



# Motivation

- Learning rich representations from unlabeled images
- Reducing the gap between supervised and self-supervised learning (SSL) for smaller models (e.g., MobileNet)
- Compressing a rich SSL model to a smaller one to enable learning at the edge for preserving privacy



## Key Ideas

- Train a high-capacity model using an off-the-shelf self-supervised method and compress it to a small model
- Each data point should have the same relationship with its neighbors in both teacher and student embeddings.

### Method

• Capture the similarity of each data point (query) to the other training data (anchors) in the teacher's embedding space f:

$$p_i(q, a, f) = \frac{exp(f(q)^{\mathsf{T}}f(a_i)/\tau)}{\sum_j exp(f(q)^{\mathsf{T}}f(a_j)/\tau)}$$

- Do the same on the student's embedding space g .
- Define the loss for a particular query point as the KL divergence between the probabilities over all anchor points under the teacher and student models:

L = KL(p(q, a, f) || p(q, a, g))

• We optimize the student by minimizing the summation of L over all images



set of anchor points.

• **Ours-2q** : Use a separate memory bank for teacher and student to decouple the embeddings. • **Ours-1q** : Use the teacher's anchor points in calculating the similarity for the student model.

student as well.

anchor points.

### Met

Sup Con Jigsa Cou Rot Dee RFI SeL MoC Our Ou

# **CompReSS: Self-Supervised Learning by Compressing Representations** Hamed Pirsiavash

#### Soroush Abbasi Koohpayegani\* Ajinkya Tejankar\* University of Maryland Baltimore County

\* Equal contribution



### Results

• Our method is better than other compression methods by a large margin across 3 different evaluation benchmarks and 2 different teacher SSL methods.

• For the first time, a self-supervised AlexNet outperforms supervised one on ImageNet classification • We reduce the gap between supervised and SSL in smaller models

• A linear classifier on our embeddings outperforms finetuning SSL methods on small ImageNet

	Ale	ResNet-50				
hod	ImageNet top-1	Places top-1	PASCA Cls. (mAP)	L VOC Det. (mAP)	Method	ImageNet top-1
. on ImageNet	56.5 (f7)	39.4 (c4)	79.9 (all)	59.1 (all)	Sup. on ImageNet	76.2 (L5)
ntext aw inting Net pCluster	$\begin{array}{c} 31.7 \ (c4) \\ 34.0 \ (c3) \\ 34.3 \ (c3) \\ 38.7 \ (c3) \\ 39.8 \ (c4) \end{array}$	$\begin{array}{c} 32.7 \ (c4) \\ 35.0 \ (c3) \\ 36.3 \ (c3) \\ 35.1 \ (c3) \\ 37.5 \ (c4) \end{array}$	65.3 (all) 67.6 (all) 67.7 (all) 73.0 (all) 73.7 (all)	51.1 (all) 53.2 (all) 51.4 (all) 54.4 (all) 55.4 (all)	SimCLR MoCo InfoMin BYOL SwAV	69.3 (L5) 71.1 (L5) 73.0 (L5) 74.3 (L5) <b>75.3</b> (L5)
Decouple a Co rs-2q rs-1q	44.3 (c5) 44.7 (c5) 45.7 (f7) 57.6 (f7) <b>59.0</b> (f7)	$\begin{array}{c} 38.6 \ ({\rm c5}) \\ 37.9 \ ({\rm c4}) \\ 36.6 \ ({\rm c4}) \\ {\bf 40.4} \ ({\rm c5}) \\ 40.3 \ ({\rm c5}) \end{array}$	74.7 (all) 77.2 (all) 71.3 (f8) <b>79.7</b> (f8) 76.2 (f8)	58.0 (all) 59.2 (all) 55.8 (all) 58.1 (all) <b>59.3</b> (all)	Ours-1q Compressed from ResNet-50 (	71.9 (L5) n SimCLR (x4)







the two distributions.

- K-means on our AlexNet embeddings (k = 1000)
- Each row is a cluster
- No cherry-picking: random images from random clusters

### Comparison with other compression methods

_								
-	Teacher Student	MoC R	SwAV ResNet-50 ResNet-18					
-		linear	NN	CA	liı	near	NN	CA
-	Teacher	70.8	57.3	34.2	75.6		60.7	27.6
-	Supervised	69.8	63.0	44.9	69.8		63.0	44.9
-	CC	61.1	51.1	25.2	6	0.8	51.0	22.8
	CRD	58.4	43.7	17.4	5	8.2	44.7	16.9
	Reg-BN	58.2	47.3	27.2	60.6		47.6	20.8
	Ours-1q	62.6	53.5	33.0	6	5.6	56.0	26.3
-								
			Method		Top-1		Top-5	
					1%	10%	1%	10%
Сс	omparison w	Supervised		25.4	56.4	48.4	80.4	
other SSI methods		All layers are fine-tuned.						
on small labelled ImageNet for ResNet-50			InstDisc		_	-	39.2	77.4
			PIRL		-	-	57.2	83.8
			SimCLR		48.3	65.6	75.5	87.8
			BYOL		53.2	68.8	78.4	89.0
			SwAV		53.9	70.2	78.5	89.9

Only the linear layer is trained. **59.7** 67.0 **82.3** 87.5 ≪ Ours-1q