

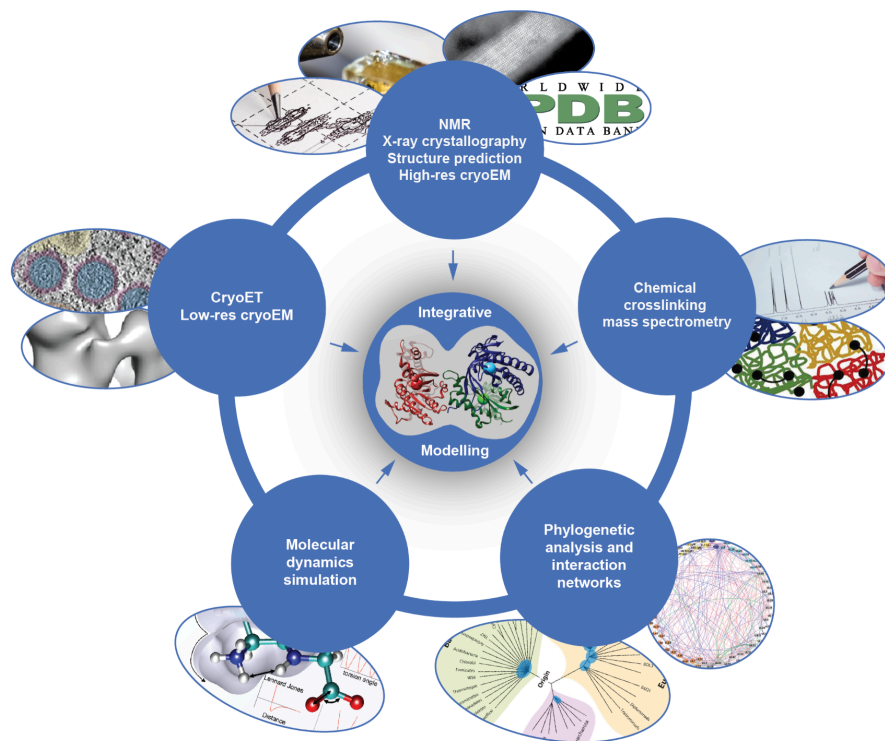
# TEMPy-ReFF + ChemEM

CCP-EM Icknield Model Building Workshop

**05.11.2024**

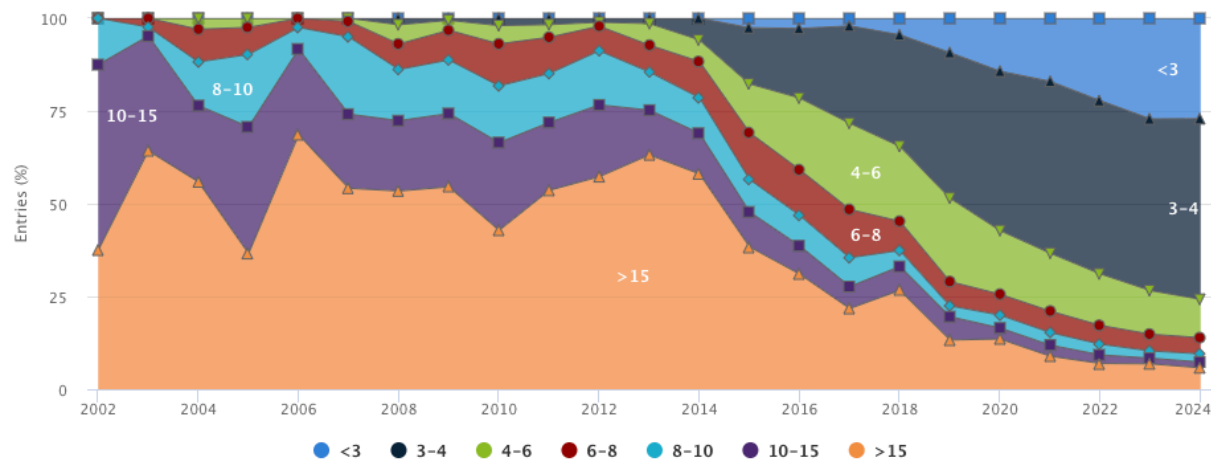
Maya Topf & Aaron Sweeney (CSSB Hamburg)





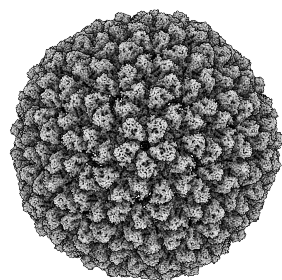
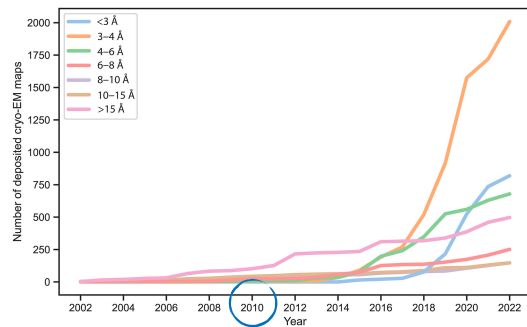
- **Develop** methods for integrative modelling of macromolecular assemblies
- **Apply** methods to model structures of viral assemblies (e.g. herpesviruses, arenaviruses)
- **Analyse** structural models to further understand their function
- **Distribute** and support new software for the structural biology community

# EMDB ENTRY RESOLUTION IN SHELLS PER YEAR

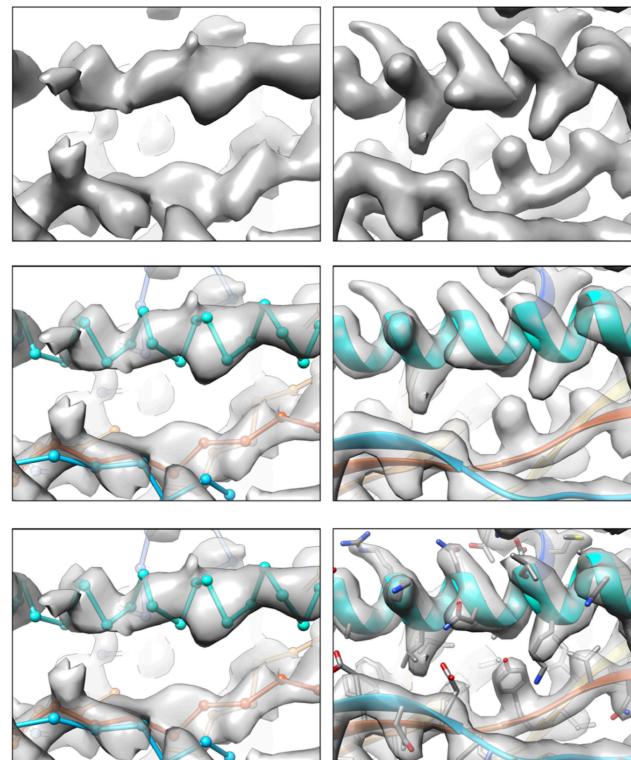


“The Resolution Revolution”  
W Kühlbrandt, *Science* 2014

# THE “RESOLUTION REVOLUTION”



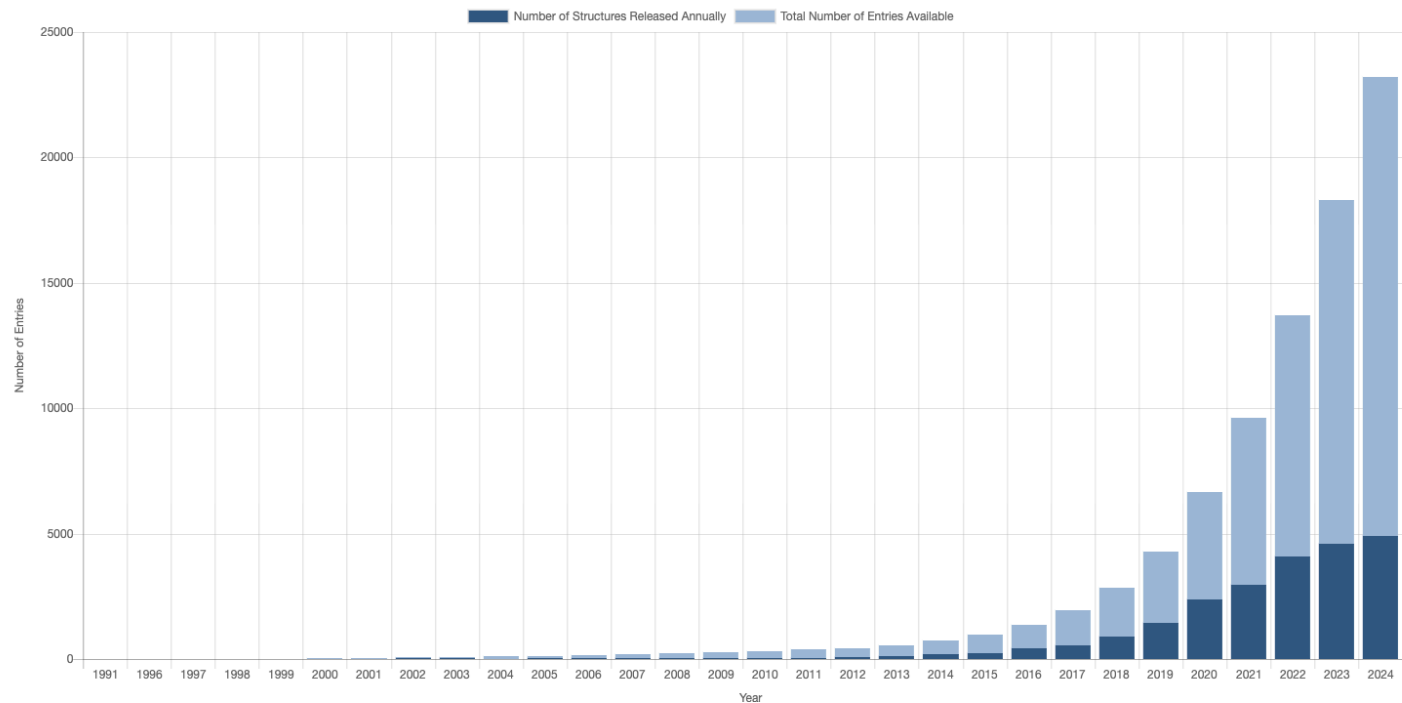
*Bacteriophage P22*  
4 Å (actually ~5 Å)  
EMDB 5137  
*Chen et al., 2011*



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

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

# NUMBER OF CRYOEM ENTRIES IN PDB



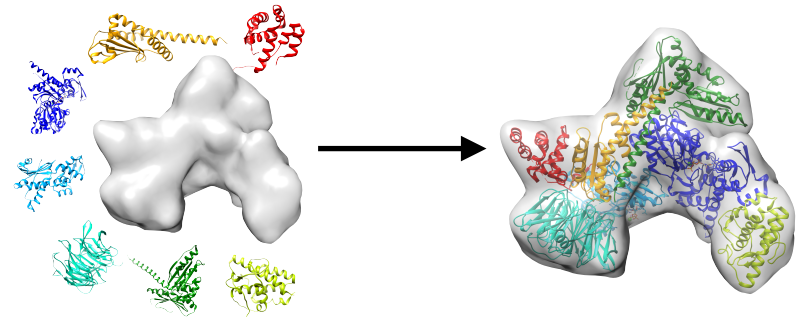
- ~10% of the entries in PDB are associated with maps in EMDB
- These days most entries are associated with either a single model or multiple models

