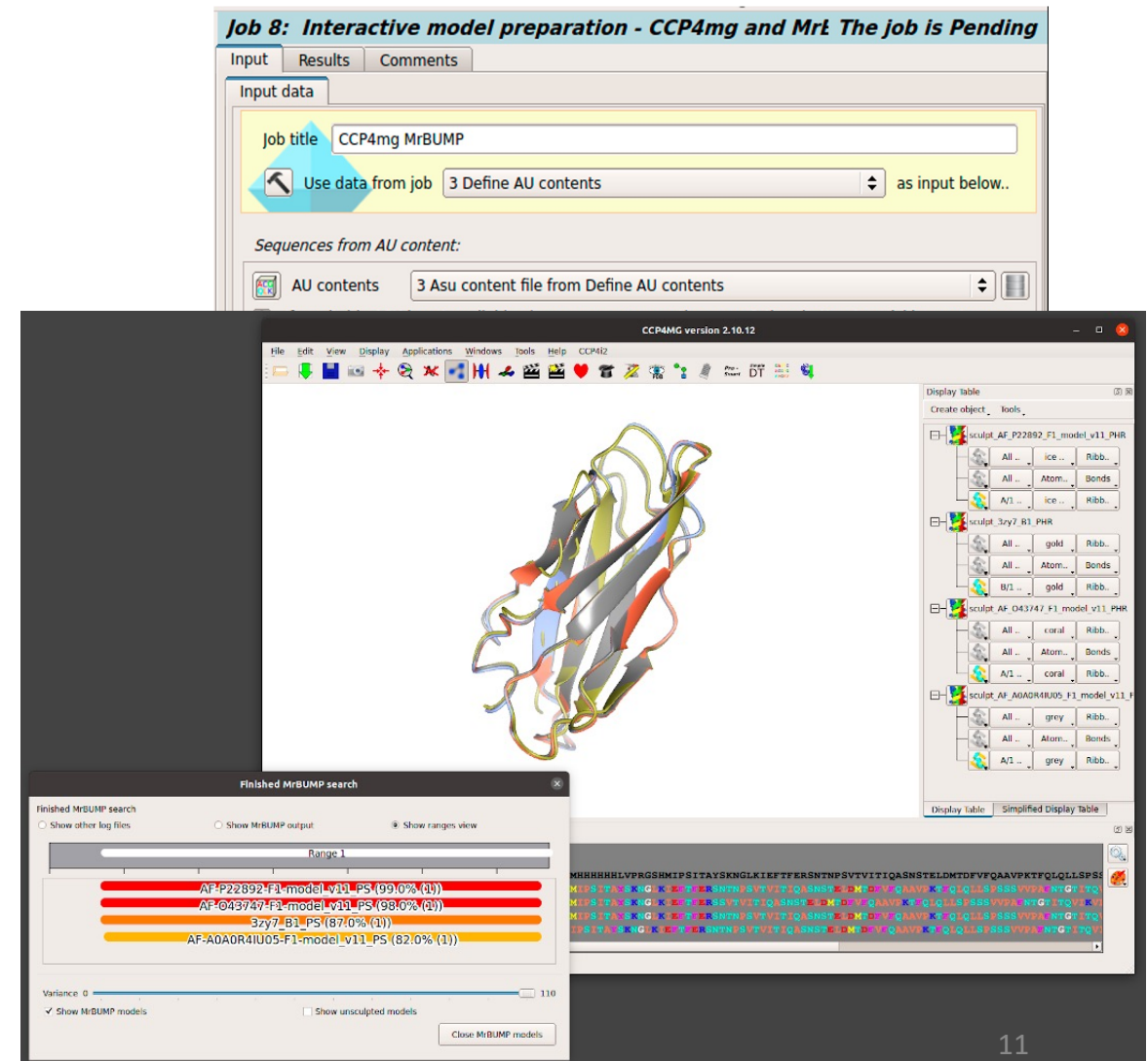
Optional Settings

☐ HHPred hhr file ..is not used

☐ Path to local PDB mirror ..is not used

Using AlphaFold database models in MrBUMP

- AFDB models can now also be used in MrBUMP via CCP4mg
- Can select threshold score so that models are retrieved without low quality regions
- Ensembles can be made from PDB and AFDB entries combined
- Slider allows interactive truncation
- Full instructions available here:
https://ccp4i2.gitlab.io/rstdocs/tutorials/alphafold/alphafold_tutorial.html



Pros and cons of AFDB

Pros

Immediately available

Likely to be a close match to your target

Large structures (up to ~2700 residues)

