

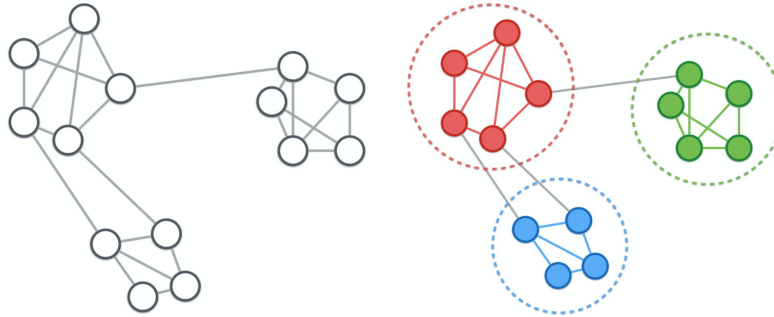
# seqSCAN: Unsupervised Classification of Proteins for New Function Discovery

Meet Barot, Vladimir Gligorijevic, Kyunghyun Cho, Richard Bonneau



# Unsupervised learning is crucial for biological data

- Most data in biology is unlabeled
  - Experiments are expensive, error-prone
  - Supervised learning algorithms are limited by this
- Most label sets are incomplete
- We need tools to get functional categories of proteins whether they have labels or not
- Discovery of these categories can enable us to infer new GO terms, correct old terms, and create entire ontologies



# Clustering with neural networks

- A neural network trained to cluster proteins would be able to be trained batchwise and have constant time inference for a new protein
  - For methods like spectral clustering, inference requires computing pairwise similarities to previous samples
- Using class activation maps, we would be able to highlight specific parts of the protein that correspond to these hypothetical classes
- Feature learning can be combined with the clustering process in neural networks

# Learning to classify images without labels (SCAN)

Step 1:  
Self-supervised  
model to learn  
feature space

$$\min_{\theta} d(\Phi_{\theta}(X_i), \Phi_{\theta}(T[X_i])).$$

Van Gansbeke, Wouter, et al. "Learning To Classify Images Without Labels." *arXiv preprint arXiv:2005.12320* (2020).

# Learning to classify images without labels (SCAN)

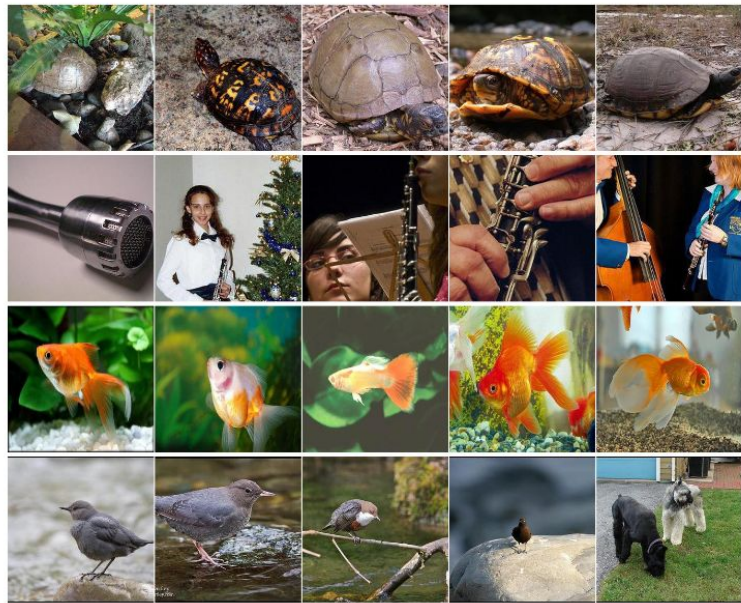


Fig. 1: Images (first column) and their nearest neighbors (other columns) [51].

