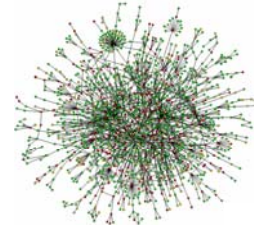# Atomic Resolution Modeling
## of
# Large Macromolecular Assemblies

### Haim J. Wolfson
### School of Computer Science
### Tel Aviv University

---

## Complexes as functional modules of the cell

**Protein-Protein interaction network**


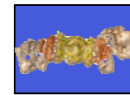
*Jeong et. al. , 2001*

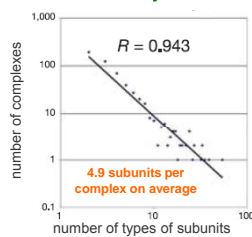**Complexes**



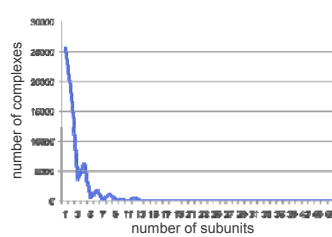ATP synthase    Virus    26S proteasome    Chaperonin   ...   Nuclear pore complex

---

## Protein complex size statistics

**distribution of complex size in yeast**     **protein data bank**



$R = 0.943$

**4.9 subunits per complex on average**

number of complexes / number of types of subunits
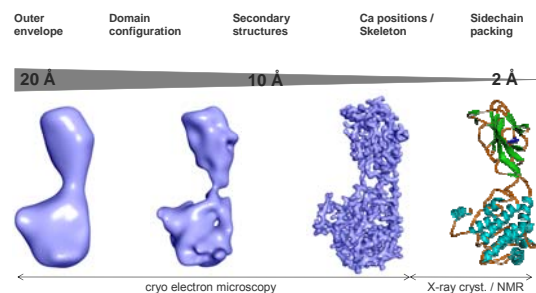
number of complexes / number of subunits

*Krogan et. al., Nature,2006*

**There are thousands of biologically relevant macromolecular complexes whose structures are yet to be characterized.**
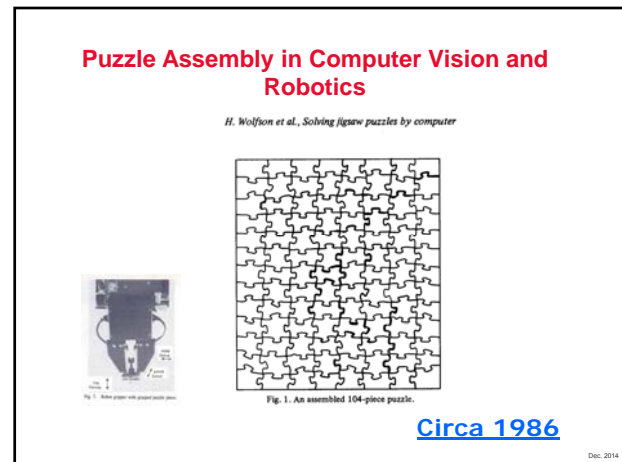
---

## Experimental techniques for Protein Structure determination

| Outer envelope | Domain configuration | Secondary structures | Ca positions / Skeleton | Sidechain packing |
|---|---|---|---|---|

20 Å      10 Å      2 Å



cryo electron microscopy     X-ray cryst. / NMR

**Use hybrid methods to bridge the resolution gaps**

## Analogy

**Multi-molecular assembly is analogous to the solution of 3D puzzles –a classical spatial Pattern Discovery task.**

High resolution data

Low resolution data



Dec. 2014

## Puzzle Assembly in Computer Vision and Robotics

*H. Wolfson et al., Solving jigsaw puzzles by computer*



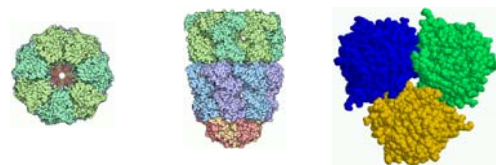Fig. 1. An assembled 104-piece puzzle.

