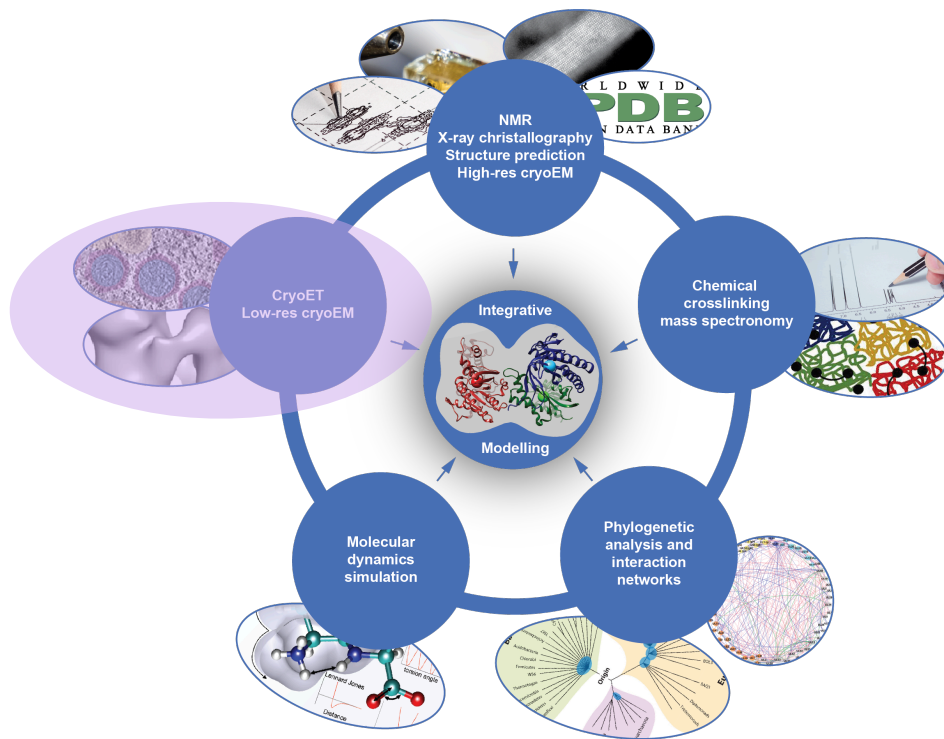


Flex-EM & TEMPy-ReFF

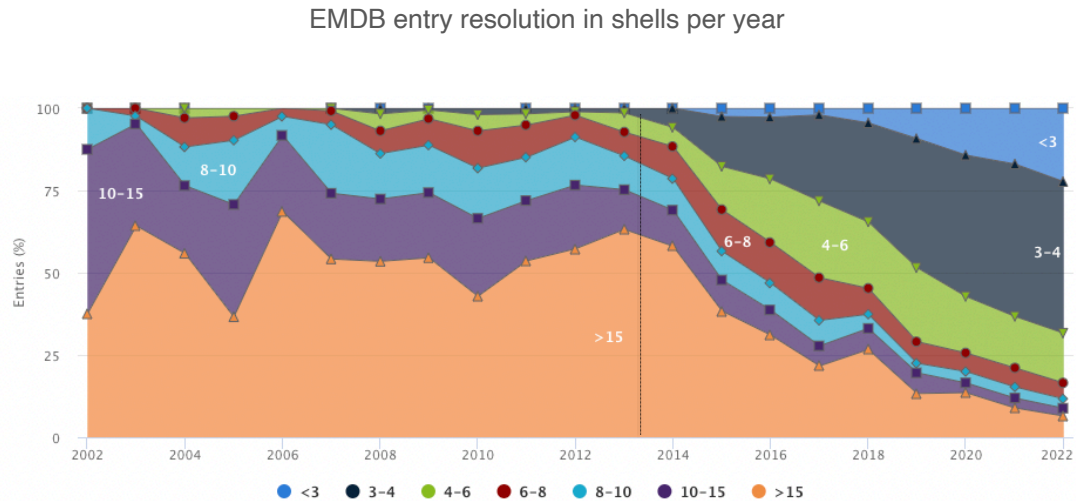
CCP-EM Icknield Model Building Workshop

11.10.2023

Maya Topf (CSSB Hamburg)



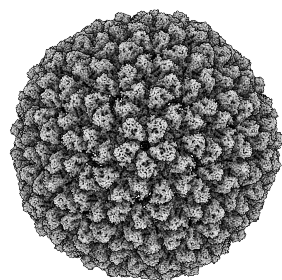
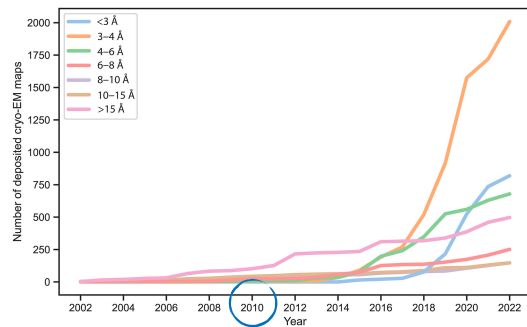
- **Develop** methods for integrative modelling of macromolecular assemblies
- **Apply** the methods to model structures of viral assemblies (e.g. herpesviruses) and other systems (e.g. bacterial pore-forming proteins)
- **Analyse** the structures to further understand their function and identify potential drug targets
- **Distribute** and support new software for the structural biology community



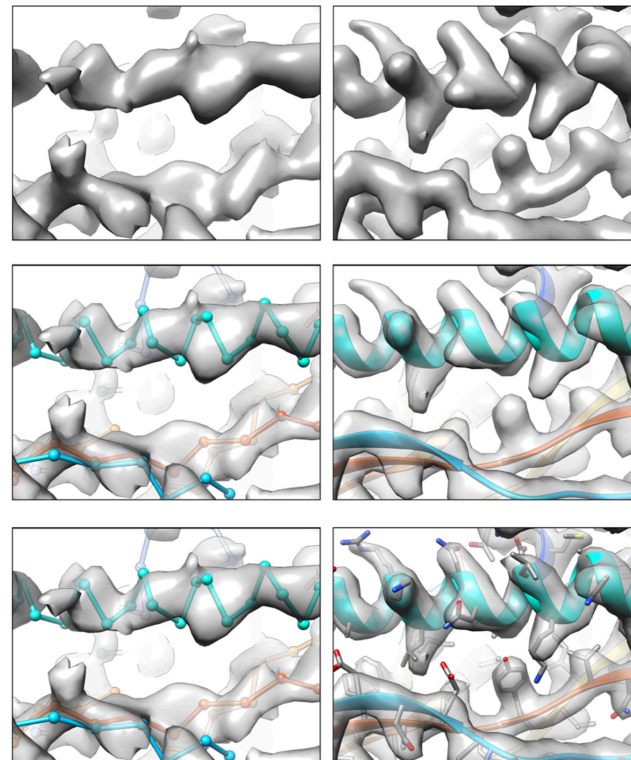
“The Resolution Revolution”
W Kühlbrandt, *Science* 2014

- More than 50% of the maps are now at better than $\sim 5 \text{ \AA}$ resolution

THE “RESOLUTION REVOLUTION”



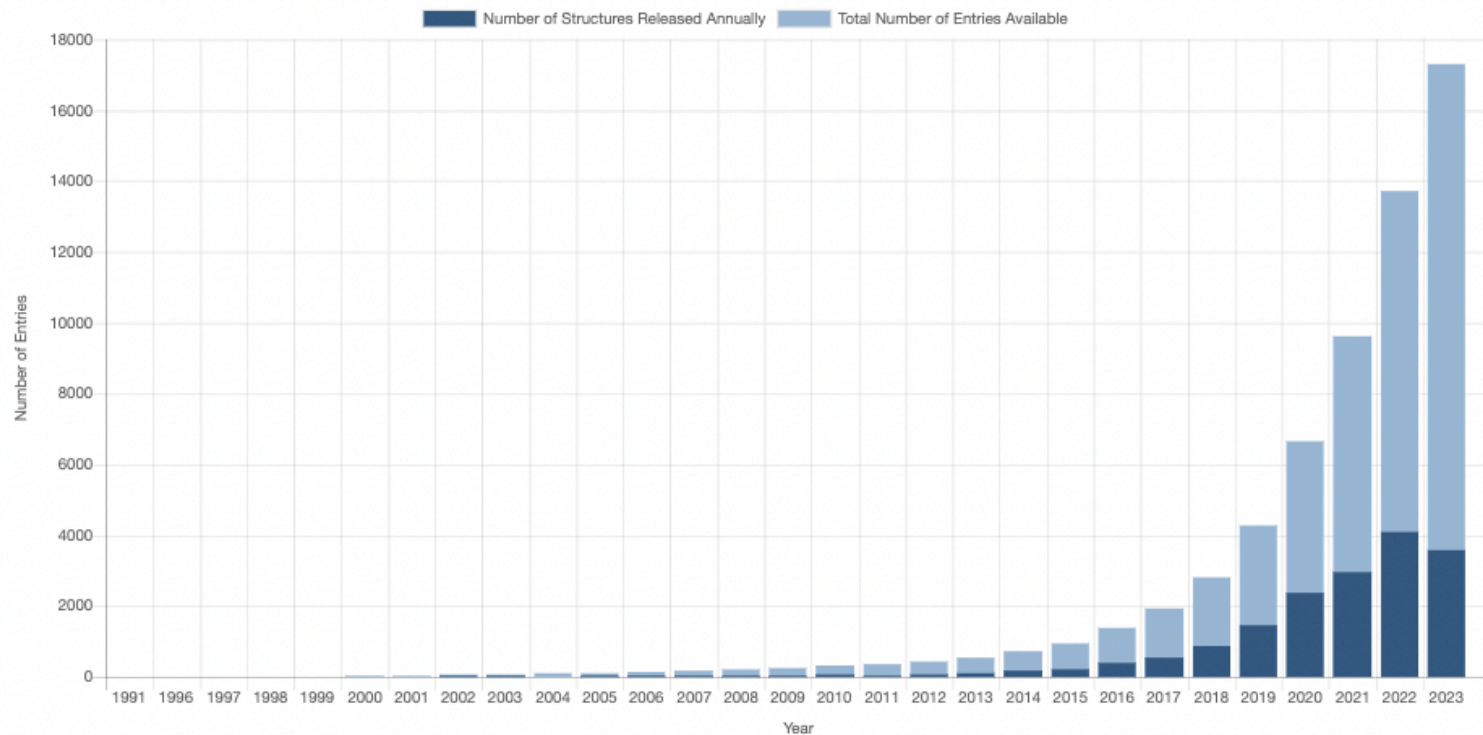
Bacteriophage P22
4 Å (actually ~5 Å)
EMDB 5137
Chen et al., 2011



Chen, 2011. *PNAS*
Non-validated C-alpha trace

Hryc, 2017. *PNAS*
de novo all-atom model

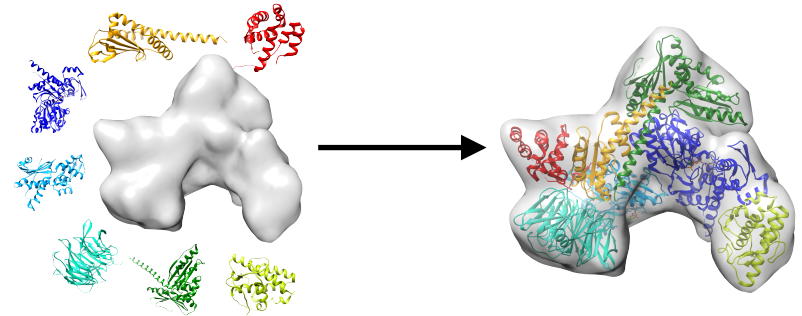
PROTEIN DATA BANK STATISTICS



- 7-8% of the structures in PDB are associated with maps in EMDb
- These days most entries are associated with either a single model or multiple models

What is it?