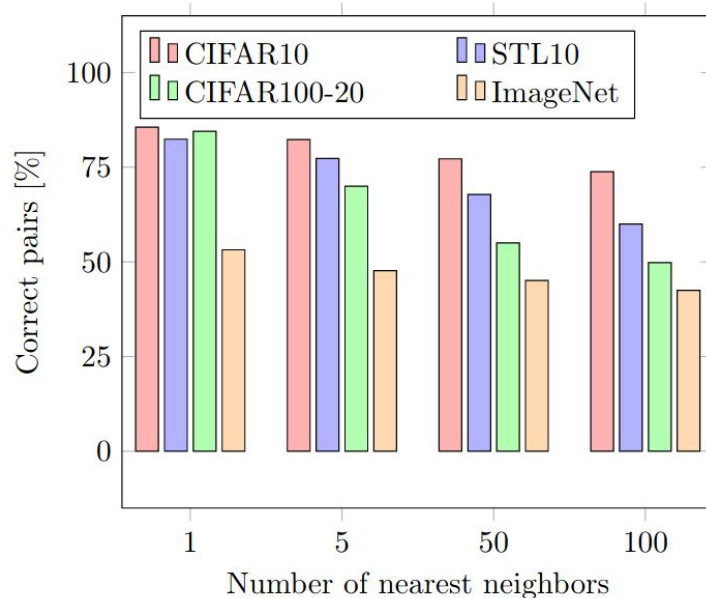


Fig. 2: Neighboring samples tend to be instances of the same semantic class.

# Learning to classify images without labels (SCAN)

Step 2: Train previous neural network, now with a softmax output, with the following loss:

$$\Lambda = -\frac{1}{|\mathcal{D}|} \sum_{X \in \mathcal{D}} \sum_{k \in \mathcal{N}_X} \log \langle \Phi_\eta(X), \Phi_\eta(k) \rangle + \lambda \sum_{c \in \mathcal{C}} \Phi'_\eta{}^c \log \Phi'_\eta{}^c,$$

$$\text{with } \Phi'_\eta{}^c = \frac{1}{|\mathcal{D}|} \sum_{X \in \mathcal{D}} \Phi_\eta^c(X).$$