## 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 small molecules and other macromolecules;
5. To help predict how new related molecules might look.
6. To simulated the dynamics of the model to better understand its function



# STRUCTURAL FEATURES AT DIFFERENT RESOLUTION LEVELS

Low

Intermediate

High

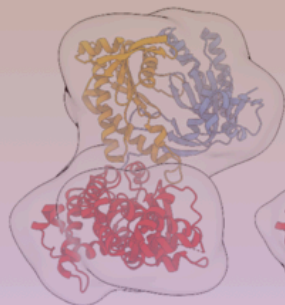
20 Å

15 Å

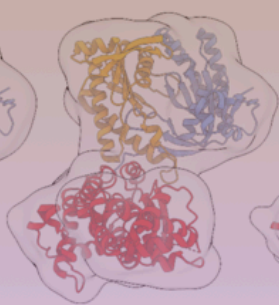
10 Å

5 Å

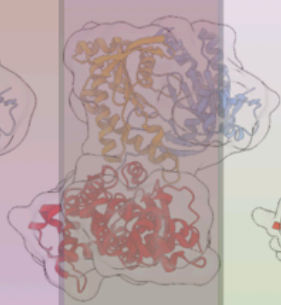
3 Å



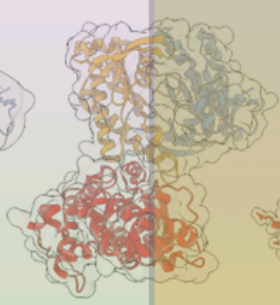
Subunit orientation  
and boundaries



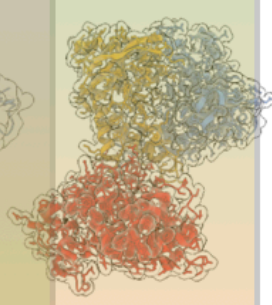
Structural elements



helices



Sheets  
Side chains



Atoms

Multi-component- fitting

Rigid fitting

Flexible refinement

*de novo* modelling

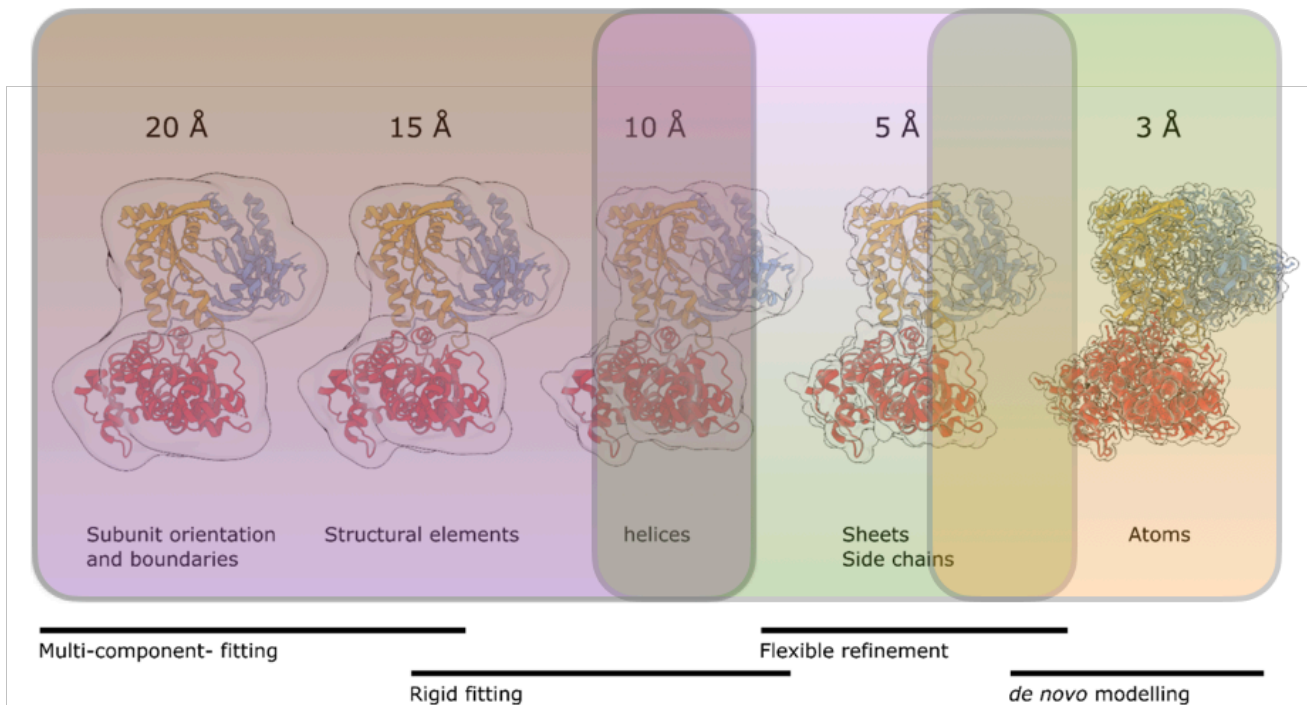
- If a complete experimental structure is available (and/or a reliable predicted model) 10 Å data may be sufficient.

- If no sequence/composition data is available even 3 Å may be challenging

# STRUCTURAL FEATURES AT DIFFERENT RESOLUTION LEVELS

TEMPy2/ $\gamma$ -TEMPy

Flex-EM/TEMPy-ReFF



- If a complete experimental structure is available (and/or a reliable predicted model), 10 Å data may be sufficient.

- If no sequence/composition data is available even 3 Å may be challenging

## TEMPy2

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

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

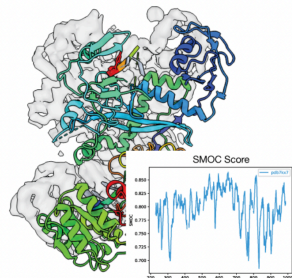
### Get TEMPy2

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

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

## TEMPy can do...

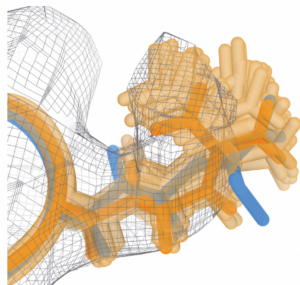
### Validation



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

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

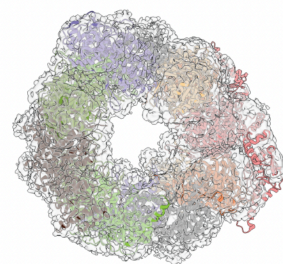
### Flexible Fitting and refinement



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 y-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 (y-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

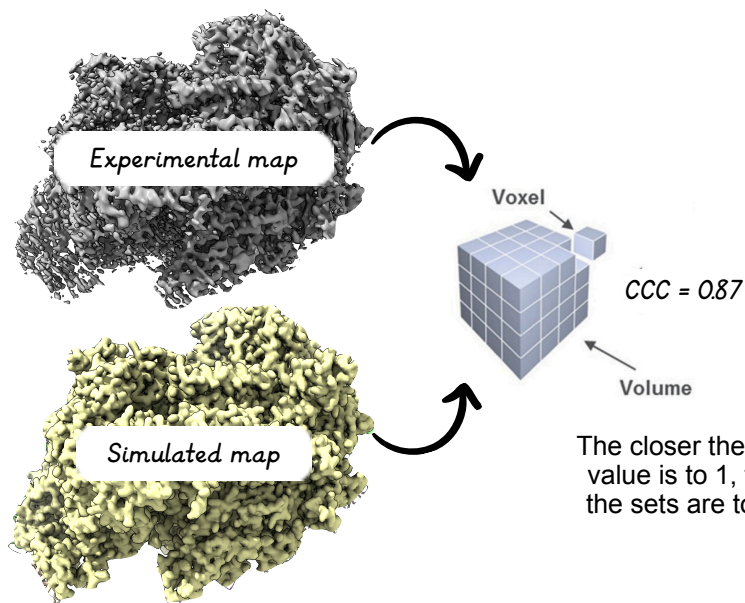
### Map comparison

- Difference maps
- FSC

### Chimera scripts

Plots  
Attribute files  
PDB files

# FITTING AN INITIAL STRUCTURE



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

- \* When possible, a known or pre-calculated model is placed and fitted in the cryo-EM map as a rigid body.
- \* Often an exhaustive six dimensional (6D) grid search of the 3 translational and 3 rotational degrees of freedom is performed to locate the highest global cross-correlation.

Cross Correlation Coefficient

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

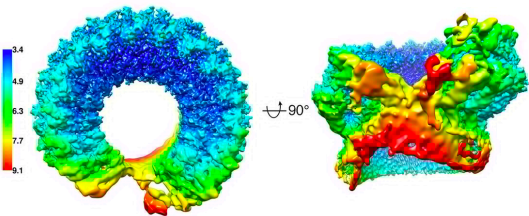
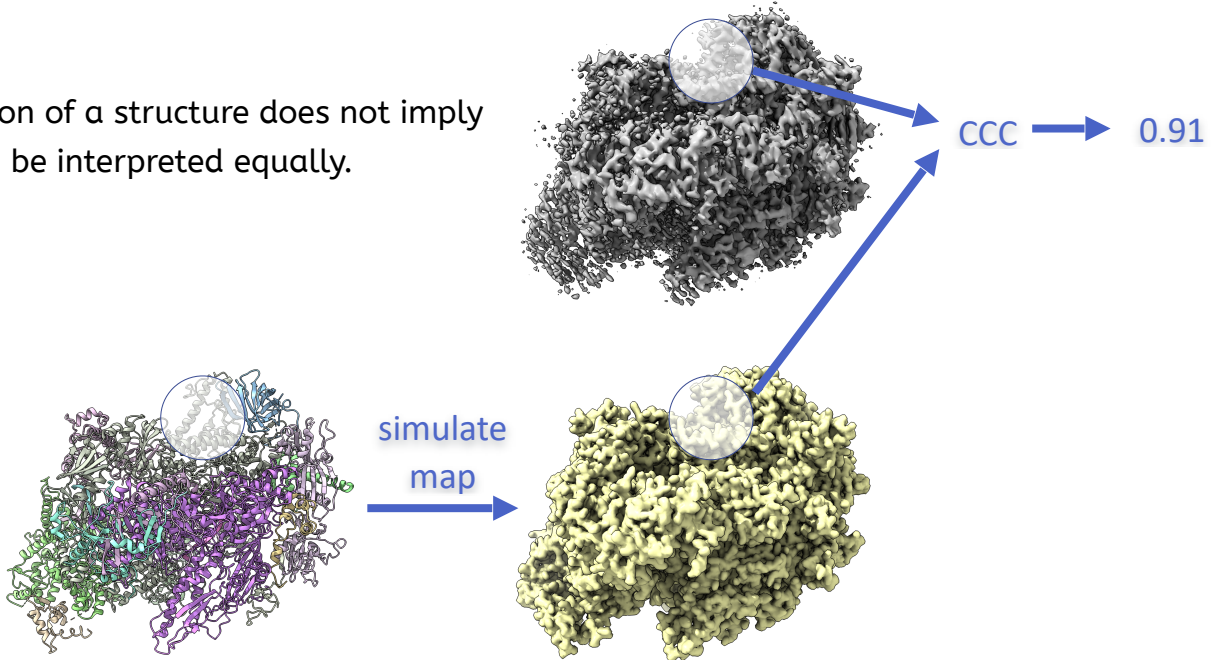
Mutual information score

$$I(X;Y) = \sum_{x \in X} \sum_{y \in Y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

