

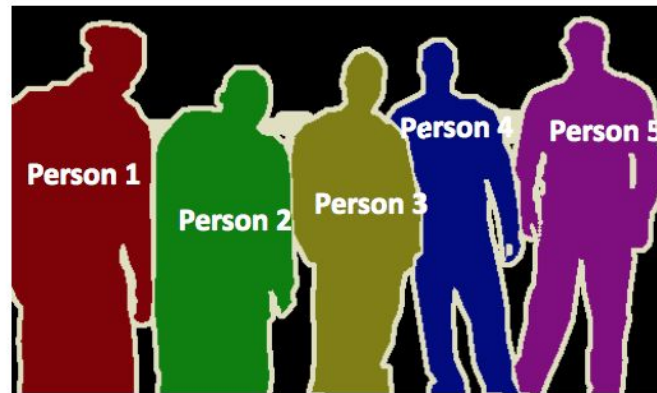


Instance Segmentation

Riley Simmons-Edler, Berthy Feng

Instance Segmentation Task

- Label each foreground pixel with object and instance
- Object detection + semantic segmentation



Instance Segmentation



In This Lecture...

- Microsoft COCO dataset
- Mask R-CNN (fully supervised)
- Mask^X R-CNN (partially supervised)

Microsoft COCO: Common Objects in Context

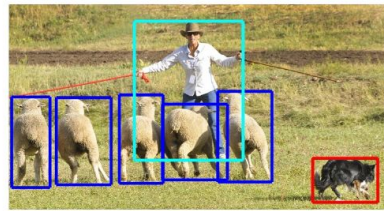
Tsung-Yi Lin, Michael Maire, Serge Belongie, et al.
“Microsoft COCO: Common Objects in Context.” arXiv,
2015.

Previous Datasets

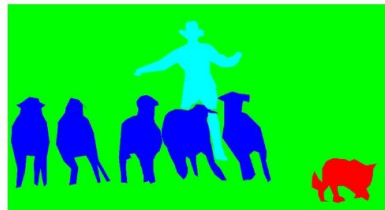
- **ImageNet**: many object categories
- **PASCAL VOC**: object detection in natural images, small number of classes
- **SUN**: labeling scene types and commonly occurring objects, but not many instances per category



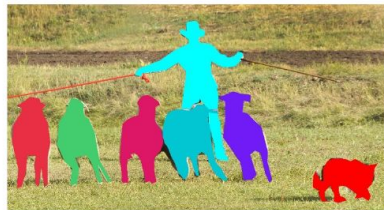
(a) Image classification



(b) Object localization



(c) Semantic segmentation



(d) This work



Goal: Push research in scene understanding

1. Detecting non-iconic views
2. Contextual reasoning between objects
3. Precise 2D localization of objects

MS COCO Dataset

Person



Dog



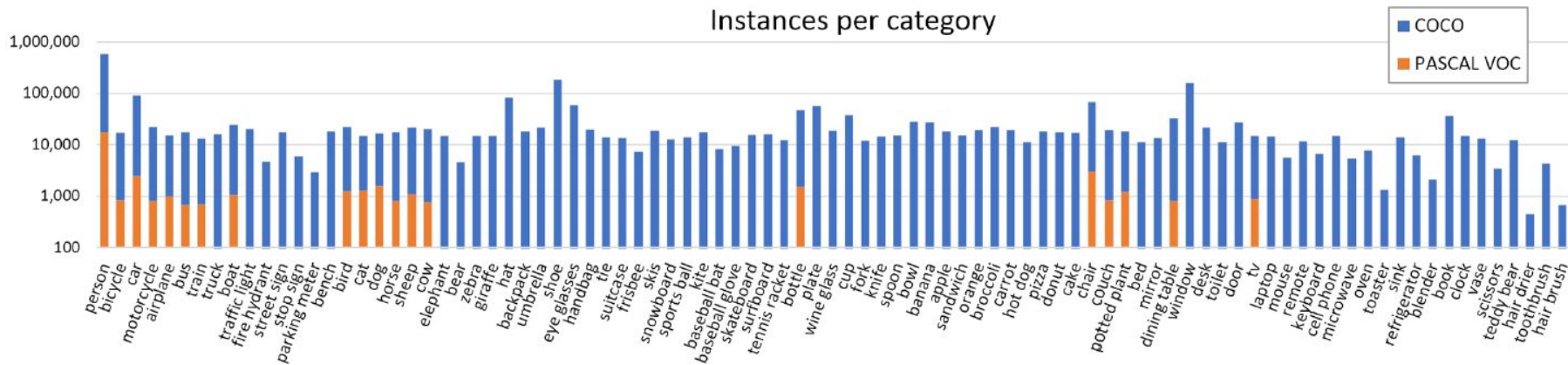
Cow



- ❖ 91 object classes
- ❖ 328,000 images
- ❖ 2.5 million labeled instances

Image Collection & Annotation

Object Categories



Non-Iconic Image Collection



(a) Iconic object images

(b) Iconic scene images

(c) Non-Iconic images

Annotation



(a) Category labeling



(b) Instance spotting



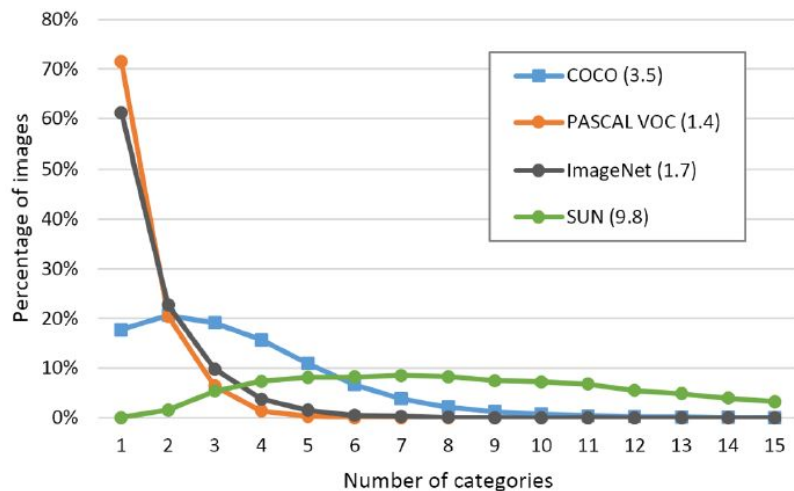
(c) Instance segmentation

Fig. 3: Our annotation pipeline is split into 3 primary tasks: (a) labeling the categories present in the image (§4.1), (b) locating and marking all instances of the labeled categories (§4.2), and (c) segmenting each object instance (§4.3).