Nice interactive output

PAE json file available

Cons

Not made with today's data

Only monomers

Only a single model, i.e. no conformational variation

Unable to vary AlphaFold2 settings to optimise modelling

Temporarily, some models have errors

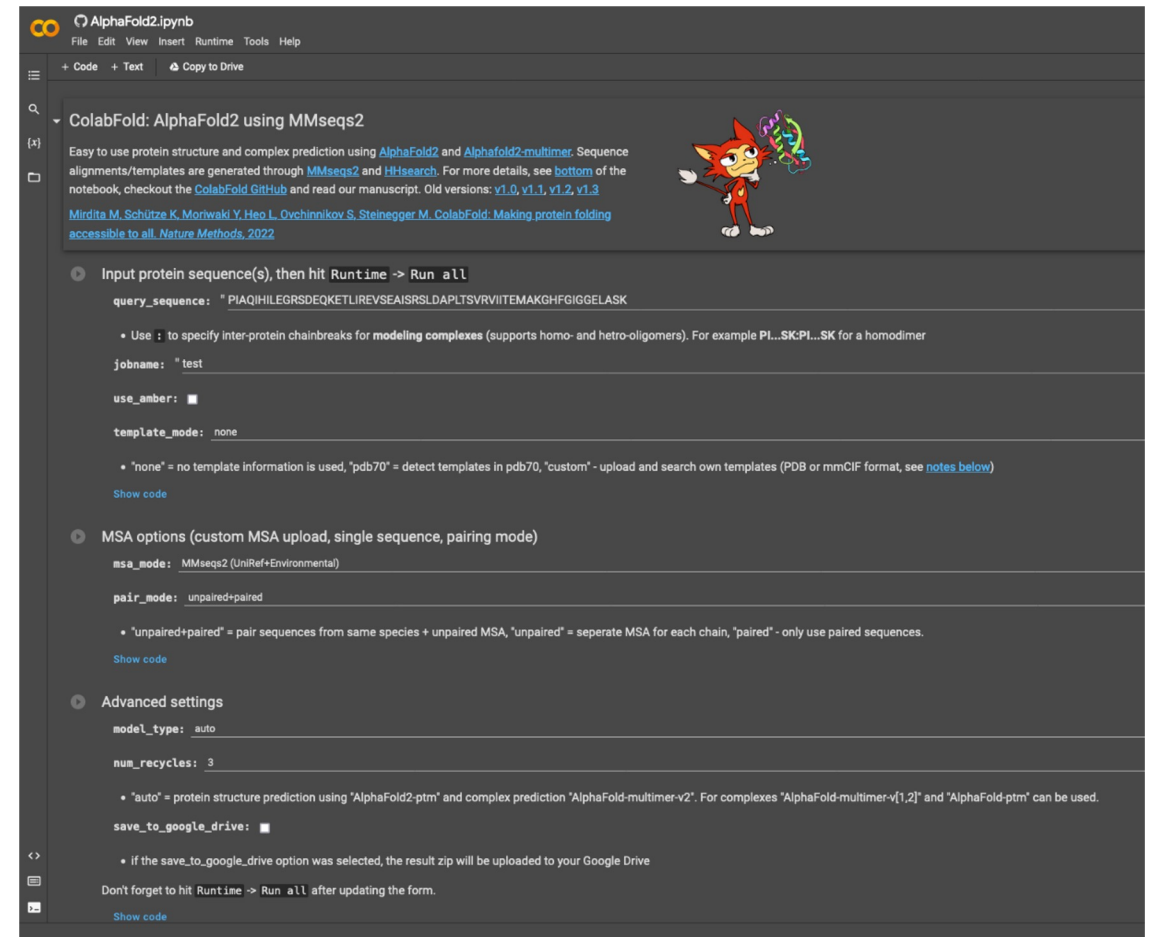
Making your own models using Colab pages

ColabFold and other Colab pages

A variety of Colab pages run AlphaFold2 on free Google resources

These differ in options for:

- MSA generation (MMSeqs2 vs JackHMMer)
- Multimer modelling
- Template usage



The screenshot shows the ColabFold AlphaFold2 using MMseqs2 notebook interface. The title bar indicates the notebook is named 'AlphaFold2.ipynb'. The main content area is titled 'ColabFold: AlphaFold2 using MMseqs2' and includes a brief description of the tool's purpose and a cartoon fox mascot. The interface is divided into several sections for configuring the protein structure prediction process:

- Input protein sequence(s), then hit Runtime -> Run all**: This section contains a text input field for the 'query_sequence' with the value 'PIAQIHILEGRSDEQKETLIREVSEAIRSLDAPLTSVRVITEMAKGHFGIGGELASK'. Below this are options for 'use_amber' (set to false) and 'template_mode' (set to 'none'). A note explains that 'none' means no template information is used, while 'pdb70' would detect templates in the PDB70 database.
- MSA options (custom MSA upload, single sequence, pairing mode)**: This section includes 'msa_mode' (set to 'MMseqs2 (UniRef+Environmental)') and 'pair_mode' (set to 'unpaired+paired'). A note explains that 'unpaired+paired' means pairing sequences from the same species with an unpaired MSA, while 'unpaired' would use separate MSAs for each chain.
- Advanced settings**: This section includes 'model_type' (set to 'auto'), 'num_recycles' (set to 3), and 'save_to_google_drive' (set to false). A note explains that 'auto' means protein structure prediction using 'AlphaFold2-ptm' and complex prediction using 'AlphaFold-multimer-v2'. The 'save_to_google_drive' option, if selected, would upload the result zip file to the user's Google Drive.

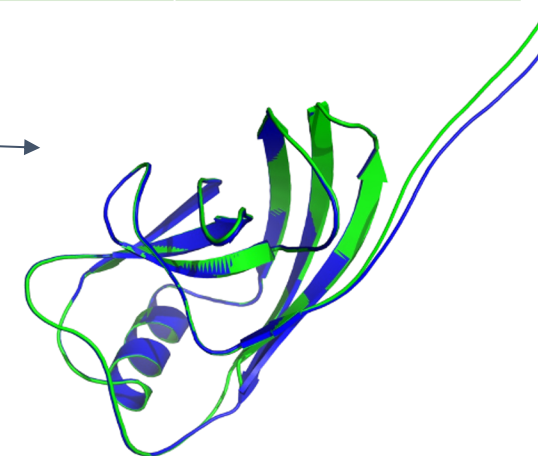
At the bottom, there is a reminder to hit 'Runtime -> Run all' after updating the form, along with a 'Show code' link.

