

13. A Bayesian approach to cryo-EM structure determination

Sjors Scheres

EMBO course 2017
Birkbeck College, London

An intuitive introduction

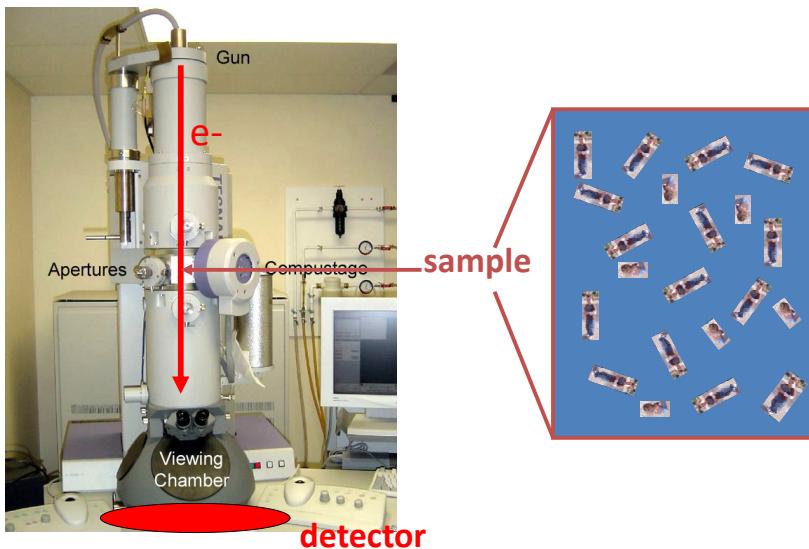
Agenda

- An intuitive introduction
- Alignment
 - Dealing with the incomplete problem
 - maxCC vs ML (real-space)
- Classification
 - Multi-reference alignment in 2D
 - and in 3D
- Fourier-space formulation
 - Regularised likelihood optimisation (Bayesian approach)

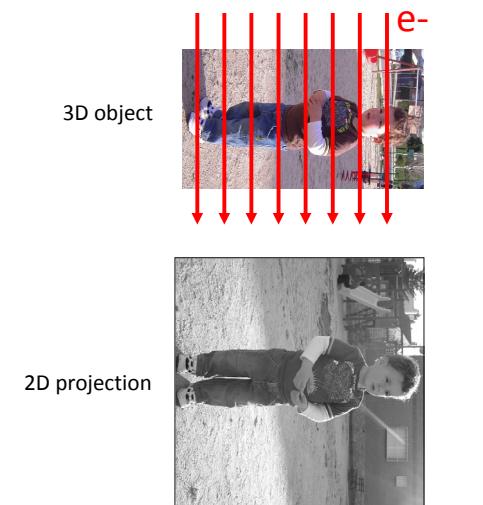
An example “protein”



Experimental setup

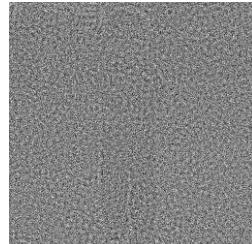


Electron microscopy imaging



Further inconveniences

- Microscope imperfections introduce artefacts
 - Contrast Transfer Function (CTF)
- Large amounts of noise