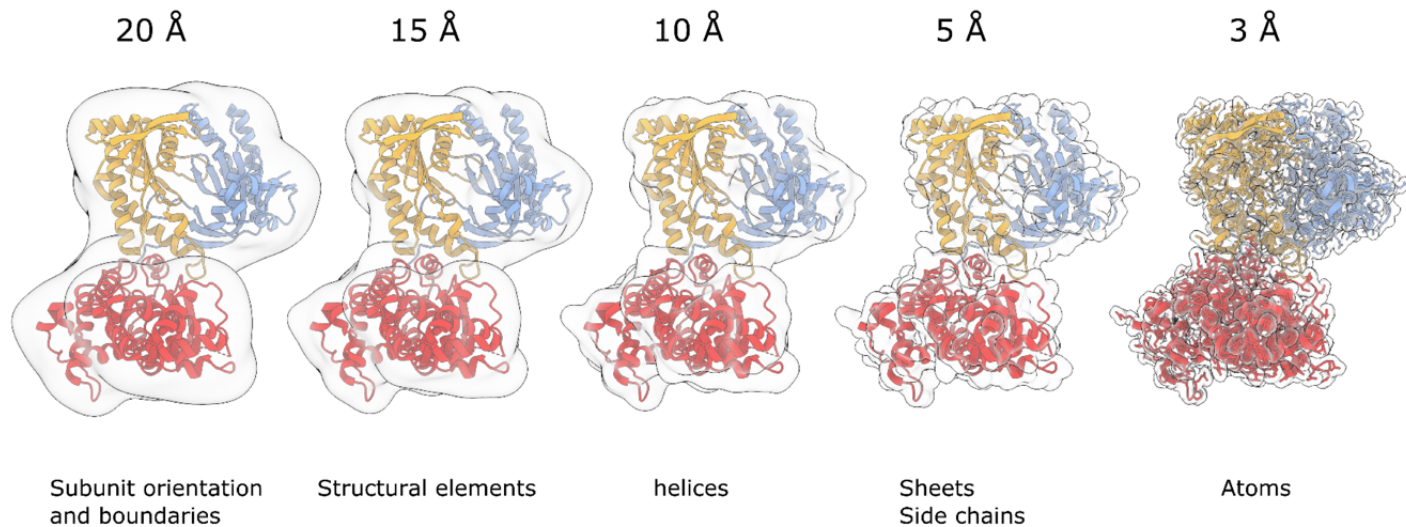
- * Generation, representation, “manipulation” of the 3D structure of biological molecules
- * A molecular model in cryoEM is a compact interpretation of a density map in light of everything known *a priori* about the structure-composition of the macromolecule of interest



Why we need it?

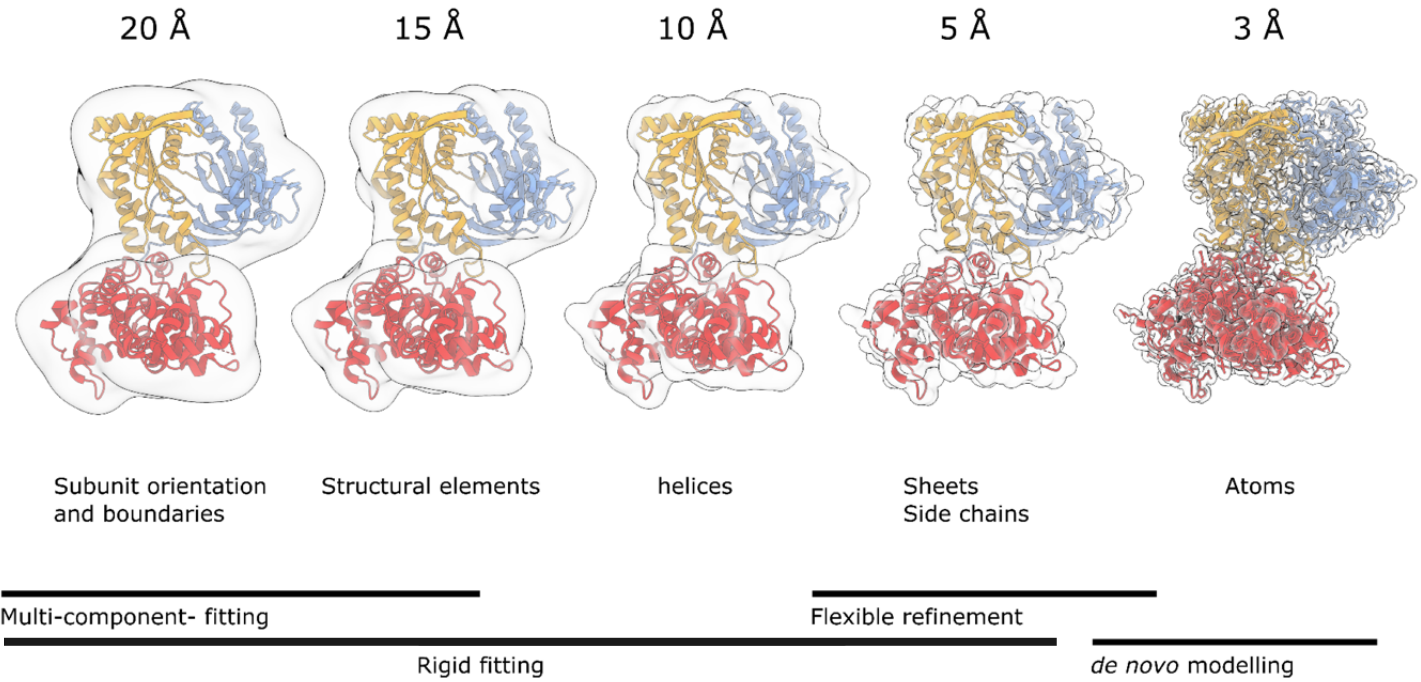
1. To get the atomic 3D structure of the molecules;
2. to know physicochemical characteristics of the molecules;
3. to compare the structure of a molecule with different molecules;
4. to visualize complexes formed between different molecules/macromolecules;
5. and to predict how new related molecules might look.

STRUCTURAL FEATURES AT DIFFERENT RESOLUTION LEVELS

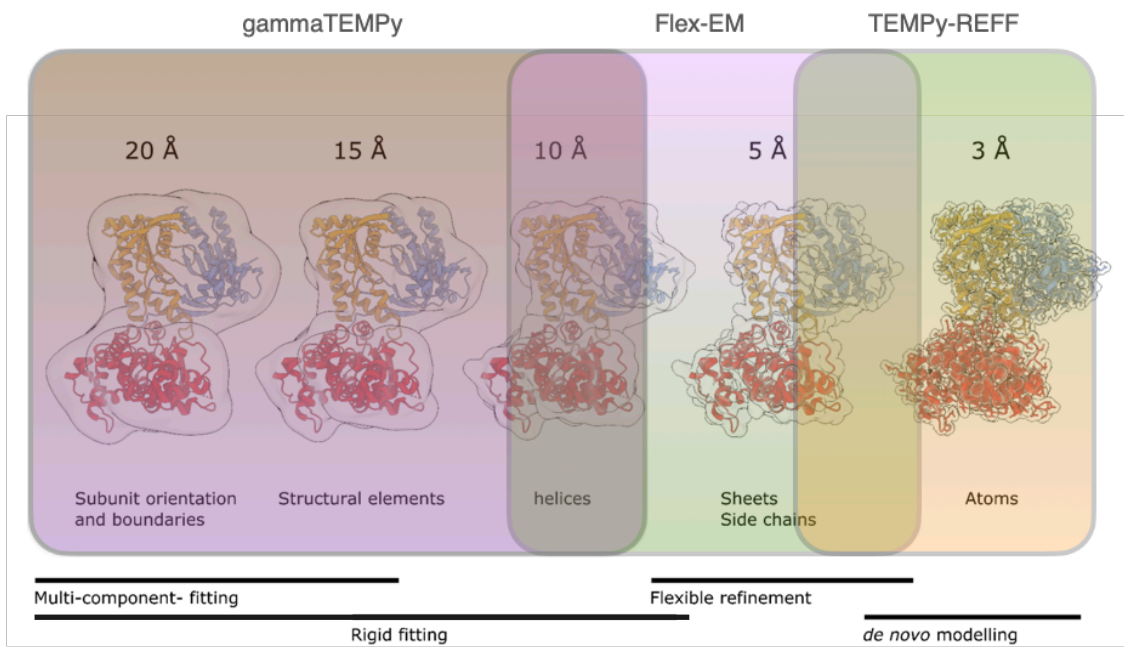
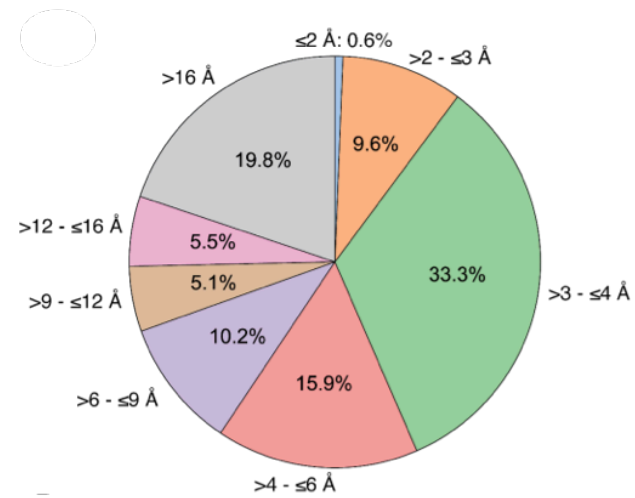


if a complete crystal structure is already available, 10Å data may be sufficient, while if no sequence/composition data is available even 3Å may not suffice

STRUCTURAL FEATURES AT DIFFERENT RESOLUTION LEVELS



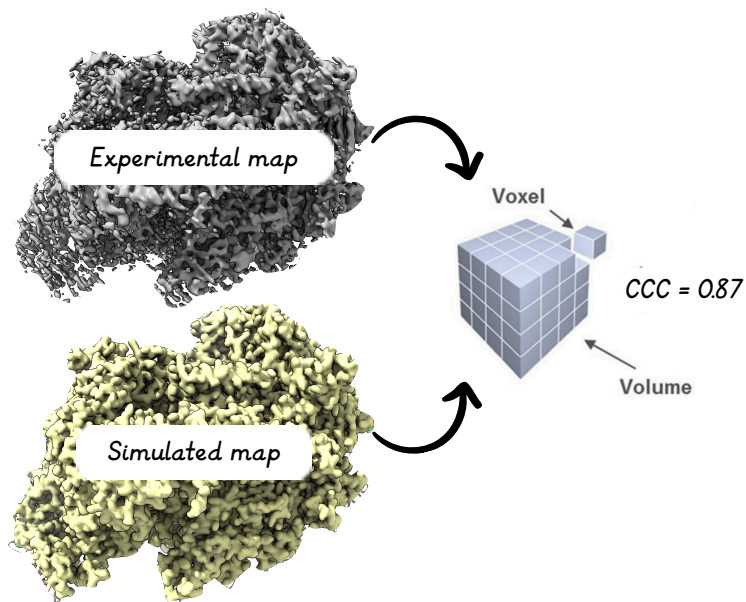
FITTING TOOLS ACROSS THE RESOLUTION SPECTRUM



* Impossible to mention all the modelling softwares developed till now...
 We just want to tell you about some of our tools.

Topf et al. *Structure*, 2008
 Pandurangan et al. *Structure* 2015
 Cargnolini et al. *Proteins* 2021
 Beton et al. *WIREs Comput Mol Sci.* 2023

FITTING AN INITIAL STRUCTURE



The closer the cross-correlation value is to 1, the more closely the sets are to being identical.

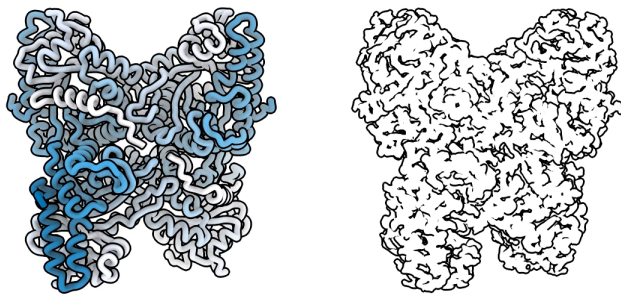
$$\rho(x, y, z) = \sum_N \frac{Z_N}{(\sigma\sqrt{2\pi})^3} e^{-\frac{(x-x_n)^2 + (y-y_n)^2 + (z-z_n)^2}{2\sigma^2}}$$

Cross Correlation Coefficient

$$\text{CCC}_{\rho, \text{lin}} = \int_{\mathbf{x}} (\rho_{\text{obs}}(\mathbf{x}) \rho_{\text{calc}}(\mathbf{x}, \mathbf{m}))^2 d^3\mathbf{x}$$

RIGID (BODY) FITTING (ALL RESOLUTIONS)

- * When possible, a known or pre-calculated model is placed and fitted in the cryo-EM map as a rigid body.
- * Often a full-exhaustive six dimensional (6D) grid search of the three translational and three rotational degrees of freedom is performed.
- * Analysis of all possible solutions is performed to locate the global cross-correlation minimum.



