Pixel Consensus Voting for Panoptic Segmentation

Haochen Wang

Ruotian Luo

Michael Maire Greg Shakhnarovich





Object detection/Instance segmentation



PASCAL VOC 2007 Bounding Box Detection



MSCOCO (2014): Bounding Box and Instance Mask

Object detection/Instance segmentation

 reason about densely enumerated bounding boxes, and then refine to instance masks.



The Mask R-CNN framework for instance segmentation

Semantic segmentation

- Label each pixel in the image with a category label
- Don't differentiate instances, only care about pixels



Semantic segmentation



Panoptic Segmentation

- Instance-aware semantic segmentation.
- Things (bicycle, dog, car, person)
- Stuff (pavement, ground, dirt, wall)



Existing methods for panoptic segmentation

Merge instance segmentation and semantic segmentation



However?

- Panoptic Segmentation removes the concept of boxes and focuses on pixels.
- Bounding boxes may not be the optimal intermediate representations to predict.
- Pixel Consensus Voting: pixels vote for object centroid
 - Training: per pixel classification
 - Inference: Generalized Hough Transform similar to Implicit Shape Model (ISM)



B. Leibe, A. Leonardis, and B. Schiele, Robust Object Detection with Interleaved Categorization and Segmentation, International Journal of Computer Vision, Vol. 77(1-3), 2008.

Pipeline

- Discretize regions around each pixel. CNN classifies the likely regions that contain centroid
- Vote aggregation: probabilities at each pixel are cast to corresponding regions through dilated deconvolution.
- Peaks in voting heatmap are detections.
- Back-projection by convolving the query filter within a peak region to get an instance mask.
- Category information provided by the parallel semantic segmentation head



Classification or Regression

- A 2d offset vector is a limited representation.
- Regression fails to capture uncertainty. Spurious peaks and false positives.
- Insight echoed by bounding box detectors: classify a proposal into anchors rather than direct regression.



Direct regression might create spurious peaks





Classification models uncertainty about regions

Classification or Regression

- Perfect centroid prediction is unnecessary. What matters is consensus.
- Tolerance for coarse prediction depends on the scale.
- Easier to learn and train.









Pixels of a large object need only a rough estimate

Pixels of a small object need to be precise

Discretization Scheme

- Discretize regions around each pixel into radially expanding cells.
- Voting mask records the ground truth label if the centroid falls into a cell.
- Full discretization has 233 cells covering 243^2 regions at 1/4 input resolution.

			-						
12	12	12	11	11	11	10	10	10	
12	12	12	11	11	11	10	10	10	
12	12	12	11	11	11	10	10	10	
13	13	13	4	3	2	9	9	9	
13	13	13	5	0	1	9	9	9	
13	13	13	6	7	8	9	9	9	
14	14	14	15	15	15	16	16	16	
14	14	14	15	15	15	16	16	16	
14	14	14	15	15	15	16	16	16	
			-						

A toy discretization of the 9 x 9 region around location (4, 4) There are 17 cells.



To assign voting label, align mask center with the current pixel, read off the label from the region that contains the centroid.



Full discretization of the 243² region around each pixe There are 233 cells.

Voting as Dilated Deconvolution

- Spreads probabilistic votes from a point to spatial locations
- Dilation allows a pixel to send its votes afar
- Fixed kernel parameters to enable voting



Voting as Dilated Deconvolution

- The previous step sends votes to points, but we are voting for regions.
- Average pooling spreads the votes evenly within each cell
- Overall voting takes 1.3ms/image over COCO val.



Smooting kernel shape [1, 1, 3, 3], no dilation





Peak detection

• Thresholding + connected component



Back-projection as convolutional filtering

- For a peak, find the pixels that favor this peak above all others
- Query filter: "inward reflection" of the voting mask. Convolve in a peak region to pick up instance mask.
- 81.8ms/image.







Voting mask: gt votes by a single pixel for possible centroids

Query filter: gt votes by surrounding pixels for a particular centroid















Network



Qualitative results

	Methods	Split	PQ	SQ	RQ	$ PQ^{th} $	SQ^{th}	RQ^{th}	$ PQ^{st} $
Mask R-CNN	PFPN [24] (1x) PFPN [24] (3x) UPSNet [62] (1x)	val val val	39.39 41.48 42.5	77.83 79.08 78.0	48.31 50.52 52.5	45.91 48.26 48.6	80.85 82.21 79.4	55.35 57.85 59.6	29.55 31.24 33.4
Single Stage Detection	SSPS [59] SSPS [59]	val test-dev	32.4 32.6	- 74.3	42.0	34.8 35.0	- 74.8	- 44.8	28.6 29.0
Proposal-Free	AdaptIS [56] DeeperLab [63] DeeperLab [63] SSAP [16] SSAP [16]	val val test-dev val test-dev	35.9 33.8 34.3 36.5 36.9	- 77.1 80.7	- 43.1 44.8	40.3 37.5 40.1	- 77.5 81.6	- 46.8 - 48.5	29.3 29.6 32.0
	Ours (1x) Ours (1x)	val test-dev	37.51 37.7	77.65 77.8	47.18 47.3	40.01 40.7	78.39 78.7	50.03 50.7	33.74 33.1

Qualitative results





PCV

UPSNet





PCV

UPSNet

Conclusion

- Instances emerge from pixel consensus on centroid locations
- Voting as dilated deconvolution
- Back-projection as convolutional filtering
- Training reduced to pixel labeling

