



# Cryo-EM Pre-Processing at Full Warp

NeCEN workshop 10/2018

# Data pre-processing, the old way



Acquisition



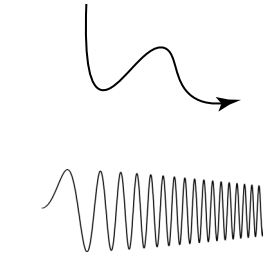
Terabytes



Storage



Days

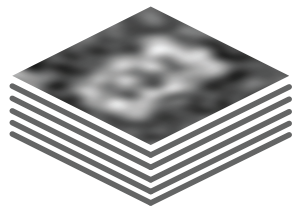


Motion correction,  
CTF estimation

All micrographs move in one chunk, one step at a time



Hours



Particles



Days



Autopicking &  
micrograph  
inspection



Hours

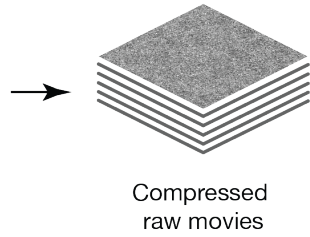
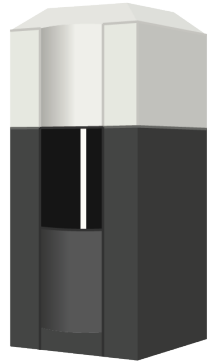


Manual picking



## Acquisition

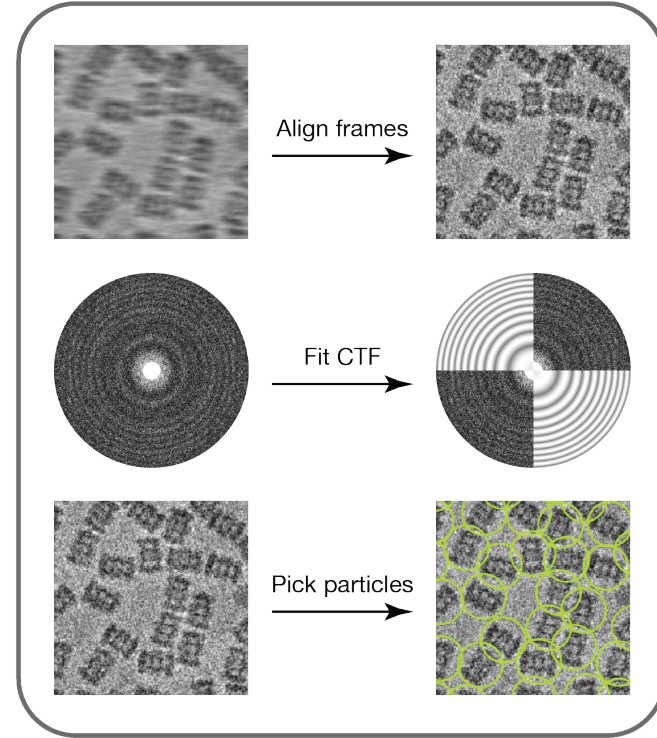
Automated in SerialEM, EPU



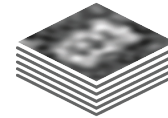
Continuously import

## Pre-processing

Automated in Warp,  $\approx 40$  s per item, results updated continuously as new data arrive