Mirdita, M., Schütze, K., Moriwaki, Y. *et al.* ColabFold: making protein folding accessible to all. *Nat Methods* (2022). <https://doi.org/10.1038/s41592-022-01488-1>

Colab comparison

Method	Monomers	Complexes	MMSeqs2	JackHMMer	Templates
AlphaFold colab	Yes	Yes	No	Yes	No
ColabFold	Yes	Yes	Yes	No	Yes
ColabFold Advanced	Yes	Yes	Yes	Yes	Yes

- MMSeqs2 runs significantly quicker than JackHMMer
 - Sequence coverage in the MSA is lower, however in the majority of cases this doesn't affect the quality of the output model
- ColabFold can accept template models uploaded by users
- Complexes are created differently:
 - AlphaFold colab and ColabFold use AlphaFold-multimer code to produce complexes.
 - ColabFold Advanced uses a residue-index mode (Similar to how RoseTTAFold makes complexes).
 - The AlphaFold-multimer code has been shown to produce the best complexes overall, however in certain instances the residue-index mode will create better complexes [1].



Making models with RoseTTAFold online server

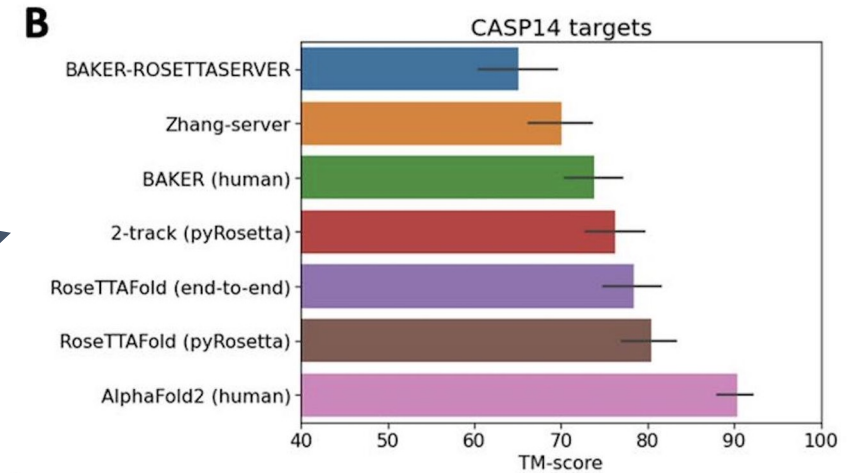
Shortly after CASP14, RoseTTAFold was released

RoseTTAFold is also a deep learning based method that can produce highly accurate models

Takes template models

RoseTTAFold is available from <https://robetta.bakerlab.org/>

(nb registration required)



Robetta Project ▾ Structure Prediction ▾

Submit a job for structure prediction
Please do not submit jobs under different user accounts. Such jobs will be removed.

Required

Target Name

Protein sequence

or upload FASTA No file selected.

Optional

RoseTTAFold ☒ CM ☐ AB ☐ Predict domains ☐

Upload MSA No file selected.

3 + 2 = Keep private ☐

POWERED BY

Exploring conformational diversity with ColabFold

The default pipeline often produces a set of similar models. If your protein has multiple conformations you may want to look for others

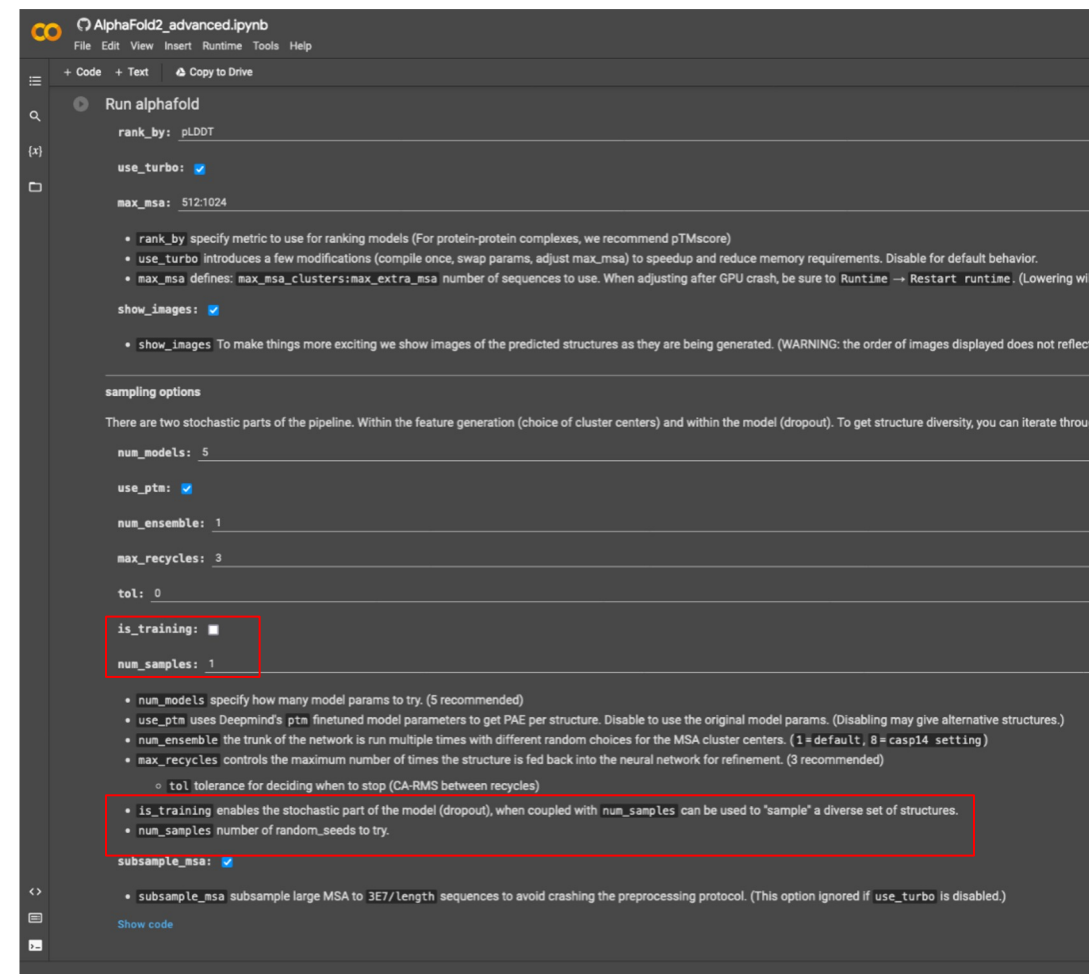
Ways to sample conformation more broadly

