

# Automatic particle picking with minimal supervision using positive-unlabeled neural networks



Massachusetts  
Institute of  
Technology

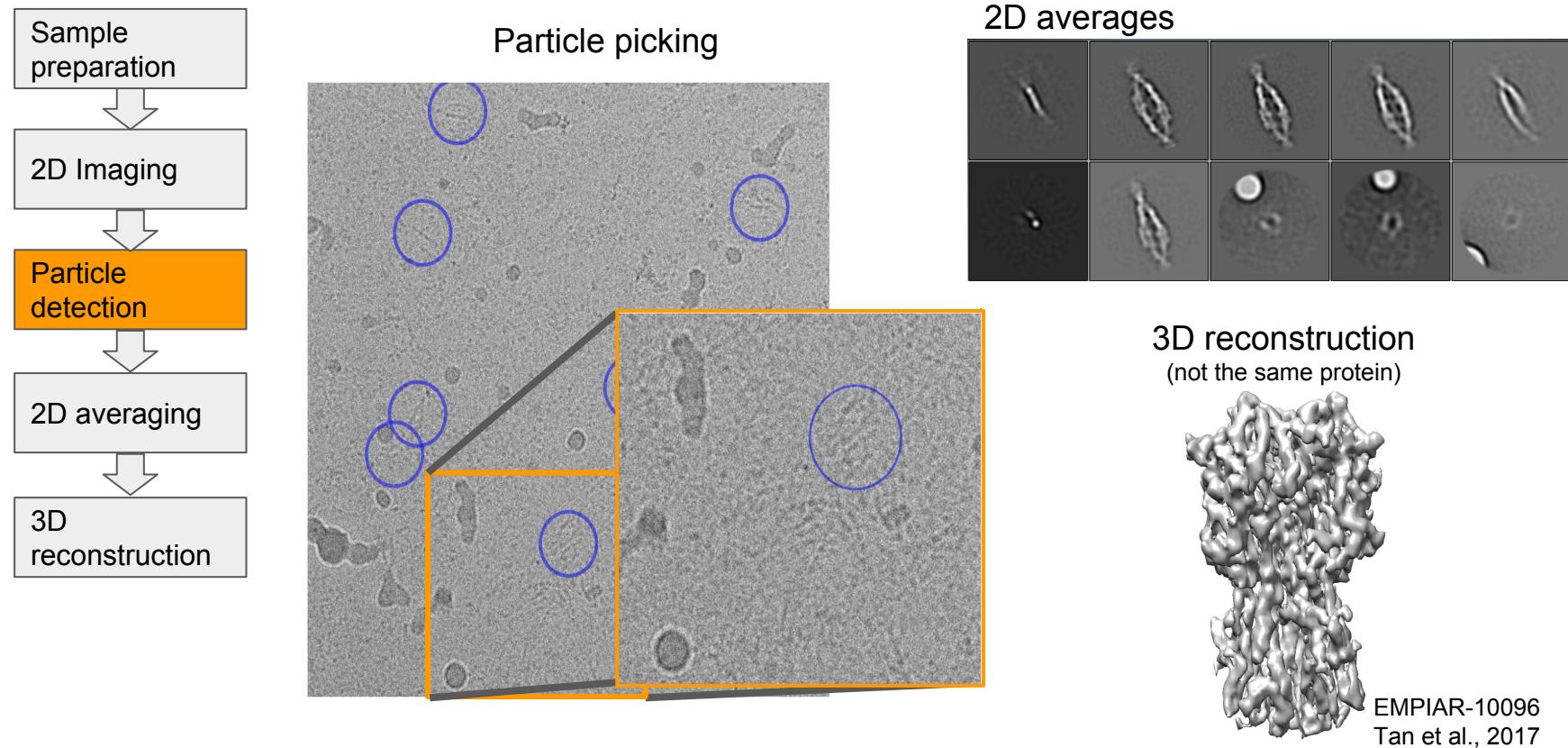
Tristan Bepler  
Berger Lab, MIT  
NYSBC Deep Learning Workshop  
April 10, 2018



SIMONS ELECTRON  
MICROSCOPY CENTER

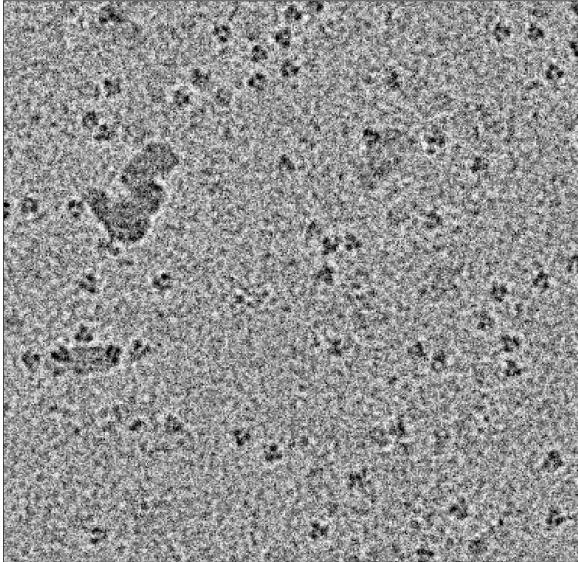


# Single-particle cryo-electron microscopy



# Particle picking

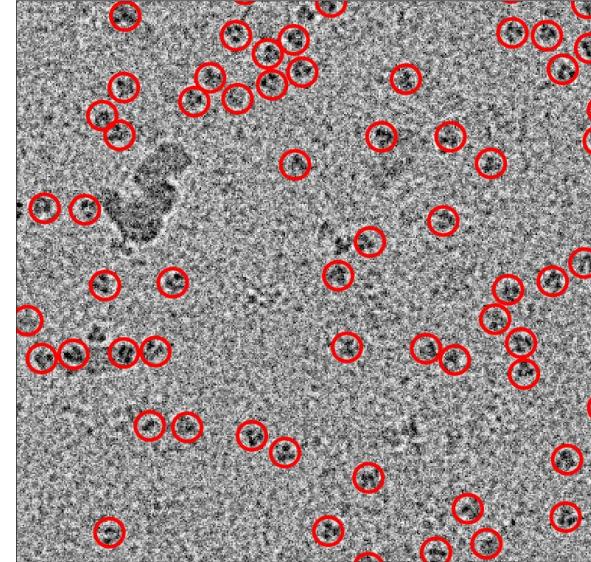
Micrograph (X)



?