Extract & export

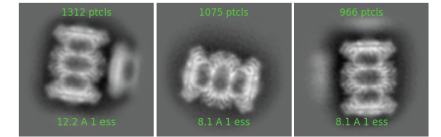


Particles,  
CTF values

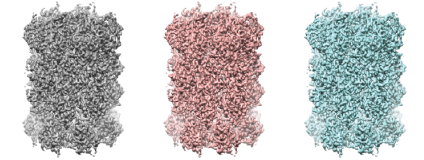
Continuously import

## Processing

Semi-automated in cryoSPARC



2D classification



3D classification, refinement

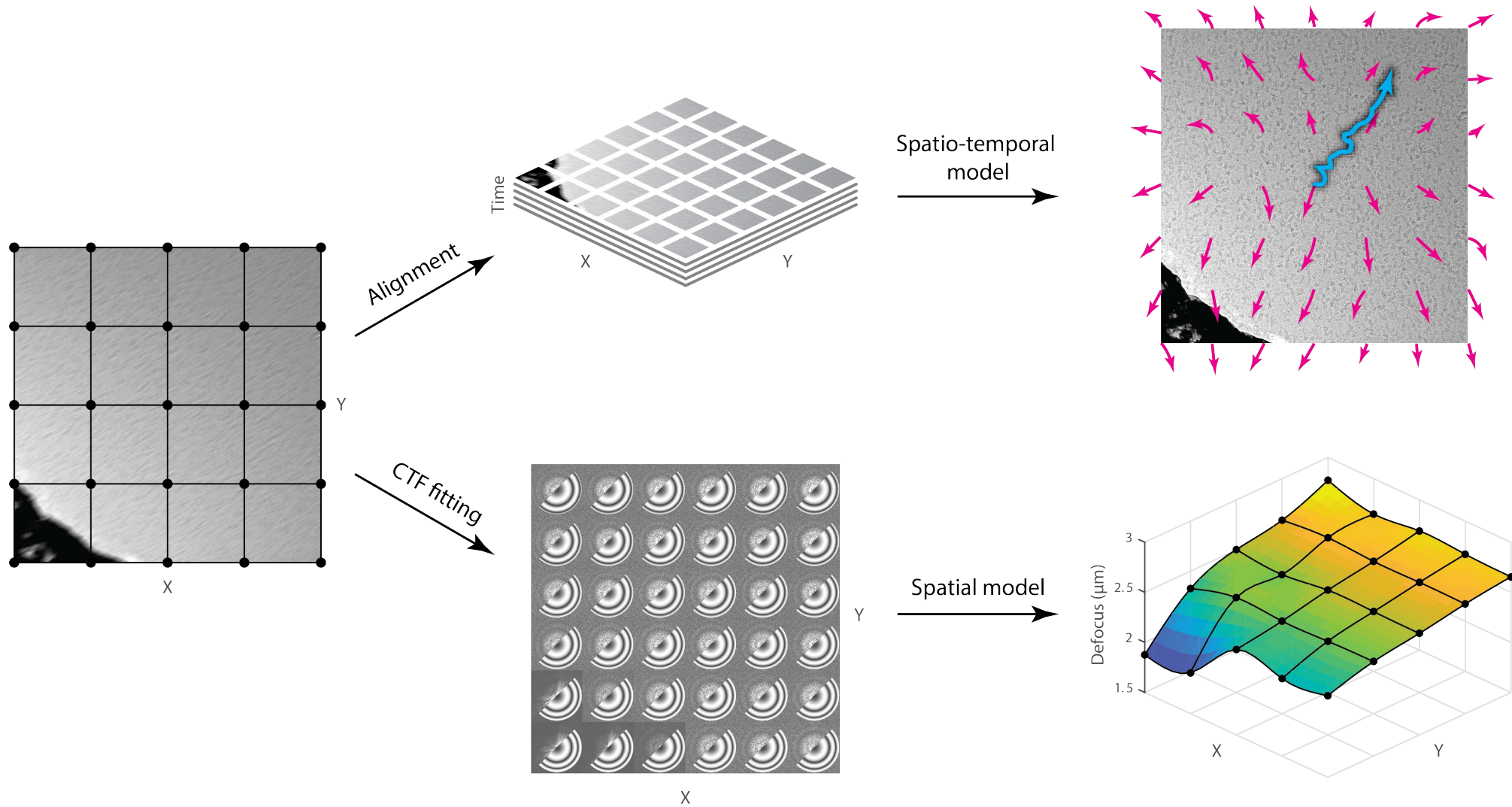


抖音

抖音号:902760646



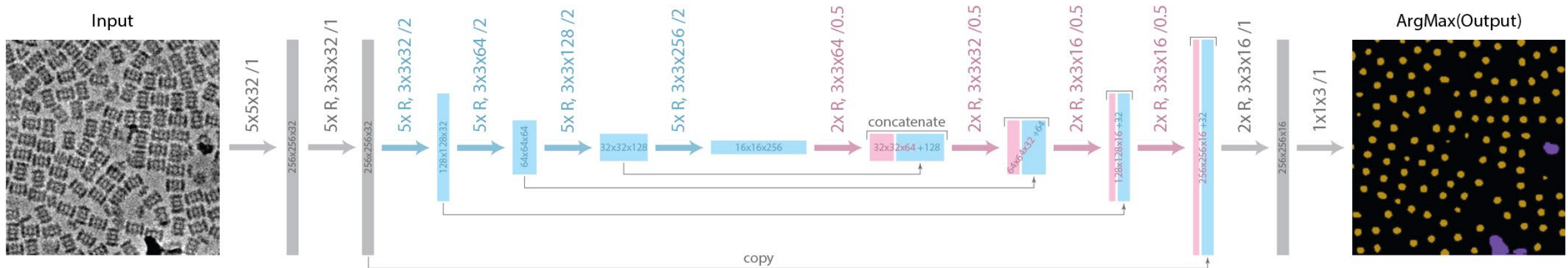
# Modeling sample deformation



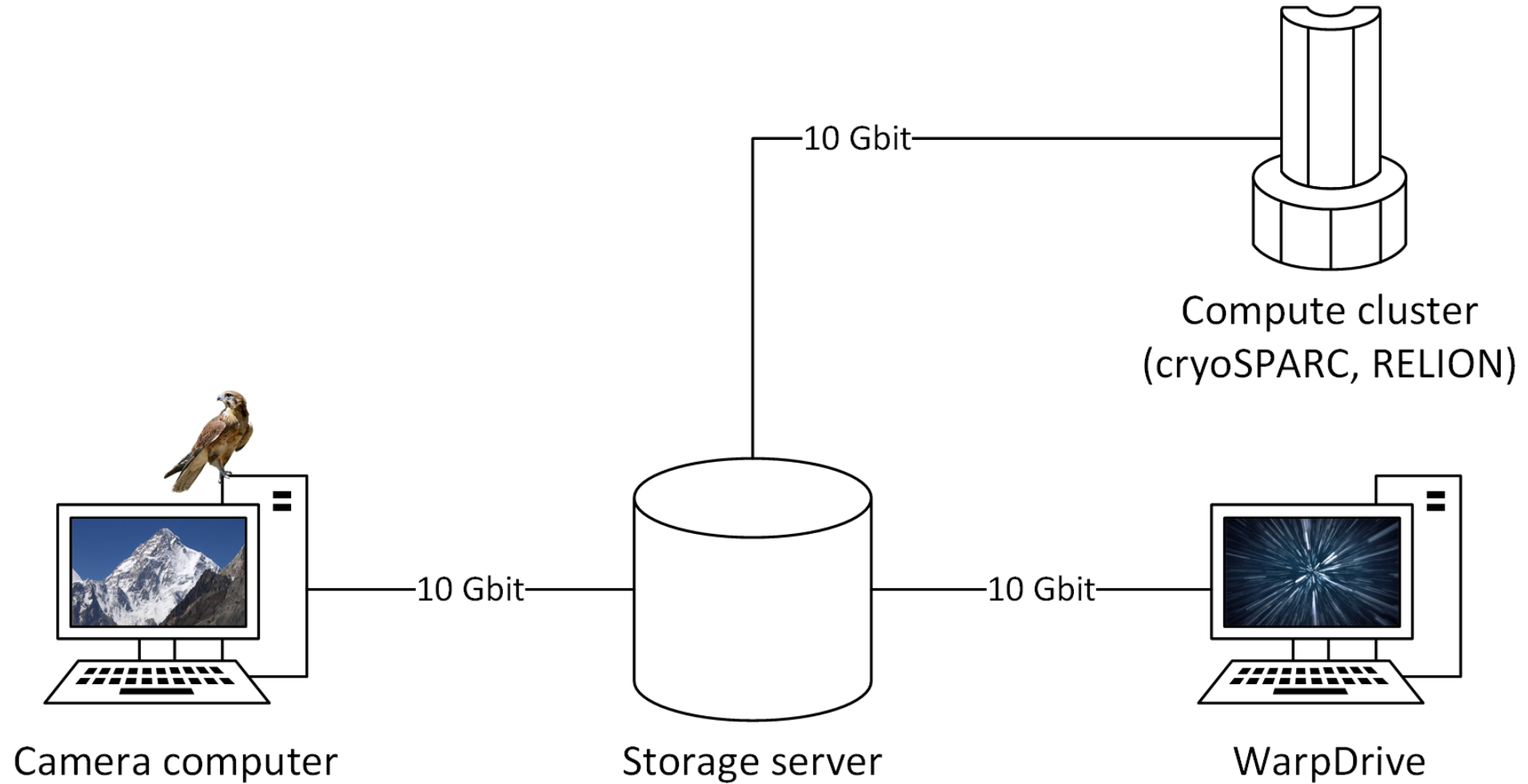