**Circa 1986**

Dec. 2014

## Additional Low Resolution Data Sources

- FRET
- Existence of di-sulfide bonds
- MasSpec (e.g.distance constraints by chemical cross linking).
- SAXS
- Interaction Data (Y2H, gene fusion, similarity with known complexes, etc.)
- and more…

Dec. 2014

## SPECIAL FREQUENT CASE:

*Structure Prediction of (cyclically) Symmetric Multi-Molecular Assemblies*



D. Schneidman-Duhovny et al., Proteins, 60, 217--223, (2005).

D. Schneidman-Duhovny et al., NAR 33 (web server issue), W363—W367, (2005).
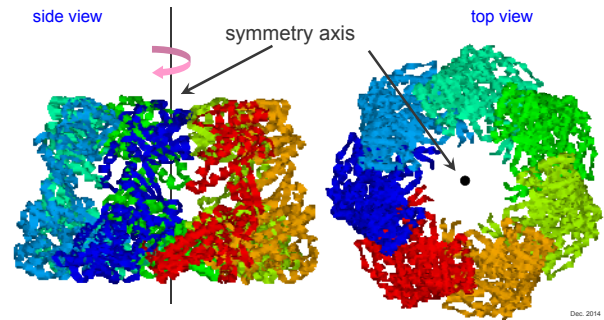
## Exploiting the Symmetry Constraints

- A trivial "naïve" approach – perform "regular" multimolecular docking and discard non-symmetric solutions.
- A more sophisticated approach – use the symmetry constraints as an integral part of the algorithm to reduce complexity and improve accuracy.
- Observation – if point A in the protein is matched after the symmetry rotation to point B, one can detect a plane to which the symmetry axis is perpendicular and its location is restricted to a known circle in that plane.
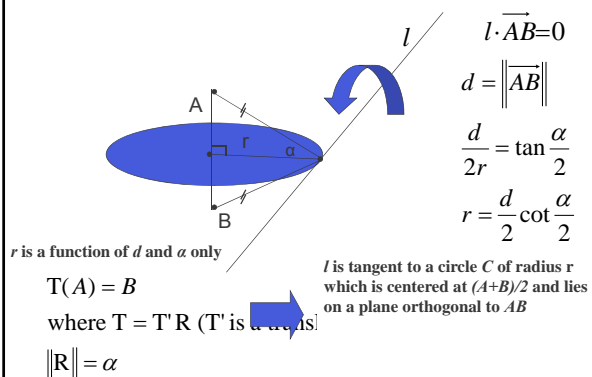
Dec. 2014

## Cyclic Symmetry

- Cyclic symmetry is defined by rotation of a single unit around an **axis**.
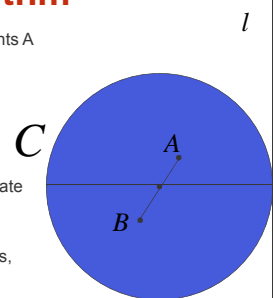- The angle is determined by a number of units **n**.



side view          symmetry axis          top view

Dec. 2014

## Geometric Analysis



$$l \cdot \overrightarrow{AB} = 0$$

$$d = \left\| \overrightarrow{AB} \right\|$$

$$\frac{d}{2r} = \tan \frac{\alpha}{2}$$

$$r = \frac{d}{2} \cot \frac{\alpha}{2}$$

*r* is a function of *d* and *α* only

$$T(A) = B$$

where T = T' R (T' is a trans)

$$\|R\| = \alpha$$

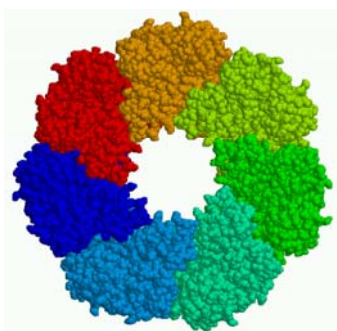*l* is tangent to a circle *C* of radius r which is centered at *(A+B)/2* and lies on a plane orthogonal to *AB*

Dec. 2014

## The Algorithm

- For each pair of matching interest points A and B
  - Calculate C$_{AB\alpha}$
    - For δ = 0 to 360-Δ step Δ
    - Calculate l$_{C\delta}$
    - Calculate T$_{l\alpha}$
    - If T is valid add T to the candidate transformation list
- Cluster transformations
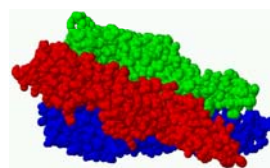- Calculate the score for transformations, which are cluster representatives



Dec. 2014

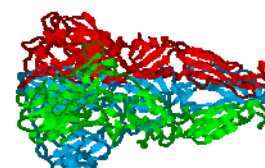Chaperon: 2.5 Å RMSD prediction for the homo-heptamer.