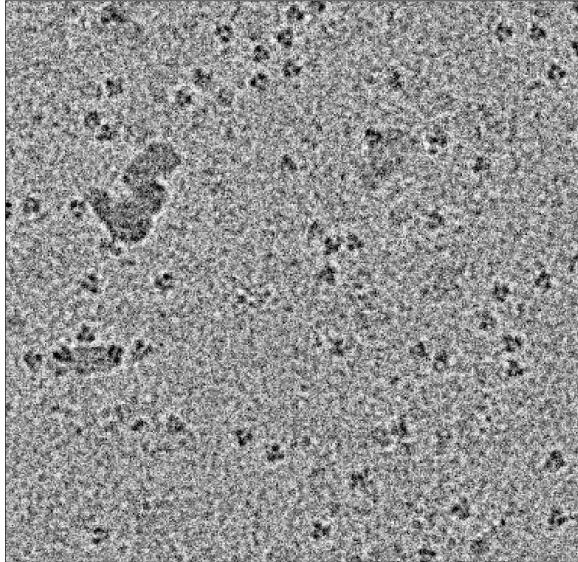


Particle coordinates

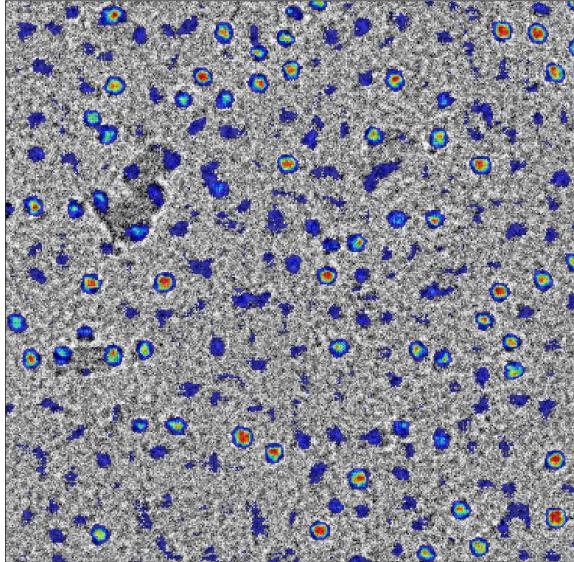


# Particle picking

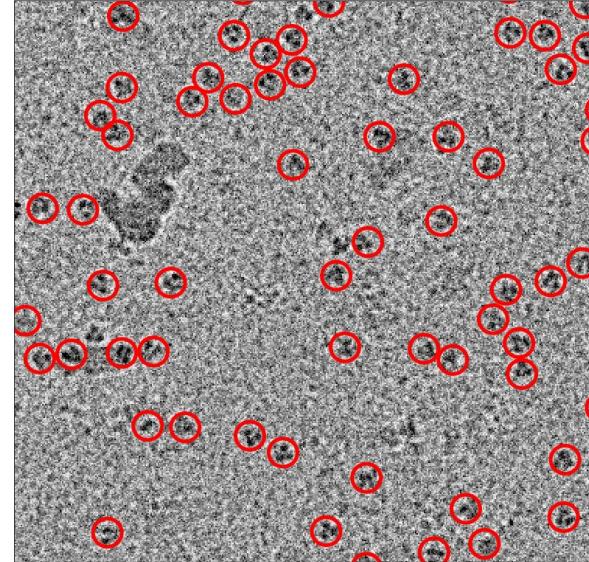
Micrograph ( $X$ )



Region scores ( $g * X$ )



Particle coordinates

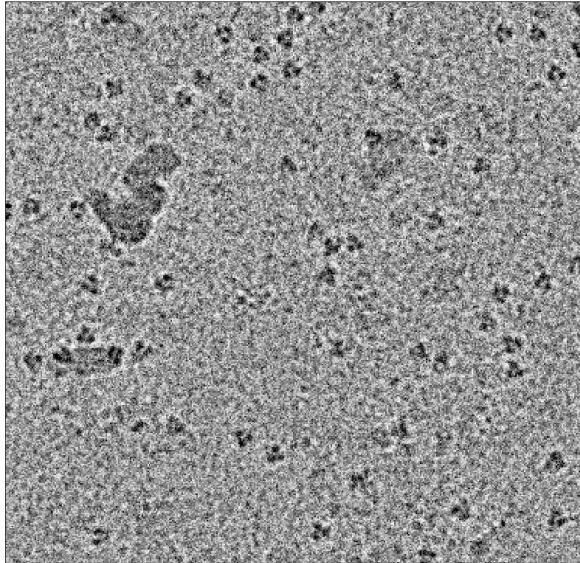


Given a scoring function,  $g$ , convolve it over the micrograph,  $X$ , to get per region scores

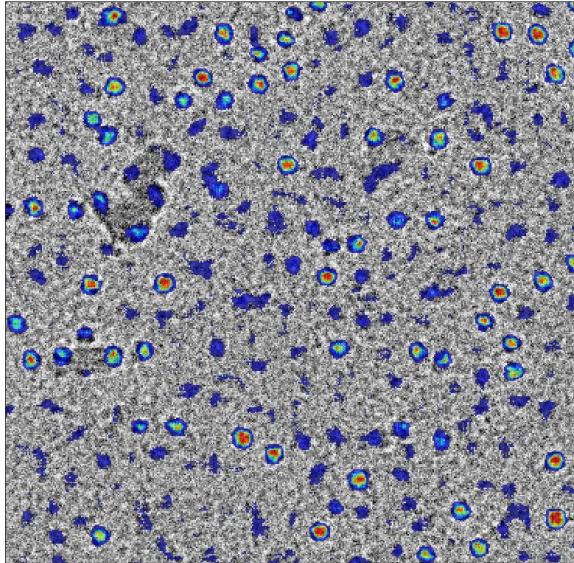
Extract coordinates by greedily selecting regions and removing nearby regions (non-maximum suppression)

# Particle picking

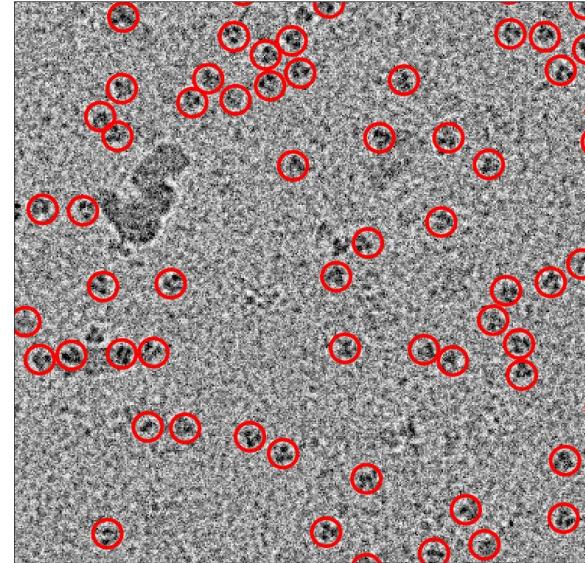
Micrograph ( $X$ )



Region scores ( $g * X$ )



Particle coordinates



What should  $g$  be?

# Convolutional neural network

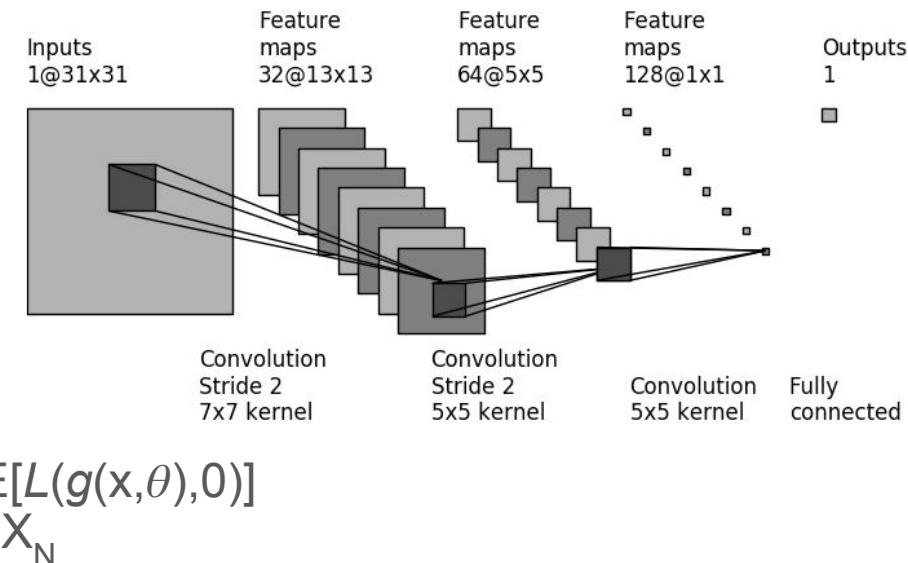
- Learn parameters,  $\theta$ , of  $g$  from data