# Deep learning-based picking: BoxNet

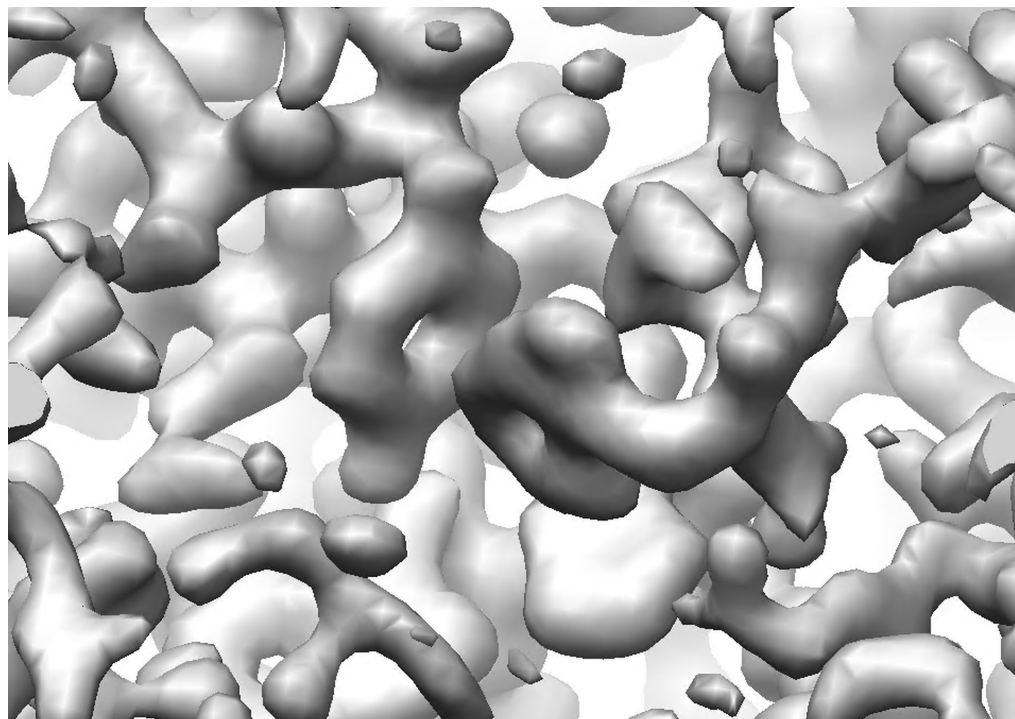
- Fully convolutional U-Net with residual blocks
- Pre-trained on 26 hand-picked data sets
- Easily retrainable within Warp
- Save and manage retrained models for projects



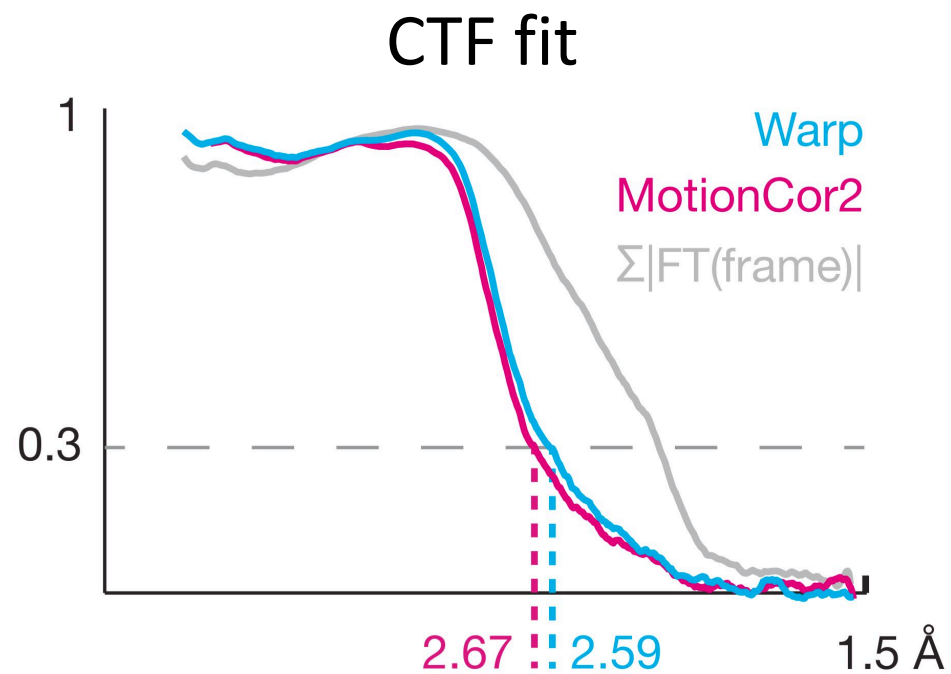
# Infrastructure



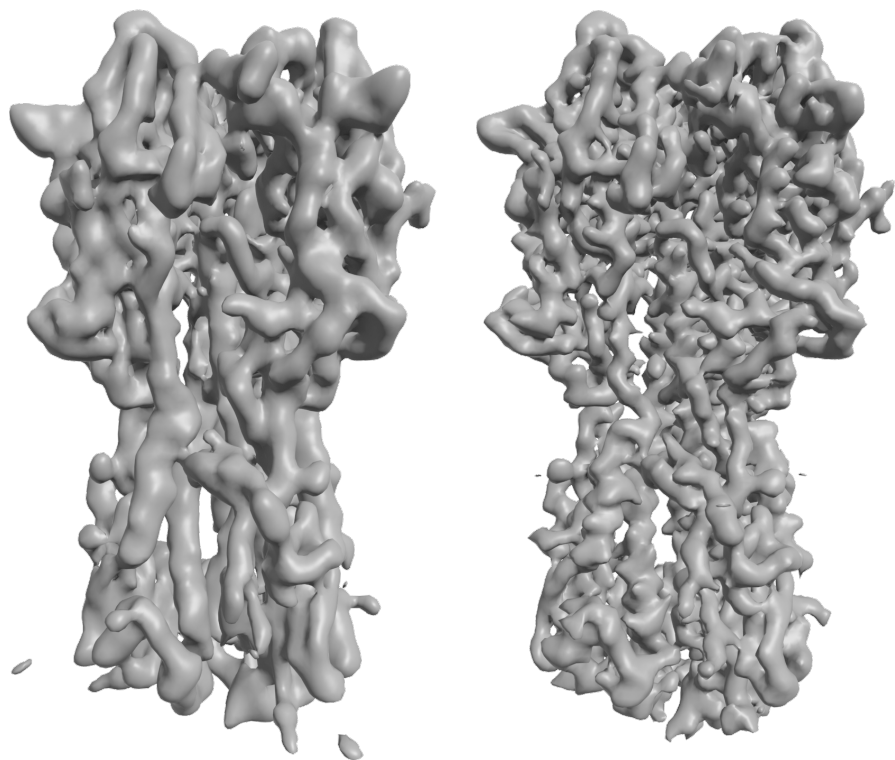
With EMPIAR-10061



Warp + RELION: 2.09 Å  
+ Beam tilt: 1.95 Å  
+ Defocus: 1.95 Å  
+ Polishing: 1.86 Å