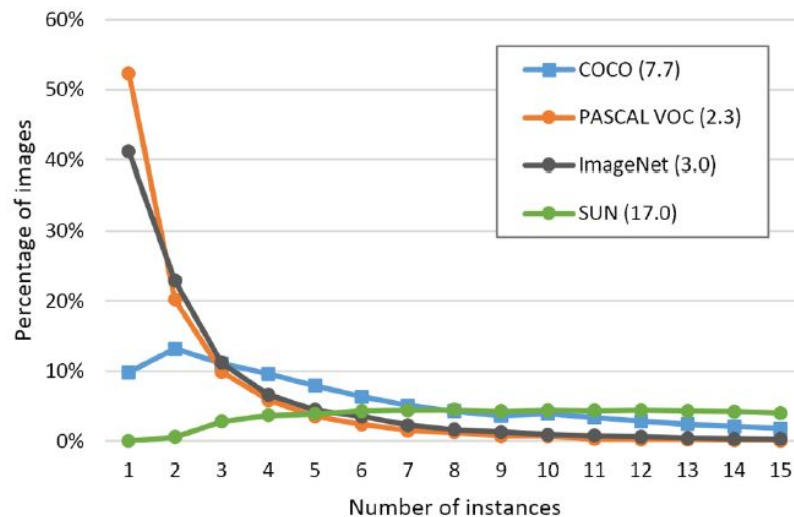
Dataset Evaluation

Statistics

Categories per image

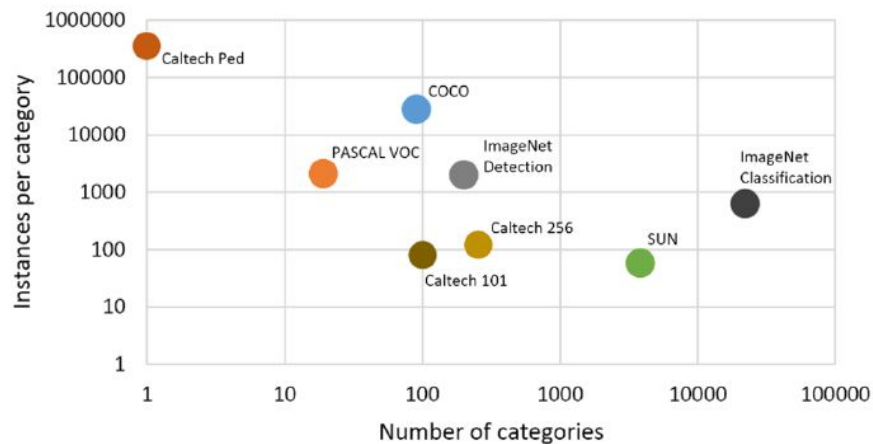


Instances per image

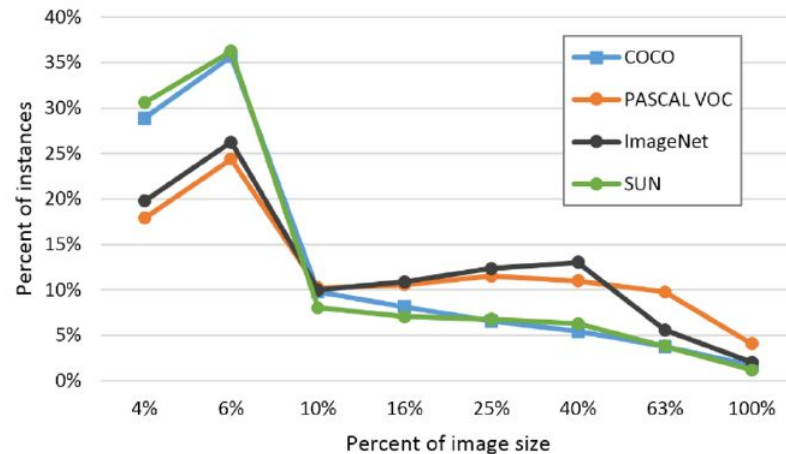


Statistics

Number of categories vs. number of instances



Instance size



COCO Detection Challenge



The COCO 2017 Detection Challenge is designed to push the state of the art in object detection forward. Teams are encouraged to compete in either (or both) of two object detection challenges: using bounding box output or object segmentation output. For full details of this task please see the [COCO Detection Challenge](#) page.

COCO Keypoint Challenge



The COCO 2017 Keypoint Challenge requires localization of person keypoints in challenging, uncontrolled conditions. The keypoint challenge involves simultaneously detecting people *and* localizing their keypoints (person locations are *not* given at test time). For full details of this task please see the [COCO Keypoints Challenge](#) page.

COCO Stuff Challenge



The COCO 2017 Stuff Segmentation Challenge is designed to push the state of the art in semantic segmentation of *stuff* classes. Whereas the COCO 2017 Detection Challenge addresses *thing* classes (person, car, elephant), this challenge focuses on *stuff* classes (grass, wall, sky). For full details of this task please see the [COCO Stuff Challenge](#) page.

COCO Places Challenges



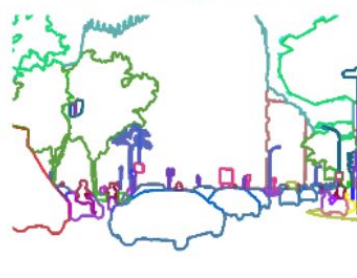
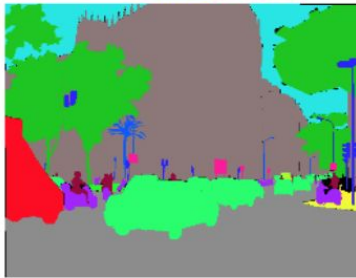
Scene Parsing



Instance Segmentation



Semantic Boundary Detection

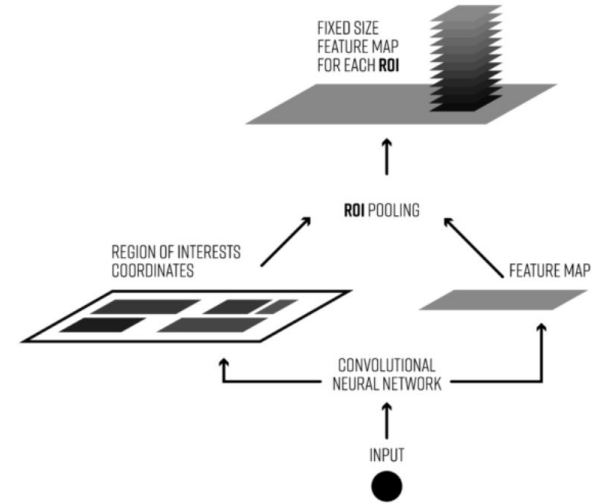
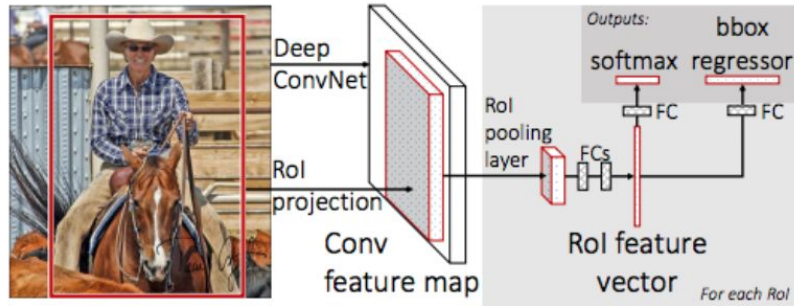


Mask R-CNN

Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. “Mask R-CNN.” ICCV, 2017.