- Use the `num_samples` and/or `is_training` options on the advanced colab page
- Deliberately make the input MSA more shallow [1] (nb: this may reduce pLDDT scores)
- Feed AF2 templates in the 'right' conformation (and maybe ignore MSA features) [2]
- Edit the input MSA to mutate the Ala residue pairs that are driving the 'wrong' conformation [3] (have to provide own MSA)

[1] <https://www.biorxiv.org/content/10.1101/2021.11.22.469536v1.full.pdf>

[2] <https://www.biorxiv.org/content/10.1101/2021.11.26.470086v1.full.pdf>

[3] <https://www.biorxiv.org/content/10.1101/2021.11.29.470469v1.full.pdf>



Pros and cons of Colab pages

Pros

You can model your exact protein

Can make complexes

Can specify template proteins

Made with latest databases

Only requires an internet connection to run

Lots of runtime options e.g. conformation diversity

Nice looking output

Cons

Can only make models up to ~1400 residues
(Depends on GPU assigned)

Advantages of making models locally

Pros and cons of local runs

Pros

- You can model your exact protein
- Can typically make much bigger models
- Can make complexes
- Made with latest databases

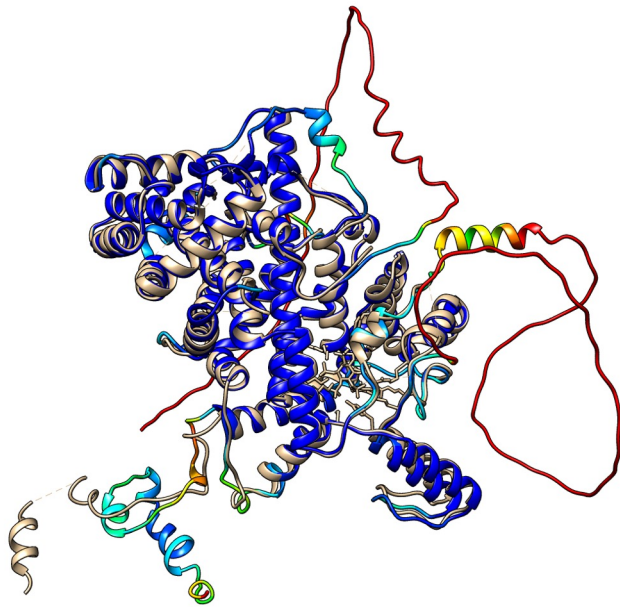
Cons

- Less intuitive to run
- Requires decent hardware
- Less informative output?

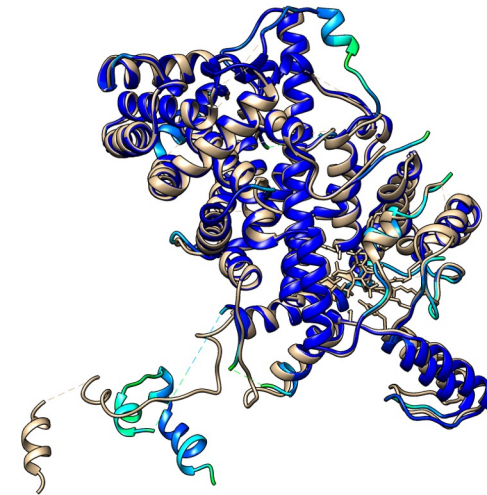
Optimising models for MR

Problems with predicted models

- The quality of a predicted model is not uniform
- Some regions will be poorly predicted
 - These regions don't provide any useful information and can interfere with crystal packing
- Model quality assessment (MQA) scores such as pLDDT can be used as a guide to remove the poorly predicted regions