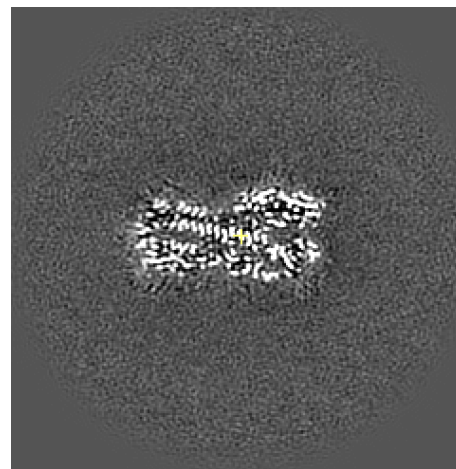
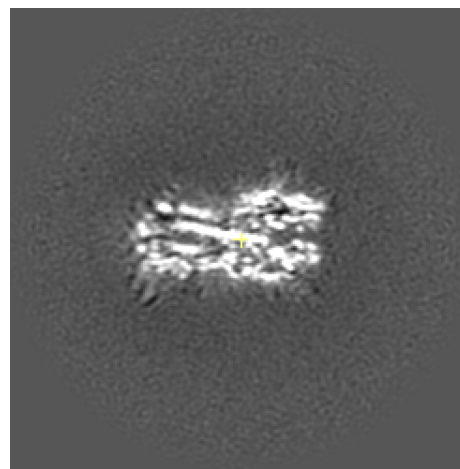


# With EMPIAR-10097

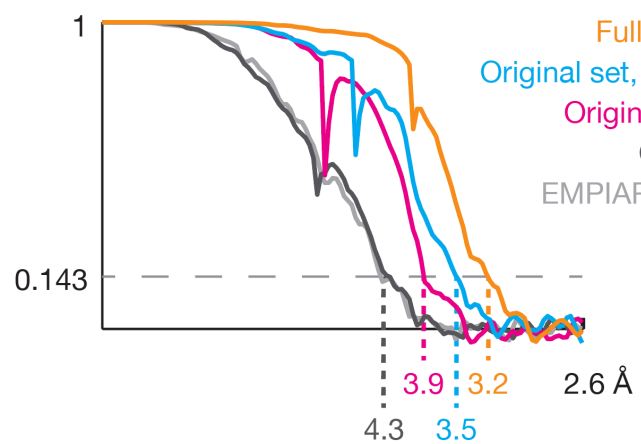


Original data,  
manual processing

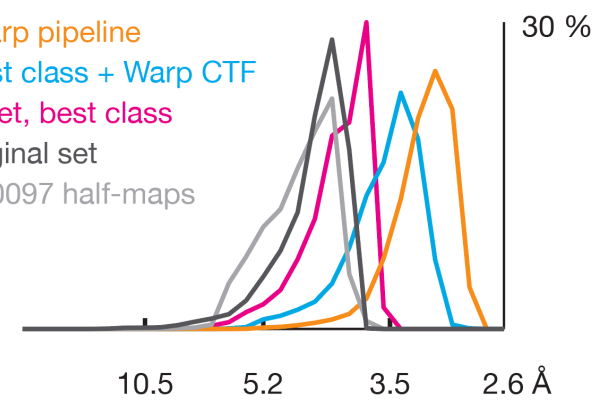
Warp pipeline,  
automated



Fourier shell correlation



Local resolution





Get it at [warpem.com](http://warpem.com)!

# Denoising with deep neural nets

I have no idea why this works.

---

# Noise2Noise: Learning Image Restoration without Clean Data

---

Jaakko Lehtinen<sup>1,2</sup> Jacob Munkberg<sup>1</sup> Jon Hasselgren<sup>1</sup> Samuli Laine<sup>1</sup> Tero Karras<sup>1</sup> Miika Aittala<sup>3</sup>  
Timo Aila<sup>1</sup>

## Abstract

We apply basic statistical reasoning to signal reconstruction by machine learning — learning to map corrupted observations to clean signals — with a simple and powerful conclusion: under certain common circumstances, it is possible to learn to restore signals without ever observing clean ones, at performance close or equal to training using clean exemplars. We show applications in photographic noise removal, denoising of synthetic Monte Carlo images, and reconstruction of MRI scans from undersampled inputs, all based on only observing corrupted data.

have been reported in several applications, including Gaussian denoising, de-JPEG, text removal (Mao et al., 2016), super-resolution (Ledig et al., 2017), colorization (Zhang et al., 2016), and image inpainting (Iizuka et al., 2017). Yet, obtaining clean training targets is often difficult or tedious. A noise-free photograph requires a long exposure; full MRI sampling is slow enough to preclude dynamic subjects, etc.

In this work, we observe that under suitable, common circumstances, we can *learn to reconstruct signals from only corrupted examples, without ever observing clean signals*, and often do this just as well as if we were using clean examples. As we show below, our conclusion is almost trivial from a statistical perspective, but in practice, it significantly eases learning signal reconstruction by lifting requirements on the availability of clean data.

## 1. Introduction

Signal reconstruction from corrupted or incomplete measurements is an important subfield of statistical data analysis

## 2. Theoretical background

Assume that we have a set of unreliable measurements

---

# Noise2Noise: Learning Image Restoration without Clean Data

---

Jaakko Lehtinen<sup>1,2</sup> Jacob Munkberg<sup>1</sup> Jon Hasselgren<sup>1</sup> Samuli Laine<sup>1</sup> Tero Karras<sup>1</sup> Miika Aittala<sup>3</sup>  
Timo Aila<sup>1</sup>

## Abstract

We apply basic statistical reasoning to the problem of signal reconstruction by machine learning — we map corrupted observations to clean ones with a simple and powerful convolutional neural network. In certain common circumstances, we learn to restore signals without clean ones, at performance close to that achieved by training using clean exemplars. We show this in photographic noise removal, denoising synthetic Monte Carlo images, and reconstructing MRI scans from undersampled input, all on only observing corrupted data.



