

A critical analysis of self-supervision, or what we can learn from a single image

Yuki M. Asano

CDT Annual Meeting Oct 2019

work with *Christian Rupprecht* and *Andrea Vedaldi* at VGG

Outline

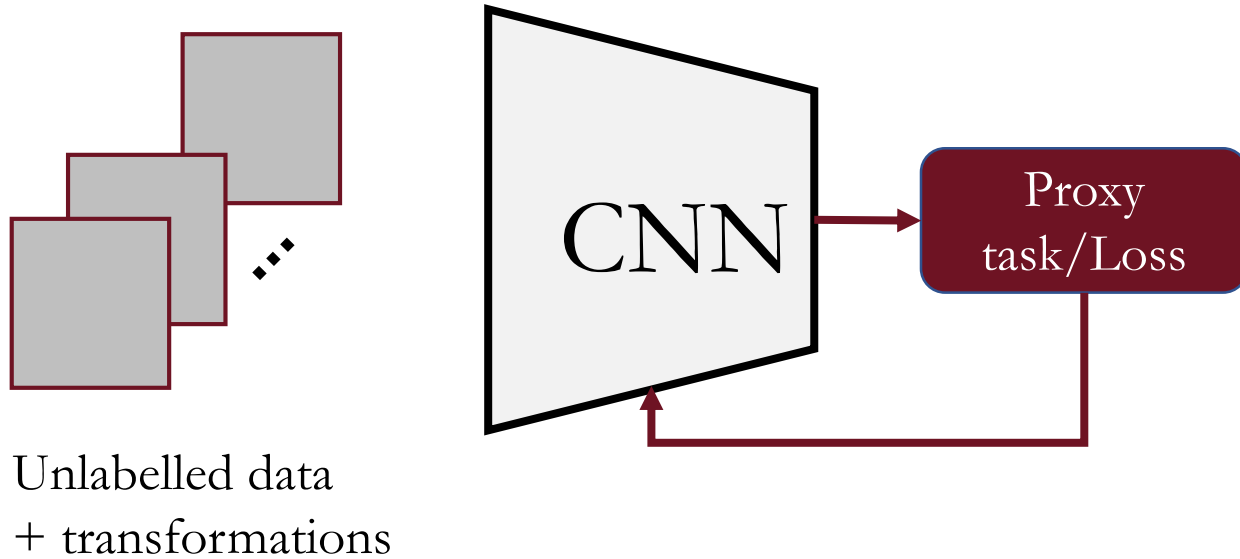
- Self-supervised learning saga
- *Or is it?*

Self-supervised learning like we do



1. Unlabelled, large collection of images
2. Train your network *without labels*
3. Use the image representations (vectors) for new tasks

Self-supervised learning like we do?



e.g. DeepCluster

- Run k-means on features
- Train classifier on k classes
- Repeat for 200 epochs

e.g. RotNet

- Create 4 classes based on rotations
- Exploits photographer bias
- Simple but works

Or colorizing images



Hypothesis

“Priors”