Positives:  $X_P$  ,  $y=1$  ,  $\pi_P$

Negatives:  $X_N$  ,  $y=0$  ,  $1 - \pi_P$

Loss function  $L$

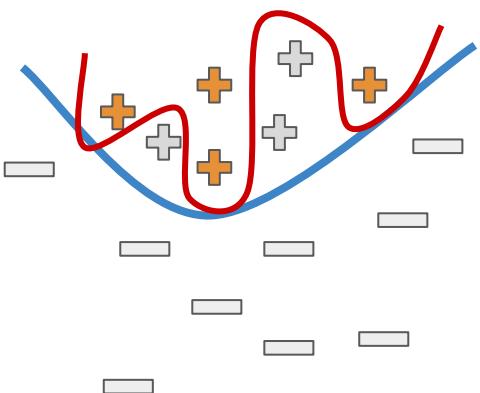
$$\operatorname{argmin}_{\theta} \pi_P E[L(g(x,\theta), 1)] + (1 - \pi_P) E[L(g(x,\theta), 0)]$$



- Problem: learning requires large amounts of labeled examples
  - Costly for a researcher to label enough particles
  - Can we learn  $\theta$  from a small amount of labeled data and the rest of the unlabeled data?

# Positive-unlabeled classification

- Learn parameters,  $\theta$ , from positive,  $X_P$ , and unlabeled,  $X_U$ , data
- Unlabeled data contains both positive and negative examples
- Loss function: find parameters that minimize this function of the training data
- Assume we know  $\pi_P$



Naive:

$$\pi_P E[L(g(x), 1)] + (1 - \pi_P) E[L(g(x), 0)]$$
$$X_P \qquad \qquad \qquad X_U$$

Unbiased estimator (du Plessis et al. 2016):

$$\pi_P E[L(g(x), 1)] - \pi_P E[L(g(x), 0)] + E[L(g(x), 0)]$$
$$X_P \qquad X_P \qquad X_U$$

Non-negative estimator (Kiryo et al. 2017):

$$\pi_P E[L(g(x), 1)] + \max\{ 0, E[L(g(x), 0)] - \pi_P E[L(g(x), 0)] \}$$
$$X_P \qquad X_U \qquad X_P$$

Generalized expectation criteria (KL) (Mann and McCallum 2010):

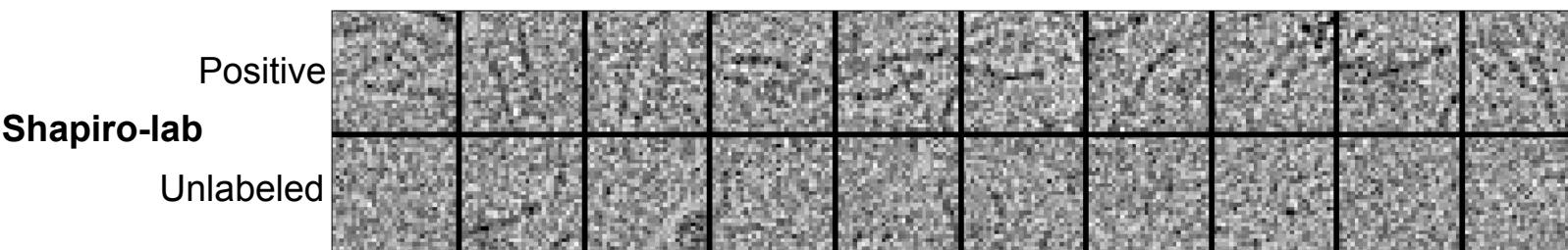
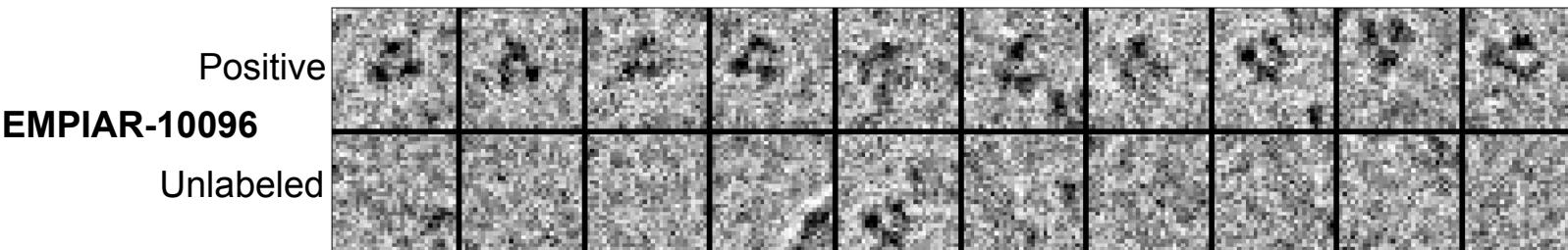
$$E[L(g(x), 1)] + \lambda \text{KL}(\pi_P, E[g(x)])$$
$$X_P \qquad X_U$$

# GE-binomial: a better GE criteria for positive-unlabeled learning with SGD

- Problem: neural network needs to be trained with **minibatch stochastic gradient descent (SGD)** - need to estimate gradient using samples of data
  - $E_{X_P}[L(g(x), 1)] + \lambda KL(\pi_P, E_{X_U}[g(x)])$
- The number of positive data points, P, in an N data point minibatch follows the **binomial distribution with probability of success  $\pi_P$** 
  - $p_k = \text{binomial}(N, \pi_P)$
- Classifier predictions,  $g(x_i)$  where  $x_i$  are unlabeled data points in the minibatch, also define a distribution over the number of positives - approximate this with a normal distribution
  - $\mu = \sum g(x_i)$  and  $\sigma^2 = \sum g(x_i)(1-g(x_i))$
  - let  $q_k$  be the discretized probability of k positives given by this distribution
- Define a **new GE criteria (GE-binomial)** using these distributions:  $\sum q_k \log(p_k)$   
 $E_{X_P}[L(g(x), 1)] + \sum q_k \log(p_k)$

# CryoEM datasets for evaluation

Dataset	Train		Test	
	Images	Particles	Images	Particles
<b>EMPIAR-10096</b>	347	100465	100	29535
<b>Shapiro-lab</b>	67	1167	20	373



# GE-binomial outperforms other positive-unlabeled learning objectives on cryoEM particle classification

Model	EMPIAR-10096			Shapiro-lab		
	10	100	1000	10	100	1167
<b>classifier</b>						
PN	0.072 ±0.029	0.187 ±0.038	-	0.012 ±0.005	0.036 ±0.012	-
NNPU	0.101 ±0.035	0.226 ±0.033	-	0.014 ±0.008	0.039 ±0.013	-
GE-KL	0.072 ±0.035	0.240 ±0.043	-	0.010 ±0.005	<b>0.062 ±0.017</b>	-
GE-binomial	<b>0.155 ±0.044</b>	<b>0.258 ±0.040</b>	<b>0.392 ±0.006</b>	<b>0.020 ±0.008</b>	0.061 ±0.010	<b>0.150 ±0.008</b>