Faster R-CNN

Fast R-CNN



Insight: Region Proposal and Detection Use Same Features

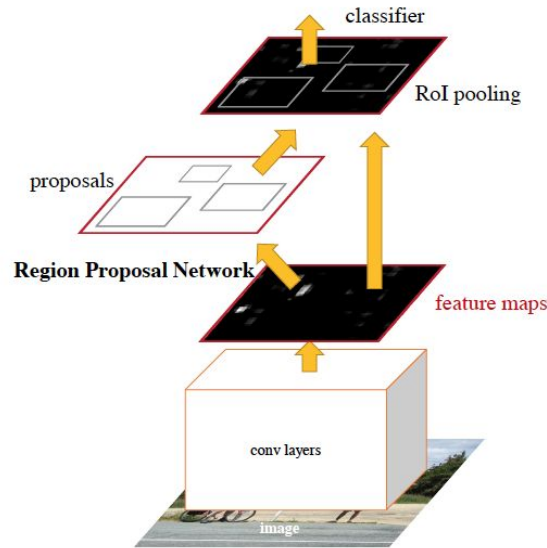
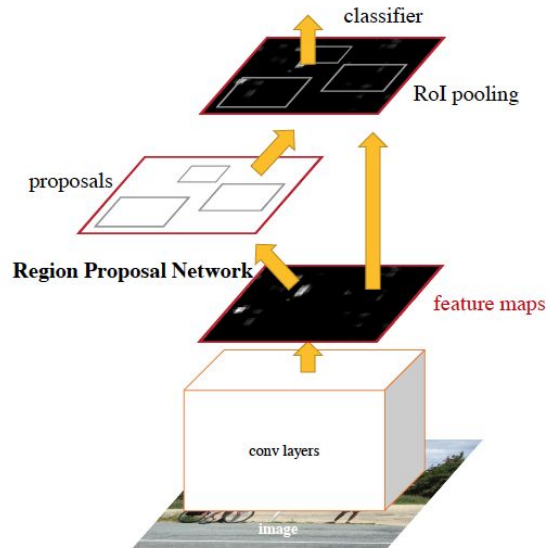
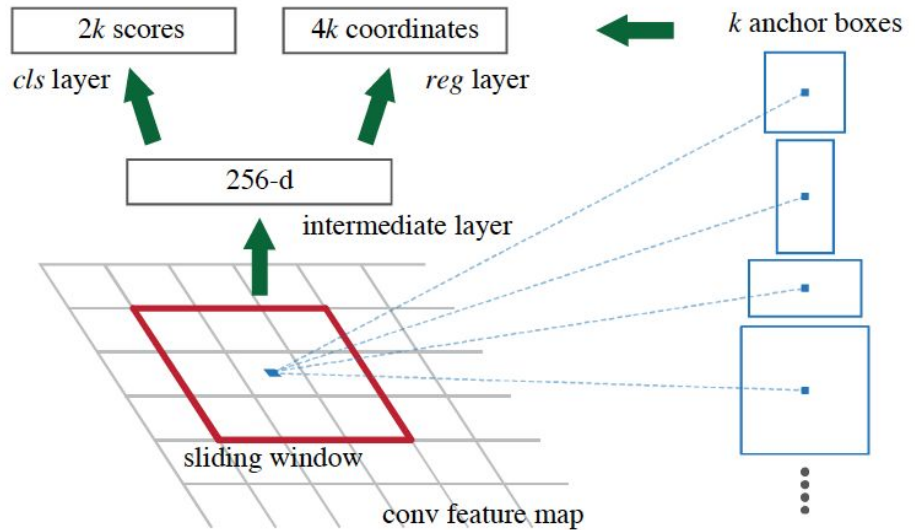


Figure 2: Faster R-CNN is a single, unified network for object detection. The RPN module serves as the 'attention' of this unified network.

