### Optimizing image processing

#### Sjors Scheres MRC Laboratory of Molecular Biology

What if one would have many noisy 2D projections of assumedly identical 3D objects in unknown orientations, and one would want to know that 3D structure?

And how strictly adhering to theory helps...

## Questions?

Conventional and ML 3D (projection matching) refinement (and the differences between them), Bayesian extension, how to avoid overfitting, how to avoid model bias, multi- reference refinement (classification).



## Cryo-EM inconveniences



## Inverse problems

### The forward model

 $X_i = \mathbf{CTF}_i \mathbf{P}_{\phi} V_k + N_i$ 

Given V and CTF, we can simulate X very well

But the other way around is more difficult!

Incompleteness

## 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



### Incomplete data problems



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

## Incomplete data problems

• Option 1: add *Y* to the model

Maximum cross-correlation / least-squares

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

• Option 2: marginalize over 
$$Y \rightarrow$$
  

$$L(\Theta) = P(X | \Theta) = \int_{Y} P(X | Y, \Theta) P(Y | \Theta) d\phi$$
Probability of X,  
regardless Y

### The maxCC approach

### Statistical data model

 $X_i = P_{\varphi} V_k$ 



## **Reference-based alignment**

• Starts from some initial guess about the structure



Compare initial guess with each experimental image



### Align and average



### Align and average



### The ML approach

### Statistical data model

 $X_i = P_{\varphi} V_k$ 



### Statistical data model

 $X_i = P_{\varphi}V_k + N_i$ white / coloured Gaussian noise **Statistical** description of the noise

## Maximum likelihood





## Incomplete data problems

• Option 1: add Y to the model

 $I(Y_{\text{the limit}} = P(X | Y_{\text{the limit}})$ Two techniques are equivalent!



$$L(\Theta) = P(X | \Theta) = \int_{Y} P(X | Y, \Theta) P(Y | \Theta) d\phi$$

Probability of X, regardless Y

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

## maxCC projection matching

- Compare X<sub>i</sub> with CTF<sub>i</sub>P<sub>φ</sub>V for all φ, and select optimal φ\* based on some similarity measure (e.g. CC)
  - Reconstruction:

terate  

$$V^{(n+1)} = \frac{\sum_{i=1}^{N} \mathbf{P}_{\phi^*}^{\mathrm{T}} \mathrm{CTF}_i X_i}{\sum_{i=1}^{N} \mathbf{P}_{\phi^*}^{\mathrm{T}} \mathrm{CTF}_i^2}$$

• Least-squares solution to V (?)

## Maximum likelihood refinement

- Calculate a probability P(X<sub>i</sub> | φ, Θ) for all φ, based on an explicit noise model (e.g Gaussian)
  - Probability-weighted angular assign

iterate

Theory says this is the best one can do (in the limit of infinitely large data sets) ht:



## Remaining issues

- 1. What to use as initial guess?
  - Local optimizer:
  - Wrong initial model -> wrong answer!
  - Model bias!

## Model bias

- common-lines models are difficult
  - 2D projections are OK
  - Their combination in 3D is not
- Better (?)
  - RCT, sub-tomogram averaging, homologous structure
- EMAN(2) better than projection matching
  - But also not guaranteed...

## Remaining issues

- 1. What to use as initial guess?
  - Wrong initial guess may lead to wrong answers!
  - Model bias!
- 2. What if multiple structures are present?
  - Cannot align against 1 reference
  - Alignment + classification problem

### **Prelim. ribosome reconstruction** 91,114 particles; 9.9 Å resolution



In collaboration with Haixiao Gao & Joachim Frank

### Seed generation



### **ML-derived classes**



(Results coincided with a supervised classification)

### BUT....

- 3D-classification is not a cure for bad data....
- Works best for few well-defined states
- Not all variability can be resolved
  - Continuous heterogeneity -> compromises
  - Many states may be tricky (expensive at least)
    - Supervised classification may be an alternative:
    - Fischer et al, Nature, 2010 (>20 states, 2M particles)
  - Ultimately a signal-to-noise ratio issue

## Remaining issues

- 1. What to use as initial guess?
  - Wrong initial guess may lead to wrong answers!
     Model bias!
- 2. What if multiple structures are present?
  - 1. Cannot align against 1 reference
  - 2. Alignment + classification problem

3. What if I do not infinite amounts of data...

### Ill-posedness

## The bad news

- The experimental data alone is not enough to determine a unique solution! (*ill-posed*)
  - Noise tends to accumulate in the reconstruction

## The bad news

- The experimental data alone is not enough to determine a unique solution!
  - Noise tends to accumulate in the reconstruction
  - Overfitting
  - Over-estimation of resolution
  - –Incorrect interpretations

## The good news

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

## 2D Wiener filter

- Assume noise is independent – with spectral power  $\sigma^2(\upsilon)$
- Assume signal is independent

   with spectral power τ<sup>2</sup>(v)
- Minimise noise in 2D average: (optimal filter)



Damp A for those Fourier components where all CTFs are zero or  $\tau^2/\sigma^2$  is small

Correct CTF AND low-pass filter!