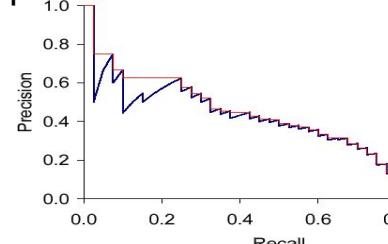
Area under the precision-recall curve on the test set for models trained with subsets of positives from the training set

PN = naive

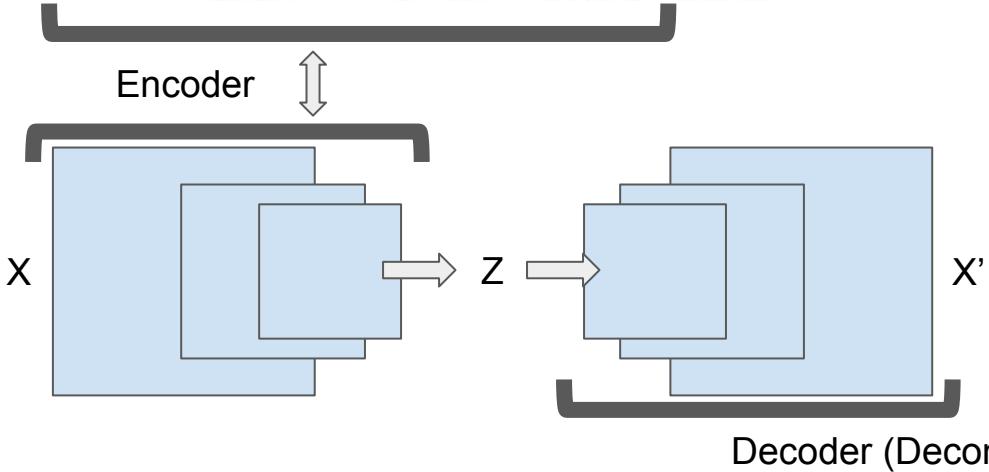
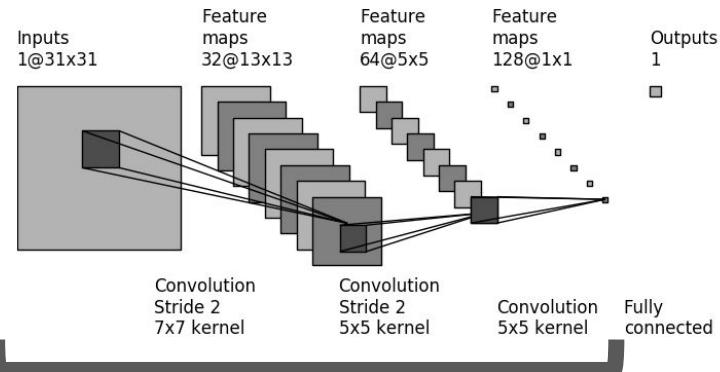
NNPU = non-negative estimator (Kiryo et al. 2017)

GE-KL = GE criteria with KL-divergence

$$\text{recall} = \text{TP}/(\text{TP} + \text{FN})$$
$$\text{precision} = \text{TP}/(\text{TP} + \text{FP})$$



# Hybrid classifier-autoencoder model



GE-binomial and autoencoder:

$$E[L(g(x), 1)] + \sum q_k \log(p_k) + a \|x - x'\|^2$$

$x_p$

GE-binomial  
loss

reconstruction  
error

# Including a decoder and reconstruction error improves generalization with few training examples

Model	EMPIAR-10096			Shapiro-lab		
	10	100	1000	10	100	1167
<b>classifier</b>						
PN	0.072 ±0.029	0.187 ±0.038	-	0.012 ±0.005	0.036 ±0.012	-
NNPU	0.101 ±0.035	0.226 ±0.033	-	0.014 ±0.008	0.039 ±0.013	-
GE-KL	0.072 ±0.035	0.240 ±0.043	-	0.010 ±0.005	<b>0.062 ±0.017</b>	-
GE-binomial	<b>0.155 ±0.044</b>	<b>0.258 ±0.040</b>	<b>0.392 ±0.006</b>	<b>0.020 ±0.008</b>	0.061 ±0.010	<b>0.150 ±0.008</b>
<b>+autoencoder</b>						
GE-binomial	0.260 ±0.016	0.324 ±0.017	0.368 ±0.013	0.029 ±0.011	0.078 ±0.018	0.120 ±0.006

PN = naive

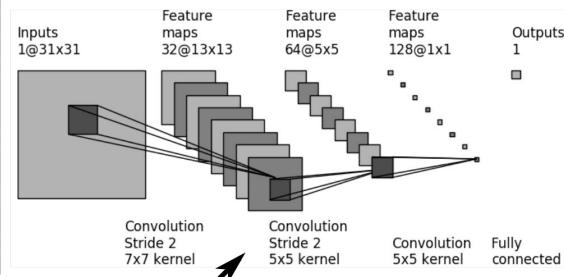
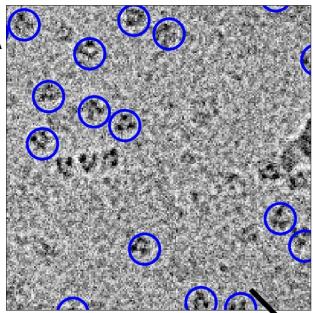
NNPU = non-negative estimator (Kiryo et al. 2017)

GE-KL = GE criteria with KL-divergence

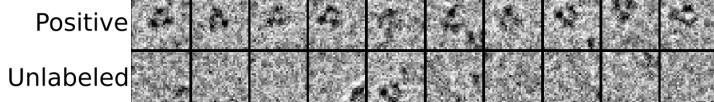
Adding a decoder ( $a = N/10$ ) can further improve classification performance when very few labeled data points are available