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

# LOCAL RESOLUTION AND LOCAL CORRELATION

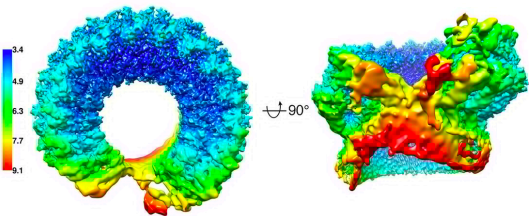
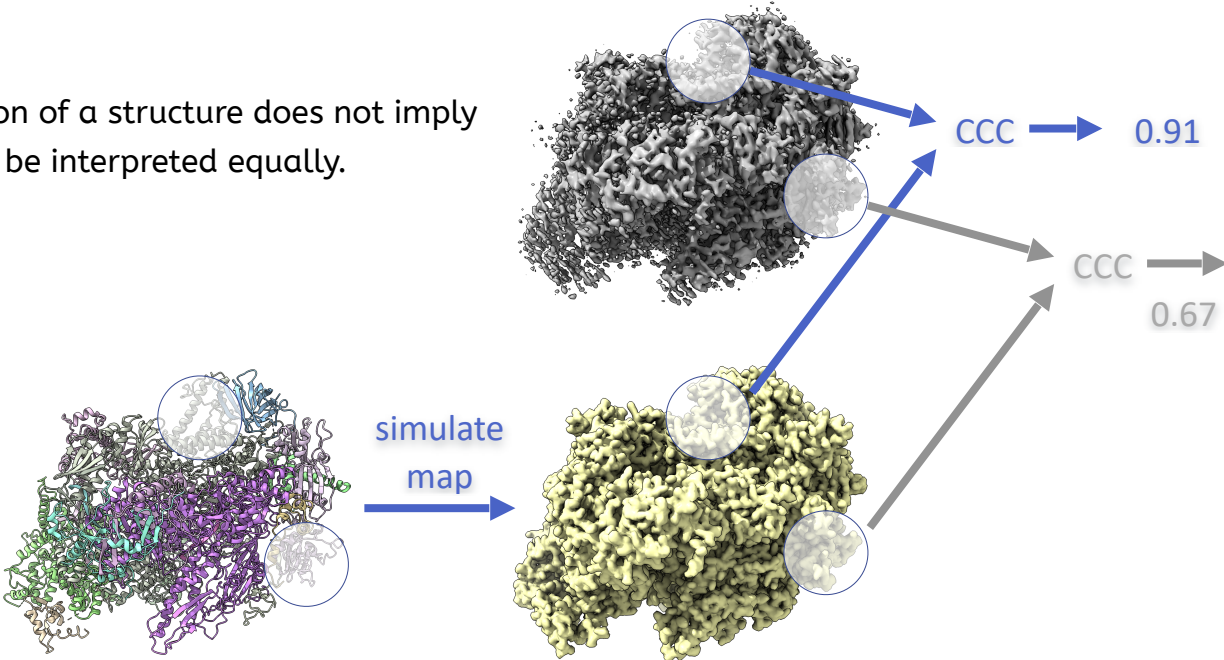
\* The overall resolution of a structure does not imply that all regions can be interpreted equally.



\* Local resolution maps are useful for estimating resolution variability and map quality in general.

# LOCAL RESOLUTION AND LOCAL CORRELATION

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

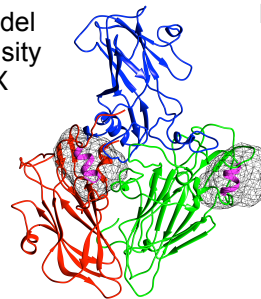


\* Local resolution maps are useful for estimating resolution variability and map quality in general.

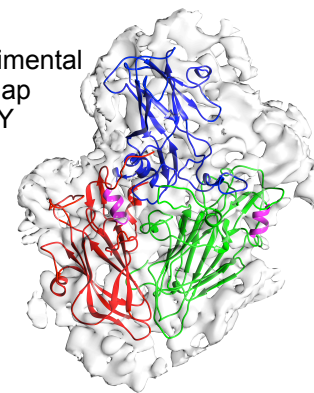
**SCCC:** **Segment**-based cross-correlation coefficient:

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

Model density  
X

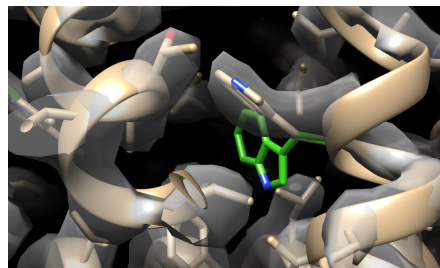


Experimental map  
Y



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

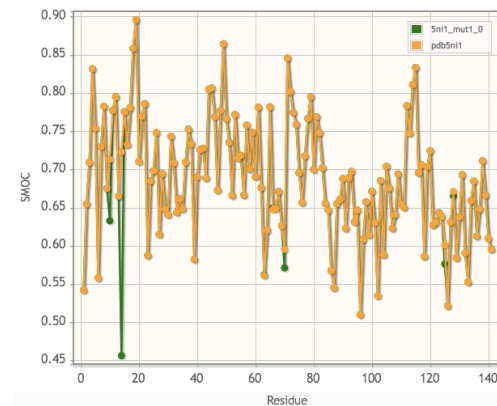
$$SMOC = \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n X_i^2} \sqrt{\sum_{i=1}^n Y_i^2}}$$