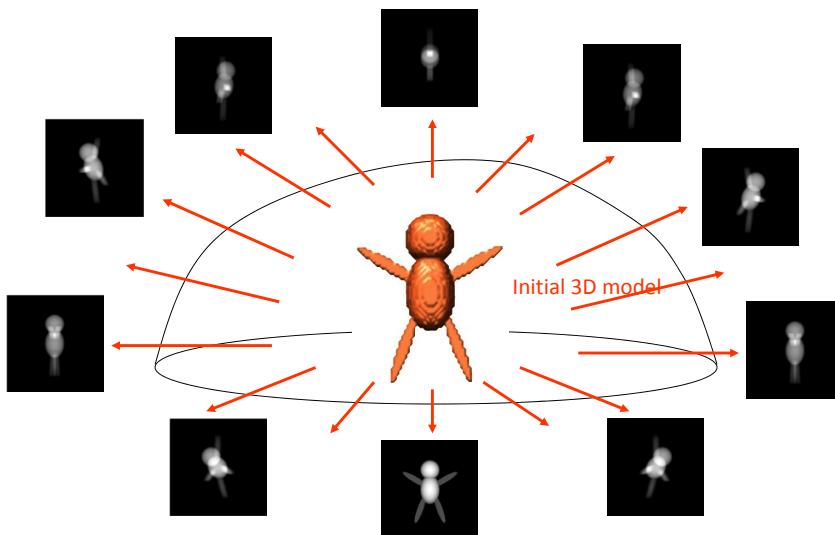


Single particle analysis

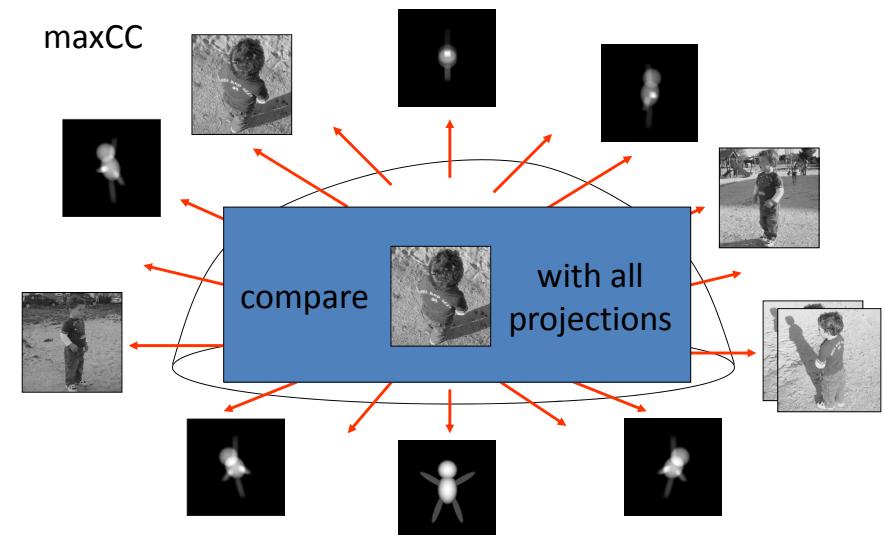
- Embedded in ice: many unknown orientations
- Two rows of eight grayscale images each, showing a young child in various poses and orientations, used for single-particle analysis.

