# Cryo-EM Map-Based Modeling

Wah Chiu ([wah@bcm.edu](wah@bcm.edu))


Grigore Pintilie  (pintilie@bcm.edu)

Matthew Baker (mbaker@bcm.edu)
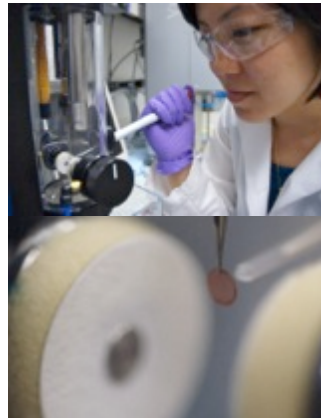

Baylor College of Medicine

NUS Cryo-EM Workshop
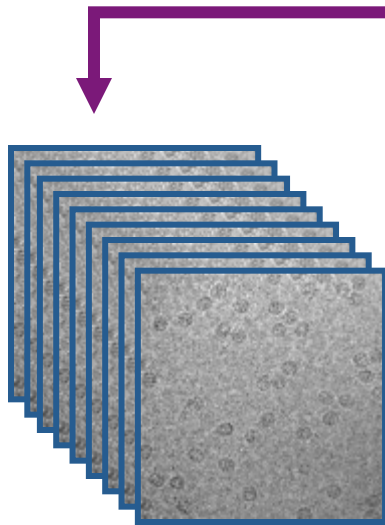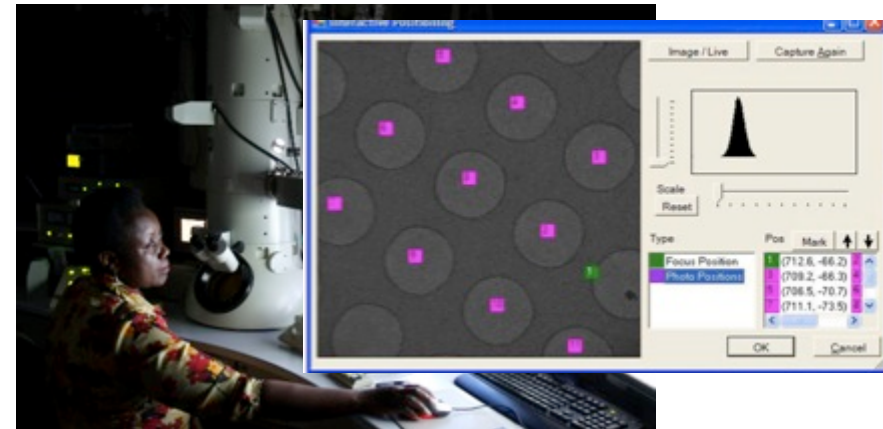
July 12 2012

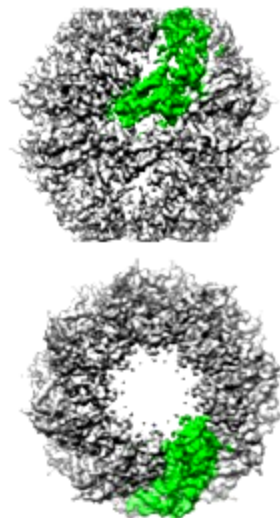# Pipeline in Single Particle Cryo-EM
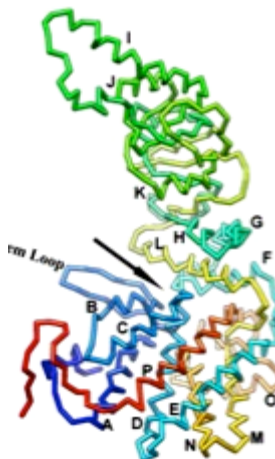


**Biochemical Preparation**

**Cryo-EM Sample Preparation**

**High Resolution Automated Data Collection**
**JADAS**

**Data Archiving & Processing**
**EMEN**

**3D Reconstruction**
**EMAN**

**Model Building & Validation**
**Gorgon**

**Structure Deposition**

# Features as a Function of Resolution



BTV VP3

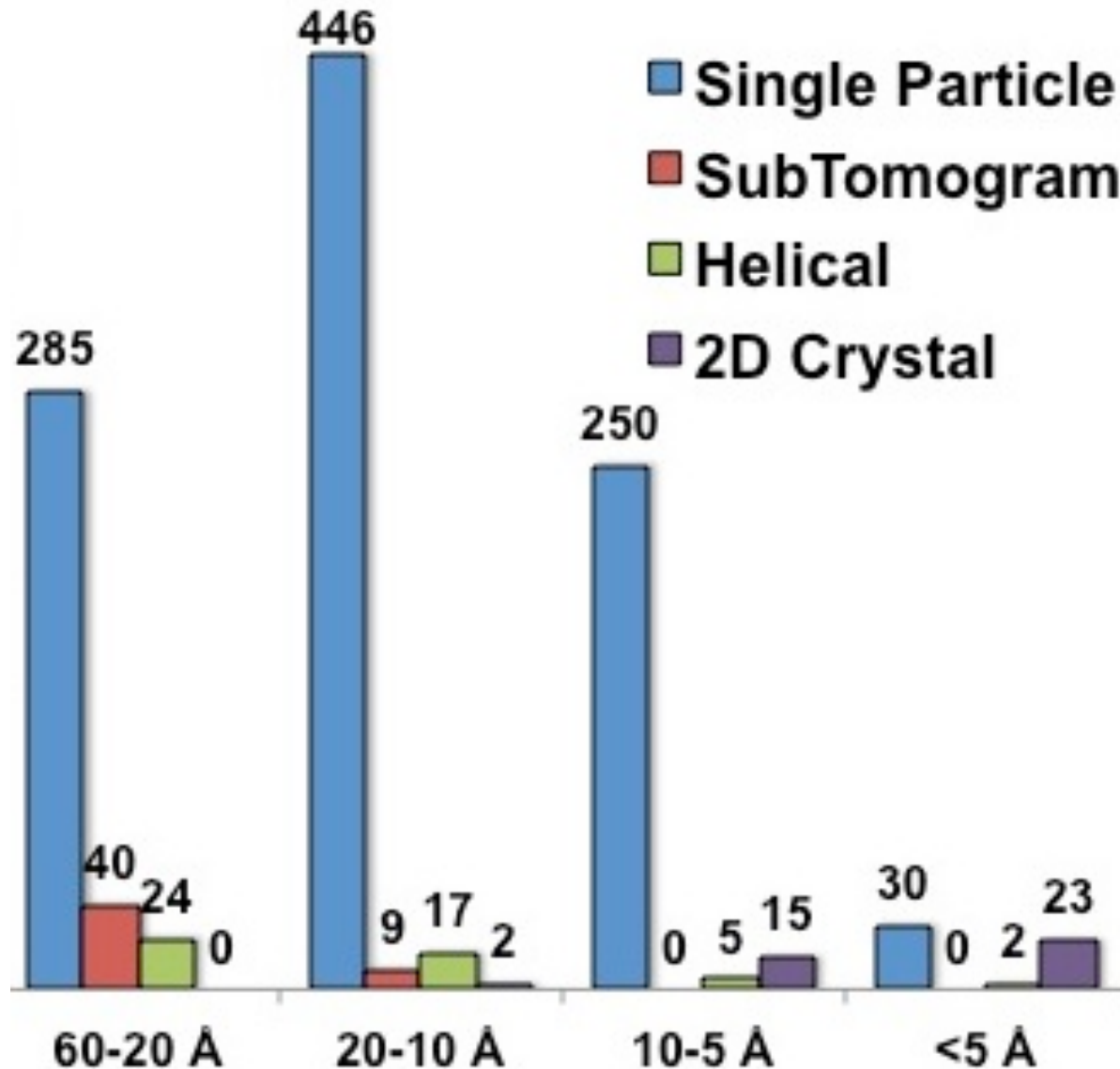| 15+ Å | 9 Å | 6 Å | <4 Å |
|---|---|---|---|
| Size Shape | Domains    α-helices β-sheets | Strands Connectivity | Sidechains |

Courtesy of Matthew Baker

# Structural Features at Different Resolutions



Courtesy of Ryan Rochat

# EMDB Deposited Maps at Different Resolutions



Courtesy of Cathy Lawson

# Cryo-EM Map-Based Modeling

- <u>Segmentation</u>

  Identify locations of molecular components in a complex

- <u>Feature Extractions</u>

  Identify quaternary and secondary structure elements

- <u>Rigid-body Docking</u>

  Dock crystal or homology model into the density map

- <u>Flexible fitting</u>

  Deform model to better fit density map (and to discover new conformations as seen by Cryo-EM)

- <u>*De Novo* modeling</u>

  Build model with no template

# References on Cryo-EM Based Modeling

- Pintilie, G. and Chiu, W. (2012). Comparison of Segger and other methods for segmentation and rigid-body docking of molecular components in Cryo-EM density maps. *Biopolymers* **97**: 742-760.

- Baker, M. L., Baker, M. R., Hryc, C. F., Ju, T. and Chiu, W. (2012). Gorgon and pathwalking: Macromolecular modeling tools for subnanometer resolution density maps. *Biopolymers* **97**:655-668.

- Baker, M. R., Rees, I., Ludtke, S. J., Chiu, W. and Baker, M. L. (2012). Constructing and validating initial Calpha models from subnanometer resolution density maps with pathwalking. *Structure* **20**: 450-463.

- Ludtke, S. J., Lawson, C. L., Kleywegt, G. J., Berman, H. and Chiu, W. (2012). The 2010 cryo-em modeling challenge. *Biopolymers* **97**: 651-654.
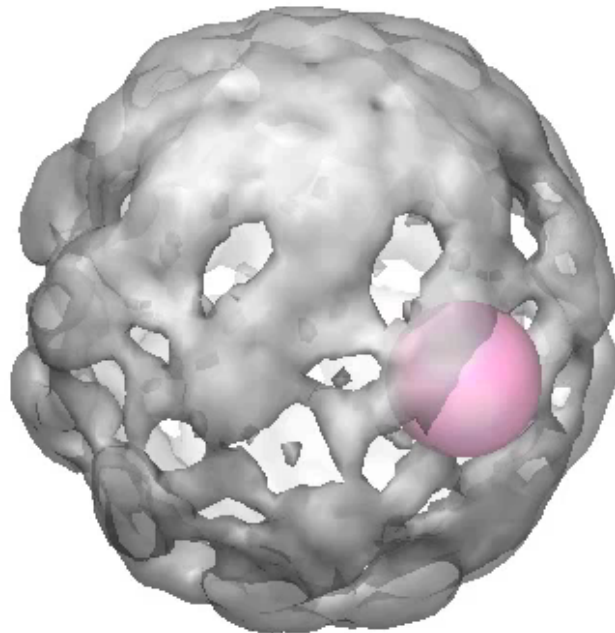
# Segmentation

- Partition map into regions or segments that correspond to individual components such as proteins, subunits, DNA, etc.