EMD-3488 (3.2Å)  
Deposited model PDB: 5N11



Agnel Joseph

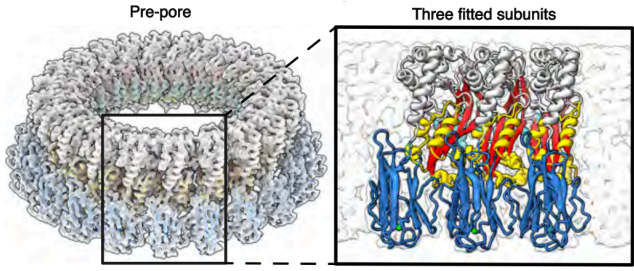


# ENTER THE "AI" ERA: GOOD INITIAL PREDICTIONS BUT...

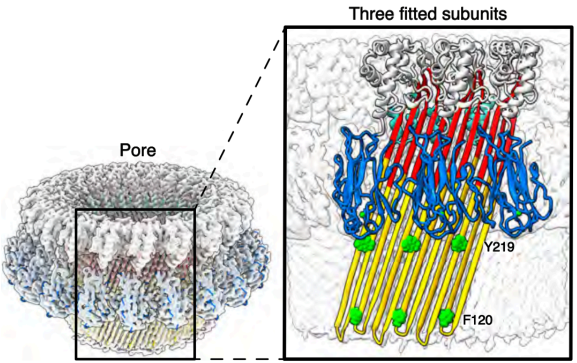


Gwen Marini

## Mpf2Ba1 pre-pore

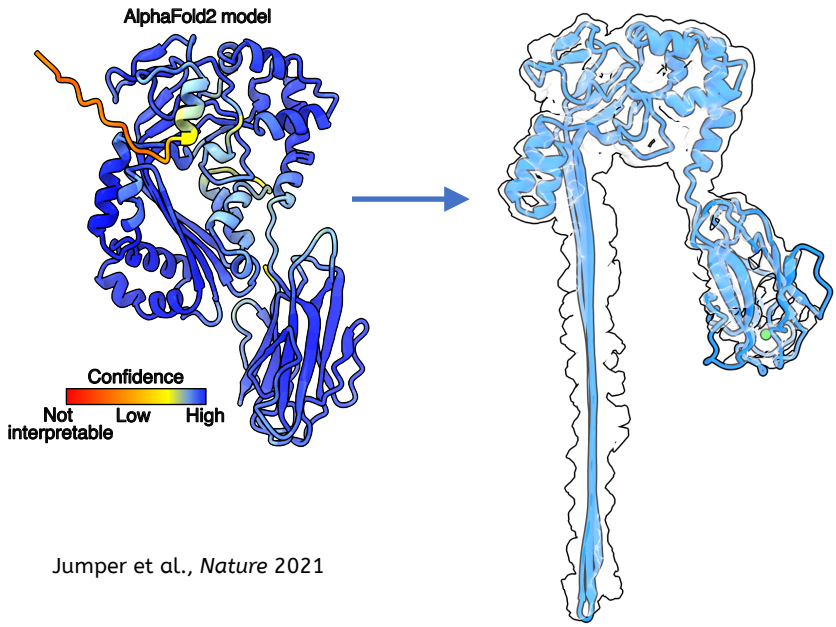


## Mpf2Ba1 pore



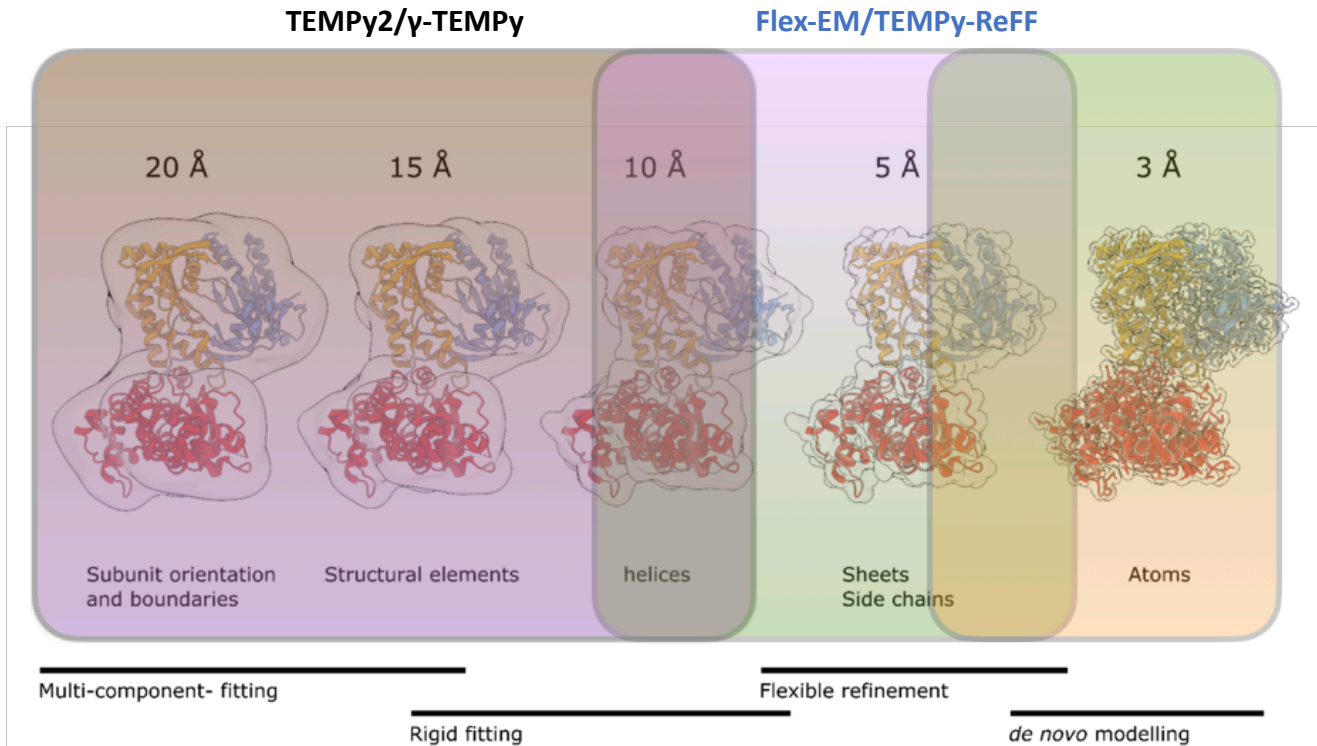
2-3 Å resolution

Marini et al., *Nat Comms* 2023



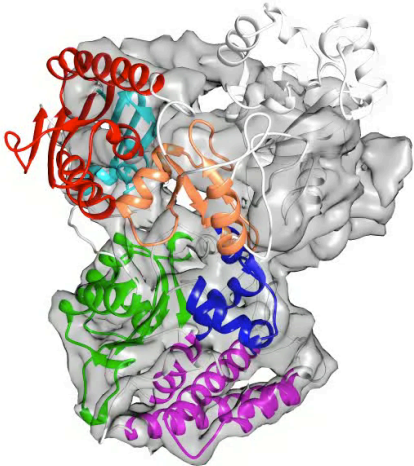
Jumper et al., *Nature* 2021

# STRUCTURAL FEATURES AT DIFFERENT RESOLUTION LEVELS



- What about higher resolutions?

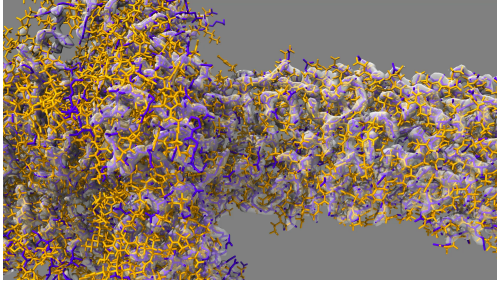
# FLEXIBLE FITTING AND REFINEMENT AT LOW TO HIGH RESOLUTIONS



sub-domains

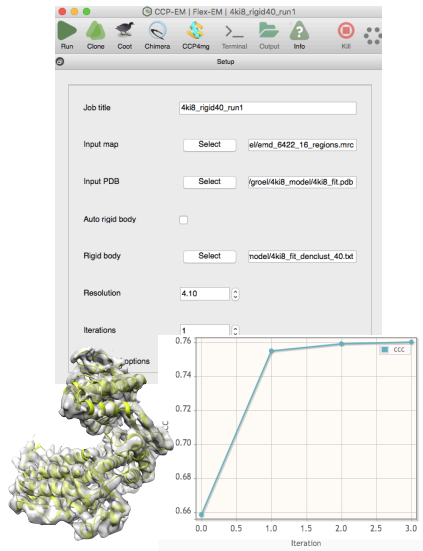
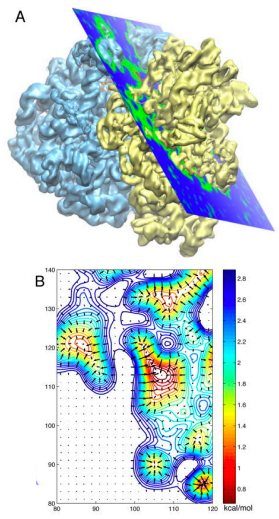


secondary structure elements



residue / all-atom

# FLEXIBLE FITTING AND REFINEMENT



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

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

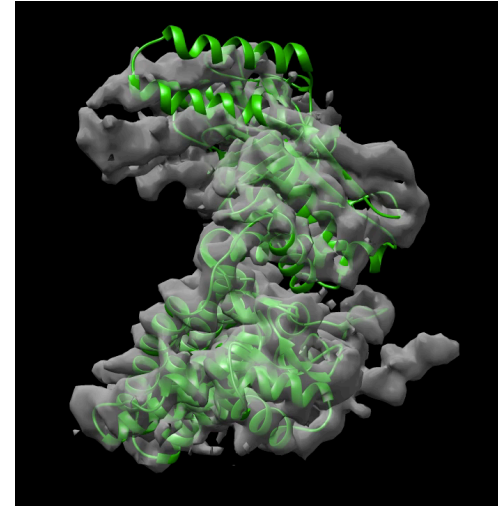
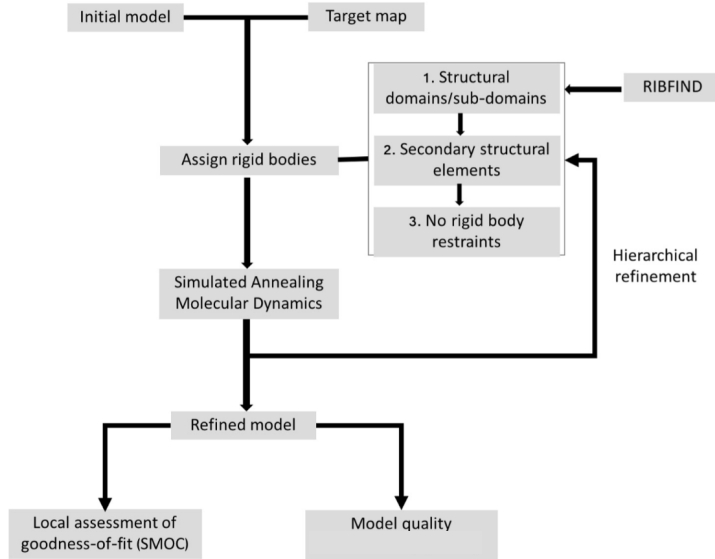
Trabuco et al., *Methods* 2008

Topf et al., *Structure* 2008;  
Joseph et al., *Methods* 2016

# FLEXIBLE FITTING AT MEDIUM RESOLUTION BY HIERARCHICAL REFINEMENT



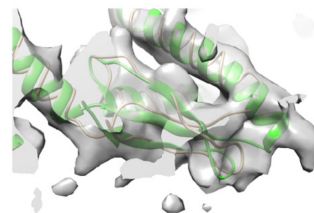
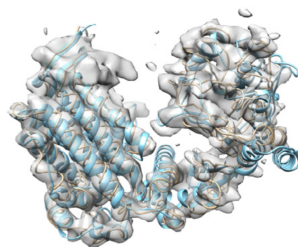
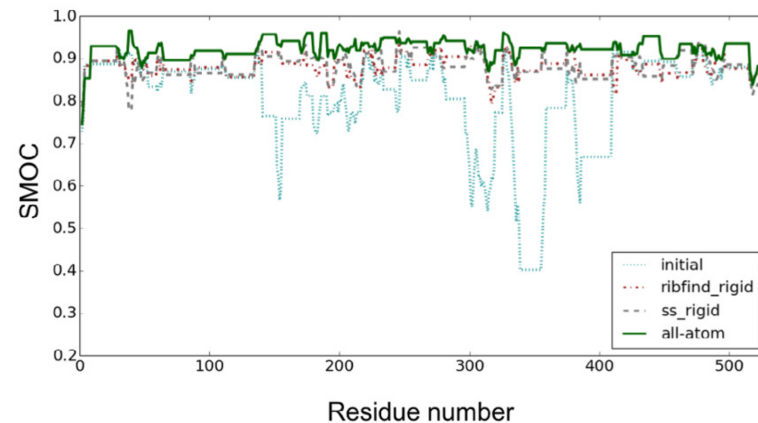
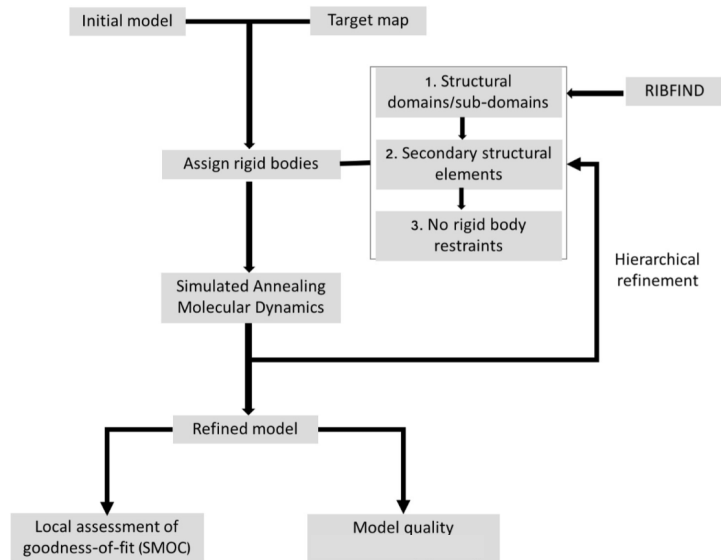
Agnel Joseph



Unliganded GroEL at 4.2 Å resolution (EMD-5001)