Dec. 2014

## CAPRI Target 10: 9.0 Å RMSD prediction for the homo-trimer of a viral coat protein
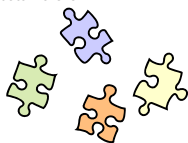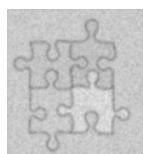


Our Prediction          Crystal Structure

Dec. 2014

### Exploit Low Resolution Info – EM, SAXS, FRET etc.

Structural models of the subunits at atomic level

Low/Medium resolution EM density map



Dec. 2014

## Previous Work

**Early work : Fitting of atomic structures to the density map by cross-correllation.**

**In essence – structural alignment at different resolutions.**
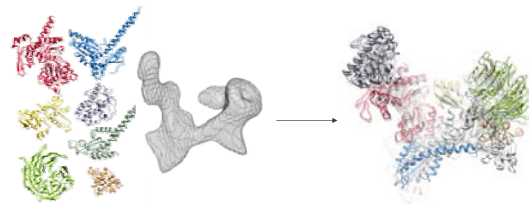
**Recent work : Hybrid Methods.**

Dec. 2014

## Publications

- W. Wriggers, R.A. Milligan, J.A. McCammon, Situs: a package for docking crystal structures into low resolution maps for electron microscopy, J. Struct. Biol. 125, (1999), 185—195.
- Z. Yang, K. Lasker, D. Schneidman-Duhovny, B. Webb, C.C. Huang, E.F. Petersen, T. D. Goddard, E.C. Meng, A. Sali, T.E. Ferrin, UCSF Chimera MODELLER, and IMP: An integrated modeling system, J. Struct. Biol. 179, (2011), 269—278.
- E. Karaca, A.S.J. Melquiond, S.J. deVries, P.L. Kastritis and A.M.J.J. Bonvin, Building Macromolecular Assemblies by Information-driven Docking : Introducing the HADDOCK MultiBody docking server, Mol. Cel. Proteomics 9, (2010), 1784—1794.
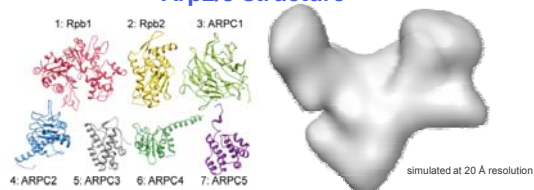
Dec. 2014

## MultiFit



**Find the placements ( translation and orientation) of atomic components in the density map of their association.**

**Lasker, Topf, Sali, Wolfson, JMB 2009**

**Lasker, Sali, Wolfson, Proteins 2010**

Dec. 2014

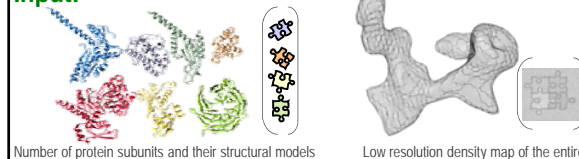## MultiFit - Example of a  Task :Assemble the Arp2/3 structure



simulated at 20 Å resolution

| component | %seq id | Cα RMSD |
|-----------|---------|---------|
| Rpb1  | 40 | 5.1  |
| Rpb2  | 48 | 2.5  |
| ARPC1 | 16 | 6.1  |
| ARPC2 | 29 | 21.4 |
| ARPC3 | 99 | 0.4  |
| ARPC4 | 29 | 14.3 |
| ARPC5 | 94 | 5.5  |

**COMPONENT STRUCTURE –**

**OUTPUT of HOMOLOGY MODELING**

Dec. 2014

## MultiFit: A geometric view

**Input:**



Number of protein subunits and their structural models          Low resolution density map of the entire assembly