...reported in several applications, including Gaussian denoising (Saito et al., 2017), denoising of JPEG images (Ledig et al., 2017), text removal (Mao et al., 2016), image inpainting (Iizuka et al., 2017). Yet, learning targets is often difficult or tedious. For example, photograph requires a long exposure; full MRI scan requires a long time to preclude dynamic subjects, etc. We observe that under suitable, common circumstances, we can learn to reconstruct signals from only corrupted observations, without ever observing clean signals, just as well as if we were using clean exemplars. In the following, our conclusion is almost trivial from a theoretical perspective, but in practice, it significantly eases learning signal reconstruction by lifting requirements on the availability of clean data.

## 1. Introduction

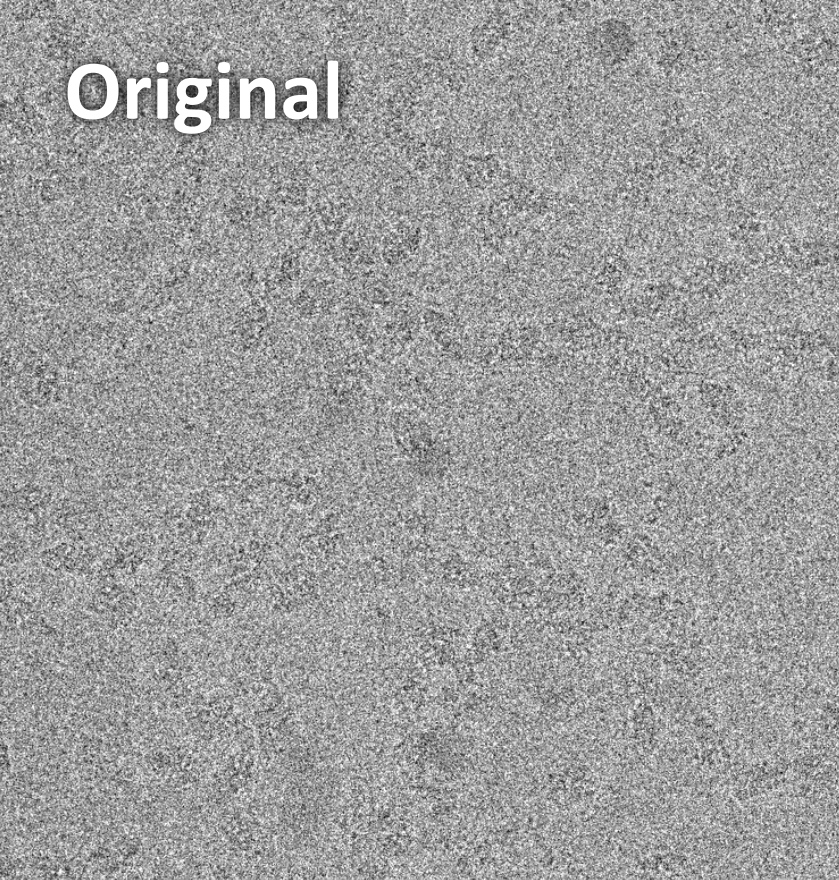
Signal reconstruction from corrupted or incomplete measurements is an important subfield of statistical data analysis

## 2. Theoretical background

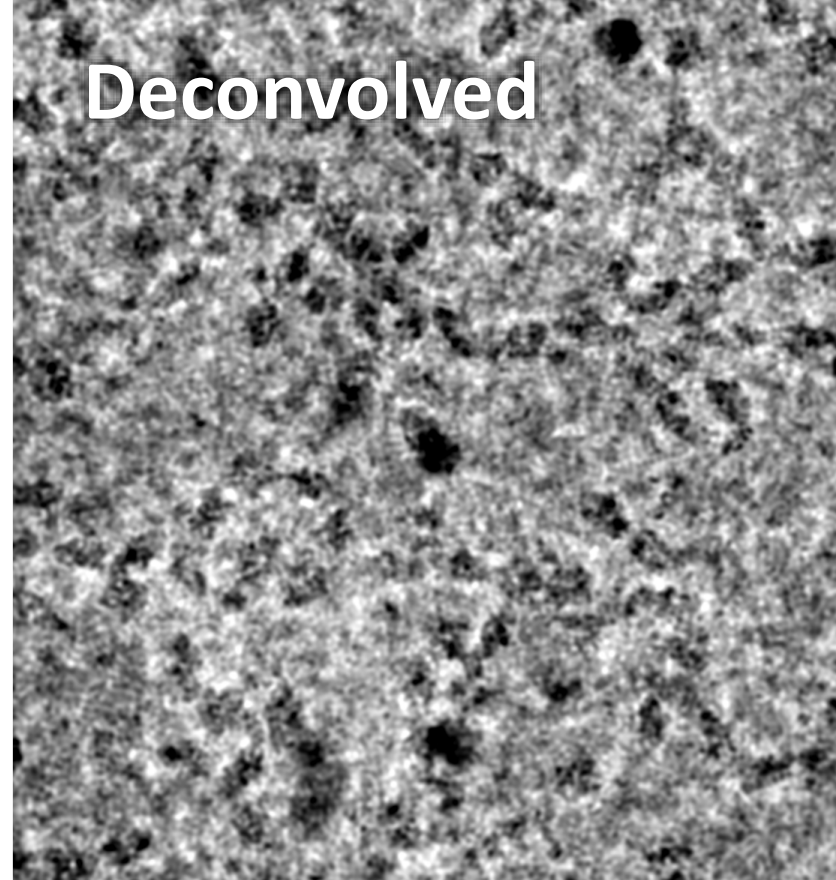
Assume that we have a set of unreliable measurements