# Learning to classify images without labels (SCAN)

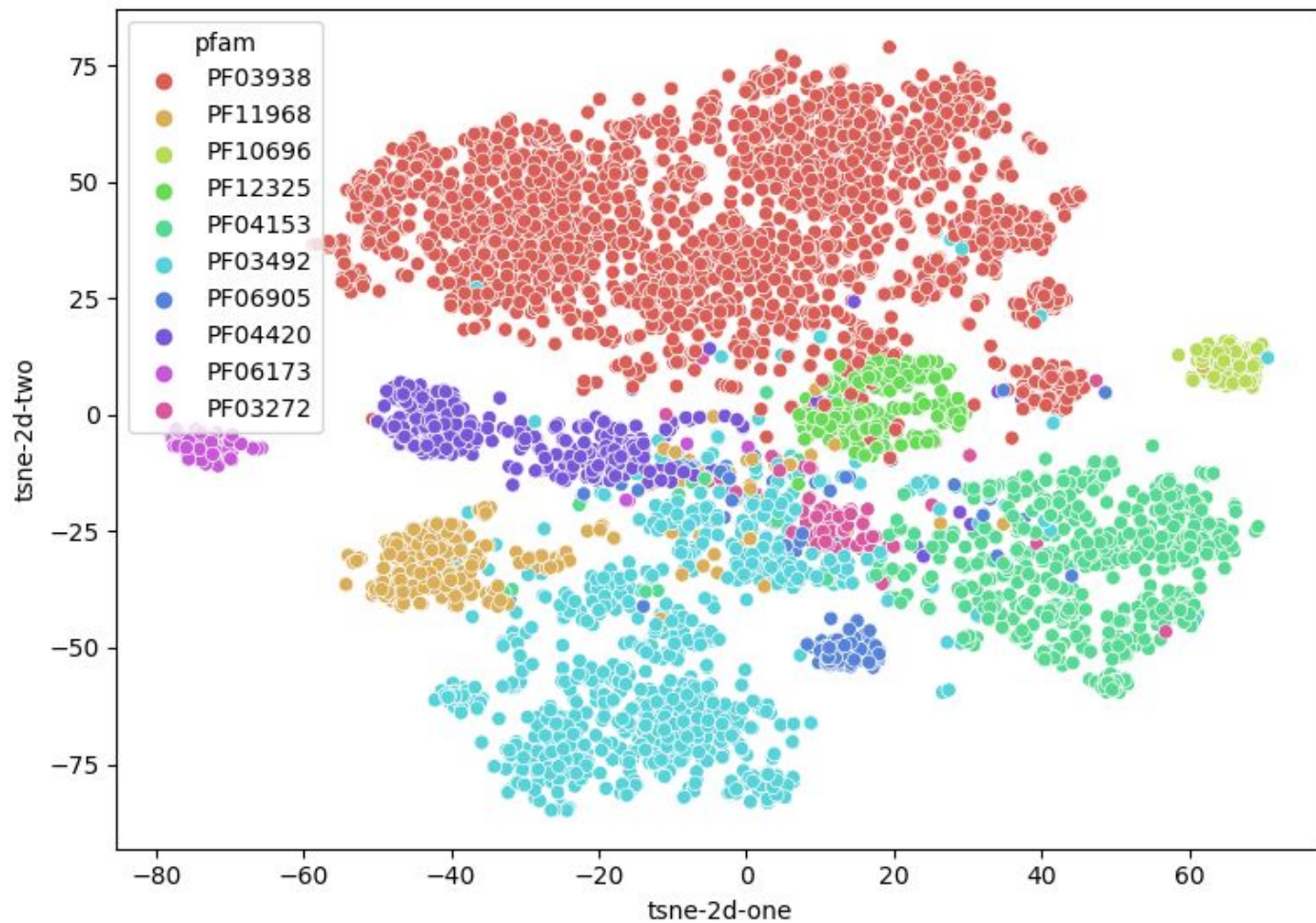
Dataset	CIFAR10			CIFAR100-20			STL10		
Metric	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
Supervised	93.8	86.2	87.0	80.0	68.0	63.2	80.6	65.9	63.1
Pretext <b>Z</b> + K-means	$65.9 \pm 5.7$	$59.8 \pm 2.0$	$50.9 \pm 3.7$	$39.5 \pm 1.9$	$40.2 \pm 1.1$	$23.9 \pm 1.1$	$65.8 \pm 5.1$	$60.4 \pm 2.5$	$50.6 \pm 4.1$
SCAN* (Avg $\pm$ Std)	$81.8 \pm 0.3$	$71.2 \pm 0.4$	$66.5 \pm 0.4$	$42.2 \pm 3.0$	$44.1 \pm 1.0$	$26.7 \pm 1.3$	$75.5 \pm 2.0$	$65.4 \pm 1.2$	$59.0 \pm 1.6$
SCAN <sup>†</sup> (Avg $\pm$ Std)	$87.6 \pm 0.4$	$78.7 \pm 0.5$	$75.8 \pm 0.7$	$45.9 \pm 2.7$	$46.8 \pm 1.3$	$30.1 \pm 2.1$	$76.7 \pm 1.9$	$68.0 \pm 1.2$	$61.6 \pm 1.8$
SCAN <sup>†</sup> (Best)	<b>88.3</b>	<b>79.7</b>	<b>77.2</b>	<b>50.7</b>	<b>48.6</b>	<b>33.3</b>	<b>80.9</b>	<b>69.8</b>	<b>64.6</b>
SCAN <sup>†</sup> (Overcluster)	$86.2 \pm 0.8$	$77.1 \pm 0.1$	$73.8 \pm 1.4$	$55.1 \pm 1.6$	$50.0 \pm 1.1$	$35.7 \pm 1.7$	$76.8 \pm 1.1$	$65.6 \pm 0.8$	$58.6 \pm 1.6$

# Pfam Experiment --- seqSCAN

- Using self-supervised sequence model to extract useful features from sequence (language model trained on 10 million protein sequences from Pfam)
- Dataset of 10 protein families, total ~6000 proteins
- Evaluate clusters obtained using Normalized Mutual Information (NMI) with respect to protein family labels

$$\text{NMI}(\Omega, \mathbf{C}) = \frac{I(\Omega; \mathbf{C})}{[H(\Omega) + H(\mathbf{C})]/2}$$