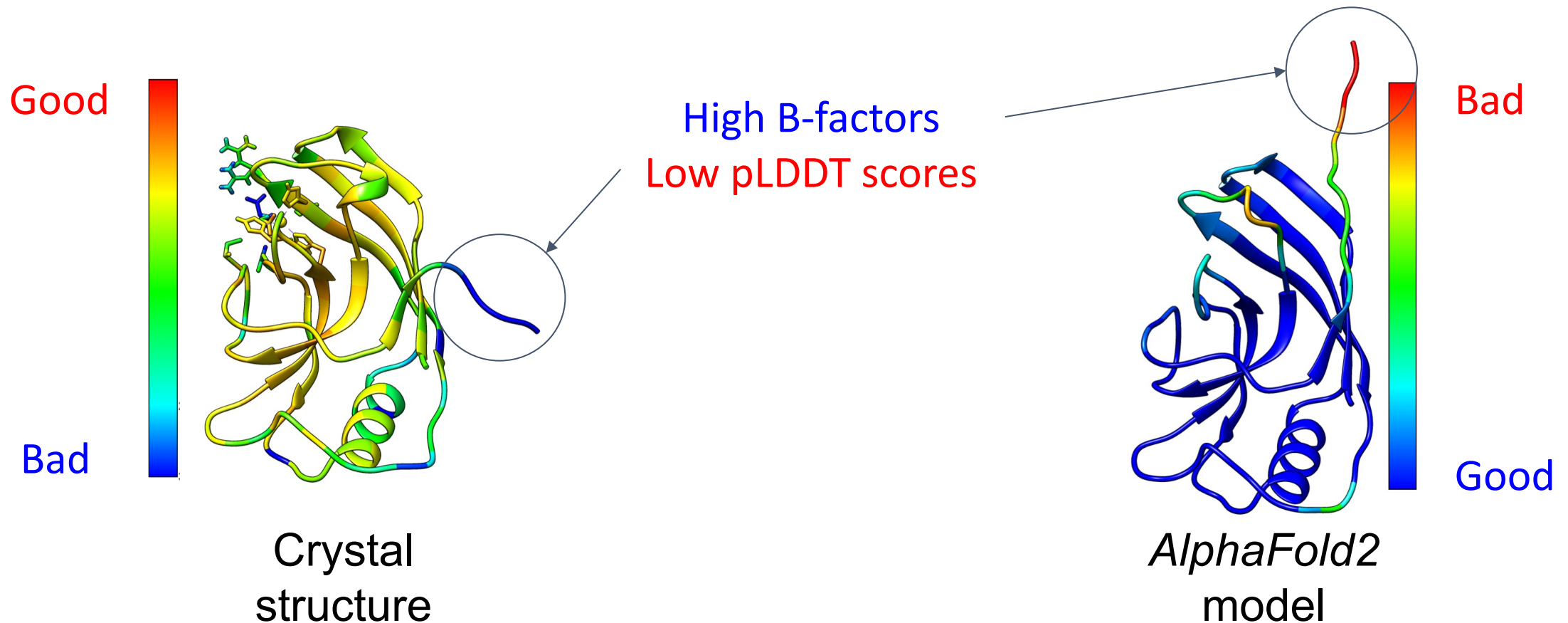
AFDB: Q5VSL9
PDB: 7K36_I



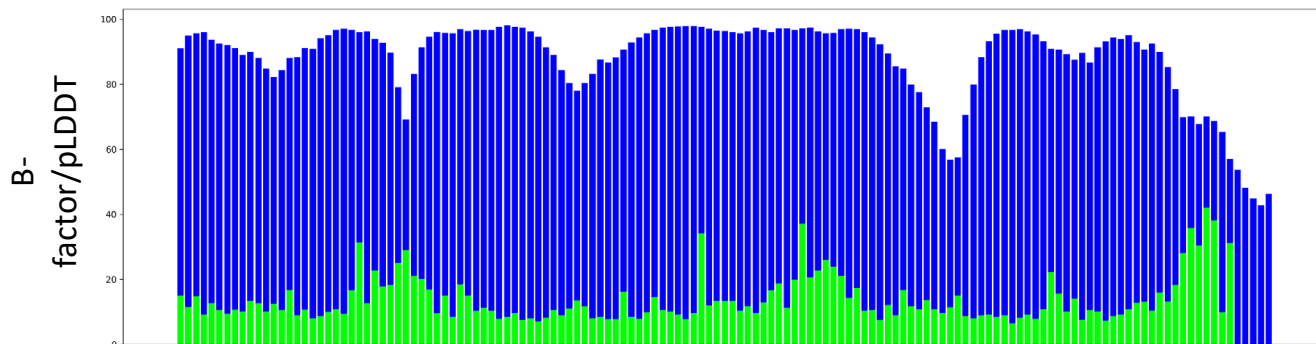
pLDDT threshold: 70

Problems with predicted models

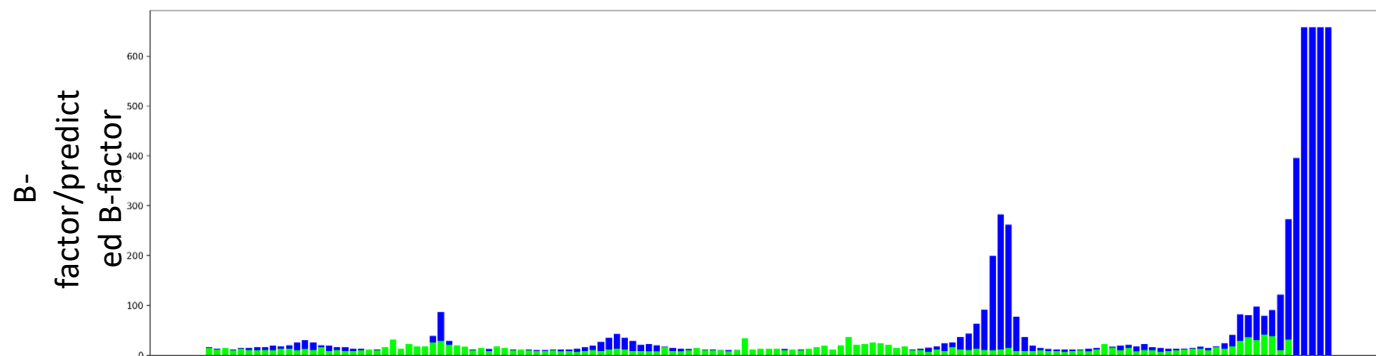
- MQA scores are often stored in the B-factor column of the PDB file
- B-factors are used weight residues in MR programs



Problems with predicted models



$$\frac{8 \times \pi^2}{3} \times \left(\frac{0.6}{pLDDT \times 0.01^3} \right)^2$$



Problems with predicted models

Many programs can now automatically remove low pLDDT regions and convert pLDDT scores to predicted B-factors, these include:

- CCP4 process predicted models
- Phenix.process_predicted_model
- MrParse (will correct for this when downloading models from the AFDB)
- MrBUMP (will correct for this when using models from the AFDB)
- Slice'N'Dice

Problems with predicted models

- Conformational differences between the predicted model and your target protein
- Individually, protein domains can be very accurately modelled, however if the orientation of these domains is incorrect relative to one another, the model may not work in molecular replacement



AlphaFold2 model
Crystal structure

SLICE'N'DICE

SLICE'N'DICE

Composed of two elements

SLICE

Slices model up into distinct structural units using clustering methods.