Faster R-CNN = RPN + Fast R-CNN



RPN = Fully Convolutional Network

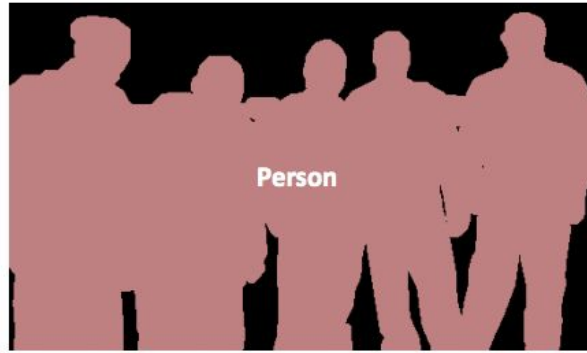


Extending to Instance Segmentation

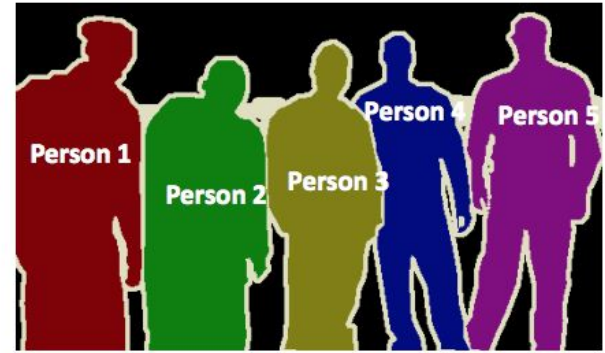
Visual Perception Problems



Object Detection



Semantic Segmentation

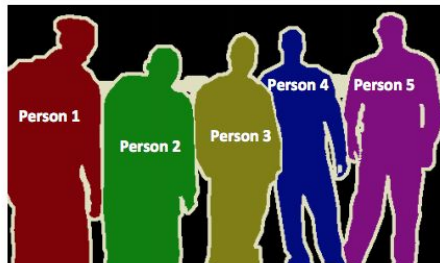


Instance Segmentation

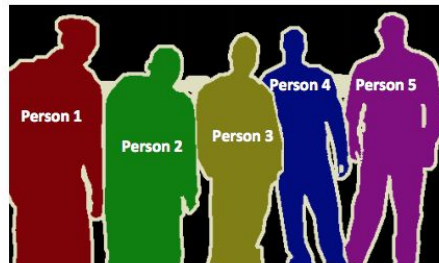


Instance Segmentation Methods

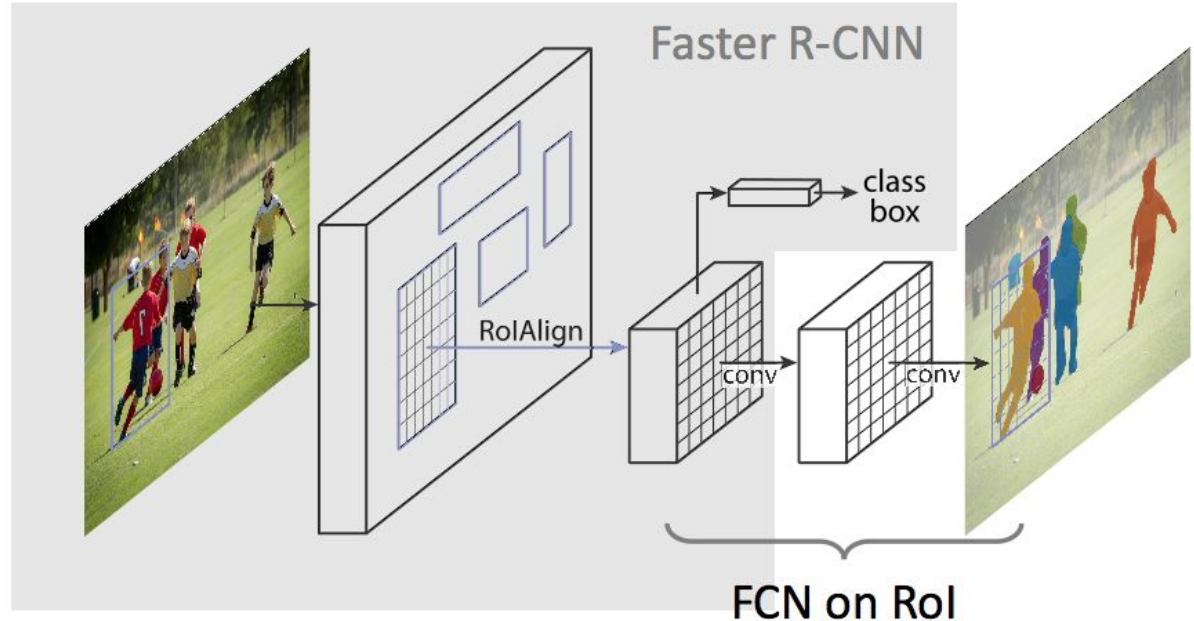
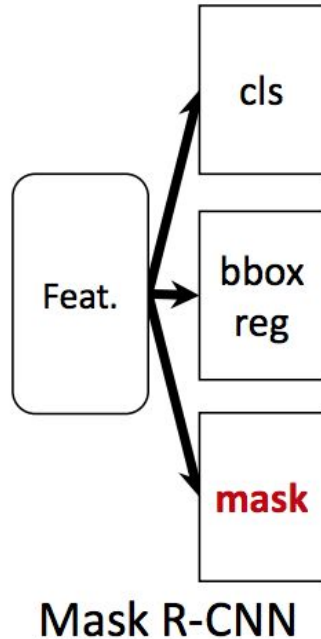
R-CNN driven



FCN driven



Insight: Mask Prediction in Parallel



RoIPool

input

0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

region proposal

0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

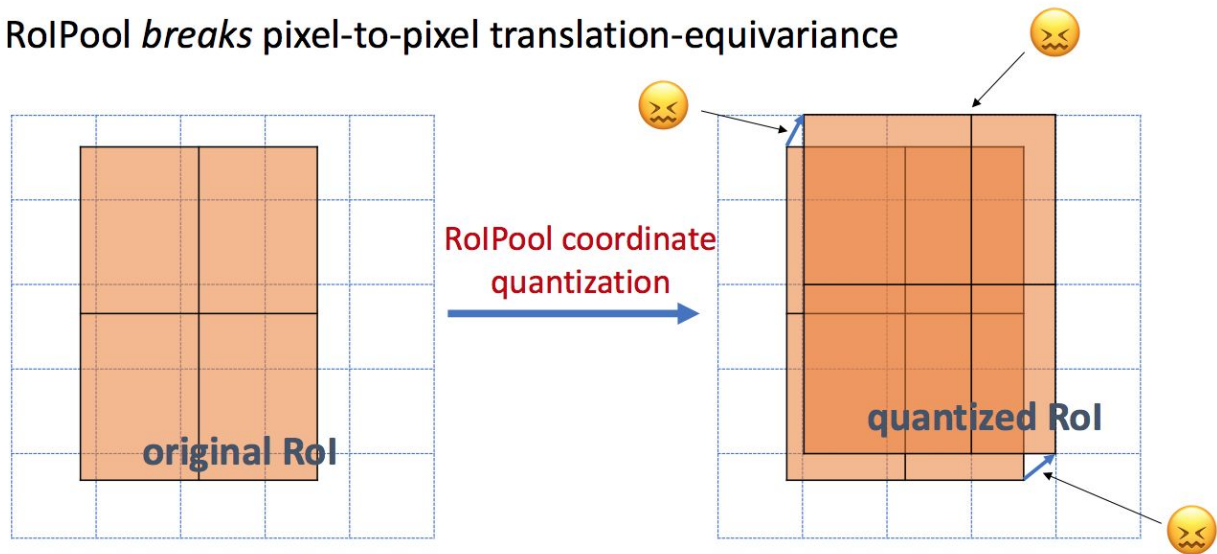
pooling sections

0.88	0.44	0.14	0.16	0.37	0.77	0.96	0.27
0.19	0.45	0.57	0.16	0.63	0.29	0.71	0.70
0.66	0.26	0.82	0.64	0.54	0.73	0.59	0.26
0.85	0.34	0.76	0.84	0.29	0.75	0.62	0.25
0.32	0.74	0.21	0.39	0.34	0.03	0.33	0.48
0.20	0.14	0.16	0.13	0.73	0.65	0.96	0.32
0.19	0.69	0.09	0.86	0.88	0.07	0.01	0.48
0.83	0.24	0.97	0.04	0.24	0.35	0.50	0.91