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

# 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)

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

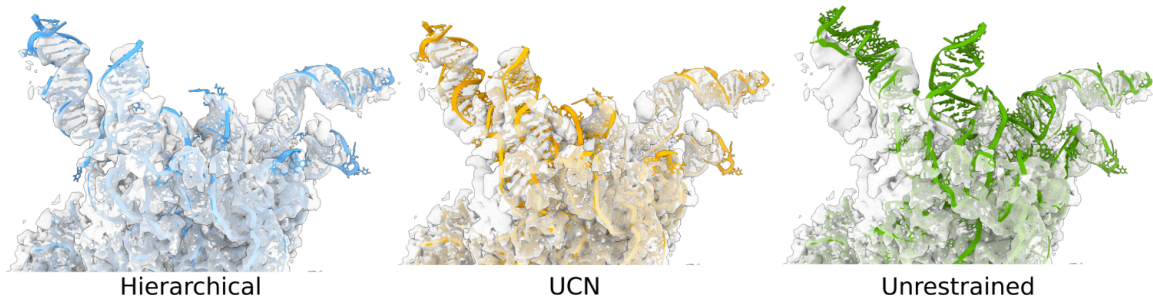
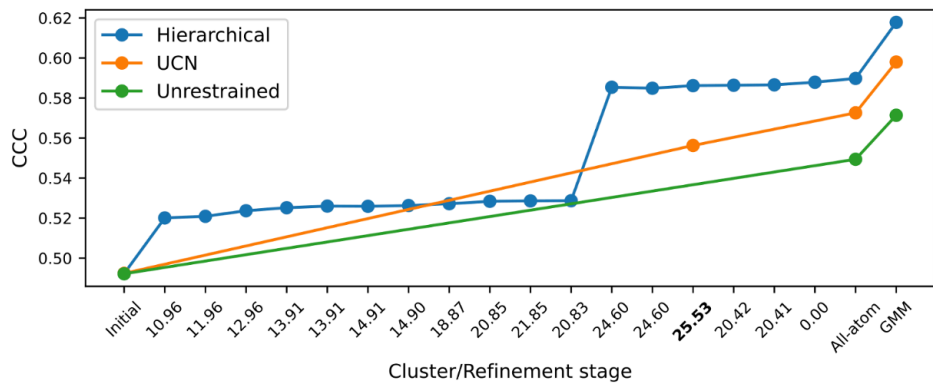
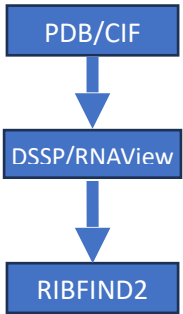


Thomas Mulvaney

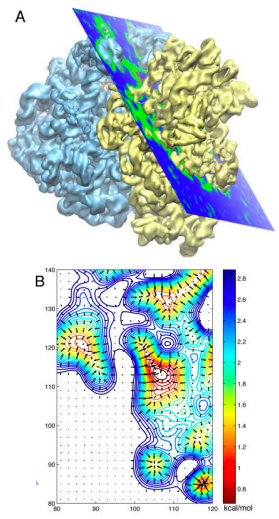


Sony Malhotra

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

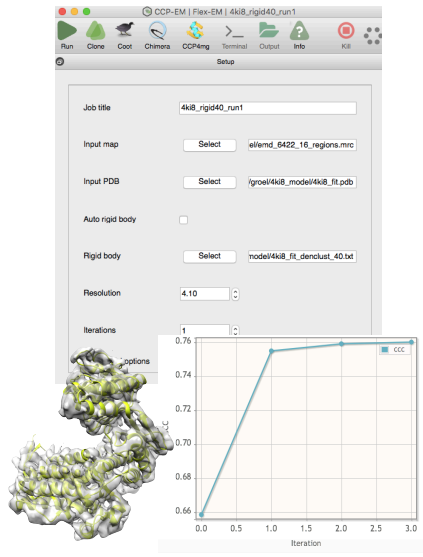


# FLEXIBLE FITTING AND REFINEMENT



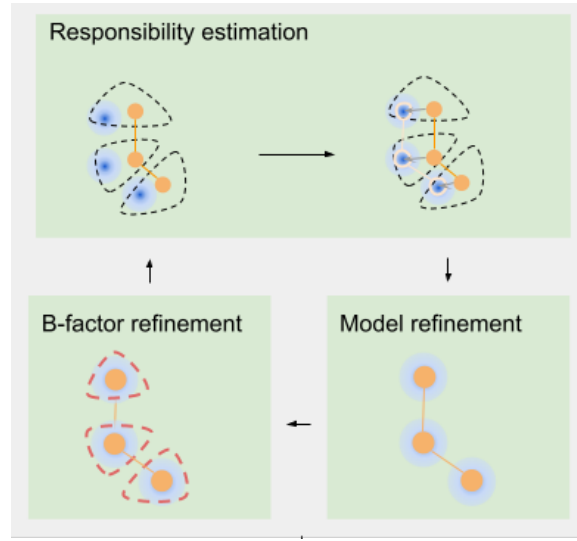
Molecule Dynamics Flexible Fitting (MDFF) uses a standard force field with a potential based on the density

Trabuco et al., *Methods* 2008



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

Topf et al., *Structure* 2008;  
Joseph et al., *Methods* 2016



TEMPy-ReFF (Responsibility-based Flexible-Fitting)

Cragolini et al., *Proteins* 2021;  
Beton\*, Mulvaney\* et al, *Nat Comm* 2024



Thomas  
Mulvaney



Joseph  
Betoni



Tristan  
Cragolini

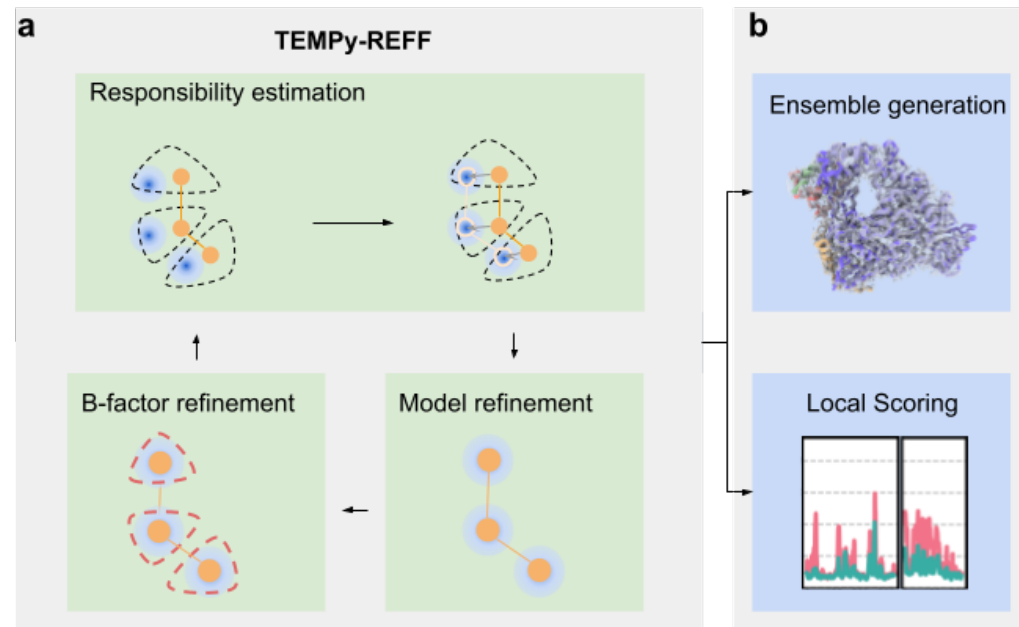
- Use an expectation-maximisation approach
- Use one gaussian per atom, and a background noise term
- Improves atomic position and variance (B-factor)

## TEMPy and OpenMM

Amber14 forcefield

Generalised Born implicit solvent model

Integration with Langevin dynamics, at 100K

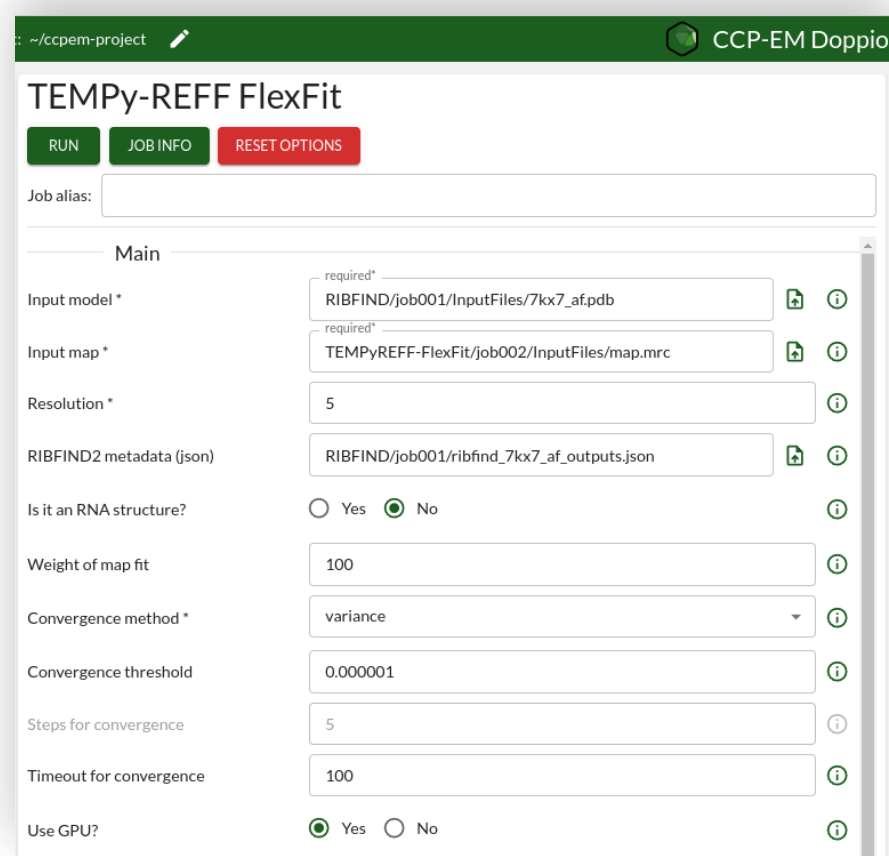


- **Flexible fitting**

- Flexible fitting of proteins intermediate resolutions using density guided force + AMBER.
- Faster – built on top of OpenMM, can handle much larger models
- Hierarchical rigid body protocol using **RIBFIND2**