- Methods:

  - Manual

    » By using "sphere eraser" in Chimera

    » By drawing boundary contours using Aviso$^{TM}$

  - Using docked model

    » Take all grid points close to atoms in a docked model

  - Semi-automated

    » **EMAN** (http://blake.bcm.edu/emanwiki/Segment3D)

    » **Segger** (http://ncmi.bcm.edu/ncmi/software/segger/docs)

    » **Volrover** (http://www.cs.utexas.edu/~bajaj/cvc/software/volrover.shtml)

# Manual Segmentation – Sphere Eraser
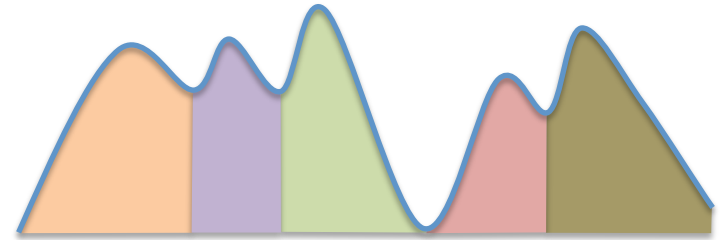
(http://www.cgl.ucsf.edu/chimera/)

- Erase parts of the map that are not in the region that corresponds to a component based on human judgment
- Very tedious, error-prone, and highly subjective
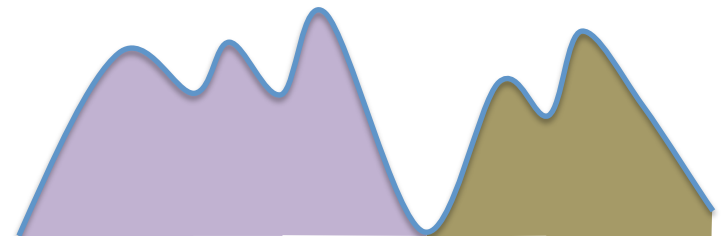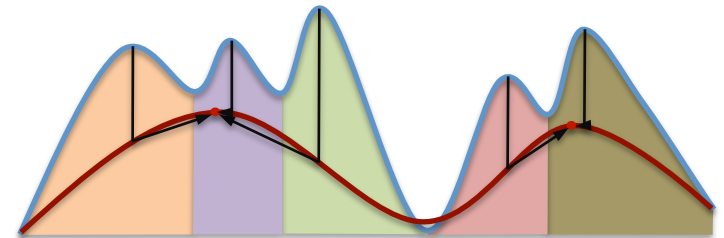
# Semi-Automatic Segmentation - Segger

## Example: 1-dimensional map
- Height proportional to density
- Watershed region produces 5 regions
- Each region corresponds to a peak in the density
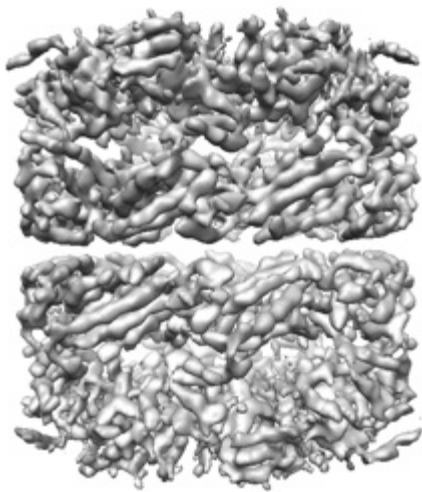
## Smoothing and Grouping
- Map is smoothed using Gaussian filter (red line)
- Regions are grouped based on which peak is reached after uphill climb (black lines)
- The result in this example is 2 regions, since there are two peaks in the smoothed map
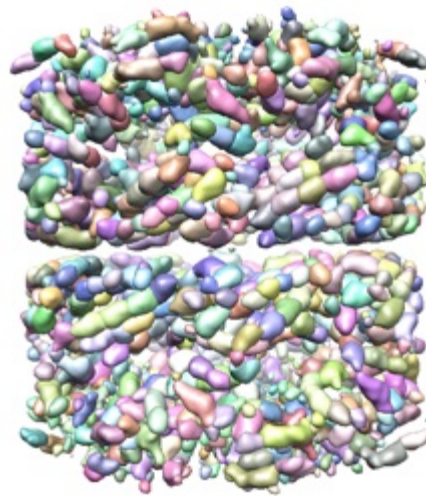


Pintilie et al (2010) *J Struct Biol* **170**(3):427-38.
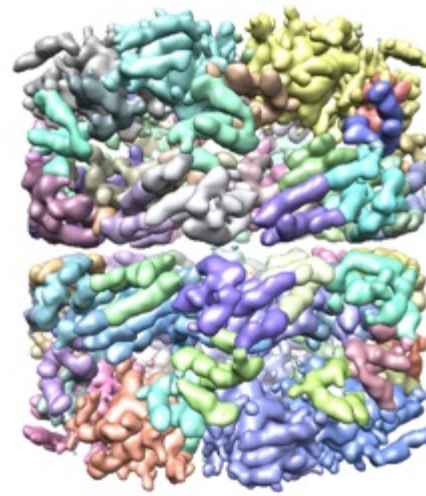
# Segmentation – Segger
# A Case that Works

- CryoEM map of GroEL @4Å resolution
- After 3 steps of size 7 (Å), 14 regions are produced; each region corresponds to a protein
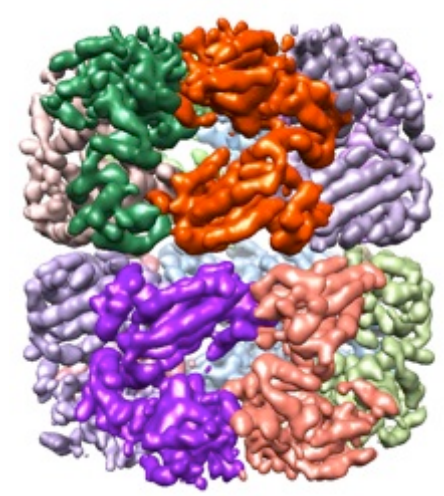


Cryo-EM map
EMDB:5001
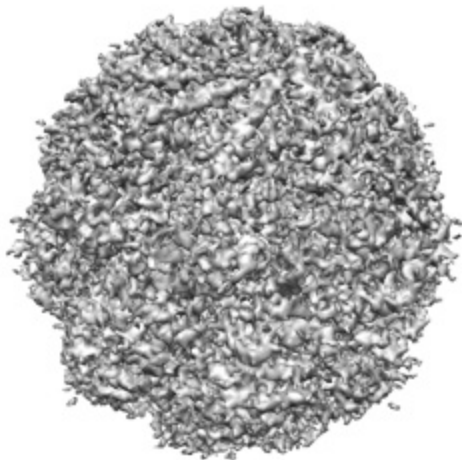GroEL @ 4Å

2124 watershed regions
(too many)

After smoothing and
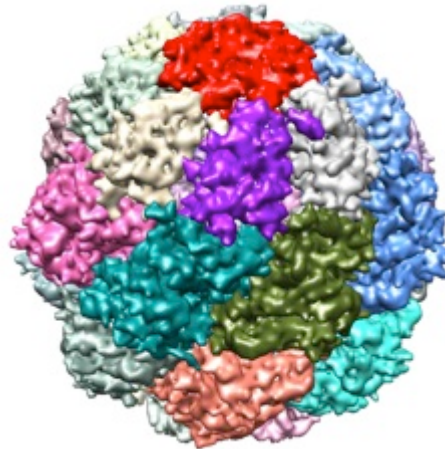grouping 1 step (42
regions)

After smoothing and
grouping 3 steps (14
regions)

Pintilie et al (2010). *J Struct Biol* **170**(3):427-38.

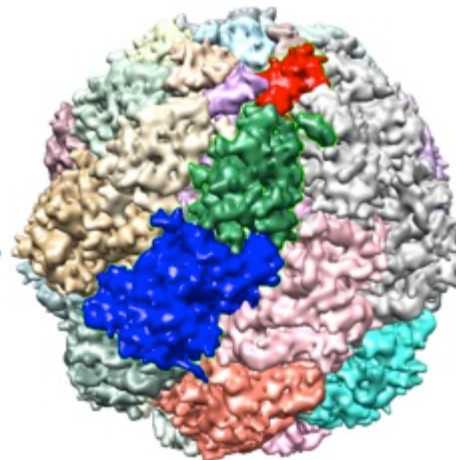# User-Assisted Semi-Automatic Segmentation - Segger

- In some cases, manual (subjective) grouping is needed as shown in the example below
- Use prior knowledge, or known homology models, for guidance
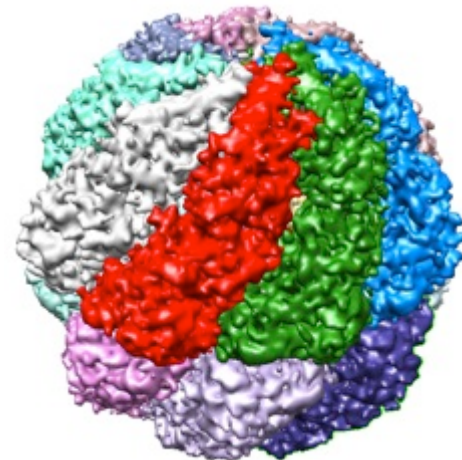


