

- **Goal:** Learning rich representations from unlabeled images



Iterative Similarity Distillation

- We can view the recent direction of SSL methods as iterative self-distillation where there is a teacher and a student.
- Our study shows that the teacher is always better in intermediate checkpoints
- We iteratively distill the teacher to the student and update the teacher as a moving average



Key Ideas

- In contrastive learning, all negatives are not equally negative
 - Relaxing binary instance classification of contrastive learning with soft classification
- Since teacher has better accuracy, use the teacher as the source of knowledge to assign soft-labels to negatives

ISD: Self-Supervised Learning by Iterative Similarity Distillation

Ajinkya Tejankar^{1,}* Soroush Abbasi Koohpayegani^{1,*} ¹University of Maryland Baltimore County * Equal contribution



Vipin Pillai¹ University Of Bern

Hamed Pirsiavash³ Paolo Favaro² ³University of California, Davis



https://github.com/UMBCvision/ISD

On ImageNet task, our method outperforms other SOTA methods in

	Batch	Epochs	Sym. Loss	Top-1	NN	20-NN					
	Size		2x FLOPS	Linear							
ResNet-50											
d	256	100	-	76.2	71.4	74.8					
	4096	200	·	69.1		-					
	4096	1000	1	69.3	-	-					
2	256	200	1	69.9	-	-					
	256	200	1	70.0	-	-					
	4096	200	1	70.6	-	-					
2	256	400	1	71.0	-	-					
2	256	800	×	71.1	57.3	61.0					
s	256	1K+130	×	71.9	63.3	66.8					
	4096	1000	1	74.3	62.8	66.9					
	4096	800	✓	75.3	-	-					
2	256	200	×	67.5							
	256	200	×	68.0	-	-					
ym	256	200	×	69.3	55.0	59.2					
	256	200	X	69.8	59.2	62.0					
ResNet-18											
d	256	100	-	69.8	63.0	67.6					
2 - 1	256	200	× ×	51.0	37.7	42.1					
ym	256	200	×	52.6	40.0	44.8					
	256	200	X	53.8	41.5	46.6					

Comparison on ImageNet with Limited Labels for ResNet-50

4	Encoha	Top-1		Top-5					
bu	Epocus	1%	10%	1%	10%				
network is fir	ne-tuned.								
vised		25.4	56.4	48.4	80.4				
	800	-	-	57.2	83.8				
	200	-	-	71.0	85.7				
LR	1000	48.3	65.6	75.5	87.8				
	800	-	-	78.2	88.7				
_	1000	53.2	68.8	78.4	89.0				
7	800	53.9	70.2	78.5	89.9				
the linear layer is trained.									
	1000	55.7	68.6	80.0	88.6				
Ress	1K+130	59.7	67.0	82.3	87.5				
	200	43.6	58.4	71.2	82.9				
asym	200	47.9	61.3	74.6	84.7				
	200	53.4	63.0	78.8	85.9				