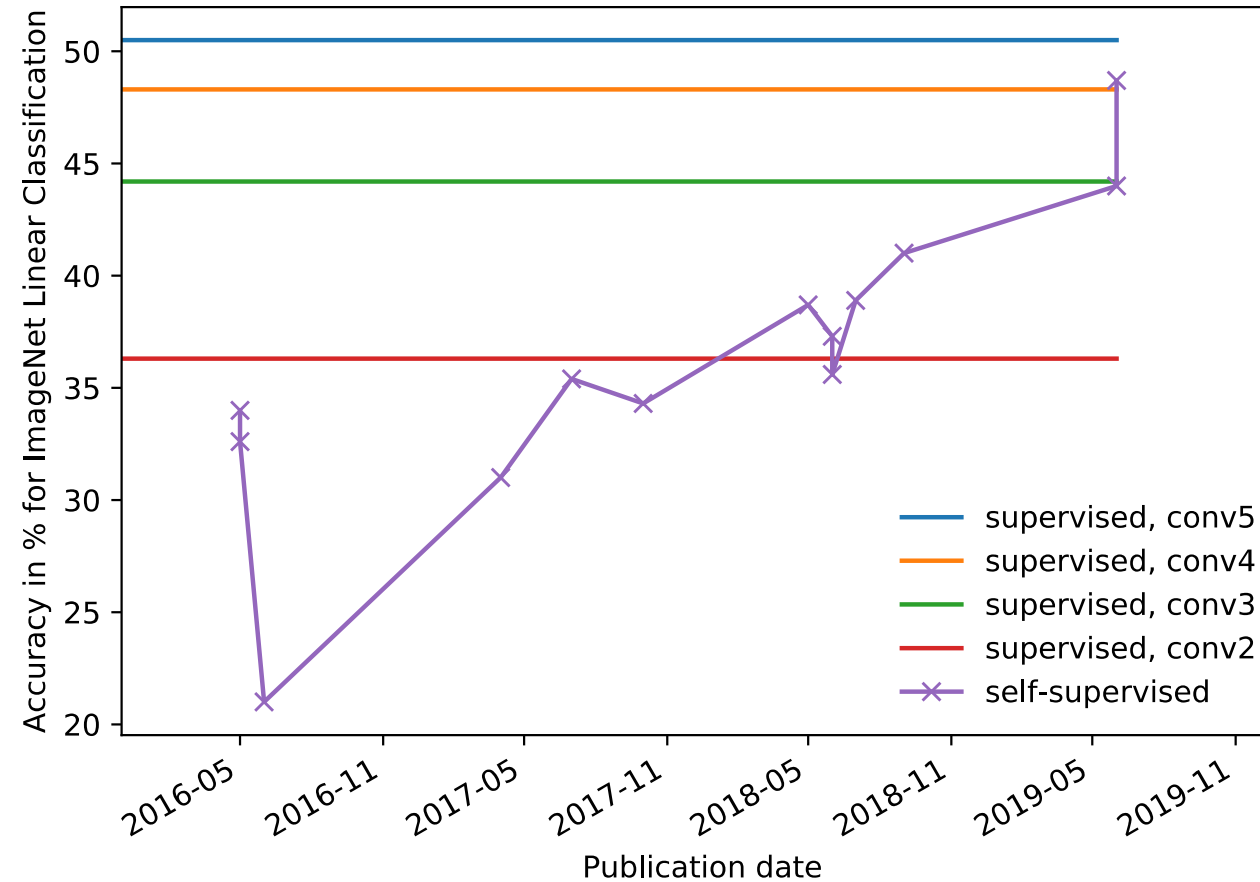
What/how humans learn



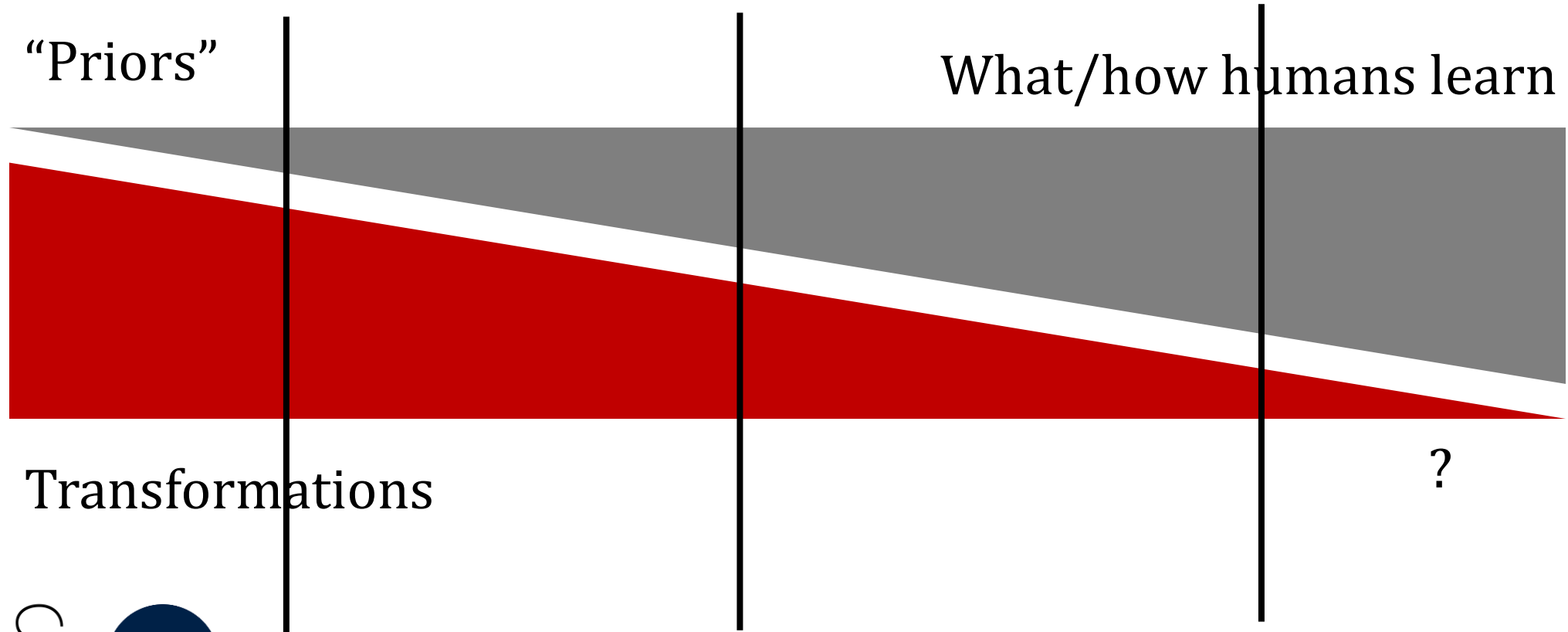
