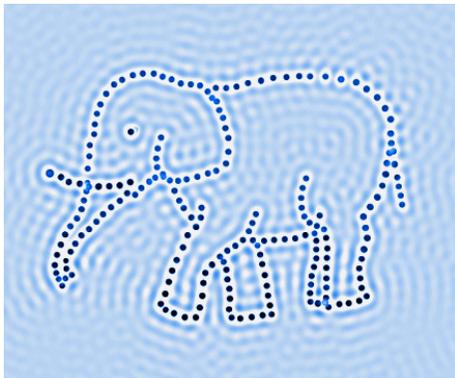


# Lecture 15

## Maya Topf

# Fitting of structures, scoring and assessment



EMBO course on image processing  
for cryo EM

10 September 2019



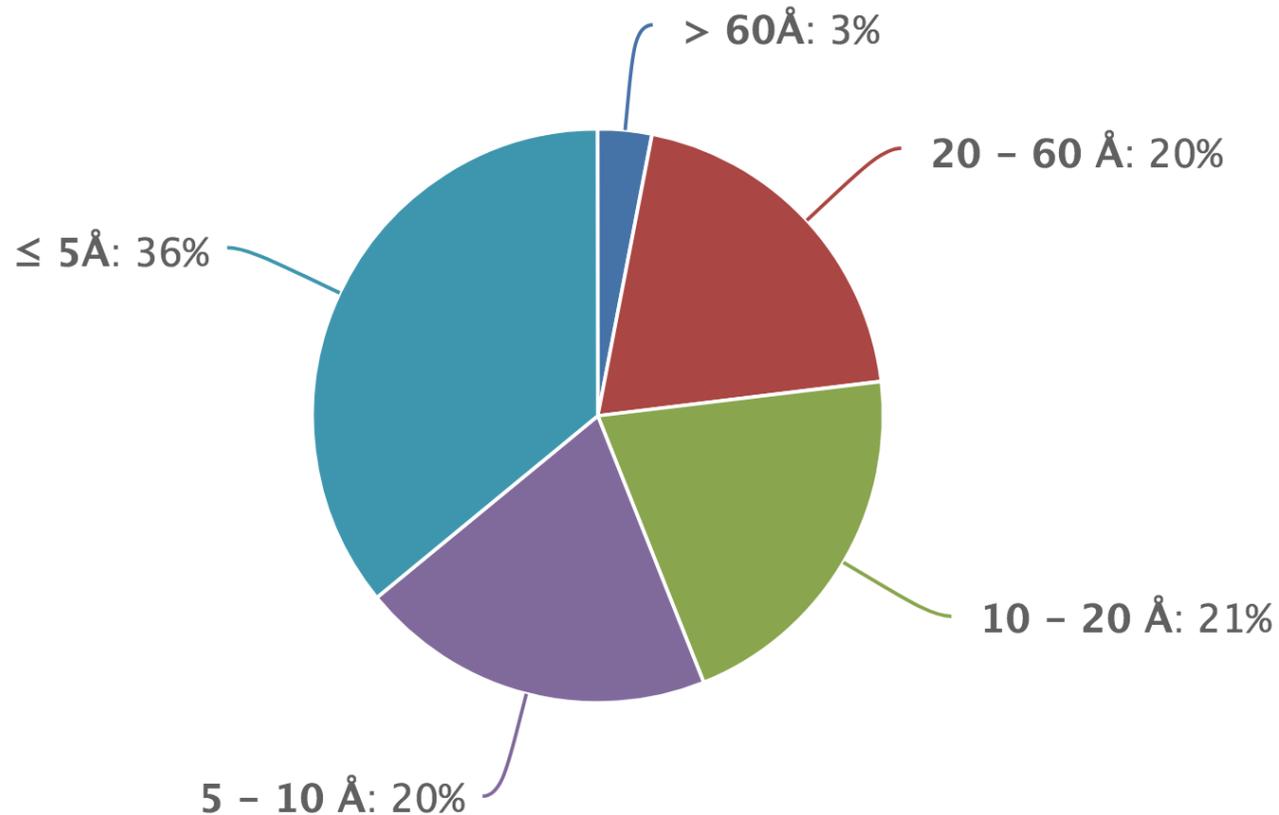
**ISMB**  
Institute of Structural  
and Molecular Biology

# Aims of this lecture:

- To understand 3D-EM density fitting and what we can achieve with it.
- To describe the different types of density fitting methods:
  - rigid fitting
  - flexible fitting
  - assembly (multiple) fitting
- To be aware of some software tools used for visualization and density fitting.

# EMDB Statistics

## Resolution distribution for released maps

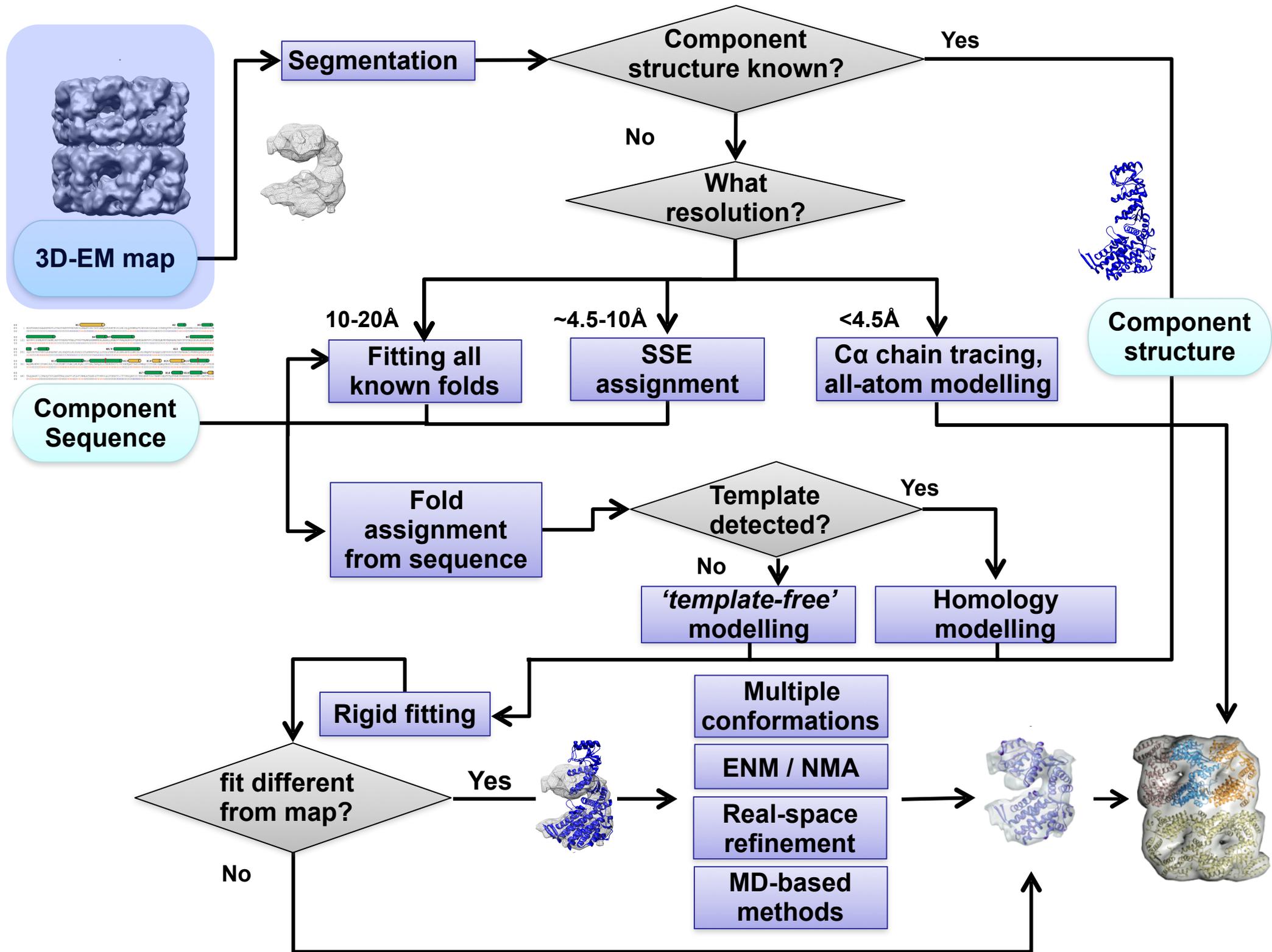


~More than half of the maps are better than ~10 Å resolution!

**Q:** How can we get more out of these 3D-EM maps?

**A:** We can use them as a constraint in model building

Model building and refinement is often  
**interactive** and **iterative**.



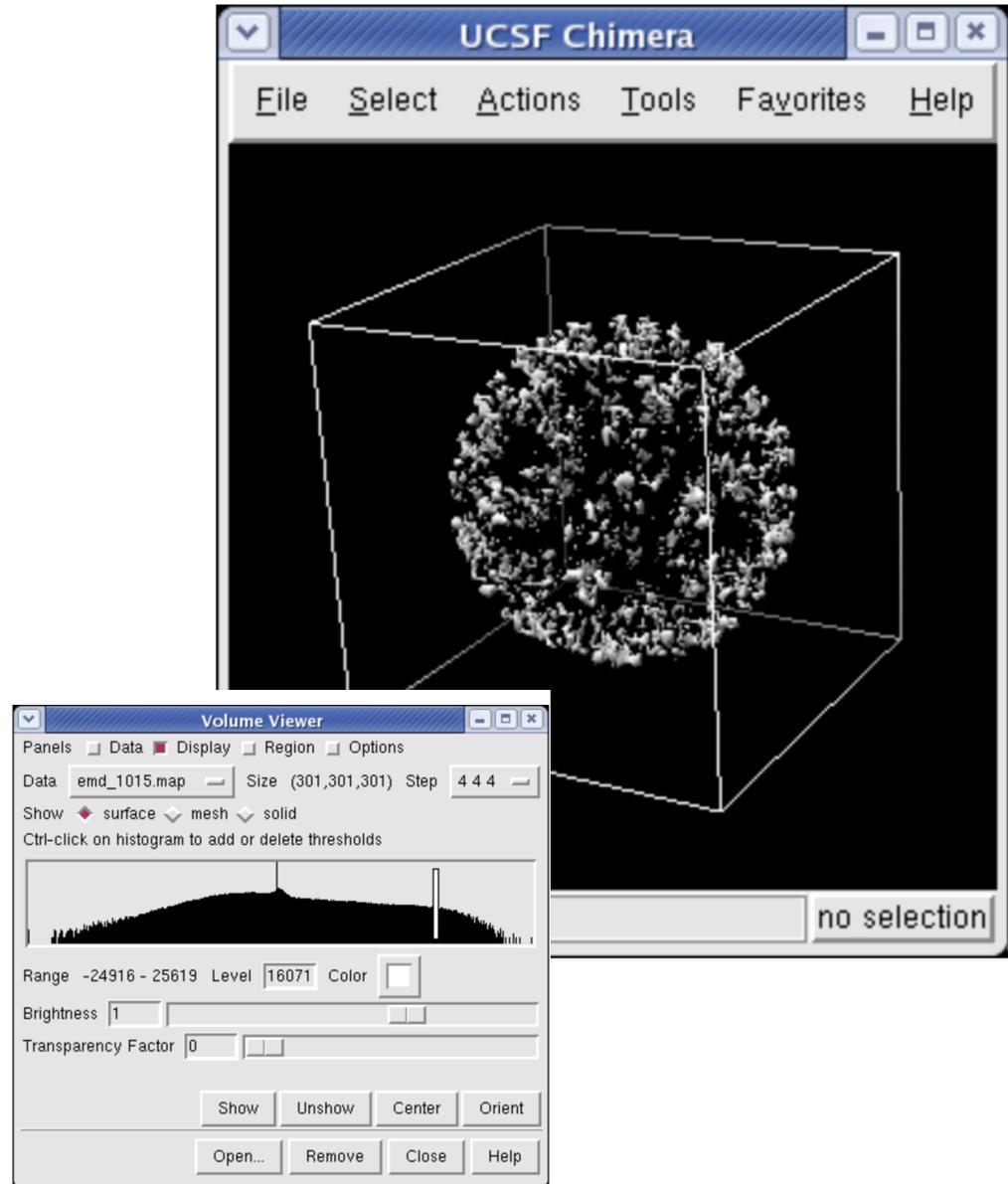
# Visualization tools

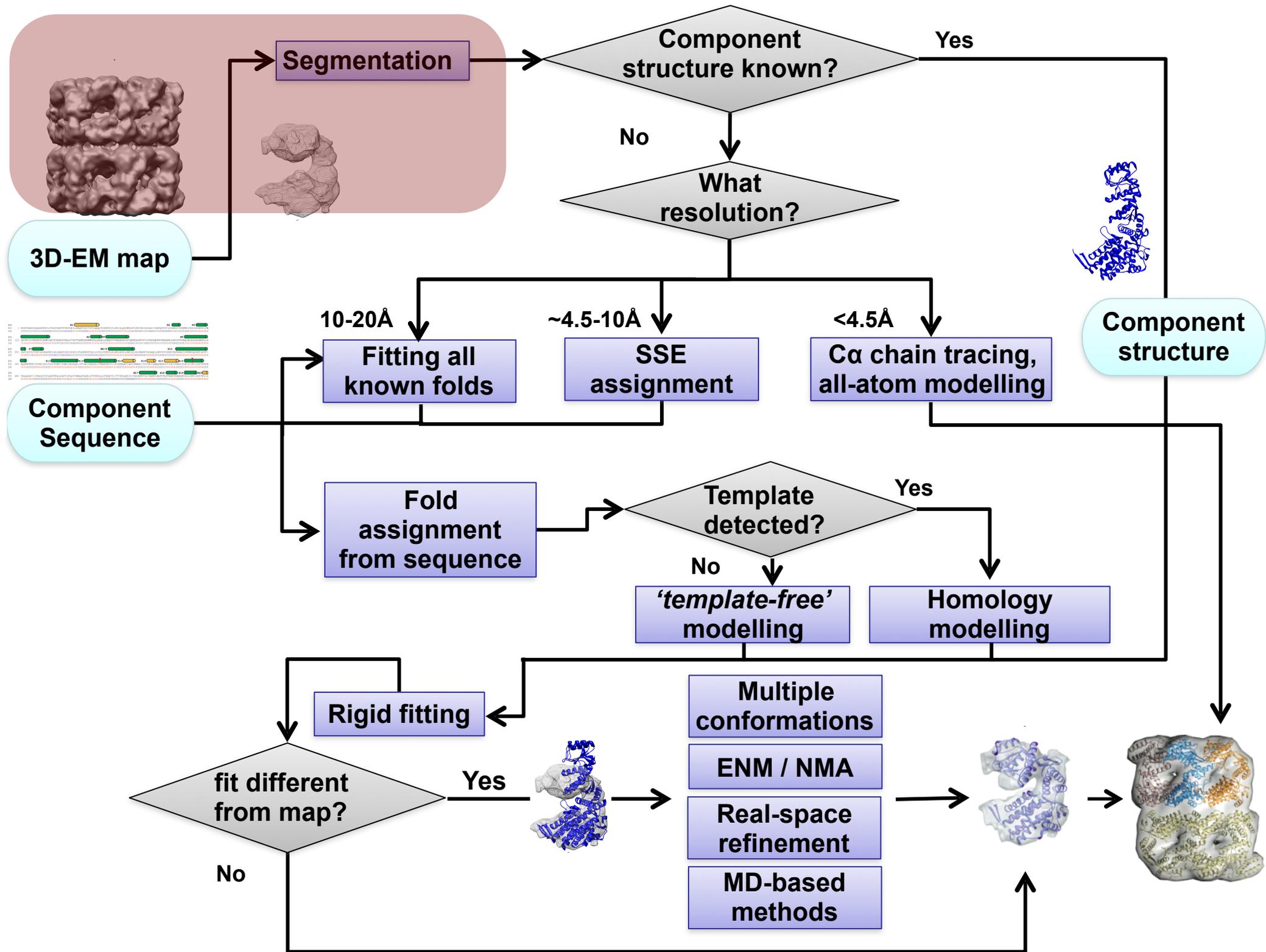
## - Academic programs:

- Chimera, ChimeraX
- Coot
- Python Molecular Viewer (PMV)
- VMD
- VolRover
- Gorgon

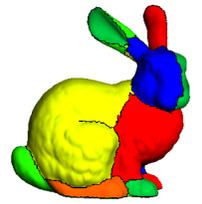
## - Commercial programs:

- PyMOL (Schrödinger)
- Amira (Thermo Fisher Scientific)