TEMPy2

[Guide](#) [Docs](#)

TEMPy2 is a Python library and set of tools for validating, fitting and refining atomic models in cryo-EM maps

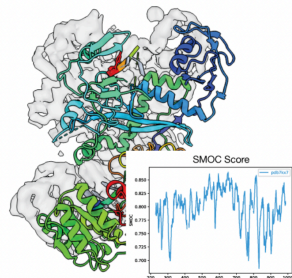
Get TEMPy2

```
$ pip install BioTEMPy==2.0.0
```

Checkout the [Quickstart](#) guide to get started.

TEMPy can do...

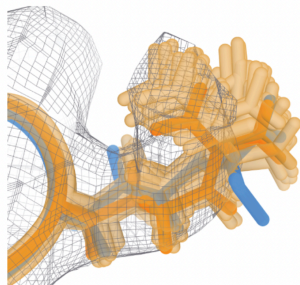
Validation



TEMPy2 offers a variety of local and global validation scores such as LoQFit, SMOC, SCCC and more.

- [Get started](#)
- [Try our tutorial](#)
- [Read more](#)

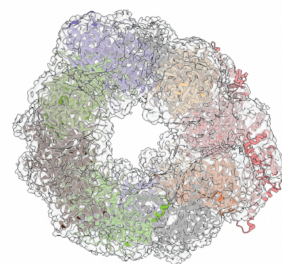
Flexible Fitting and refinement



TEMPy-REFF offers B-factor refinement and can produce ensembles which better explain underlying dynamics.

- [Get started](#)
- [Try our tutorial](#)
- [Read more](#)

Assembly Fitting



Build models into low resolution density with γ -TEMPy using a genetic algorithm.

- [Get started](#)
- [Try our tutorial](#)
- [Read more](#)

Map processing

- Transformation (rotation/translation)
- Filters

Fitting

- Local random and exhaustive search
- Multicomponent fitting (γ -TEMPy)
- Refinement (TEMPy-REFF)

Model processing

- Transformation (rotation/translation)
- Model-to-map
- Ensemble generation
- Clustering

Global Scoring

- Density based
- Surface based

Local scoring

- SCCC
- SMOC
- LoQfit

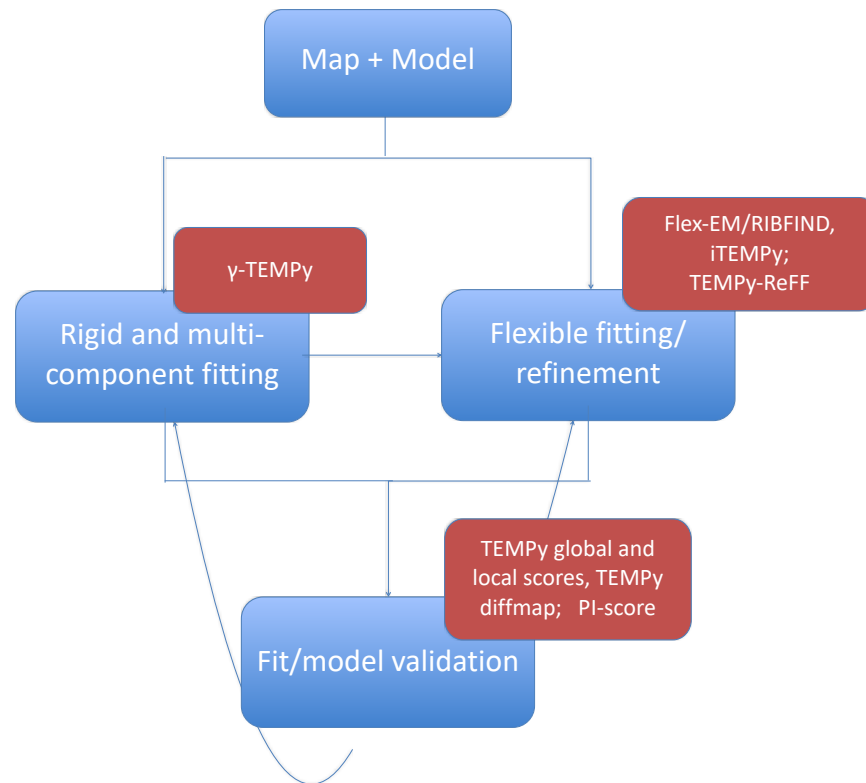
Map comparison

- Difference maps
- FSC

Chimera scripts

Plots
Attribute files
PDB files

TEMPY2 TOOLS FOR MODELLING STRUCTURES FROM CRYO-EM MAPS

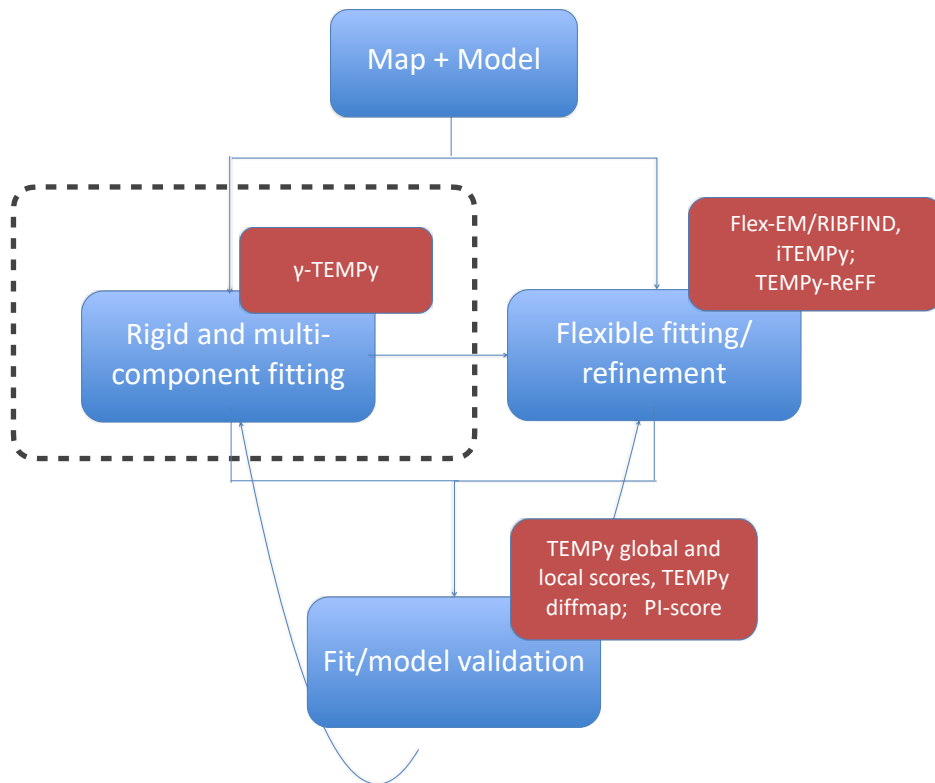
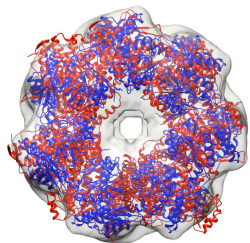


Topf et al. *Structure*, 2008
Vasishthan et al. *J Struct Biol* 2011
Farabella et al. *J App Cryst* 2015
Joseph et al. *Methods* 2016
Cragolini et al. *Acta Cryst D* 2021
Malhotra et al. *Nat Comm* 2021
Cargnolini et al. *Proteins* 2021

<https://topf-lab.gitlab.io/tempy>

TEMPy2: NumPy and SciPy, Gemmi, Matplotlib, OpenMM

TEMPY2 TOOLS FOR MODELLING STRUCTURES FROM CRYO-EM MAPS

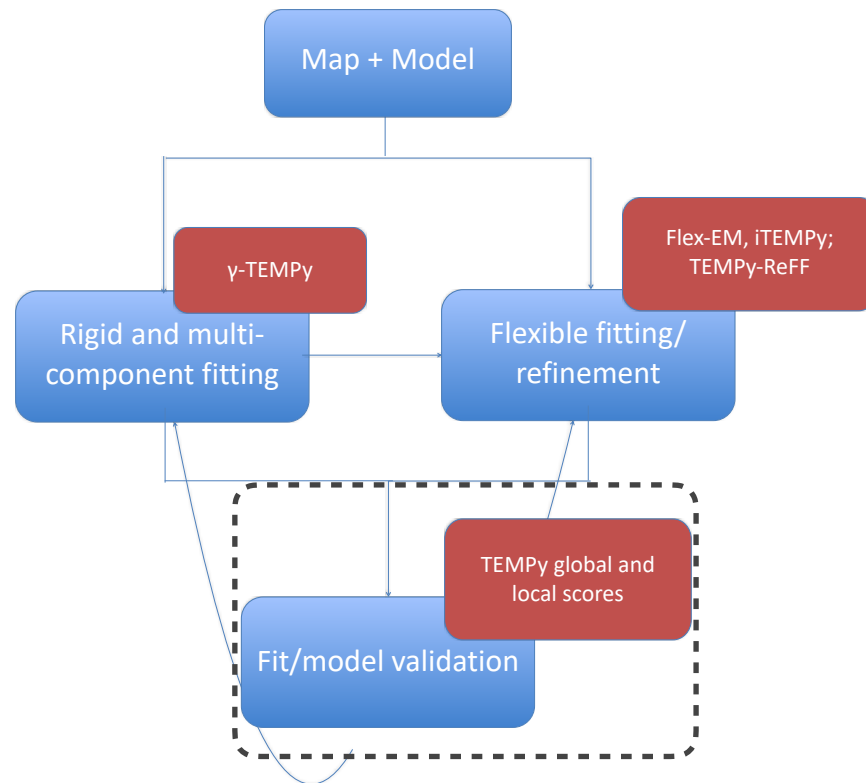


Topf et al. *Structure*, 2008
Vasishthan et al. *J Struct Biol* 2011
Farabella et al. *J App Cryst* 2015
Joseph et al. *Methods* 2016
Cragolini et al. *Acta Cryst D* 2021
Malhotra et al. *Nat Comm* 2021
Cargolini et al. *Proteins* 2021

<https://topf-lab.gitlab.io/tempy>

TEMPy2: NumPy and SciPy, Gemmi, Matplotlib, OpenMM

TEMPY2 TOOLS FOR MODELLING STRUCTURES FROM CRYO-EM MAPS



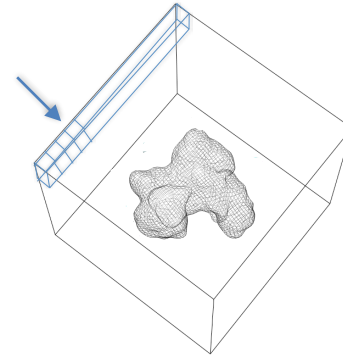
Topf et al. *Structure*, 2008
Vasishthan et al. *J Struct Biol* 2011
Farabella et al. *J App Cryst* 2015
Joseph et al. *Methods* 2016
Cragolini et al. *Acta Cryst D* 2021
Malhotra et al. *Nat Comm* 2021
Cargolini et al. *Proteins* 2021

<https://topf-lab.gitlab.io/tempy>

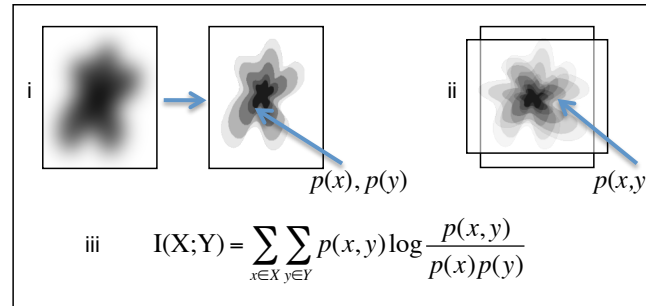
TEMPy2: NumPy and SciPy, Gemmi, Matplotlib, OpenMM