Cryo-EM map
EMDB:5137
Mm-cpn @ 4.3Å

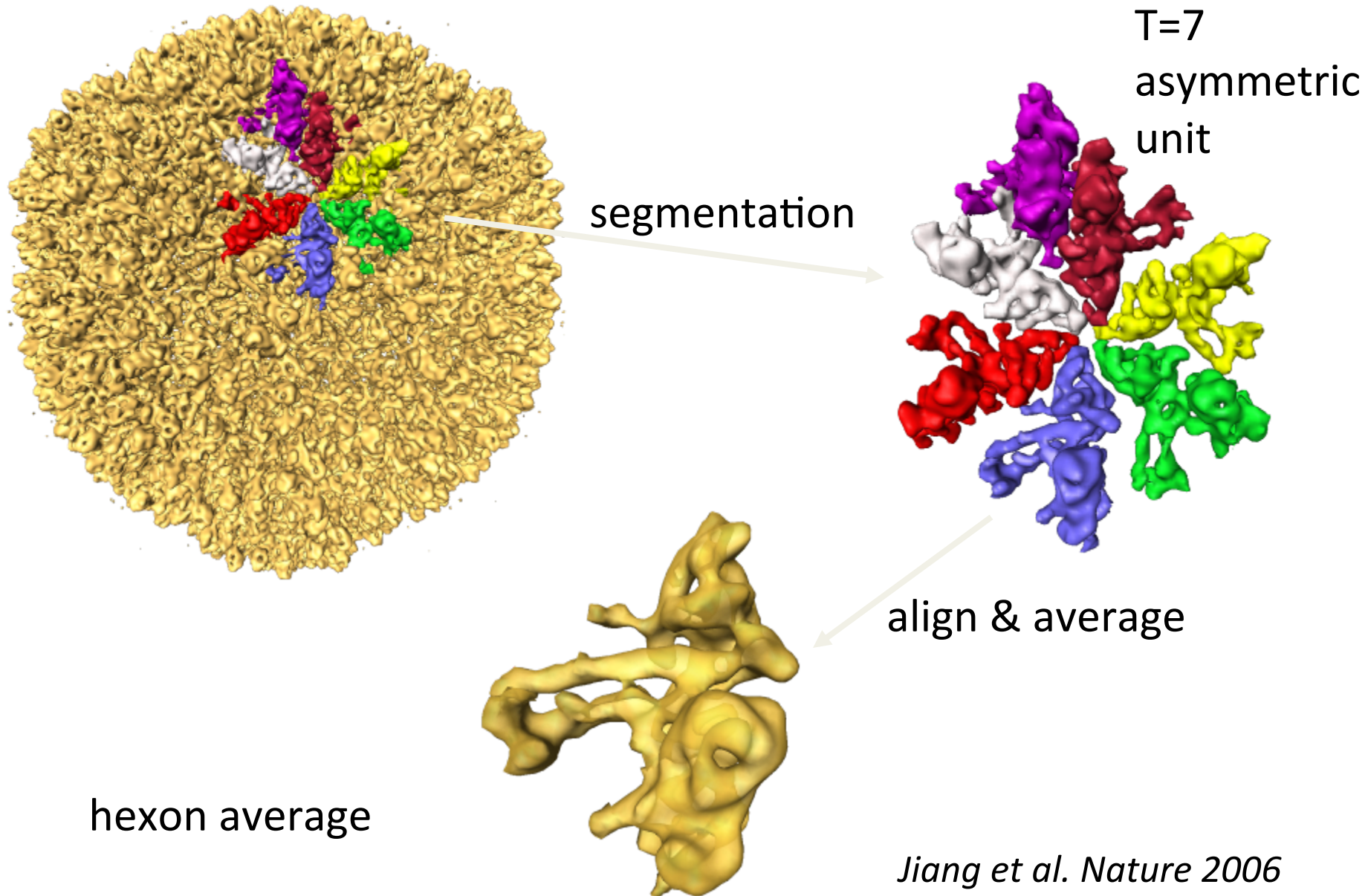Regions after smoothing and grouping – note that red region is incorrect as it spans multiple proteins

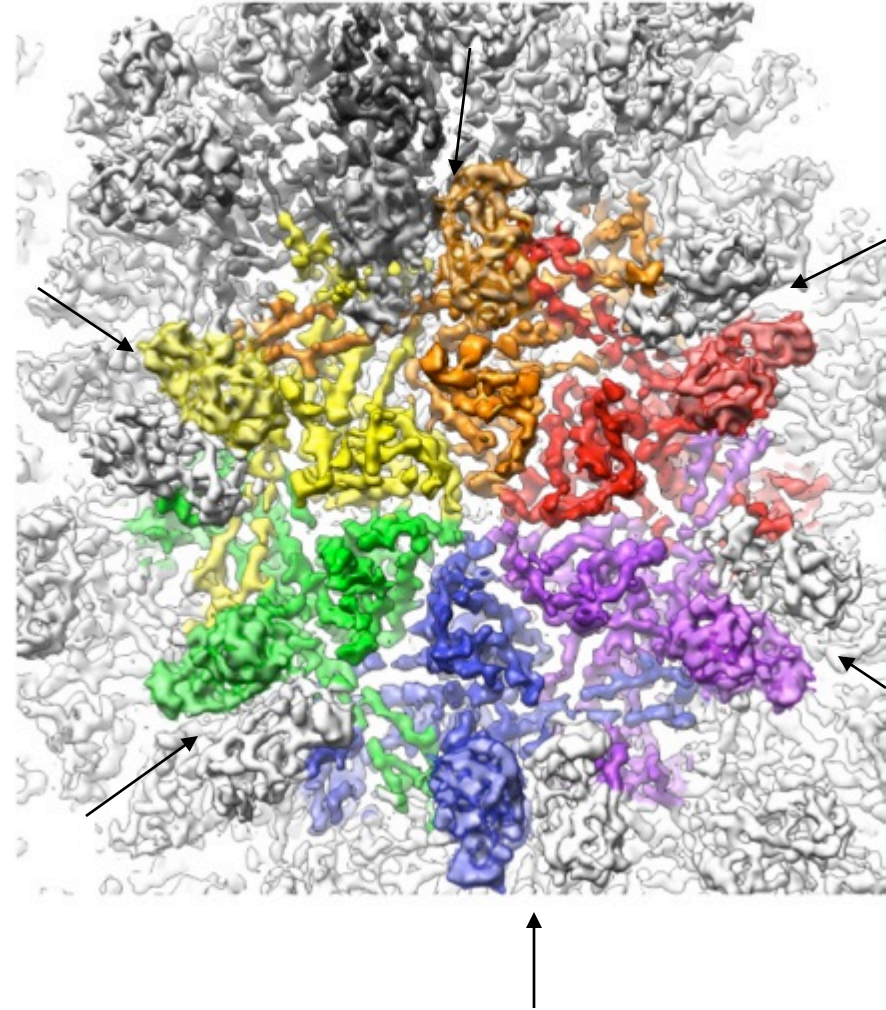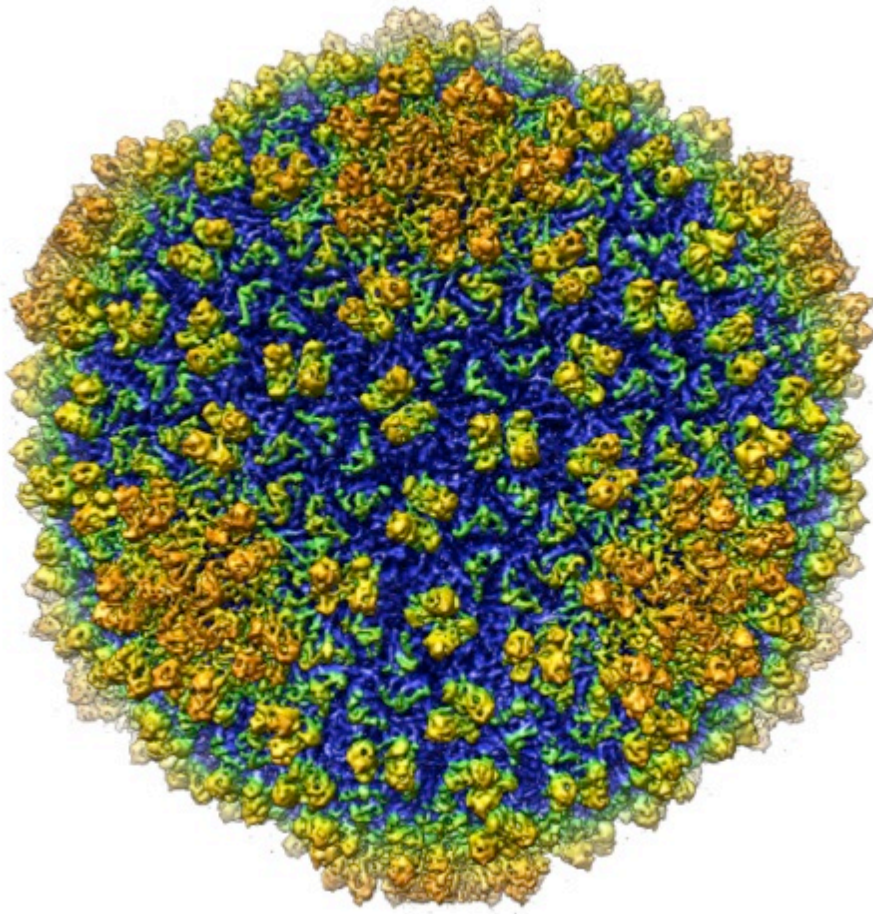Regions can be selected and grouped interactively (the blue, green and red regions are selected above)

The correct segmentation, after manual grouping, contains 16 regions

# Epsilon15 Phage at 9.5 Å Resolution



T=7 asymmetric unit

segmentation

align & average

hexon average

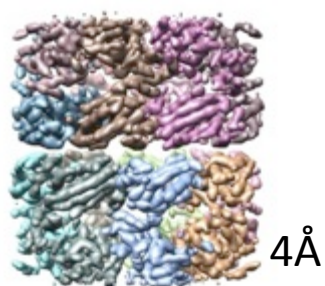*Jiang et al. Nature 2006*

# Epsilon15 Phage at 4.5 Å Resolution



Jiang et al (2008) *Nature* **451**: 1130-4.
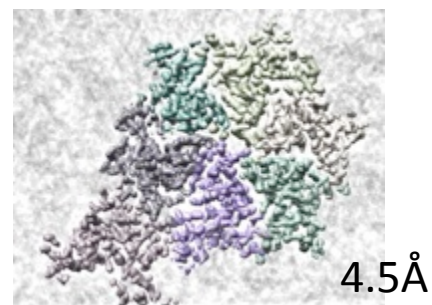
# 2010 Cryo-EM Challenge
## Segger Segmentation Results

# Segmentation Summary

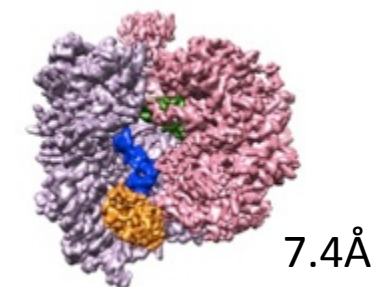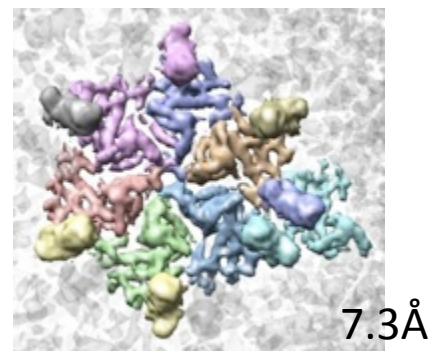- **Segmentation is a difficult process. Why?**
  - Proteins have complex 3D shapes and contacts with other proteins
  - Complexes can consist of large numbers of proteins
  - Contact boundaries between different proteins can be as dense as intra-protein contacts; hence they are hard to identify accurately

- **In general:**
  - It tends to be easier in density maps where components are well-separated (e.g. GroEL)
  - It tends to be hard in density maps where there are many contacts between proteins, and proteins have long, narrow segments that contact adjacent proteins (e.g. in phages)
  - Knowing what the protein looks like (e.g. from homology model) helps a great deal

# Feature Extraction

- Quaternary Structure
- Domains
- Helices and sheets

# Herpes Simplex Virus-1 Capsid

- 8.5 Å resolution cryo-EM map (Zhou, *Science* 2000)

- 4 structural proteins
  - VP5: major capsid protein forms pentons and hexons
  - VP26: binds only VP5 at hexon specific positions
  - VP23 and VP19C form triplexes between hexons and pentons

# 8.5 Å Map of HSV-1 Capsid



P

C   E

Manually identified helices

Zhou et al (2000) *Science* **288**: 877-880.