**Goal: Determine the assembly configuration** *optimizing*

$$S = \textbf{docking} + \textbf{Structural alignment} + \textbf{docking}$$

Geometric complementarity          Fitting score          Envelope penetration

Structural accuracy

Find the placements ( translation and orientation) of atomic components in the density map that minimizes the scoring function

Dec. 2014

## Few representative reasons for the difficulty of multiple fitting

- Scoring
  - Cross-correlation measure alone is not always sufficient to place a component in the map.
  - Cross-correlation score does not check for geometric complementary between interacting components.
  - Docking alone is problematic, since the accuracy of docking methods depends on the accuracy of the individual atomic structures

    *Solution: use a scoring function that considers fitting and geometric complementarity simultaneously*

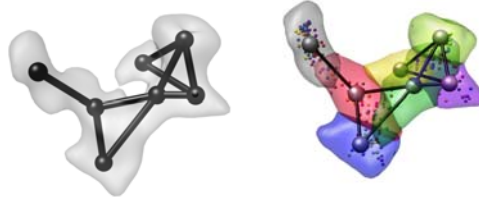| Pair of components | Pairwise docking rank |
|---|---|
| ARP3/ARPC2 | 12185 |
| ARP3/ARPC3 | 854 |
| ARP3/ARPC4 | 5888 |
| ARPC1/ARP2 | 4663 |
| ARPC1/ARPC5 | 5504 |

- Optimization
  - Sequential fitting or sequential pairwise docking may not result in the right configuration in the general case.
  - Enumerating all possible configurations of components of large assemblies is too expensive

## Focus the subunit placement search around anchor points

- **anchor graph:** a low-resolution description of the assembly.
  - **nodes:** points in 3D that approximate the centroid positions of the assembly components.
  - **edges:** between nodes that are close in space.
- The anchor graph was constructed using a Gaussian Mixture Model segmentation of the density map.
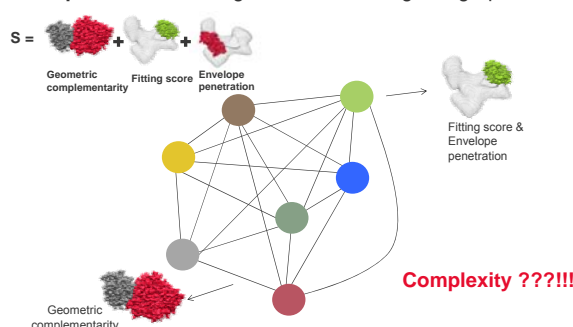


