

Instance segmentation

CV3DST | Prof. Leal-Taixé

Semantic segmentation

Label every pixel, including the background (sky, grass, road)



Do not differentiate between the pixels coming from instances of the same class

Instance segmentation

Label every pixel, including the background (sky, grass, road)



Do not differentiate between the pixels coming from instances of the same class

Do not label pixels coming from uncountable objects (sky, grass, road)



Differentiate between the pixels coming from instances of the same class

Instance segmentation methods

FCN-based

Proposal-based



Instance segmentation methods

FCN-based





FCN-based methods



A semantic map...

We already know how to obtain this!

Why FCN-based?

• Fully Convolutional Networks for Semantic Segmentation



Long, Shelhamer, Darrell - Fully Convolutional Networks for Semantic Segmentation, CVPR 2015, PAMI 2016

FCN-based methods

• X. Liang et al. "Proposal-free Network for Instancelevel Object Segmentation". Arxiv 2015

• A. Kirillov et al. "InstanceCut: from Edges to Instances with MultiCut". CVPR 2017

• M. Bai and R. Urtasun "Deep Watershed Transform for Instance Segmentation ". CVPR 2017

Instances through clustering



CV3DST | Prof. Leal-Taixé

A. Kirillov et al. "InstanceCut: from Edges to Instances with MultiCut". CVPR 2017

Instance segmentation methods

FCN-based

Proposal-based



Proposal-based methods

Bounding boxes.....

We already know how to obtain those!



Proposal-based methods

- B. Hariharan et al. "Simultaneous Detection and Segmentation". ECCV 2014
 - Follow-up work: B. Hariharan et al. "Hypercolumns for Object Segmentation and Fine-grained Localization ". CVPR 2015
- Dai et al. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". CVPR 2016
 - Previous work: Dai et al. "Convolutional Feature Masking for Joint Object and Stuff Segmentation". CVPR 2015

SDS

SDS: Simultaneous Detection and Segmentation



MNC

• MNC: Multi-task network cascades







Mask R-CNN

What is Mask-RCNN?

• Starting from the Faster R-CNN architecture



What is Mask-RCNN?

• Faster R-CNN + FCN for segmentation



Region Proposal Network

What is Mask-RCNN?

• Faster R-CNN + FCN for segmentation



Mask loss = binary cross entropy per pixel for the k segmentation classes

He at al. "Mask R-CNN" ICCV 2017

Mask R-CNN



Detection vs. segmentation

• Detection: for object classification, you require invariant representations



Translation equivariance: wherever the penguin is in the image, I still want to have "penguin" as my classification output

Detection vs. segmentation

- Detection: for object classification, you require
 invariant representations
- Segmentation: you require equivariant representations
 - Translated object \rightarrow Translated mask
 - Scaled object \rightarrow scaled mask
 - For semantic segmentation, small objects are less important (less pixels), but for instance segmentation, all objects (no matter the size) are equally important

Mask-RCNN: operations

• What operations are equivariant?



Features extraction = convolutional layers → equivariant

Segmentation head is a fully convolutional network \rightarrow equivariant

Mask-RCNN: operations

RolAlign

class box

• What operations are equivariant?

Fully connected layers and global pooling layers give invariance!



Segmentation head is a fully convolutional network \rightarrow equivariant

Recall: Rol pooling

Region of Interest Pooling: for every proposal



Recall: Rol pooling

Not suitable to extract pixel-wise precise masks

• Let us look at sizes



Mask-RCNN: operations

Make all operations equivariant

Fully connected layers and global pooling layers give invariance!



RolAlign

• Erase quantization effects



Chose 48.75







CV3DST | Prof. Leal-Taixé



CV3DST | Prof. Leal-Taixé



Mask R-CNN: extended for joints



Model a keypoint's location as a one-hot mask, and adopt Mask R-CNN to predict K masks (which are in the end only 1 pixel), one for each of K keypoint types (e.g., left shoulder, right elbow). This demonstrates the flexibility of Mask R-CNN.

Improving Mask-RCNN

• One problem with Mask R-CNN is that the mask quality score is computed as the confidence score for the bounding box

The only way the "instance" is evaluated is through the box loss



Recall the mask loss just evaluates if the pixels have the correct semantic class, not the correct instance!

Both instances have the same class = person

Mask IoU head


Mask confidence score



Typically, Mask scoring R-CNN gives lower confidence scores than Mask R-CNN, which corresponds to masks not being perfect (IoU < 1.0).

This tiny modification achieves SOTA results.

πп

Is one-stage vs twostage also applicable to masks?

One-stage vs two-stage detectors

Faster R-CNN



S × S grid on input

YOLO

Slower, but has higher performance

Faster, but has lower performance

One-stage vs two-stage instance segmenters

Mask R-CNN

YOLACT





Slower, but has higher performance

Faster, but has lower performance

CV3DST | Prof. Leal-Taixé

YOLO with masks?