# SSEHunter of a Cryo-EM Density



Helix

Sheet

Skeleton

6.8 Å resolution cryoEM density map
*Z.H. Zhou et al, 2001, Nature Struct Biol.*

X-ray model

M. Baker et al. (2007) *Structure* **15**:7-19

# Skeleton: Feature and Topology

- Compact geometric representation of a volume

- Feature preserving
  – Sheets are represented as flat surfaces
  – Helices and loops are represented as curves

- Topology preserving
  – Maintains density connectivity while minimizing number of branches and breaks

Baker, M.L., Ju, T. and Chiu, W. (2007) *Structure*, **15**:7-19.

# 9.5 Å Cryo-EM Map of RyR1 (2.2 MDa)



500 Å

Ludtke et al (2005) *Structure* **13**: 1203-11.

# Identified SSE in RyR1 Cryo-EM Map



**S6** **S5** **S4** **S3** **S2** **S1**

**central part**

**S7**

**CY region**

**TM region**

41 α-helices
(**8**+**23**+**5**+**5**)

7 β sheets

Serysheva et al (2008) *Proc Natl Acad Sci U S A* **105**: 9610-5.

# References on Fitting methods

- Wriggers W, Milligan RA, McCammon JA. Situs: A Package for Docking Crystal Structures into Low-Resolution Maps from Electron Microscopy. *J Struct Biol* 1999;125:185–195. [PubMed: 10222274]

- Roseman AM. Docking structures of domains into maps from cryo-electron microscopy using local correlation. *Acta Crystallogr D Biol Crystallogr* 2000;56:1332–40. [PubMed: 10998630]

- Rossmann MG, Bernal R, Pletnev SV. Combining Electron Microscopic with X-Ray Crystallographic Structures. *J Struct Biol* 2001;136:190–200. [PubMed: 12051899]

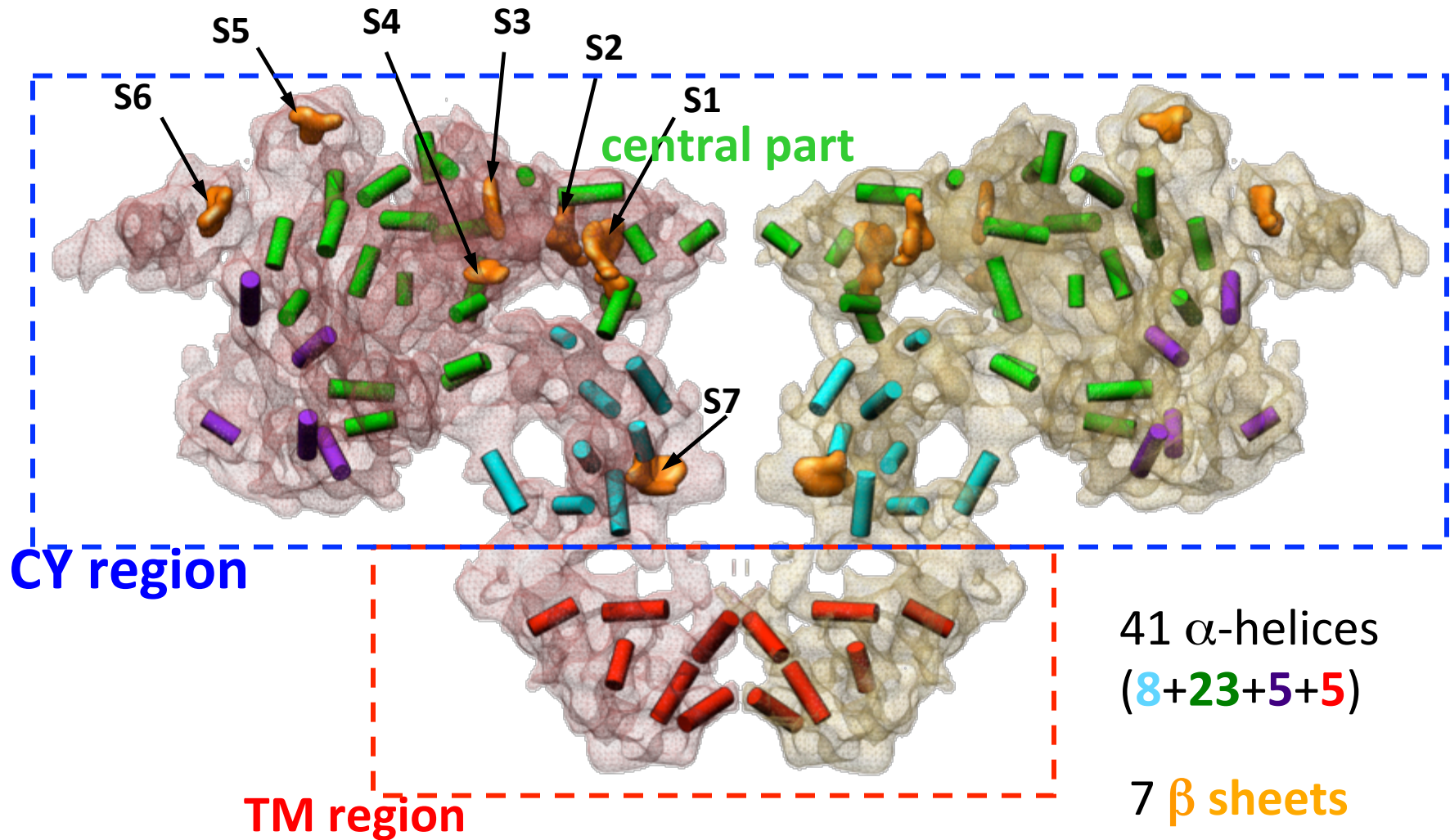- Jiang W, Baker M, Ludtke S, Chiu W. Bridging the information gap:Computational tools for intermediate resolution structure interpretation. *J Mol Biol* 2001;308:1033–1044. [PubMed: 11352589]

- Topf, M., Baker, M. L., John B., Chiu, W. and Sali, A. (2005). Structure characterization of components of protein assemblies by comparative modeling and electron cryo-microscopy. *J Struct Biol* **149**: 191-203.

- Pintilie, GD, Zhang, J, Goddard, TD, Chiu, W, and Gossard, DC (2010). Quantitative analysis of cryo-EM density map segmentation by watershed and scale-space filtering and fitting of structures by alignment to regions. *J Struct Biol* **170**:427-38.
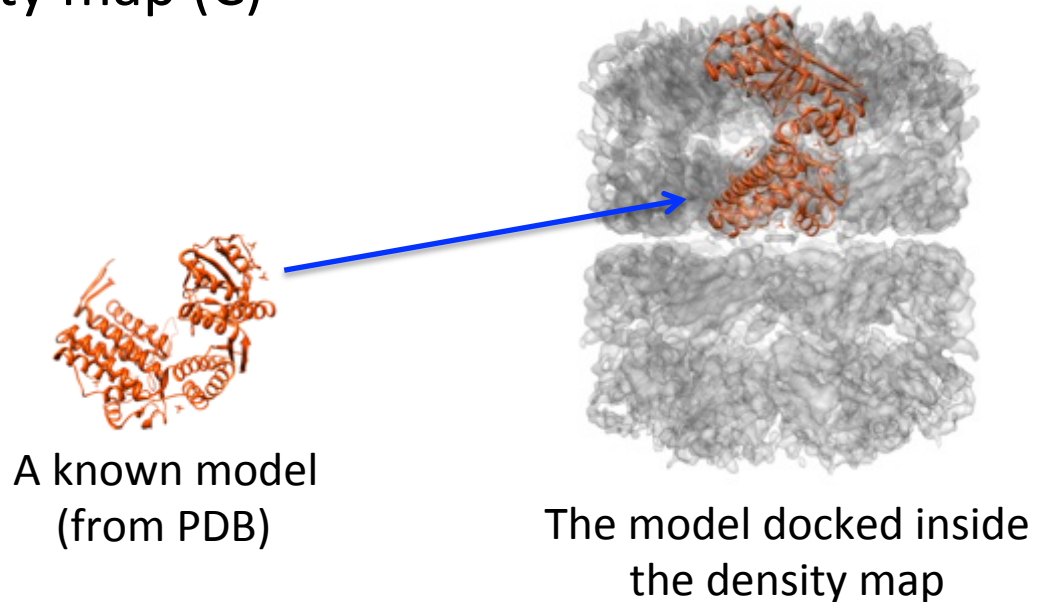
# Rigid-body Docking

- What is rigid-body docking:
  - Finding the translation and orientation that takes a known model and places it inside a density map so that it accurately overlaps the same component as it occurs in the density map
- Quantitatively speaking, we try to find the translation and orientation that maximizes the cross-correlation score between:
  - a simulated density map of the structure (S),
  - and the cryo-em density map (C)

$$cc = \frac{\displaystyle\sum_{i,j,k \in (n_1 \times n_2 \times n_3)} d^C d^S_{i,j,k}}{\displaystyle\sum_{i,j,k \in (n_1 \times n_2 \times n_3)} d^C \sum_{i,j,k \in (n_1 \times n_2 \times n_3)} d^S_{i,j,k}}$$