The anchor graph    Sampling of subunit centroids at anchor graph pts

## Reduce the multiple fitting problem to optimization of a subunit location and orientation graph

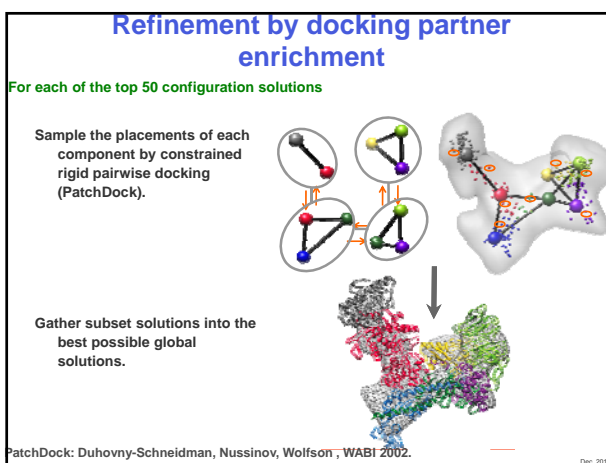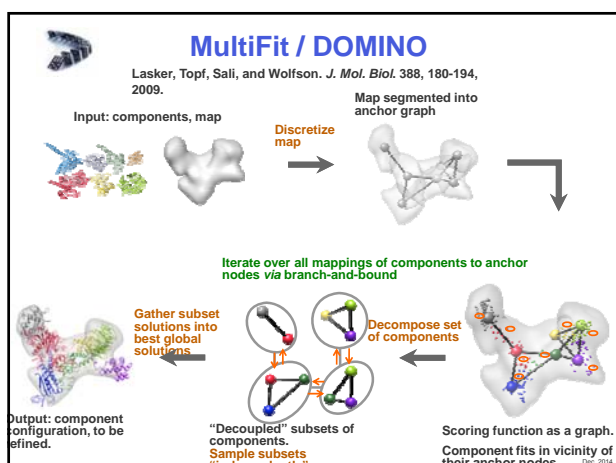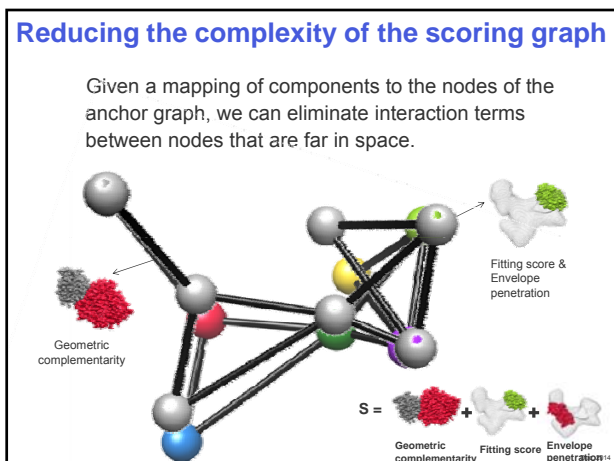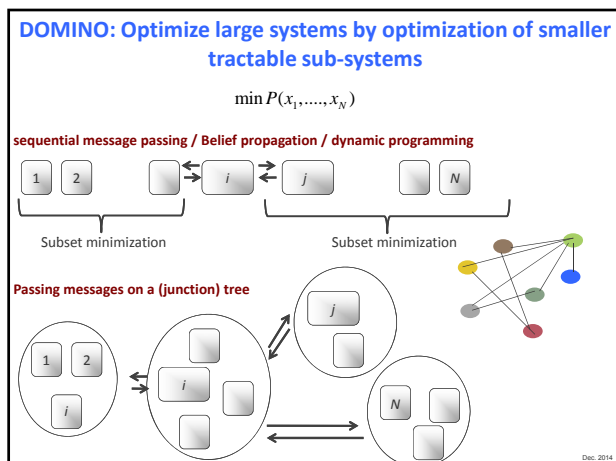**1. Represent** the scoring function as a weighted graph.