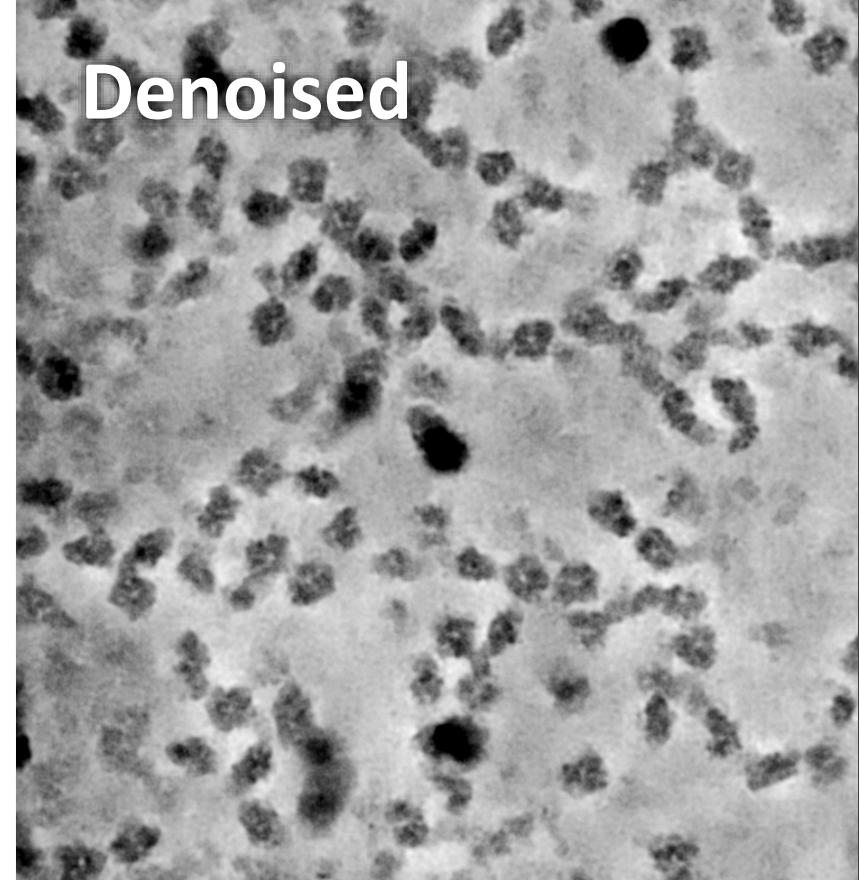
Original



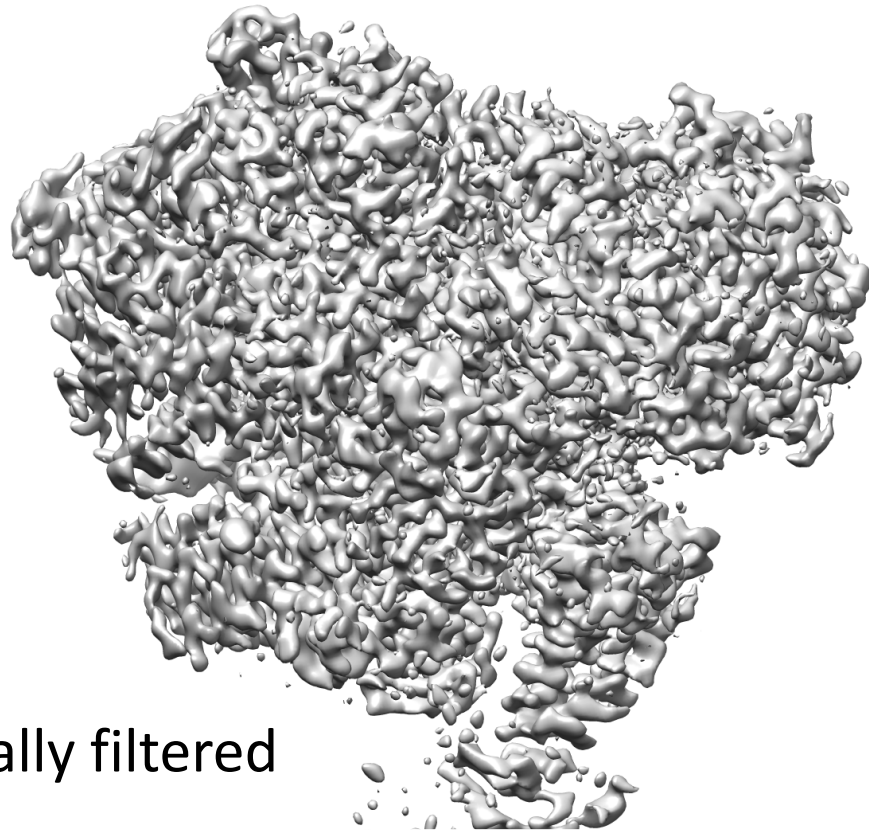
Deconvolved



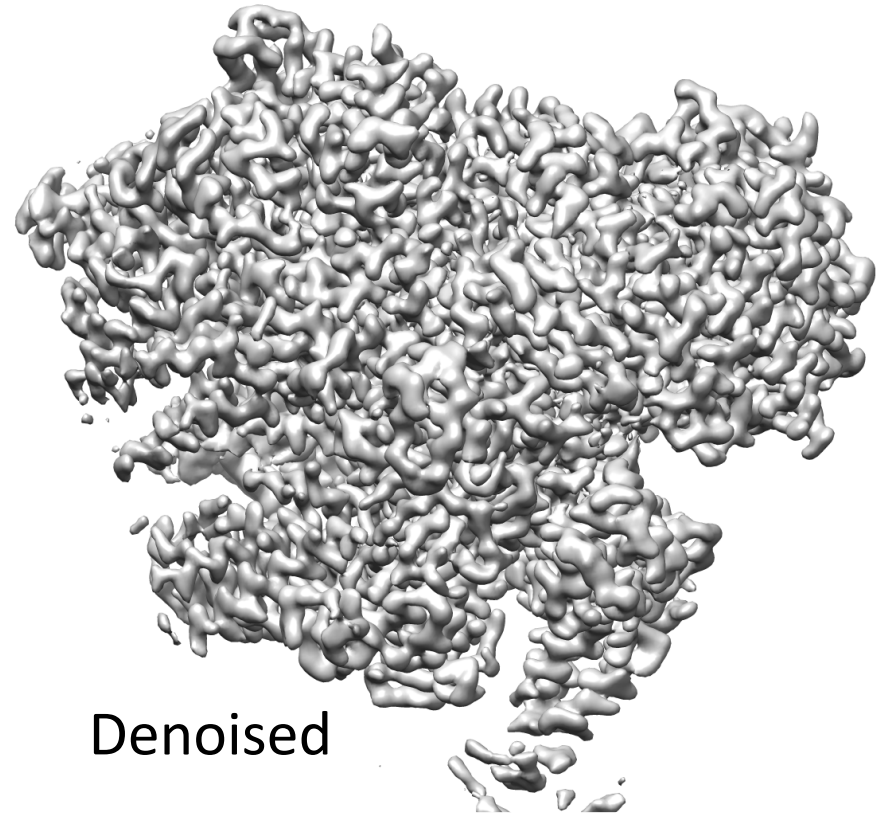
Denoised



Denoising: 2D (0.8  $\mu\text{m}$  defocus)