# Topaz particle picking pipeline

Train classifier with positive and unlabeled micrograph regions

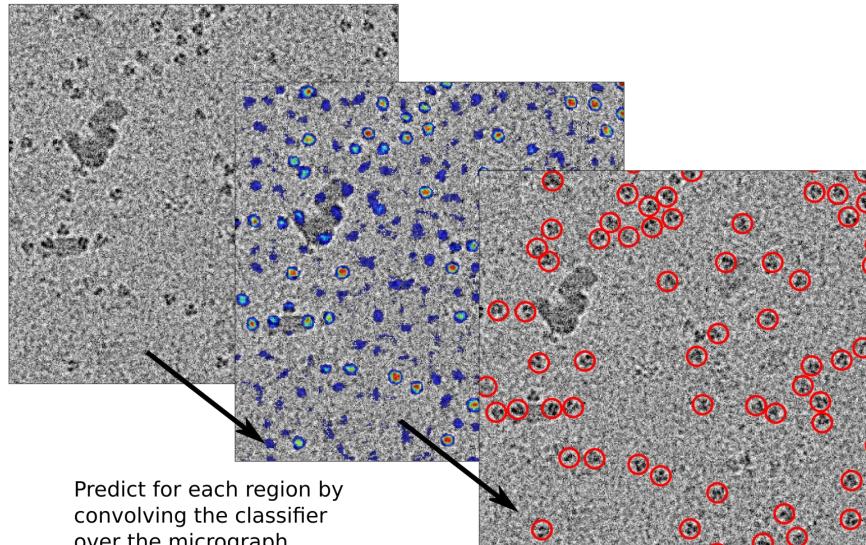


Region classes are taken from particle labels



Train CNN classifier with positive and unlabeled examples

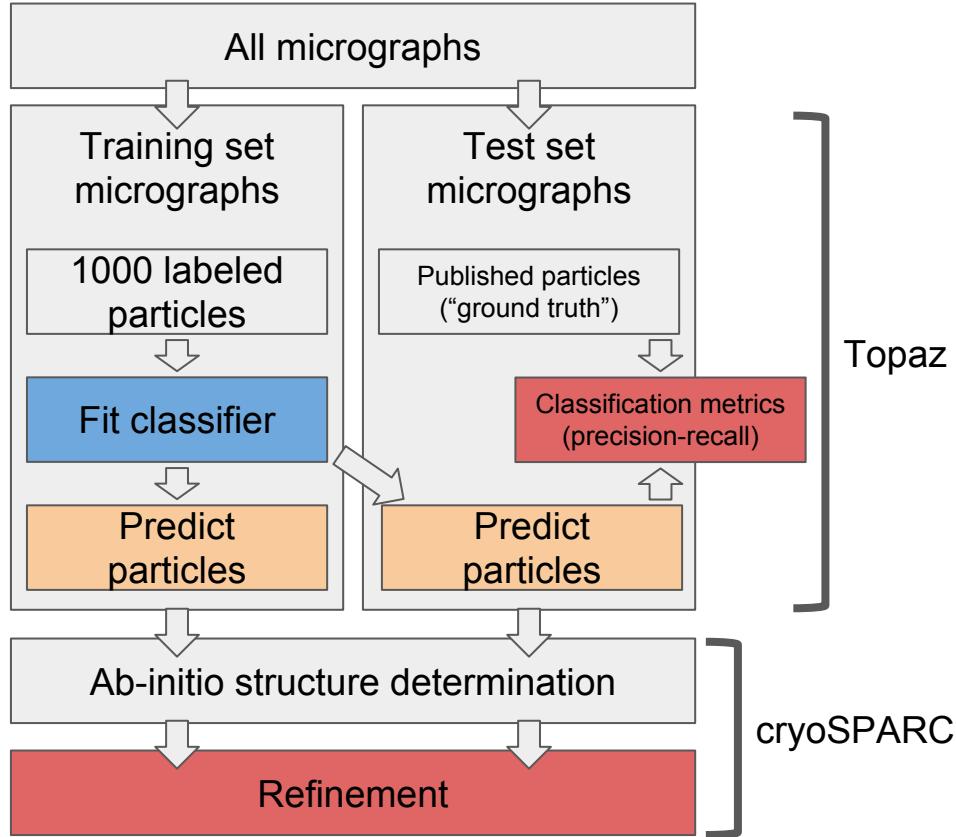
Score micrograph regions and extract predicted particle coordinates



Extract coordinates from region predictions using non-maximum suppression

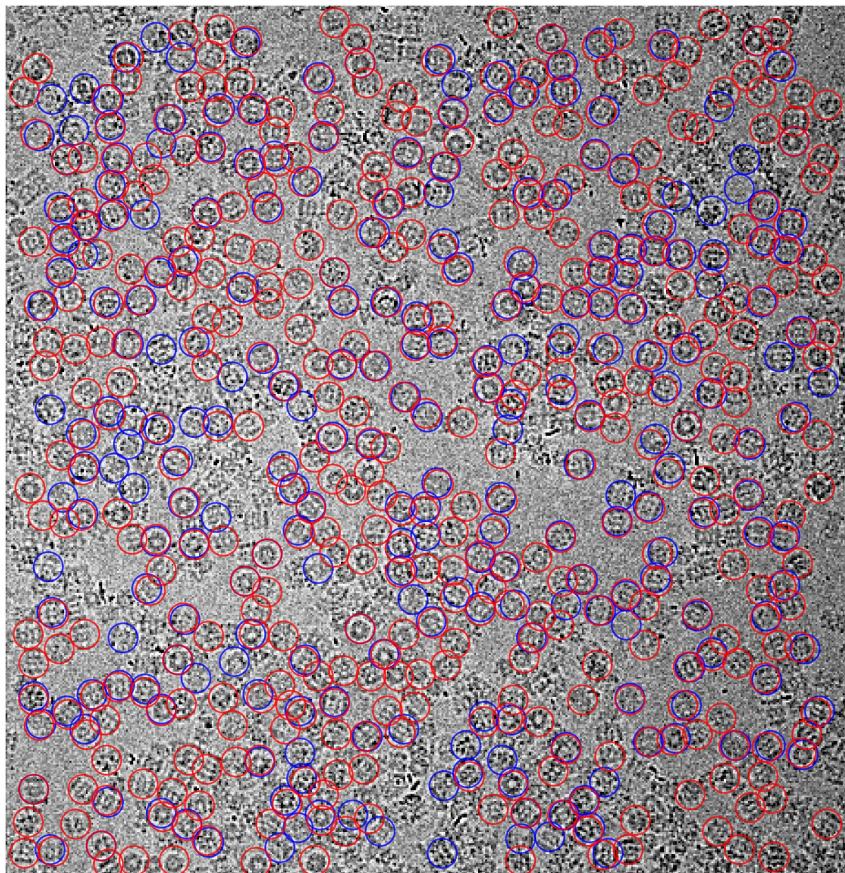
# Structure determination with predicted particles

- 2 new datasets: EMPIAR-10025 (T20S proteasome) and EMPIAR-10028 (80S ribosome)
- 20% of micrographs held out for testing particle detection
- 1000 labeled training particles
- Predicted particles selected at decreasing score thresholds
- Ab-initio structure determination and refinement performed with each particle set with cryoSPARC
- *No post-processing of predicted particles (no 2D/3D classification)*