**Complexity ???!!!**
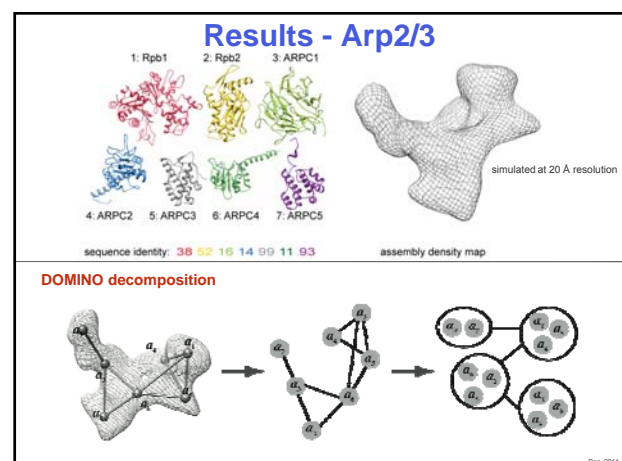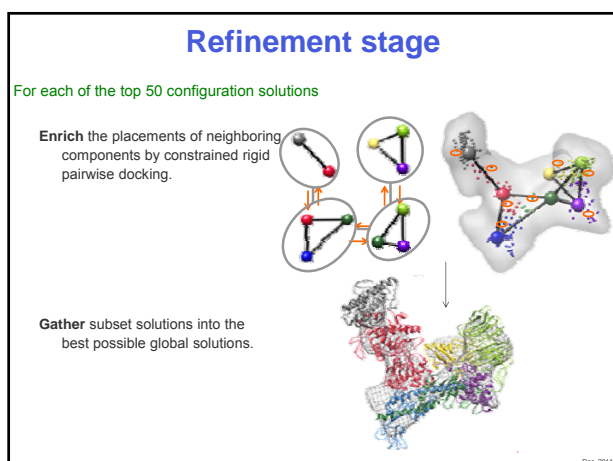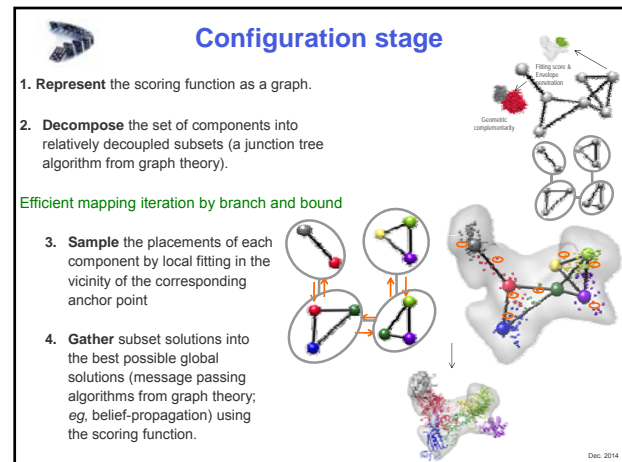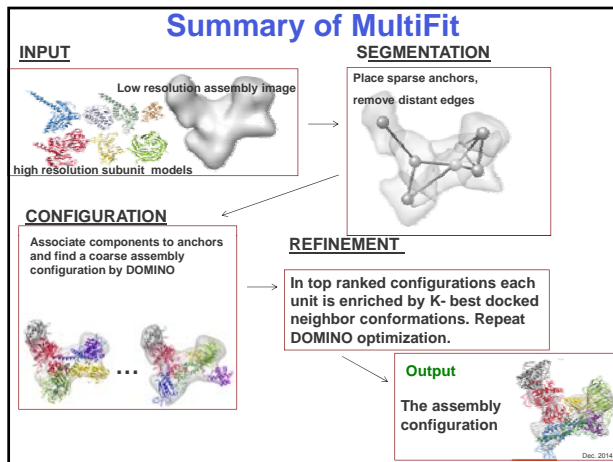
## Graphical Models

- Use a belief propagation type algorithm to detect the optimal solution.
- Apply the algorithm both in the placement stage and orientation refinement stages.
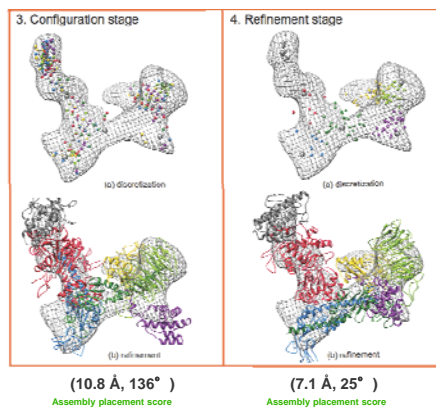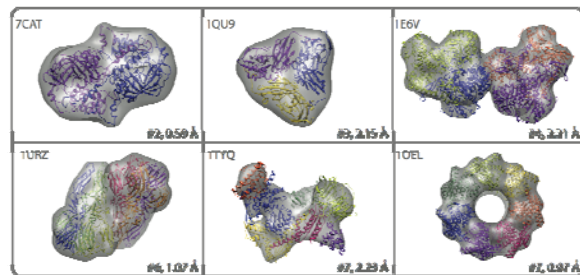- Utilise the Junction Graph structure.

## DOMINO: Optimize large systems by optimization of smaller tractable sub-systems

$$\min P(x_1,....,x_N)$$

**sequential message passing / Belief propagation / dynamic programming**



Subset minimization       Subset minimization

**Passing messages on a (junction) tree**



Dec. 2014

## Reducing the complexity of the scoring graph

Given a mapping of components to the nodes of the anchor graph, we can eliminate interaction terms between nodes that are far in space.



Fitting score & Envelope penetration

Geometric complementarity

$$S = \text{Geometric complementarity} + \text{Fitting score} + \text{Envelope penetration}$$

Dec. 2014

## MultiFit / DOMINO

Lasker, Topf, Sali, and Wolfson. *J. Mol. Biol.* 388, 180-194, 2009.

**Input: components, map**

**Discretize map**

**Map segmented into anchor graph**



**Iterate over all mappings of components to anchor nodes *via* branch-and-bound**

**Gather subset solutions into best global solutions**

**Decompose set of components**

Output: component configuration, to be refined.