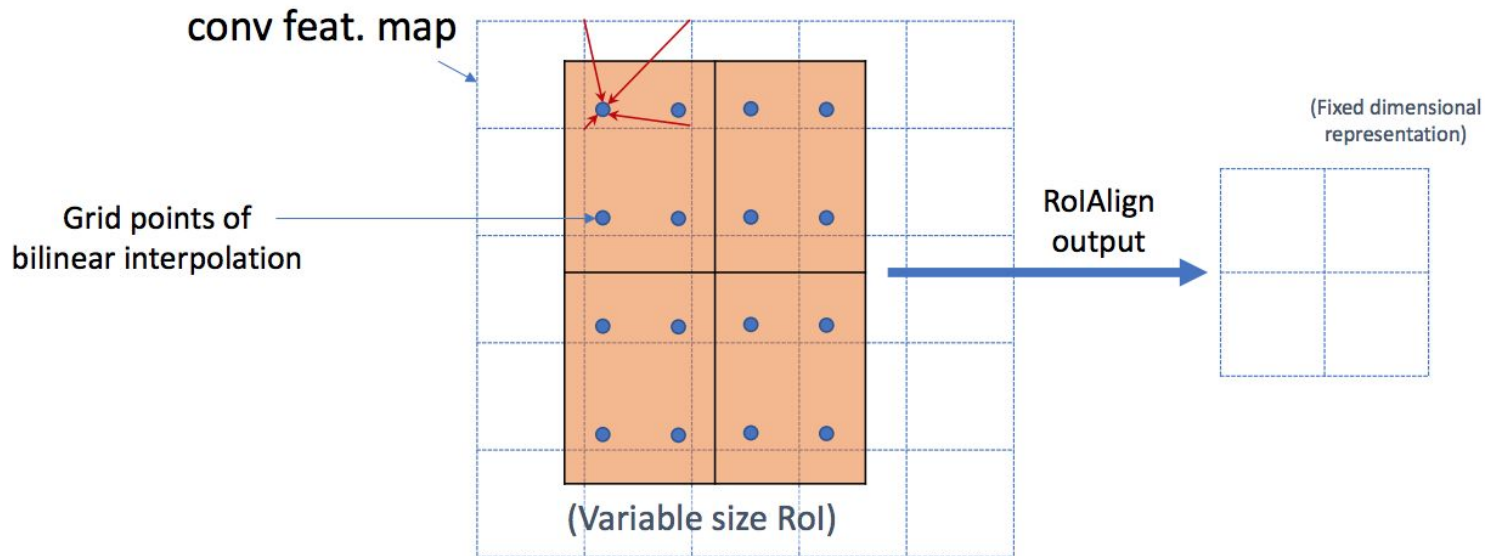
0.85	0.84
0.97	0.96

RoIPool

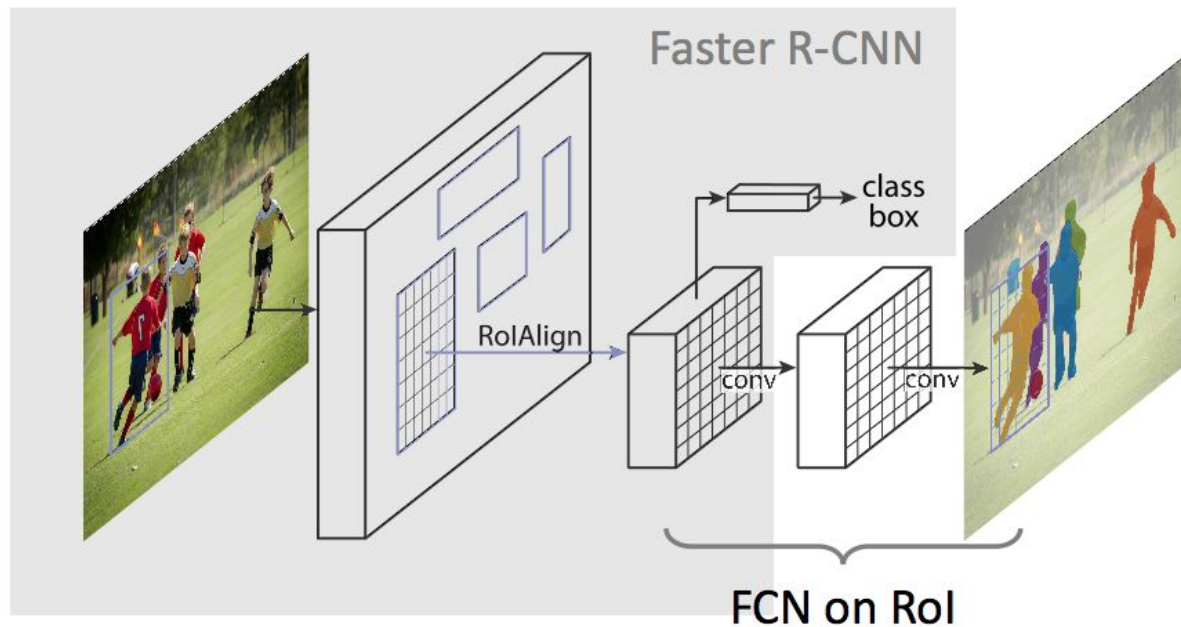
- RoIPool *breaks* pixel-to-pixel translation-equivariance



RoIAlign



Mask R-CNN



Mask R-CNN Results

Examples

- Mask AP = 35.7

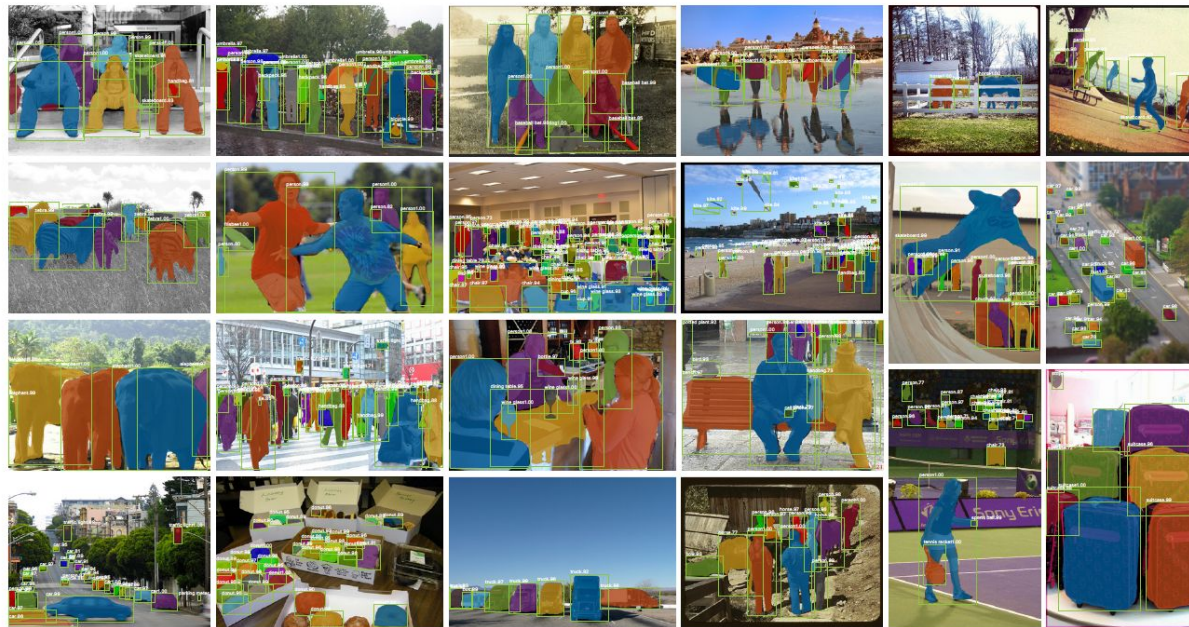


Figure 5. More results of **Mask R-CNN** on COCO test images, using ResNet-101-FPN and running at 5 fps, with 35.7 mask AP (Table 1).



Comparisons

	backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
MNC [10]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [26] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [26] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
Mask R-CNN	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
Mask R-CNN	ResNet-101-FPN	35.7	58.0	37.8	15.5	38.1	52.4
Mask R-CNN	ResNeXt-101-FPN	37.1	60.0	39.4	16.9	39.9	53.5

Comparisons

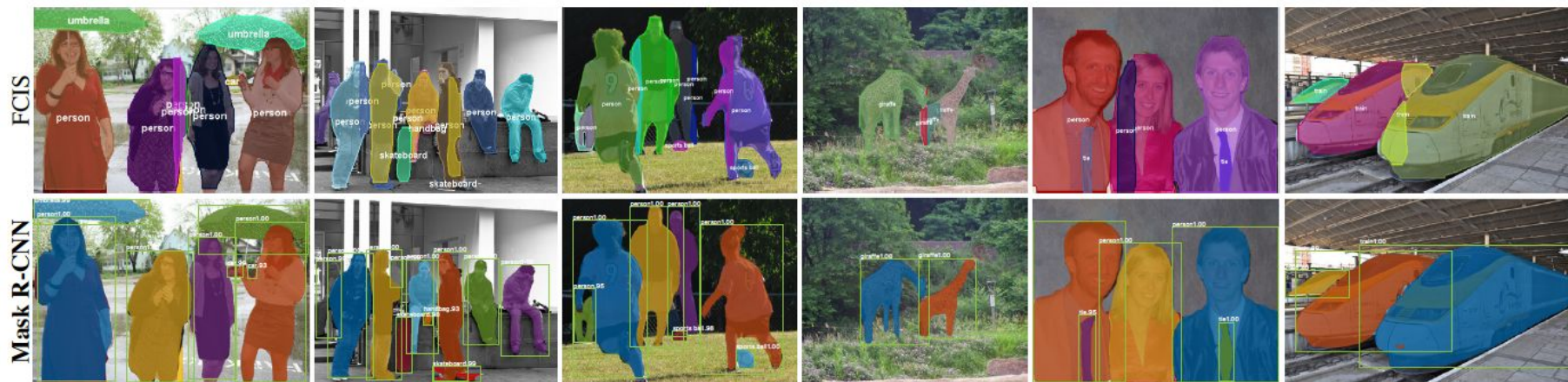


Figure 6. FCIS+++ [26] (top) vs. Mask R-CNN (bottom, ResNet-101-FPN). FCIS exhibits systematic artifacts on overlapping objects.

