

# Fast algorithms for *in situ* single protein matching in Cryo tomography

Manas Rachh<sup>1</sup>, Alex Barnett<sup>1</sup>, Leslie Greengard<sup>1</sup>, Nikolaus Grigorieff<sup>2</sup>

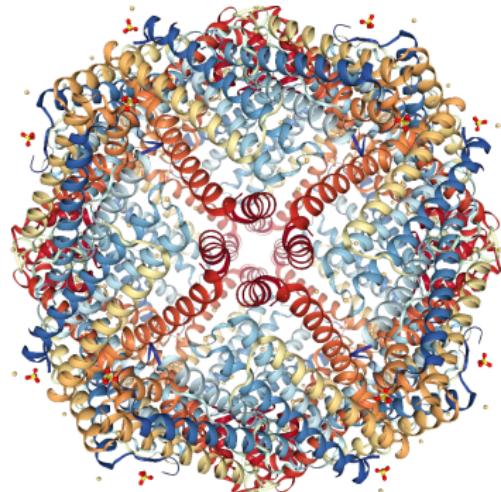
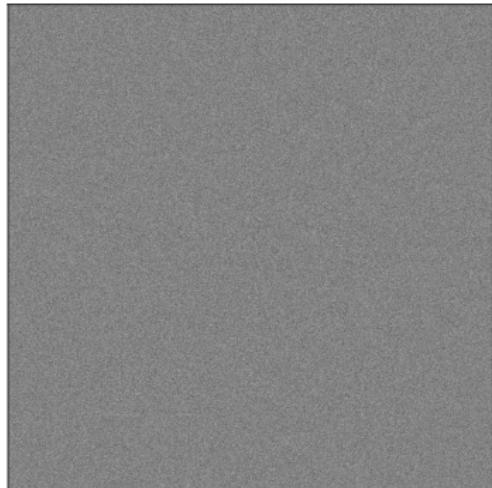
Nov 28, 2018

<sup>1</sup>Center for Computational Mathematics, Flatiron Institute, Simons Foundation

<sup>2</sup>Howard Hughes Medical Institute

## Problem statement

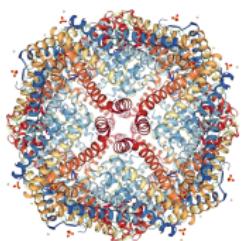
**Task:** Find orientation and location of proteins in cryo-tomograms using template matching



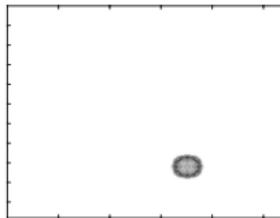
(left) A simulated tomogram with  $\text{SNR} = 0.025$ , (right) Apoferritin molecule (source: rscb pdb)

# Maximum likelihood noise model - I

$$R \in \mathbf{SO}(3)$$

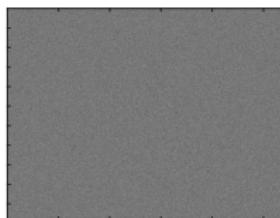


$$T_R(\mathbf{x}) = \int V(\tilde{x}_R, \tilde{y}_R, \tilde{z}_R) dz$$



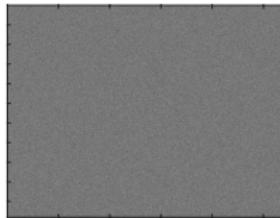
$$T_R(\mathbf{x} + \mathbf{x}_0)$$

+



$$\eta(\mathbf{x})$$

=



$$I(\mathbf{x})$$



## Maximum likelihood noise model - II

$$I(\mathbf{x}) = T_R(\mathbf{x} + \mathbf{x}_0) + \eta(\mathbf{x}), \quad \eta \sim \mathcal{N}(0, \sigma^2), \text{ iid}$$

$$p_{\text{lik}}(I|\mathbf{x}_0, R) = \prod_{\mathbf{x}} \frac{1}{2\pi\sigma^2} \exp(-|T_R(\mathbf{x} + \mathbf{x}_0) - I(\mathbf{x})|^2/(2\sigma^2))$$

$$\hat{\mathbf{x}}_0, \hat{R} = \operatorname{argmin}_{\mathbf{x}_0, R} (-\log p_{\text{lik}}) = |I|^2 + |T_R|^2 - 2 \langle I, T_R(\cdot + \mathbf{x}_0) \rangle,$$

where

$$\langle I, T_R(\cdot + \mathbf{x}_0) \rangle = \int I(\mathbf{x}) T_R(\mathbf{x} + \mathbf{x}_0)$$