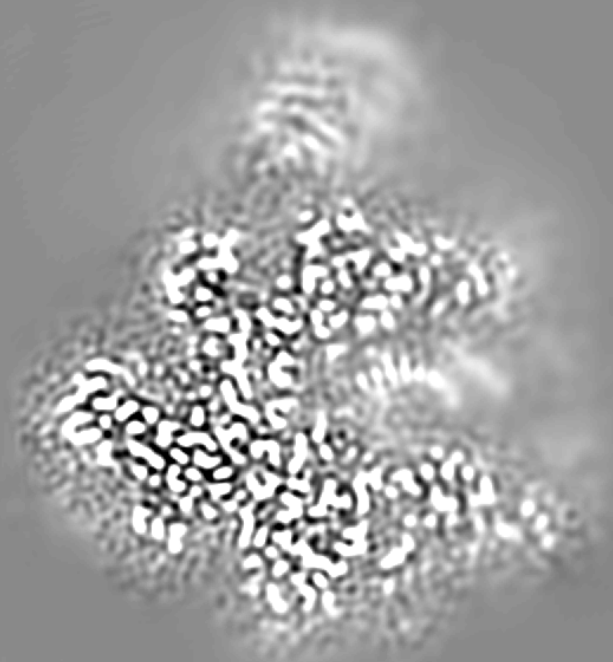


Locally filtered

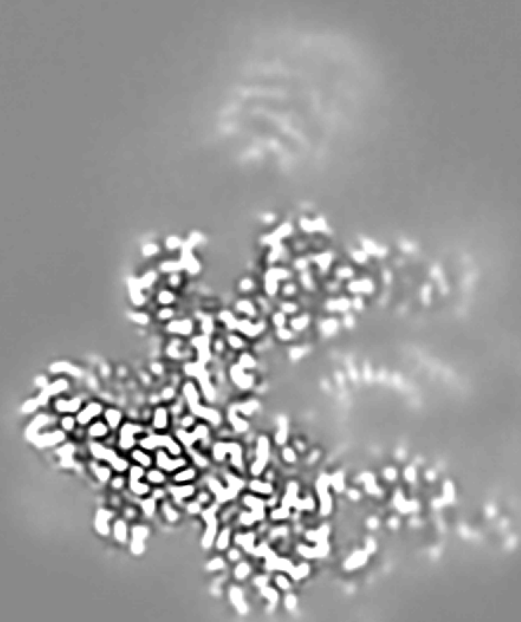


Denoised

Denoising: 3D half-maps



Locally filtered



Denoised

Denoising: 3D half-maps





Denoising: 3D *in situ* tomograms



Where do we want to go?  
Tomography.



Ben Engel



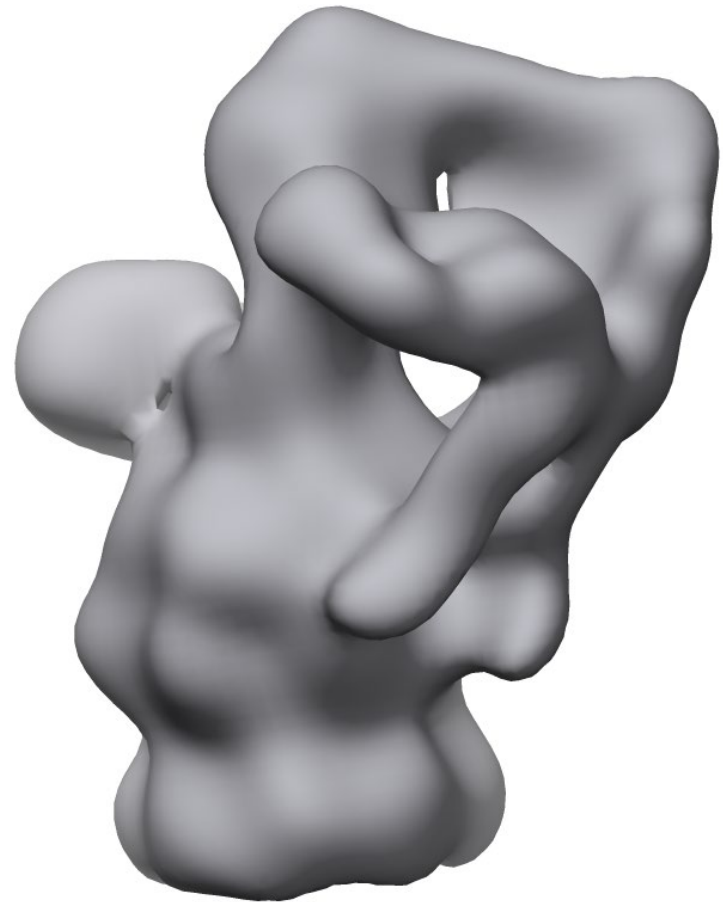
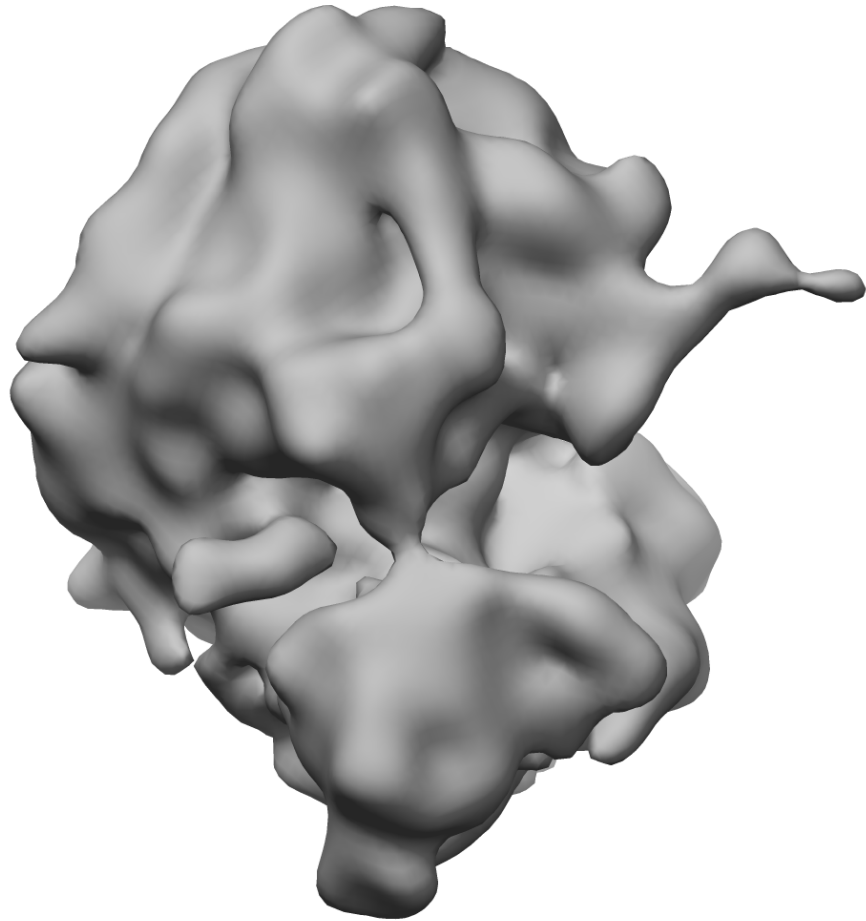