 - Combine all 2D projections into a 3D reconstruction

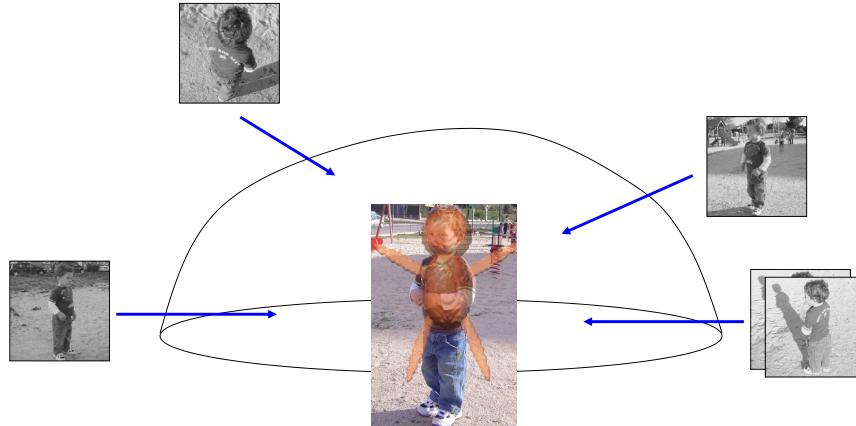
Projection matching



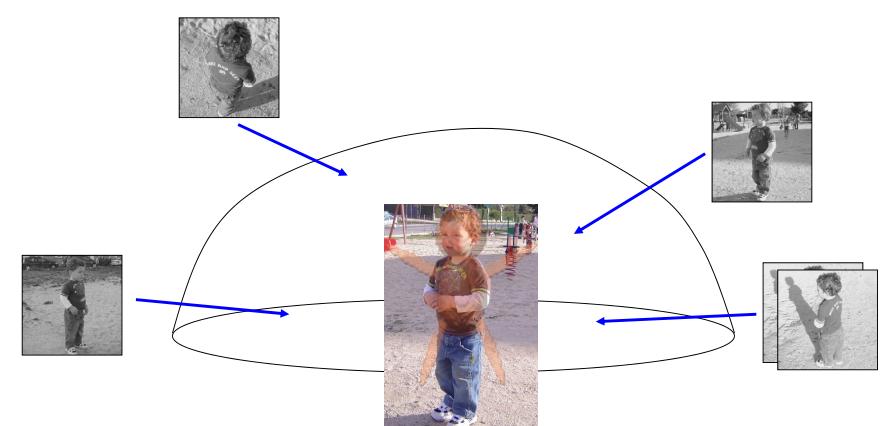
Projection matching



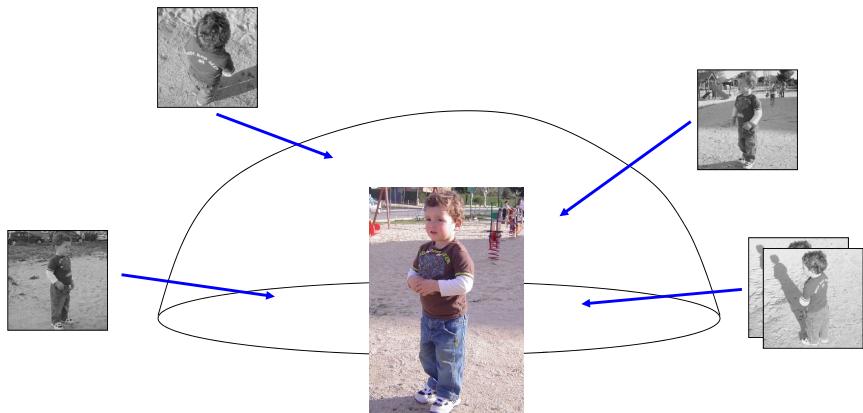
3D reconstruction



Iterative refinement



Iterative refinement



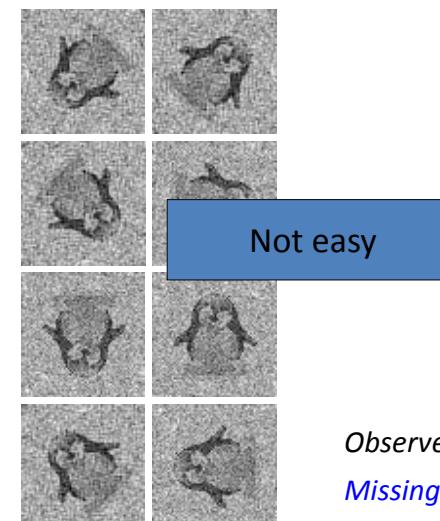
Alignment

Or how to 'match' projections

Incomplete data problems

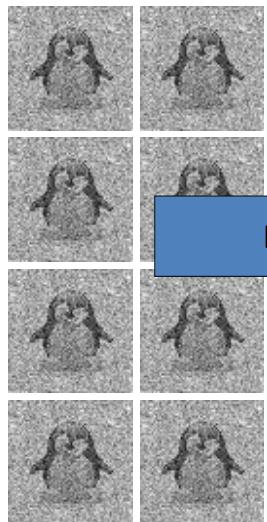
- Part of the data was not observed experimentally
 - Orientations
 - Class assignments
- Difficult to solve!
 - Iterative methods?
- Complete data problem would be very easy to solve
- (Another famous one: the phase problem in XRD)

Incomplete data problems



Observed data (X): images
Missing data (Y): orientations

Complete data problems



white Gaussian noise

$$L(\Theta) = P(X | \Theta)$$

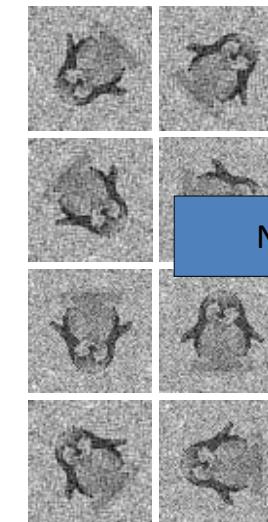
Easy!



$$\hat{A}^{MLE} = \frac{1}{N} \sum_{i=1}^N X_i$$

Observed data (X): images

Incomplete data problems



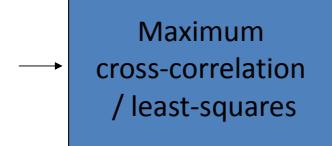
Not easy

Observed data (X): images

Missing data (Y): orientations

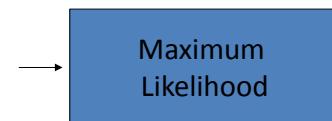
Incomplete data problems

- Option 1: add Y to the model



$$L(Y, \Theta) = P(X | Y, \Theta)$$

- Option 2: marginalize over Y



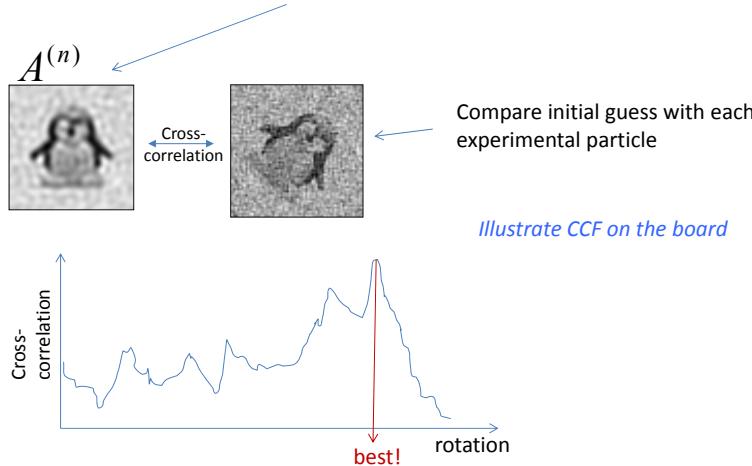
$$L(\Theta) = P(X | \Theta) = \int_Y P(X | Y, \Theta) P(Y | \Theta) d\phi$$