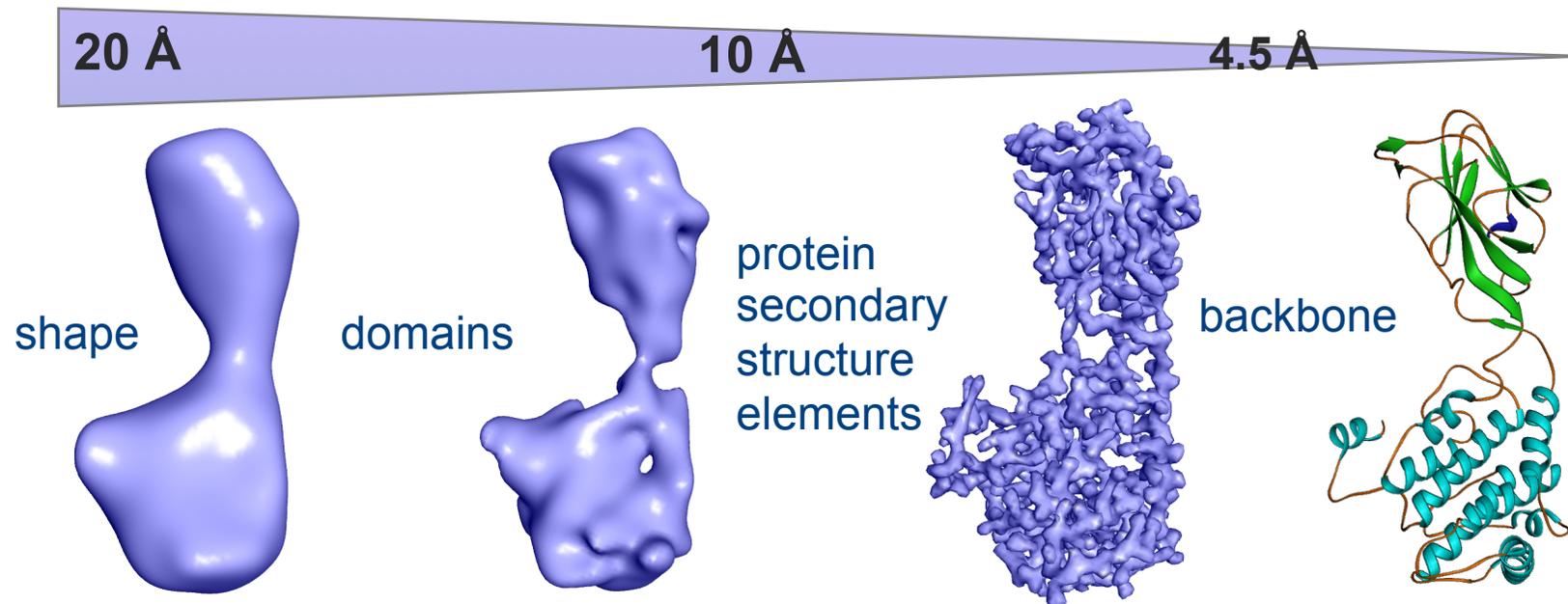




# Segmentation

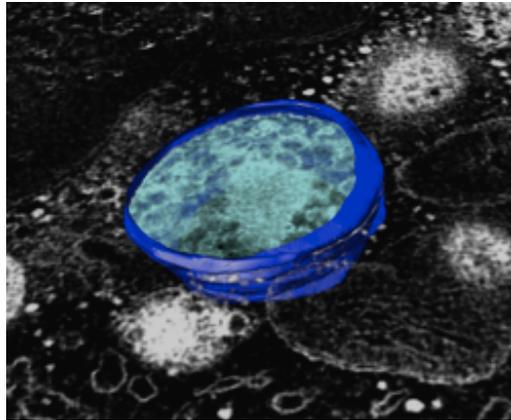


- Identify boundaries between 3D regions that represent structural components in the context of *structural*, *biochemical* and *bioinformatics* knowledge.
- The identified boundaries can be useful in detecting the positions of known component structures in the map.
- The size of the segmented components is related to the map resolution.

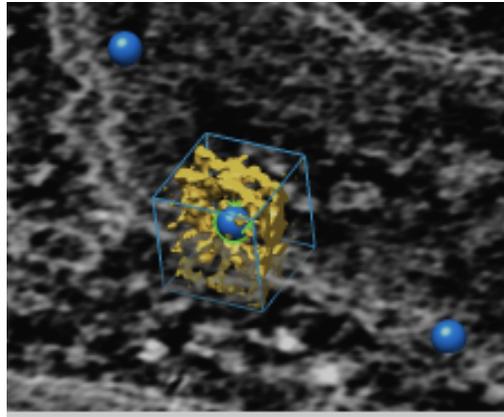


# Segmentation

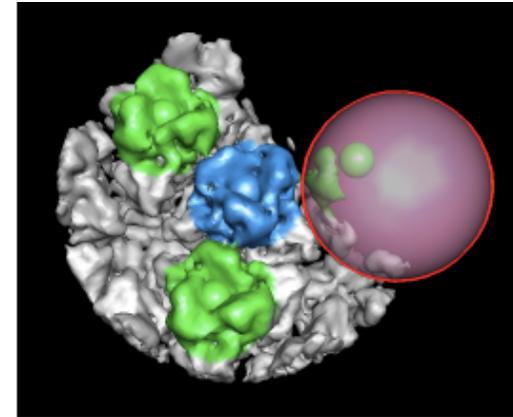
## - Manual segmentation



Mask



Box around marker/atoms

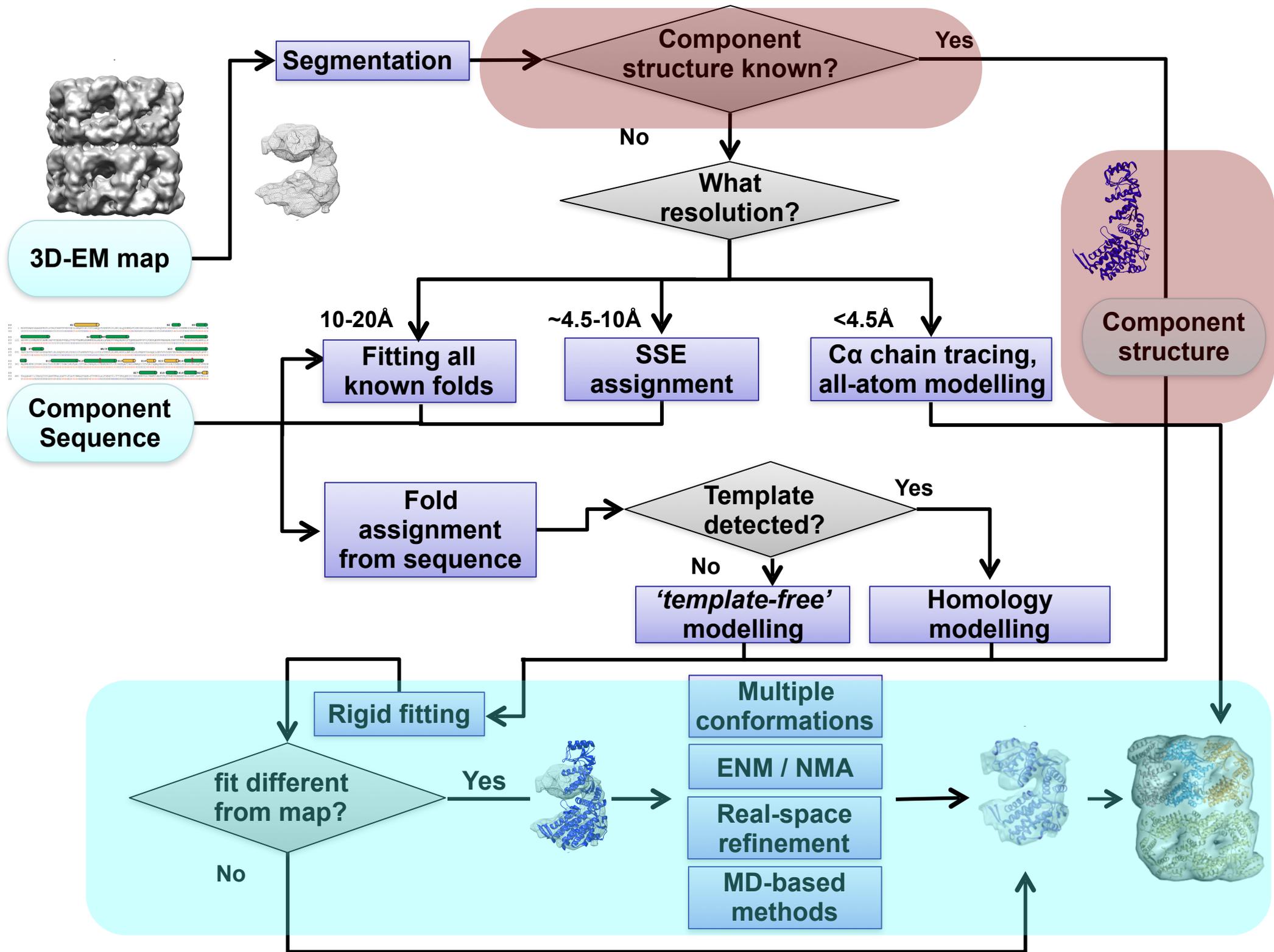


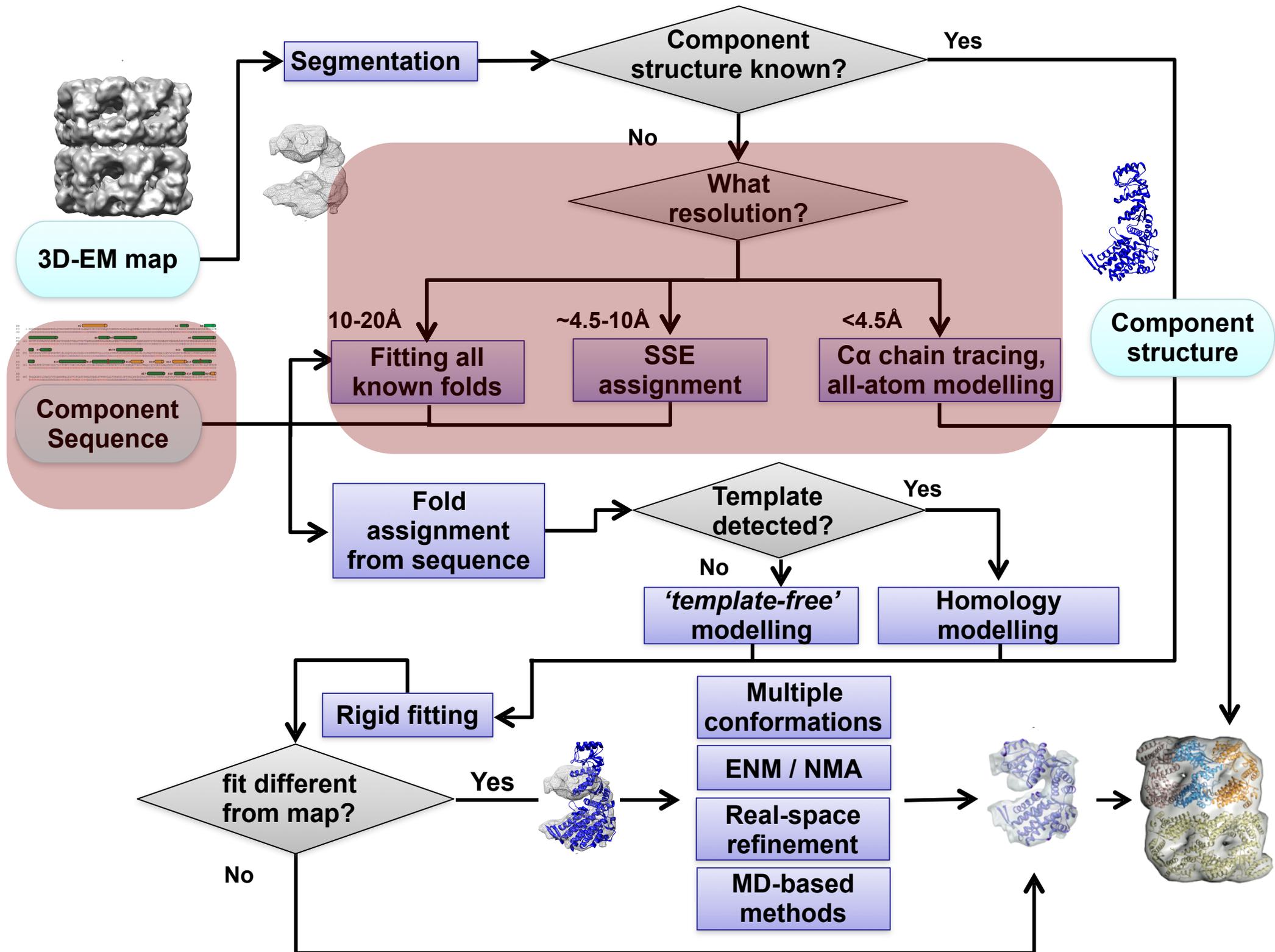
Hand erasing

- **Automated segmentation:** based on density alone, with or without the use of symmetry information. (e.g. in VolRover, Segger, Amira, IMOD)

## - Knowledge-based segmentation:

- Antibody labelling; gold clusters; subunit/domain deletion (difference mapping).
- Recognition of structural components - **density fitting**.

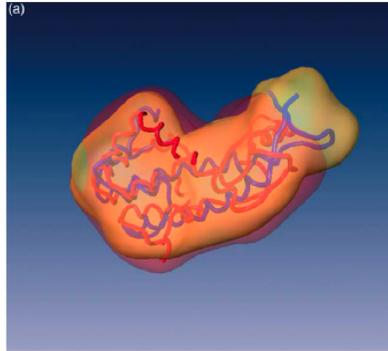




# Fold recognition from density

< ~10-20 Å: Fit domains from a non-redundant protein domain database (e.g. CATH);

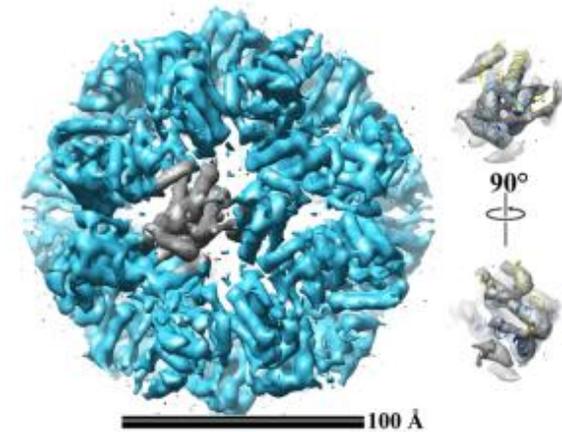
12 Å



Fitting of a domain from 1.20.1060.10 (mainly alpha) into 1.10.530.10 (mainly-alpha).