Transformations

?

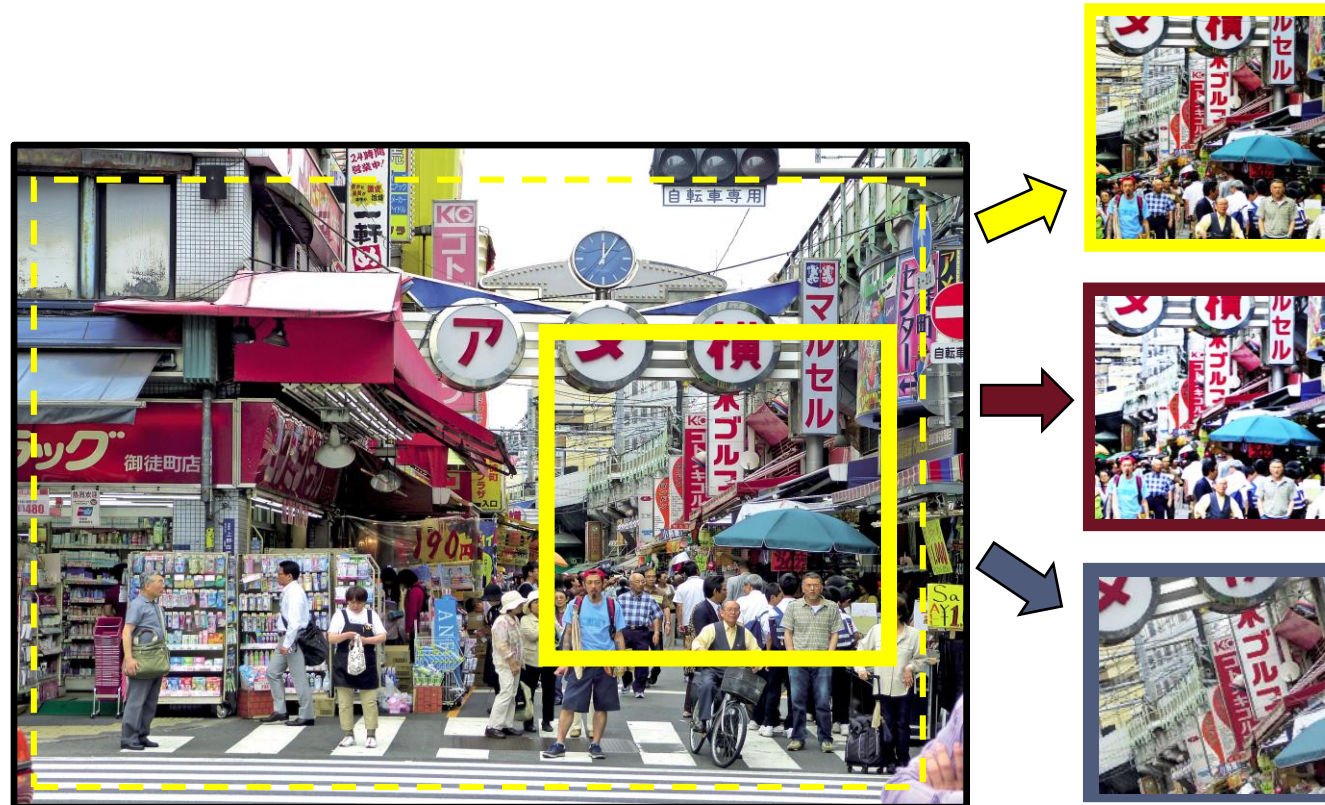
Getting there, but not quite yet



Where are we?