↓
Probability of X ,
regardless Y

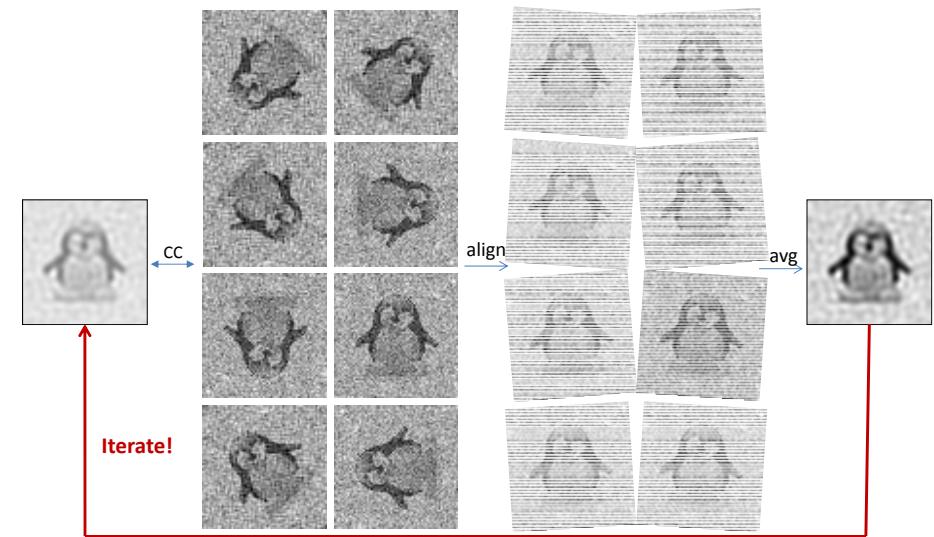
The maxCC approach

Reference-based alignment

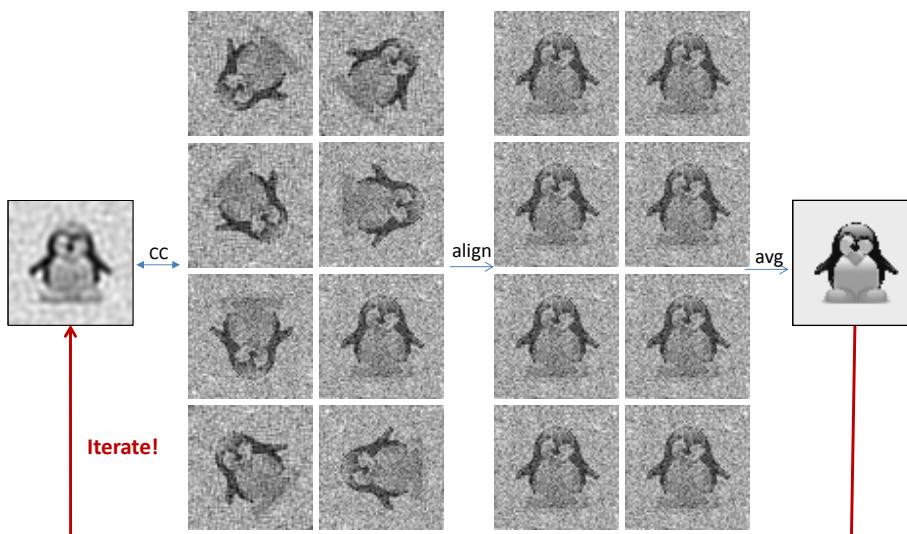
- Starts from some initial guess about the structure



Align and average

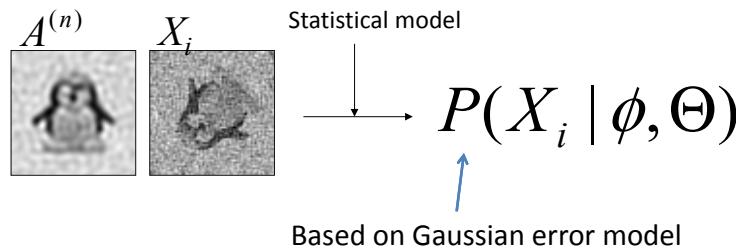


Align and average



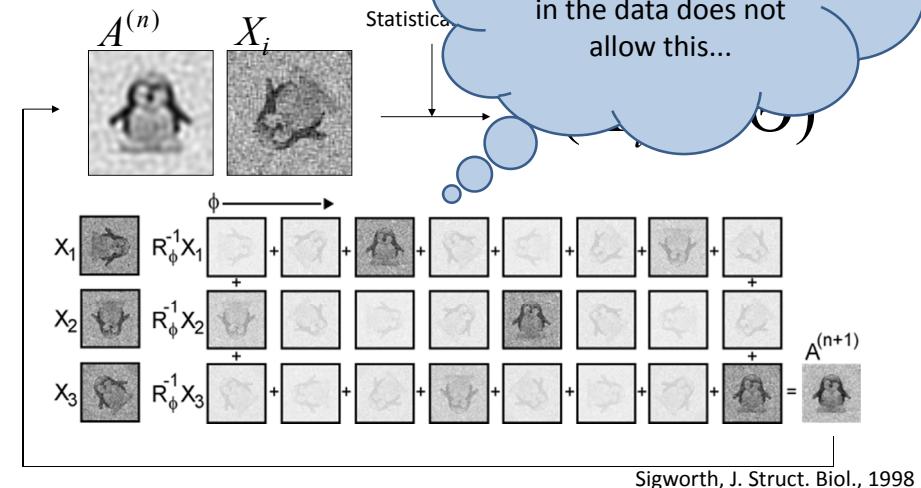
The ML approach

Maximum likelihood



$$P(X_i | \phi, \Theta) = \prod_{j=1}^J \frac{1}{2\pi\sigma^2} \exp\left(\frac{(\mathbf{P}_\phi V_j - X_{ij})^2}{-2\sigma^2}\right)$$