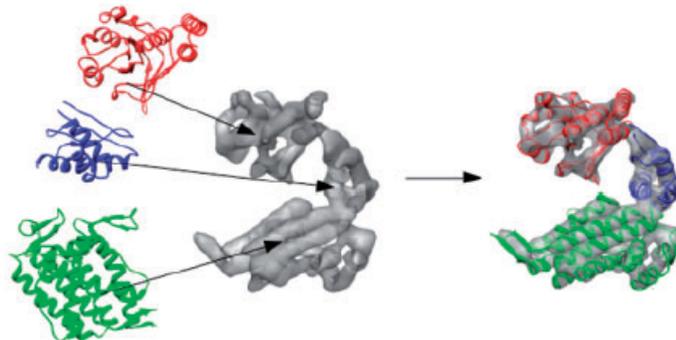
Velazquez-Muriel et al. *JMB* 2005

7 Å



Detection of bacteriophage Lambda

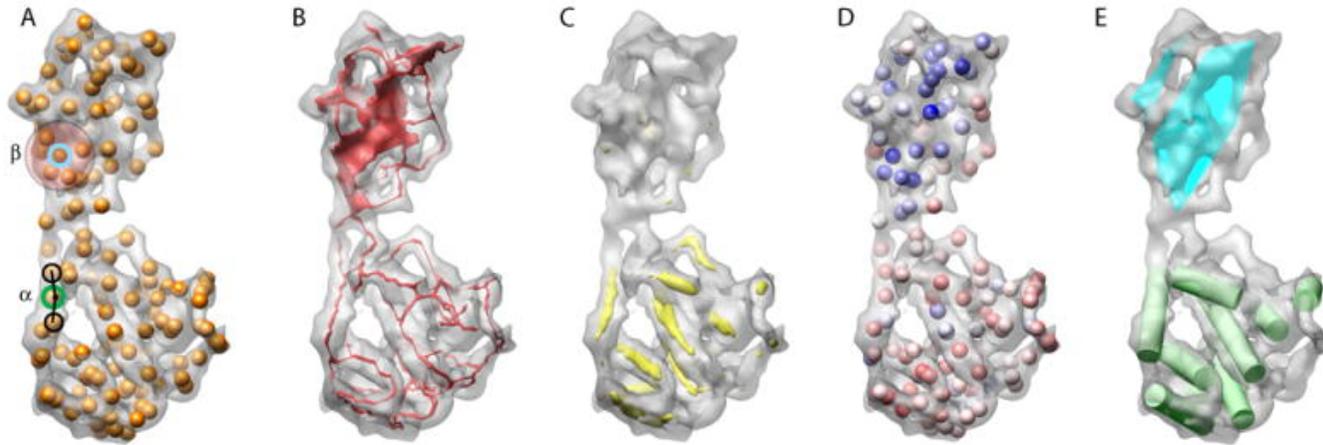
Khayat et al. *JSB* 2010



FOLD-EM: Saha et al. *Bioinformatics* 2012

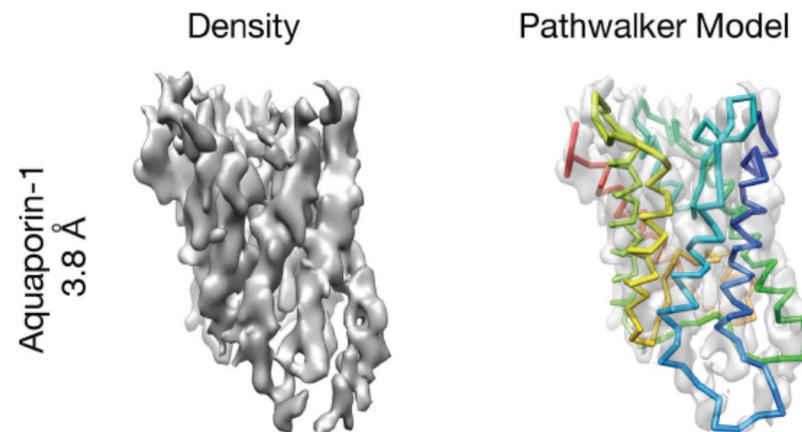
# Fold recognition from density

~4.5-10 Å: secondary structure element detection



Baker et al. *Structure* 2007

4.5 Å and better: *de novo* C $\alpha$  tracing and model building



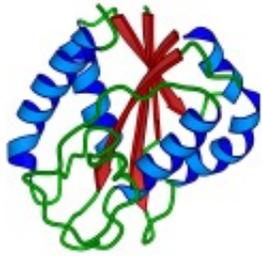
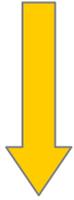
Baker et al. *Structure* 2012

**Programs:** SSEhunter (Gorgon), SSETracer, Ematch, Pathwalker, **Coot**, **Buccaneer**, EM-fold, Rosetta (sequence information), Phenix autobuild, ARP/wARP, MAINMAST



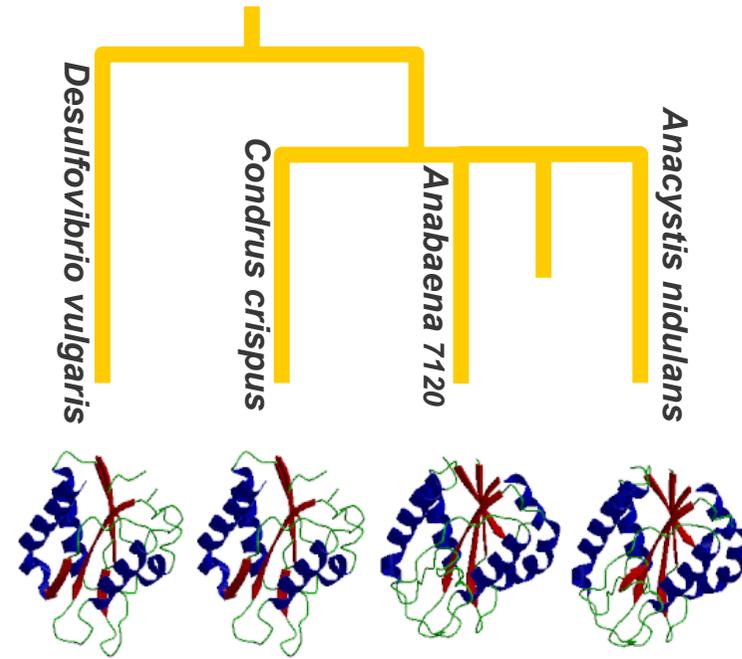
# Fold recognition from sequence

GFCHIKAYTRLIMVG...



Template-free

*Ab initio (de novo)* prediction  
Fragment Assembly  
Evolutionary Couplings



Template-based

Threading  
Comparative (Homology) Modelling

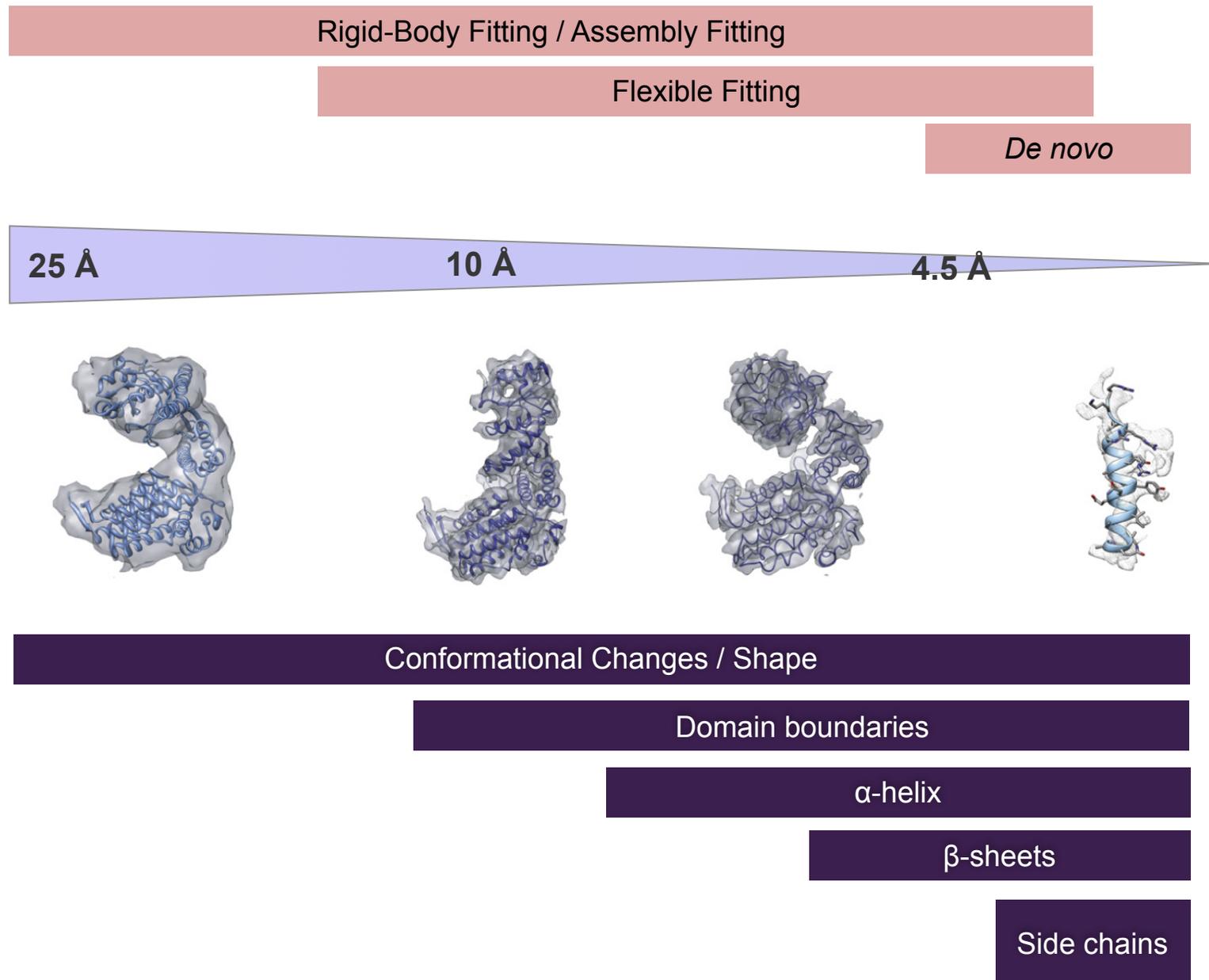
Programs: MODELLER, SWISS-MODEL, Phyre2, RaptorX, I-TASSER, Rosetta, EVfold...



# Density fitting

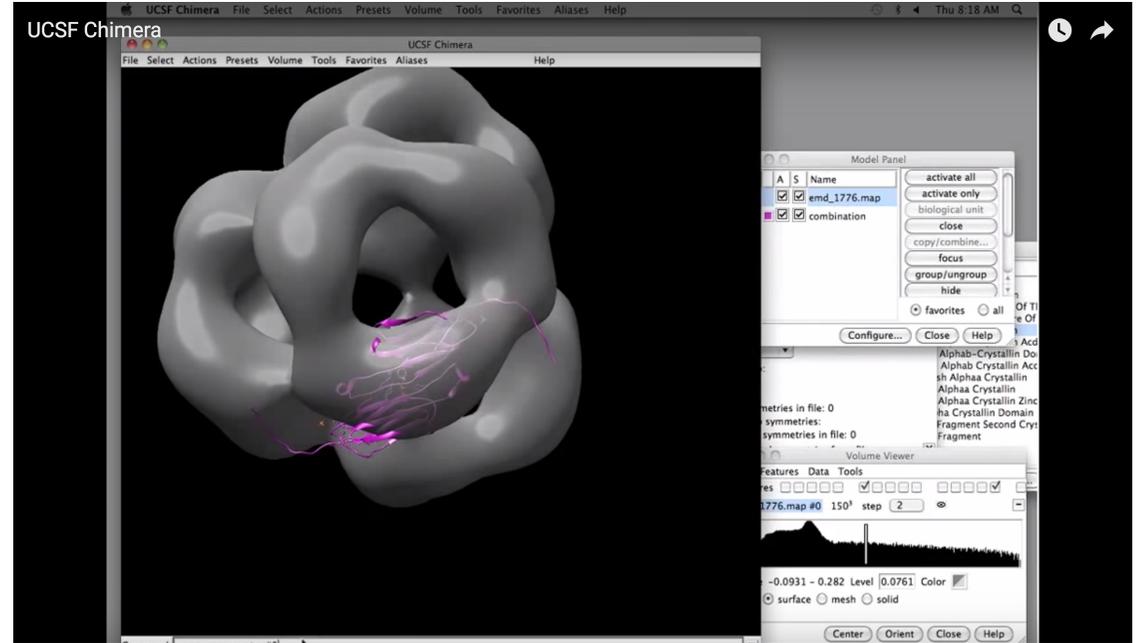
Find the best “match” between the atomic model and the 3D-EM density map

# Density fitting



# Manual fitting

Fitting an atomic structure within the envelope (an isocontour) of the density using visualisation programs.



## Pros:

- Human brain is efficient in certain pattern recognition tasks.
- Immediate feedback and intelligent choices by the user.
- Often good for the initial placement of the component in the map.

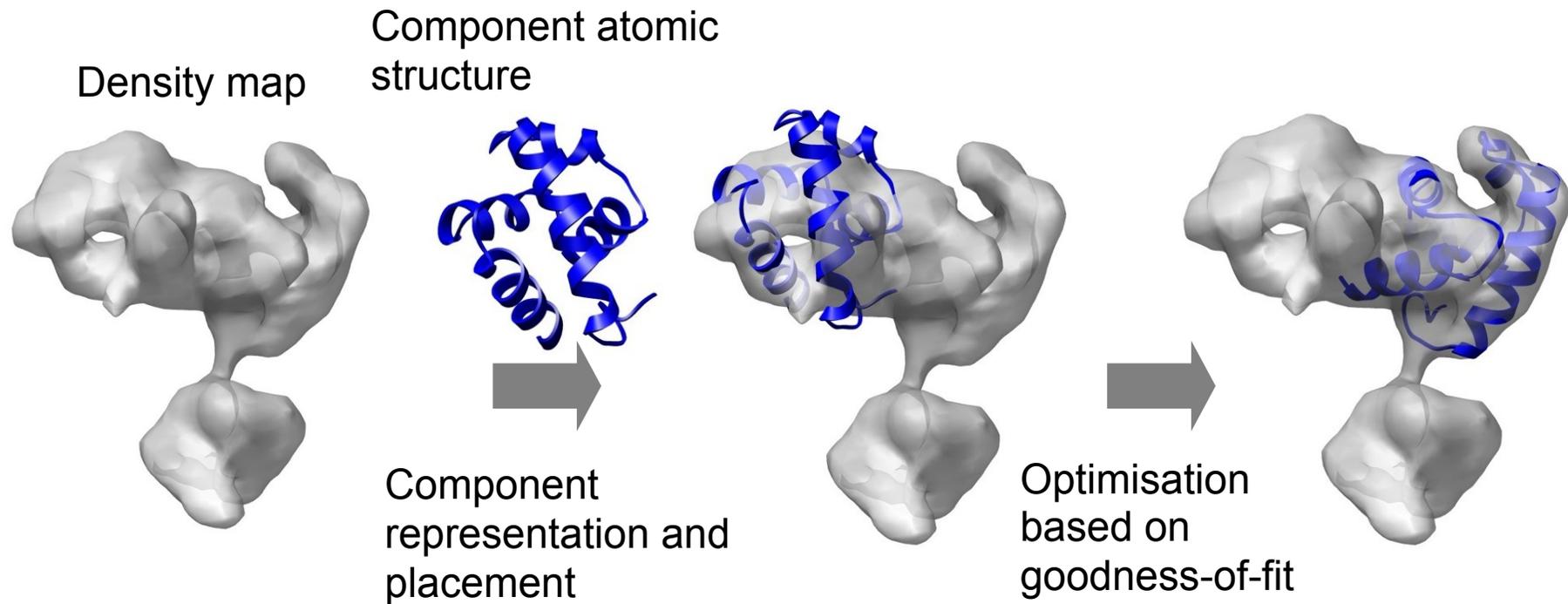
## Cons:

- High level of subjectivity may lead to error, especially if the map does not have sufficient distinctive features for an unambiguous placement of the component.
- Depends on contour level.
- Conformational rearrangements cannot be modelled (misfits and steric clashes).

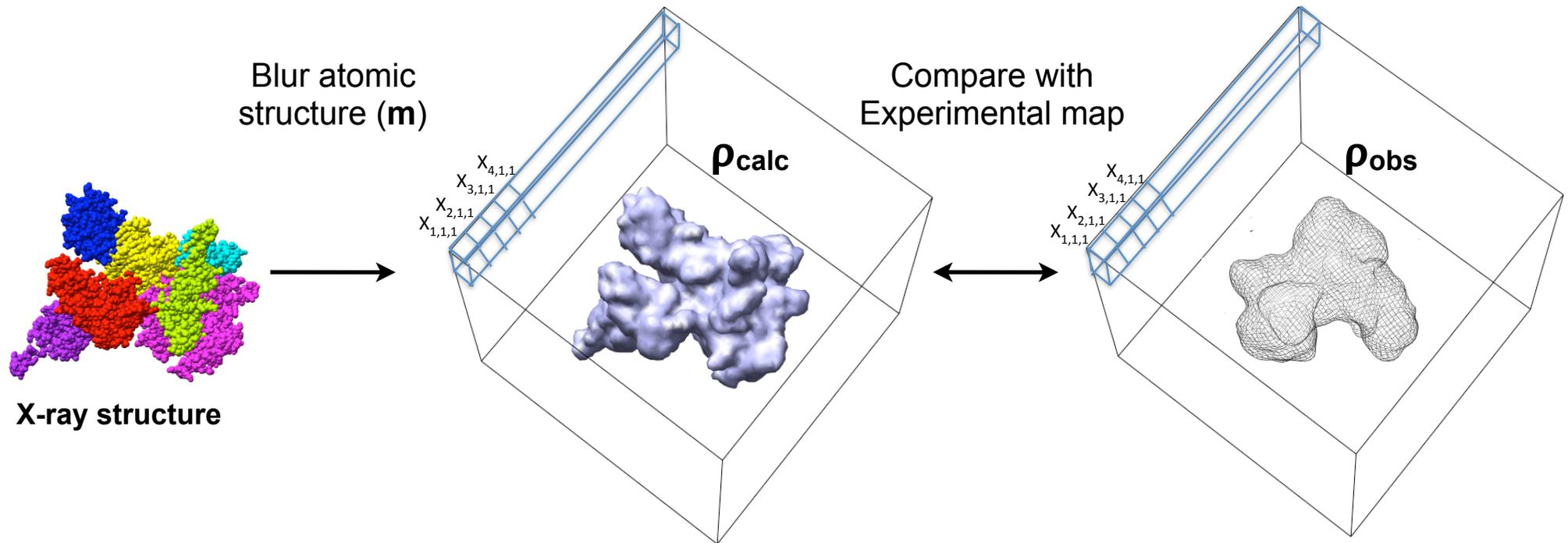
# Automated fitting

All automated fitting methods require:

1. a way of representing both the structure and the density map (**representation**).
2. a way of measuring the goodness-of-fit (**scoring**).
3. a method of finding the best fit (**optimisation**).