Cross-correlation score
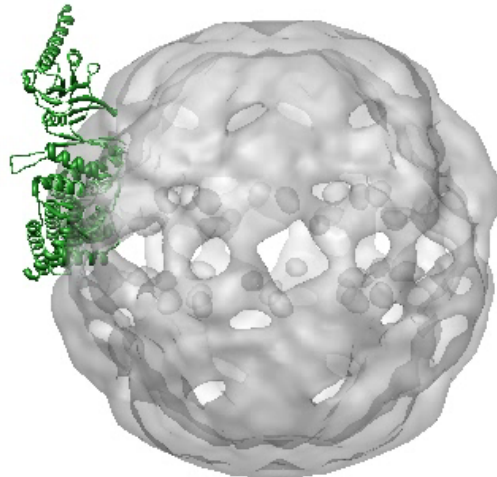
$d^S_{i,j,k}$  Density in S at grid point i,j,k

$d^C$  (Interpolated) Density in C at grid point i,j,k in S



A known model
(from PDB)



The model docked inside
the density map
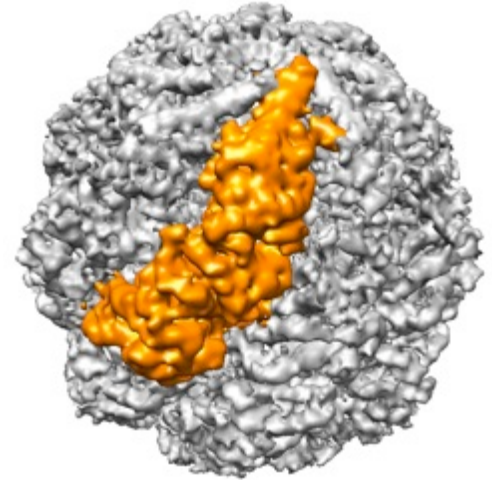
# Rigid-body Docking I

Exhaustive search:

- Used by Situs, Foldhunter, ADP-EM

- The structure is placed at evenly-spaced positions and orientations

- The cross-correlation score is computed for each position/orientation

- This can take a long time in large maps (hours), but it can be accelerated using the Fourier transform, or by using multiple processors
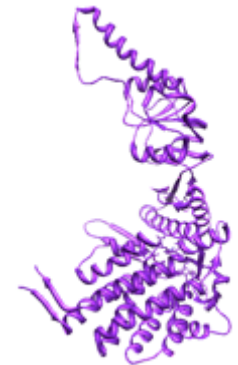
# Rigid-body Docking II

## Segger:

- Performs rigid-body docking by aligning a model to a segmented region

- Faster than exhaustive search, however user-guidance is required

- Steps:
  1. Segment map
  2. Choose a region and a model
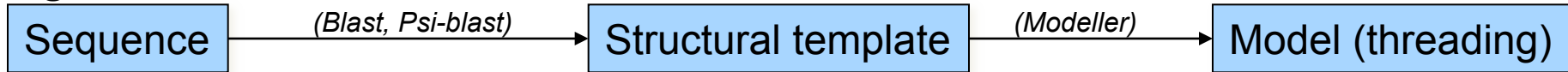  3. Align model to region using one of two methods
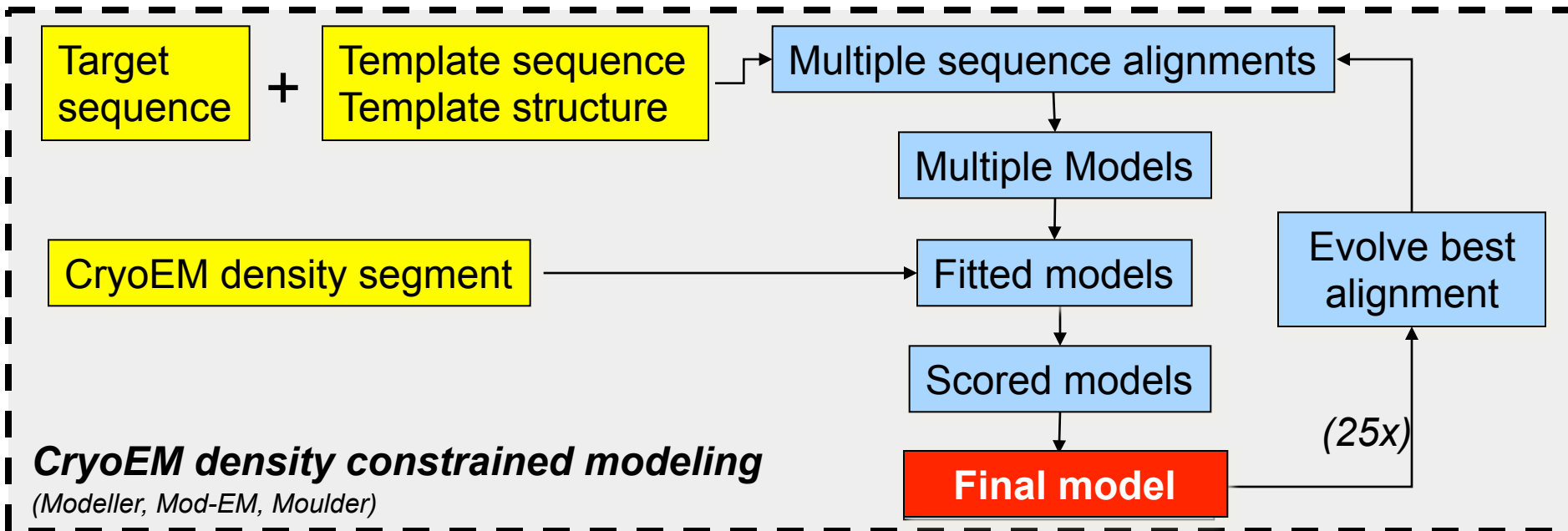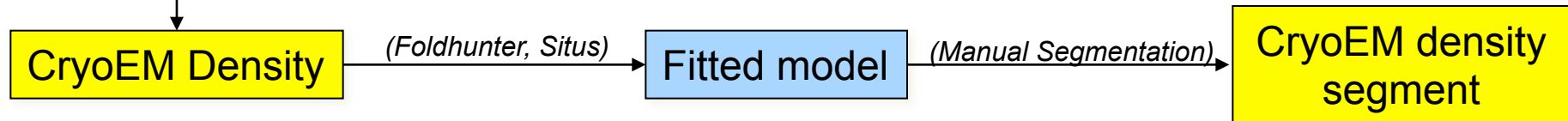


Mm-cpn @4.2Å
EMDB:5137



PDB:3kfb, chain A

# CryoEM Restrained Comparative Modeling

**Target identification**

| Sequence | *(Blast, Psi-blast)* → | Structural template | *(Modeller)* → | Model (threading) |

**Model localization**

| CryoEM Density | *(Foldhunter, Situs)* → | Fitted model | *(Manual Segmentation)* → | CryoEM density segment |

**Target sequence** + **Template sequence Template structure** → Multiple sequence alignments

Multiple Models

CryoEM density segment → Fitted models

Scored models

Evolve best alignment

*(25x)*

**Final model**

*CryoEM density constrained modeling*
*(Modeller, Mod-EM, Moulder)*

Topf et al, *JMB* (2006)

# CryoEM Restrained Modeling
# of the N-terminal of RyR1



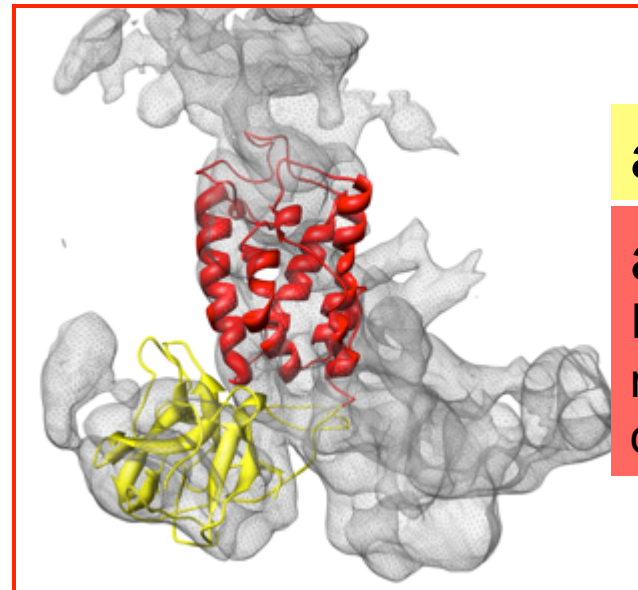N terminal Antibody binding

SSEhunter results:

**α-helices**

β-sheets

**aa 215-420**

**aa 418-564**
IP3-binding region (1N4K) of IP3 Receptor

**aa 12-207**
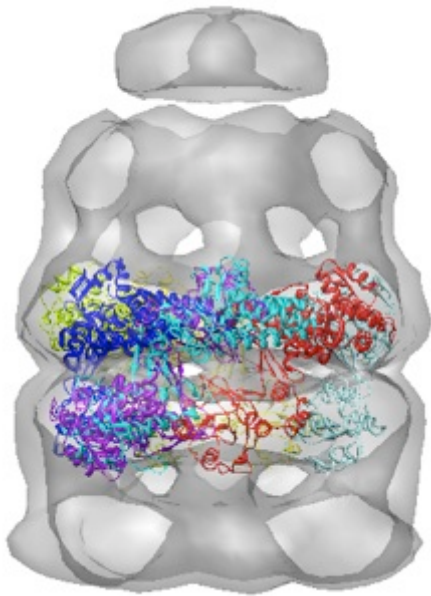IP3R-supressor binding region (1XZZ)

# Rigid-body Docking

- Which alignment, produced by exhaustive search, or by Segger, or by other methods, is the right fit?

- It is typical to pick the fit with the highest score

  ▪ Which score to use?

    ○ Cross-correlation tends to be the most reliable
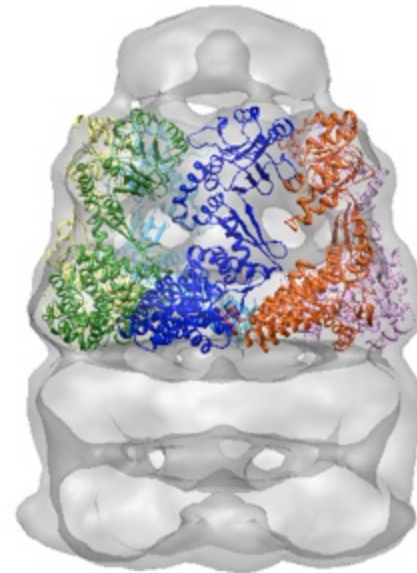    However this can fail, especially at lower resolutions

Vasishtan and Topf (2011). *J Struct Biol* 174:333-343.