Maximum likelihood



Incomplete data problems

- Option 1: add Y to the model

$$L(Y, \Theta) = P(X | Y, \Theta)$$

Maximum cross-correlation

- Option 2: marginalize over Y

$$L(\Theta) = P(X | \Theta) = \int_Y P(X | Y, \Theta) P(Y | \Theta) d\phi$$

Probability of X , regardless Y

Maximum Likelihood

Incomplete data problems

In the limit of **noiseless data** the Two techniques are equivalent!

Maximum cross-correlation

Maximum Likelihood

Many software packages now use ML:
cryoSPARC, SPARX/SPHIRE, FREALIGN,
XMIPP, RELION

Read more? See *Methods in Enzymology*, 482 (2010)

Classification

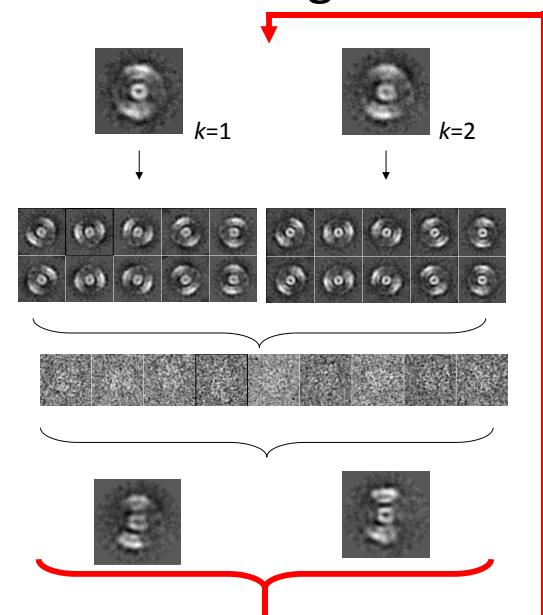
The 2D multi-reference algorithm

estimates for K
2D objects

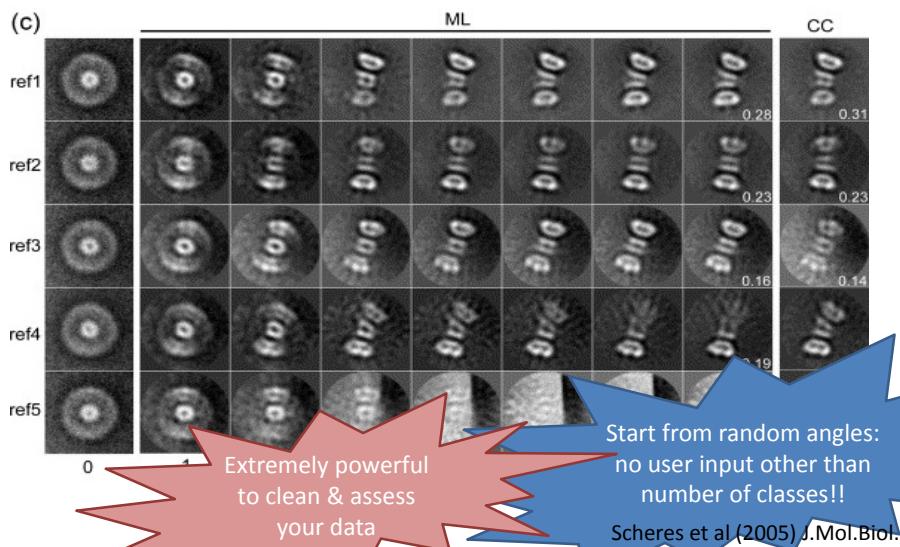
sampled rotations 360°

for each image, calculate all
 $P(\text{image}_i | k, \text{rot})$

calculate new 2D average
as *probability weighted averages*

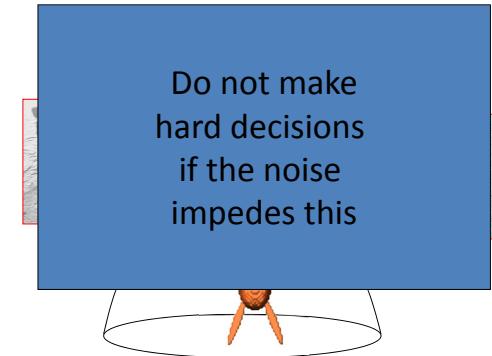


Reference-free 2D class averaging



3D alignment & classification

3D ML refinement



“Probability-weighted angular assignment”