- Cross-correlation coefficient (CCC)

$$CCC = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$



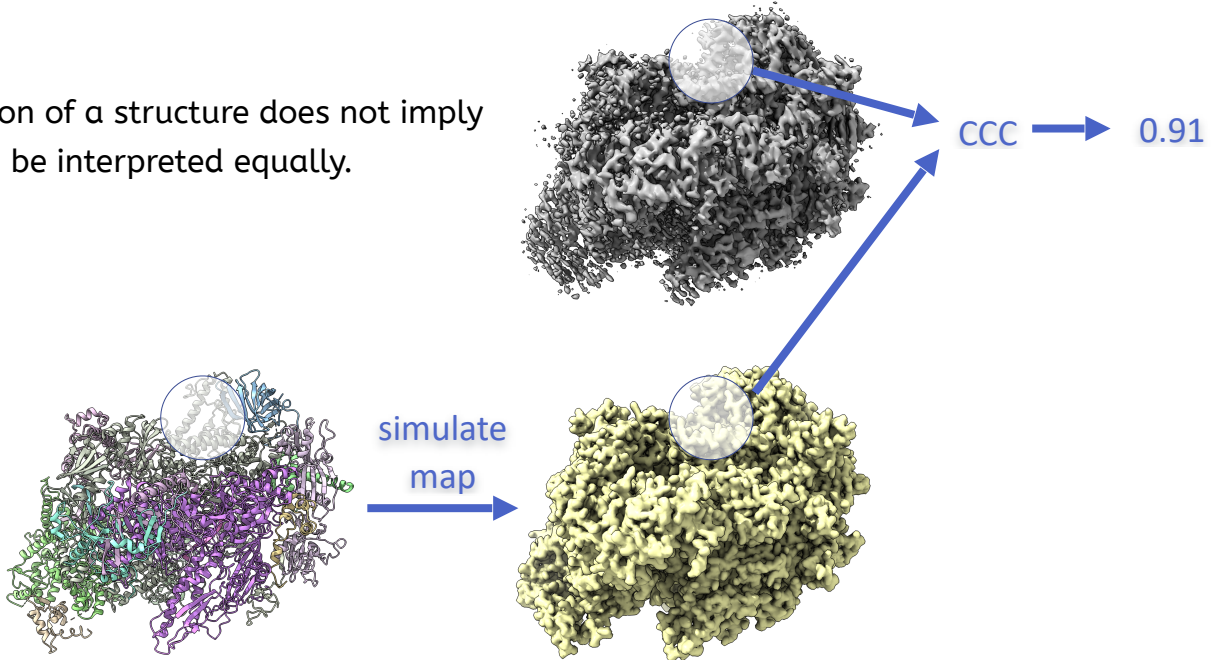
- Mutual information-based score (MI)



Useful at intermediate resolutions; noisy maps; less sensitive to relative intensity levels

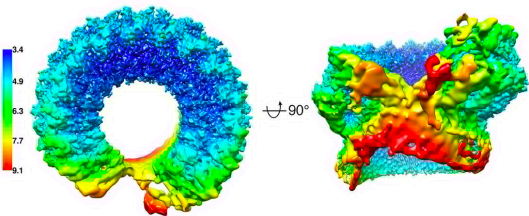
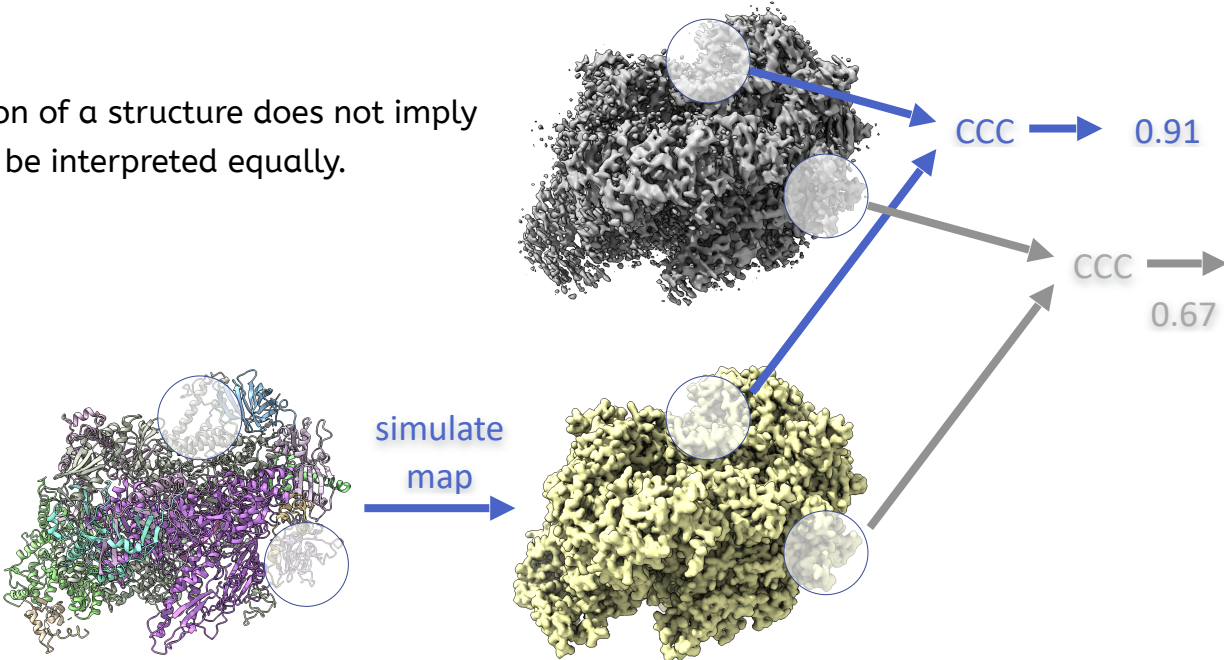
LOCAL RESOLUTION AND LOCAL CORRELATION

* The overall resolution of a structure does not imply that all regions can be interpreted equally.



LOCAL RESOLUTION AND LOCAL CORRELATION

* The overall resolution of a structure does not imply that all regions can be interpreted equally.



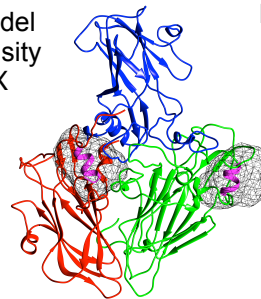
* Local resolution maps are useful for estimating resolution variability, but visual inspection is always essential to assess map quality.

LOCAL SCORING FOR MEDIUM-TO-HIGH RESOLUTION MAPS

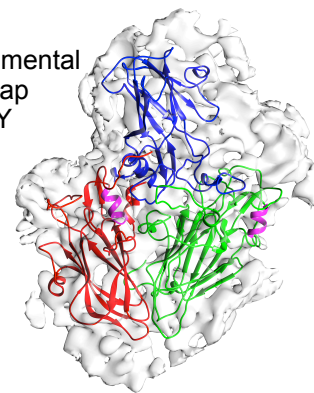
SCCC: Segment-based cross-correlation coefficient:

$$\text{SCCC} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Model density
X



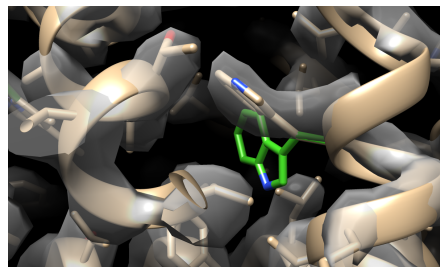
Experimental map
Y



Farabella et al. *J Appl. Cryst.* 2015

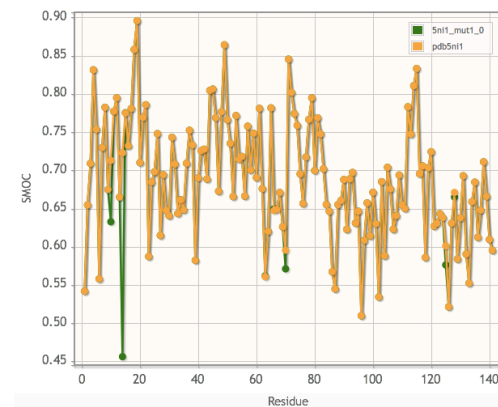
SMOC: An overlap coefficient is calculated over voxels covered by each **Residue** (and the local neighbourhood):

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



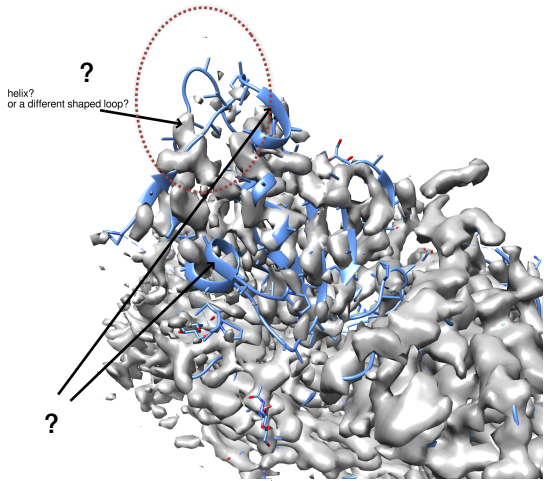
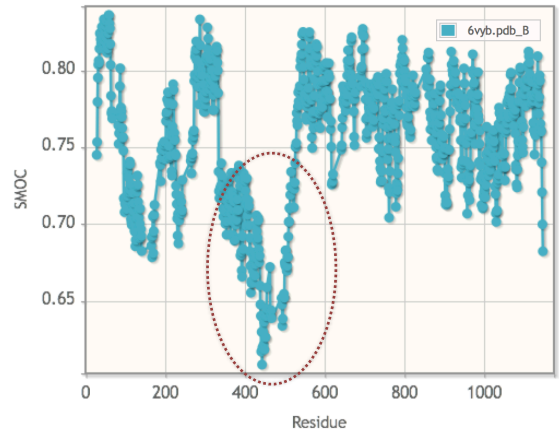
EMD-3488 (3.2Å)
Deposited model PDB: 5NI1

Joseph et al. *Methods* 2016



EXAMPLE

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

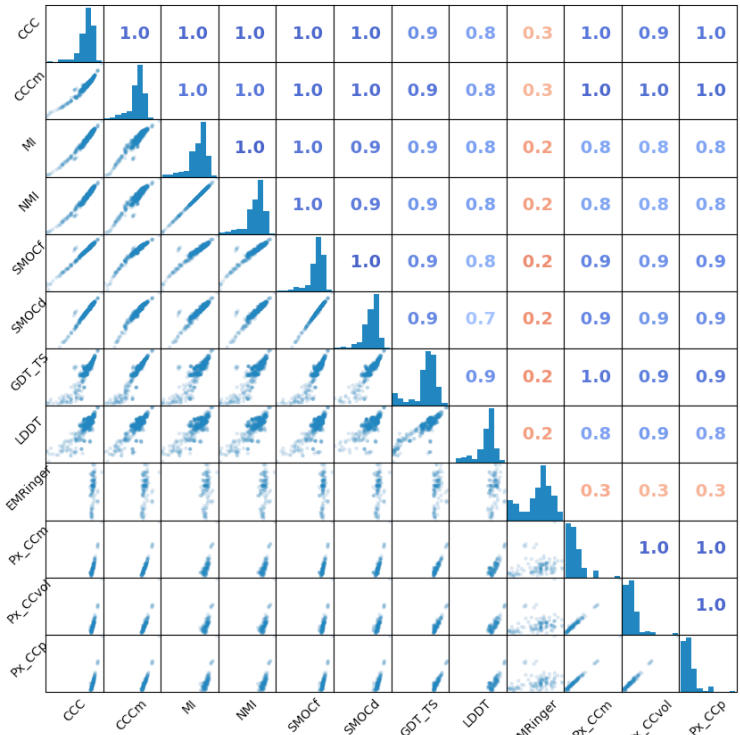
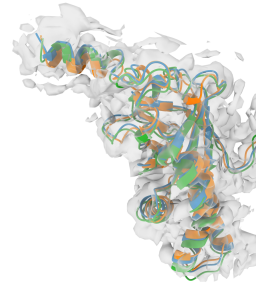
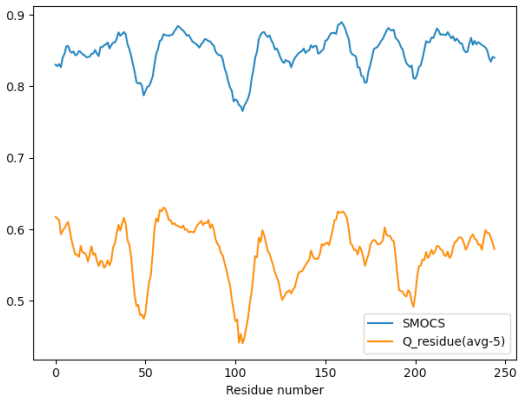


SARS-CoV-2 spike ectodomain (3.2 Å resolution)