“Learn” from one image... using multiple transformations



Learned first convolutional layer – from one image

Method, Image A

Method, Image B

BiGAN

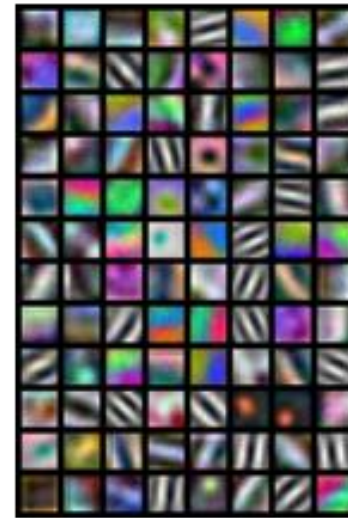
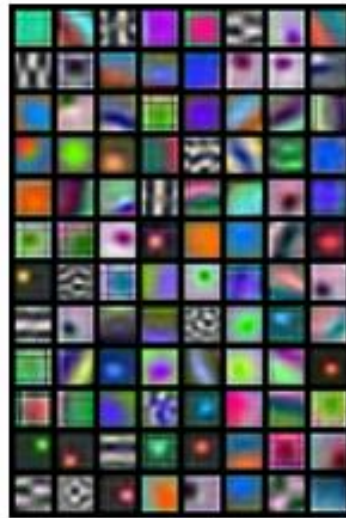
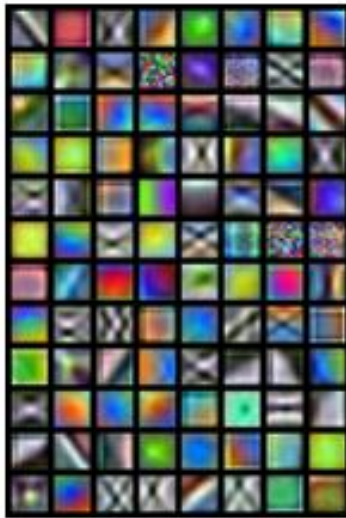
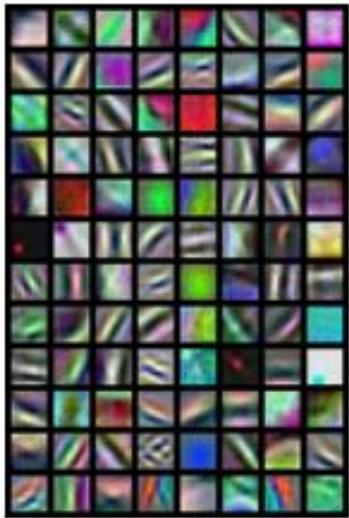
RotNet