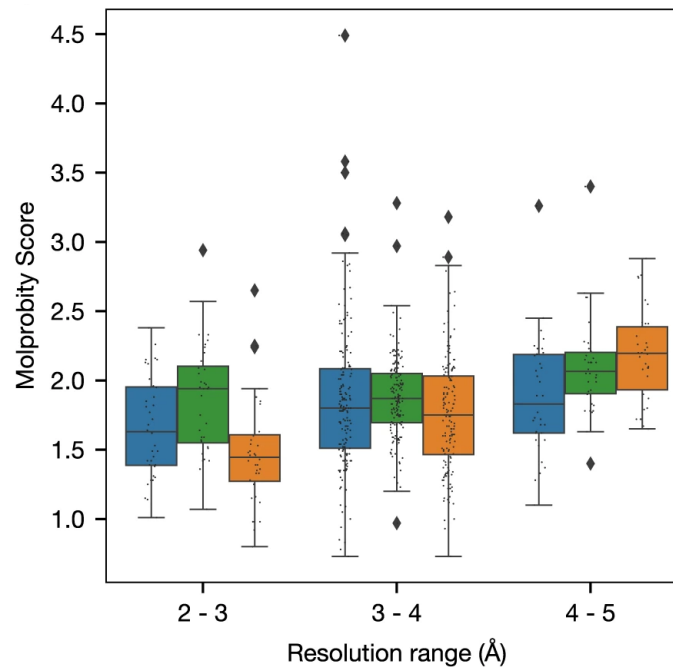
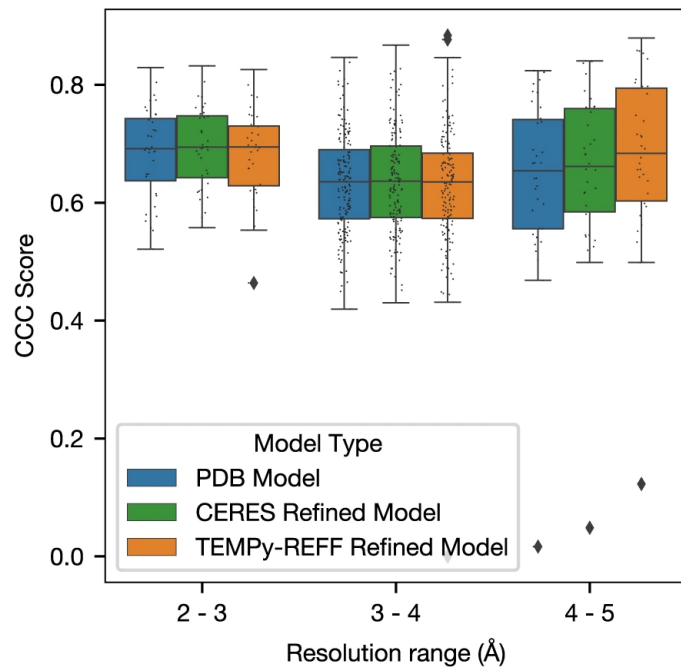
- **Refinement**

- Optimises fit of GMM to data using molecular dynamics with OpenMM and AMBER.

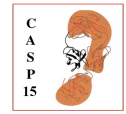
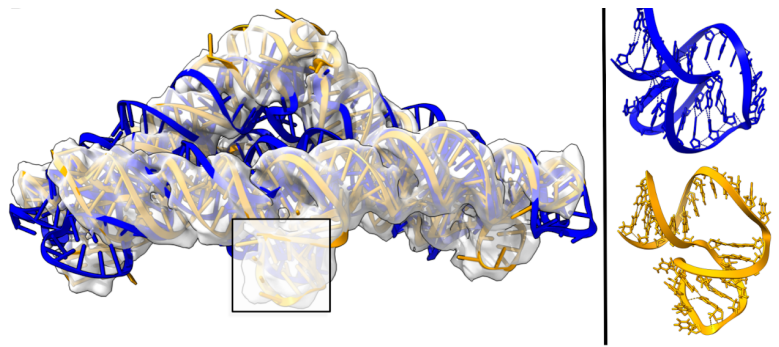
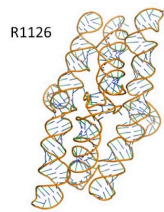
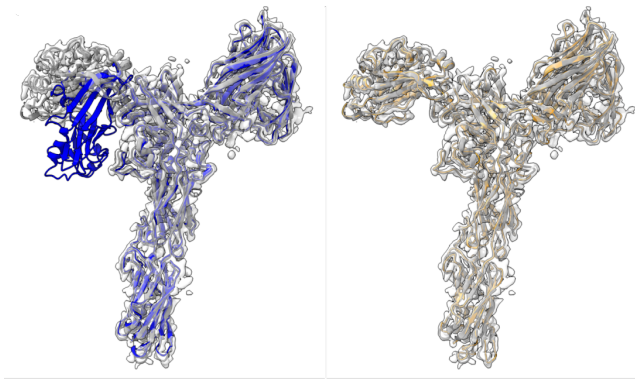
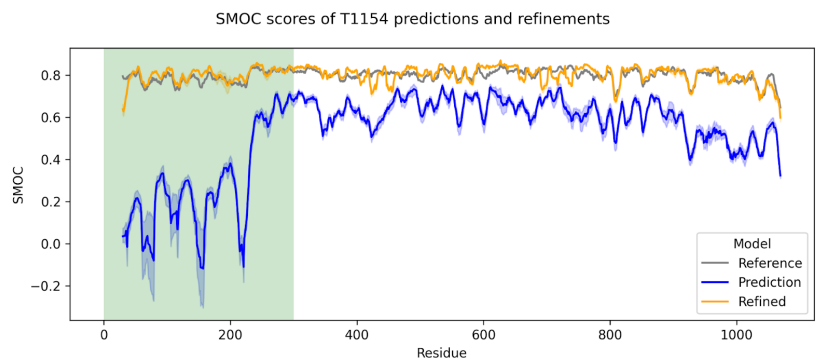
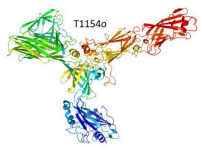


Fancy GUI in CCP-EM Doppio

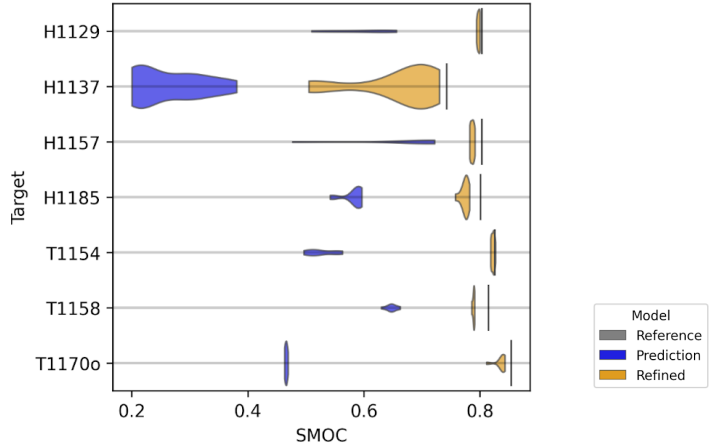
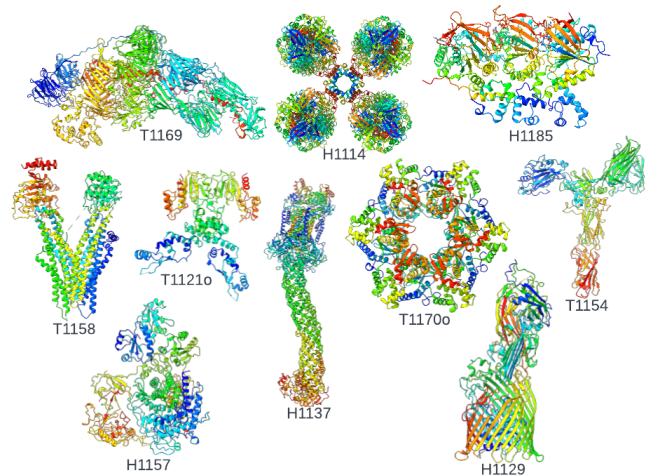
# COMPARISON AGAINST 229 CERES MODELS



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



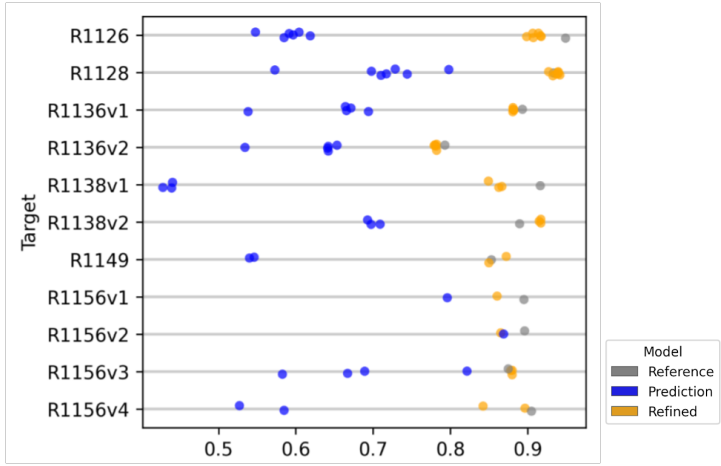
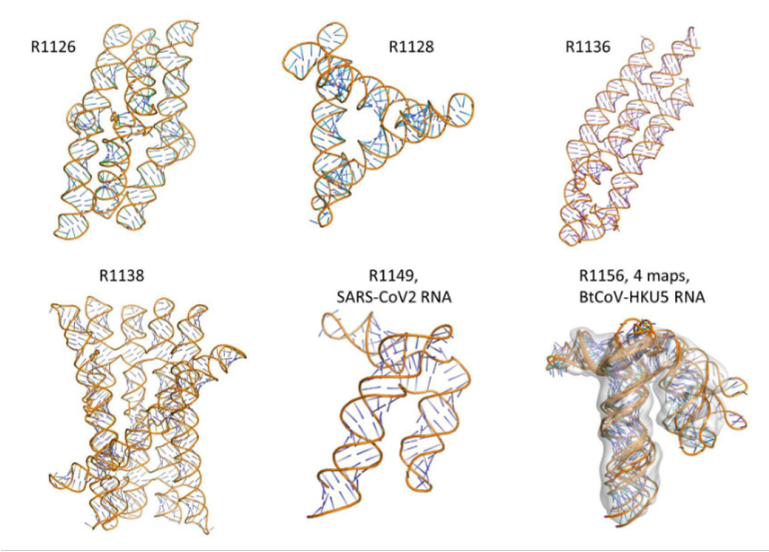
# TEMPY-REFF REFINEMENT OF CASP15 MODELS (PROTEIN COMPLEXES AND RNA)



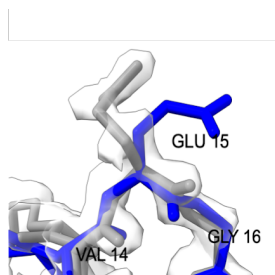
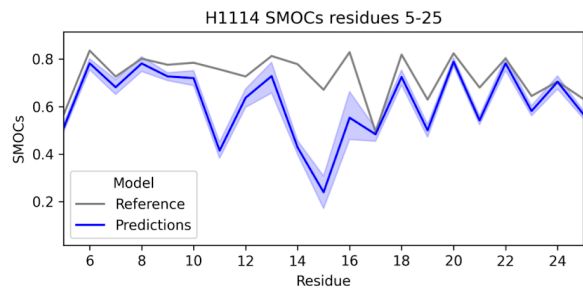
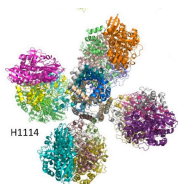
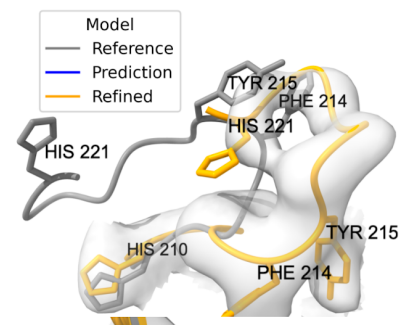
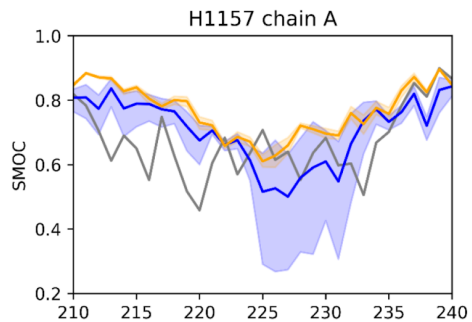
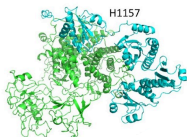
# TEMPY-REFF REFINEMENT OF CASP15 MODELS (PROTEIN COMPLEXES AND RNA)