# Representation and scoring



$$R_{\rho, \text{lsq}} = \int_{\mathbf{x}} (\rho_{\text{obs}}(\mathbf{x}) - \lambda_{\rho} \rho_{\text{calc}}(\mathbf{x}, \mathbf{m}))^2 d^3 \mathbf{x}$$

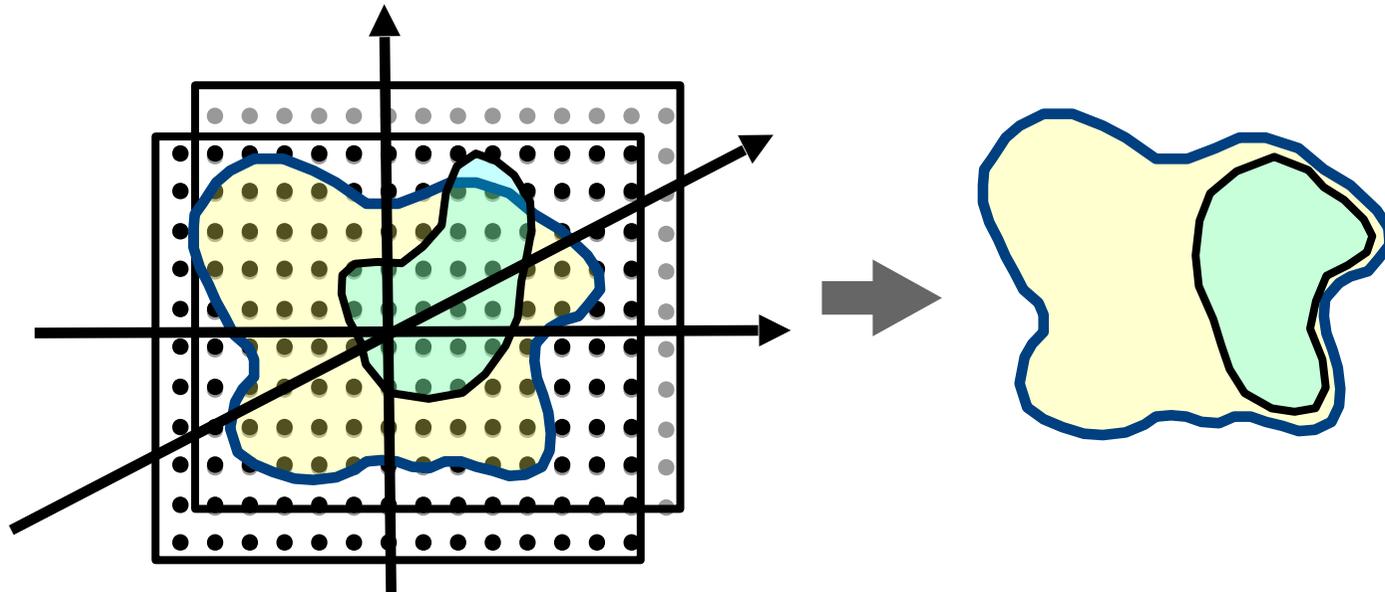
Cross Correlation Coefficient

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

# Exhaustive search

**Pros:** Get the global solution in respect to a given scoring function.

**Cons:** The search in real space is too large for most scores (very expensive).



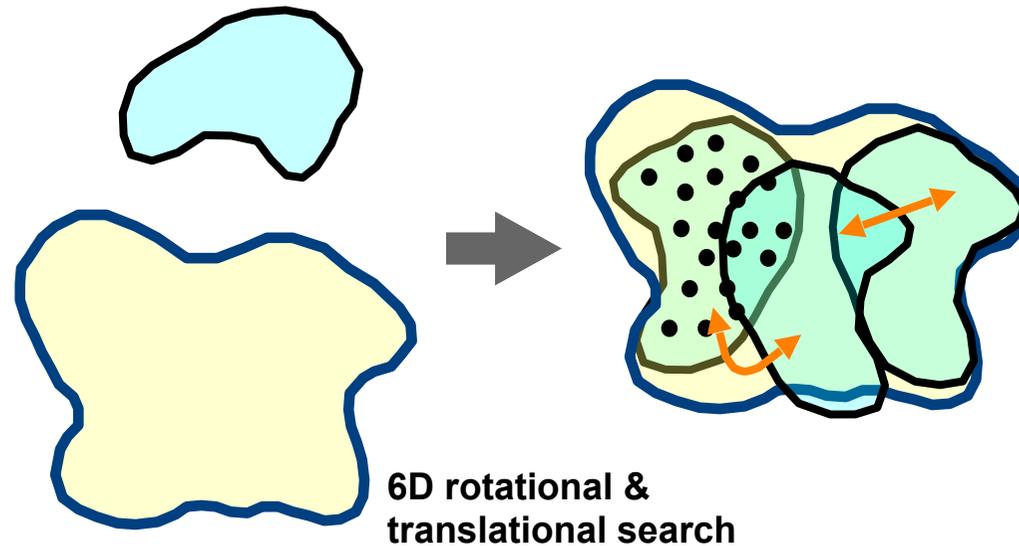
- **Acceleration:** FFT (translational moves); Spherical harmonics (rotational moves)  
COLORES, DOCKEM, ADP-EM, PowerFit, gEMfitter (GPU acceleration)...

- **Local fitting** - Search exhaustively a given sub-region in the map (Mod-EM, Chimera)

# Stochastic/random and gradient methods

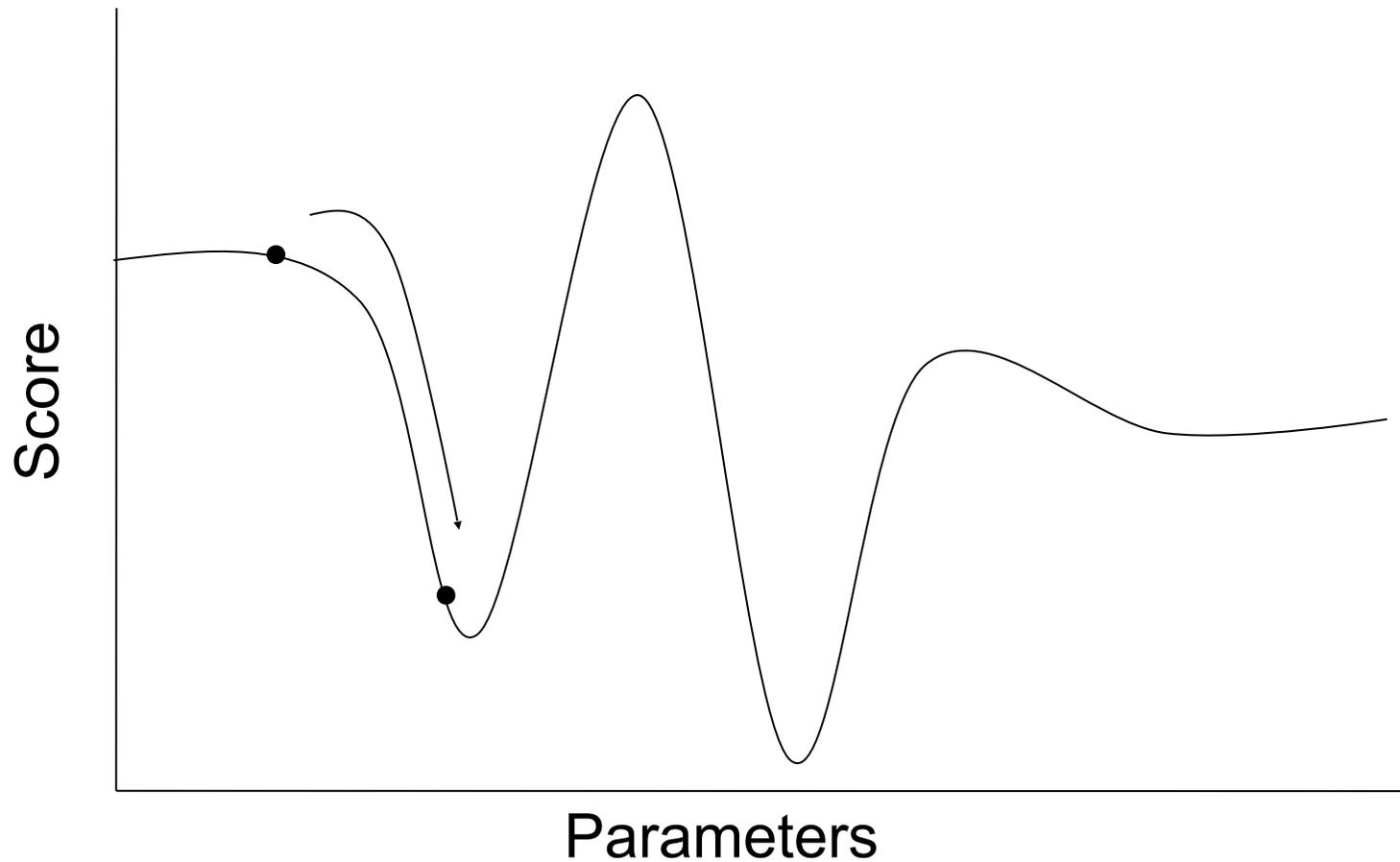
**Pros:** Fast; easy to implement different scoring functions.

**Cons:** The model can be “trapped” in a local minimum



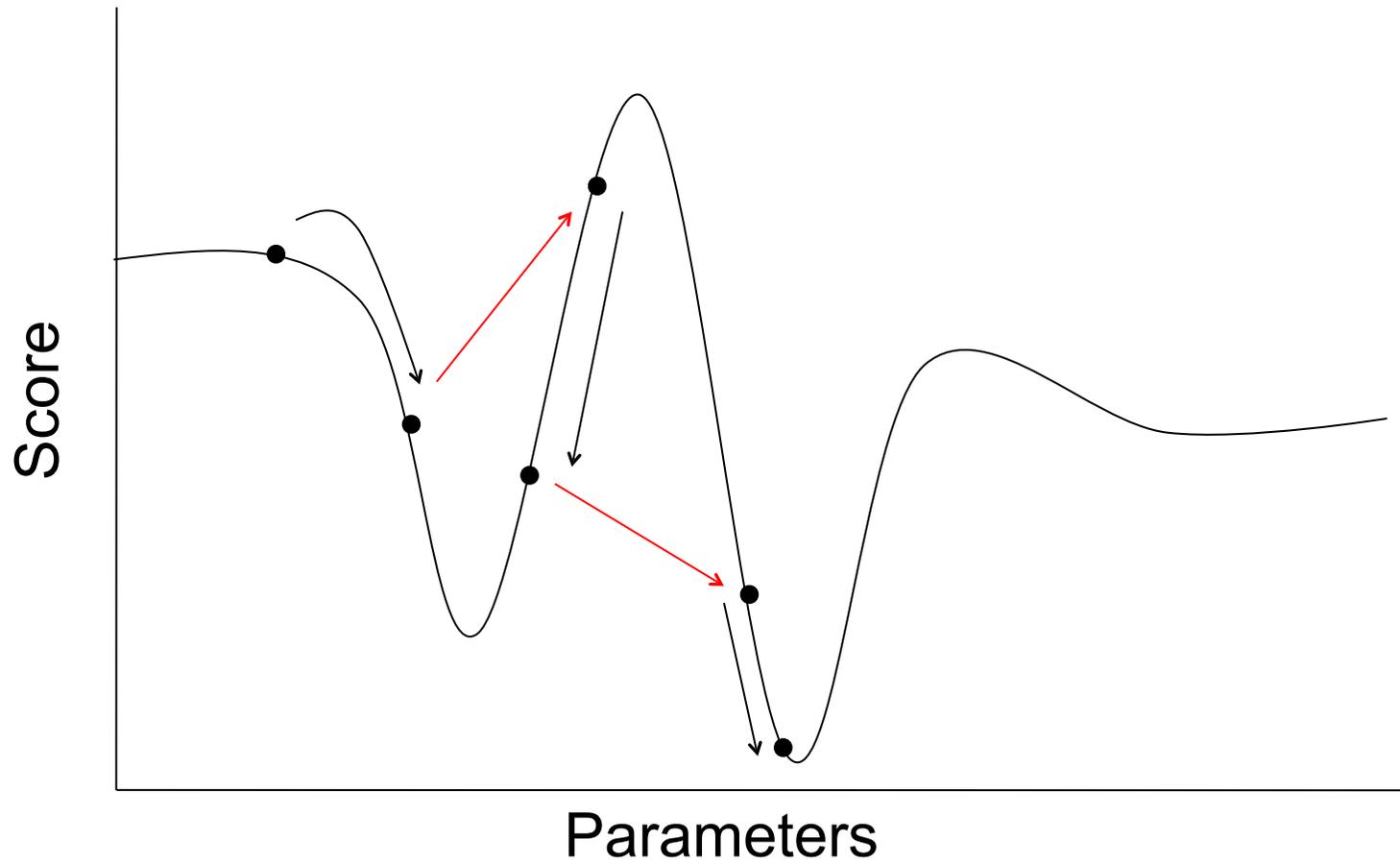
**Programs:** Mod-EM, Chimera, Rosetta, IMP, HADDOCK, GMfit (gaussian approximation)

# Gradient-based methods



Optimisation follows steepest gradient

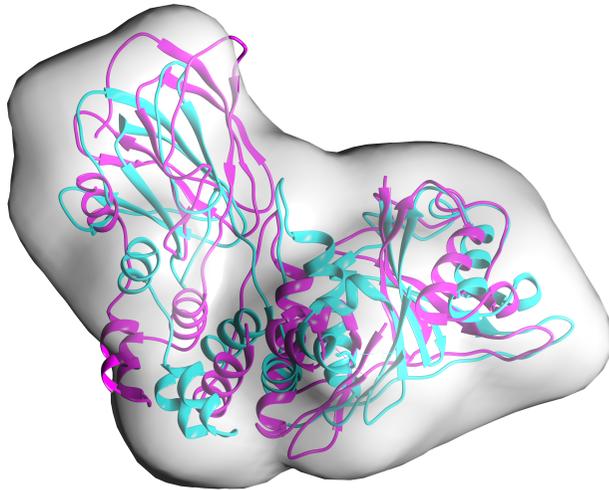
# Stochastic methods



Example: simulated annealing the optimisation follows a gradient method, with 'jumps' to avoid local minima