Application: Human Pose Estimation



Figure 7. Keypoint detection results on COCO test using Mask R-CNN (ResNet-50-FPN), with person segmentation masks predicted from the same model. This model has a keypoint AP of 63.1 and runs at 5 fps.



Mask R-CNN Recap

- Add parallel mask prediction head to Faster-RCNN
- RoIAlign allows for precise localization
- Mask R-CNN improves on AP of previous state-of-the-art, can be applied in human pose estimation

Learning to Segment Every Thing

Ronghang Hu, Piotr Dollar, Kaiming He, Trevor Darrell, and Ross Girshick. “Learning to Segment Every Thing.” arXiv, 2017.

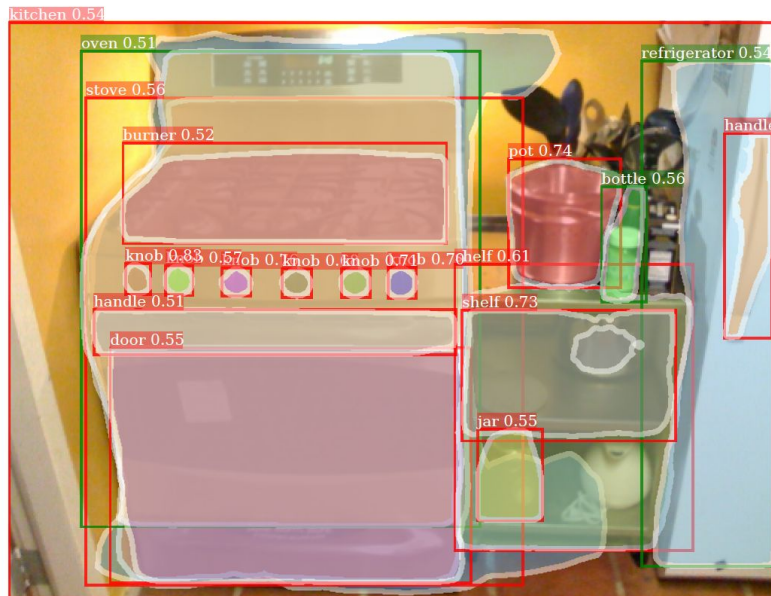
Partially Supervised Model

Motivation for a Partially Supervised Model

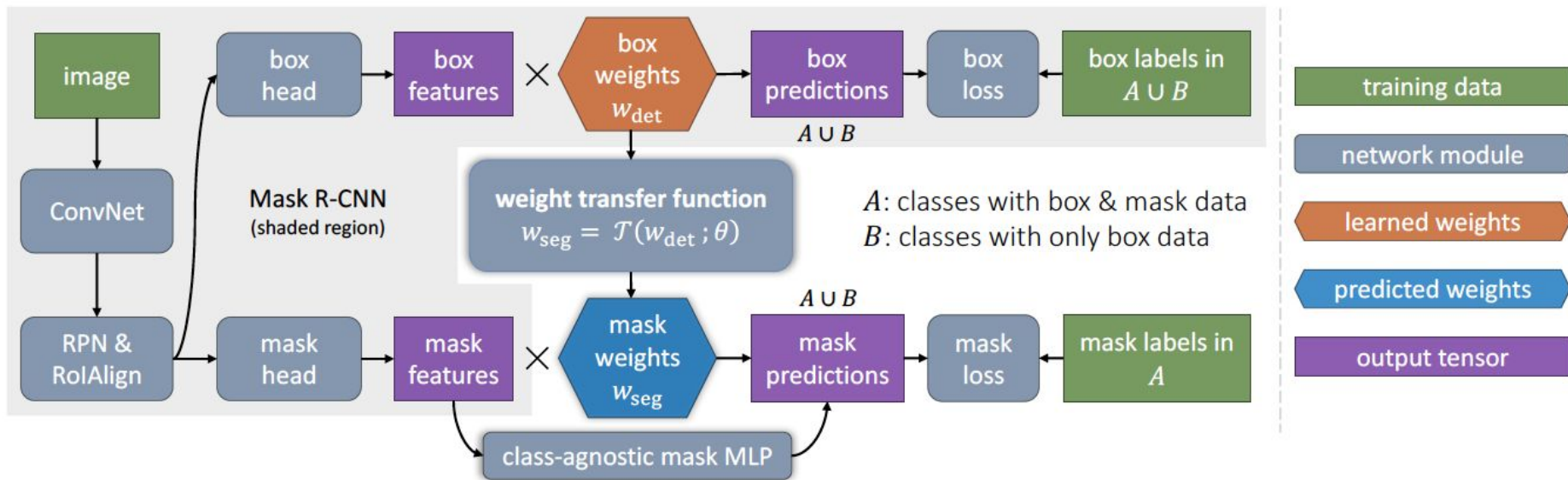
A = set of object categories with complete mask annotations

B = set of object categories with only bounding boxes (no segmentation annotations)

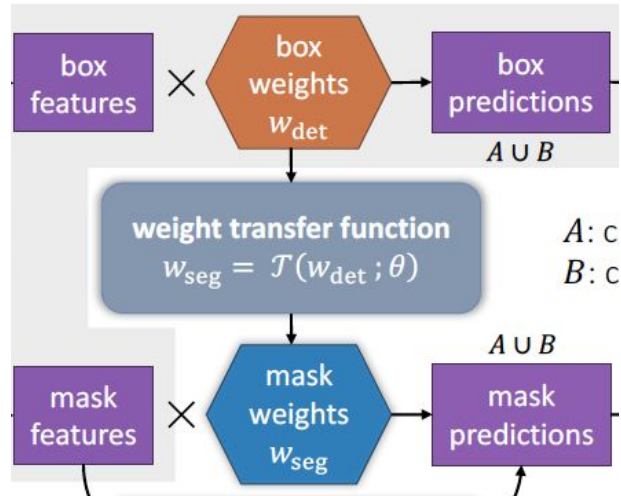
How can we know $C = A \cup B$?



