

Positive-unlabeled convolutional neural networks for particle picking in cryo-electron micrographs

Tristan Bepler^{1,2}, Andrew Morin⁶, Alex J. Noble³, Julia Brasch⁴, Lawrence Shapiro^{4,5},
and Bonnie Berger^{2,6,*}

¹ Computational and Systems Biology, MIT, Cambridge, MA, USA

² Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, MA, USA

³ National Resource for Automated Molecular Microscopy, Simons Electron Microscopy Center, New York Structural Biology Center, NY, USA

⁴ Department of Biochemistry and Molecular Biophysics, Columbia University, NY, USA

⁵ Mortimer B. Zuckerman Mind Brain Behavior Institute, NY, USA

⁶ Department of Mathematics, MIT, Cambridge, MA, USA

* Correspondence: bab@csail.mit.edu

Abstract

Cryo-electron microscopy (cryoEM) is fast becoming the preferred method for protein structure determination. Particle picking is a significant bottleneck in the solving of protein structures from single particle cryoEM. Hand labeling sufficient numbers of particles can take months of effort and current computationally based approaches are often ineffective. Here, we frame particle picking as a positive-unlabeled classification problem in which we seek to learn a convolutional neural network (CNN) to classify micrograph regions as particle or background from a small number of labeled positive examples and many unlabeled examples. However, model fitting with very few labeled data points is a challenging machine learning problem. To address this, we develop a novel objective function, GE-binomial, for learning model parameters in this context. This objective uses a newly-formulated generalized expectation (GE) criteria to learn effectively from unlabeled data when using minibatched stochastic gradient descent optimizers. On a high-quality publicly available cryoEM data set, we show that CNNs trained with this objective classify particles accurately with very few positive training examples. Using 1000 randomly sampled particles (out of 100k total) as references, EMAN2's byRef method achieves 33% precision at 90% recall. With the same 1000 labeled training particles, we improve this result by roughly 40% to 46% precision at 90% recall. Remarkably, we achieve 41% precision with 1/10th the number of labeled particles and still reach 34% precision with only 1/100th the number of labeled particles at 90% recall. At all numbers of labeled particles, we improve substantially over EMAN2's area under the precision-recall curve (AUPR). Our relative performance increase is even greater on a difficult unpublished dataset supplied by the Shapiro lab. Furthermore, we show that incorporating an autoencoder improves generalization when very few labeled data points are available. We also compare our GE-binomial method with other positive-unlabeled learning methods never before applied to particle picking. We expect our particle picking tool, Topaz, based on CNNs trained with GE-binomial, to be an essential component of single particle cryoEM analysis and our GE-binomial objective function to be widely applicable to positive-unlabeled classification problems.

Introduction

Structure determination with cryoEM involves reconstructing a 3D molecule from 2D projections, often requiring tens to hundreds of thousands of experimental projections, or particles, to be obtained. Locating these particles in cryoEM micrographs (Figure 1), referred to as particle picking, is a major bottleneck in the current protein structure determination pipeline. This pipeline generally consists of sample and EM grid preparation, imaging, particle picking, and eventually structure determination¹. Labelling a sufficient number of particles to determine a high resolution structure can require months of effort. A number of methods exist for automating this process²⁻⁴. Template-based approaches, for example, search micrographs for patterns matching particle ‘exemplars’. However, the usefulness of these tools can be limited due to high false positive rates. Recent high profile controversies illustrate that initializing template-based methods from homologous structures can easily identify sufficient false positives to give erroneous 3D structures⁵⁻⁷. Therefore, it is still necessary for particle picking methods to use particle examples from the data set of interest and hand-picking remains the gold standard.

Recently, a number of particle picking methods have been proposed using convolutional neural network classifiers trained to classify micrograph regions as particle or background using positive and negative examples^{8,9}. However, these approaches typically require researchers to label a large number of regions for training-- a non-trivial task. Moreover, the diverse characteristics of negative data make it difficult to manually label a representative set of negative examples (Figure 1). To overcome the problems inherent in template-based methods and fully supervised neural network methods, we newly frame this problem as a semi-supervised classification problem in which we seek to learn a model to classify micrograph regions as particle or background based on a small number of labeled regions and the large number of unlabeled regions in the micrographs collected. Our approach enables us to effectively pick particles using a very small number of positive examples, a challenging machine learning problem. Furthermore, it is common practice for cryoEM researchers to label only positive examples, as is regularly done to generate picking templates, allowing our method to fit easily into current analysis pipelines (e.g. Appion¹⁰).

Positive-unlabeled learning is a problem in which we seek to learn a classifier that can separate positive from negative data points where the training data only contains labeled positive and unlabeled examples. This problem has been studied extensively in the context of binary classifiers for information retrieval¹¹⁻¹⁴. Two of the most successful approaches are the generalized expectation criterion proposed by Mann and McCallum¹⁵, typically using the KL-divergence (GE-KL), and recent work from du Plessis et al. framing the problem as cost sensitive learning between positive and unlabeled data¹³ to which Kiryo et al. have proposed an improvement suited for use with neural network classifiers (NNPU)¹⁴.

Here, we present two major advances: 1) Topaz, a pipeline for particle picking using classifiers trained with positive and unlabeled data (Figure 2); and 2) GE-binomial, a general objective function for learning classifier parameters from positive and unlabeled data with a novel generalized expectation criteria. We choose CNNs as the specific classifier for our particle picking pipeline, because they are state-of-the-art models for image classification and object detection^{16,17}. We show that with a very small number of data points, our method gives superior particle predictions when compared with the widely-used EMAN2 byRef particle picker², which also uses only positive examples. Furthermore, we show that augmenting the classifier with a decoder to form a hybrid classifier+autoencoder improves model performance when only a small number of labeled positives are available. Lastly, we compare against other positive-unlabeled learning methods (including GE with KL-divergence and NNPU), none of which have been applied to particle picking, and demonstrate empirically that GE-binomial improves over other positive-unlabeled objective functions when training CNNs to classify micrograph regions.

On a high-quality publicly available cryoEM dataset consisting of 100k particles¹⁸, we show that CNNs trained with GE-binomial on 1000 randomly sampled particles achieve 46% precision at 90% recall as compared to EMAN2 byRef which achieves 33% precision at 90% recall using the same 1000 particles as references. With 100 labeled positives (10x less labeled data), our method achieves 41% precision and with only 10 labeled positives (100x less labeled particles than EMAN2), our method still reaches 34% precision at 90% recall. In addition, our models have substantially higher AUPRs at all numbers of particles. On a second, unpublished, dataset that confounds current particle pickers (e.g. Gautomatch, DoG picker, FindEM template picker, and RELION DoG picker, as included in Appion¹⁰), we achieve even greater comparative success. This dataset contains a particle that is particularly difficult for other pickers to locate due to its thin, elongated shape (see appendix figure S1) -- demonstrating that our method is particularly well suited for picking challenging proteins.

We make our source code freely available for academic use at (<https://github.com/tbepler/topaz>). It runs efficiently on a single GPU machine. Topaz is currently being integrated into Appion¹⁰ and, we hope, other cryoEM software suites in the near future.

Methods

Datasets

EMPIAR-10096¹⁹ and an unpublished dataset provided by the Shapiro lab (Shapiro-lab dataset) were used for model evaluation. For both datasets, frame aligned micrographs were used. The particle annotations for EMPIAR-10096 were taken from Tan et al.¹⁸. For the Shapiro-lab dataset, particles were hand labeled by a cryoEM expert. One hundred (out of 447) micrographs were randomly held out from EMPIAR-10096 as the test set. For the Shapiro-lab dataset, 20 (out of 87) micrographs were randomly held out. Table 1 reports the number of images and number of particles in the train and test splits for each dataset. Example particles and unlabeled regions from EMPIAR-10096 are included in Figure 2. Examples from the Shapiro-lab dataset can be found in appendix figure S1.

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

** Table 1: Summary of datasets with the number of micrographs and particles in the training and test splits. **

Existing positive-unlabeled classification methods

Here we summarize existing positive unlabeled classification methods to which we later compare but note that none of these have previously been used for particle picking. As notation, we use π to denote the positive class prior, P , the training data labeled as positive, U the unlabeled training data, g , the classifier with parameters we wish to learn, and L , the loss function we wish to minimize.

A naive approach to this problem is to assume that π is small and treat unlabeled data as negative for purposes of learning the model parameters. We refer to this as the positive-negative (PN) learning method which minimizes the standard classification objective with unlabeled data considered as negatives

$$\pi E_P[L(g(x), 1)] + (1 - \pi)E_U[L(g(x), 0)]$$

Recently, du Plessis et al. frame the positive-unlabeled learning problem as an instance of cost sensitive classification and propose an unbiased estimator of the true misclassification risk¹³. Kiryo et al.

have since proposed a modification in which they bound the objective function to be non-negative, making it more appropriate for use with neural networks¹⁴. This objective function is

$$\pi E_P[L(g(x), 1)] + \max\{0, E_U[L(g(x), 0)] - \pi E_P[L(g(x), 0)]\}$$

We refer to this objective as the ‘NNPU’ method throughout this work.

An alternative approach to learning from positive and unlabeled data, proposed by Mann et al., is to impose a penalty on the expectation of the model over the unlabeled data while still minimizing the standard classification loss on the labeled data. They refer to this penalty as a generalized expectation (GE) criteria¹⁵. A common choice for this penalty is the KL-divergence between the expectation of the model predictions and π , weighted by a slack term λ . This gives the following objective

$$E_P[L(g(x), 1)] + \lambda KL(\pi \| E_U[g(x)])$$

We call this objective ‘GE-KL.’

Because we wish to use a minibatched stochastic gradient descent optimizer, the expectation over the unlabeled data needs to be estimated at each step from a minibatch of data. These noisy KL-divergence estimates are biased and tend to drive the parameters towards $g(x) = \pi$ for all data points in the unlabeled set.

GE-binomial: a novel positive-unlabeled classification method

We propose an alternative to the KL-divergence GE criteria that attempts to reduce the bias caused by estimating the penalty from minibatches of unlabeled data. We recognized that, given a sample of N unlabeled data points in a minibatch, the number of these that are positive follows a Binomial distribution where the probability of drawing a positive data point is given by π , the positive class prior. Furthermore, the model predictions, $g(x)$, for these data points define a distribution over the number of positive data points in the minibatch according to the current classifier. This distribution can be approximated by the normal distribution with $\mu = \sum_{i=0}^N g(x_i)$ and $\sigma^2 = \sum_{i=0}^N g(x_i)(1 - g(x_i))$. Let q_k denote the probability that k of N unlabeled data points are positive given by the normal approximation from the $g(x)$ ’s and let p_k denote the probability that k unlabeled data points are positive given by the binomial distribution prior, then we propose to use the cross entropy of q from p as the GE criteria. This gives the GE-binomial objective function:

$$E_P[L(g(x), 1)] + \sum_{k=0}^N q_k \log p_k$$

In general, for all of these positive-unlabeled learning methods, π is unknown and should be determined by cross validation.

Particle extraction with non-maximum suppression

To extract predicted particle coordinates from the per window model predictions, we adopt a greedy version of the widely used non-maximum suppression algorithm for object detection²⁰. Given a fixed radius r , we iteratively choose coordinates starting from the maximum window score such that those coordinates are at least r distance apart. Here, we use a radius of 10 pixels for EMPIAR-10096 and a radius of 15 pixels for the Shapiro-lab dataset.

In order to calculate the precision-recall curve of a set of predicted particle coordinates with scores, we first match these coordinates to the ground truth coordinates by finding the one-to-one assignment of predicted coordinates to ground truth coordinates that minimizes the Euclidean distance between them up to a maximum distance threshold. This threshold is chosen to be the approximate semi-major axis of each particle to ensure that predictions are not assigned to distant particles. This corresponds to a maximum distance threshold of 7 pixels for EMPIAR-10096 and 15 pixels for the Shapiro-lab dataset. The precision-recall curve is then calculated from these assignments.

Results

Overview of GE-binomial objective and particle picking pipeline

We break particle picking from positive examples down into two main stages (Figure 2). The first stage is region classification using a CNN. The second is particle coordinate extraction from per region predictions using non-maximum suppression.

In order to learn a classifier from positive and unlabeled regions, we propose a novel GE criteria based on the observation that, given positive class prior π and a minibatch containing N unlabeled data points, the number of positive data points in a sample of unlabeled data is distributed binomially with N trials and π chance of success. We can then approximate the distribution over the number of positive data points given by the classifier by a normal distribution with parameters as a function of the predicted probability that each data point is a positive. We denote the probability that k data points are positive given this distribution as q_k and the probability that k data points are positive given the binomial prior as p_k . Treating the cross entropy between this distribution and the binomial prior as the GE criteria gives our GE-binomial learning objective (see *GE-binomial: a novel-positive-unlabeled classification method* for more details):

$$E_P[L(g(x), 1)] + \sum_{k=0}^N q_k \log p_k$$

Furthermore, we find that augmenting the CNN classifier with a decoder and the objective with a corresponding reconstruction error term to form a hybrid classifier+autoencoder leads to improved generalization. We use squared loss as the reconstruction error. Given a reconstruction x' and reconstruction error weight c , the full objective with reconstruction loss becomes

$$E_P[L(g(x), 1)] + \sum_{k=0}^N q_k \log p_k + cE[\|x - x'\|_2^2]$$

Because we wish to weight the autoencoder less when more labeled data points are available, we set c to be $10/P$ where P is the number of labeled particles.

We provide a particle picking software package, Topaz, implementing this pipeline for use by the research community (<https://github.com/tbepler/topaz>).

CNNs trained with GE-binomial greatly outperform EMAN2’s byRef particle picker

We compare CNN classifier+autoencoders trained using our GE-binomial method (see supplementary methods for architecture and training details) against EMAN2’s byRef picker (EMAN version 2.2) and show that our method gives more accurate predictions using fewer labeled particles on two cryoEM datasets by comparing the AUPR on heldout micrographs (see *Datasets*). In this and all subsequent analyses, micrographs were downsampled and normalized as described in the supplementary methods. For EMPIAR-10096, we trained our models with a subset of positive particles randomly sampled from the training set, segment and extract particle predictions for test set images using non-maximum suppression with a suppression radius of 10 pixels. We repeated this process 50 times for 10 random positives and 10 times for 100 random positives. EMAN2 byRef was run using 1000 randomly sampled positives from the training set as references with a box size of 32 and particle diameter of 15. Because the EMAN2 byRef picker does not output confidences for each predicted particle, we ran it with thresholds of 0, 2, 4, and 6 and used these to estimate the AUPR. To compare using 1000 labeled particles, we also trained CNNs with GE-binomial 10 times on the same particles used as references for EMAN2. In order to calculate classification metrics, we matched each predicted particle coordinate to the nearest ground truth particle coordinate within 7 pixels such that no two predicted particles are assigned to the same ground truth particle. We chose a 7 pixel cutoff because it corresponds to the approximate particle radius in the EMPIAR-10096 images.

We report mean classification metrics for particle detection with each of these methods in Table 2. With only 10 positive labeled data points, our GE-binomial model achieves a superior area-under the precision-recall curve (AUPR) and equivalent precision at 0.9 recall levels to EMAN2 with 1000 labeled positives. Our models perform even better when trained with 100 and 1000 particles. On the same 1000 particles used for EMAN2, our models achieve roughly a 40% relative increase in precision at 90% recall. Appendix figure S2 depicts predictions for these methods on an example micrograph.

Number of labeled particles	Method	AUPR	Precision at recall			
			0.25	0.5	0.75	0.9
10	Topaz [classifier+autoencoder, GE-binomial]	0.453	0.529	0.489	0.421	0.338
100	Topaz [classifier+autoencoder, GE-binomial]	0.545	0.633	0.577	0.497	0.405
1000	Topaz [classifier, GE-binomial]	0.616	0.701	0.642	0.556	0.465
	EMAN2 byRef	0.397	-	-	-	0.329

** Table 2: Particle prediction on the EMPIAR-10096 test set. We report the AUPR and precision at various recall levels averaged over GE-binomial CNNs with autoencoders trained on either 50 replicates of 10 or 10 replicates of 100 randomly-sampled labeled particles. We also report the average metrics for 10 models trained on a single set of 1000 randomly sampled particles. These results are compared against AUPRs and precisions for EMAN2’s byRef picker using those same 1000 particles as references. At the thresholds used for EMAN2, precisions at 0.25, 0.5, and 0.75 recall could not be calculated. Boldfaced results indicate top performance in each category.**

We repeated this experiment with a second, particularly challenging dataset provided by the Shapiro lab. For this dataset, we trained models on 10 randomly sampled sets of 10 positive examples from the training set, 10 randomly sampled sets of 100 positive examples from the training set, and 10 times on the full training set of 1167 particles. Predicted particles were extracted using a suppression radius of 15 pixels and matched to ground truth particles with a maximum distance of 15 pixels. EMAN2 byRef picker was run with a box size of 48 and particle size of 31 for this dataset and at the same thresholds as for EMPIAR-10096. Table 3 shows results for particle extraction on the test set averaged over each set of samples compared with predictions from EMAN2 byRef using all 1167 particles in the training set; predictions on an example micrograph can be found in appendix figure S3.

We find that for this dataset our CNN+autoencoders trained using the GE-binomial objective dramatically outperform EMAN2’s byRef particle picker for all numbers of labeled examples. We note that the test set is less completely labeled than for EMPIAR-10096, causing AUPR and precisions to be lower. Manual examination of the predictions shows that many particles predicted by Topaz appear to be unlabeled positives in the test set. Furthermore, this dataset is significantly more difficult to pick due to the thin, elongated shape of the particles (see appendix figure S1) -- demonstrating that our model is particularly well suited for picking challenging proteins.

Number of labeled particles	Method	AUPR	Precision at recall			
			0.25	0.5	0.75	0.9
10	Topaz [classifier+autoencoder, GE-binomial]	0.142	0.203	0.119	0.066	0.044
100	Topaz [classifier+autoencoder, GE-binomial]	0.293	0.464	0.227	0.105	0.062
1167	Topaz [classifier, GE-binomial]	0.415	0.660	0.368	0.159	0.082
	EMAN2 byRef	0.011	0.076 [†]	-	-	-

**Table 3: Particle predictions on the Shapiro-lab test set. We report the AUPR and precision at various recall levels averaged over hybrid classifier+autoencoder models trained with the GE-binomial objective on either 10 replicates of 10 or 10 replicates of 100 particles randomly sampled from the training set. We also report the averages for 10 models trained on the full 1167 particle training set. These results are compared against AUPRs and precisions for EMAN2’s byRef picker given all 1167 particles in the training set as references. [†]EMAN2 does not reach recall 0.25 at any threshold for this dataset. We report the precision at the highest recall achieved, 0.18. **

Including an autoencoder improves performance with few labels

To assess the importance of augmenting the classifier with an autoencoder, we trained models using the GE-binomial objective with and without decoder components. These models were trained on the same sets of 10, 100, and 1000 training particles from EMPIAR-10096 or 10, 100, and all 1167 training particles from the Shapiro-lab dataset while using the same number of replicates as described in the previous section. The mean particle detection metrics on the test sets for these models are reported in Table 4. We observe that incorporating an autoencoder into the model improves the ability of these models to generalize when the number of labeled data points is small. However, for both datasets, classifiers without autoencoders achieve better performance around 1000 labeled particles. Although we report particle detection results here, the region classification results, found in appendix Table 1, show the same trend.

Number of labeled particles	Model	AUPR	Precision at recall			
			0.25	0.5	0.75	0.9
EMPIAR-10096						
10	classifier	0.335	0.476	0.372	0.200	0.091
	classifier+autoencoder	0.452	0.529	0.489	0.420	0.337
100	classifier	0.488	0.609	0.520	0.395	0.248
	classifier+autoencoder	0.545	0.633	0.577	0.497	0.405
1000	classifier	0.616	0.701	0.642	0.556	0.465
	classifier+autoencoder	0.598	0.685	0.626	0.542	0.454
Shapiro-lab						
10	classifier	0.105	0.131	0.075	0.048	0.038
	classifier+autoencoder	0.142	0.203	0.119	0.066	0.044
100	classifier	0.249	0.357	0.162	0.080	0.051
	classifier+autoencoder	0.293	0.464	0.227	0.105	0.062
1167	classifier	0.415	0.660	0.368	0.159	0.082
	classifier+autoencoder	0.378	0.599	0.307	0.140	0.079

** Table 4: Comparison of models trained using GE-binomial with and without autoencoder components on various numbers of labeled positive data points for each dataset. AUPR and precision at various recall levels for each test set averaged over the replicates are reported.**

GE-binomial outperforms other positive-unlabeled learning methods

We next examine CNN classifiers trained for particle detection using our proposed GE-binomial method and three other approaches: PN, GE-KL, and NNPU (see *Existing positive-unlabeled classification methods*). We trained CNN classifiers on 10 and 100 randomly sampled positive data points from the EMPIAR-10096 and Shapiro-lab training. As before, particle extraction results are reported on the test

sets averaged over 50 samples of 10 and 10 samples of 100 positive training data points for EMPIAR-10096 and 10 sampled of 10 and 10 samples of 100 positive training data points for Shapiro-lab in Table 5. Appendix Table 1 reports the region classification results for these experiments.

These results indicate that models trained using our GE-binomial objective function produce better test set predictions in nearly every case. For EMPIAR-10096 with 10 and 100 labeled positives, the difference in the mean AUPR between GE-binomial and GE-KL trained models is statistically significant according to a paired t-test (p-value 1.9×10^{-14} for 10 and p-value 0.0042 for 100 labeled particles). For the Shapiro-lab dataset, we do not find a significant difference in the mean AUPR between GE-binomial and GE-KL methods, but GE-binomial does significantly improve the mean AUPR over the next best method, NNPU (p-value 0.016 for 10 and p-value 0.00017 for 100 labeled particles). We hypothesize that GE-KL achieves similar results to GE-binomial on the Shapiro-lab dataset, because the Shapiro-lab dataset contains a smaller fraction of positives than EMPIAR-10096 (π estimated to be 0.01 vs. 0.02 for EMPIAR-10096) and minibatch estimates of GE-KL have less bias when this value is small. Furthermore, this suggests that GE-binomial will provide even larger benefit for datasets with a larger fraction of positives.

Number of labeled particles	Method	AUPR	Precision at recall			
			0.25	0.5	0.75	0.9
EMPIAR-10096						
10	PN	0.178 \pm 0.062	0.260 \pm 0.093	0.157 \pm 0.071	0.092 \pm 0.038	0.070 \pm 0.015
	NNPU	0.228 \pm 0.065	0.348 \pm 0.096	0.220 \pm 0.087	0.108 \pm 0.043	0.070 \pm 0.014
	GE-KL	0.222 \pm 0.084	0.378 \pm 0.160	0.179 \pm 0.130	0.083 \pm 0.041	0.065 \pm 0.005
	GE-binomial	0.335 \pm0.072	0.476 \pm0.092	0.372 \pm0.098	0.200 \pm0.082	0.091 \pm0.034
100	PN	0.353 \pm 0.057	0.457 \pm 0.061	0.368 \pm 0.070	0.263 \pm 0.065	0.167 \pm 0.045
	NNPU	0.401 \pm 0.060	0.500 \pm 0.069	0.432 \pm 0.067	0.329 \pm 0.060	0.215 \pm 0.042
	GE-KL	0.464 \pm 0.062	0.590 \pm 0.064	0.495 \pm 0.071	0.362 \pm 0.077	0.221 \pm 0.060
	GE-binomial	0.488 \pm0.057	0.609 \pm0.060	0.520 \pm0.063	0.395 \pm0.067	0.248 \pm0.061
Shapiro-lab						
10	PN	0.069 \pm 0.024	0.079 \pm 0.024	0.059 \pm 0.015	0.044 \pm 0.008	0.036 \pm 0.004
	NNPU	0.075 \pm 0.031	0.089 \pm 0.037	0.061 \pm 0.018	0.044 \pm 0.009	0.037 \pm 0.004
	GE-KL	0.086 \pm 0.036	0.105 \pm 0.044	0.049 \pm 0.012	0.038 \pm 0.004	0.034 \pm 0.002
	GE-binomial	0.105 \pm0.036	0.131 \pm0.049	0.075 \pm0.020	0.048 \pm0.009	0.038 \pm0.004
100	PN	0.153 \pm 0.042	0.197 \pm 0.061	0.121 \pm 0.027	0.070 \pm 0.009	0.049 \pm 0.004
	NNPU	0.162 \pm 0.053	0.223 \pm 0.089	0.121 \pm 0.031	0.073 \pm 0.011	0.050 \pm 0.005
	GE-KL	0.250 \pm0.050	0.369 \pm0.105	0.161 \pm 0.034	0.080 \pm 0.013	0.052 \pm0.006
	GE-binomial	0.249 \pm 0.033	0.357 \pm 0.074	0.162 \pm0.030	0.080 \pm0.007	0.051 \pm 0.004

Table 5: Comparison of particle predictions from CNNs trained with four different positive-unlabeled objective functions with 10 or 100 labeled particles for each dataset. Mean and standard deviations of the AUPR and precision at a variety of recall levels are reported for 50 replicates of 10 positives and 10 replicates of 100 positives from EMPIAR-10096 and 10 replicates of 10 positives and 10 replicates of 100 positives from the Shapiro-lab dataset.

Our positive-unlabeled CNNs are time and memory efficient

Not only does our positive-unlabeled CNN particle picking method dramatically outperform current approaches using less labeled data, it also runs in a practical amount of time. Model training takes \sim 4-6 hours with a single Nvidia K40 or 980ti GPU. Prediction is much faster, taking only minutes to segment and and extract particle predictions from 100 micrographs. In contrast, EMAN2’s byRef method took a hours to extract particles from 100 micrographs given 1000 particle references with a 4-core Intel i7 940

CPU. Additionally, we attempted to use EMAN2’s local search method on the same machine but stopped when the particle extraction had not finished after more than 5 days.

Discussion

CryoEM is revolutionizing structural biology with widespread applications ranging from basic biology to the understanding of disease linked proteins and development of novel therapeutics. Fully realizing the promise of cryoEM and achieving rapid turnaround from imaging to structure determination requires state-of-the-art computational methods. To this end, we have presented GE-binomial, a generalized expectation criteria for learning classifier parameters from positive and unlabeled examples. This objective function explicitly models the number of positive examples in minibatches of unlabeled data as being drawn from a binomial distribution and penalizes the cross entropy between the distribution over number of positives given the model predictions and this binomial prior. We have shown empirically on two cryoEM datasets that convolutional neural networks trained using this objective achieve superior particle detection results to CNNs trained with three other positive-unlabeled objectives. We have demonstrated that augmenting the CNN classifier with an autoencoder can further improve performance with very few labeled data points. We also compared our full particle picking pipeline, Topaz, with EMAN2’s template-based particle picker and found that our method gave improved particle predictions even with fewer labeled particles. We note that the results reported here are underestimates of the true model performance, because the ground truth labeling is not complete, leading to higher apparent false positive rates.

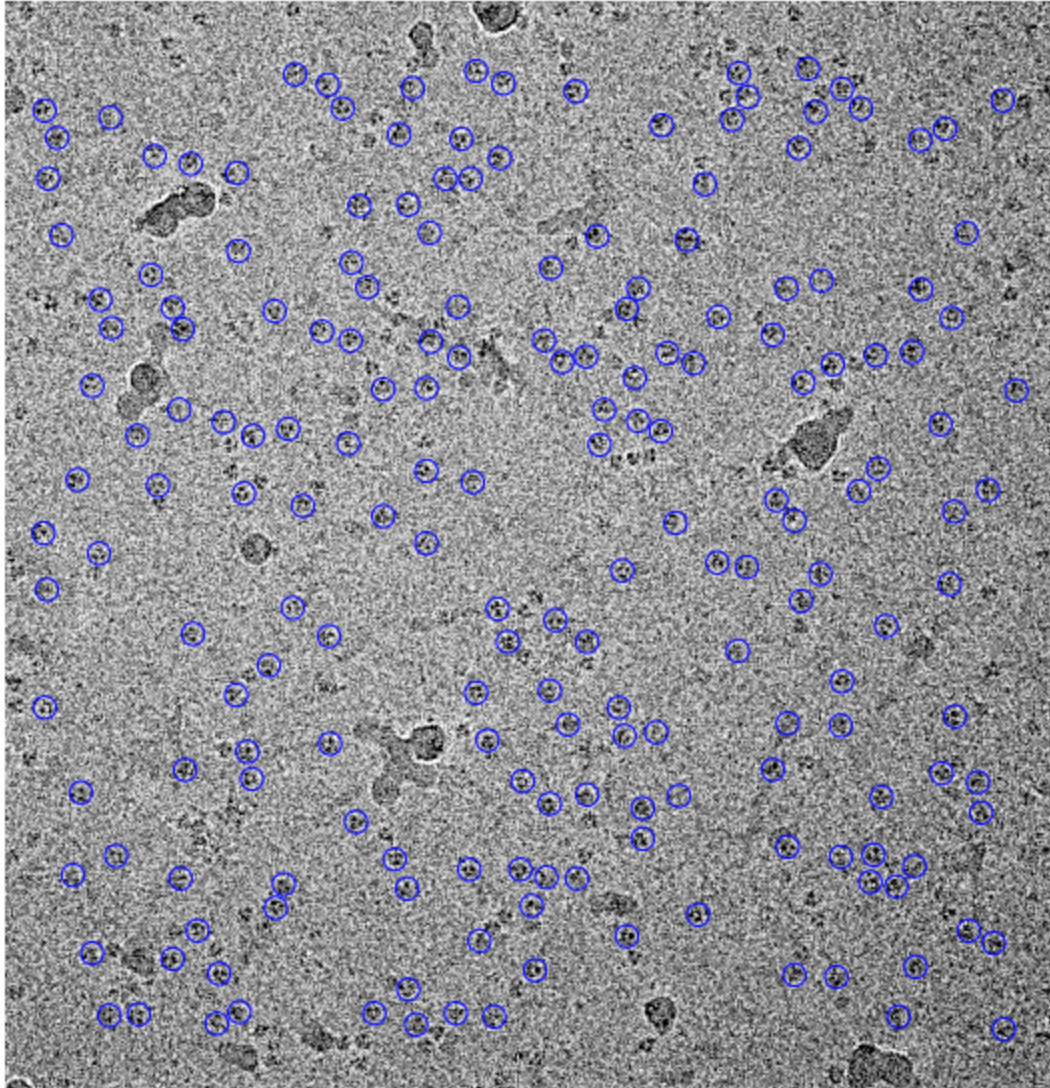
Although we used a simple CNN architecture with reasonable default hyperparameters and showed that it performed well on these datasets, any alternative model architectures that can be trained with gradient descent can use our GE-binomial objective to learn from positive and unlabeled data. In particular, L2 or dropout regularization can improve generalization. These hyperparameter choices should ideally be made by cross validation for new datasets. The only hyperparameter introduced by our objective function, and other positive-unlabeled objectives that we consider, is the unknown positive class prior. Although this parameter should also be chosen by cross validation, we observed that our results were relatively insensitive to its choice. Finally, it is straightforward to extend the GE-binomial objective to include labeled negatives by taking the expectation of the loss over all labeled data in the first term.

Topaz requires researchers to label far fewer particles to achieve high quality predictions. It performs well independently of particle shape, opening automated picking to a wide selection of proteins previously too difficult to locate computationally. In addition, our pipeline is computationally efficient--training in a few hours on a single GPU and producing predictions for hundreds of micrographs in minutes. Furthermore, once a model is trained for a specific particle, it can be reused for new imaging runs of the same particle. In the future, these methods will help to expedite structure determination by cryoEM.

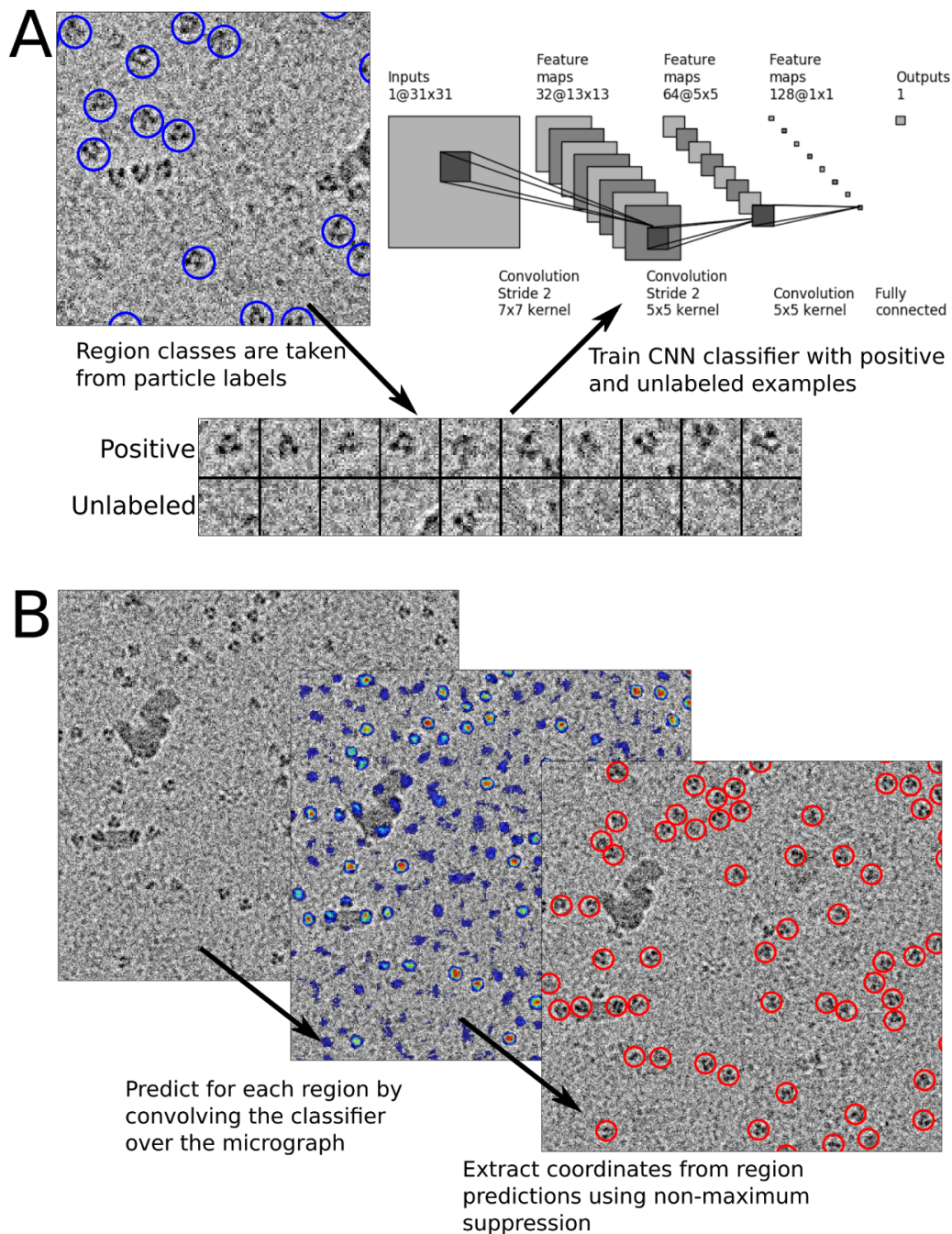
Acknowledgements

We would like to thank Tommi Jaakkola, members of SEMC at NYSBC, and members of the Berger lab for useful discussion. T.B. and A.M. were partially supported by NIH grant R01-GM081871 (to B.B.). J.B. was supported by NIH grant R01-MH1148175 (to L.S.).

Some of this work was performed at the Simons Electron Microscopy Center and National Resource for Automated Molecular Microscopy located at the New York Structural Biology Center, supported by grants from the Simons Foundation (349247), NYSTAR, and the NIH National Institute of General Medical Sciences (GM103310) with additional support from Agouron Institute [Grant Number: F00316] and NIH S10 OD019994-01.



** Figure 1: Example micrograph with particle annotations (blue circles) from EMPIAR-10096. This micrograph demonstrates the diversity of negatives present in this type of data, including noisy background and varieties of ice chunks. **



**** Figure 2: Topaz particle picking pipeline using CNNs trained with positive and unlabeled data. A)** Given a set of labeled particles, a CNN is trained to classify positive and negative regions using particle locations as positive regions and all other regions as unlabeled. Labeled particles from EMPIAR-10096 are indicated by blue circles and a few positive and unlabeled regions are depicted. **B)** Once the CNN classifier is trained, particles are predicted in two steps. First, the classifier is convolved over each micrograph to give predictions for each region of each micrograph. Then, coordinates are extracted from the region predictions using non-maximum suppression. The left image shows a raw micrograph from EMPIAR-10096. The middle image depicts the micrograph with overlaid region predictions [blue = low confidence, red = high confidence]. The right image indicates predicted particles after using non-maximum suppression on the region predictions. **

Literature Cited

1. Cheng, Y., Grigorieff, N., Penczek, P. A. & Walz, T. A primer to single-particle cryo-electron microscopy. *Cell* **161**, 438–449 (2015).
2. Tang, G. *et al.* EMAN2: an extensible image processing suite for electron microscopy. *J. Struct. Biol.* **157**, 38–46 (2007).
3. Voss, N. R., Yoshioka, C. K., Radermacher, M., Potter, C. S. & Carragher, B. DoG Picker and TiltPicker: software tools to facilitate particle selection in single particle electron microscopy. *J. Struct. Biol.* **166**, 205–213 (2009).
4. Scheres, S. H. W. Semi-automated selection of cryo-EM particles in RELION-1.3. *J. Struct. Biol.* **189**, 114–122 (2015).
5. Borrell, B. Rift widens over structure of HIV’s molecular anchor. *Nature News* (2013). doi:10.1038/nature.2013.14071
6. Bartesaghi, A., Merk, A., Borgnia, M. J., Milne, J. L. S. & Subramaniam, S. Prefusion structure of trimeric HIV-1 envelope glycoprotein determined by cryo-electron microscopy. *Nat. Struct. Mol. Biol.* **20**, 1352–1357 (2013).
7. Mao, Y. *et al.* Molecular architecture of the uncleaved HIV-1 envelope glycoprotein trimer. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 12438–12443 (2013).
8. Wang, F. *et al.* DeepPicker: A deep learning approach for fully automated particle picking in cryo-EM. *J. Struct. Biol.* **195**, 325–336 (2016).
9. Zhu, Y., Ouyang, Q. & Mao, Y. A deep convolutional neural network approach to single-particle recognition in cryo-electron microscopy. *BMC Bioinformatics* **18**, 348 (2017).
10. Lander, G. C. *et al.* Appion: an integrated, database-driven pipeline to facilitate EM image processing. *J. Struct. Biol.* **166**, 95–102 (2009).
11. Elkan, C. & Noto, K. Learning Classifiers from Only Positive and Unlabeled Data. in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* 213–220 (ACM, 2008).
12. Blanchard, G., Lee, G. & Scott, C. Semi-Supervised Novelty Detection. *J. Mach. Learn. Res.* **11**, 2973–3009 (2010).
13. du Plessis, M., Niu, G. & Sugiyama, M. Convex Formulation for Learning from Positive and Unlabeled Data. in *International Conference on Machine Learning* 1386–1394 (2015).
14. Kiryo, R., Niu, G., du Plessis, M. C. & Sugiyama, M. Positive-Unlabeled Learning with Non-Negative Risk Estimator. *arXiv [cs.LG]* (2017).
15. Mann, G. S. & McCallum, A. Generalized Expectation Criteria for Semi-Supervised Learning with Weakly Labeled Data. *J. Mach. Learn. Res.* **11**, 955–984 (2010).
16. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
17. Ren, S., He, K., Girshick, R. & Sun, J. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**, 1137–1149 (2017).
18. Tan, Y. Z. *et al.* Addressing preferred specimen orientation in single-particle cryo-EM through tilting. *Nat. Methods* **14**, 793–796 (2017).
19. Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a public archive for raw electron microscopy image data. *Nat. Methods* **13**, 387–388 (2016).
20. Girshick, R., Donahue, J., Darrell, T. & Malik, J. Rich feature hierarchies for accurate object detection and semantic segmentation. in *Proceedings of the IEEE conference on computer vision and pattern recognition* 580–587 (2014).
21. Ioffe, S. & Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. in *International Conference on Machine Learning* 448–456 (2015).
22. Radford, A., Metz, L. & Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434* (2015).
23. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv [cs.LG]* (2014).
24. Paszke, A., Gross, S. & Chintala, S. PyTorch. (2017).

Appendix

Supplementary Methods

Image preprocessing

EMPIAR-10096 images were downsampled 4x and Shapiro lab images were downsampled 8x by cropping in Fourier space.

Images were then normalized using a per-image scaled two component Gaussian mixture model. Given K images, each pixel is modeled as being drawn from a two component Gaussian mixture model, parameterized by π , the mixing parameter, μ_0 , σ_0 , μ_1 , and σ_1 , the means and standard deviations of the Gaussian distributions, with a scalar multiplier for each image, $\alpha_{1...K}$. Let $x_{i,j,k}$ be the value of the pixel at position i,j in image k , it is distributed according to

$$z_{i,j,k} \sim \text{Bernoulli}(\pi)$$
$$x_{i,j,k} | z_{i,j,k} \sim \text{Gaussian}(\alpha_k \mu_{z_{i,j,k}}, (\alpha_k \sigma_{z_{i,j,k}})^2)$$

where $z_{i,j,k}$ is a random variable denoting the component membership of the pixel. The maximum likelihood values of the parameters π , μ_0 , μ_1 , σ_0 , σ_1 and $\alpha_{1...K}$ are found by expectation-maximization for each data set. Then, the pixels are normalized by first dividing by the image scaling factor and then standardizing to the dominant mixture component. Let μ' , σ' be μ_0 , σ_0 if $\pi < 0.5$ and μ_1 , σ_1 otherwise, then the normalized pixel values $x'_{i,j,k}$ are given by

$$x'_{i,j,k} = (\frac{x_{i,j,k}}{\alpha_k} - \mu')/\sigma'$$

Convolutional neural network (CNN) structure

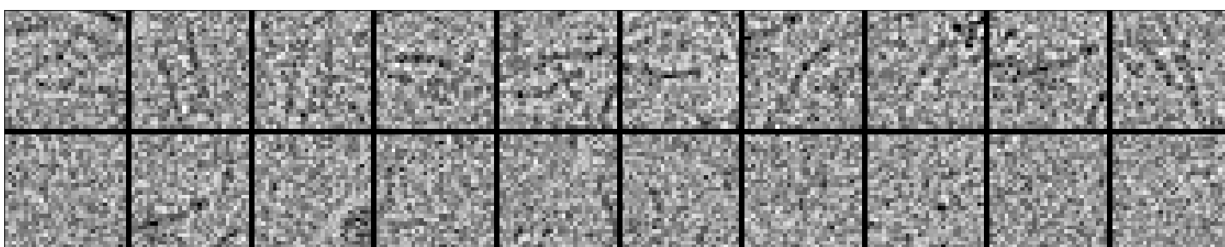
We use a simple three-layer convolutional neural network with striding, batch normalization²¹, and parametric rectified linear units (PReLU) as the classifier in this work. The model is organized as 32 conv7x7 filters with batch normalization and PReLU, stride by 2, 64 conv5x5 filters with batch normalization and PReLU, stride by 2, 128 conv5x5 filters with batch normalization and PReLU, and a final fully connected layer with a single output.

When augmenting with an autoencoder, we use a decoder structure similar to that of DCGAN²². The d -dimensional representation output by the final convolutional layer of the classifier network is projected to a small spatial dimension but large feature dimension representation. This is repeatedly projected into larger spatial dimension and smaller feature dimension representations until the final output is of the original input image size. Specifically, this model is structured as repeated transpose convolutions with batch normalization and leaky ReLU activations. Let z be the representation output by the final convolutional layer of the classifier and X' be the image reconstruction given by the decoder, the decoder structure is $z \rightarrow$ transpose conv4x4 128-d, batch normalization, leaky ReLU \rightarrow transpose conv4x4 64-d, stride 2, batch normalization, leaky ReLU \rightarrow transpose conv4x4 32-d, stride 2, batch normalization, leaky ReLU \rightarrow transpose conv3x3 1-d, stride 2 $\rightarrow X'$.

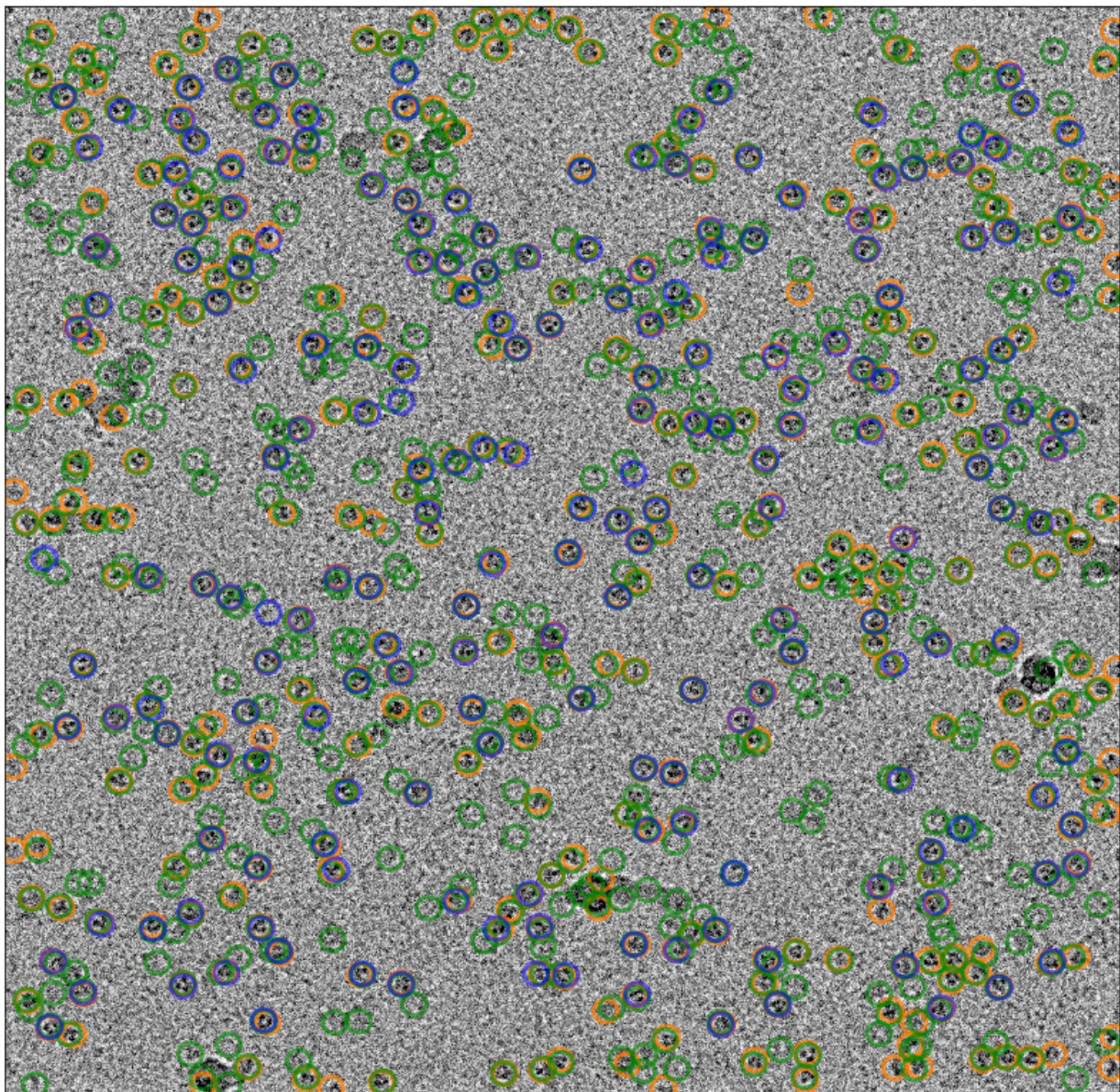
Implementation details

Image regions of size 31x31 were considered positively labeled if the center of the region fell within 4 pixels of a labeled particle coordinate for both datasets. For EMPIAR-10096, the dataset was reasonably completely labeled, so π was set to 0.02, slightly larger than the observed fraction of positive regions in the dataset, 0.016. For the Shapiro-lab dataset, π was chosen to be 0.01 by 5-fold cross validation on the training set from possible values of 0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, and 0.1. When using the GE-KL objective, we set $\lambda = 100$.

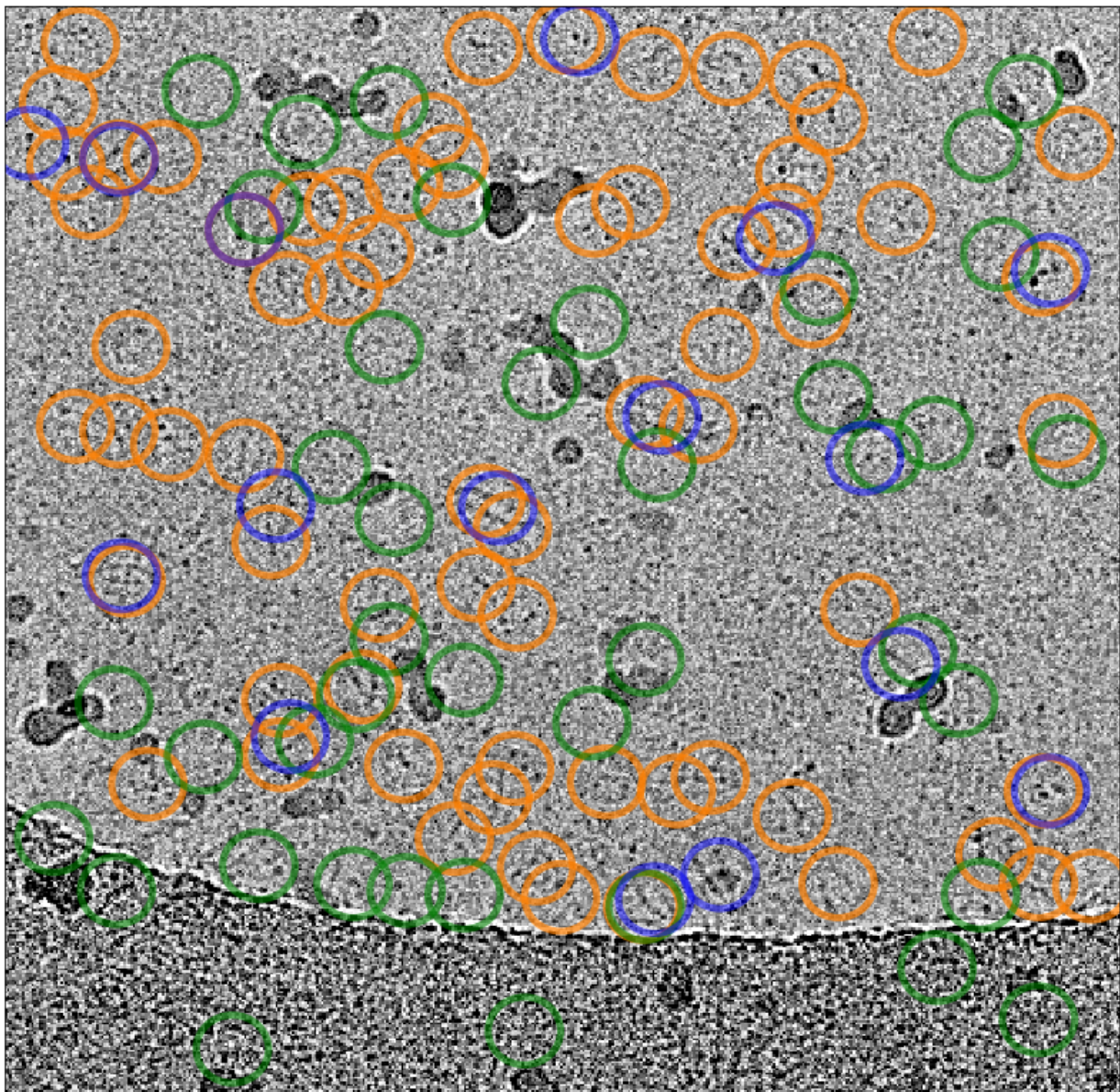
All models were trained for 50,000 iterations using ADAM²³ with a minibatch size of 256 and a learning rate of 0.0002. Each minibatch was composed of 16 positively labeled and 240 unlabeled image patches randomly sampled from the training micrographs with random rotation. We used cross entropy loss for all experiments. Models were implemented with PyTorch²⁴.



** Supplementary figure 1: example particles (top row) and unlabeled regions (bottom row) from the Shapiro-lab dataset. **



** Supplementary figure 2: Particle predictions on an example micrograph from the EMPIAR-10096 test set. **Orange:** CNN+autoencoder trained with GE-binomial on 100 training particles. Predictions are shown for a threshold giving 90% recall. **Green:** EMAN2 predictions from 1000 reference particles with a threshold of 2. **Blue:** particle labels from EMPIAR. **



** Supplementary figure 3: Particle predictions on an example micrograph from the Shapiro-lab test set. **Orange:** CNN+autoencoder trained with GE-binomial on 100 training particles. Predictions are shown for a threshold giving 75% recall. **Green:** EMAN2 predictions from 1000 reference particles with a threshold of 0. **Blue:** expert hand-labeled particles. **

Model	EMPIAR-10096			Shapiro-lab		
	10	100	1000	10	100	1167
classifier						
PN	0.072 \pm 0.029	0.187 \pm 0.038	-	0.012 \pm 0.005	0.036 \pm 0.012	-
NNPU	0.101 \pm 0.035	0.226 \pm 0.033	-	0.014 \pm 0.008	0.039 \pm 0.013	-
GE-KL	0.072 \pm 0.035	0.240 \pm 0.043	-	0.010 \pm 0.005	0.062 \pm0.017	-
GE-binomial	0.155 \pm0.044	0.258 \pm0.040	0.392 \pm0.006	0.020 \pm0.008	0.061 \pm 0.010	0.150 \pm0.008
+autoencoder						
GE-binomial	0.260 \pm 0.016	0.324 \pm 0.017	0.368 \pm 0.013	0.029 \pm 0.011	0.078 \pm 0.018	0.120 \pm 0.006

** Supplementary table 1: AUPR of region classification on test sets for EMPIAR-10096 and Shapiro-lab datasets for CNN classifiers trained with four different objective functions and a CNN+autoencoder trained with GE-binomial. Means and standard deviations are reported for 50 replicates of 10 labeled particles, 10 replicates of 100 labeled particles, and 10 models trained on the same 1000 particles for EMPIAR-10096 and for 10 replicates of 10 labeled particles, 10 replicates of 100 labeled particles, and for 10 models trained on all 1167 training particles for the Shapiro-lab dataset. **