# Challenges in Rigid-body Docking





**GroEL @23Å resolution**
Alignments with highest cross-correlation are incorrect – they overlap the middle part of the map, which has higher densities, and hence give higher cross-correlation scores

The correct docking results are as shown. Can we compute other scores, or compute a confidence level in our docking results?

# Rigid-body Docking

Q: How confident can we be  that the alignment with the highest score is the correct fit?

A: Statistically, if the highest score is higher than scores of other alignments, this means the fit is significant.

- o *z-score:* indicates how much higher the score for the best alignment is compared to the mean score of the other alignments
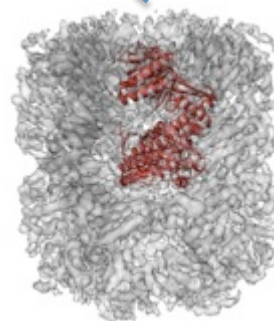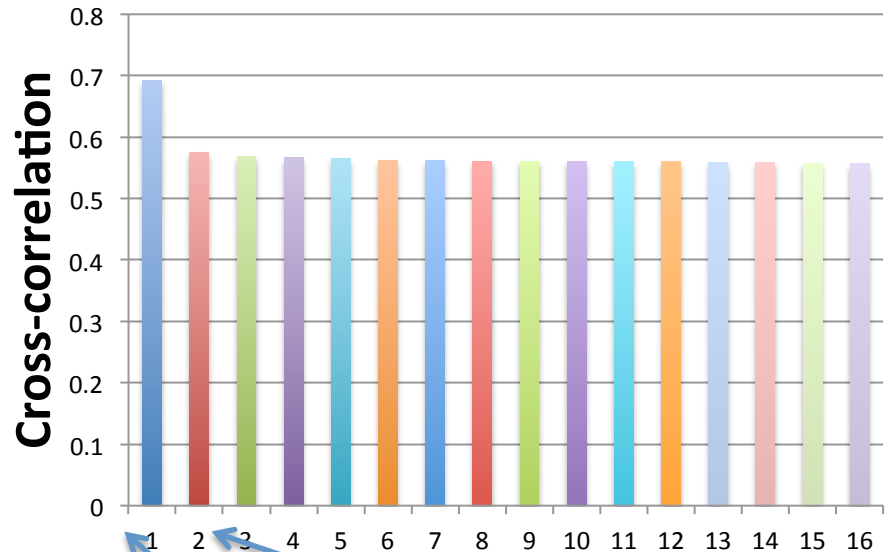
# Score Rigid-body Docking

- *z-score:*
  - How much higher is the best score compared to the mean of the other scores?
  - More precisely: how many standard deviations is the highest score above the mean?

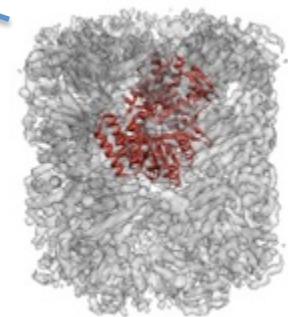$$z - score = \frac{S_1 - mean(S_{2..n})}{stdev(S_{2..n})}$$

$S_1$  Top score

$avg(S_{2..n})$  Mean of scores 2..n

$stdev(S_{2..n})$  Standard deviation of scores 2..n
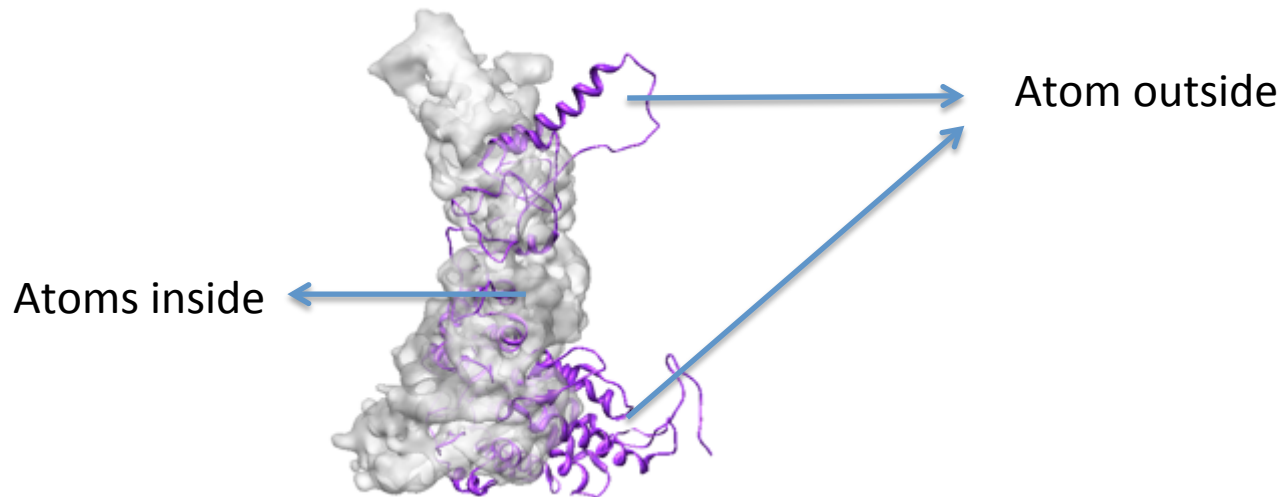


Correct fit
(highest CC score)

Incorrect fit
(significantly lower CC score)

# Score Rigid-body Docking

- Other scores can be computed, to evaluate how confident we can be about a rigid-body docking result:

  - Atom inclusion:
    - How many atoms are inside the observed density?
    - If many atoms are outside, then the model may not be the best match for the map, or it may be incorrectly docked.

  - Density occupancy:
    - How many grid points with high density values are occupied by atoms
    - If many areas inside the density maps are un-occupied, then the docked models may not fully explain the experimental density map.

  - Clashes with other docked models, or with symmetric copies
    - Do multiple docked models clash
    - If there are many clashes, this could signify an unreliable docking result.
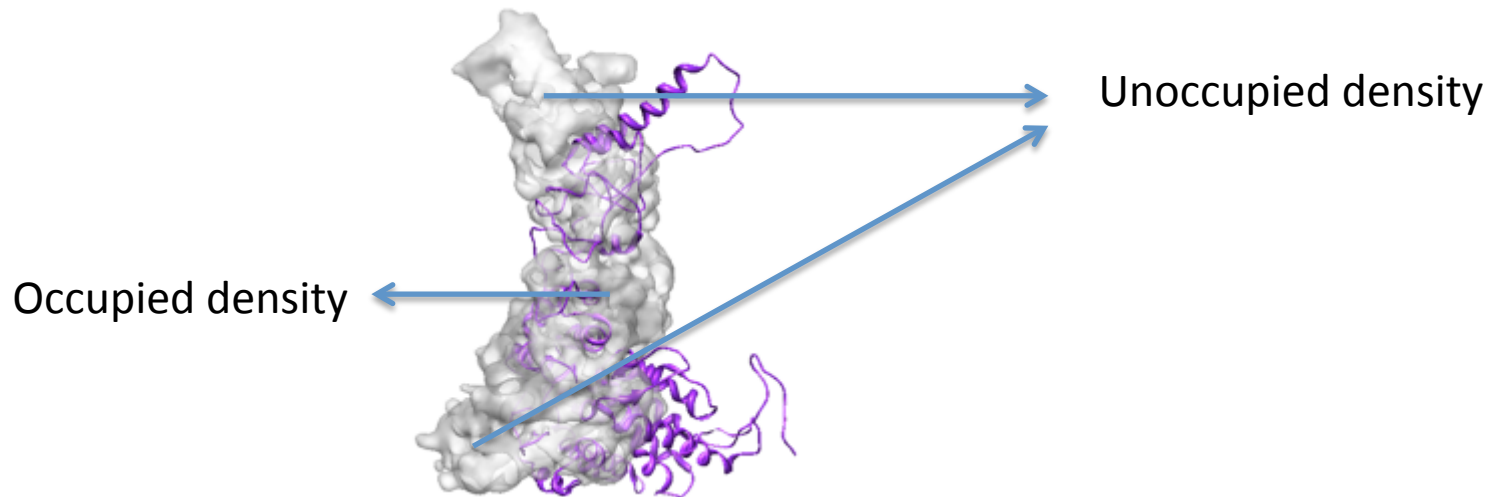
# Rigid-body Docking

- Other scores can be computed, to evaluate how confident we can be about a rigid-body docking result:

  - Atom inclusion:
    - How many atoms are inside the observed density?
    - If many atoms are outside, then the model may not be the best match for the map, or it may be incorrectly docked.

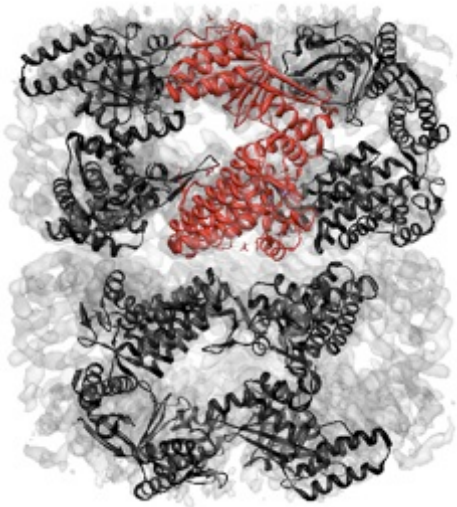

Atom outside

Atoms inside

# Rigid-body Docking

- Other scores can be computed, to evaluate how confident we can be about a rigid-body docking result:

  - Density occupancy:
    - How many grid points with high density values are occupied by atoms
    - If many areas inside the density maps are unoccupied, then the docked models may not fully explain the experimental density map.
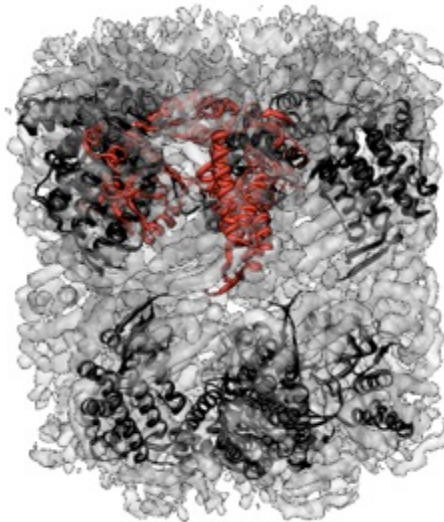


Unoccupied density

Occupied density

# Rigid-body Docking

- Other scores can be computed, to evaluate how confident we can be about a rigid-body docking result:

  - Clashes with other docked models, or with symmetric copies
    - Do multiple docked models or symmetric copies clash?
    - If there are many clashes, this could signify an unreliable docking result.