Assume:  $|T_R|^2 \approx \text{const}$

$$\hat{\mathbf{x}}_0, \hat{R} = \operatorname{argmin}_{\mathbf{x}_0, R} (-\log p_{\text{lik}}) = \operatorname{argmax}_{\mathbf{x}_0, R} \langle I, T_R(\cdot + \mathbf{x}_0) \rangle$$

Maximum likelihood estimator is the maximizer of the cross-correlation

## Computational cost

- $N$  - number of pixels in image  $I \sim 1080^2$  for numerical experiments
- $M$  - number of pixels in template  $T \sim 360^2$  for apoferritin at  $0.965\text{\AA}$  resolution
- $N_r$  - number of rotations in  $\text{SO}(3) \sim 2 \times 10^6$  for  $\Delta\theta = 1.7^\circ$
- Fourier slice theorem for template generation  $O(N_r \cdot M)$

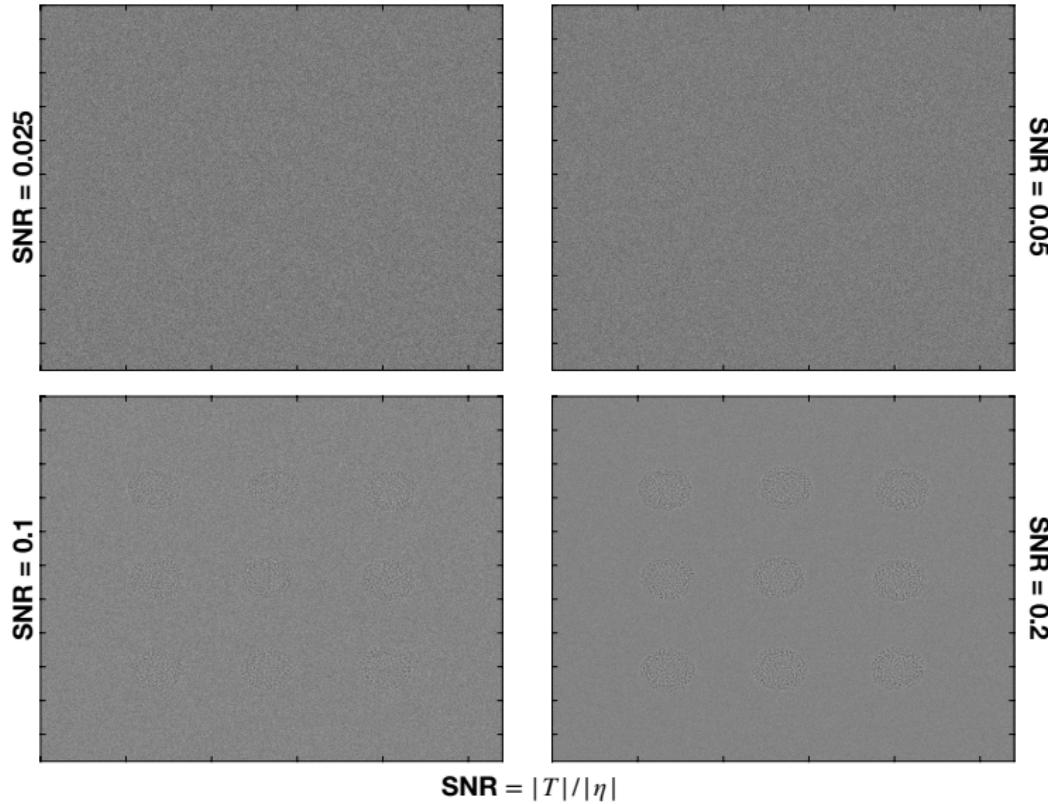
$$T_R(\mathbf{x}) = \int V(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}}) dz$$