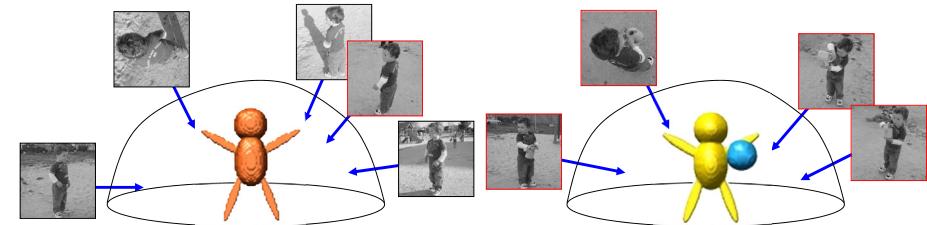
Structural heterogeneity



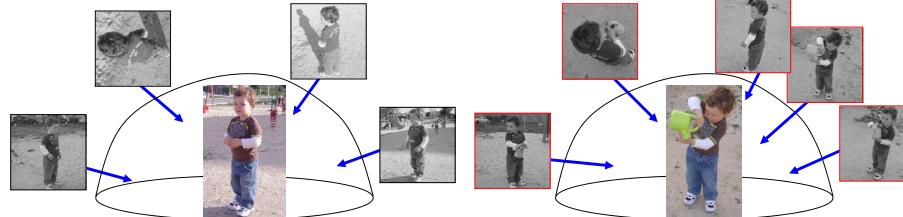
Initial model

- Expectation-Maximisation is a local optimizer!
 - Gets stuck in nearest (local) minimum
- Bad model in -> bad model out!!!
 - Much less of a problem with high-resolution data
- Stochastic methods may reach global minimum
 - Stochastic Hill Climbing (SIMPLE)
 - Stochastic Gradient Descent (cryoSPARC & RELION)

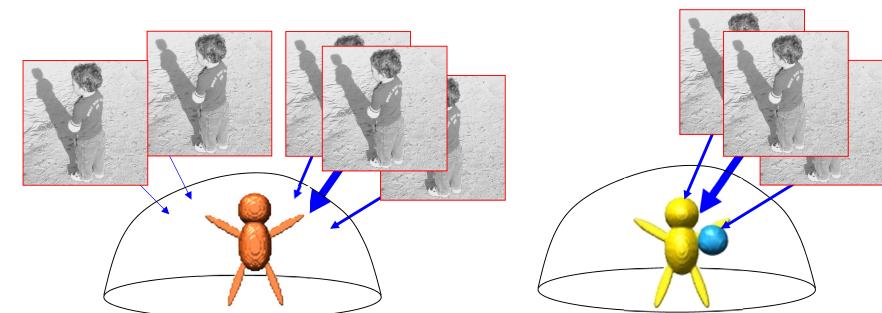
Multi-reference refinement



Multi-reference refinement



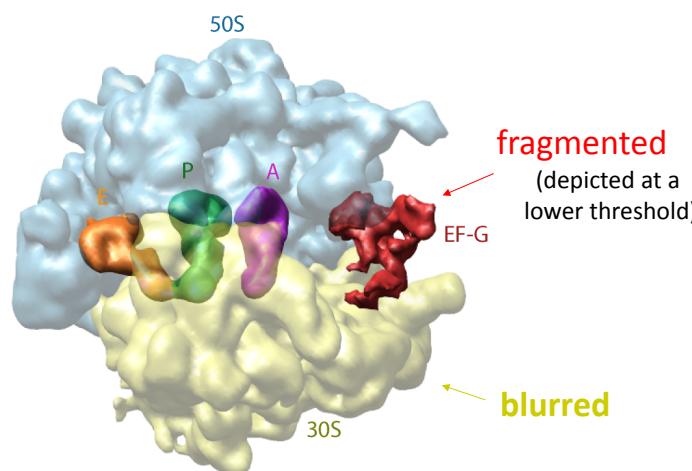
ML3D classification



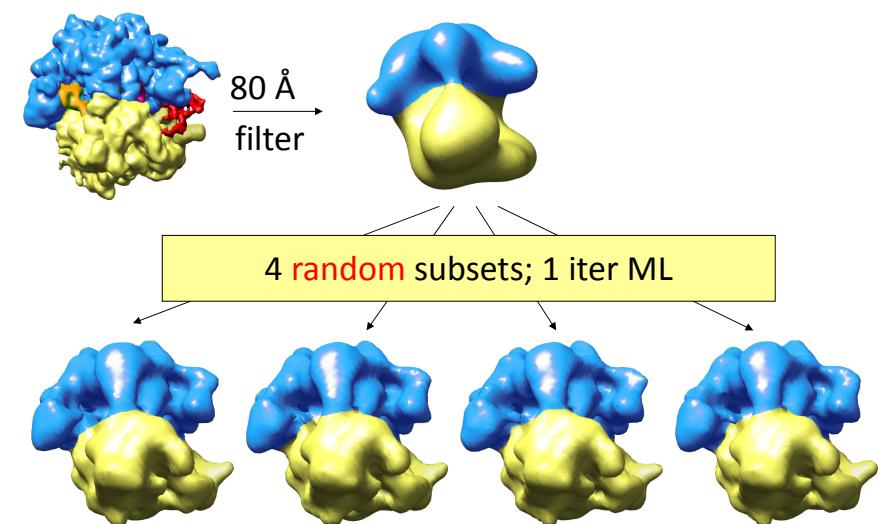
“Probability-weighted angular assignment”