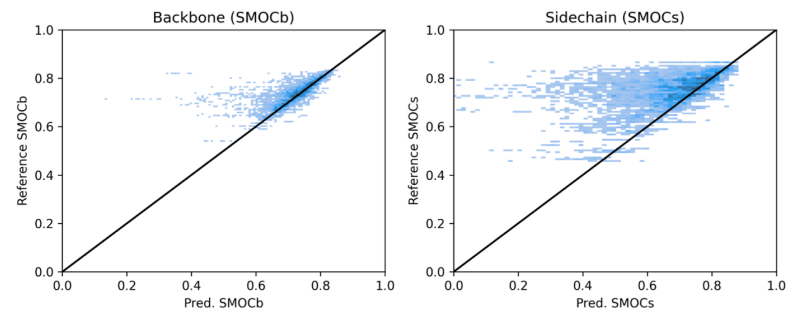
Rachael Kretsch,  
Das Lab, Stanford



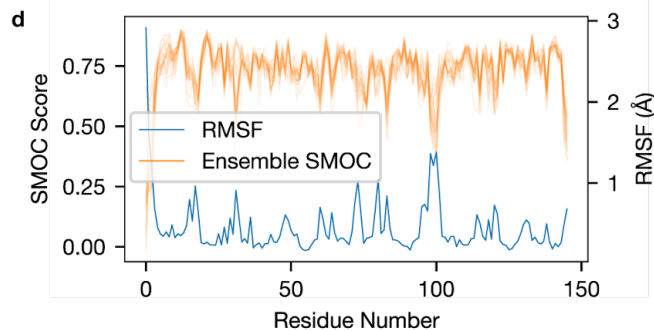
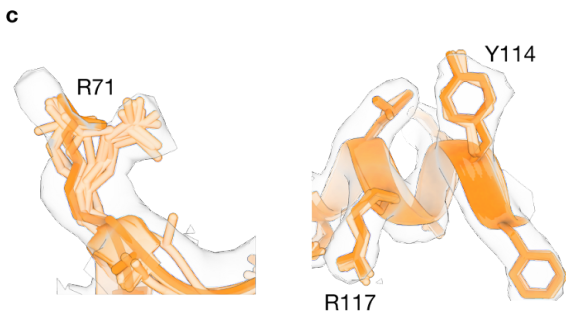
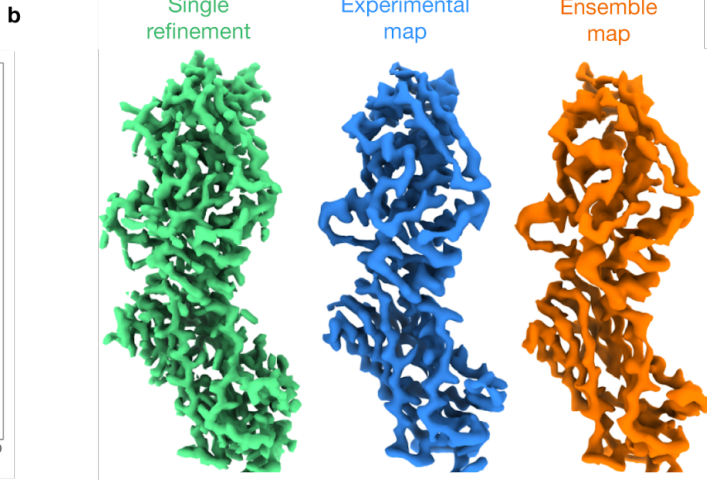
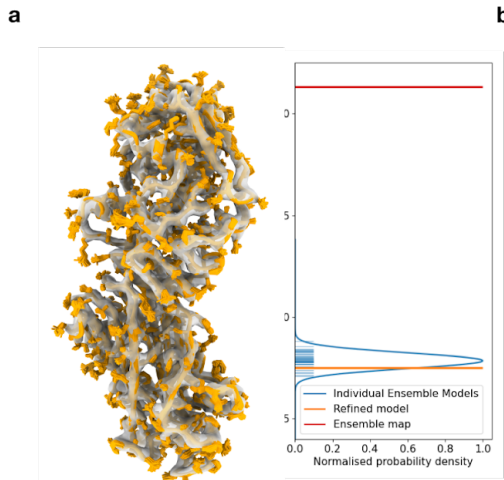
# EXPERIMENTAL VALIDATION AND REFINEMENT BY CRYO-EM (CASP15)



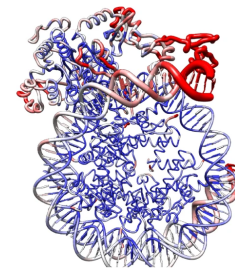
SMOc score distributions for H1114 chain A



# TEMPY-REFF CAN GENERATE ENSEMBLES

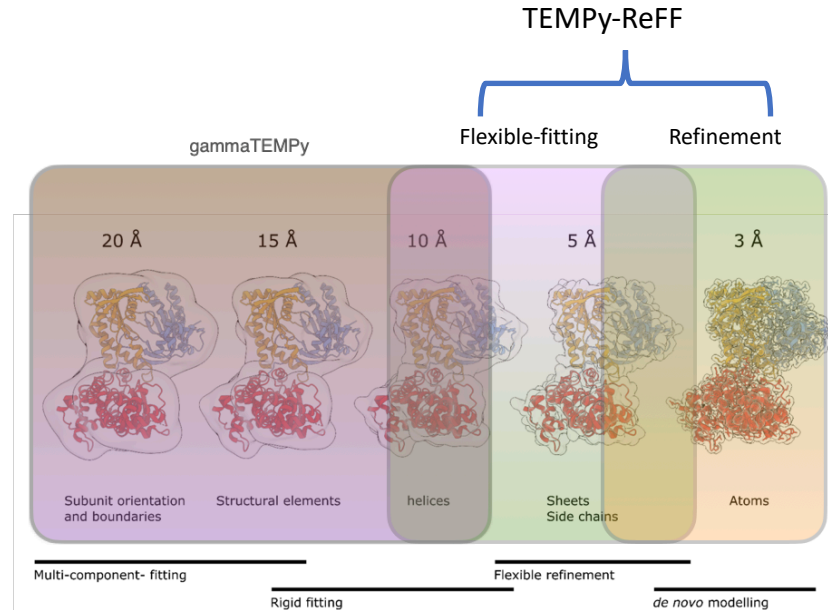


Color: TEMPy-ReFF Bfactor  
Width: ResMap



## SUMMARY (TEMPY-REFF)

- 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
- **TEMPy-ReFF** is useful across a range of resolutions.
- **RIBFIND2** can generate rigid bodies for Protein and RNA.
- **TEMPy-ReFF** ensembles can model uncertainty/heterogeneity.
- **CASP15** predictions were good starting points for models for flexible fitting and refinement.



- **More attention is needed for small molecules...**

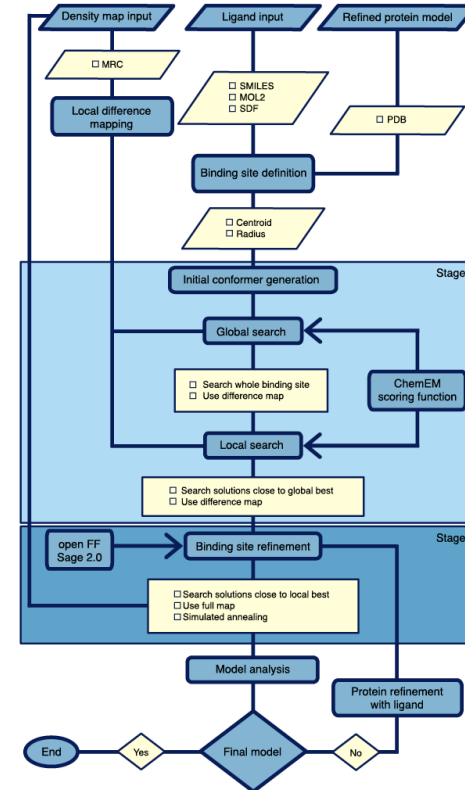
# ChemEM: flexible docking of small molecules in Cryo-EM structures

Uses Molecular docking, Molecular Dynamics and Density Difference mapping.

Fits ligands and refines binding site residues.

Accurate ligand fitting at 2.2-5.6 Å resolution.

Can fit multiple ligands in parallel to a single site.



## ChemEM: Ligand Placement and Refinement Workflow

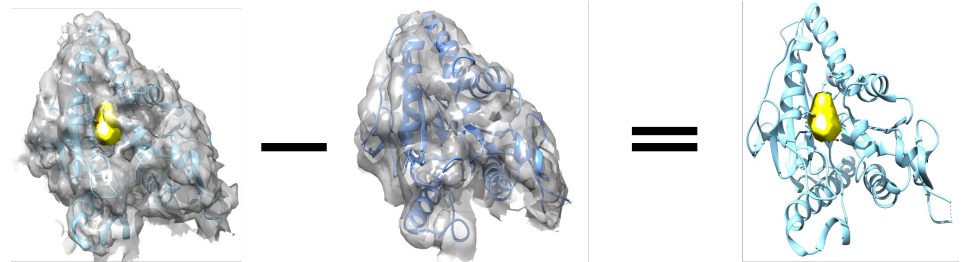
### • Stage 1: Initial Ligand Docking

- o Segment ligand density from the electron density map.
- o Perform molecular docking using ligand density as a 'restraint'.
- o Score solutions by combining:
  - **Docking Score** for physicochemical reasonableness.
  - **Mutual Information Score** for alignment with density.