Clustering input data:

- Predicted Aligned Error (PAE)
- C-alpha coordinates

Clustering methods fall into two categories:

- Manual clustering methods - number of clusters to produce must be specified
- Automatic clustering methods - number of clusters is automatically determined by the method

DICE

Performs automated Molecular replacement or Map fitting using the 'sliced' up model

Molecular Replacement uses:

- Phaser
- Molrep

Map Fitting uses:

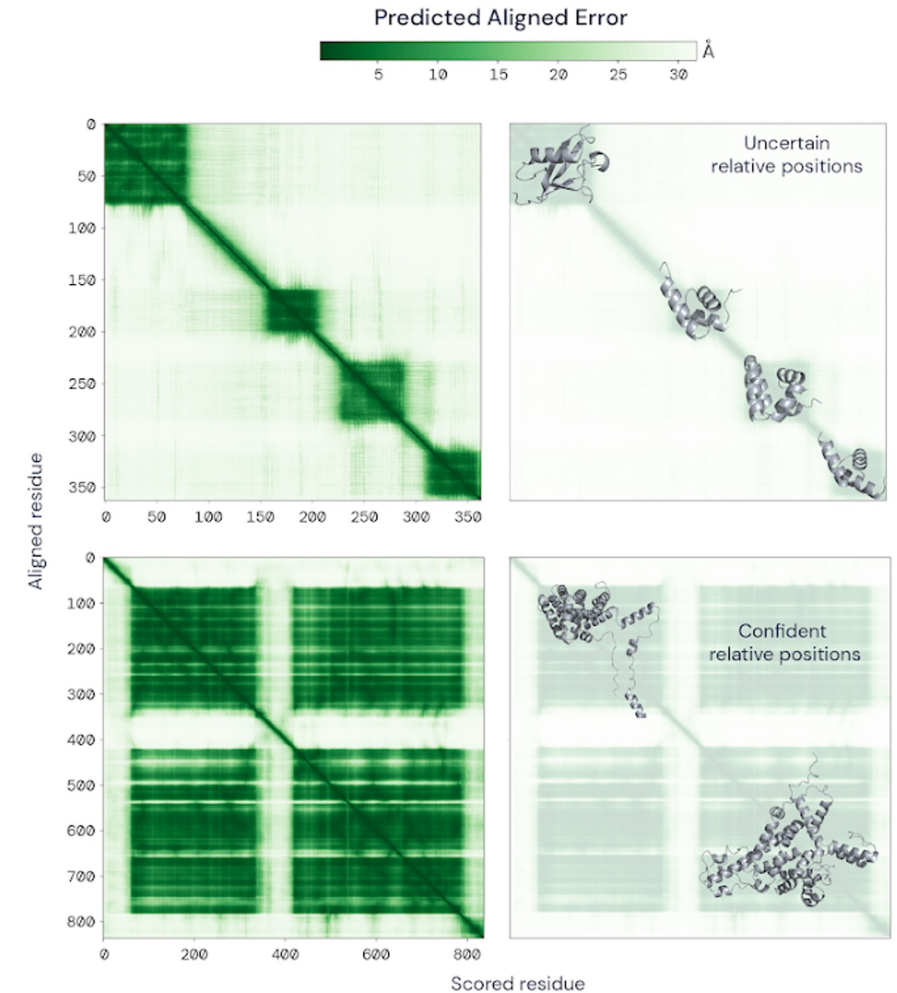
- Molrep
- Powerfit

Clustering Algorithms

Clustering on PAE

AlphaFold2 produces a metric called Predicted Aligned Error (PAE) that measure the confidence of residues relative to one another

Tristan Croll has created a tool in CCTBX that uses the information in the PAE file to identify individual structural units in an AlphaFold2 model

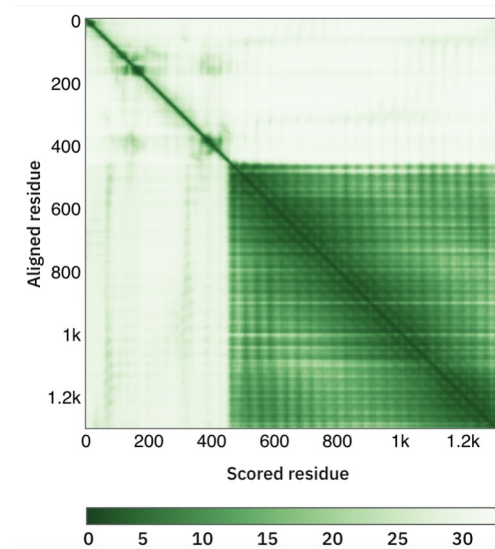
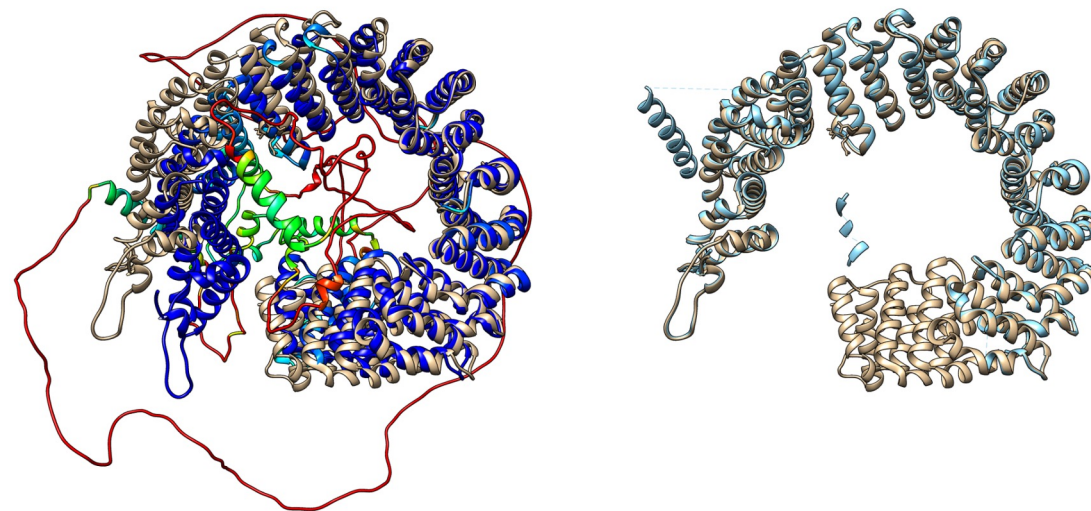