Prelim. ribosome reconstruction

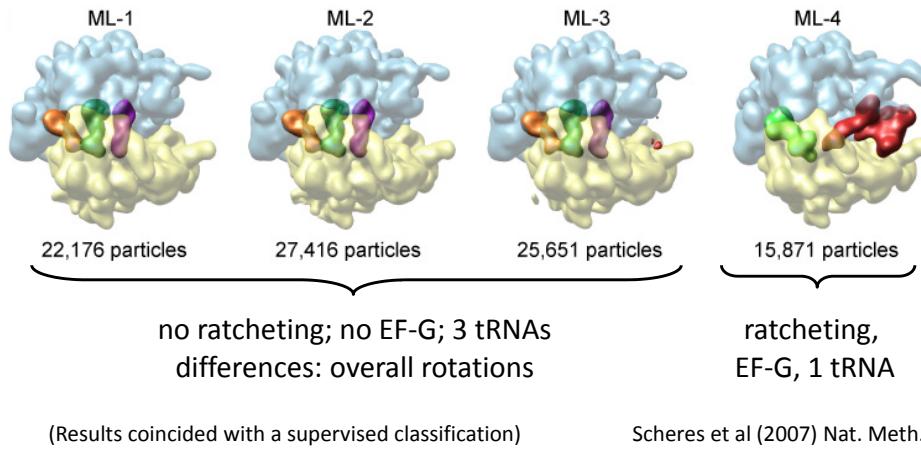
91,114 particles; 9.9 Å resolution



Seed generation

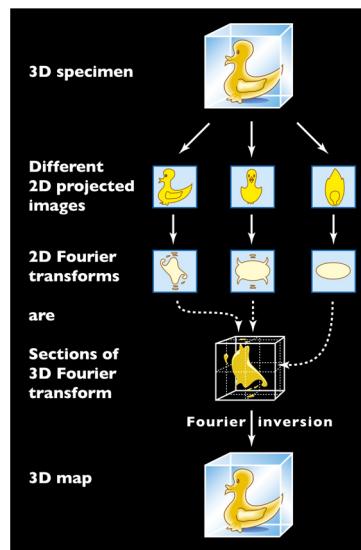


ML-derived classes

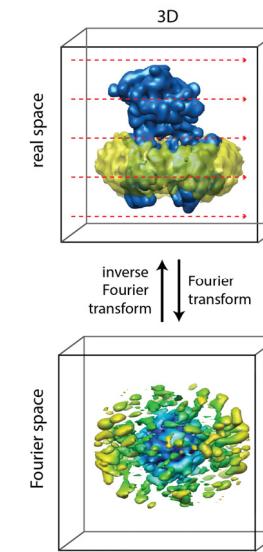


Fourier-space formulation

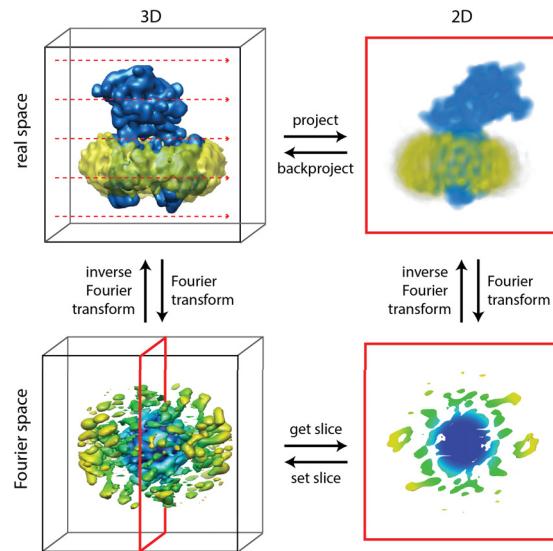
Projection-slice theorem



Projection-slice theorem



Projection slice theorem



Data model

- Real-space
 - Fourier space
- $$X_i = \text{CTF}_i \otimes \mathbf{P}_\phi V_k + N_i \quad X_i = \text{CTF}_i \mathbf{P}_\phi V_k + N_i$$
- Convolute w/ CTF
 - \mathbf{P}_ϕ implements integrals
 - N_i describes white noise
 - Multiply w/ CTF
 - \mathbf{P}_ϕ takes a slice
 - N_i describes coloured noise

Regularised Likelihood

Maximum-likelihood estimators

- The best one can do...
- ...in the limit of *infinitely large data sets*
- But my data set is limited in size, right?!
 - Even with Krios, K3 & EPU!