CASP14 ASSESSMENT OF MODELS IN RESPECT TO CRYO-EM MAP

Nearly all scores exhibit a relatively high degree of correlation and this is true across all targets

Strong correlations with GDT_TS



T1092-D2

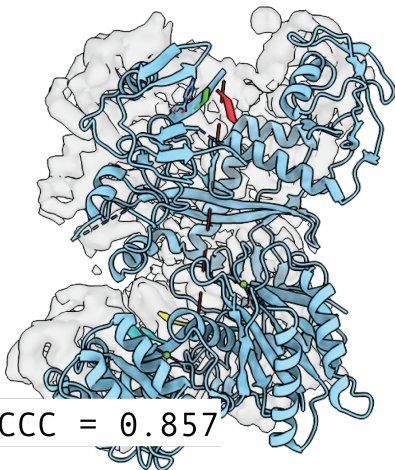
AR9 RNA polymerase

TEMPY2 SCORING FUNCTIONS

- There are many different scoring functions in TEMPy, which can be grouped into:

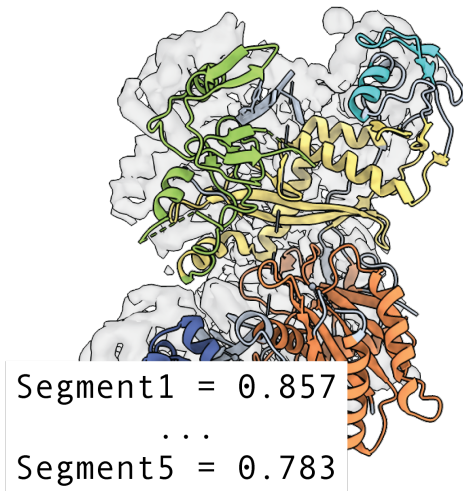
Global Scores

(one number per model)



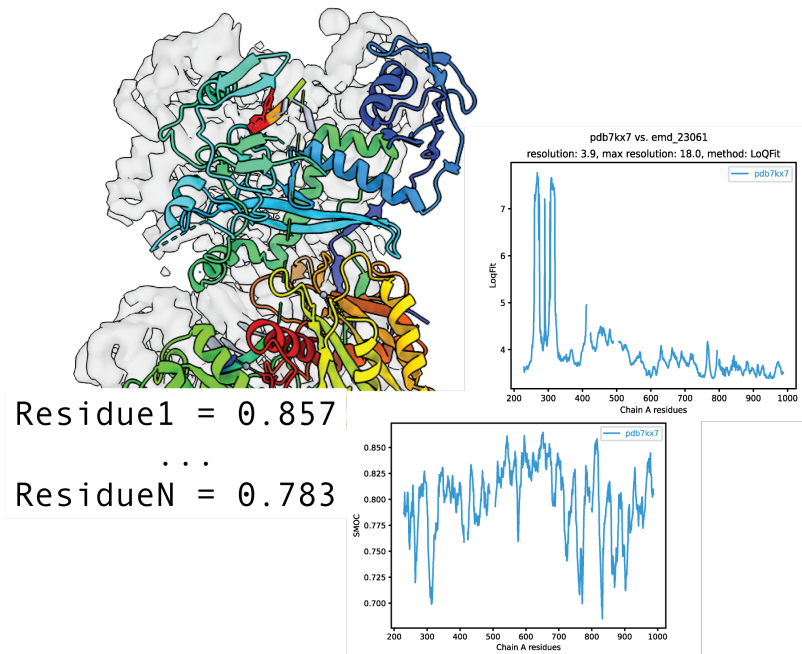
Segmented Scores

(one number per segment)

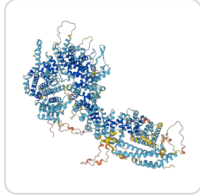


Residue Scores

(one number per residue)

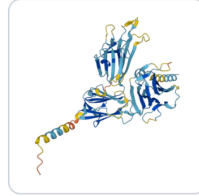


Enter the “AI” era



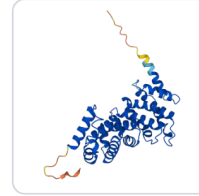
**Nuclear pore complex
protein Nup205**

Part of a large complex that
acts as a gateway in and out
of the cell nucleus



**Gametocyte surface
protein P45/48**

From the malaria parasite;
a candidate protein for
including in vaccines



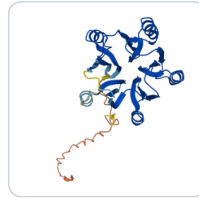
**CCR4-NOT transcription
complex subunit 9**

Regulates an important
cellular process (the rate
of mRNA degradation)



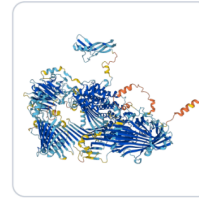
Ice nucleation protein

Bacterial protein that can trigger
ice formation at relatively high
temperatures, causing frost
damage to plants



F20H23.2 protein

Plant protein; represents
a potential new structural
superfamily unlike anything
seen before

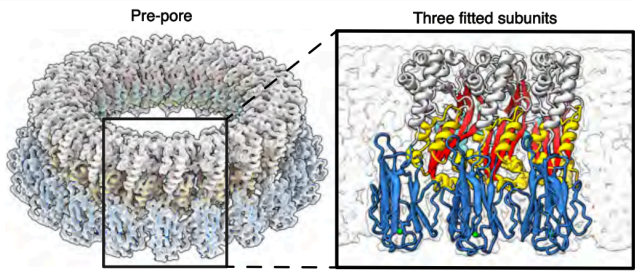


Vitellogenin

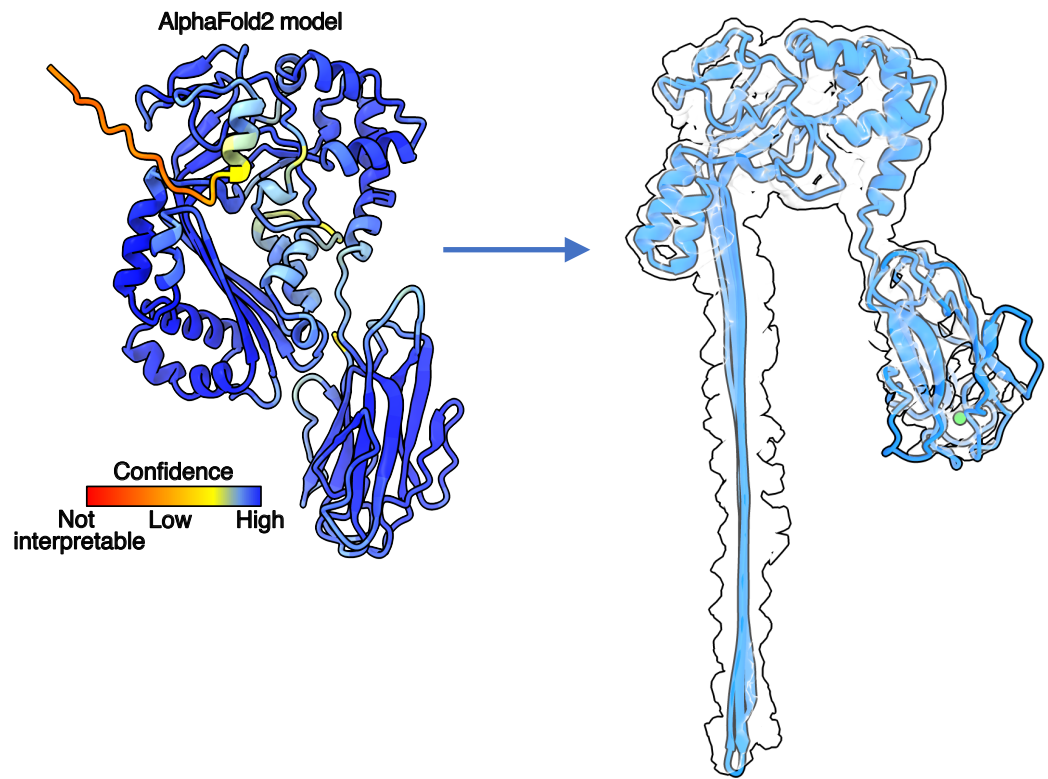
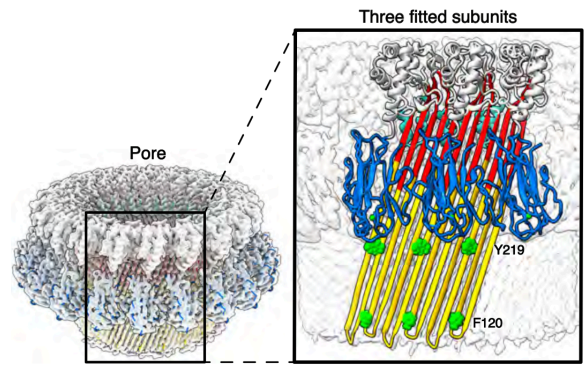
Involved in the immune
system of egg-laying animals
including honeybees

GOOD INITIAL PREDICTIONS FOR PROTEIN ASSEMBLIES BUT...

Mpf2Ba1 pre-pore

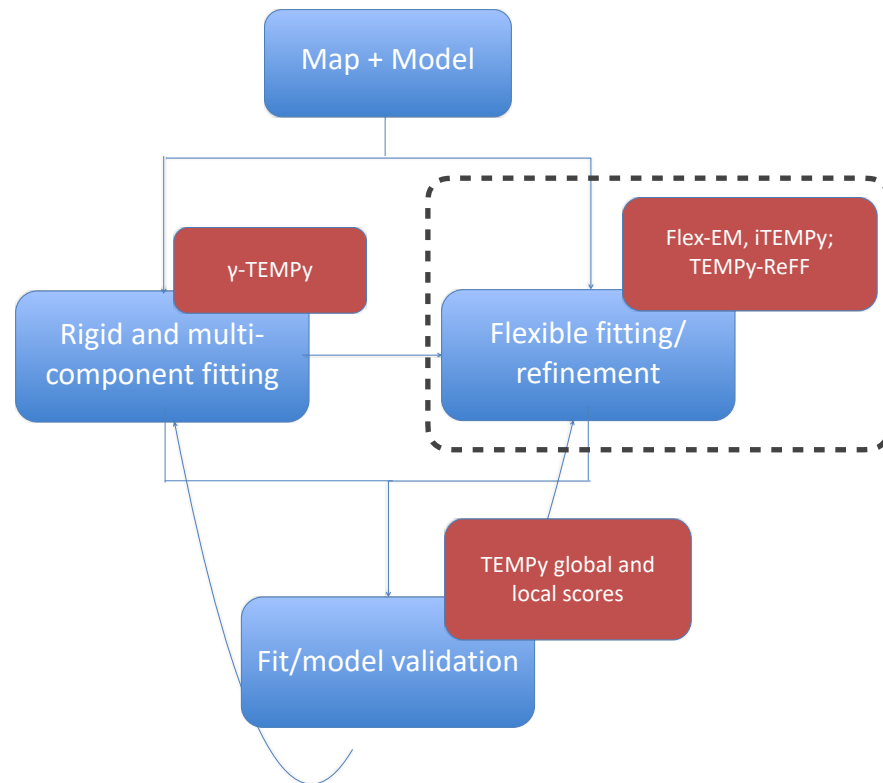


Mpf2Ba1 pore

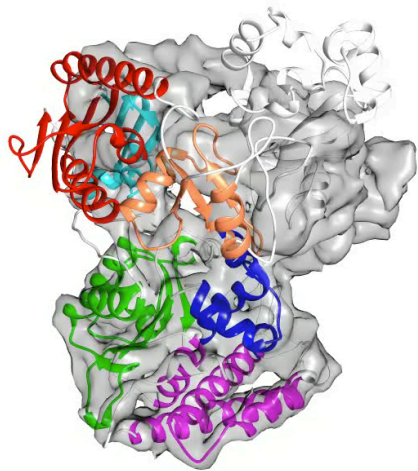


Using validation approaches, we can *iteratively* improve the fit to the density

TEMPY2 TOOLS FOR MODELLING STRUCTURES FROM CRYO-EM MAPS



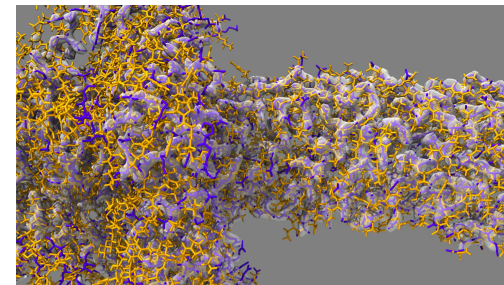
FLEX-EM REFINEMENT AT MEDIUM TO LOW RESOLUTION



sub-domains



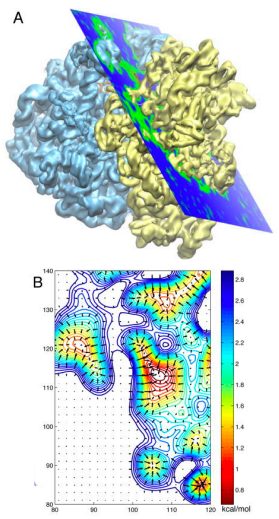
secondary structure elements



residue / all-atom

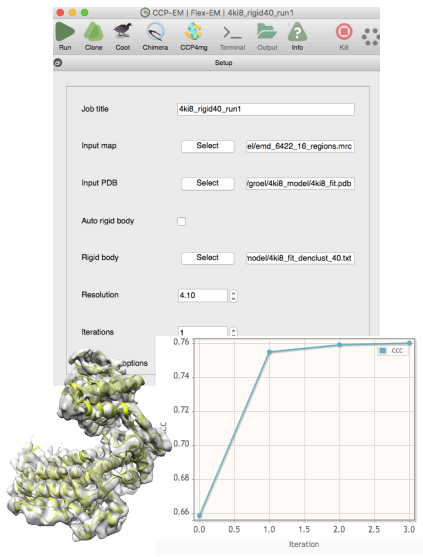
FLEXIBLE FITTING AND REFINEMENT

Different approaches for flexible fitting have been developed over the years.