DeepCluster

BiGAN

RotNet

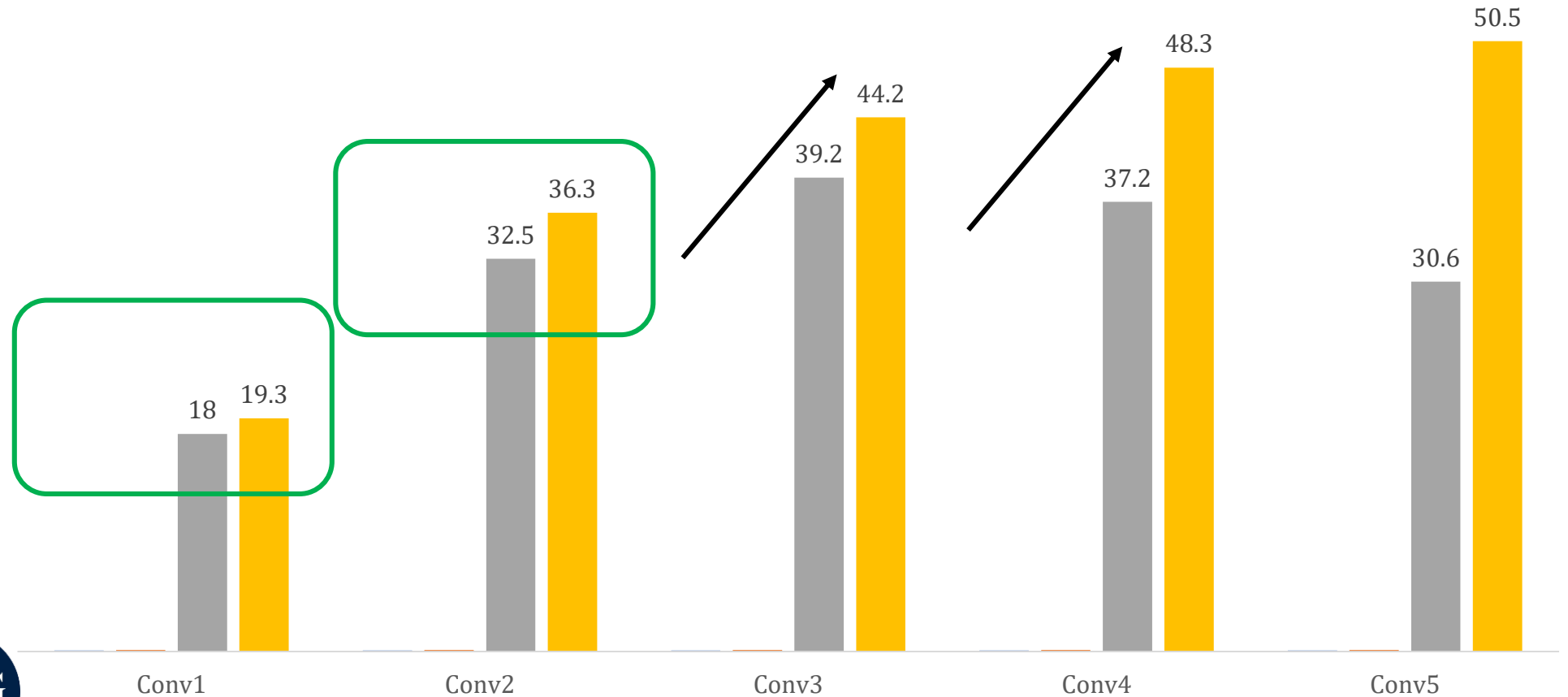
DeepCluster



Performance

Comparison of random, DeepCluster (1 & 1M images) and supervised

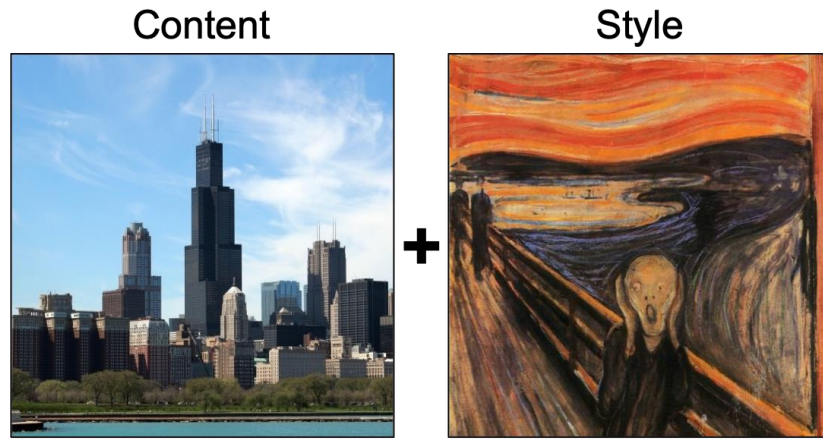
■ Random ■ 1-image ■ 1M images ■ Supervised



Conclusion

1. Early layers of deep networks contain limited information about natural images
2. These can be learned through self-supervision or supervised learning
3. Notably, only one **image + transformations** are necessary for this
4. Much space to go the *right* direction in self-supervised learning

Style transfer with a 1-image trained CNN



Appendix

Table 4: **Finetuning experiments** The pre-trained model’s conv1 and conv2 are left frozen and only the higher levels are re-trained using ImageNet LSVRC-12 training set. Accuracy is averaged over 10 crops.

	c1	c2	c3	c4	c5
Full sup.	19.3	36.3	44.2	48.3	50.5
BiGAN, A	22.5	37.6	44.2	47.6	48.3
RotNet, A	22.0	38.2	44.8	49.2	51.8
DeepCluster, A	21.8	35.9	43.6	48.8	50.4

		CIFAR-10			
		conv1	conv2	conv3	conv4
(a)	Fully sup.	66.5	70.1	72.4	75.9
(b)	Random feat.	57.8	55.5	54.2	47.3
(c)	No aug.	57.9	56.2	54.2	47.8
(d)	Jitter	58.9	58.0	57.0	49.8
(e)	Rotation	61.4	58.8	56.1	47.5
(f)	Scale	<u>67.9</u>	<u>69.3</u>	<u>67.9</u>	<u>59.1</u>
(g)	Rot. & jitter	64.9	63.6	61.0	53.4
(h)	Rot. & scale	67.6	69.9	68.0	60.7
(i)	Jitter & scale	<u>68.1</u>	<u>71.3</u>	<u>69.5</u>	<u>62.4</u>
(j)	All	68.1	72.3	70.8	63.5