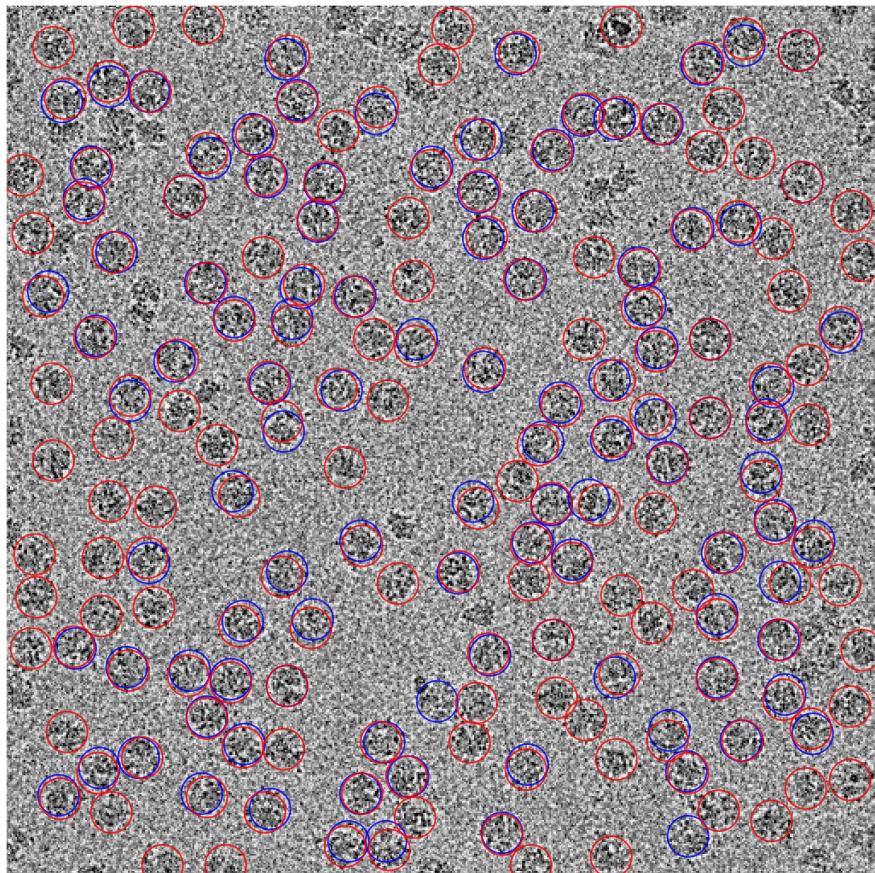


Example test set micrographs show extra predicted particles are true particles

EMPIAR-10025



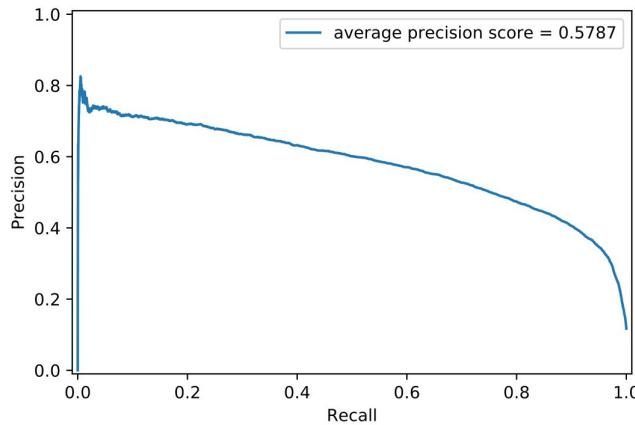
EMPIAR-10028



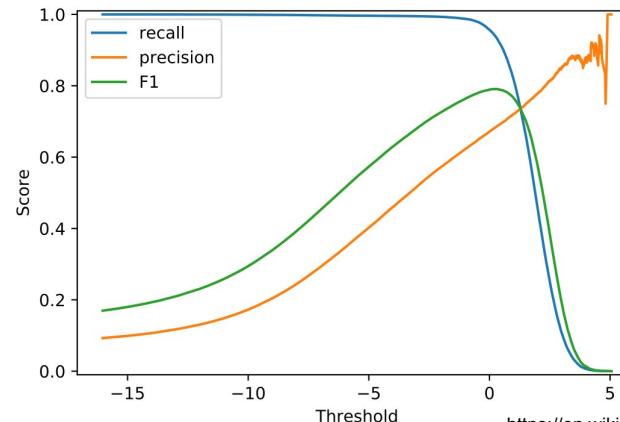
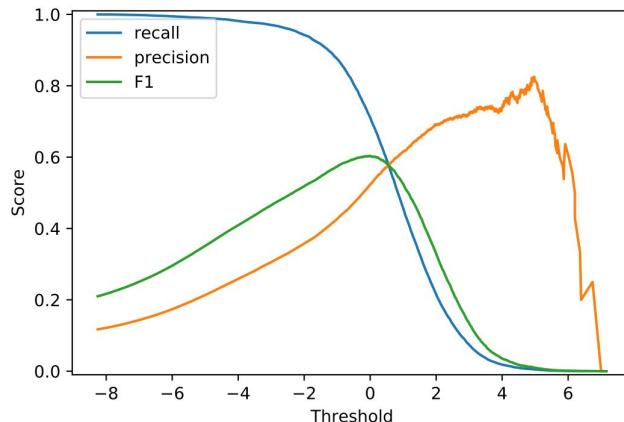
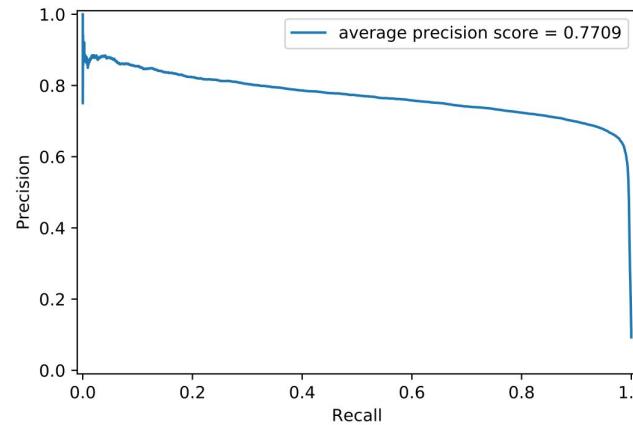
Red: predicted particles, Blue: published particles

# Models detect test set particles with good average precision scores

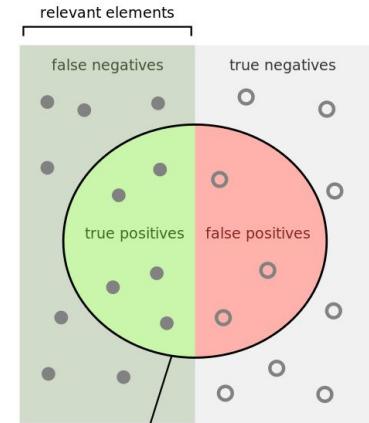
EMPIAR-10025



EMPIAR-10028



$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$
$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

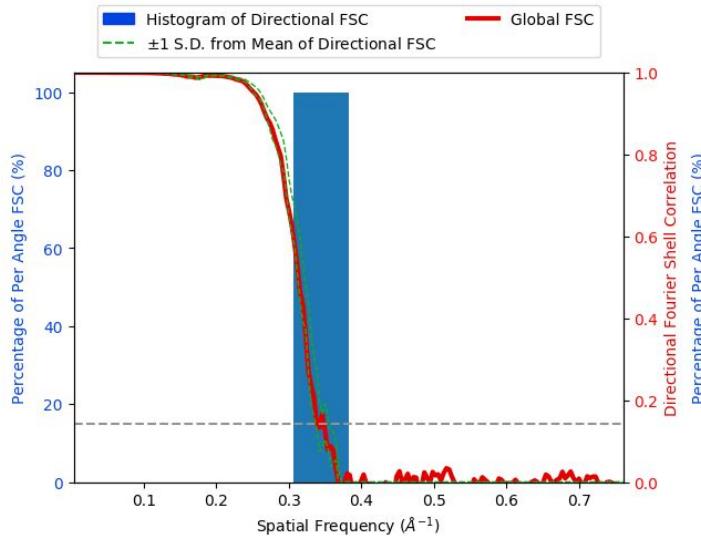


$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$
$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