"Decoupled" subsets of components. Sample subsets "independently"

Scoring function as a graph. Component fits in vicinity of their anchor nodes.

Dec. 2014

## Refinement by docking partner enrichment

**For each of the top 50 configuration solutions**

Sample the placements of each component by constrained rigid pairwise docking (PatchDock).



Gather subset solutions into the best possible global solutions.

PatchDock: Duhovny-Schneidman, Nussinov, Wolfson, WABI 2002.

Dec. 2014

Summary of MultiFit



Configuration stage



Refinement stage



Results - Arp2/3

**Arp2/3 Example: Optimization stages**



**Benchmark results**

density maps simulated to 20Å
no proteomics data was used as input
Best model within the top 10 models

Lasker, Sali and Wolfson. *Proteins*, 78, 3205-3211, 2010



**2011 EM Modeling challenge**



**2011 EM modeling challenge: GroEL**

## 2011 EM modeling challenge: MmCpn

**model on map**



|  | mmcpn opened | mmcpn closed |
|---|---|---|
| resolution (Å) | 8 | 4.3 |
| cross-correlation | 0.9 (0.94) | 0.78 (0.81) |
| $C_\alpha$-RMSD to reference | 1.7 | 0.8 |

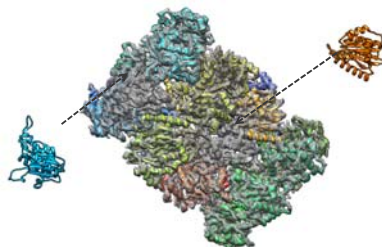Dec. 2014

## 3D-MOSAIC



**D. Cohen, N. Amir, H.J. Wolfson - submitted**

Dec. 2014

## New Multimolecular Assembly Method: 3D-Mosaic

- Capitalizes on the steady improvement in EM map resolution to sub-nanometer accuracy.
- Fits _simultaneously_ numerous atomic resolution subunits into intermediate res
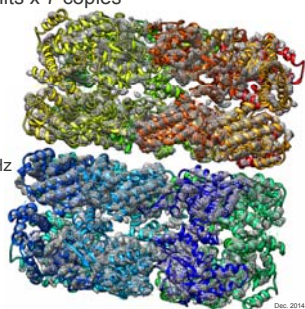


Dec. 2014

## Advantages of 3D-Mosaic

- Requires _no prior segmentation_ of the EM map.
- Handles _"missing"_ subunits.
- _Highly efficient_ handling of a large number of _multiple structurally homologous copies_ of complex subunits.
- Efficient new method for integrative _simultaneous_ modeling of large multi-molecular assemblies by formulating the optimization task as an Integer Linear Program (ILP).
- Incorporates both EM and X-link information into the same framework.

_D. Cohen, N. Amir, H.J. Wolfson, 3D-MOSAIC: An efficient method for integrative modeling of large multimolecular assemblies, (to be submitted)._
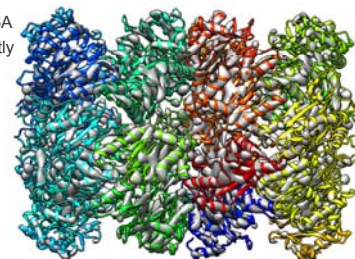
Dec. 2014

### Results - GroEL

- Protein chaperonin important for proper protein folding
- 14 subunits, 2 unique subunits x 7 copies
- @4.2A resolution :
  - RMSD of solution : 2.5A
  - All units placed correctly
  - Run time :
    - Placement: 10min
    - Optimization: 15sec
  - Measured on 12 core, 3.06GHz
    Ubuntu 12.04 machine



Dec. 2014

### Results : 20S Proteasome – experimental map

- Breakdown of proteins
- 28 subunits, 2 unique subunits x 14 copies
- @6.8A resolution :
  - RMSD of solution : 1.5A
  - All units placed correctly
  - Run time :
    - Placement: 2-4min
    - Optimization: 1min



Dec. 2014

### Current Major Challenge

**Modeling a multimolecular assembly from sequence data alone by threading the sequences on the EM structural scaffold.**

Dec. 2014