Transfer Learning



Weight Transfer Function



$$w_{\text{seg}}^c = \mathcal{T}(w_{\text{det}}^c; \theta)$$

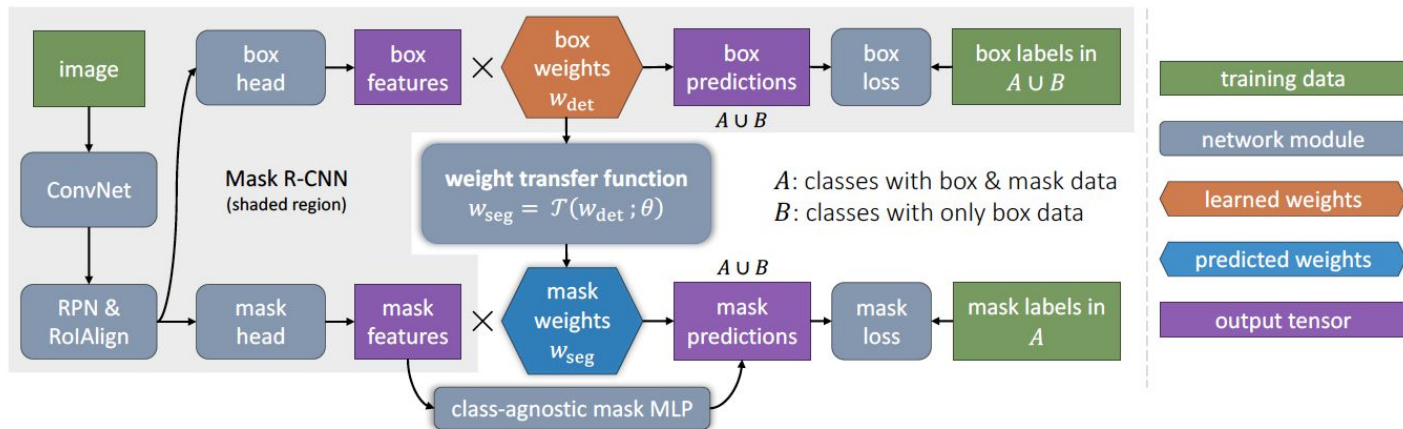
$$w_{\text{det}}^c = w_{\text{box}}^c$$

$$w_{\text{det}}^c = w_{\text{cls}}^c$$

$$w_{\text{det}}^c = [w_{\text{cls}}^c, w_{\text{box}}^c]$$

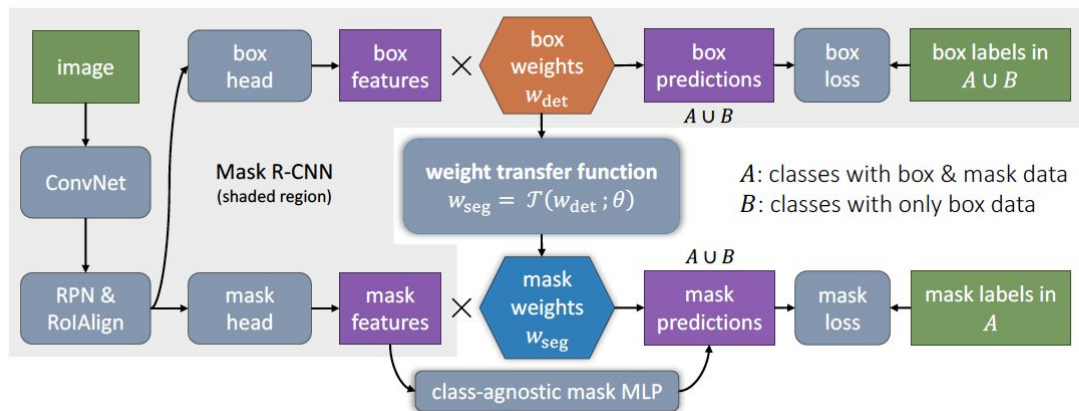
Training

- Train bounding box head using standard box detection losses on all classes in $A \cup B$
- Train mask head, weight transfer function using mask loss on classes in A



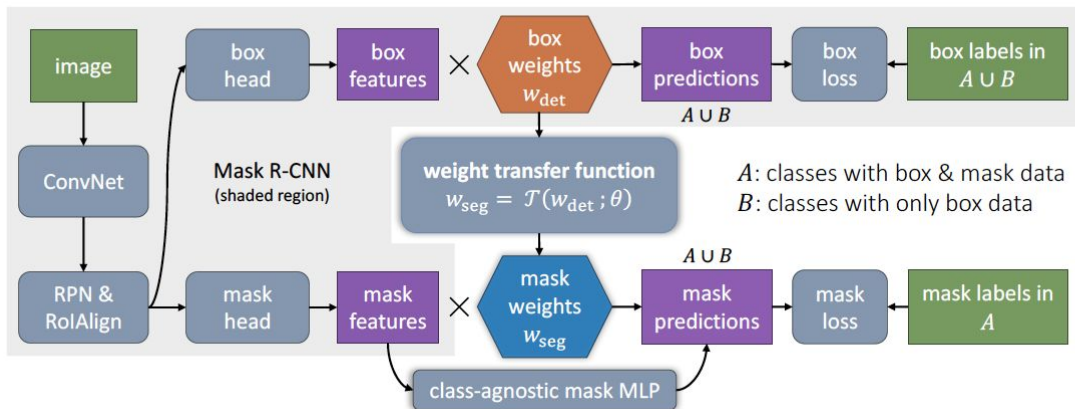
Stage-Wise Training

1. Detection training
 2. Segmentation training
- Train detection once and then fine-tune weight transfer function
 - Inferior performance



End-to-End Joint Training

- Jointly train detection head and mask head end-to-end
- Want detection weights to stay constant between A and B



$$w_{seg}^c = \mathcal{T}(\text{stop_grad}(w_{det}^c); \theta)$$

End-to-End Training Better

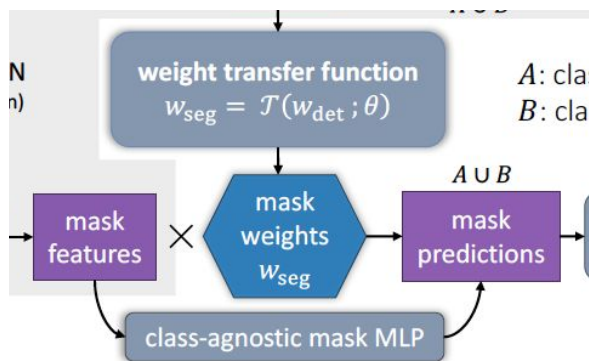
method	training	stop grad	voc \rightarrow non-voc		non-voc \rightarrow voc	
		on w_{det}	AP on B	AP on A	AP on B	AP on A
class-agnostic	sw	n/a	14.2	34.4	21.5	30.7
transfer	sw	n/a	20.2	35.2	26.0	31.2
class-agnostic	e2e	n/a	19.2	36.8	23.9	32.5
transfer	e2e		20.2	37.7	24.8	33.2
transfer	e2e	✓	22.2	37.6	27.6	33.1

(d) **Ablation on the training strategy.** We try both stage-wise (‘sw’) and end-to-end (‘e2e’) training (see §3.2), and whether to stop gradient from \mathcal{T} to w_{det} . End-to-end training improves the results and it is crucial to stop gradient on w_{det} .

Mask Prediction

Baseline: Class-agnostic FCN mask prediction

Extension: FCN+MLP mask head:



method	voc → non-voc		non-voc → voc	
	AP on B	AP on A	AP on B	AP on A
class-agnostic	14.2	34.4	21.5	30.7
class-agnostic+MLP	17.1	35.1	22.8	31.3
transfer	20.2	35.2	26.0	31.2
transfer+MLP	21.3	35.4	26.6	31.4

(c) **Impact of the MLP mask branch.** Adding the class-agnostic MLP mask branch (see §3.4) improves the performance of classes in set B for both the class-agnostic baseline and our weight transfer approach.