Clustering on C-alphas

PAE is not always the best way to cluster models

SciKitLearn clustering algorithms that use the XYZ coordinates of C-alpha atoms can:

- Cluster predicted models from other sources (e.g. RoseTTAFold)
- Identify structural units when PAE is uninformative



MX examples

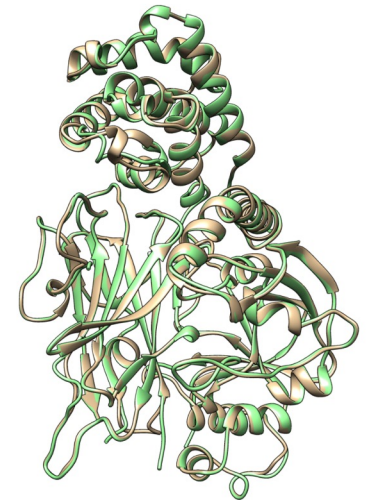
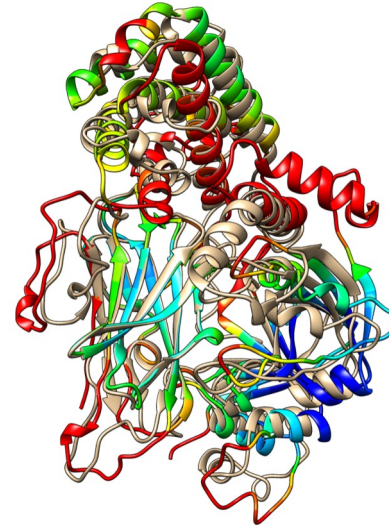
Example 1 - 7OA7

- Crystal Structure solved by SAD
- Closest hit in the PDB, 12% seq ID
- AlphaFold2 model has avg pLDDT of 85.61
- Model was unable to solve the structure by MR
- Split into two using the Birch algorithm in Slice'N'Dice (C-alpha based)
- Phaser was able to solve the structure extremely easily (LLG: 1339, TFZ: 35.6)



Example 2 - 7RB4

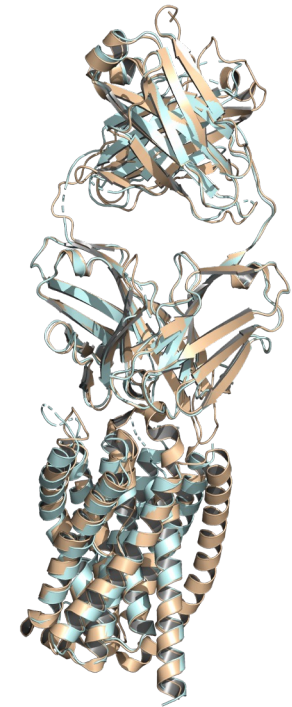
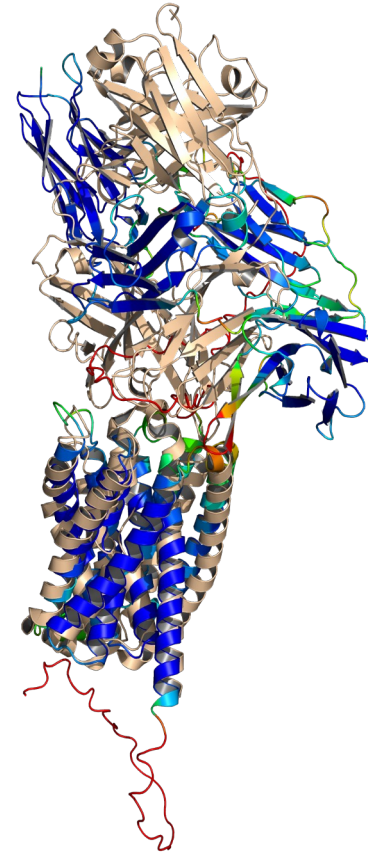
- Crystal Structure solved by SAD
- Closest hit in the PDB, 26% seq ID
- Only 26 sequences found in the MSA
- AlphaFold2 model has avg pLDDT of 61.02
- Model was unable to solve the structure by MR
- Split into three using the Birch algorithm in Slice'N'Dice (C-alpha based) with pLDDT threshold of 50
- Phaser was able to solve the structure (LLG: 114, TFZ: 13.2)
- Unable to solve with pLDDT threshold of 70