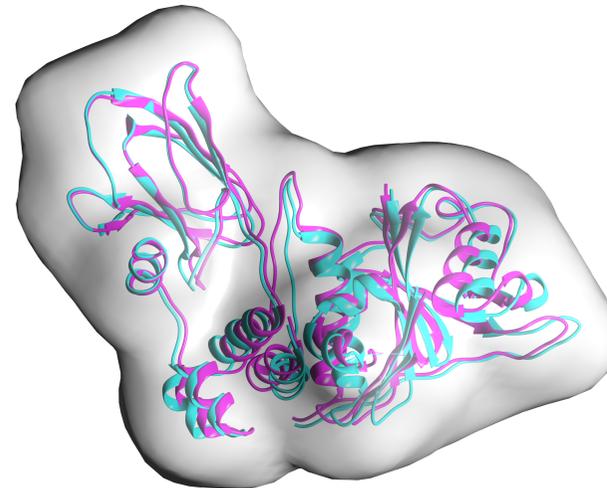
# Refinement at intermediate resolution

Before refinement



**C $\alpha$  RMSD from native: 7.5 Å**

After refinement



**C $\alpha$  RMSD from native: 2.1 Å**

1VCB, 10 Å resolution

native

best predicted fit

# Problems of density fitting

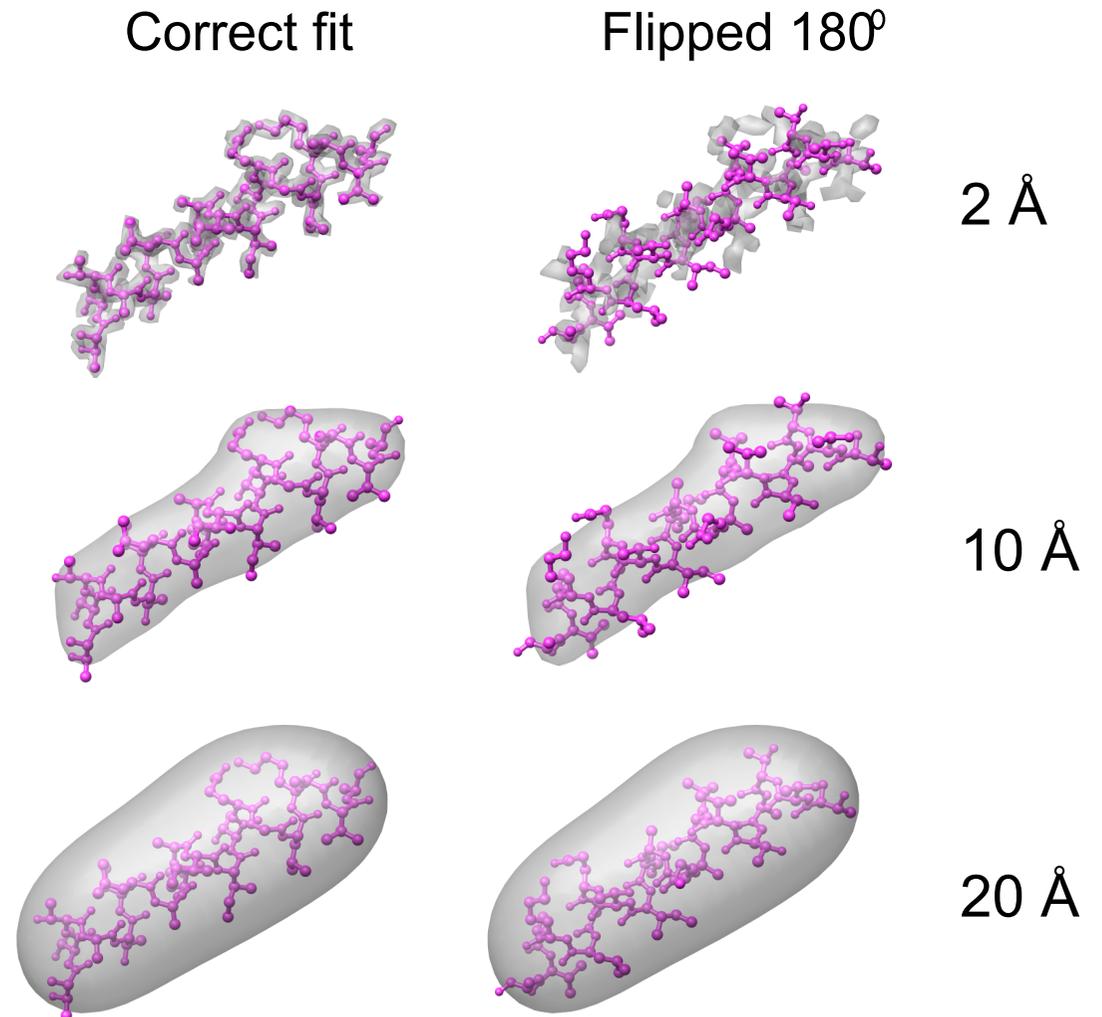
# i. Limitations of resolution

## Problems:

- At low resolution: many local optima with similar numerical values.
- Local resolution, noise, scaling, filtering, masking.
- Blurring of the atomic structure.

## Solutions:

- (Improve your resolution!)
- Improve scoring for goodness-of-fit.
- Coarse-graining (change representation)
- Fit/model validation

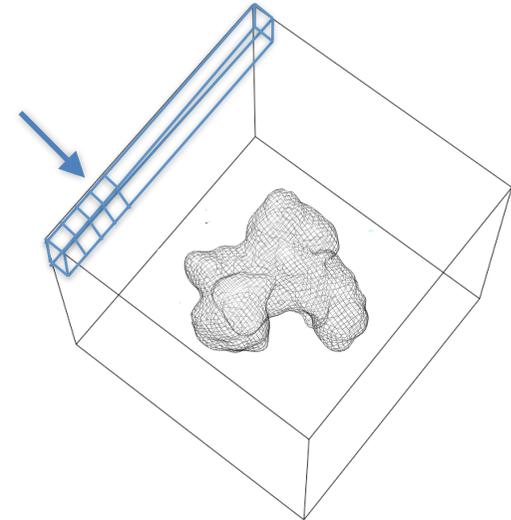


# Density-based scoring functions

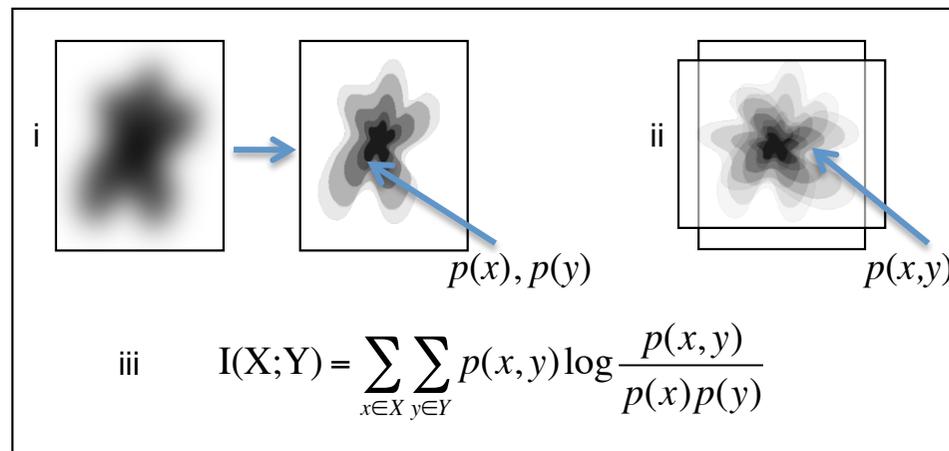
TEMPy: <http://tempy.ismb.lon.ac.uk/>

- Cross-correlation coefficient (CCC)

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



- Mutual information-based score (MI)



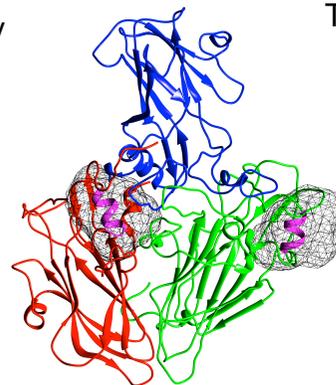
Useful at intermediate resolutions; noisy maps;  
less sensitive to relative intensity levels

# Local scoring

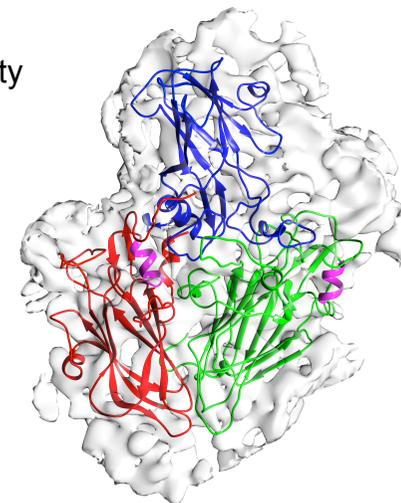
$$SCCC = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

scores local segments in structure

Probe density  
X



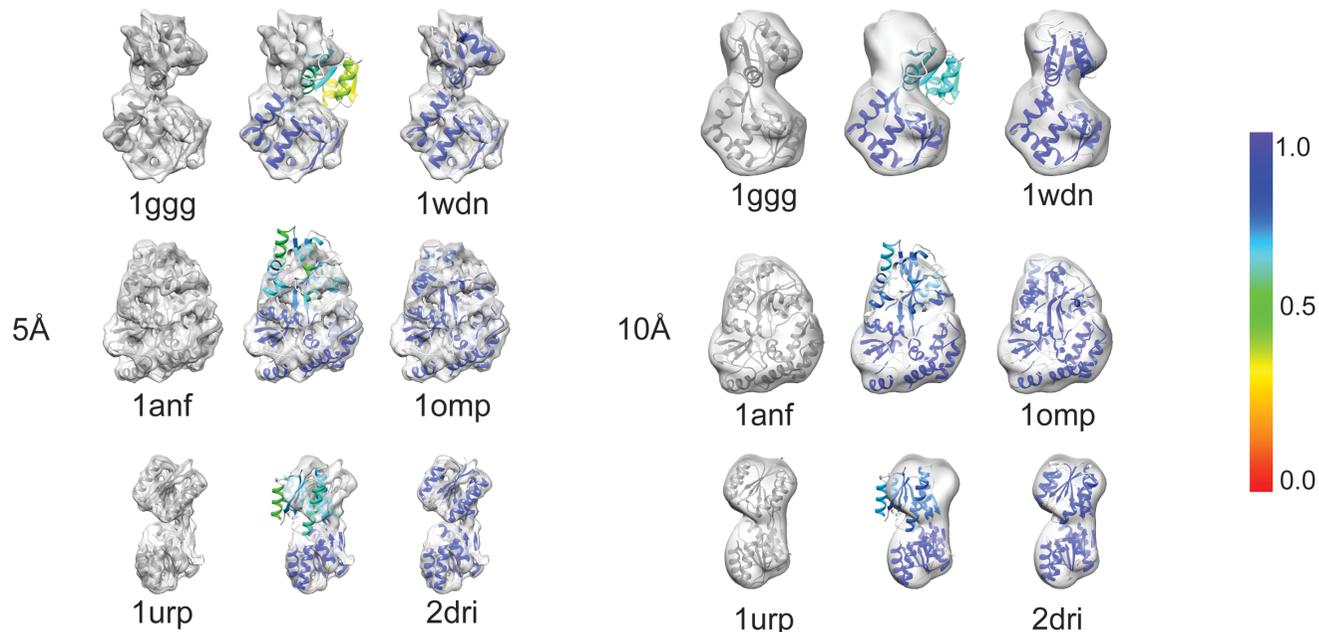
Target density  
Y



Roseman, *Acta Crystallogr D* 2000; Pandurangan *et al.*, *J Struct Biol* 2014

TEMPy + Chimera attribute files

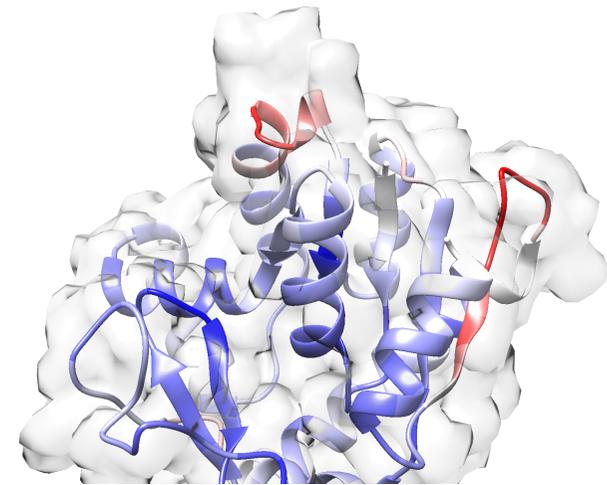
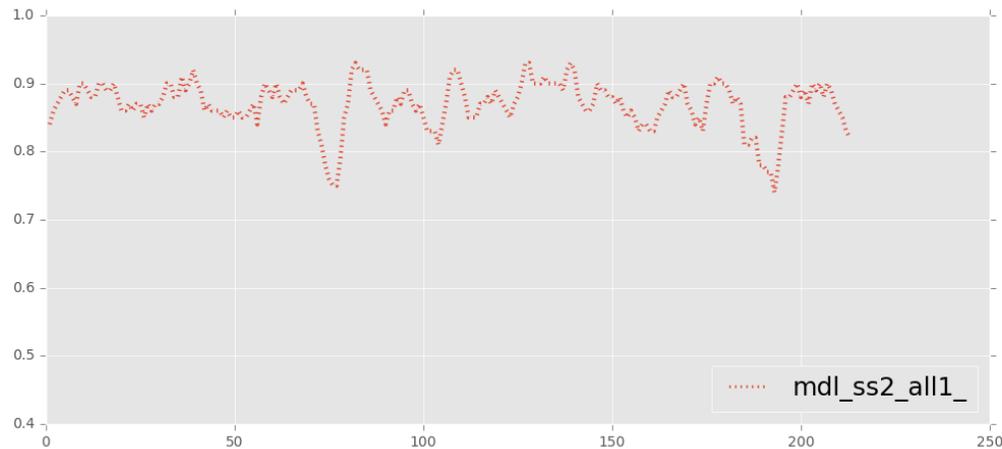
Useful for calculating CCC on any defined local segment



# Local scoring

- Segment-based manders' overlap coefficient (SMOC):

Calculated on overlapping segments along the sequence and assigned to central residue so that each residue has a score.



$$SMOC = \frac{\sum_{i \in vox\_sr} \rho_i^{EM} \times \rho_i^P}{\sqrt{\sum_{i \in vox\_sr} (\rho_i^{EM})^2 \times \sum_{i \in vox\_sr} (\rho_i^P)^2}}$$

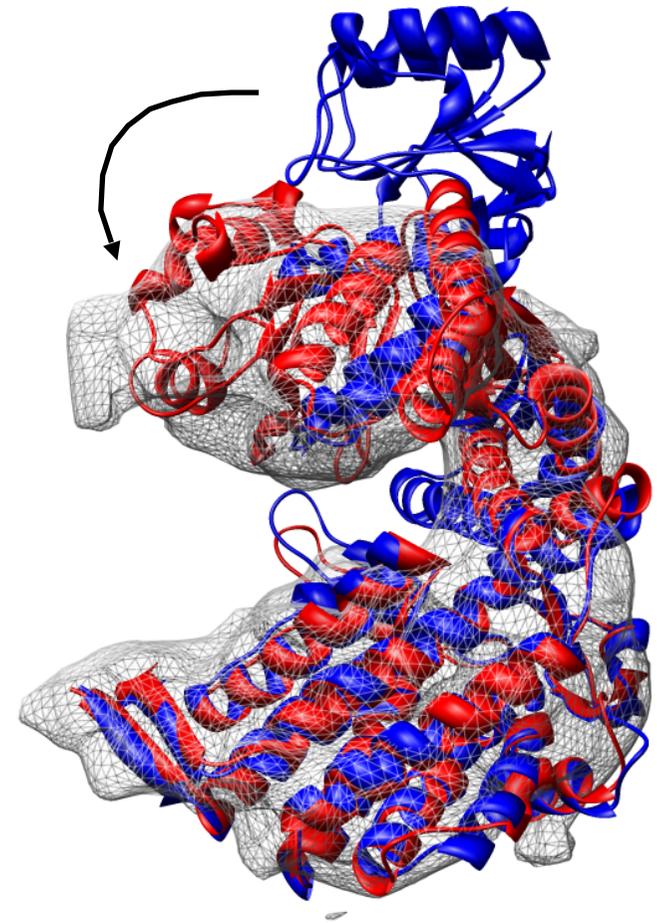
Useful to calculate local fit per residue (segment)

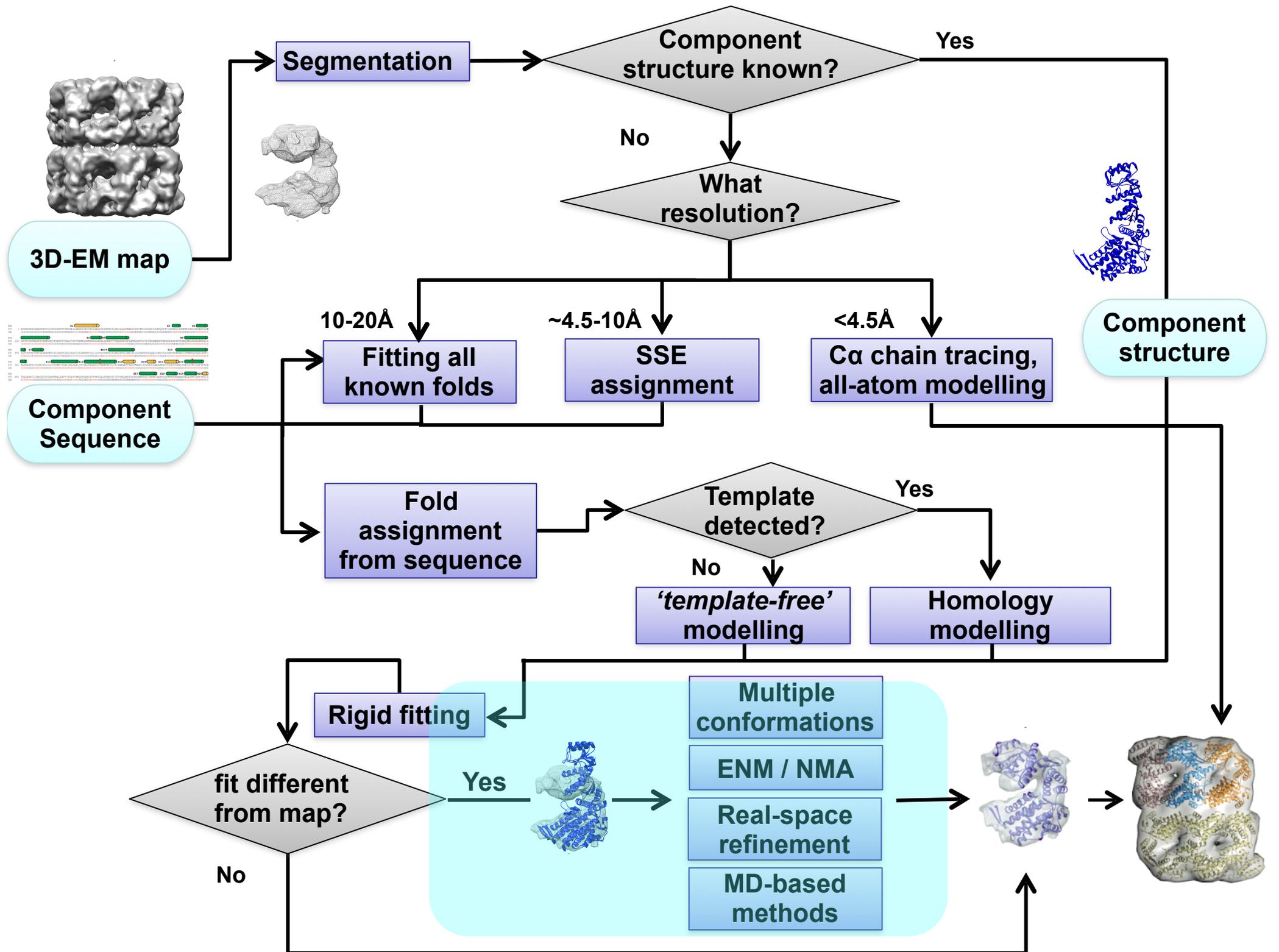
## ii. Conformational variability

**Problem:** Conformations observed by 3D EM often deviate from the conformations of the atomic models we fit.

- Dynamics.
- Crystal packing effects.
- Errors in structure prediction.