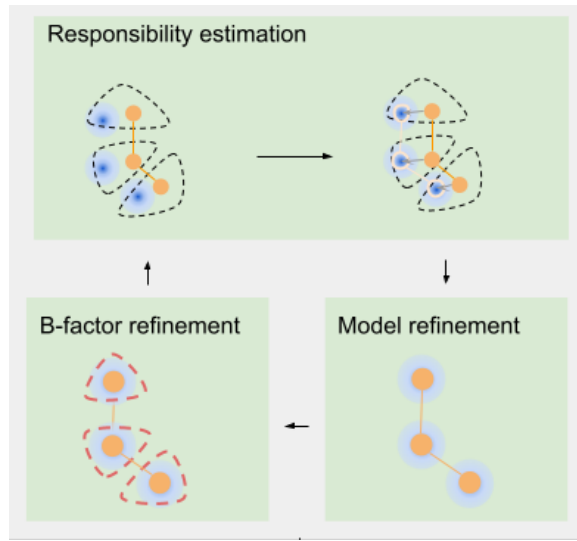
Molecule Dynamics Flexible Fitting (MDFF) uses a standard force field with a potential based on the density

Trabuco et al., *Methods* 2008



Flex-EM improves the correlation between the model and the cryo-EM map

Topf et al., 2008
Joseph et al., 2016



TEMPy-ReFF (Responsibility-based Flexible-Fitting) utilizes a gaussian mixture model with expectation-maximisation to improve the model in the map

Cragolini et al., 2021, *BioRxiv* 2022
Malhotra et al., 2023

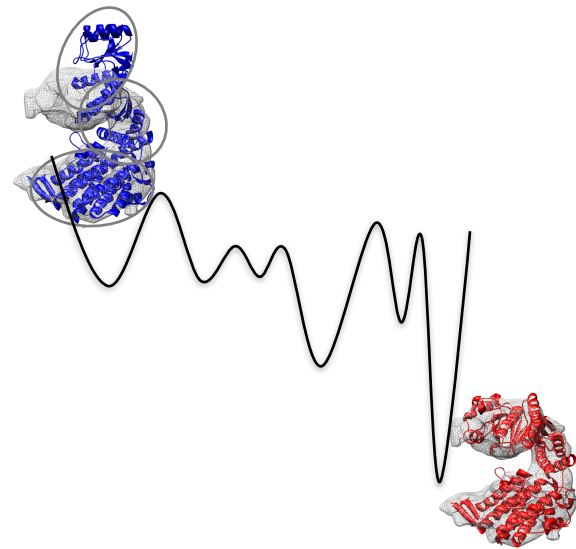
- During the refinement, the atoms are displaced in the direction that maximizes their cross-correlation with the cryoEM density map (E^{CC}) and minimizes the violations of the stereochemical (E^{SC}) and non-bonded contacts (E^{NB}):

$$E = w_1 * E^{CC}(P) + w_2 * E^{SC}(P) + w_3 * E^{NB}(P)$$

- Optimisation is performed on rigid bodies (b) by energy minimisation and simulated annealing molecular dynamics:

$$\vec{F}(b_l) = - \sum_{j \in Atom(b_l)} \frac{\partial E(b_l)}{\partial \vec{r}_j}$$

- Atoms in a rigid body move together during the course of refinement.



TEMPY-REFF: REFINEMENT USING MIXTURE MODELS

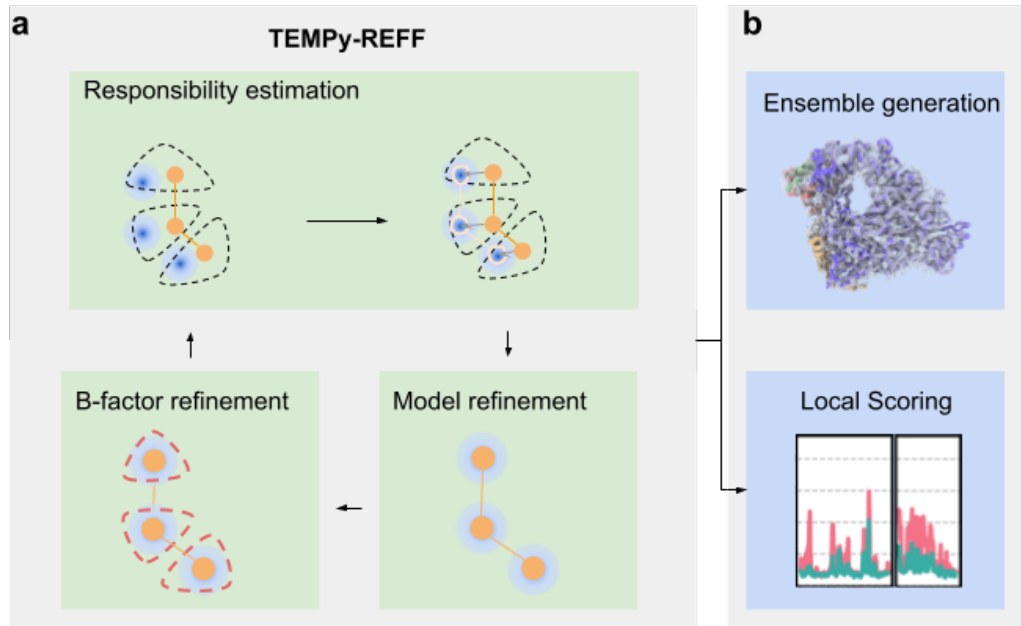
- Use an expectation-maximisation approach
- Use one gaussian per atom, and a background noise term
- Improves atomic position and variance (B-factor)

TEMPy and OpenMM

Amber14 forcefield

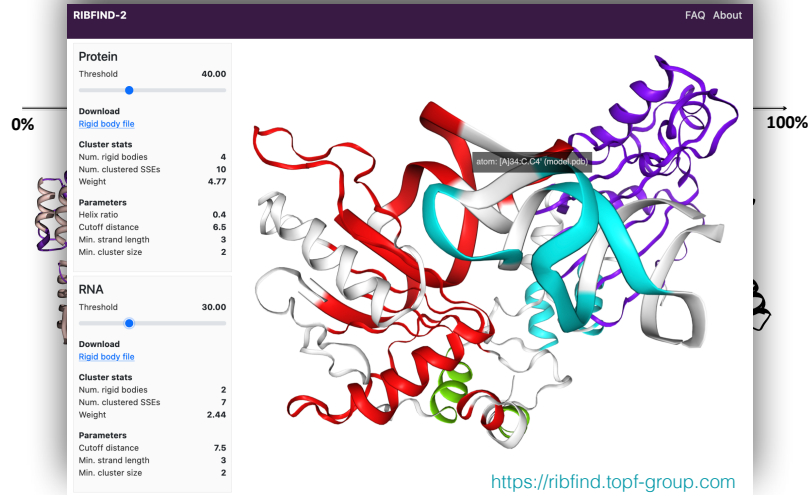
Generalised Born implicit solvent model

Integration with Langevin dynamics, at 100K



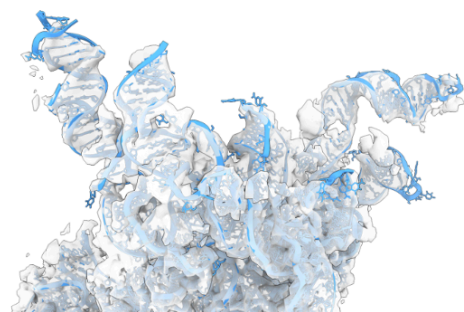
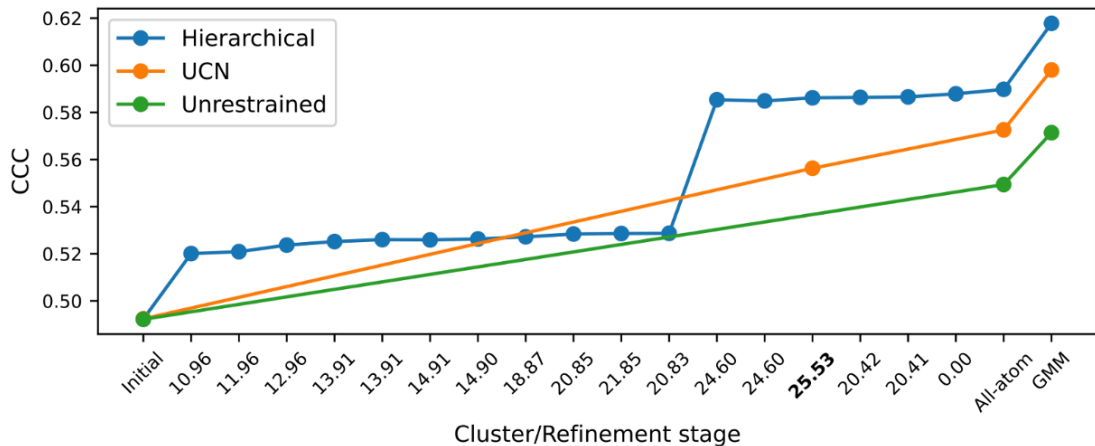
RIBFIND2 - IDENTIFY RIGID BODIES IN PROTEIN AND RNA STRUCTURES

- When the resolution of density map is insufficient to fit smaller entities like individual residues or atoms, we use RIBFIND-defined rigid bodies.
- Clustering of secondary structures together (SSEs computed using DSSP or RNAView).
- Allows faster large body movements in the initial stages or refinement.