Few clashes between
symmetric copies:
correct dock

Many clashes between
symmetric copies:
incorrect dock

Red ribbon:
   docked model

Black ribbon:
   symmetric copies

# Rigid-body Docking Summary

- Exhaustive search
  - Can take a long time
  - Thoroughly searches for alignments of the model inside the map
  - Can produce unreliable results if a map is very heterogeneous (i.e. densities are higher in some regions - this makes the CC-score higher in those regions even if the model does not fit well there)
- Segger
  - Allows docking by aligning known models to segmented regions
  - Faster than exhaustive search
  - Allows more control of where the model is docked
  - Requires more prior-knowledge about where the model might be found inside the map
  - Allows cross-validation between segmentation and docking results
- The cross-correlation score is widely used as the score that determines the best fit
- Other scores can be used to assess confidence in fit (z-score, atom inclusion, density occupancy, etc.)

# Flexible Fitting

- Typically, start by rigidly-docking a model
- Then assess whether the fit could be better
  - The conformation of the model comes from homology or X-ray crystallography
    - But, the conformation seen in the cryoEM map may be different
  - Molecules are inherently flexible – they can adopt different conformations; it is possible that the cryoEM density map shows a different conformation
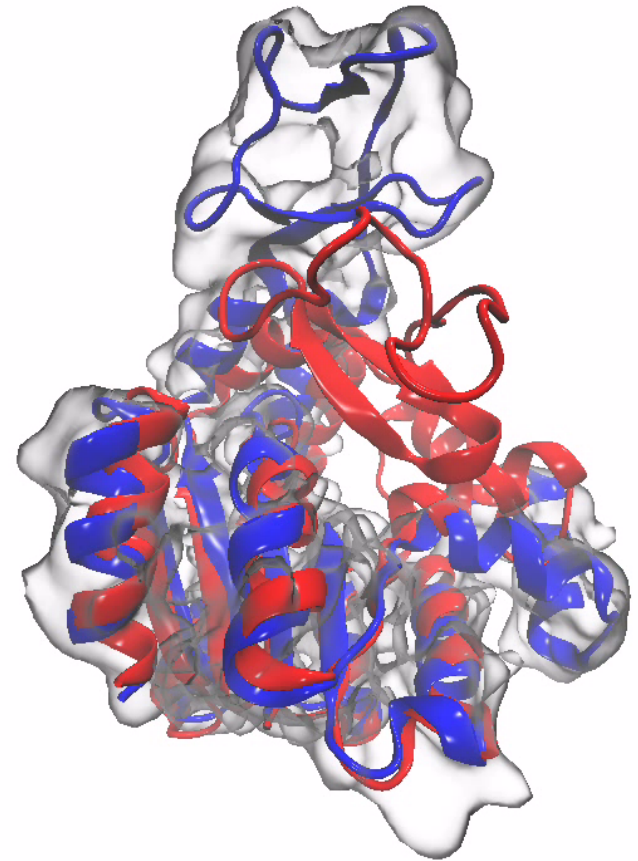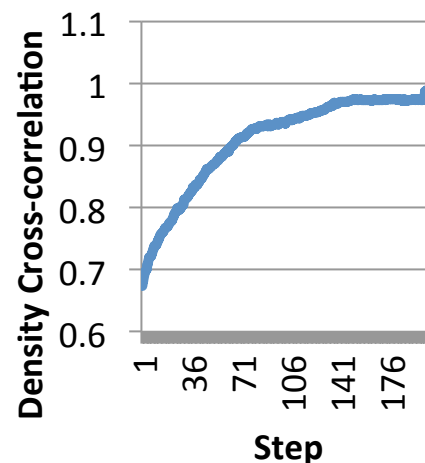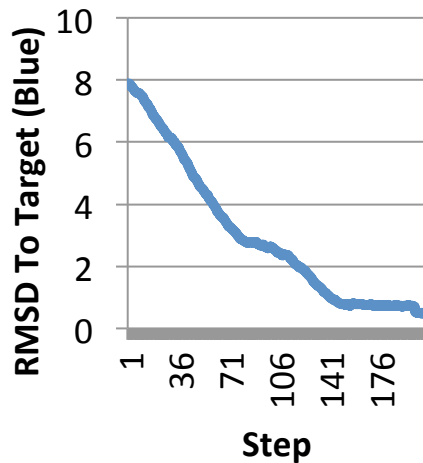
# Flexible Fitting

- How to do flexible fitting?
  - Assume atoms, or different subunits of the model (e.g. secondary structures), can move individually
  - Move each atom/secondary structure towards higher densities (in gradient direction)
  - Goals:
    - Increase fitting score (cross-correlation, atom inclusion, density occupancy, clashes, etc.)
    - Maintain a good structure – good bonds, good angles, good dihedrals, etc.

# Flexible Fitting Methods

- MDFF (Molecular Dynamics Flexible Fitting)
  - Klaus Schulten lab
  - http://www.ks.uiuc.edu/Research/mdff/
- Direx
  - Gunnar Schröder lab
  - http://www.schroderlab.org/
- Flex-EM
  - Andrej Sali lab
  - http://salilab.org/Flex-EM/
- Rosetta
  - David Baker lab
  - http://www.rosettacommons.org/software/

# Molecular Dynamics Flexible Fitting (MDFF)

- Apply force to each atom in the model to the high density of the map
- Apply full force field to keep good geometry (e.g. bonds, angles, dihedrals), and complementary charge interactions
- Use extra forces and bonds to maintain Secondary structures; H-bonds; Chirality; Cis-peptide bonds

# Flex-EM

- Also uses molecular dynamics, moving atoms to high density

- Secondary structures are rigid bodies

- Simulated annealing (temperature is cycled between low and high values)

# Deformable Elastic Network (DEN)

The network adapts itself to refine only those degree of freedom for which the experiment provides information.

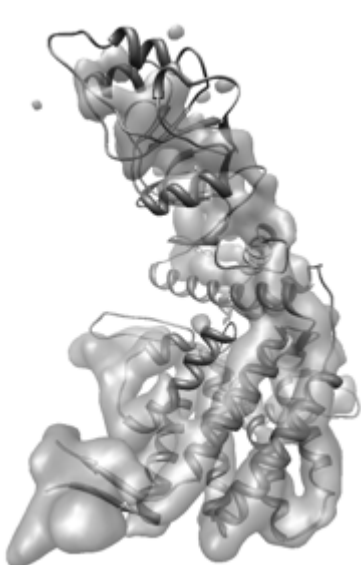. All-atom description
. Correct stereochemistry
. No atom clashes

Schroeder et al (2007) *Structure*

# Evaluation of Flexible Fitting Methods

- Does the model fit the density map better after flexible fitting?
  - Compute model-to-density scores
    - Density cross-correlation
    - Atom inclusion
    - Occupancy of high-density areas
    - Local density cross-correlation
- The model changes during flexible fitting; so we must also re-assess the quality of the model
  - Hence, we also calculate model-alone scores
    - Calculated using Molprobity
      - http://molprobity.biochem.duke.edu
    - Some of the scores it computes:
      - Clashes between atoms
      - Rotamer quality
      - Ramachandran outliers/favored
      - Bad bonds/angles

# Flexible Fitting Test Case
## Mm-Cpn @ 8Å (EMD-5140)



PDB:3KFE
(crystal structure
closed state
rigid-body docked)

MDFF

Flex-EM

Direx

Rosetta
(from cryoEM
challenge)

# Modeling Δlid Mm-Cpn Open State



Deformable Elastic Network (DEN)
Schroeder, Brunger and Levitt., *Structure*, 2007

# Validate the 8 Å Resolution Cryo-EM Built Model for Δlid Mm-Cpn in the Apo State with 2 Modeling Tools


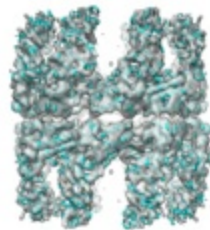
Overal Cα RMSD: 2.8Å
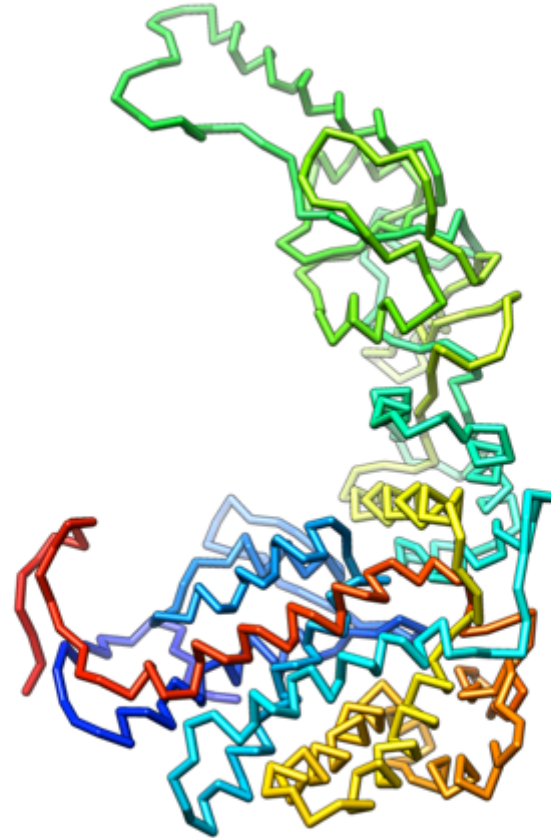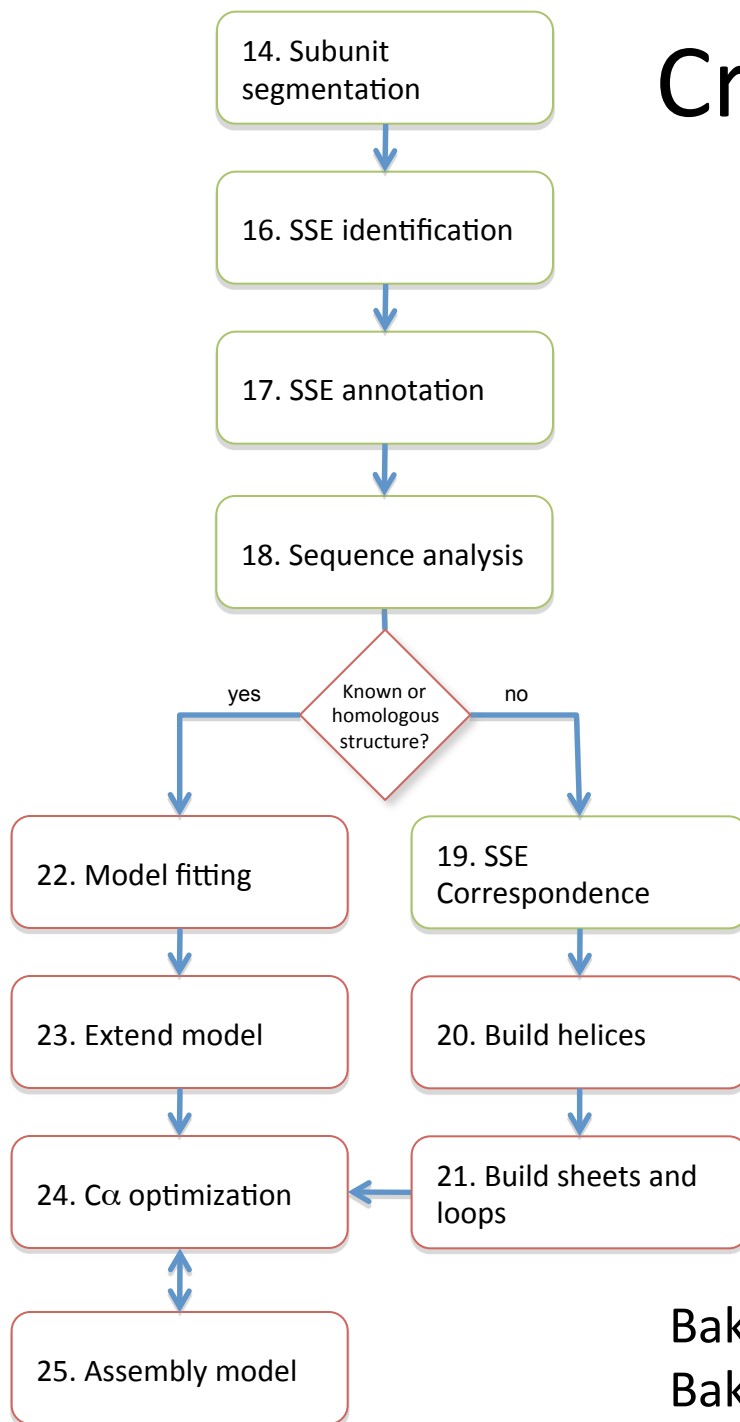
API Cα RMSD: 4.3Å