**Solution:** change the conformation of the atomic model during the fitting process — **flexible fitting**.





# Model refinement

Without any restraints a model may fit well with a high score in near-atomic-to-low resolution density: “perfectly overfitted model” (e.g. Faulkner et al. 2013)

The resulting model however will not have standard protein geometry:

backbone torsions: phi/psi (Ramachandran space), peptide planarity, chirality (trans/cis), bond lengths and angles, side chain torsions / rotamers

Refinement methods try to maintain standard geometry while fitting the model in density. These geometry restraints reduce the degrees of freedom (sampling space).

# Approaches to refinement

- Elastic Network Model (ENM)
- Normal Mode Analysis (NMA): A collection of harmonic oscillators; those with low frequency and large amplitude motions often correlate with experimentally observed conformational changes.
- Geometry-based conformational sampling using harmonic restraints

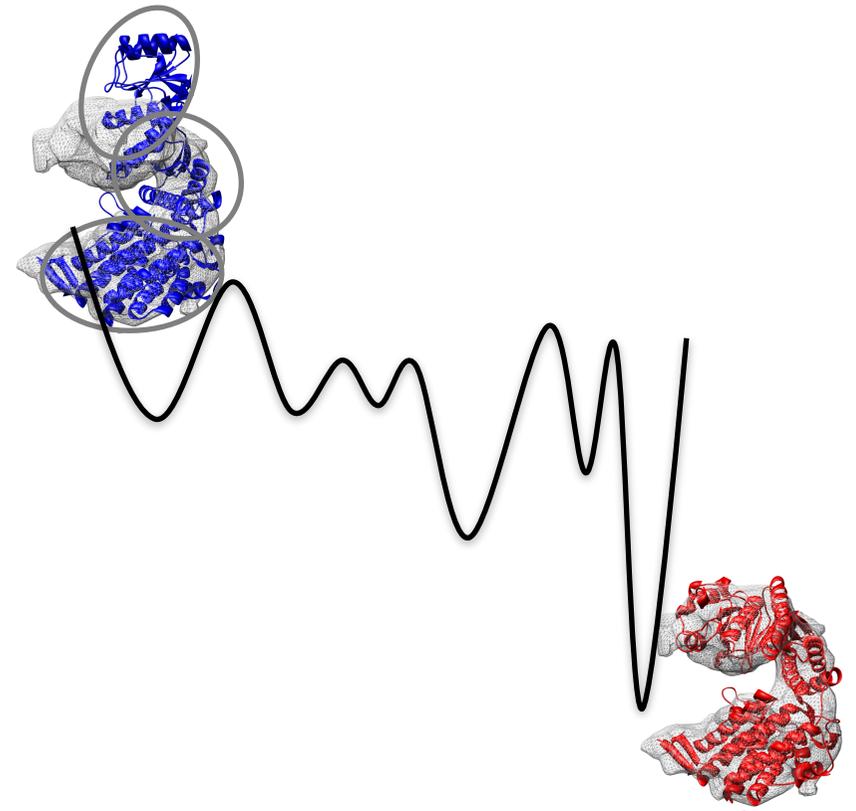
# Real-space refinement - Flex-EM

- The fit of the probe structure is optimised simultaneously with the stereo-chemical properties by the minimisation of a scoring function, such as:

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

- Optimisation is performed on **rigid bodies (b)** by energy minimisation and molecular dynamics.

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



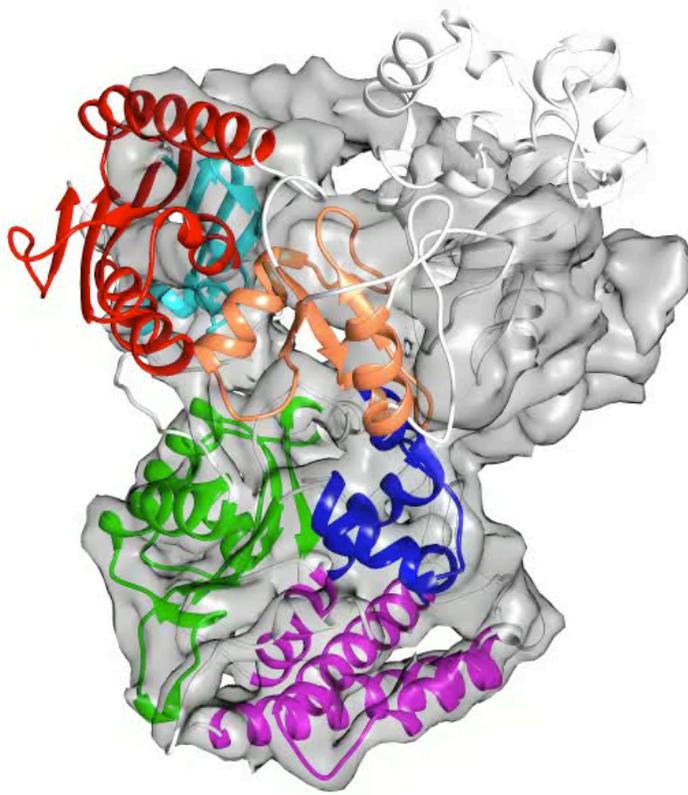
# Rigid-body restraints

A cluster of atoms that form a compact structural segment through a network of contacts can be restrained :

- when the resolution of density map is insufficient to fit smaller entities like individual residues or atoms.
- to allow faster large body movements in the initial stages or refinement

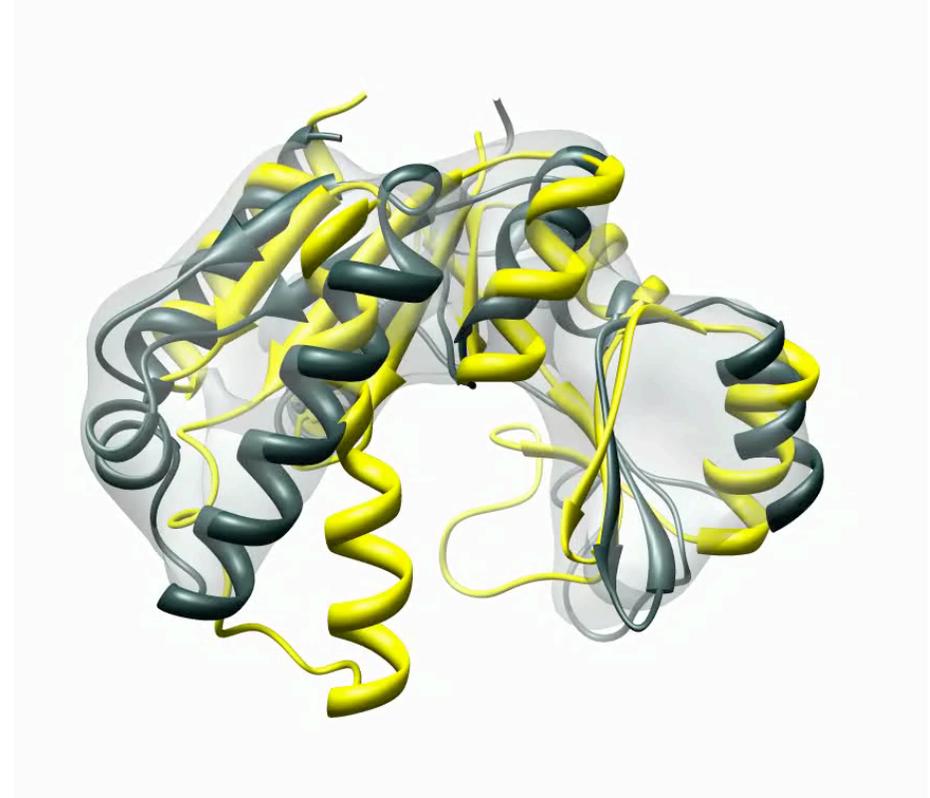
**Flex-EM** can use **RIBFIND** cluster segments based on secondary structure contacts. Long range distance restraints can be also added using MODELLER

# Refinement at intermediate resolutions



Rigid bodies: sub-domains

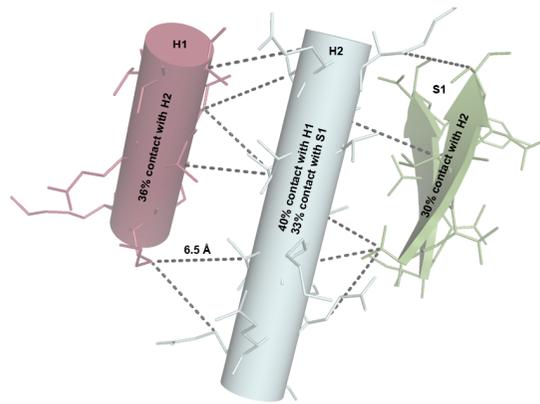
5 Å resolution



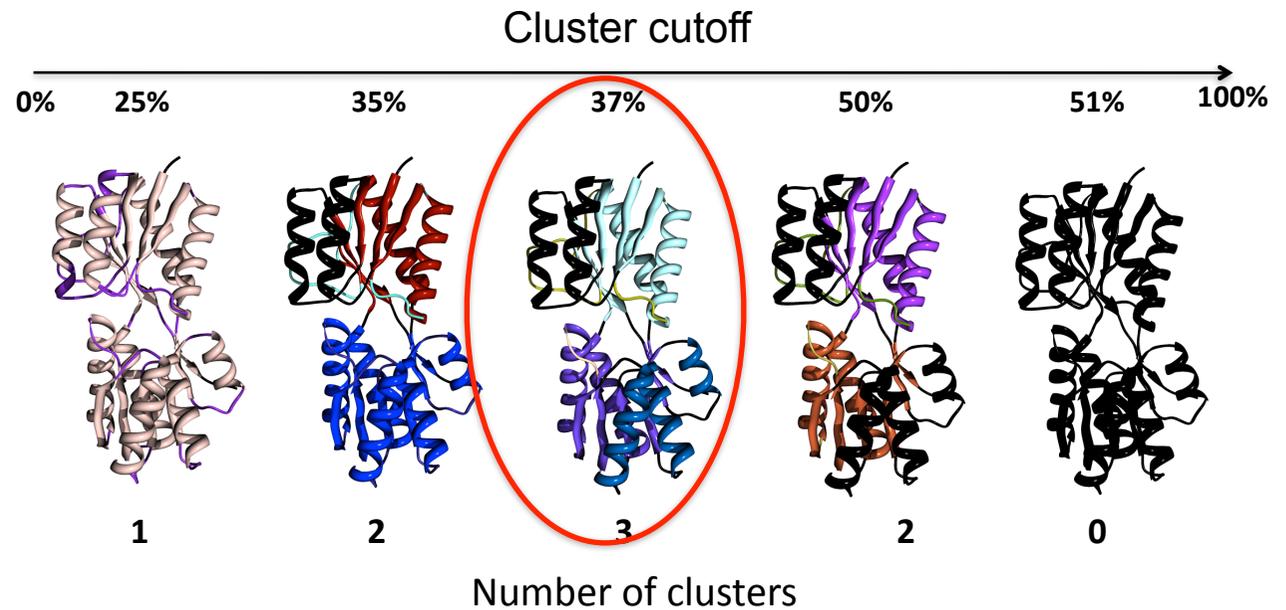
secondary structure elements

10 Å resolution

# RIBFIND: identify sets of rigid bodies



SSE-based clustering



RIBFIND: finding rigid bodies in protein structures

**Resources**  
[HOME](#)  
[Topf\\_Group](#)  
[DSSP\\_server](#)  
**Documentation**  
[Download](#)  
[RIBFIND\\_enquires](#)

Total no. of SSEs: 11  
No. of clusters: 3  
No. of SSEs in cluster: 9

Cutoff value : 36

On/Off Density map

Map threshold : 0.3

[Download rigid body file](#)

[Download all results](#)

[View log file](#)

<http://ribfind.ismb.lon.ac.uk/>

# Refinement with Flex-EM/RIBFIND

RIBFIND rigid bodies

Initial

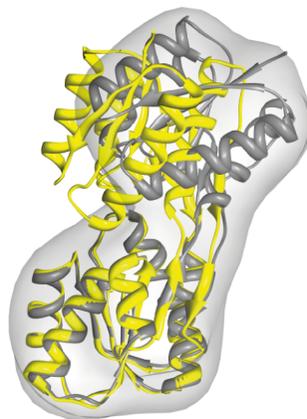
Final non-clustered

Final clustered

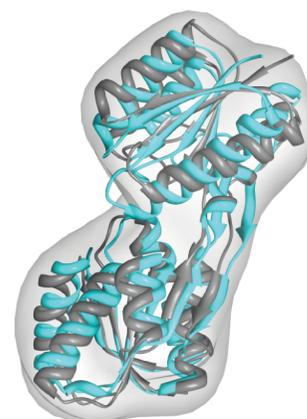
2driA, 10 Å



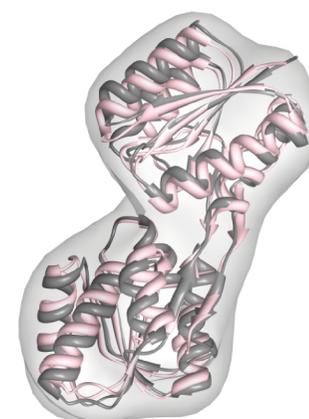
RMSD from target:



7.55 Å

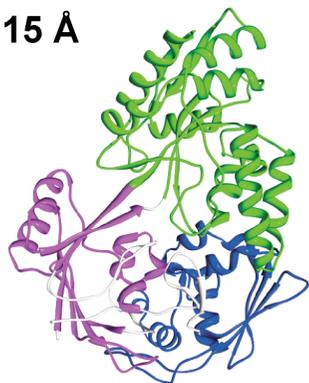


3.04 Å

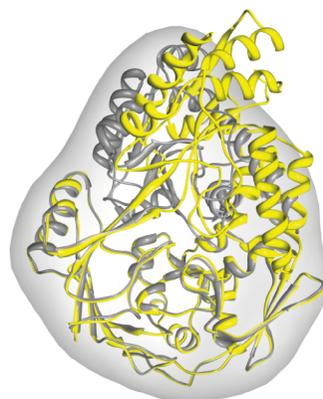


1.71 Å

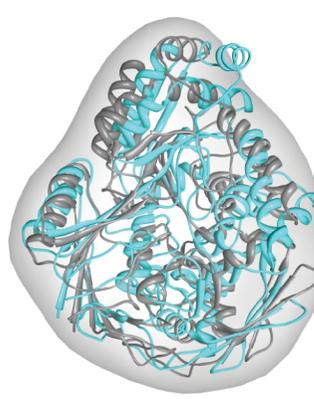
1dpeA, 15 Å



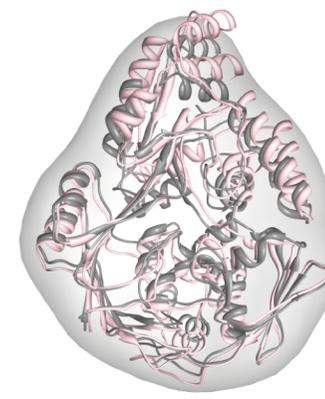
RMSD from target:



12.28 Å



7.00 Å

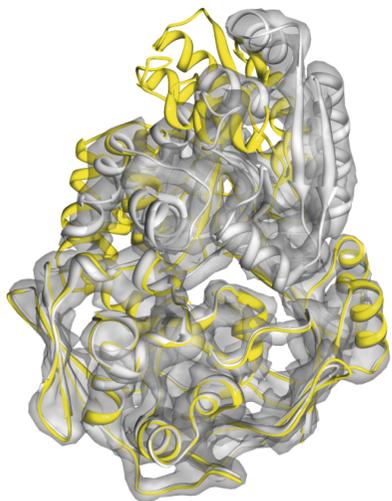


4.17 Å

# Hierarchical refinement

1dpe, 5 Å

Initial



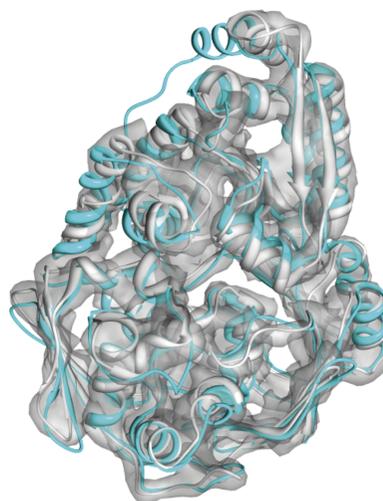
**RMSD from native:**

**12.28 Å**

**CCC:**

**0.814**

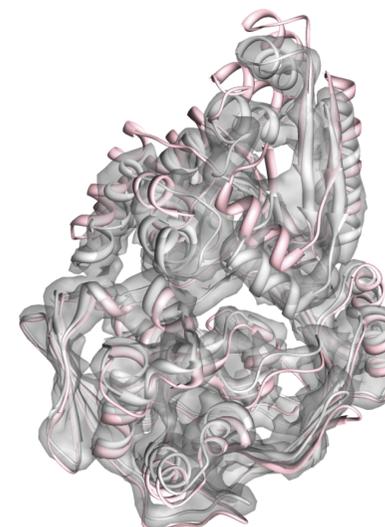
Final un-clustered



**3.69 Å**

**0.891**

Final clustered



**3.92 Å**

**0.867**

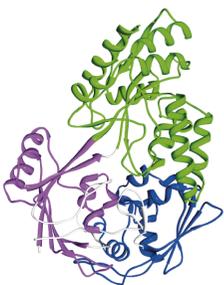
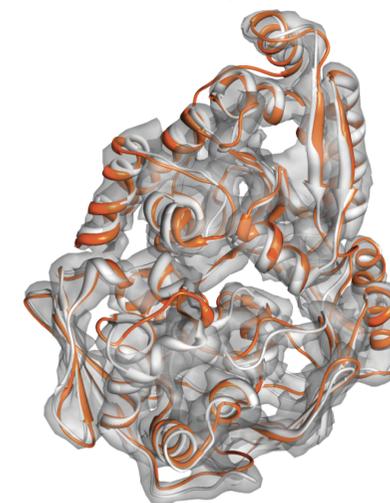
un-cluster



Final two-stage  
refinement

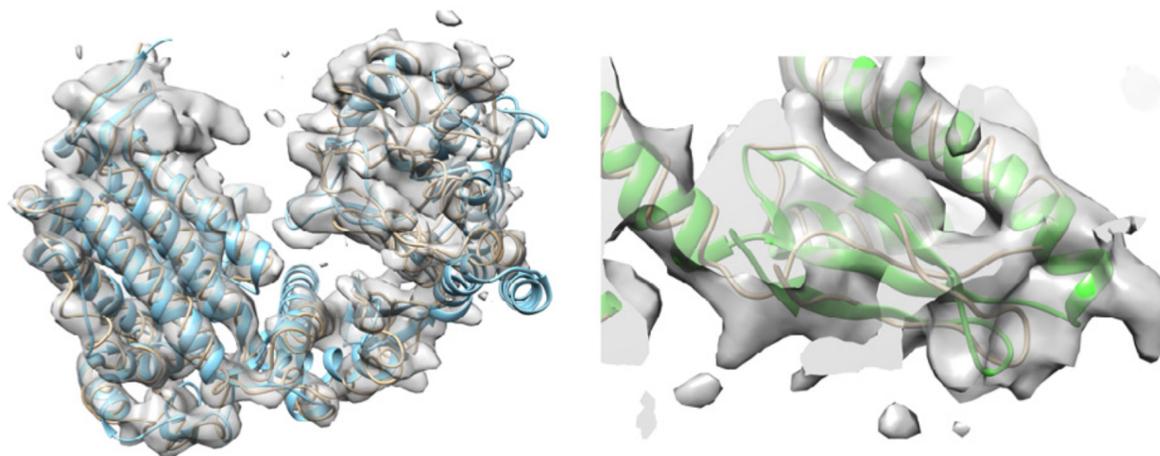
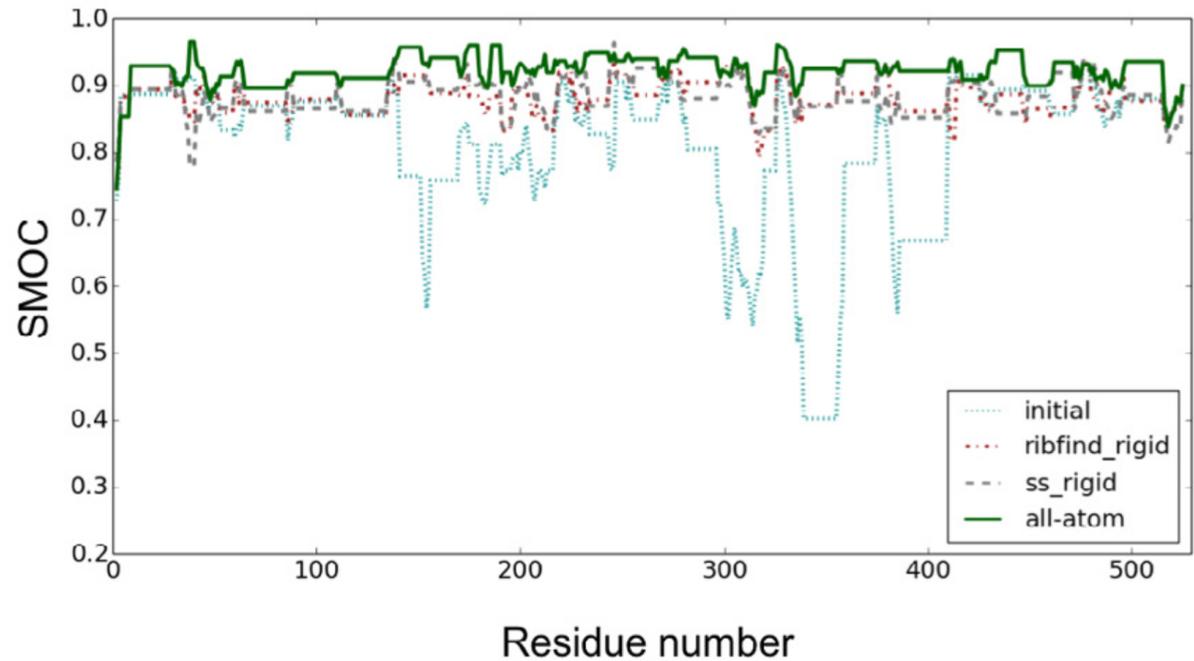
**RMSD: 2.12 Å**

**CCC: 0.92**



# SMOC plots for subnanometer resolution

ADP-bound GroEL (PDB: 4KI8) refined in the density of the unliganded form of GroEL solved at 4.2 Å resolution (EMD-5001).



# Refinement methods

MDFF: Molecular Dynamics (*Trabuco et al. 2008; Singharoy et al. 2016*)

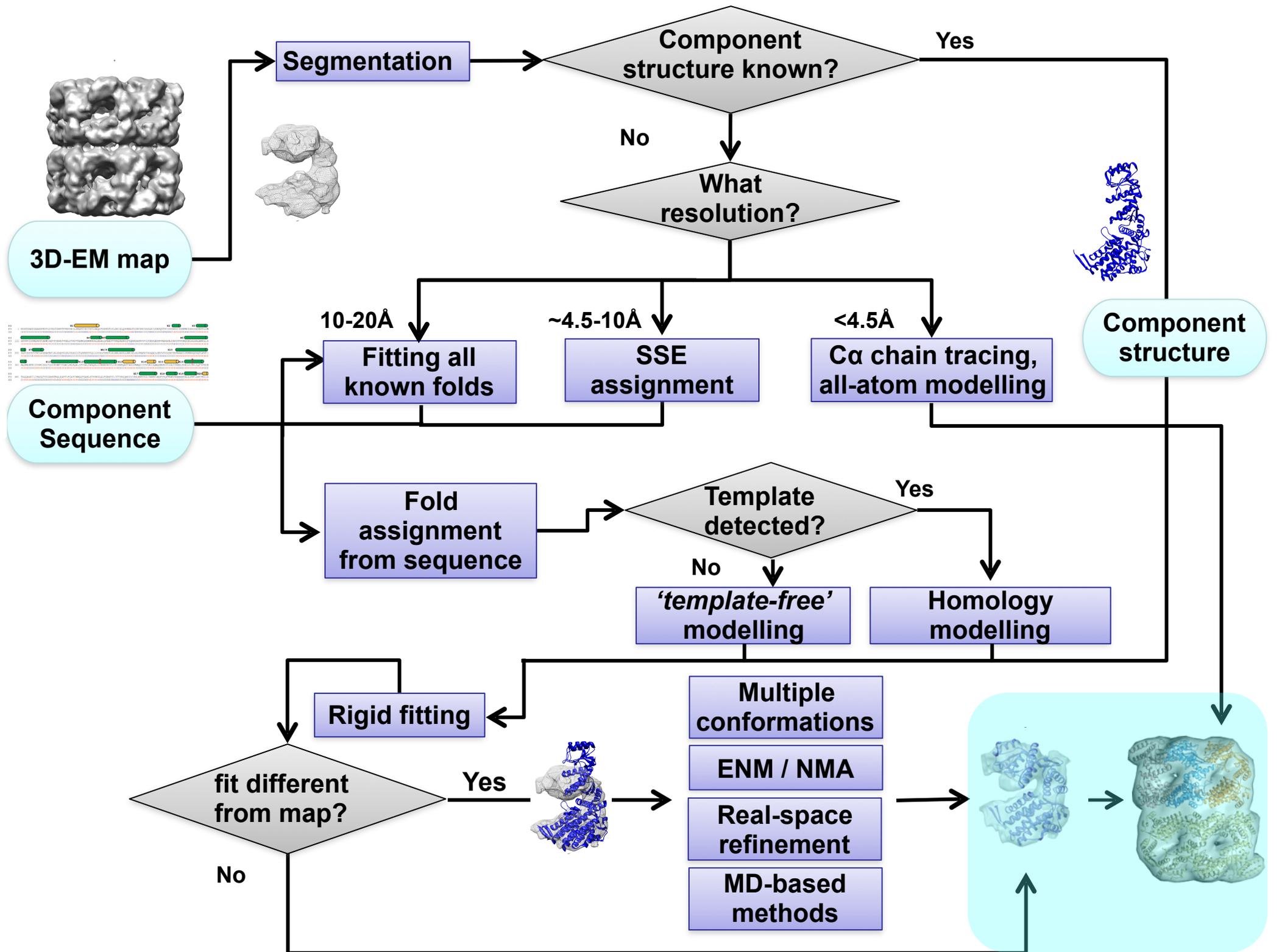
Direx, NMFF, iMODFIT: Normal modes and geometric constraints (*Wang and Schroder 2012; Tama et al. 2004; Blanco and Chacon 2013*)

Rosetta: Monte-Carlo/stochastic (*Wang et al 2016; DiMaio et al. 2015*)

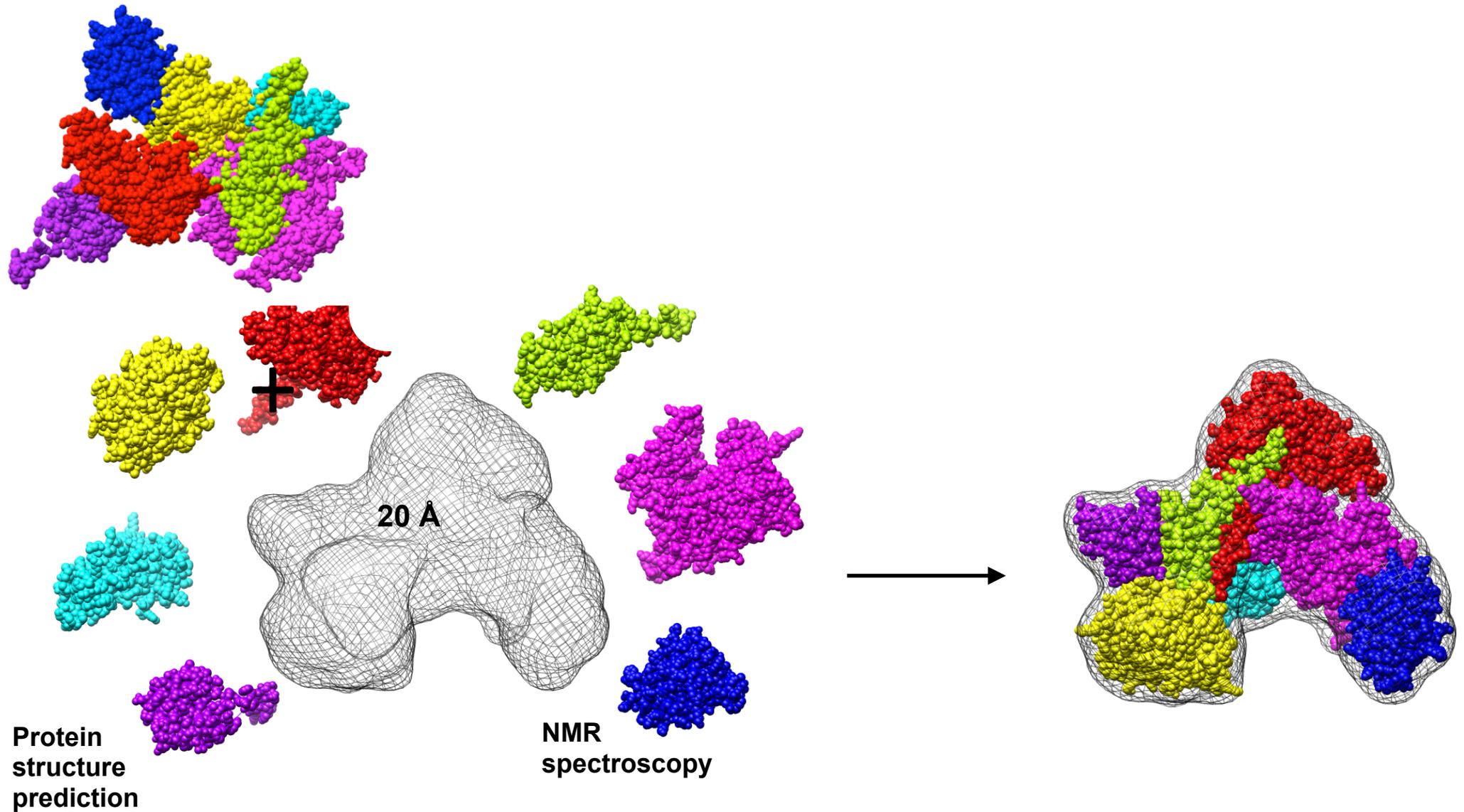
Refmac: Maximum likelihood (*Murshudov 2011; Brown et al. 2015, Nicholls et al. 2018*)

Coot: Interactive/stochastic/exhaustive/gradients (*Emsley et al. 2010; Brown et al. 2015*)