- FFTs for computing cross correlations  $O(N_r \cdot N)$

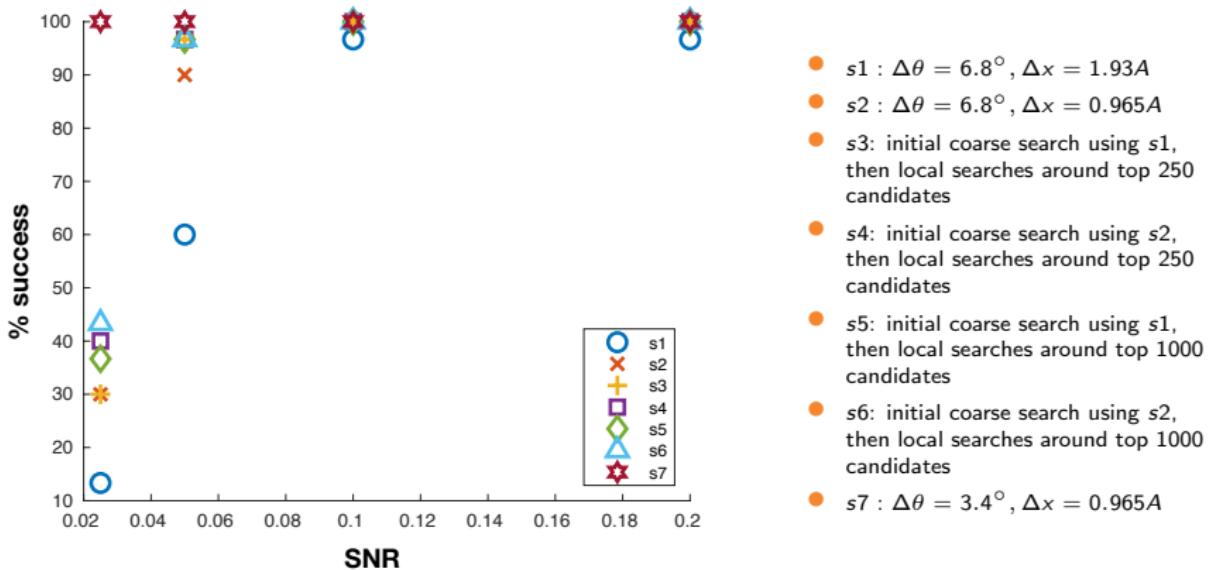
$$\langle I, T_R(\cdot + \mathbf{x}_0) \rangle \quad \text{for all } \mathbf{x}_0$$

- 100 CPU hours per template per image
- **Proposed acceleration: Frequency marching**  
Low resolution search followed by local searches at higher resolution

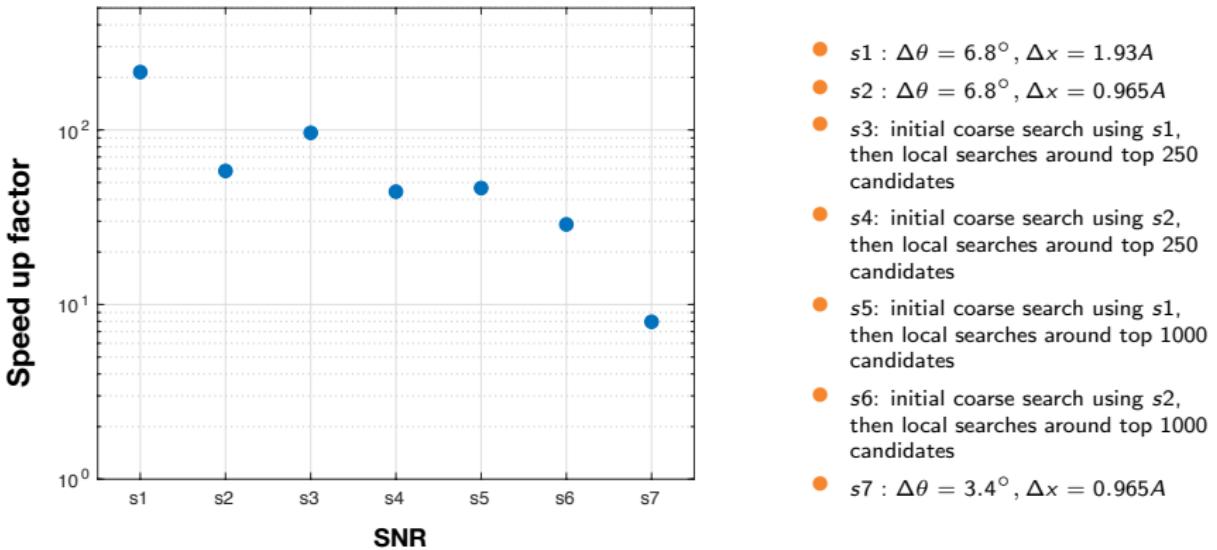
## Typical images



# Preliminary results - Accuracy



## Preliminary results - Speed



## Ongoing work

- Faster template generation using NUFFT
- Optimization of parameters for local searches
- Exploiting symmetries of proteins
- Better maximum likelihood estimators
- Performance in crowded protein environments
- Using tilt series information available to improve SNR
- Automated and robust pipeline for simultaneous detection of collection of proteins

## Ongoing work

- Faster template generation using NUFFT
  - Optimization of parameters for local searches
  - Exploiting symmetries of proteins
  - Better maximum likelihood estimators
  - Performance in crowded protein environments
  - Using tilt series information available to improve SNR
  - Automated and robust pipeline for simultaneous detection of collection of proteins
- 
- Questions?

Thank you