

Scene-centric data are easier to collect, and contain more information per image.



But effective pre-training on them requires more than instance discrimination. We endorse **object-level** contrastive learning.

Then, how to find the objects w/o supervision?

Hand-crafted object-proposal methods like saliency, selectivesearch, k-means clustering may limit the upper bound of learned representations.

We propose *learnable semantic grouping* for online object discovery.



- First learn semantic prototypes from the whole dataset, where each prototype can represent a semantic class (e.g., cat, dog).
- Then assign a nearest-neighbor prototype to each pixel.
- Pixels with the same pseudo-label forms a group (object).

Priors for objectness

- Prior 1: geometric-covariance and photometric-invariance
- Prior 2: small prototype number
- Prior 3: meaningful grouping. i.e., avoiding collapse

Table 1: Details of the	he data
Dateset	#Img
ImageNet-1K [15] COCO [48] COCO+ [48]	1.281 1181 2411

SlotCon: Self-Supervised Visual Representation Learning with Semantic Grouping

Xin Wen¹ Bingchen Zhao² ¹The University of Hong Kong

Solving object discovery & representation learning jointly



Online Object discovery

- Rand. init. the prototypes
- Train w/ pixel-level deep clustering
- Pixels w/ same assign. form an object

sets used for pre-training. #Obj./Img. #Class 10001.77.380 N/A N/A

New SOTA on scene-centric pre-training

Table 2: Main transfer results with COCO pre-training. We report the results in COCO [48] object detection, COCO instance segmentation, and semantic segmentation in Cityscapes [13], PASCAL VOC [20] and ADE20K [86]. Compared with other image-, pixel-, and object-level self-supervised learning methods, our method shows consistent improvements over different tasks without leveraging multi-crop [6] and objectness priors. (†: re-impl. w/ official weights; ‡: full re-impl.)

Method	Epochs	Multi crop	i Obj. Prior	COCO detection		COCO segmentation			Semantic seg. (mIoU)			
				$\overline{AP^{b}}$	AP_{50}^{b}	AP_{75}^b	$\overline{AP^m}$	AP_{50}^m	AP_{75}^m	City.	VOC	ADE
random init.	-	×	×	32.8	50.9	35.3	29.9	47.9	32.0	65.3	39.5	29.4
Image-level app	Image-level approaches											
MoCo v2 [‡] [9]	800	×	×	38.5	58.1	42.1	34.8	55.3	37.3	73.8	69.2	36.2
Revisit. [†] [63]	800	1	×	40.1	60.2	43.6	36.3	57.3	38.9	75.3	70.6	37.0
Pixel-level approaches												
Self-EMD [49]	800	×	×	39.3	60.1	42.8	-	-	-	-	-	-
DenseCL [†] [68]	800	×	×	39.6	59.3	43.3	35.7	56.5	38.4	75.8	71.6	37.1
PixPro [‡] [75]	800	×	×	40.5	60.5	44.0	36.6	57.8	39.0	75.2	72.0	38.3
Object / Group-level approaches												
DetCon [†] [35]	1000	×	1	39.8	59.5	43.5	35.9	56.4	38.7	76.1	70.2	38.1
ORL [†] [74]	800	1	1	40.3	60.2	44.4	36.3	57.3	38.9	75.6	70.9	36.7
Ours (SlotCon)	800	X	×	41.0	61.1	45.0	37.0	58.3	39.8	76.2	71.6	39.0

Eqv. Performance to ImageNet with only ¹/₅ data

Table 3: Pushing the limit of scene-centric pre-training. Our method further sees a notable gain in all tasks with extended COCO+ data, showing the great potential of scene-centric pre-training.

Method Dataset Ep		Epochs	COCO detection			COCO segmentation			Semantic seg. (mIoU)		
	2		AP ^b	AP_{50}^b	AP_{75}^b	AP^m	AP_{50}^m	AP_{75}^m	City.	VOC	ADE
SlotCon SlotCon SlotCon	COCO ImageNet ImageNet	800 100 200	41.0 41.4 41.8	61.1 61.6 62.2	45.0 45.6 45.7	37.0 37.2 37.8	58.3 58.5 59.1	39.8 39.9 40.7	76.2 75.4 76.3	71.6 73.1 75.0	39.0 38.6 38.8
ORL [74] SlotCon	COCO+ COCO+	800 800	40.6 41.8	60.8 62.2	44.5 45.8	36.7 37.8	57.9 59.4	39.3 40.6	- 76.5	73.9	39.2

Anlin Zheng^{1,3} Xiangyu Zhang³ Xiaojuan Qi¹ ²University of Edinburgh ³MEGVII Technology

Object-level contrastive learning

 Pool obj-level rep.(slots) from each view • Slots w/o assigned pixels are filtered out Contrastive learning between slots

Evaluation on unsupervised semantic segmentation

Table 5: Main results in COCO-Stuff unsupervised semantic segmentation.

Method	mIoU	pAcc
MaskContrast [64]	8.86	23.03
PiCIE + H. [40]	14.36	49.99
SegDiscover [38]	14.34	56.53
<i>Ours</i> (SlotCon)	18.26	42.36

The model managed to discover visual concepts



Figure 3: Examples of visual concepts discovered by SlotCon from the COCO val2017 split. Each column shows the top 5 segments retrieved with the same prototype, marked with reddish masks or arrows. Our method can discover visual concepts across various scenarios and semantic granularities regardless of small object size and occlusion. (best viewed in color)

Though it can be biased toward occupying classes, e.g., persons



Set the number of protos close to real class number

K	COCO	City	VOC	ADI
128	40.7	76.4	71.9	38.5
256	41.0	76.2	71.6	39.0
512	40.9	75.6	71.6	38.9
1024	40.7	75.8	70.9	39. 1











Geometric augs are crucial for sem. grouping