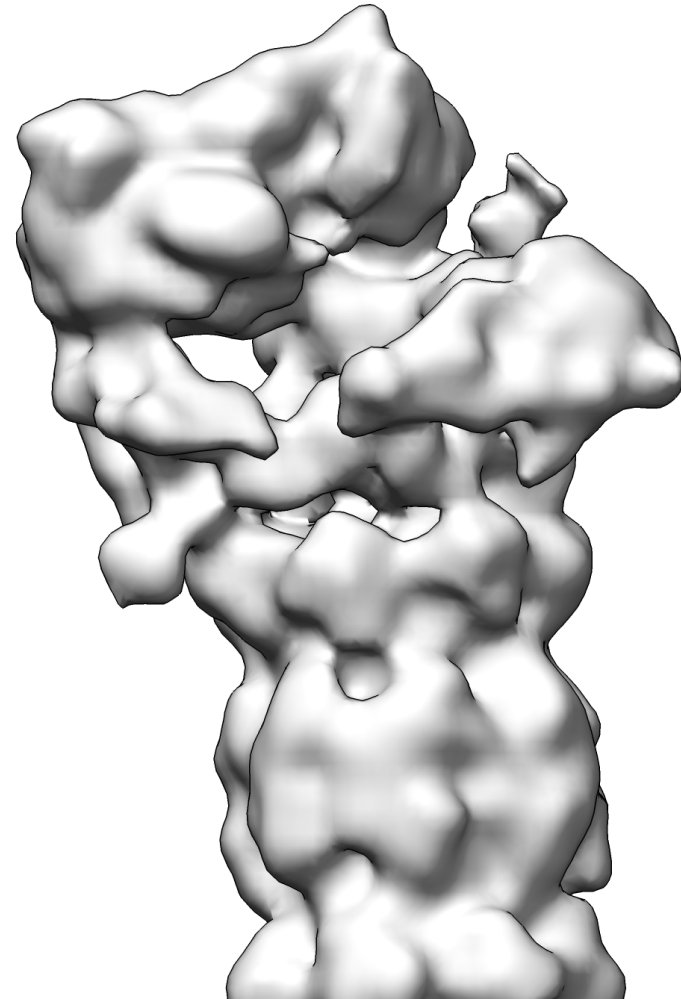
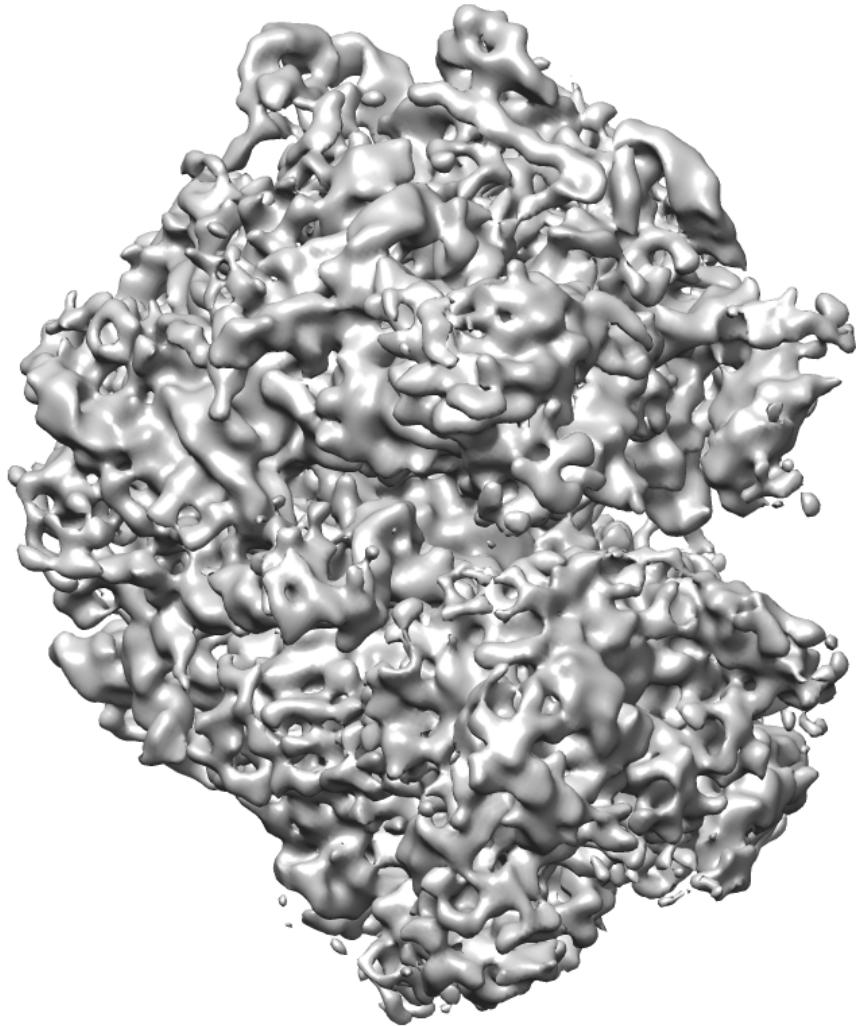


Previous results for cryo-FIBed *in situ*





Similar *in situ* data with Warp



# Implications for data processing and sharing

- Every tomogram needs as many identified particles as possible
- Every particle species needs as many copies as possible
- Every *in situ* tomogram contains less than 1 particle of interest per lab
- ... and 10 000+ particles of interest for other labs
- Everyone's resolution increases as more particles are added
- No single lab/facility will be able to produce enough data



Thanks to:  
Patrick Cramer

Data from:  
EMPIAR  
Ben Engel  
Carrie Bernecky

[warpem.com](http://warpem.com)  
[github.com/dtegunov](https://github.com/dtegunov)  
[twitter.com/dtegunov](https://twitter.com/dtegunov)