# 2.8 Å reconstruction of EMPIAR-10025 (without dose weighting)

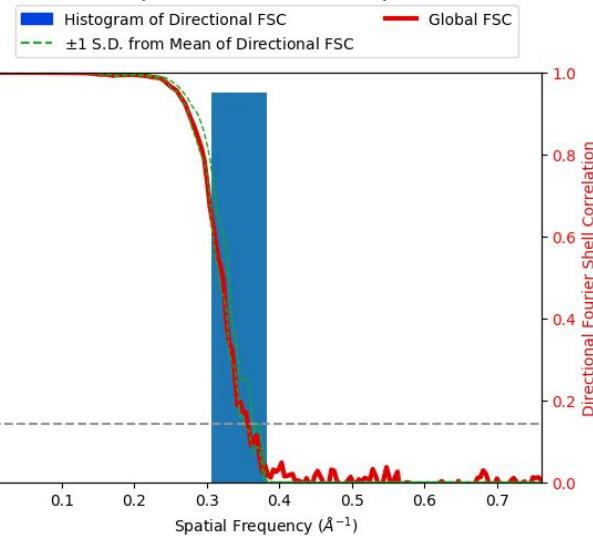
Predicted 3D structure

Published particle set

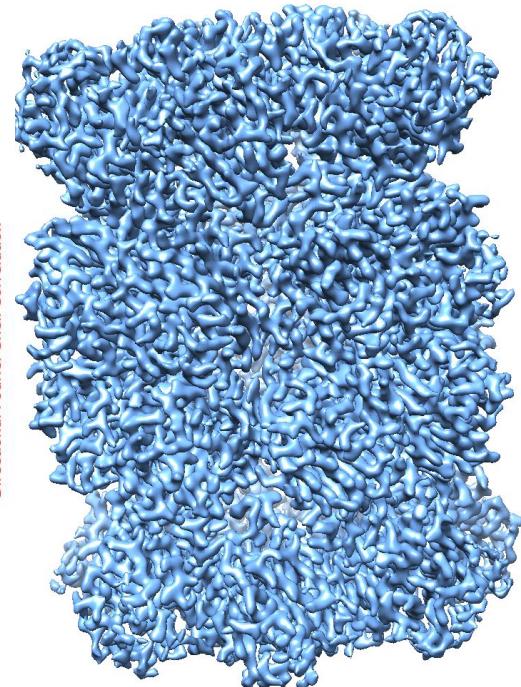


Sphericity = 0.988  
Global resolution = 2.99 Å  
Number of particles = 49954

Predicted particle set  
(best threshold)

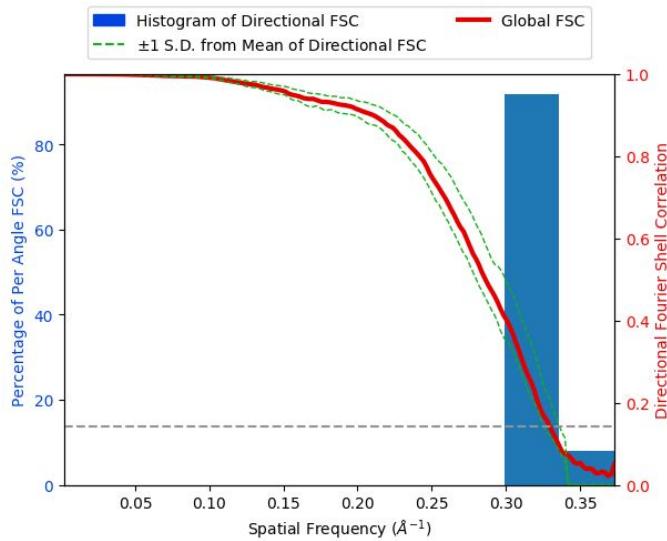


Sphericity = 0.982  
Global resolution = 2.83 Å  
Number of particles = 160658

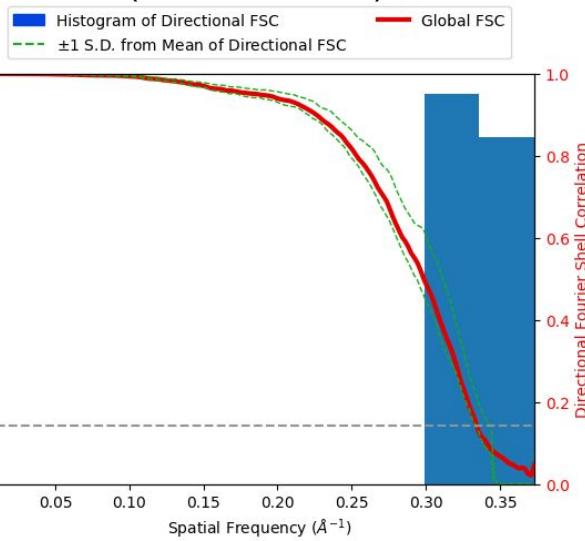


# 3.0 Å reconstruction of EMPIAR-10028

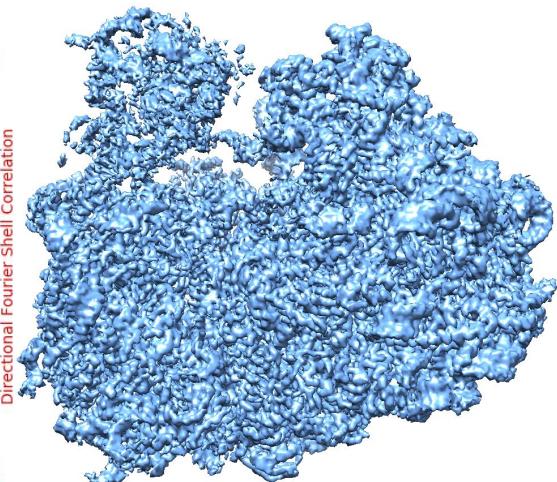
Published particle set



Predicted particle set  
(best threshold)