INT Cα RMSD: 2.3Å

EQU Cα RMSD:1.7Å

Rosetta-built Model

Direx-built Model

# Data Statistics for Mm-Cpn + ATP/AlFx

- 110,000x effective camera magnification
- 637 CCD frames (Gatan 4k CCD)
- ~29,926 particle images used for final 3-D reconstruction
- Density map determined to 4.3 Å resolution based on the 0.5 criterion of Fourier Shell Correlation

# Cryo-EM Structure Analysis

14. Subunit segmentation

16. SSE identification

17. SSE annotation

18. Sequence analysis

Known or homologous structure?

yes → 22. Model fitting

no → 19. SSE Correspondence

22. Model fitting

23. Extend model

24. Cα optimization

19. SSE Correspondence

20. Build helices

21. Build sheets and loops

25. Assembly model



Baker et al. (2010) *Nature Protocols* **5**: 1697-1708
Baker et al. (2010) *Methods in Enzymology* **483**: 1-29

# Mm–cpn
# at 4.3 Å

# closed state

Mm-Cpn @4.3 Å

*SSE detection*

*segmentation*

**constrained C-alpha homology model**

**MM-CPN subunit**

*skeletonization*

**MM-CPN primary sequence**

VLPENMKRYMGRDAQRMNILAGRIIAETVRSTLGPKGMDKMLVDDLGD
VVVTNDGVTILREMSVEHPAAKMLIEVAKTQEKEVGDGTTTAVVVAGELL
RKAEELLDQNVHPTIVVKGYQAAAQKAQELLKTIACEVGAQDKEILTKIAM
TSITGKGAEKAKEKLAEIIVEAVSAVVDDEGKVDKDLIKIEKKSGASIDDTELI
KGVLVDKERVSAQMPKKVTDAKIALLNCAIEIKETETDAEIRITDPAKLMEFI
EQEEKMLKDMVAEIKASGANVLFCQKGIDDLAQHYLAKEGIVAARRVKKS
DMEKLAKATGANVIAAIAALSAQDLGDAGLVEERKISGDSMIFVEECKHP
KAVTMLIRGTTEHVIEEVARAVDDAVGVVGCTIEDGRIVSGGGSTEVELS
MKLREYAEGISGREQLAVRAFADALEVIPRTLAENAGLDAIEILVKVRAAHA
SNGNKCAGLNVFTGAVEDMCENGVVEPLRVKTQAIQSAAESTEMLLRID
DVIAAE

*homology modeling*

**homology model**

**MM-CPN density map**

**C-alpha assembly model**

**scaled MM-CPN density map**

**model structure factors**

*calculate structure factors*

*re-scale density map*

*model assembly*

*model refinement*

A

C$\alpha$ deviation

C$\beta$ deviation
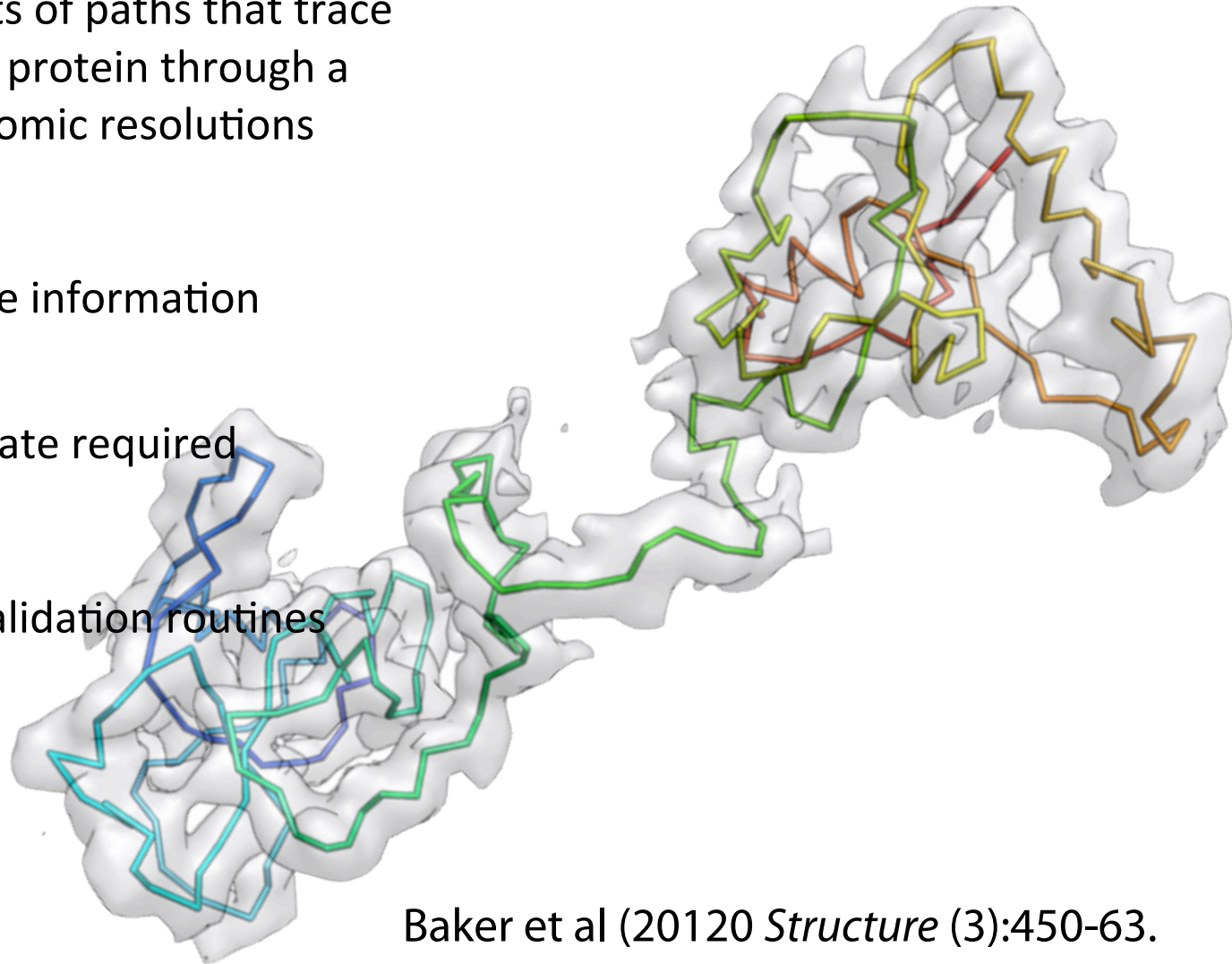
B

B-factor

C

# Cryo-EM Structure Analysis

26. Turn Cα model into a density map

↓

27. Re-scale density map

↓

28. Cα to atomic model

↓

29. Model completion

↓

30. Sidechain registration

↓

31. Model optimization

↓

32. Fix outliers

↓

33. Fit model to density map

Baker et al. (2010) *Nature Protocols* **5**: 1697-1708
Baker et al. (2010) *Methods in Enzymology* **483**: 1-29

# Automatic de Novo Modeling: Pathwalking

Goal: Find a path or sets of paths that trace the complete path of a protein through a density map at near-atomic resolutions

- No SSEs required

- No explicit sequence information required

- No structural template required

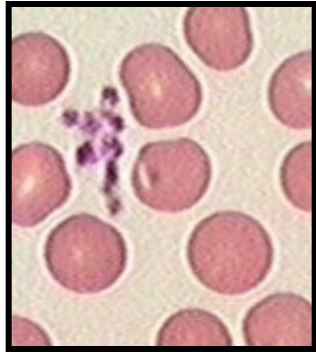- Automated

- Optimization and validation routines

Baker et al (20120 *Structure* (3):450-63.

# Flexible Fitting – Summary

- MDFF
    - Con: Slower (since it moves each atom individually, using a full-force field); however it can be run on multiple processors
    - Pro: Good model-alone scores
- Flex-EM
    - Con: Secondary structures are kept rigid, so they are less able to deform to fit the density
    - Pro: Good model-alone score
- Direx
    - Pro: Fast, since full-force field is not used. Increases model-density scores the most.
    - Con: Poor model-alone scores, particularly in lower-resolution maps
- Rosetta
    - Pro: uses information from known crystal structures (fragment library) or initial de novo bulit; hence it gives good model-alone scores
    - Con: takes longer to change structure significantly, since only small parts are modified at one time

# Structural Biology from Man to Atoms