

2

Scaling up Instance Segmentation using Approximately Localized Phrases

Meta Al

Karan Desai¹

Justin Johnson ^{1,2} Ishan Misra² ¹ University of Michigan ² Meta Al

Training data with ALPs

ALPs: natural language phrases paired with boxes derived from mouse scribbles.



Model: FCOS-MO

FCOS-style model with a deep Hourglass mask head and a CLIP-based open-vocabulary classifier.



FCOS uses box centers for localization supervision during training instead of box area or edges as used by Mask R-CNN. This makes it less sensitive to noise in imprecise ALP boxes.

- Training: alternate batches of masks (base classes) and ALPs (novel classes). - Box quantization: Training augmentation to randomly expand boxes of ALPs. Increases the coverage of underlying object in the box, provides less noisy training supervision.

Motivation and Goal

Labeled masks and boxes are expensive to collect.

Learn to segment **novel object classes** without using any training masks or labeled boxes for them.

Training instance segmentation models with:

Mask annotations (base classes)



Laurens van der Maaten²



Main Idea

Approximately Localized Phrases (novel classes)

4

Quantitative Results

Main experiments: Train Mask R-CNN* (baseline) and FCOS-MO models using COCO masks and different forms of supervision for novel classes, for learning to segment 300 classes from Open Images (60 base, 240 novel).

Evaluation: Open Images test AP_{50} for base and novel classes.

Model	Novel class supervision	AP ₅₀ -base
Open-vocab Mask R-CNN*	none	62.8
ALP-supervised Mask R-CNN*	ALPs	56.7
ALP-supervised FCOS-MO	ALPs	60.0
Box-supervised Mask R-CNN* (oracle)	Labeled boxes	61.8

Observation 1: ALP-supervised models outperform an open-vocab baseline. **Observation 2:** When trained with ALPs, FCOS-MO outperforms Mask R-CNN* indicating its effectiveness in handling noisy supervision of ALPs.

... see our paper for evaluations on COCO subsets and ablation studies.

Qualitative Results

Predictions on Open Images from ALP-supervised FCOS-MO.

Accurate labels with accurate masks:



Duplicate inaccurate labe accurate mask:



Inaccurate label with a fairly accurate mask:







