

Deep learning based structural pattern mining in tomograms -- several exploratory studies

Min Xu

Computational Biology Department
School of Computer Science
Carnegie Mellon University

Systematic detection of macromolecular structures in cellular tomograms

Structural pattern mining / in silico purification:
template-free detection of macromolecular structures

Challenges

- Imaging limits
 - Missing data (missing wedge effect)
 - Low signal-to-noise ratio
- High structural content complexity
 - Macromolecule structure highly diverse
 - High molecular crowding level
- Big data
 - Hundreds of tomograms
 - Millions of macromolecules

Deep learning

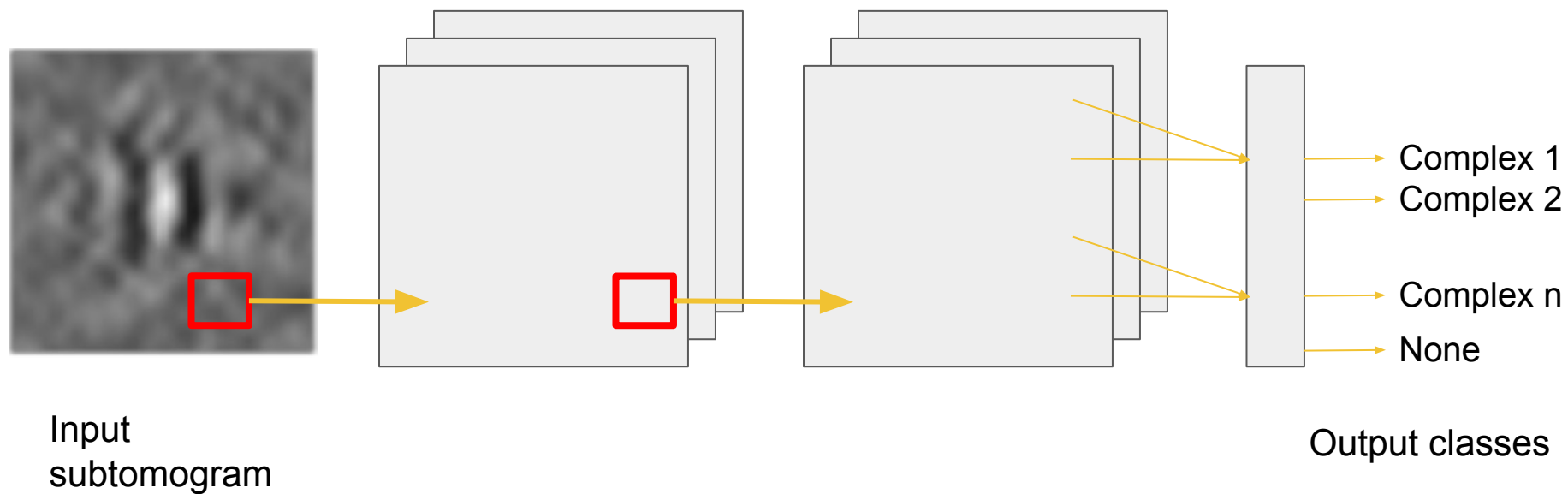
- Convolutional neural network (CNN) performs regression using large amount of parameters
 - Multiple layers + nonlinearity → exponential increase of flexibility for approximation of arbitrary mapping between input and output
 - Convolution: parameter sharing & local connectivity → increases efficiency by taking advantage of composition structure in images
 - Stacked convolutional layers → learning of image feature hierarchy
- Linear scalability respect of training sample number → learn from big data
 - Back-propagation training, easy to implement and parallelize on GPU
- Dropout → improved generalization ability

Exploratory projects

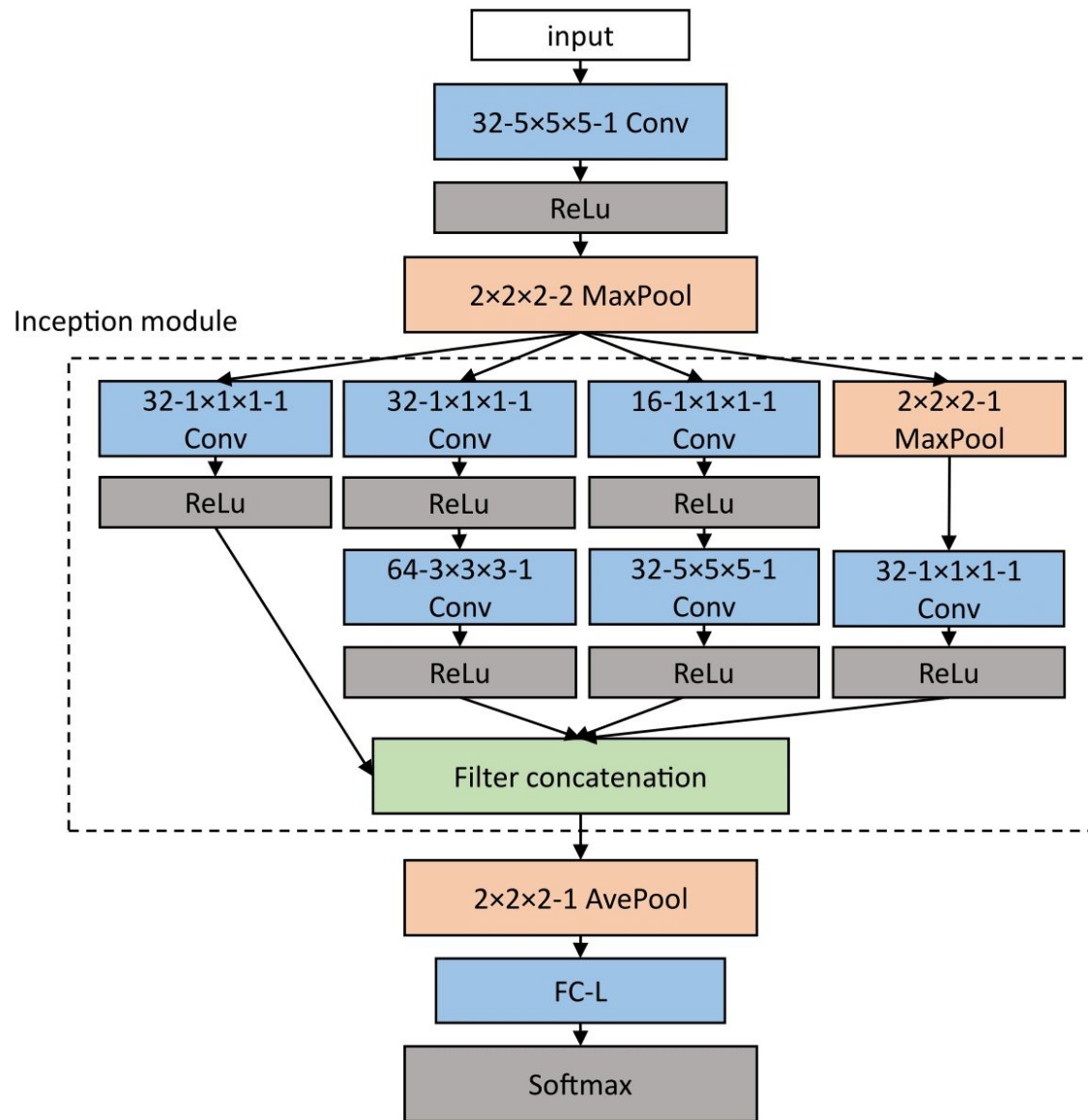
1. Macromolecule structure classification and subdivision
2. Autoencoder based pattern detection
3. Subtomogram segmentation

Supervised subtomogram classification

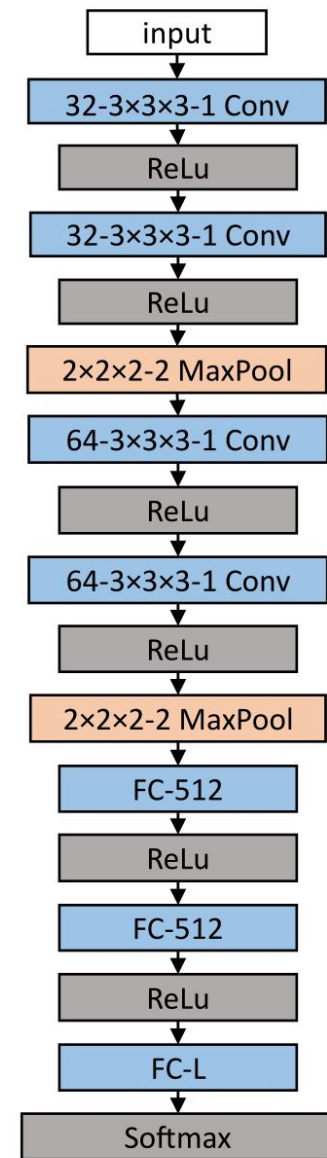
Supervised subtomogram classification



CNN classification models



(a) Inception3D network

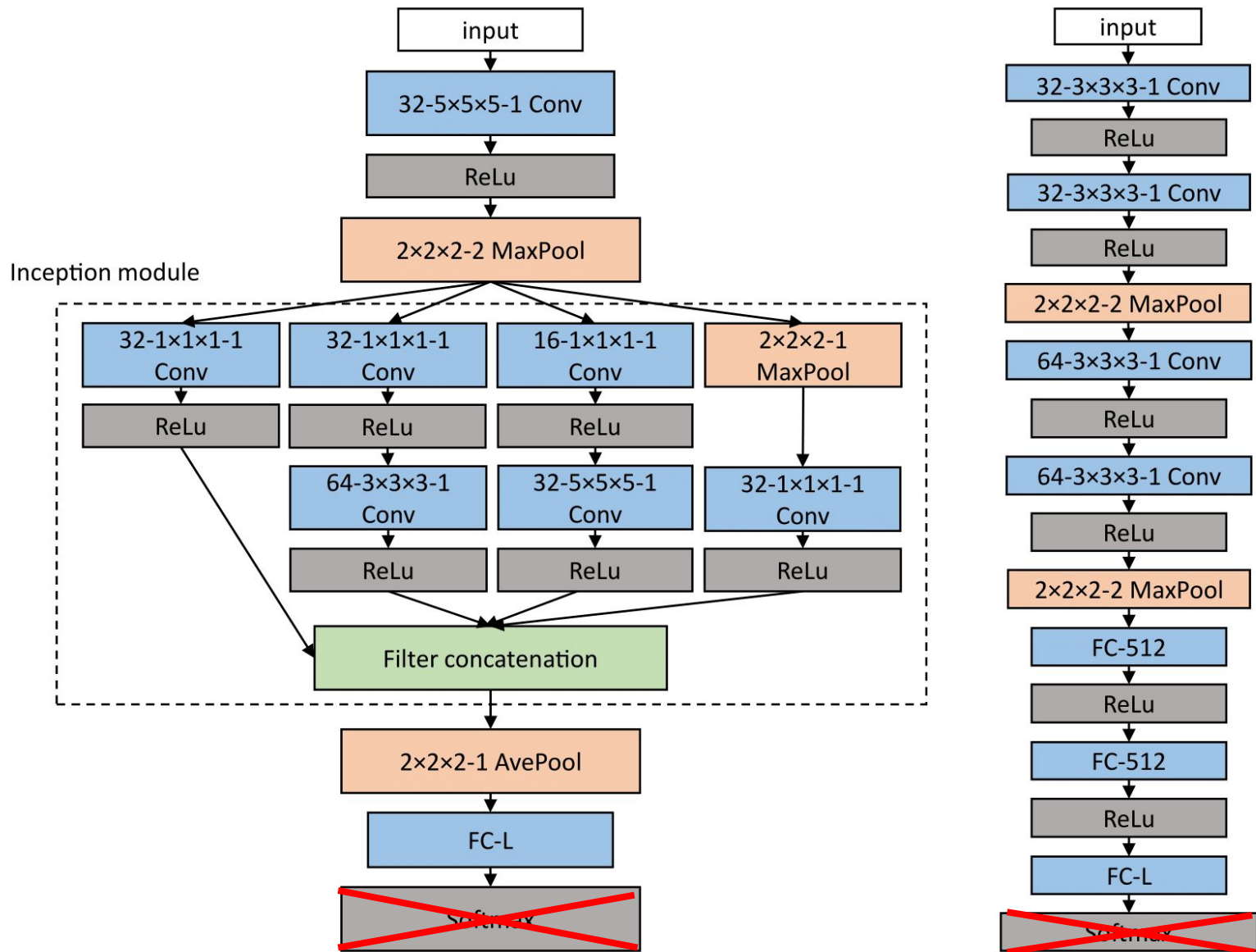


(b) DSRF3D network

Performance

- Classification accuracy significantly better than Rotational Invariant Features + Support Vector Machines
- Once trained, classifying 1M subtomograms take < 2 hours on a single GPU

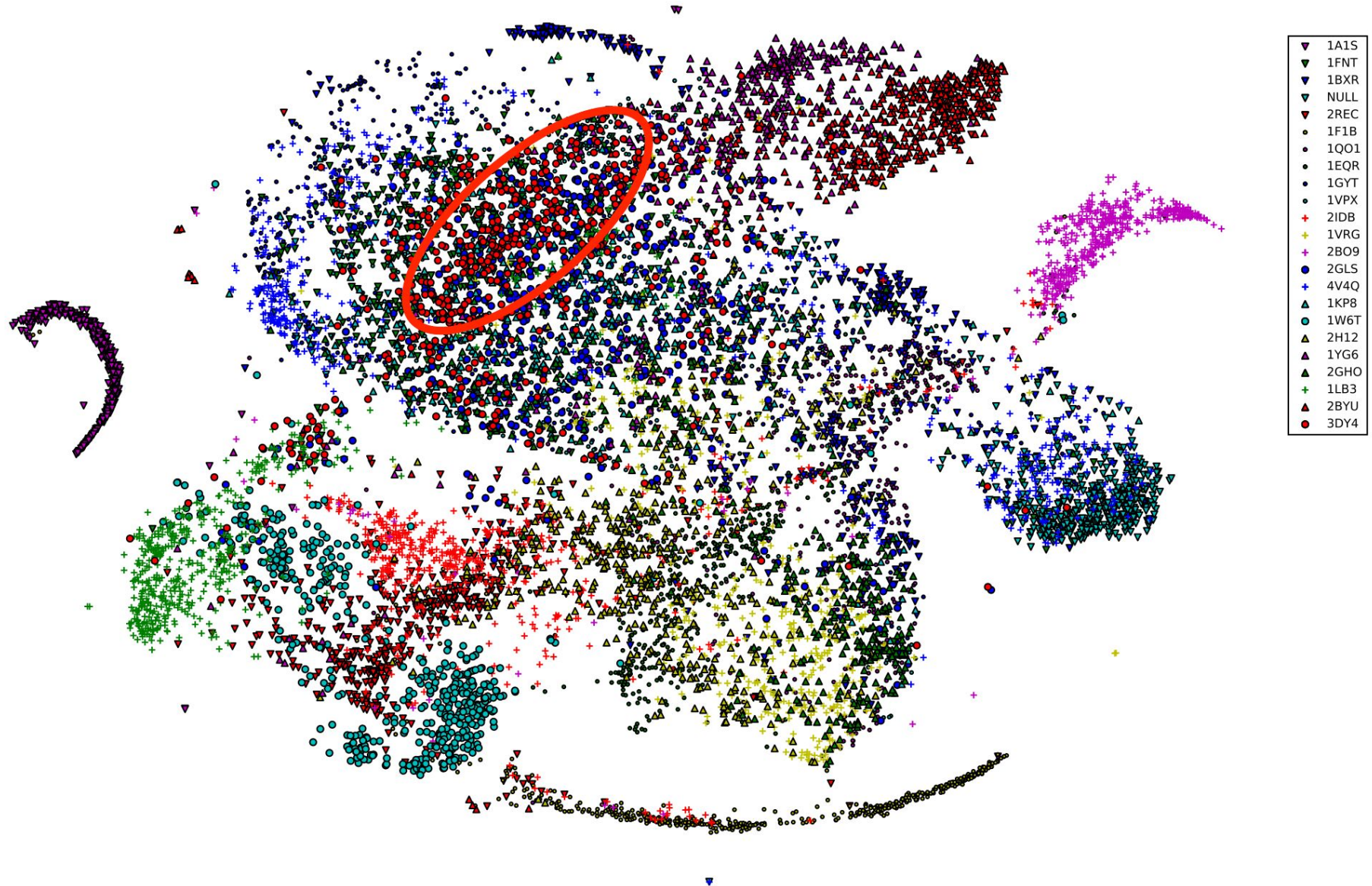
Supervised structural feature extraction



Supervised structural feature extraction

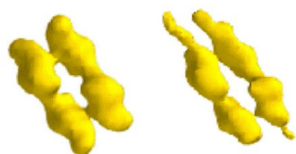
- Continuous representation of the likelihood of the class assignments
- Project the input subtomogram into a low dimensional structural feature space spanned by the training classes
- Invariant to
 - Rigid transformations
 - Missing wedge effects

Detection of new structures

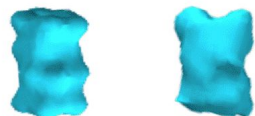


Detection of new structures: leave-one-out test

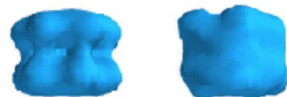
Successfully recovered



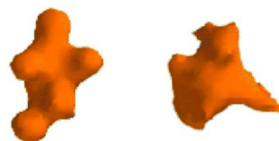
1BXR (5.1)



1VPX (5.9)



2GLS (5.1)



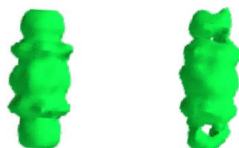
1EQR (6.9)



1VRG (6.4)



2H12 (5.3)



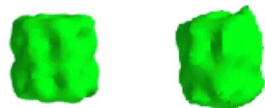
1FNT (4.2)



1W6T (4.3)



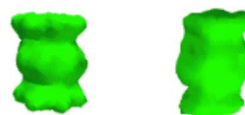
2IDB (4.5)



1KP8 (5.0)



2BO9 (4.7)



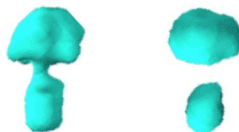
3DY4 (4.2)



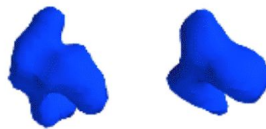
1LB3 (4.8)



2BYU (3.7)



1Q01 (5.9)



2GHO (5.7)

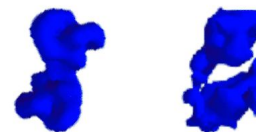
Unsuccessfully recovered



1A1S (38.6)



1F1B (38.6)



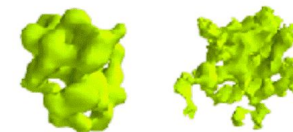
1GYT (7.6)



1YG6 (38.6)

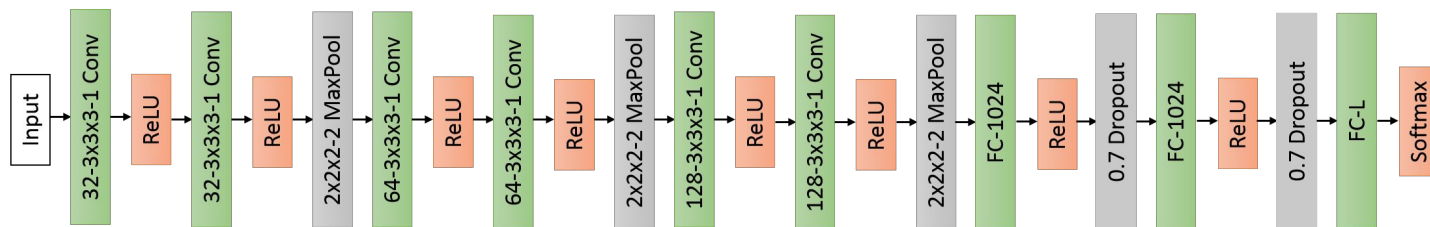


2REC (38.6)

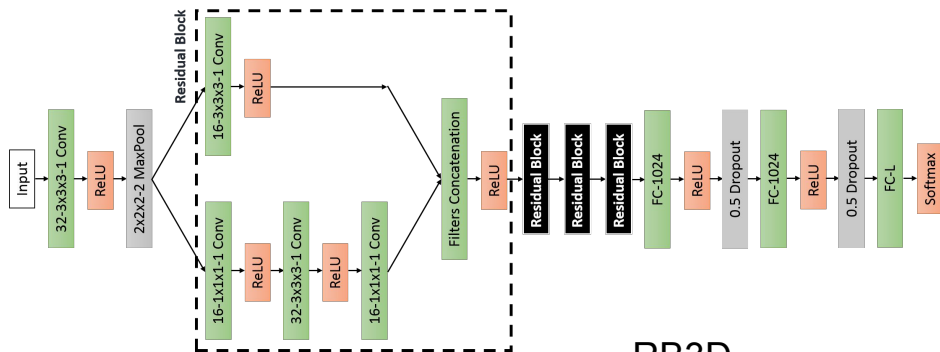


4V4Q (13.5)

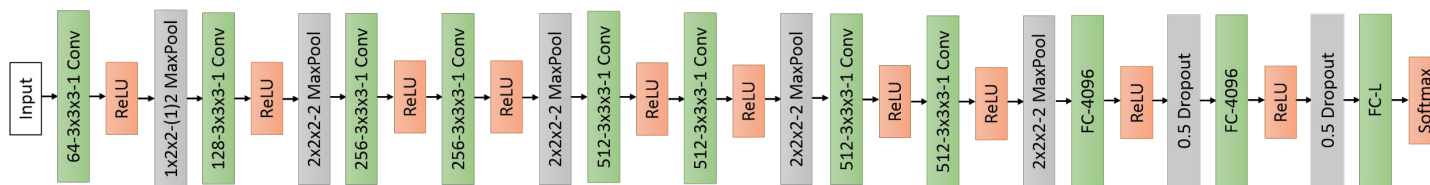
Improvements: deeper models for improved accuracy



DSRF3D-v2



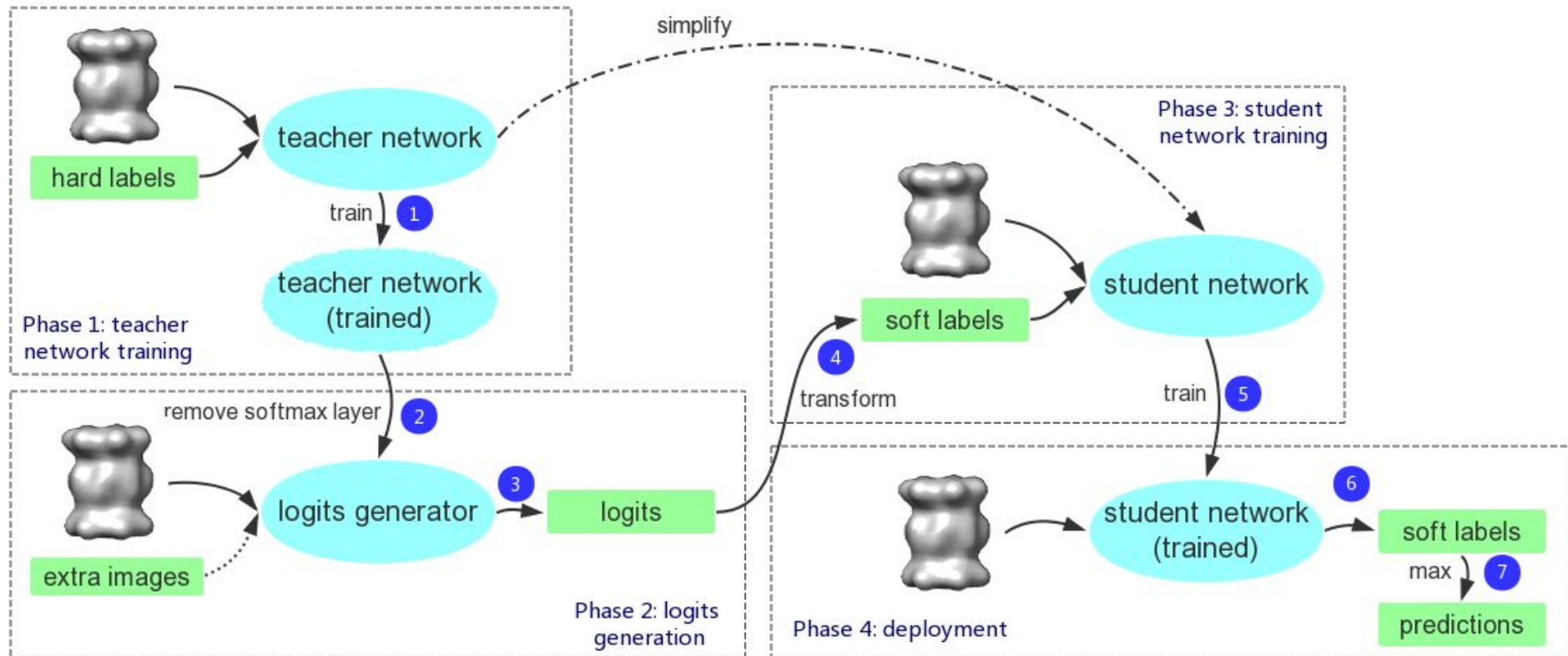
RB3D



CB3D

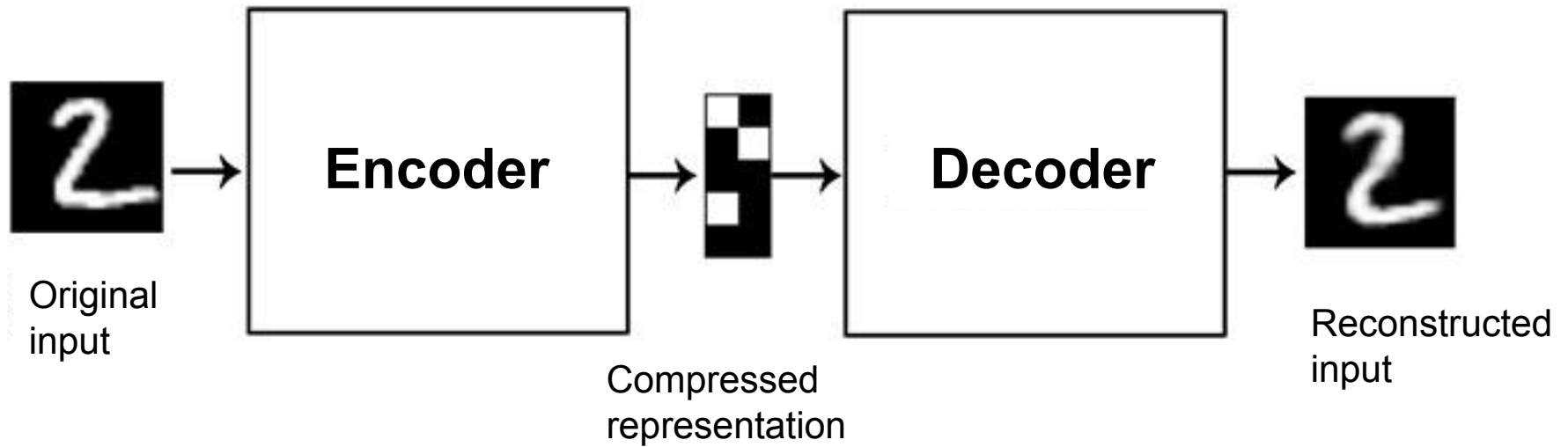
(Best performance)

Improvements: model compression for increased speed

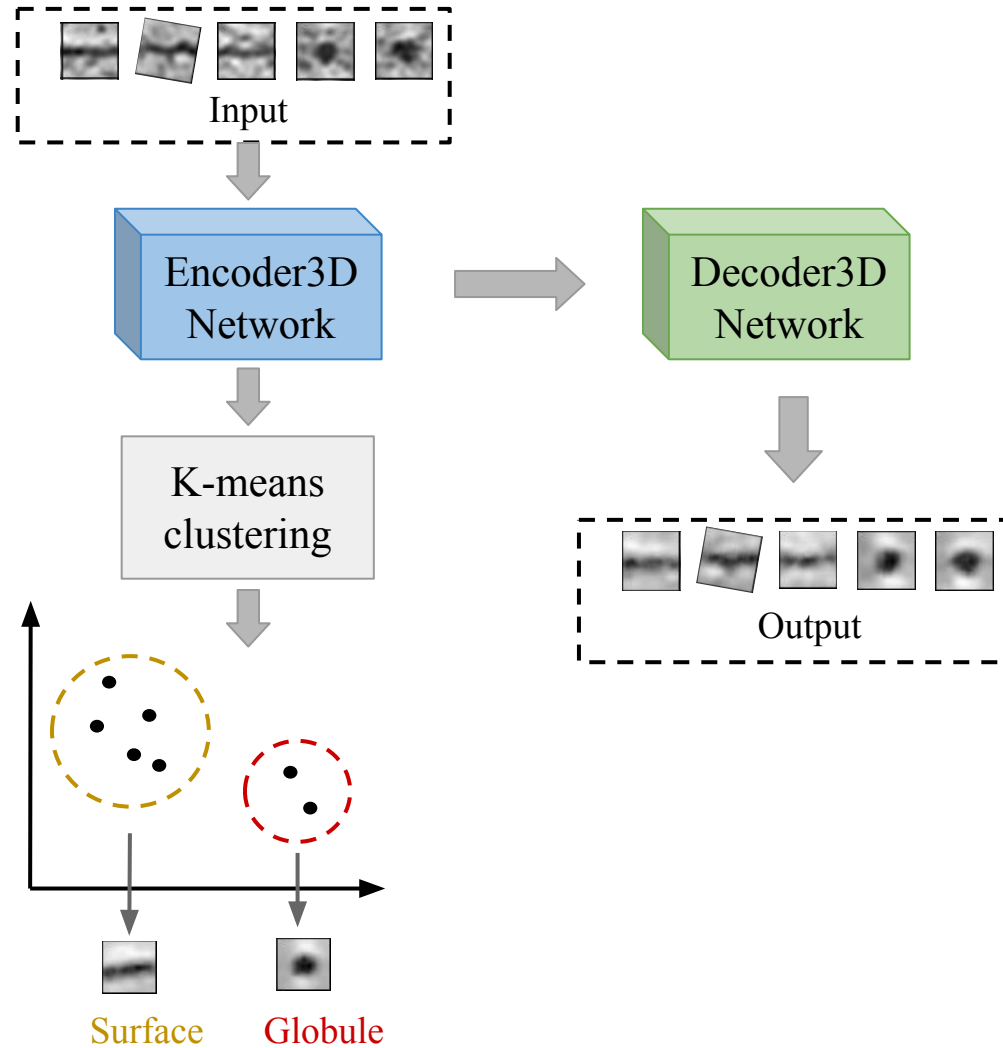


Autoencoder based pattern detection

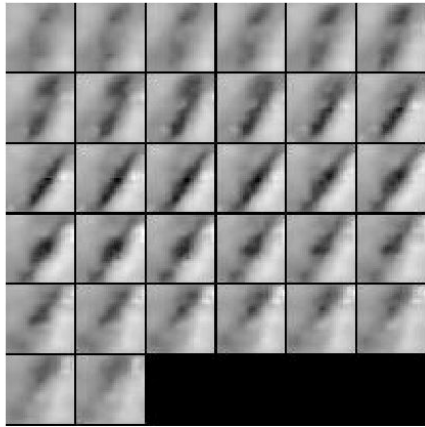
Autoencoder based pattern detection



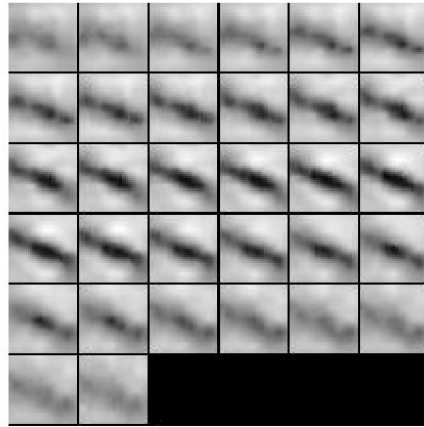
Autoencoder based pattern detection



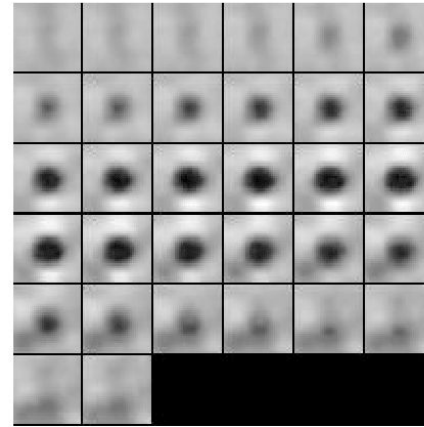
Autoencoder based pattern detection



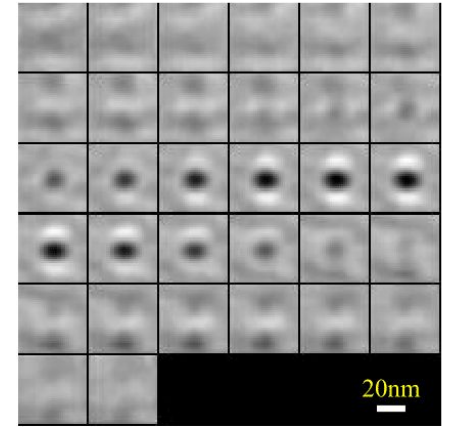
Surface patch



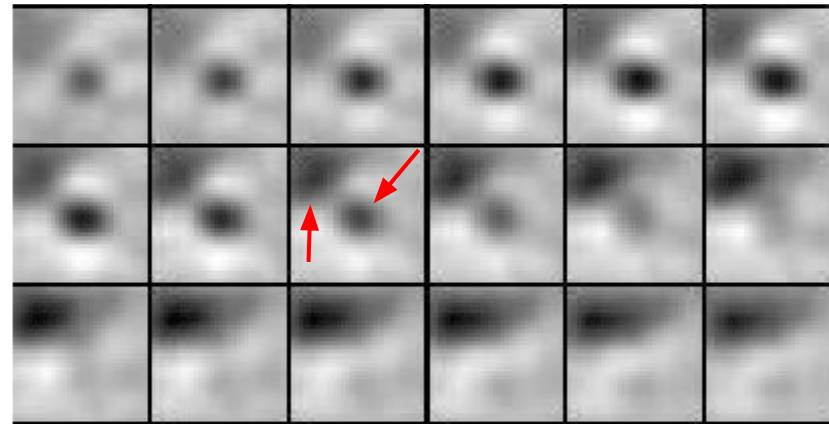
Surface patch



Large globule



Small globule



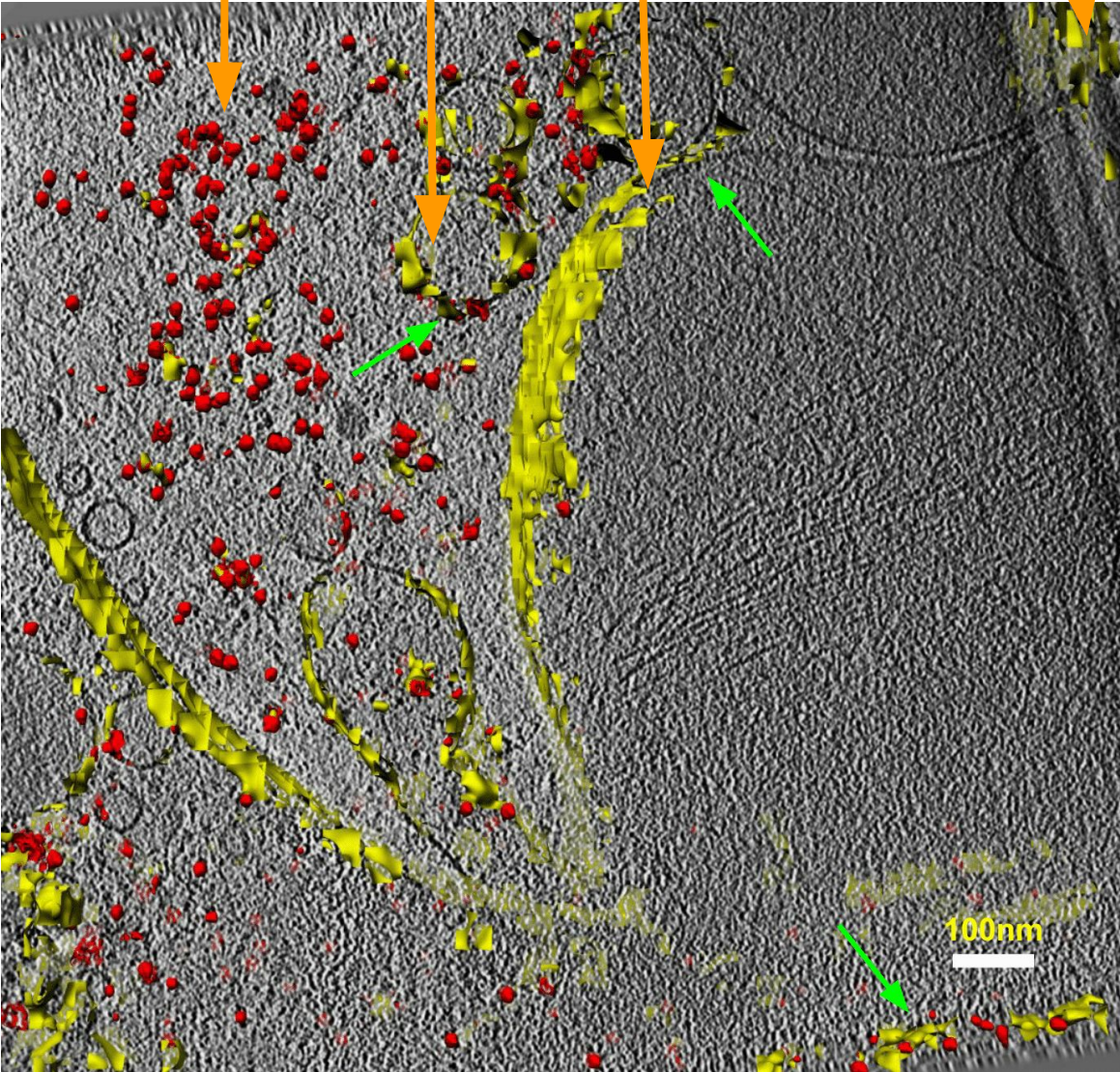
Interaction between cellular components

Embedding of detected patterns

Ribosome like macromolecule

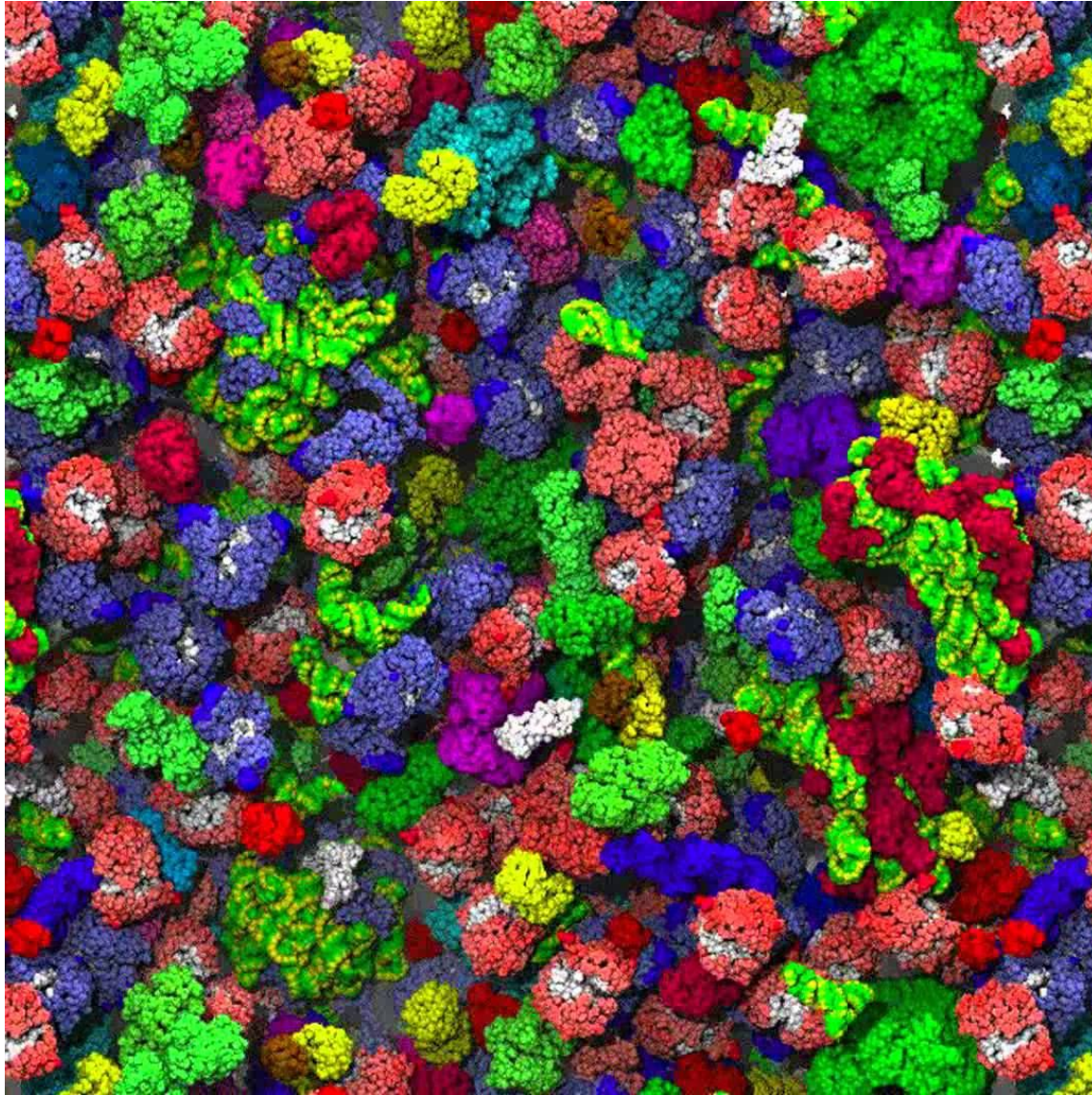
Vesicular membrane edge

Tomogram boundary

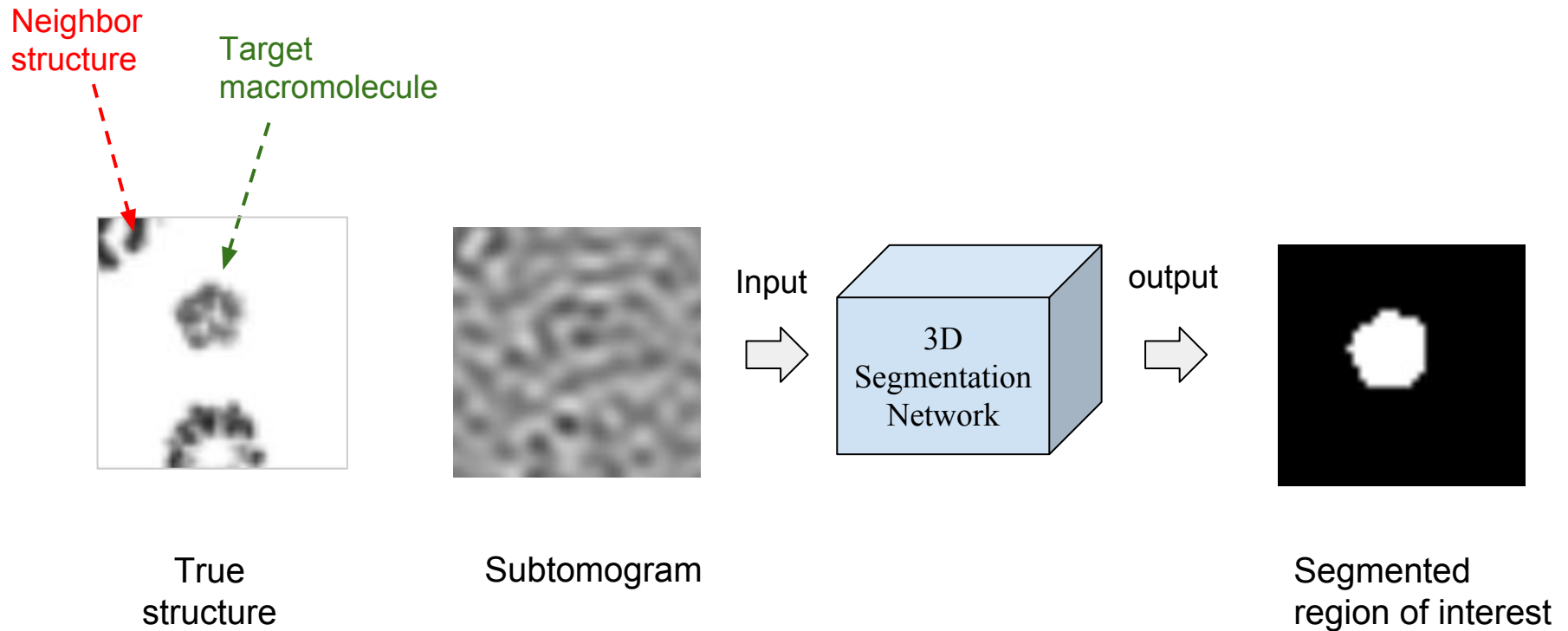


Subtomogram segmentation

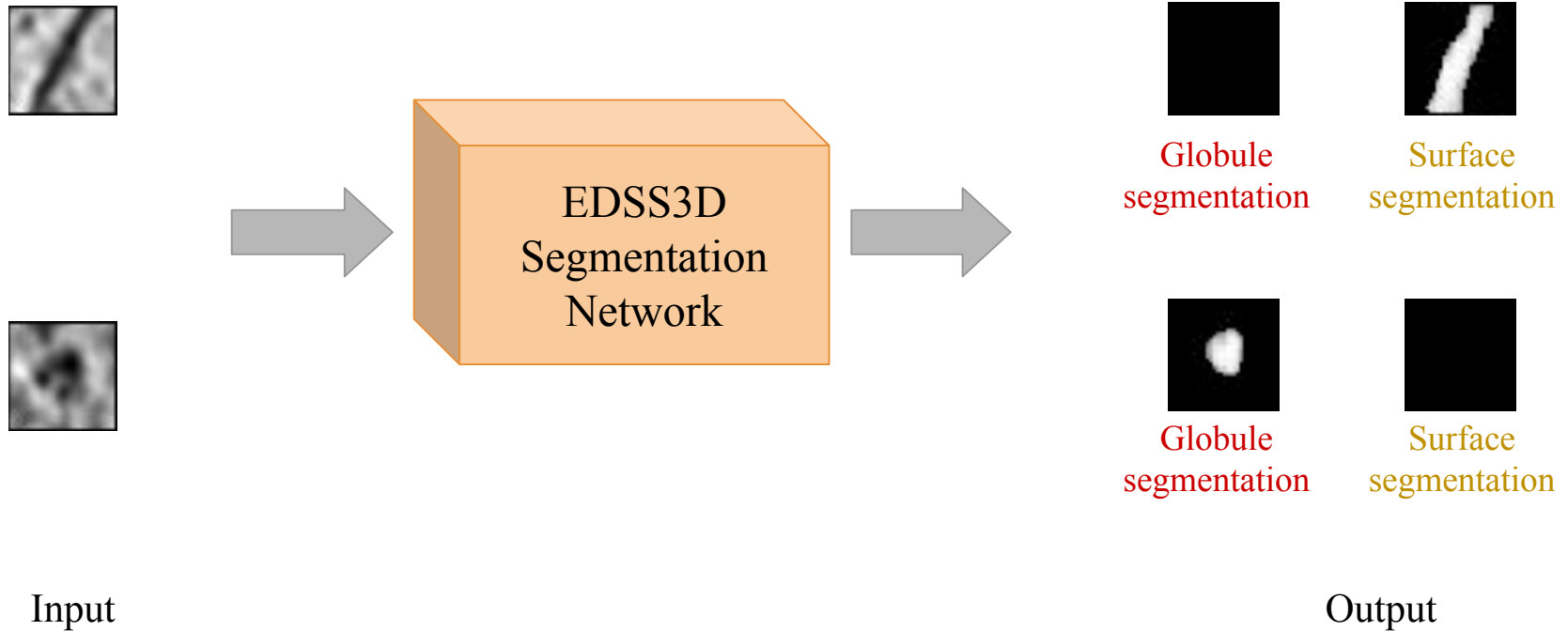
Motivation: molecular crowding



Voxelwise binary classification based segmentation

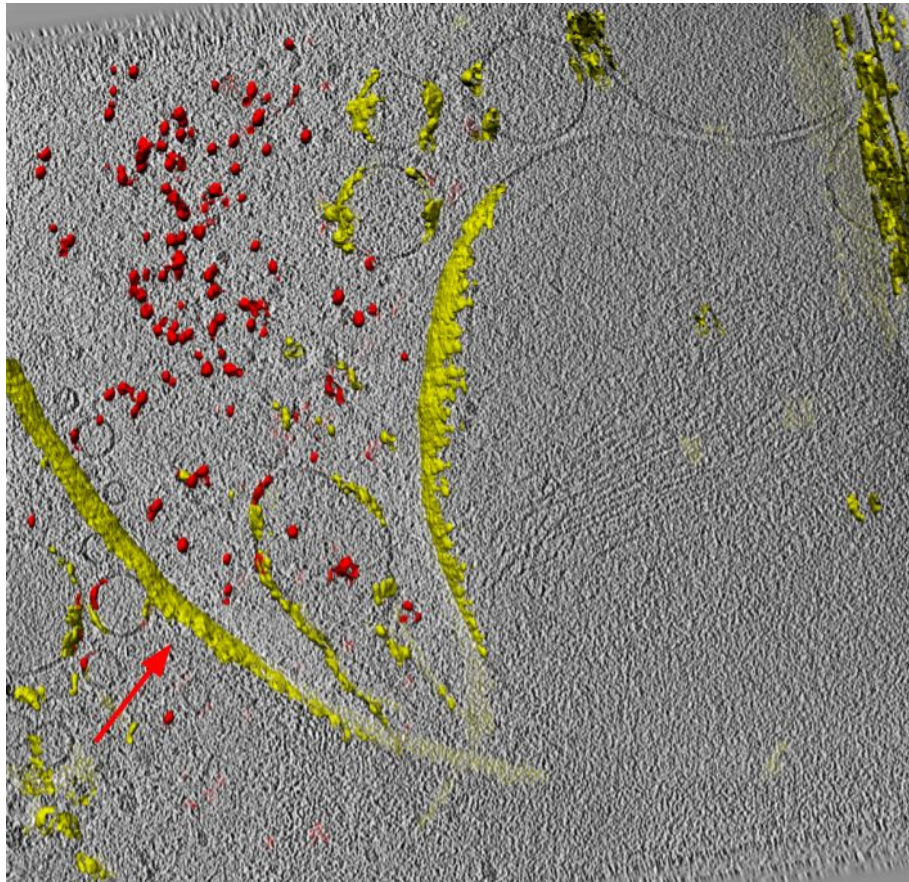


Voxelwise multiclass classification based segmentation

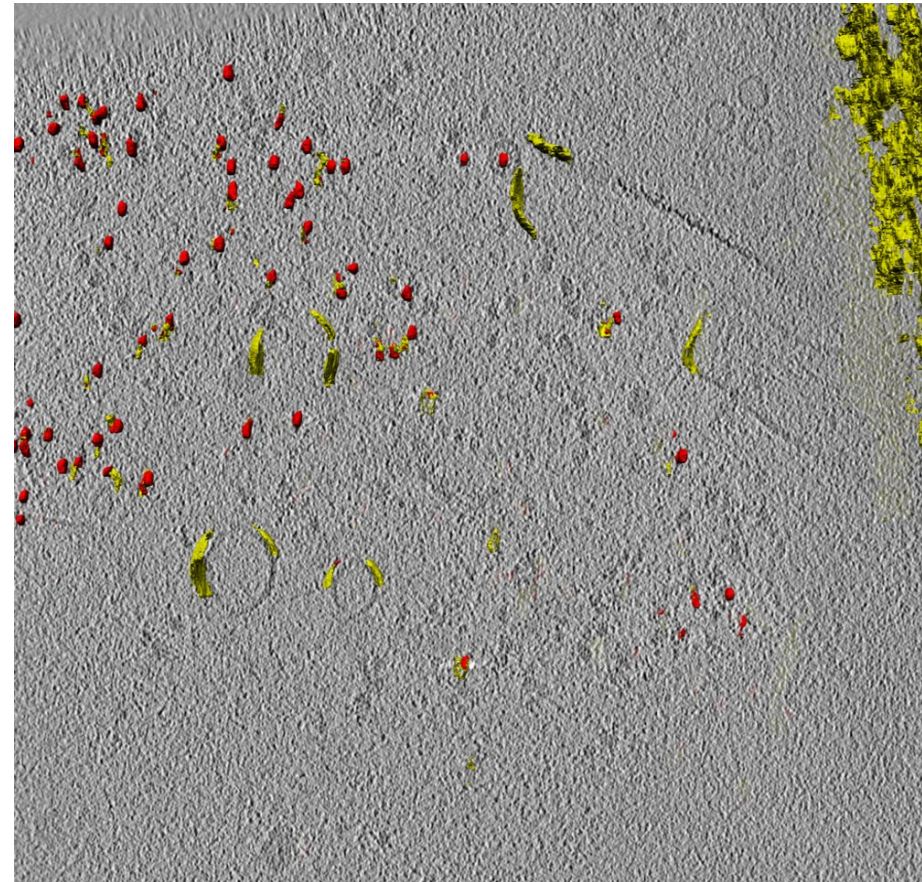


Weakly supervised segmentation

Training tomogram



Testing tomogram



Autoencoder
training



Segmentor
training



Segmenter
prediction

Summary

- Convolutional neural networks are potentially powerful tools for structural pattern mining
- Substantial further works needed to make supervised deep learning practically useful
 - Construction of good training data
 - Optimization of network models
 - Reduction of supervision

Thank you