The bad news

- The experimental data alone is not enough to determine a unique solution!
- There are many noisy reconstructions that describe the data equally well...
- Danger of incorrect interpretation...

The good news

- By incorporating external information, a different problem may be solved for which a unique solution does exist!
- Regularisation
 - Conventional regularisation approaches
 - Wiener filtering
 - Low-pass filtering

A Bayesian view on regularization

$$P(\Theta | X) = \frac{P(X | \Theta)P(\Theta)}{P(X)}$$

$$\text{Posterior} = \frac{\text{Likelihood} * \text{Prior}}{\text{Evidence}}$$

Regularised likelihood optimisation

Likelihood

- Assume noise is Gaussian and independent
 - in Fourier space
 - with spectral power $\sigma^2(v)$: *coloured noise*

$$P(X_i | k, \phi, \Theta) = \prod_{j=1}^J \frac{1}{2\pi\sigma_{ij}} \exp\left(-\frac{\|X_{ij} - \text{CTF}_{ij}(\mathbf{P}_\phi V_k)_j\|^2}{2\sigma_{ij}^2}\right)$$

Prior

- Assume signal is Gaussian and independent
 - in Fourier space
 - Limited power $\tau^2(v)$: *smoothness in real space!*

$$P(\Theta) = \prod_l \frac{1}{2\pi\tau_{kl}} \exp \left\{ \frac{\|V_{kl}\|^2}{-2\tau_{kl}^2} \right\}$$

3D Wiener filter

$$V^{(n+1)} = \frac{\sum_{i=1}^N \int \Gamma_{i\phi}^{(n)} \mathbf{P}_\phi^\top \frac{\text{CTF}_i}{\sigma_i^{2(n)}} X_i d\phi}{\sum_{i=1}^N \int \Gamma_{i\phi}^{(n)} \mathbf{P}_\phi^\top \frac{\text{CTF}_i^2}{\sigma_i^{2(n)}} d\phi + \frac{1}{\tau^{2(n)}}$$

- Calculates SSNR(v) (as a 3D function)
- Handles uneven orientational distribution
- Handles astigmatic CTFs & CTF envelope
- Corrects CTF & low-pass filter
- Optimal linear filter*

WITHOUT ARBITRARINESS!

Expectation maximization

$$V^{(n+1)} = \frac{\sum_{i=1}^N \int \Gamma_{i\phi}^{(n)} \mathbf{P}_\phi^\top \frac{\text{CTF}_i}{\sigma_i^{2(n)}} X_i d\phi}{\sum_{i=1}^N \int \Gamma_{i\phi}^{(n)} \mathbf{P}_\phi^\top \frac{\text{CTF}_i^2}{\sigma_i^{2(n)}} d\phi + \frac{1}{\tau^{2(n)}}} \longrightarrow \text{Wiener (optimal) filter for CTF-corrected 3D reconstruction / 2D class averages}$$

$$\sigma_i^{2(n+1)} = \frac{1}{2} \int \Gamma_{i\phi}^{(n)} \|X_i - \text{CTF}_i \mathbf{P}_\phi V^{(n)}\|^2 d\phi \longrightarrow \text{Estimate resolution-dependent power of noise from the data}$$

$$\tau^{2(n+1)} = \frac{1}{2} \|V^{(n)}\|^2 \longrightarrow \text{Estimate resolution-dependent power of signal from the data}$$

$$\Gamma_{i\phi}^{(n)} = \frac{P(X_i | \phi, \Theta^{(n)}) P(\phi | \Theta^{(n)})}{\int \limits_{\phi'} P(X_i | \phi', \Theta^{(n)}) P(\phi' | \Theta^{(n)}) d\phi'}$$

Recapitulating

- Alignment & classification are incomplete problems
 - Best dealt with by marginalisation (ML)
- 2D and 3D problems are very similar
- Fourier-space is most convenient
 - CTF multiplication
 - Slices instead of line integral projections
 - Coloured noise-model
 - Regularised Likelihood function -> 'optimal' filters

Further Reading

- Penczek, Fundamentals of Three-Dimensional Reconstruction from Projections, *Methods in Enzymology*, , **482** (2010) p 1
- Penczek, Image restoration in cryo-electron microscopy, *Methods in Enzymology*, , **482** (2010) p 35
- Sigworth, Doerschuk, Carazo & Scheres, An Introduction to Maximum-Likelihood Methods in Cryo-EM, *Methods in Enzymology*, **482** (2010) p 263
- Scheres, Classification of Structural Heterogeneity by Maximum-Likelihood Methods, *Methods in Enzymology*, **482** (2010) p 295
- Scheres, Processing of Structurally Heterogeneous Cryo-EM Data in RELION, *Methods in Enzymology*, **579** (2016) p 125
- www2.mrc-lmb.cam.ac.uk/relion (tutorial & Wiki pages)