# Using SCAN loss to cluster proteins using learned features

- Train a single-layer model on the learned features with softmax output with the SCAN loss function:

$$\Lambda = -\frac{1}{|\mathcal{D}|} \sum_{X \in \mathcal{D}} \sum_{k \in \mathcal{N}_X} \log \langle \Phi_\eta(X), \Phi_\eta(k) \rangle + \lambda \sum_{c \in \mathcal{C}} \Phi_\eta'^c \log \Phi_\eta'^c,$$

$$\text{with } \Phi_\eta'^c = \frac{1}{|\mathcal{D}|} \sum_{X \in \mathcal{D}} \Phi_\eta^c(X).$$

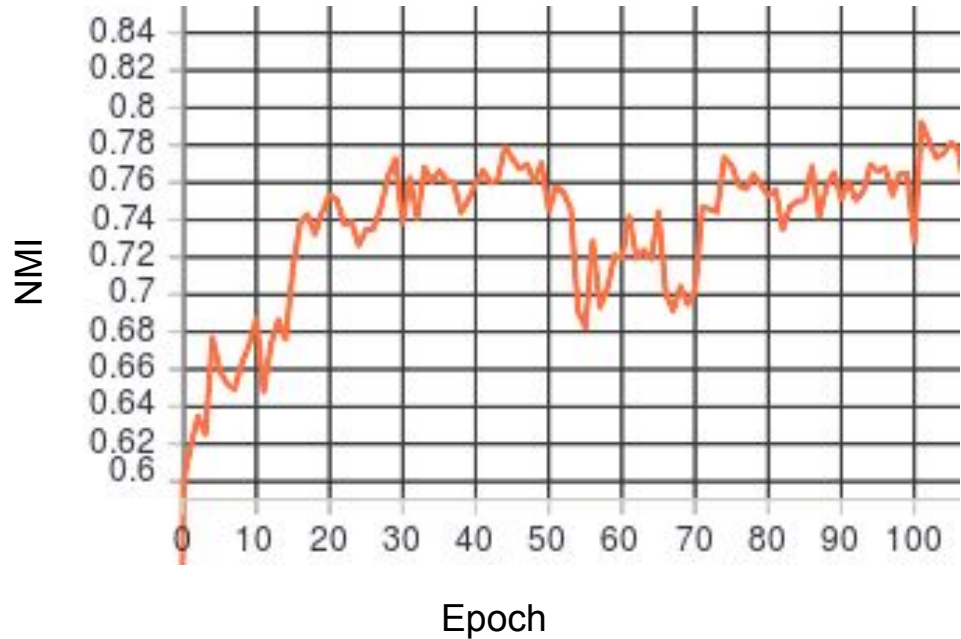
## NMI during “SeqSCAN” Training

Baseline:

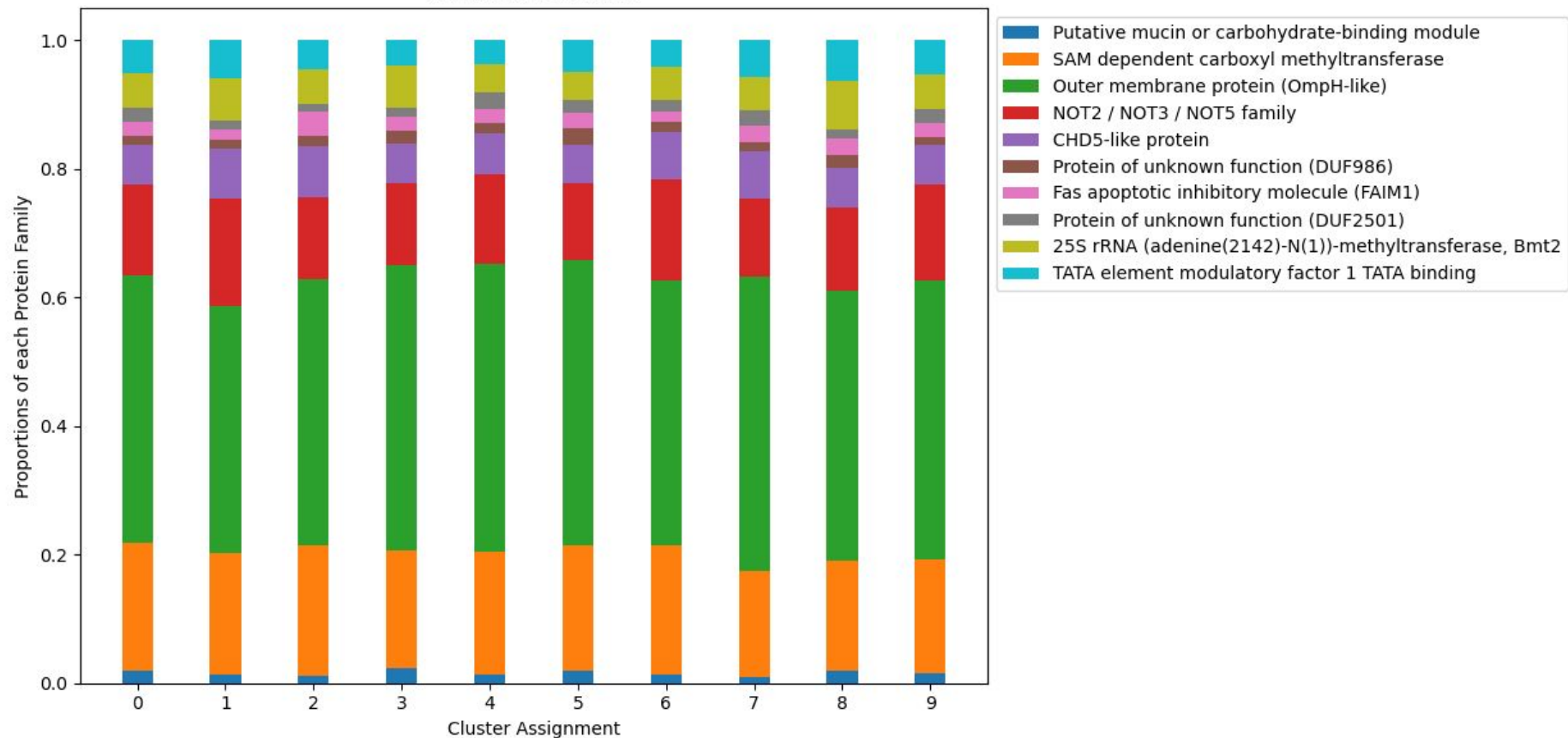
PCA+K-means on  
sequence features

- NMI: 0.61

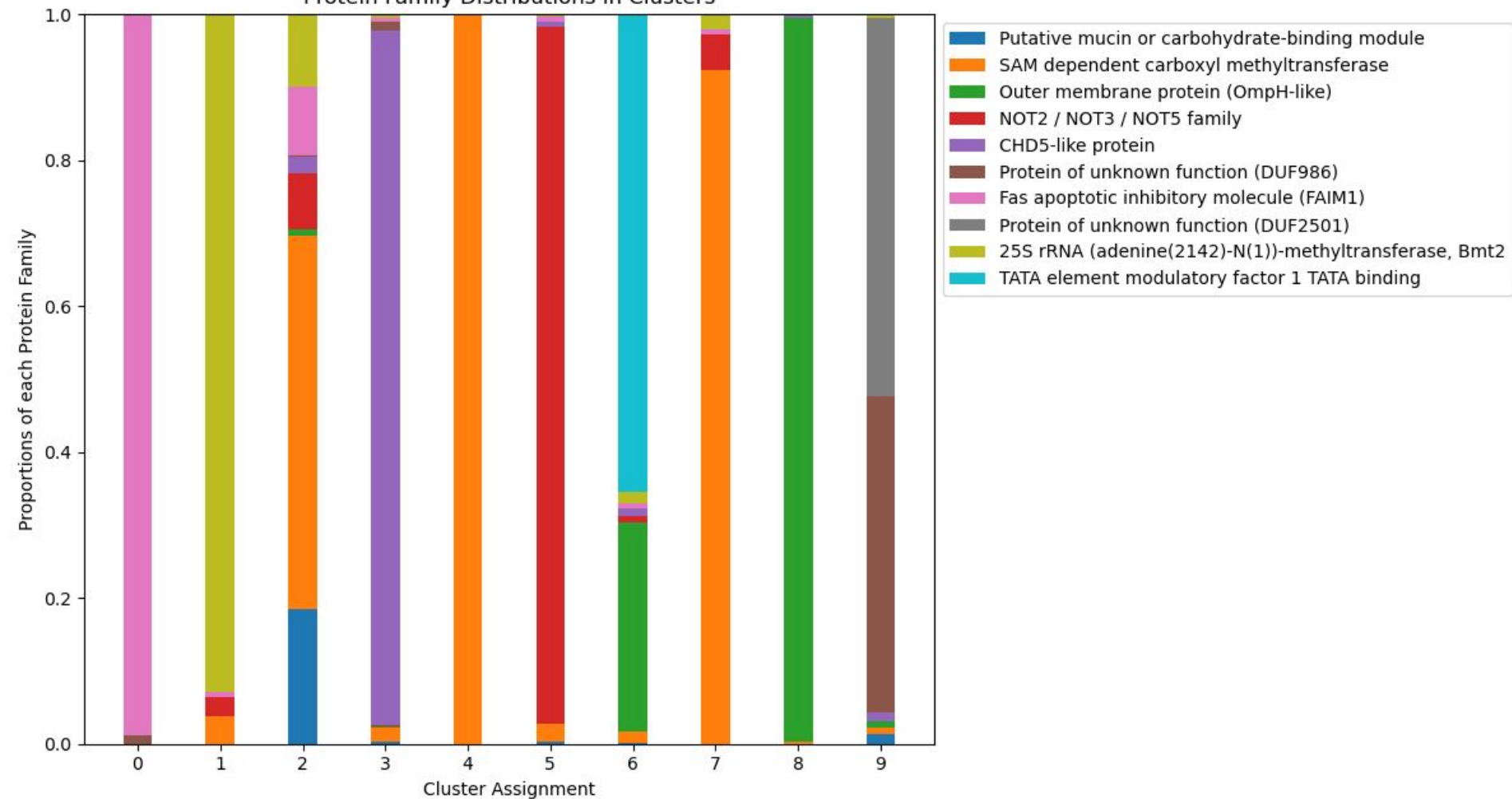
$$\text{NMI}(\Omega, \mathbf{C}) = \frac{I(\Omega; \mathbf{C})}{[H(\Omega) + H(\mathbf{C})]/2}$$



Random Predictions

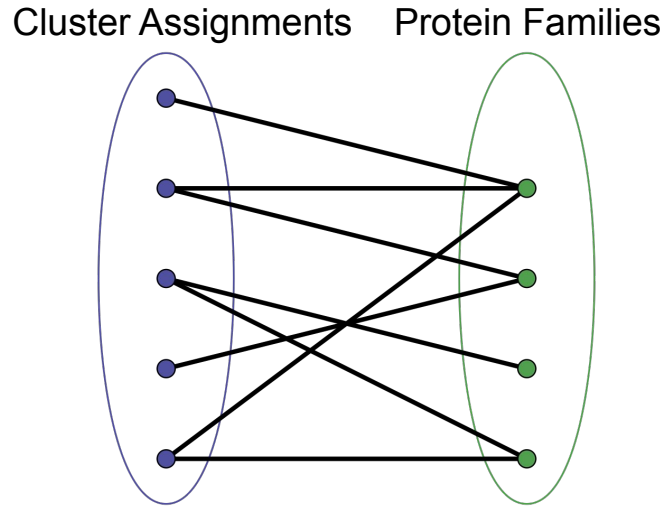


Protein Family Distributions in Clusters



# Matching cluster assignments with labels

- Bipartite matching: maximizing accuracy of the cluster assignments with respect to their protein families







## Scaling up to PfamA

- Training clustering model on 16 million proteins (15k Pfam families)
- Test on 1.8 million proteins (13k families)

	PCA+K-means	seqSCAN
Training NMI	0.499	<b>0.516</b>
Test NMI	0.523	<b>0.541</b>



Vladimir Gligorijević  
Research Scientist



Richard Bonneau



The Bonneau lab!



Kyunghyun Cho