"Boxes are stupid anyway though, I'm probably a true believer in masks except I can't get YOLO to learn them."

– Joseph Redmon, YOLOv3



YOLACT*

*You Only Look At CoefficienTs

CV3DST | Prof. Leal-Taixé

42





1) Generate mask prototypes



1) Generate mask prototypes



1) Generate mask prototypes

YOLACT: backbone



YOLACT: protonet



YOLACT: protonet

• Fully convolutional network



Similar to the mask branch in Masr R-CNN.

However, no loss function is applied on this stage.

YOLACT: mask coefficients

Predict a coefficient for every predicted mask.



YOLACT: mask coefficients



The network is similar but shallower than RetinaNet

YOLACT: mask assembly

- 1. Do a linear combination between the mask coefficients and the mask prototypes.
- 2. Predict the mask as $M = \sigma(PC^T)$ where P is a (HxWxK) matrix of prototype masks, C is a (NxK) matrix of mask coefficients surviving NMS, and σ is a nonlinearity.



YOLACT: loss function



Cross-entropy between the assembled masks and the ground truth, in addition to the standard losses (regression for the bounding box, and classification for the class of the object/mask).

YOLACT: qualitative results



YOLACT: qualitative results



For large objects, the quality of the masks is even better than those of twostage detectors

So, which segmenter to use?



YOLACT: improvements

- A specially designed version of NMS, in order to make the procedure faster.
- An auxiliary semantic segmentation loss function performed on the final features of the FPN. The module is not used during the inference stage.
- D. Boyla et al. "YOLACT++: Better real-time instance segmentation". <u>arXiv:1912.06218</u> 2019



Semantic segmentation



Instance segmentation



+

Semantic segmentation



+

FCN-like

Instance segmentation



Mask R-CNN

Semantic segmentation



FCN-like

Instance segmentation



+

Mask R-CNN

Panoptic segmentation



=

UPSNet

CV3DST | Prof. Leal-Taixé



It gives labels to uncountable objects called "stuff" (sky, road, etc), similar to FCN-like netsorks.

It differentiates between pixels coming from different distances of the same class (countable objects) called "things" (cars, pedestrians, etc).



Problem: some pixels might get classified as stuff from FCN network, while at the same time being classified as instances of some class from Mask R-CNN (conflicting results)!



Solution: Parametric-free panoptic head which combines the information from the FCN and Mask R-CNN, giving final predictions.

Network architecture



Network architecture



CV3DST | Prof. Leal-Taixé

The semantic head



New: deformable convolutions!

Recall: Dilated (atrous) convolutions 2D



(a) the dilation parameter is 1, and each element produced by this filter has reception field of 3×3. (b) the dilation parameter is 2, and each element produced by it has reception field of 7x7. (c) the dilation parameter is 4, and each element produced by it has reception field of 15x15.

Deformable convolutions



(a) Conventional Convolution, (b) Deformable Convolution, (c) Special Case of Deformable Convolution with Scaling, (d) Special Case of Deformable Convolution with Rotation

Deformable convolutions: generalization of dilated convolutions when you learn the offset

Deformable convolutions



Regular convolution

$$\mathbf{y}(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{x}(\mathbf{p}_0 + \mathbf{p}_n)$$

Deformable convolution

$$\mathbf{y}(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{x}(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n)$$

where $\Delta \mathbf{p}_n$ is generated by a sibling branch of regular convolution

Deformable convolutions



The deformable convolution will pick the values at different locations for convolutions conditioned on the input image of the feature maps.

The Panoptic head

Mask logits from the instance head

Object logits coming from the semantic head (e.g., car)

Stuff logits coming from the semantic head (e.g., sky)


The Panoptic head

Mask logits from the instance head

Object logits coming from the semantic head (e.g., car)

Stuff logits coming from the semantic head (e.g., sky)



The Panoptic head



Perform softmax over the panoptic logits. If the maximum value falls into the first stuff channels, then it belongs to one of the stuff classes. Otherwise the index of the maximum value tells us the instance ID the pixel belongs to.

Read the details on how to use the unknown class (hint: think about the exam...)

Panoptic quality



 SQ: Segmentation Quality = how close the predicted segments are to the ground truth segment (does not take into account bad predictions!)

Panoptic quality



 RQ: Recognition Quality = just like for detection, we want to know if we are missing any instances (FN) or we are predicting more instances (FP).

Panoptic quality

• As in detection, we have to "match ground truth and predictions. In this case we have segment matching.



• Segment is matched if IoU>0.5. No pixel can belong to two predicted segments.

Panoptic segmentation: qualitative



CV3DST | Prof. Leal-Taixé

Panoptic segmentation: qualitative



Next lectures

Next Friday: Video object segmentation

• Remember to work on the competition!