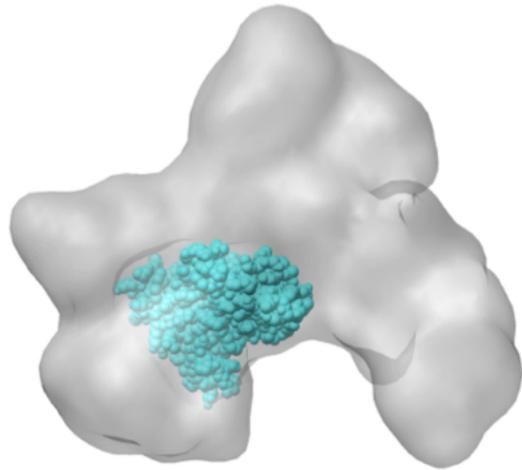
Phenix: Gradient/Simulated annealing MD/exhaustive (*Afonine et al. 2012, 2018*)



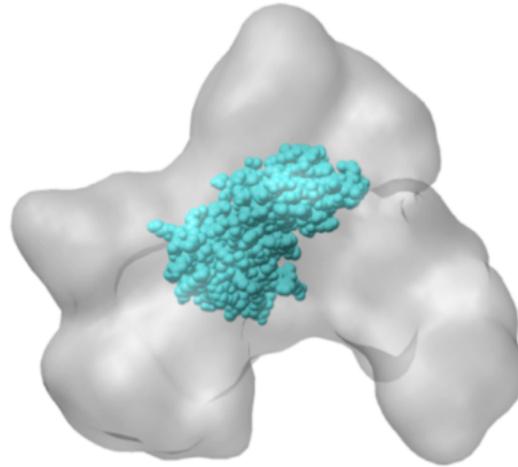
### iii. Assembly (multi-component) fitting



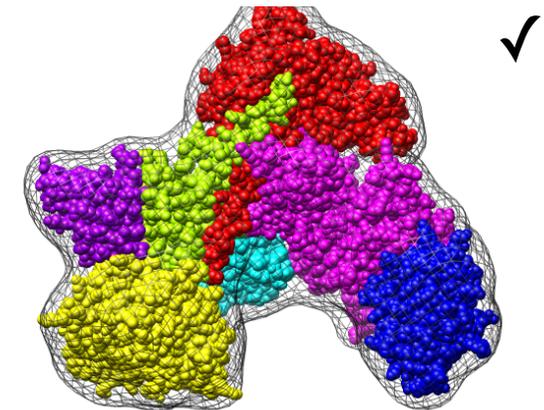
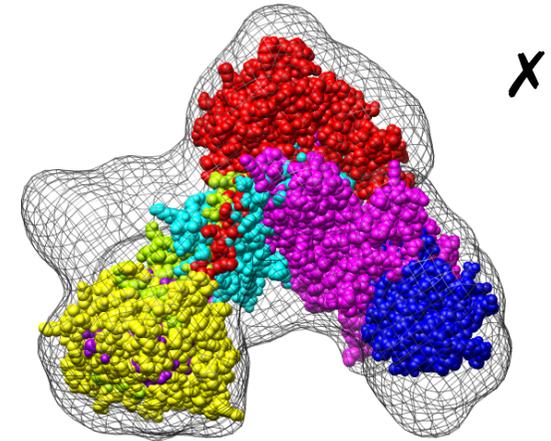
# Sequential fitting



20 Å correct position



20 Å best scoring position



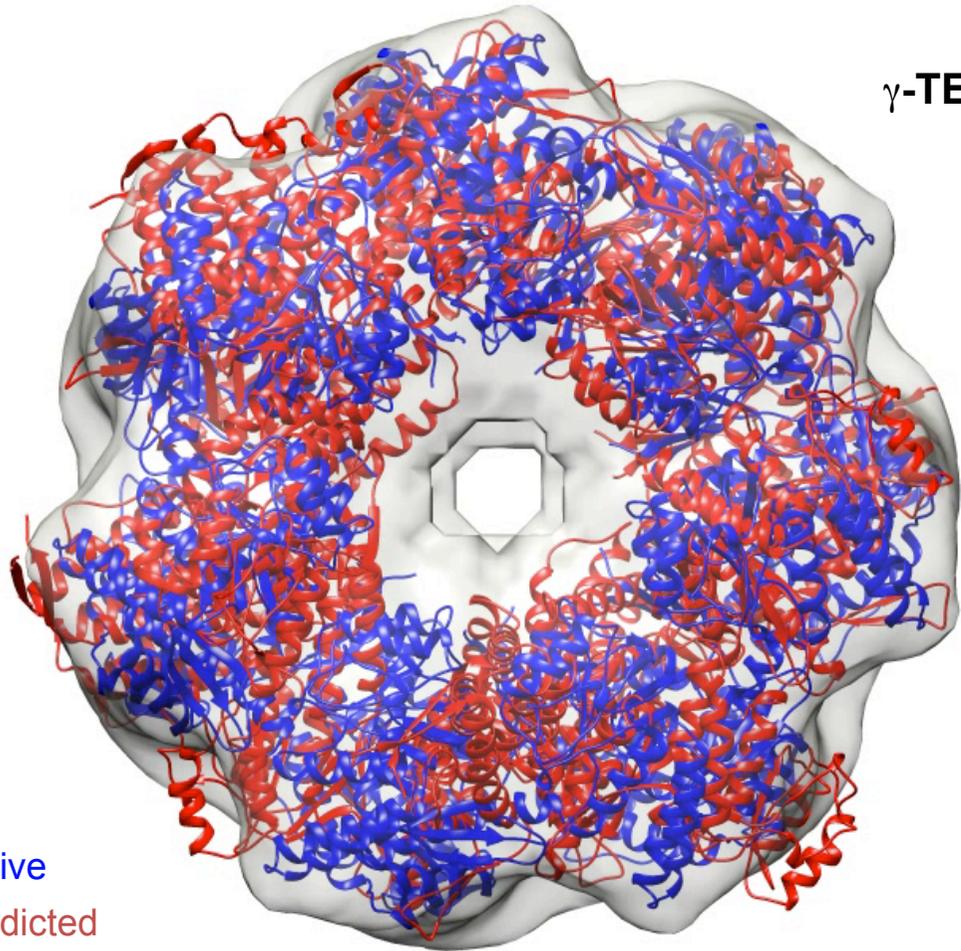
**Problem:** Components may migrate toward the centre of the map or to a different local maxima.

## Solution:

- Core-weighted cross correlation (Wu et al, Zundert & Bonvin 2015)
- Simultaneous (assembly) fitting
- Multiple scores (additional constraints - integrative modelling).

# Multi-component (assembly) fitting

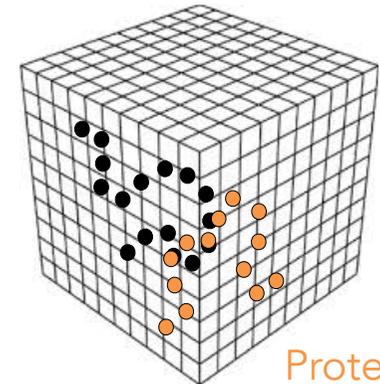
$\gamma$ -TEMPy



Native  
Predicted

1GRU, EMD-1046 (23.5Å)

Protein 1

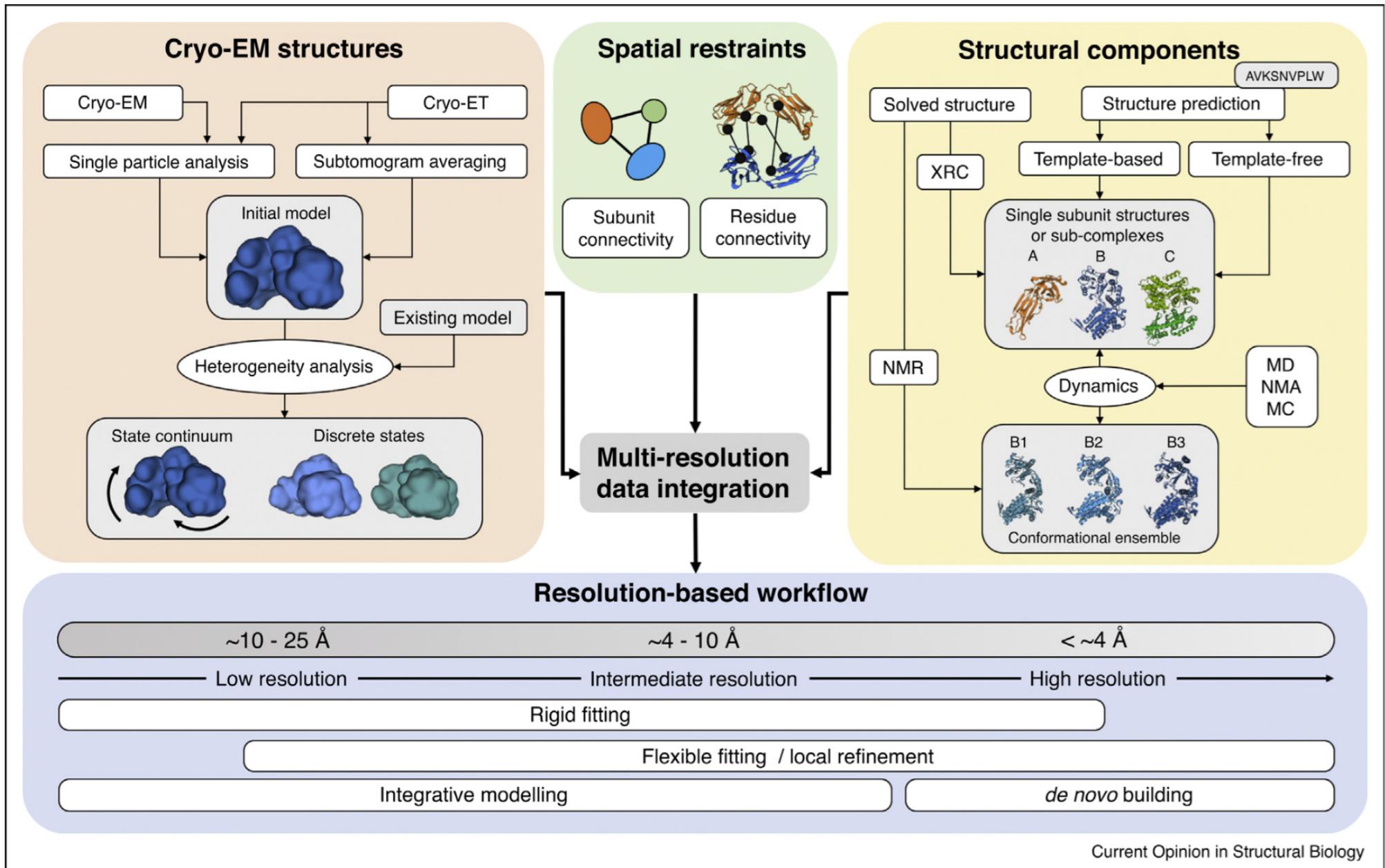


Protein 2

$$PS = \sum_{i=1}^M (\text{Vol}_{\text{Overlap}} / (\text{Vol}_{\text{Component1}} + \text{Vol}_{\text{Component2}}))$$

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

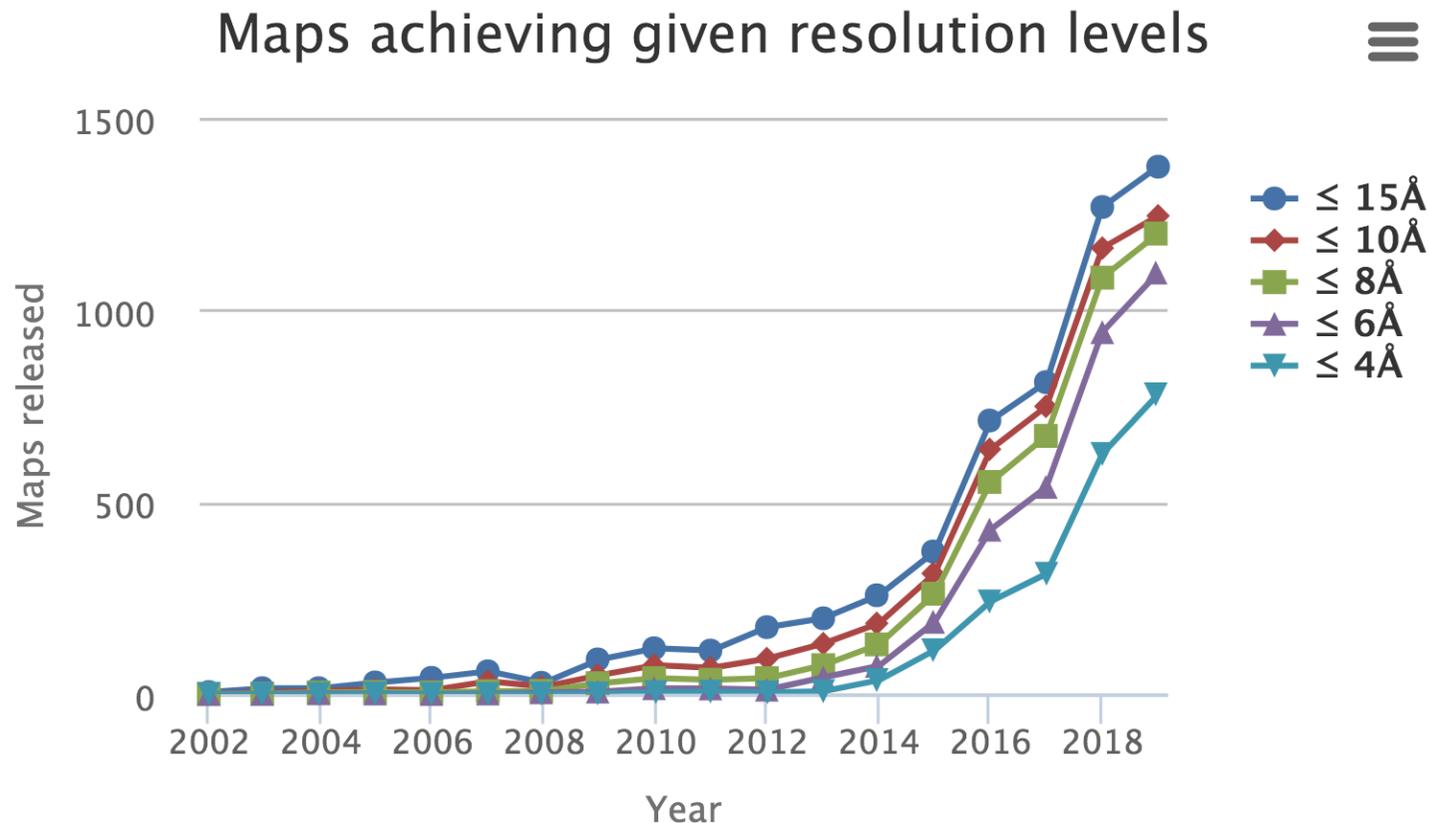
# Determining structures of macromolecular complexes using cryo-EM



# Fit / Model assessment and Validation

# Model Validation

~9000 maps in EMDB.  
~3683 fits in PDB.

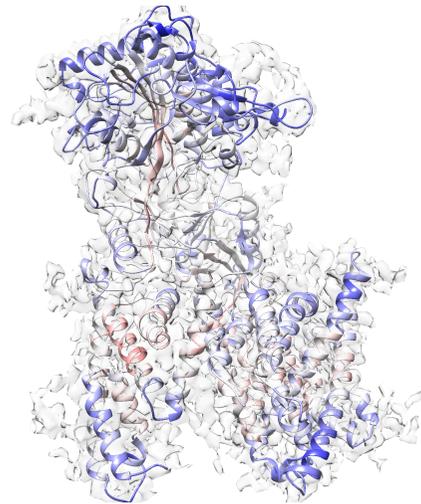


**Models are fitted across resolutions but so far there are no systematic assessment and validation pipeline across resolutions. As the structural and biological interpretation is based on the atomic models, the importance of validation is increasingly being realised.**

# Model fit

# Model geometry

peptide planarity  
backbone torsions (Ramachandran)  
bond lengths  
bond angles  
side chain rotamers



Molprobit: <http://molprobit.biochem.duke.edu/>  
What check: <http://swift.cmbi.ru.nl/gv/whatcheck/>  
PROCHECK: <http://www.ebi.ac.uk/thornton-srv/software/PROCHECK/>

# Model Validation

Many tools available for assessment and validation, but there is no systematic pipeline

## Goodness of fit

TEMPy  
Coot/Refmac  
Phenix  
EMringer  
...

## Secondary structure

MolProbity, Coot, Qmean...  
Pspired, ...

## Validation

Cross-validation:  
Half map (Refmac, Rosetta)  
Ensemble assessment with  
multiple scores (TEMPy)  
Resolution shells (Direx)

## Model geometry

Molprobity  
Coot  
What-check  
...

## Tertiary structure

Verify-3D, ProQ3, Prosa,  
DOPE (MODELLER), ModFold, ..

## Experimental validation

mutations, cross-links, ...