Results

Examples

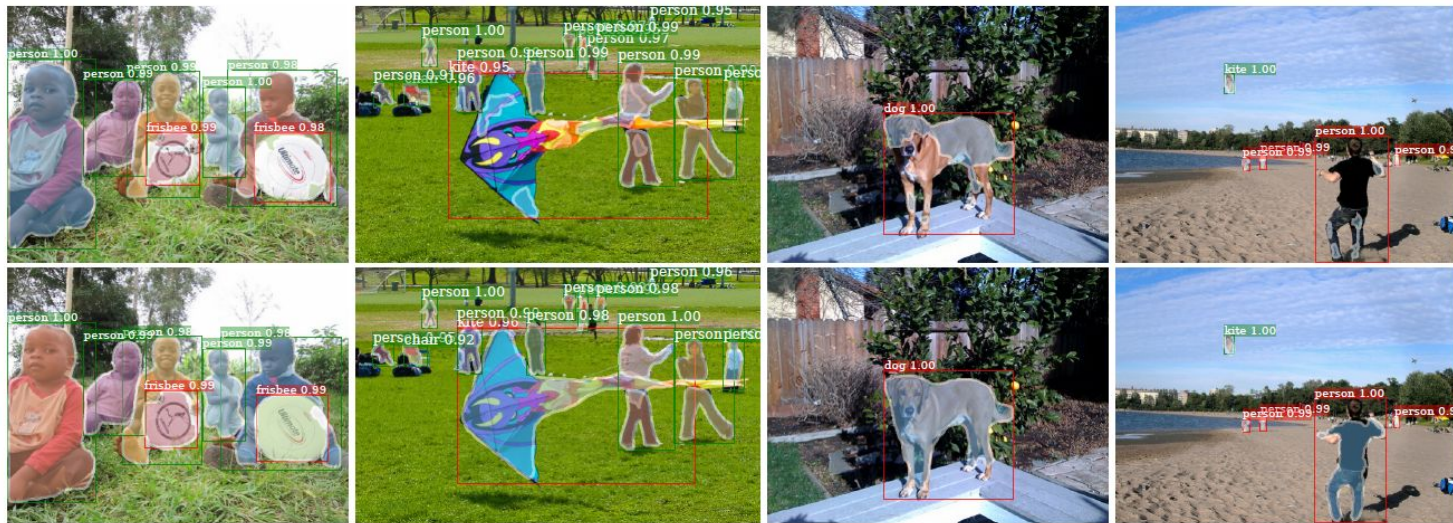


Figure 4. Mask predictions from the class-agnostic baseline (top row) vs. our Mask^X R-CNN approach (bottom row). Green boxes are classes in set A while the red boxes are classes in set B . The left 2 columns are $A = \{\text{voc}\}$ and the right 2 columns are $A = \{\text{non-voc}\}$.

Comparisons

backbone	method	voc \rightarrow non-voc: test on $B = \{\text{non-voc}\}$						non-voc \rightarrow voc: test on $B = \{\text{voc}\}$					
		AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
R-50-FPN	class-agnostic (baseline)	19.2	36.4	18.4	11.5	23.3	24.4	23.9	42.9	23.5	11.6	24.3	33.7
	Mask ^X R-CNN (ours)	23.7	43.1	23.5	12.4	27.6	32.9	28.9	52.2	28.6	12.1	29.0	40.6
	fully supervised (oracle)	33.0	53.7	35.0	15.1	37.0	49.9	37.5	63.1	38.9	15.1	36.0	53.1
R-101-FPN	class-agnostic (baseline)	18.5	34.8	18.1	11.3	23.4	21.7	24.7	43.5	24.9	11.4	25.7	35.1
	Mask ^X R-CNN (ours)	23.8	42.9	23.5	12.7	28.1	33.5	29.5	52.4	29.7	13.4	30.2	41.0
	fully supervised (oracle)	34.4	55.2	36.3	15.5	39.0	52.6	39.1	64.5	41.4	16.3	38.1	55.1

Table 2. **End-to-end training of Mask^X R-CNN.** As in Table 1, we use ‘cls+box, 2-layer, LeakyReLU’ implementation of \mathcal{T} and add the MLP mask branch (‘transfer+MLP’), and follow the same evaluation protocol. We also report AP₅₀ and AP₇₅ (average precision evaluated at 0.5 and 0.75 IoU threshold respectively), and AP over small (AP_S), medium (AP_M), and large (AP_L) objects. Our method significantly outperforms the baseline on those classes in set B without mask training data for both ResNet-50-FPN and ResNet-101-FPN backbones.

Segmenting Everything

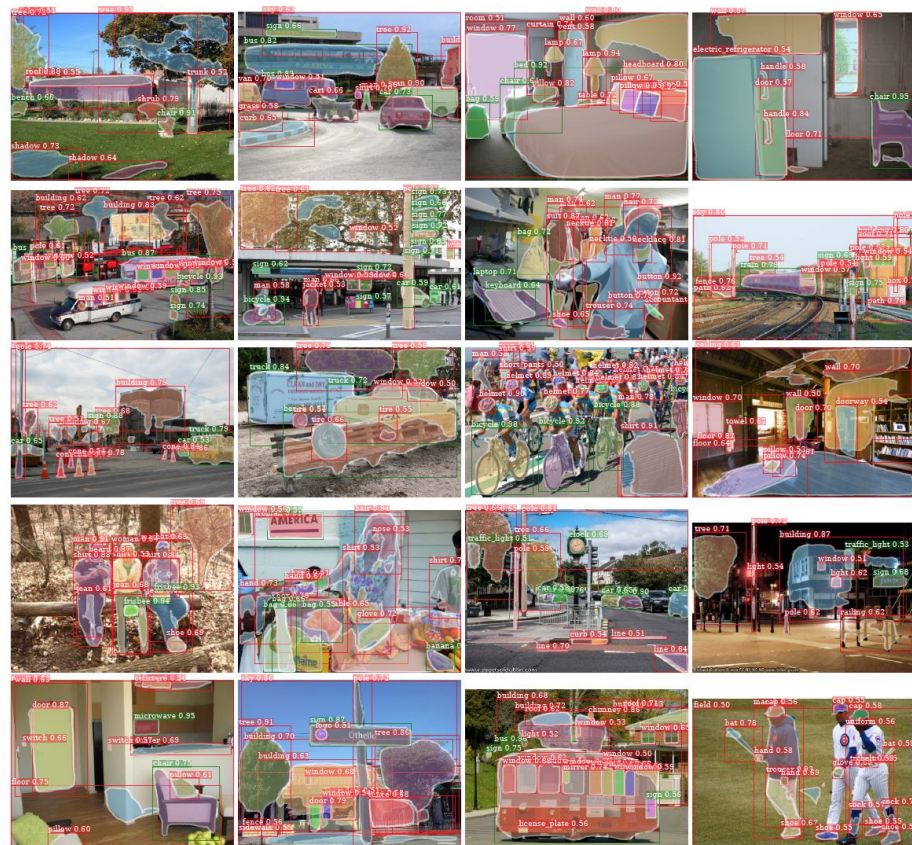


Figure 5. Example mask predictions from our Mask^X R-CNN on 3000 classes in Visual Genome. The green boxes are the 80 classes that overlap with COCO (set *A* with mask training data) while the red boxes are the remaining 2920 classes not in COCO (set *B* without mask training data). It can be seen that our model generates reasonable mask predictions on many classes in set *B*. See §5 for details.