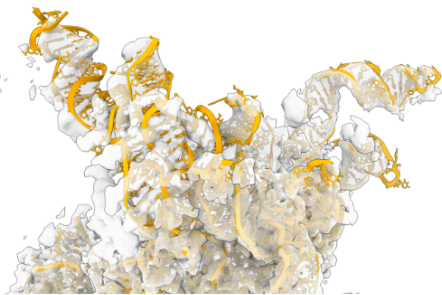


RIBFIND2 - HELPS TO AVOID GETTING "STUCK" IN LOCAL OPTIMA

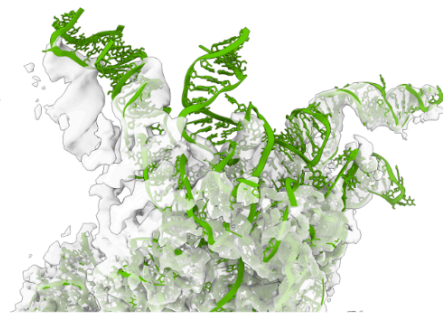
- human SSU processome
- pre-A1 into state post-A1
- EMD-23938, 2.7 Å resolution



Hierarchical

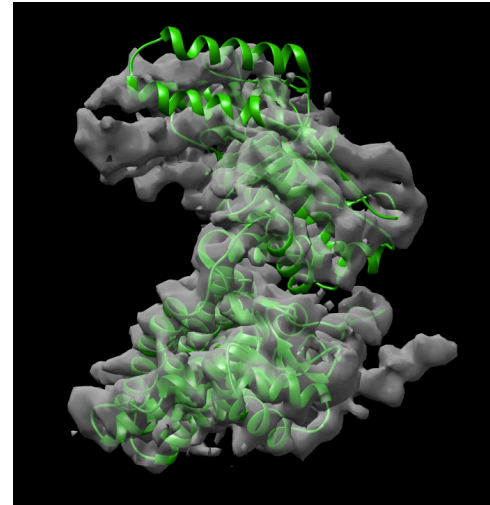
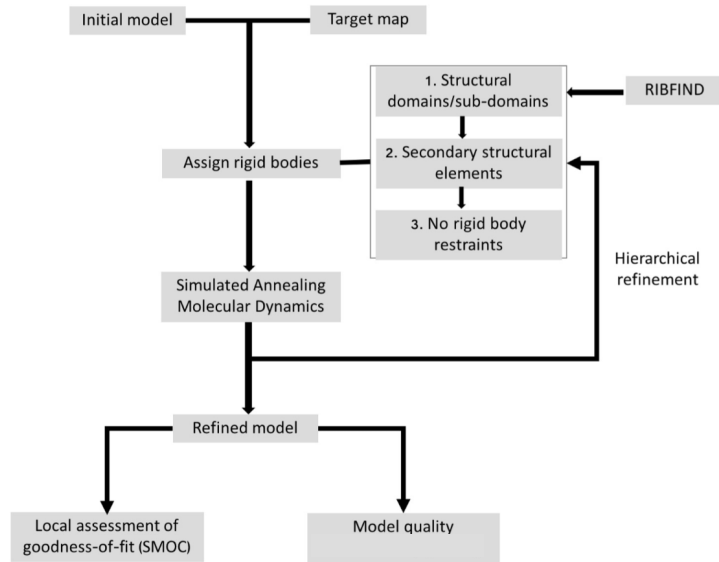


UCN



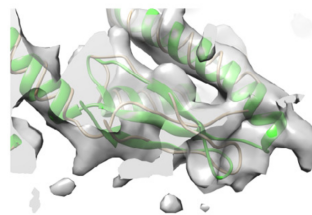
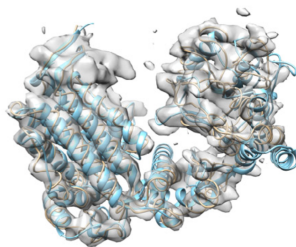
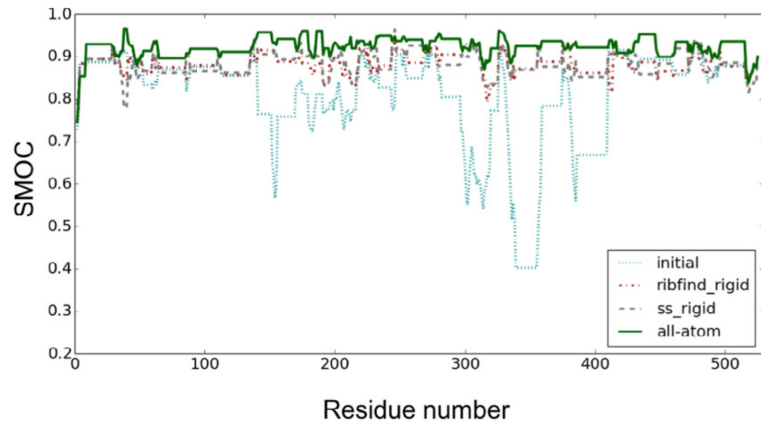
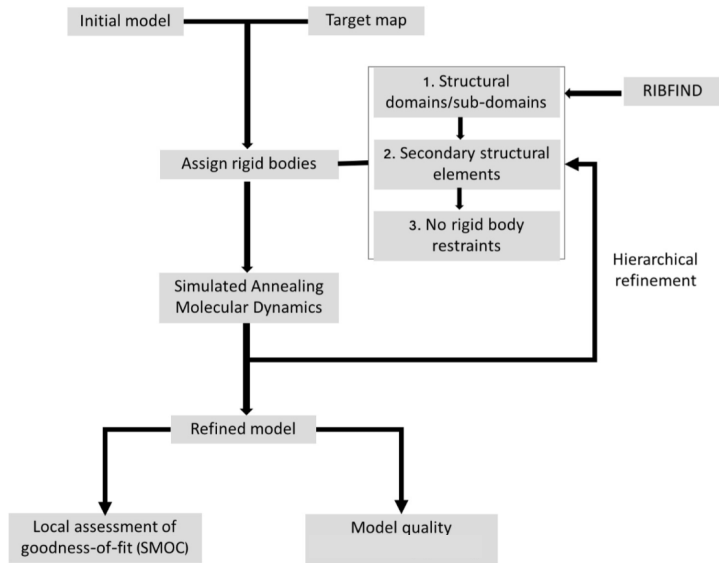
Unrestrained

FLEXIBLE FITTING AT MEDIUM RESOLUTION BY HIERARCHICAL REFINEMENT



Unliganded GroEL at 4.2 Å resolution (EMD-5001)

FLEXIBLE FITTING AT MEDIUM RESOLUTION BY HIERARCHICAL REFINEMENT



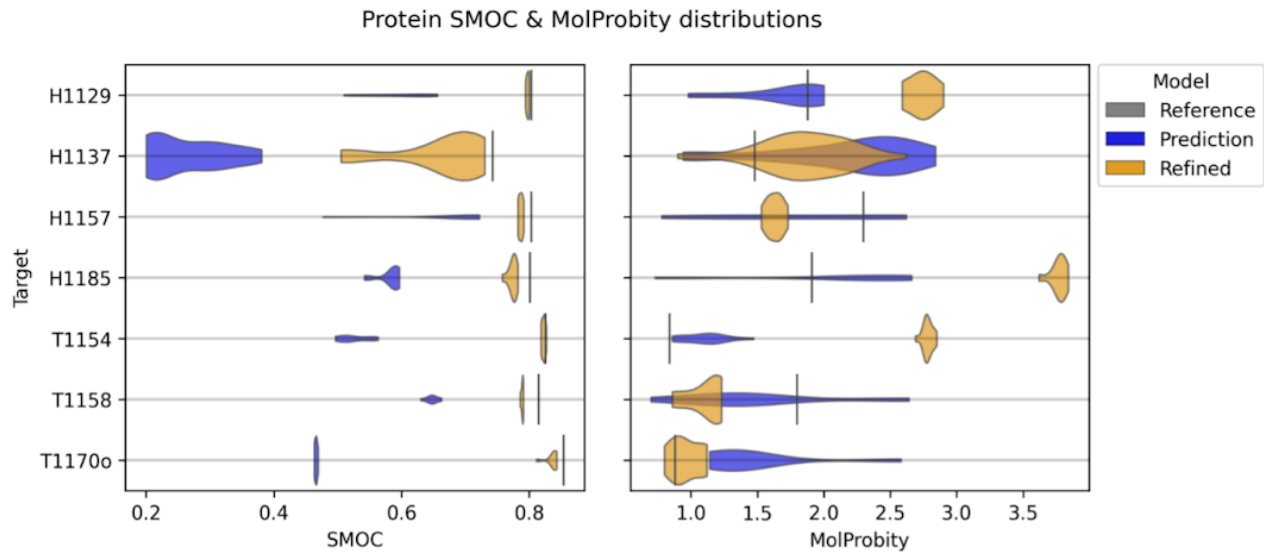
Unliganded GroEL at 4.2 Å resolution (EMD-5001)

Initial model: ADP-bound GroEL (PDB: 4KI8)

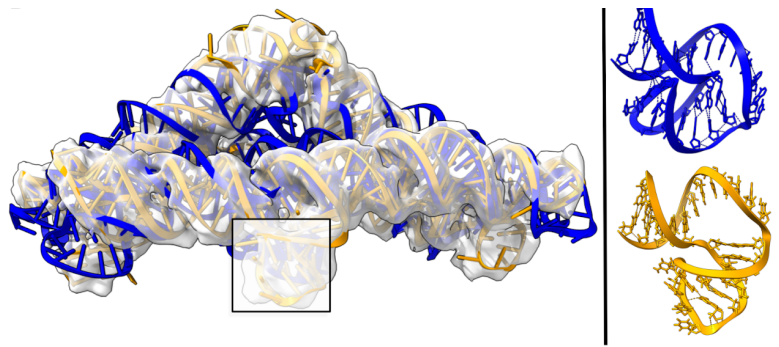
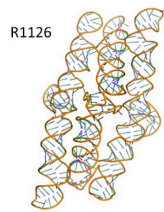
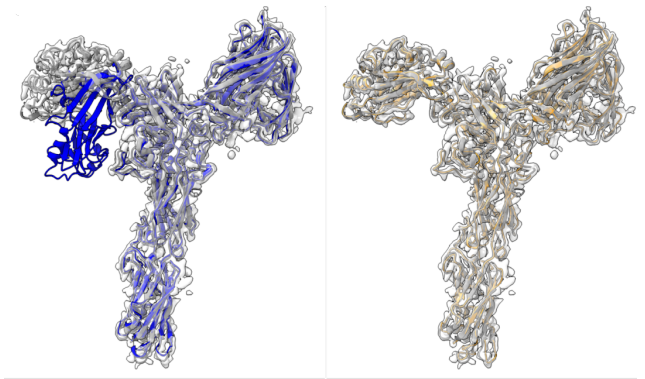
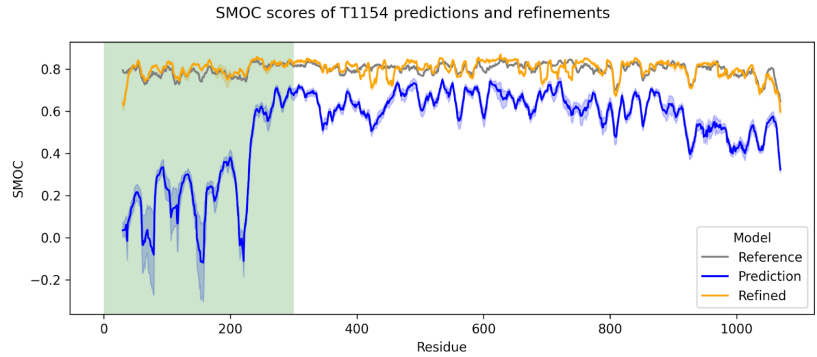
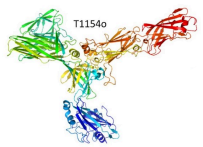
Refined model

Deposited model (PDB: 3CAU)

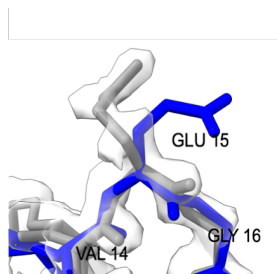
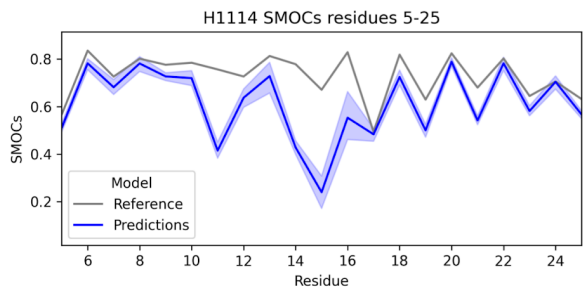
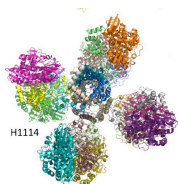
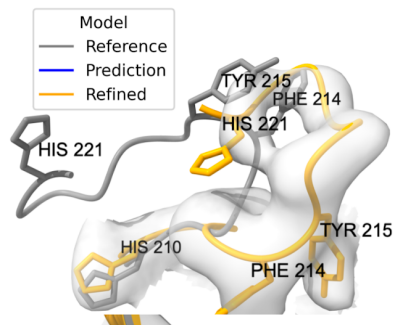
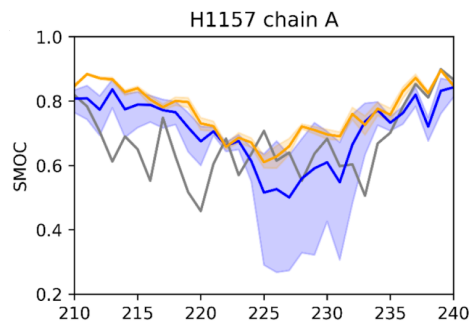
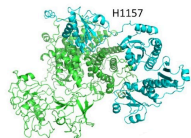
TEMPY-REFF REFINEMENT OF CASP15 MODELS IN CRYO-EM MAPS