**ChemDock Score:** Includes terms for Bonding, Pi-Pi stacking and hydrophobic interactions

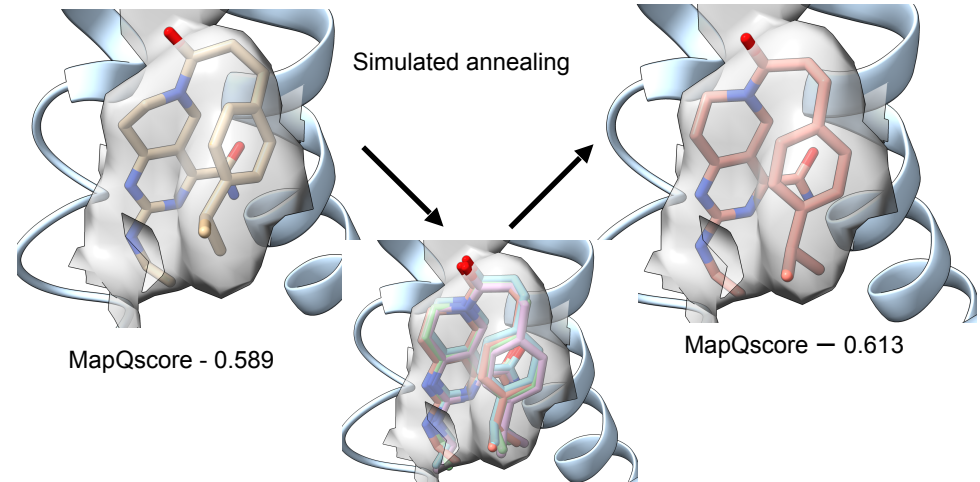
$$\text{Docking Score} = \text{Bias} + w1 * \text{HB}_{dist} + w2 * \text{HB}_{angle} + w3 * \text{Vina}_{g1} + w4 * \text{Vina}_{g2} + w5 * \text{Vina}_{steric} + w6 * \text{LogP}_{HP1} + w7 * \text{Aromatic}_{dist} + w8 * \text{Aromatic}_{angle} + w9 * \text{Intra}_{vg1} + w10 * \text{Intra}_{vg2} + w11 * \text{Intra}_{vsteric}$$

**Segmenting Ligand density:** Currently using a density difference map method.

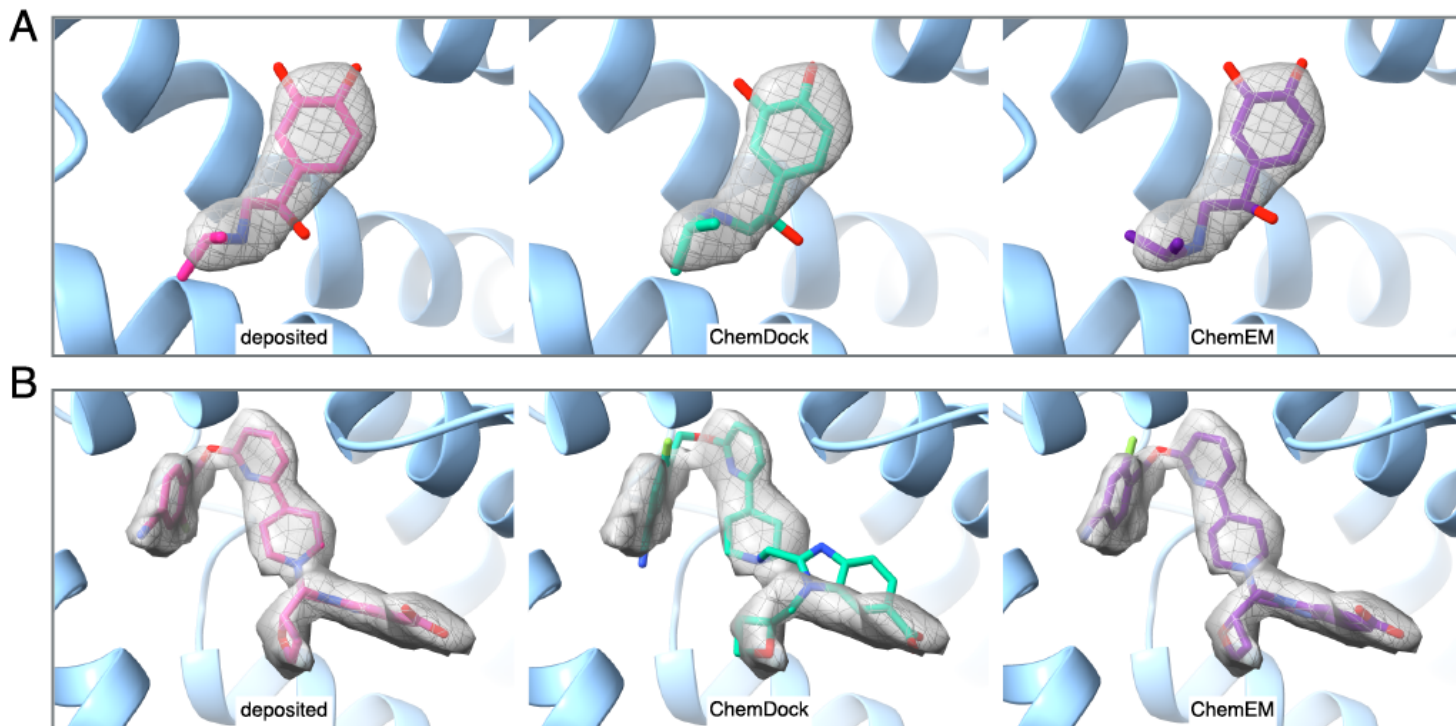


## Stage 2: Molecular Dynamics Refinement

- Select top solutions from Stage 1.
- Refine ligand placement using molecular dynamics:
  - **Amber Parameters** for protein atoms
  - **OpenFF Parameters** for ligand molecules.



# Using ChemEM : Examples



# Using ChemEM : Examples

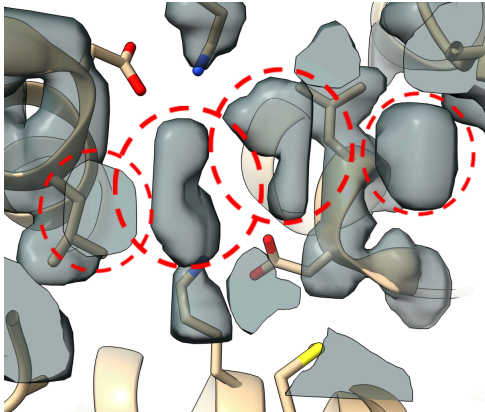
- ChemEM can be downloaded from
- It runs in terminal and all data is passed in from a configuration file.
- The enable stages block defines what protocols to run
- The rest defines the data and output locations
- ChemEM is also capable of fitting multiple ligand to the same site.

```
#ChemEM config file

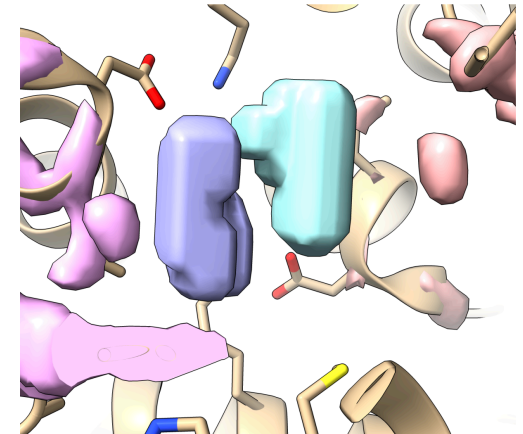
#enable stages
pre_process = 1
pre_process_split_density = 1
auto_split_point = 1
auto_split_zone = 0
fitting = 1
dock_only = 0
post_process = 0
rescore = 0

#data and working directory
protein = ~/test_data/6tti_protein.pdb
ligand = CC(=O)Nc1cc(no1)C(F)(F)F
ligand = CC(=O)Nc1cc(no1)C(F)(F)F
ligand = CS(=O)C
ligand = CS(=O)C
densmap = ~/maps/6tti.mrc
resolution = 2.5
map_contour = 0.07
centroid = ( 43.908, 54.415, 49.755)
output = ~/multi_fit
```

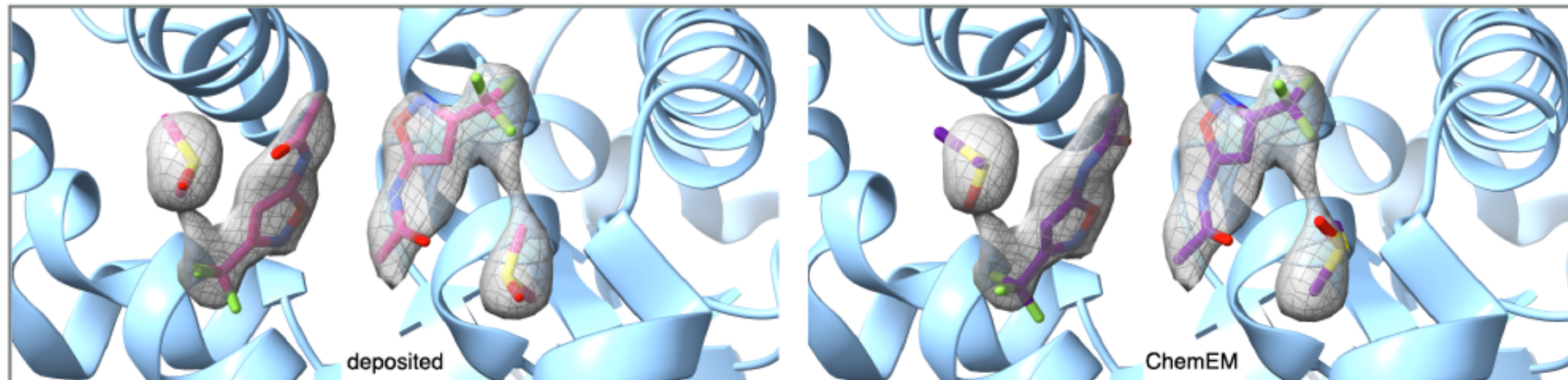
- The density is segmented by the difference mapping technique. This yields a single density map (That has a lot of extra density)
- The results can be improved and calculations sped up by specifying multiple centroid using this initial map. And assigning a ligand to each 'density'.



```
...  
centroid = ( 43.908, 54.415, 49.755)  
centroid = ( 39.049, 58.586, 49.760)  
centroid = ( 40.625, 62.434, 52.641)  
centroid = ( 42.391, 50.558, 52.642)  
...
```



# Using ChemEM : Multi-ligand fitting



- **Recent improvement to ChemEM include:**

- Ions can be included to the docking calculation
- Previously known ligand positions can be included
- Protonation state assignment
- Improved post-processing algorithms enable selective refinement of specific ligand atoms.

- **Coming soon:**

- Better handling of water molecules and binding site solvation.
- Lead optimization tools to predict the impact of R-group modifications on ligand affinity.

**Topf group (current)**

Sanjana Nair  
Thomas Mulvaney  
Karen Manalastas  
Luca Genz  
Aaron Sweeney  
Laetitia Adeler-Ohde  
Matthias Pfeifer  
Birgit Märtens  
Mauro Maiorca  
Guendalina Marini  
Annika Rammelt  
Pasquale Lamagna  
Natan Nagar  
Fabian Hausmann  
Gabriele Diana



**Previous members:**

Sony Malhotra (STFC)  
Joseph Beton (Biontech)  
Agnel Joseph (STFC)  
Tristan Cragolini (Birkbeck)

**CCP-EM (STFC, UK)**

Tom Burnley  
Agnel Joseph  
Sony Malhotra  
Martyn Winn

**DAS group (Stanford, CA)**

Rhiju Das  
Rachael Kretsch

**CASP (UC Davis, CA)**

Andriy Kryshtafovych

**Birkbeck (UK)**

Helen Saibil