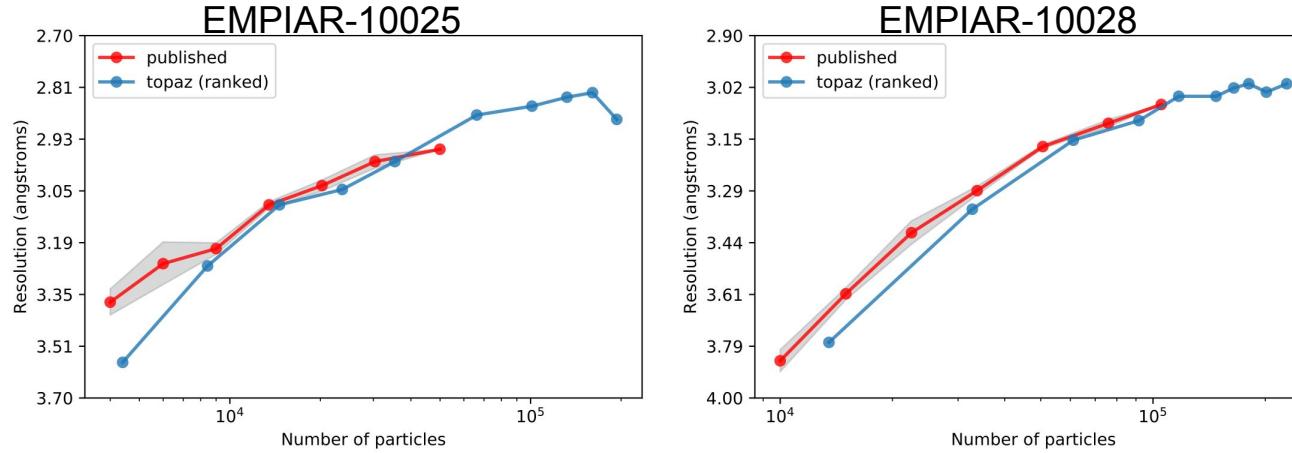
Predicted 3D structure



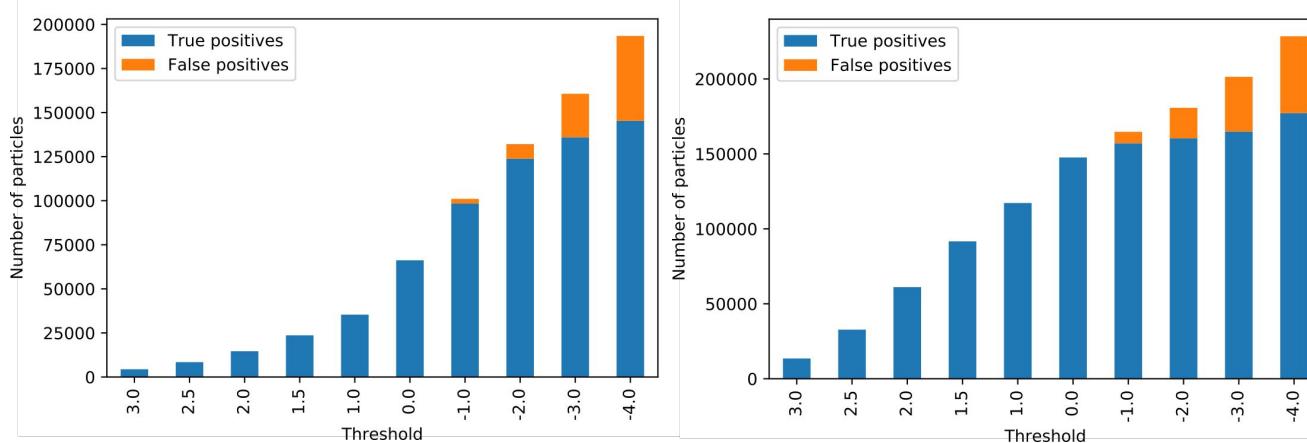
\*\* EM-map challenge best resolution 3.10 Å

# Predicted particles are well-ranked

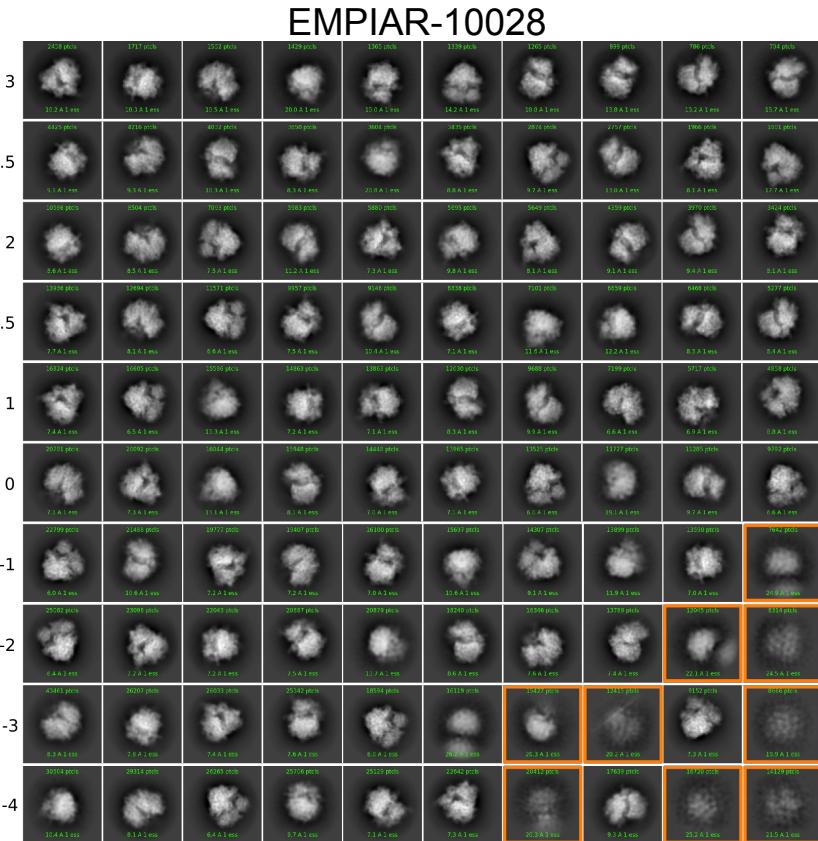
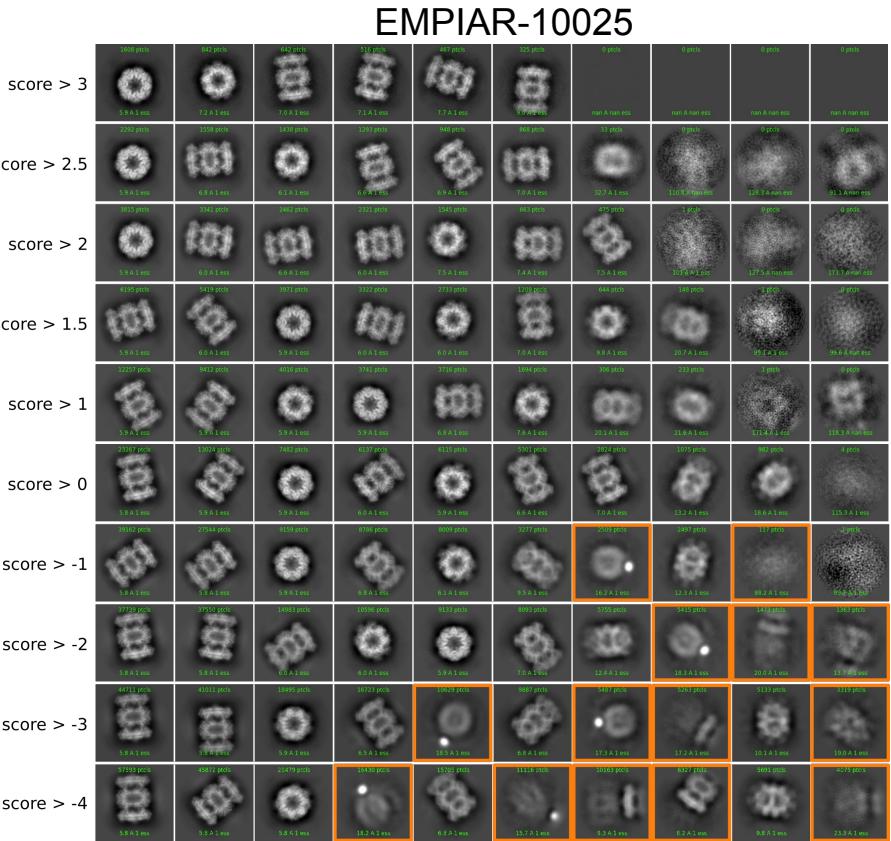
3d reconstruction  
resolution



2d class  
quantification



# 2d class averages with decreasing particle threshold



Lower threshold -> Increasing number of particles

# Summary

1. We proposed the **GE-binomial loss function** and showed that neural network classifiers trained to minimize this loss on positive and unlabeled micrograph regions **outperform classifiers trained with other positive-unlabeled learning objective functions** on 2 challenging cryoEM datasets (Shapiro-lab, EMPIAR-10096)
2. We showed that creating a joint training scheme in which the classifier is trained together with a decoder to form a **hybrid classifier+autoencoder can further improve performance** when few labeled data points are available
3. We developed an object detection pipeline for picking particles using classifiers trained from positive and unlabeled examples.
4. We showed that particles predicted by Topaz (with only 1000 labeled training examples and no postprocessing) give **state-of-the-art reconstructions** on 2 additional datasets (EMPIAR-10025, EMPIAR-10028)

Topaz - our implementation of this particle picking pipeline - is available at <https://github.com/tbepler/topaz>

Manuscript in preparation - preprint can be found at <https://arxiv.org/abs/1803.08207>

# Acknowledgements

Massachusetts Institute of Technology

- **Bonnie Berger**
- **Andrew Morin**
- Hoon Cho
- Tommi Jaakkola

New York Structural Biology Center

- **Alex Noble**
- Bridget Carragher
- Clint Potter

Columbia University

- **Julia Brasch**
- **Larry Shapiro**