### <u>3D Wiener filter</u>

CHAPTER ONE

#### FUNDAMENTALS OF THREE-DIMENSIONAL RECONSTRUCTION FROM PROJECTIONS

Pawel A. Penczek

Reconstruction methods based on backprojection and direct Fourier inversion methods require the implementation of a form of Wiener filter, which schematically is written as (see Chapter 2):

$$D = \frac{\sum_{n} \text{CTF}_{n} \text{SSNR}_{n} G_{n}}{\sum_{n} \text{CTF}_{n}^{2} \text{SSNR}_{n} + 1}.$$
(1.27)

The summation in the numerator can be realized as a backprojection of the Fourier transforms of (n - 1)D projections multiplied by their respective CTFs and SSNRs, so the result is nD. However, it is far from obvious how the summation in the denominator can be realized such that the result would have the intended meaning after the division is performed.

Meth. Enzym. (2010)

## 3D Wiener filter

- Same assumptions
- Plus (often):  $\frac{\tau^{2}(v)}{\sigma^{2}(v)} = \text{SSNR}(v) = 1/C \text{ BUT THIS IS NOT TRUE!!!!}$

Low-pass filtering effect is lost!

$$V^{(n+1)} = \frac{\sum_{i=1}^{N} \mathbf{P}_{\phi^*}^{\mathrm{T}} \mathrm{CTF}_i X_i}{\sum_{i=1}^{N} \mathbf{P}_{\phi^*}^{\mathrm{T}} \mathrm{CTF}_i^2 + C}$$
 "Wiener constant"

# "Arbitrary" low-pass filters

- Many different ones exist
  - choose shapes, effective resolution, width, etc.



### A Bayesian view on regularization



Posterior = Likelihood \* Prior Evidence

### **Maximum A Posteriori estimation**

## Likelihood

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

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

## Prior

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

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

### **Expectation** maximization



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

$$\Gamma_{i\phi}^{(n)} = \frac{P(X_i \mid \phi, \Theta^{(n)}) P(\phi \mid \Theta^{(n)})}{\int_{\phi'} P(X_i \mid \phi', \Theta^{(n)}) P(\phi' \mid \Theta^{(n)}) d\phi'}$$

## 3D Wiener filter



- Calculates SSNR( $\upsilon$ ) (as a 3D function)
- Handles uneven orientational distribution
- Handles anisotropic CTFs & CTF er
- Corrects CTF & low-pass
- Optimal linear filter

WITHOUT ARBITRARINESS!

## Recapitulating...

- Inverse problem: needs iterating
- Incomplete problem: needs marginalizing
- Ill-posed problem: needs regularizing

- Bayesian approach:
  - Does all 3 things in optimizing a single function!
  - "Learns" optimal parameters from the data
  - No *ad-hoc* parameters to tune by the user

## **Preventing overfitting**

A little detour...

Scheres & Chen (2012) Nature Methods

### The pitfalls of undetected overfitting

simulated

- 20k simulated GroEL particles
- Conventional projection matching



## **Overfitting-free refinement**



easy to script into many packages...

### Only lower resolution data drive alignment



## Experimental data

- 5,053 GroEL particles\*
- 50,330  $\beta$ -galactosidase particles
- 5,403 hepatitis B capsid particles\*\*
- High-resolution crystal structures!

*kindly provided by NCMI/Steven Ludtke kindly provided by Tony Crowther*

### GroEL



### Hepatitis B capsid



## $\beta$ -galactosidase



## Conclusions

- Overfitting may be avoided without loss of resolution
  - Gold-standard FSCs between 2 independent models
- In the absence of overfitting
  - Higher-resolutions may be obtained
  - Maps are clean and easy to interpret, fit, etc.
  - FSC=0.143 is a reliable resolution estimate

### Back to the statistical approach

### Gold-standard FSC in the Bayesian approach

- Refine two models independently
- At each iteration: calculate  $\tau^{2}(\upsilon)$  based on  $\text{FSC}_{\text{gold}}$