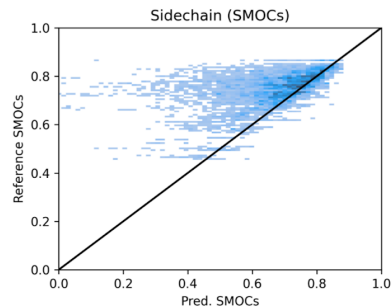
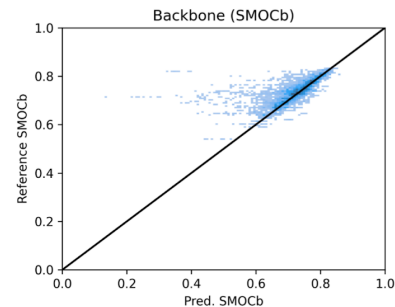
TEMPY-REFF REFINEMENT OF CASP15 MODELS (PROTEIN COMPLEXES AND RNA)



EXPERIMENTAL VALIDATION AND REFINEMENT BY CRYO-EM (CASP15)

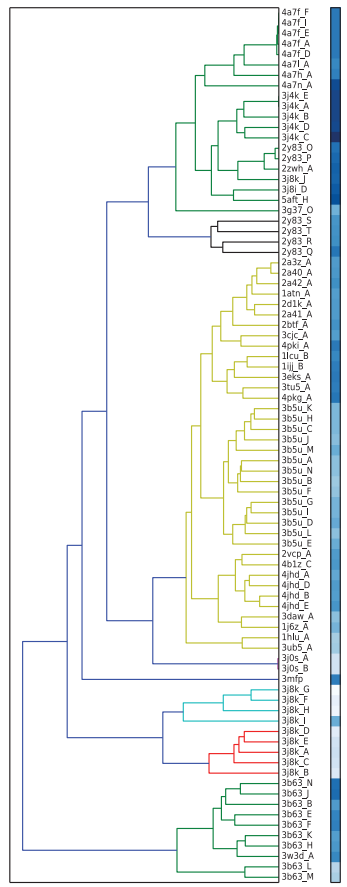
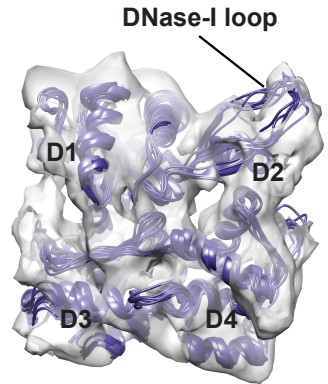


SMOC score distributions for H1114 chain A



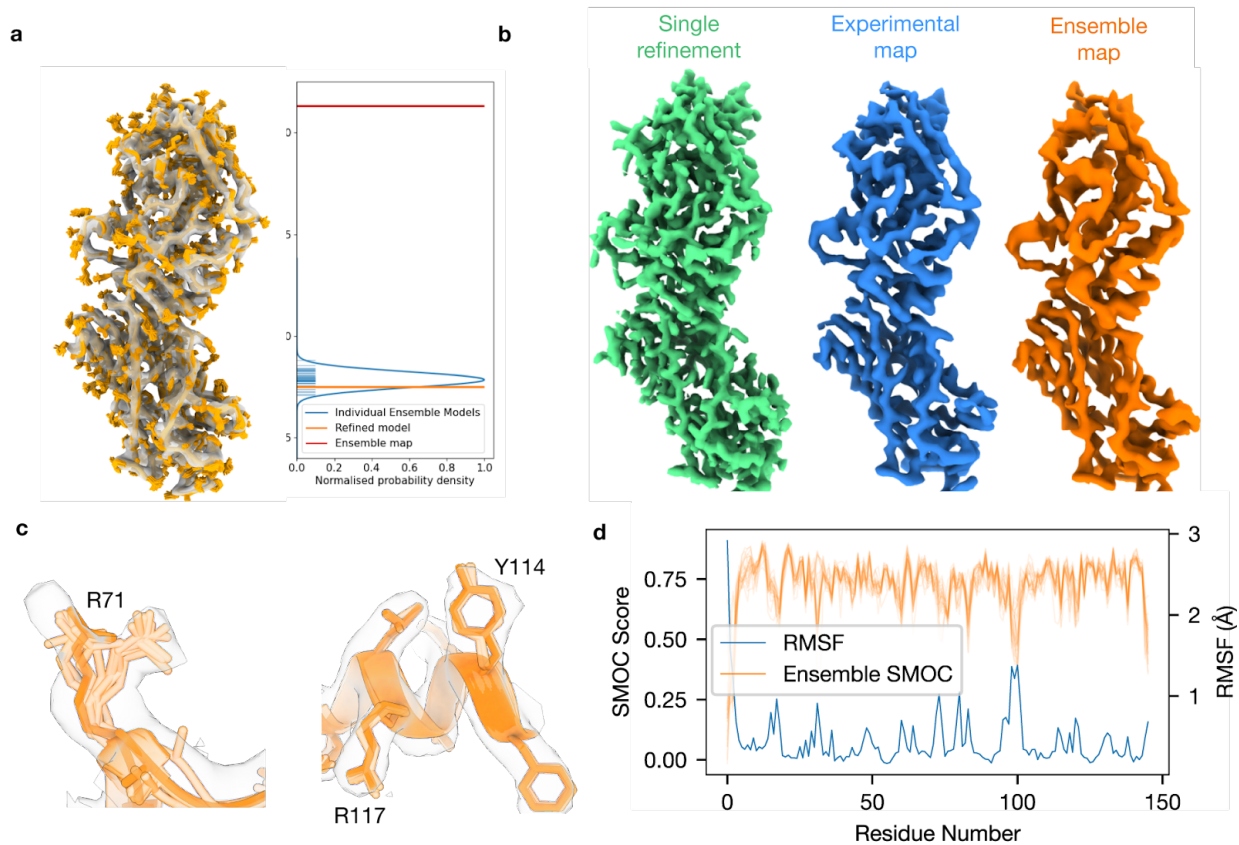
- The quality of the fit to the density often varies in different parts of the cryoEM map
- Proteins are dynamic
- Some of these parts may be better represented by multiple models

ENSEMBLE OF FITS



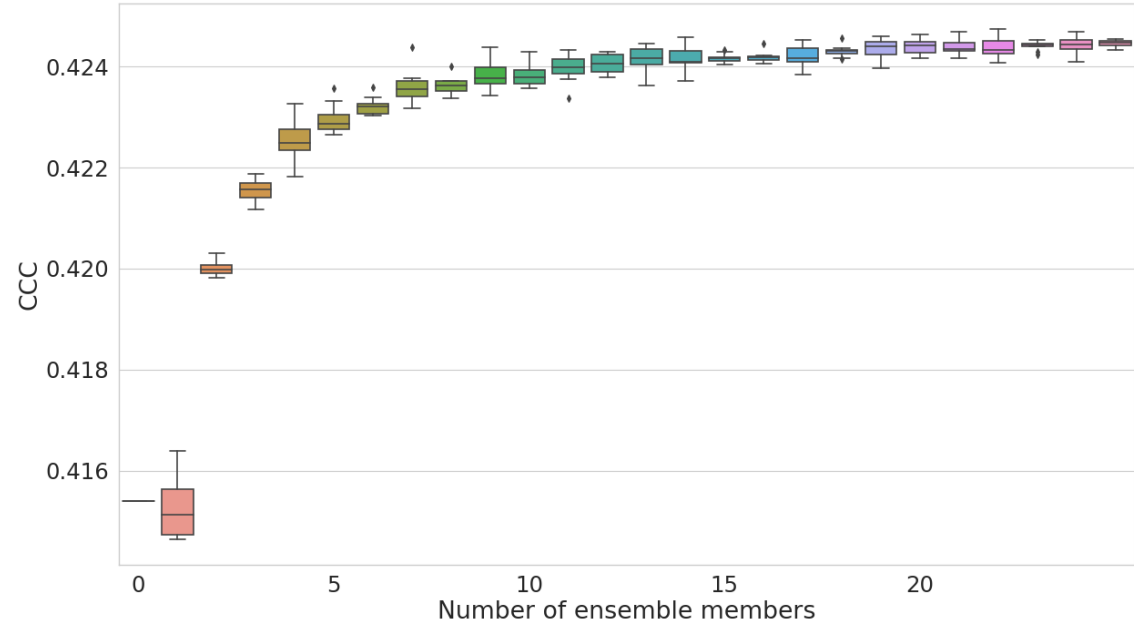
EMD: 1990 (8.9 Å resolution)
84 actin subunits from PDB

TEMPY-REFF CAN GENERATE ENSEMBLES



ENSEMBLE SIZE

CCC vs ensemble size for emd 9361



SOFTWARE: CCP-EM IMPLEMENTATION OF FLEX-EM AND TEMPY

CCP-EM | Flex-EM | 4ki8_rigid40_run1

Run Clone Coot Chimera CCP4mg Terminal Output Info Kill

Setup

Job title: 4ki8_rigid40_run1

Input map: Select e/emd_6422_16_regions.mrc

Input PDB: Select /groel/4ki8_model/4ki8_fit.pdb

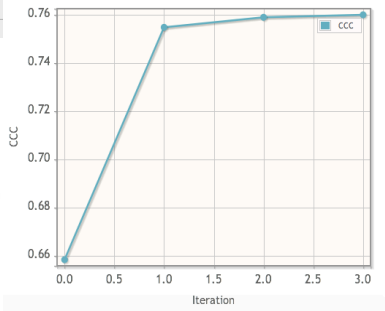
Auto rigid body:

Rigid body: Select model/4ki8_fit_denclust_40.txt

Resolution: 4.10

Iterations: 1

Extended options



CCP-EM | TEMPY-LocScore

Run Clone Coot Chimera CCP4mg PyMOL Terminal Output Info Kill

Setup

Job title: None

Input map: Select s/CCP4_meeting_israel/flexem_data/emd_6422_16_regions.mrc

Map resolution: 4.20

Metric Selection

SMOC score:

SCCC score:

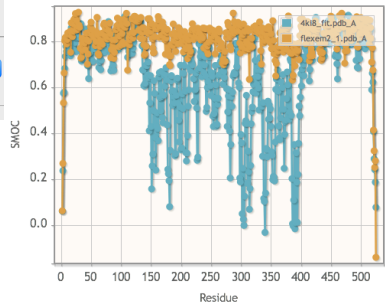
Input PDB(s):

PDB: Select s/CCP4_meeting_israel/flexem_data/4ki8_fit.pdb + -

Input rigid body: Select pobox/courses/CCP4_meeting_israel/flexem_data/flexem1_1.pdb

Advanced options

Residue neighborhood: Fragment



CCP-EM | TEMPYGlobScore | Untitled | TEMPYGlobScore_286

Run Clone Coot Chimera CCP4mg Terminal Output Info Kill

Setup Pipeline Launcher Results

Setup

Job title: Untitled

Input map: Select roel/emd_6422_16_regions.mrc

Map resolution: 4.10

Input PDB(s):

PDB: Select /roel/4ki8_fit.pdb + -

Add PDB

TEMPY scores

Atomic model	overlap_map	overlap_model	correlation	local_correlation	local_mi	ccc_ov	mi_ov
4ki8_fit.pdb	0.927	0.126	0.336	0.305	0.092	0.242	0.067
flexem1_1.pdb	0.964	0.132	0.535	0.273	0.076	0.273	0.076
flexem2_1.pdb	0.964	0.134	0.558	0.245	0.066	0.273	0.076

Global scores (CCC, MI...)

THANKS...

Topf group (current)

Sanjana Nair
Rebecca Brooker
Luca Genz
Thomas Mulvaney
Joseph Beton
Karen Manalastas
Aaron Sweeney
Laetitia Adeler-Ohde
Matthias Pfeifer
Birgit Märtens
Mauro Maiorca
Guendalina Marini
Manaz Kaleel
Aylin del Moral
Maryam Nikooei



previous members:

Sony Malhotra (STFC - CCP-EM)
Agnel Joseph (STFC - CCP-EM)
Arun Prasad Pandurangan (Cambridge)

CASP

Andriy Kryshtafovych

Birkbeck

Helen Saibil

STFC

Martyn Winn
Tom Burnley
Colin Palmer



Bundesministerium
für Bildung
und Forschung