Cryo-EM example

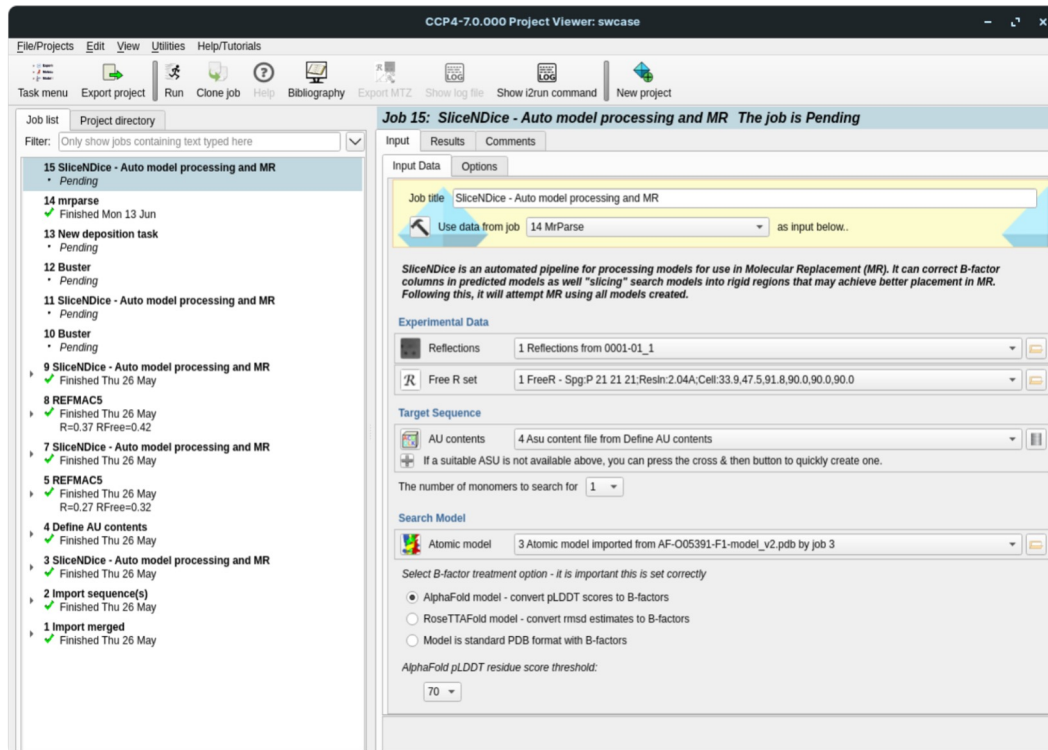
Example 3 - 7FCI or EMD-31526

- NTCP in complex with YN69083 Fab
- Resolution: 3.30 Å map
- Initial AlphaFold2 model very different to PDB (RMSD: 5.1 Å)
- Split into four using the Birch algorithm in Slice'N'Dice (C-alpha based)
- Each slice placed using Molrep
- Placed structure was much more similar to PDB (RMSD: 0.6 Å)

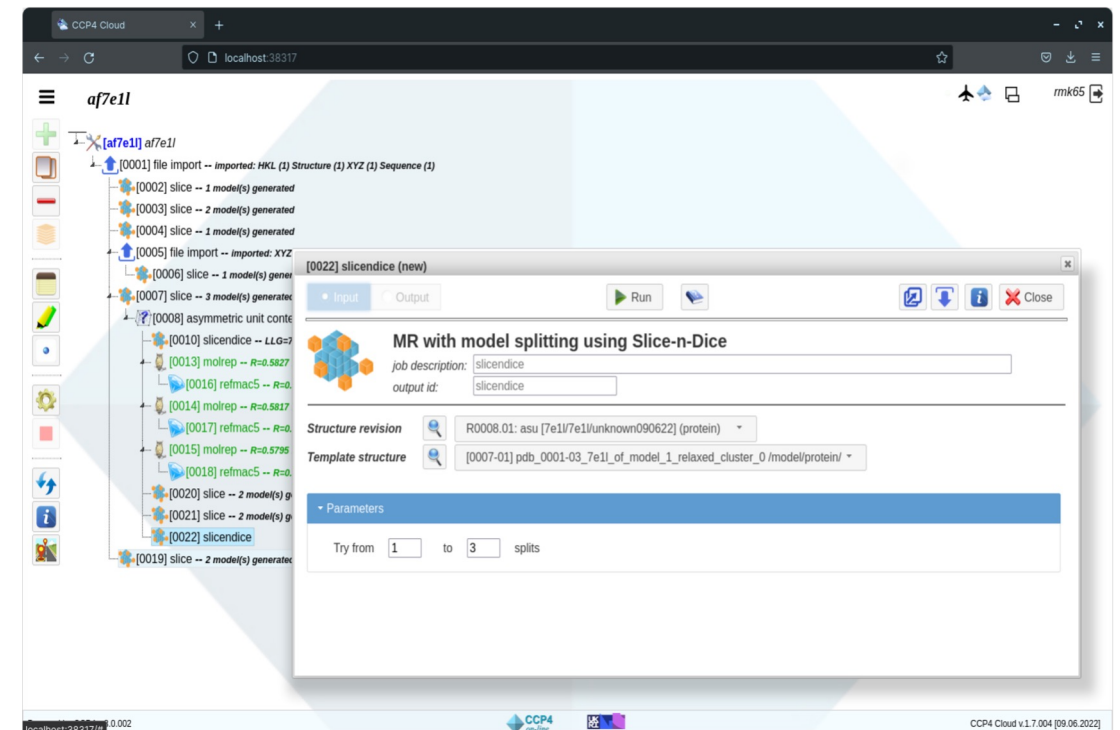


Availability

CCP4i2



CCP4cloud



Acknowledgements

Daniel Rigden

Ronan Keegan

Jens Thomas

Luc Elliott

Eugene Krissinel

Kyle Stevenson

Tom Burnley

Slice'N'Dice: Maximising the value of predicted models for structural biologists
Adam J. Simpkin, Luc G. Elliott, Kyle Stevenson, Eugene Krissinel, Daniel J. Rigden,
Ronan M. Keegan
bioRxiv 2022.06.30.497974; doi: <https://doi.org/10.1101/2022.06.30.497974>

hasimpk@liverpool.ac.uk
Slack @Adam Simpkin