## **RE**gularised **LI**kelihood **O**ptimisatio**N**

#### http://www2.mrc-Imb.cam.ac.uk/relion



Page Discussion

Read Edit View history 🔻

Go Search

[edit]

#### **Running RELION**

Using the GUI

#### Navigation

Main page Community portal

Toolbox

What links here Related changes Upload file Special pages Printable version Permanent link

ile Run type: 3D reconstruct	
	tion 🗘 Start new run 🗘
O CTF Optimisation Sampling	Running
Number of MPI procs:	8 - []
Number of threads:	8 ?
Submit to queue?	Yes ¢?
Queue name:	openmpi_8 ?
Queue submit command:	qsub ?
Standard submission script:	res/app/relion/gui/qsub.csh ? Browse
	Print command Run!

RELION may be used to perform different tasks (run types). The following run-types may be selected from the drop-down menu at the top of the GUI:

- 2D averaging: calculate reference-free 2D class averages
- 3D reconstruction: perform 3D (multi/single-reference) refinements

### Some results

Tom Walz: test new programs on old data!

# Classify structural variability

- Standard data set (i.e. used by many groups...)
  - 10,000 70S ribosomes (50% +EFG; 50% -EFG)
  - MAP-refinement K=4



8 hrs on 64 CPUs

## 3D auto-refine results

	$\beta$ -galactosidase	$\operatorname{groEL}$	hepatitis B	rotavirus
Sample characteristics				
Size (MDa)	0.45	0.8	4	60
Symmetry	D2	D7	Ι	Ι
Microscopy settings				
Microscope	FEI Polara G2	Jeol 3000SFF	Hitachi HF2000	FEI Tecnai F30
Voltage (kV)	80	300	200	300
Defocus range $(\mu m)$	1.2 - 2.7	1.9 - 3.2	1.0 - 2.0	1.2 - 2.9
Detector	Kodak SO163	Kodak SO163	Kodak SO163	Kodak SO163
Data characteristics				
Image size (pixel <sup>2</sup> )	$100 \times 100$	$128 \times 128$	$220 \times 220$	$400 \times 400^a$
Pixel size (Å)	2.93	2.12	2.00	2.40
Nr. particles	50,330	5,053	5,403	3,700
RELION parameters				
Particle mask diameter (Å)	200	205	400	785
Initial low-pass filter (Å)	60	60	50	40
Initial angular sampling (°)	7.5	7.5	3.7	3.7
Local scarches from (°)	1.8	1.8	0.5	0.5
Initial offset range (pixel)	6	6	6	6
Initial offset step (pixel)	1	1	1	1
RELION results				
Wall-clock time (hr)	13.6	2.0	8.2	41.5
Reported resolution (Å)	9.8	8.2	7.3	5.6
Resolution vs X-ray (Å)	10.1	8.4	7.3	$4.4^{b}$
Previous results				
Refinement program	$XMIPP^{c}$	$EMAN2^{d}$	MRC	$FREALIGN^{e}$
Reported resolution (Å)	13.9	8.4	7.4	$\approx 6$
Resolution vs X-ray (Å)	12.7	8.7	7.5	$4.4^{b}$

### 3D auto-refine results



## More exciting RELION results

• DNA-origami object @ 11.5 Å resolution

– See poster (Xiao-chen Bai)

## Conclusions

- 3D-EM reconstruction is ill-posed, incomplete inverse problem
   Needs: regularization, marginalization and iteration
- Initial model generation & classification remain problematic in some projects
- Overfitting may be avoided w/o loss of reconstruction quality
   Use gold-standard FSCs, or high-res limited refinement!
- Bayesian framework provides a firm theoretical basis for 3D-EM
  - Learns optimal parameters from the data
  - Very little user input -> objective and easy-to-use
  - Excellent quality reconstructions

## Acknowledgements

- HepB data
  - Tony Crowther
  - Greg McMullan
- GroEL data
  - Steven Ludtke
- 70S Ribosome data
  - Haixiao Gao
  - Joachim Frank
- β-galactosidase data
  - Shaoxia Chen
  - Richard Henderson
- Rotavirus data
  - James Chen
  - Niko Grigorieff

- 80S ribosome data
  - Xiaochen Bai
  - Israel Sanchez
  - Venki Ramakrishnan
- Computing
  - Jake Grimmett
  - Toby Darling
- Some code in RELION
  - Xmipp (Carazo et al.)
  - Bsoft (Heymann et al.)
- Discussions
  - LMB colleagues
- Funding

ARC Laboratory of Molecular Biology