

Goal

Incorporate top-down information, feedback and contextual information in Faster R-CNN

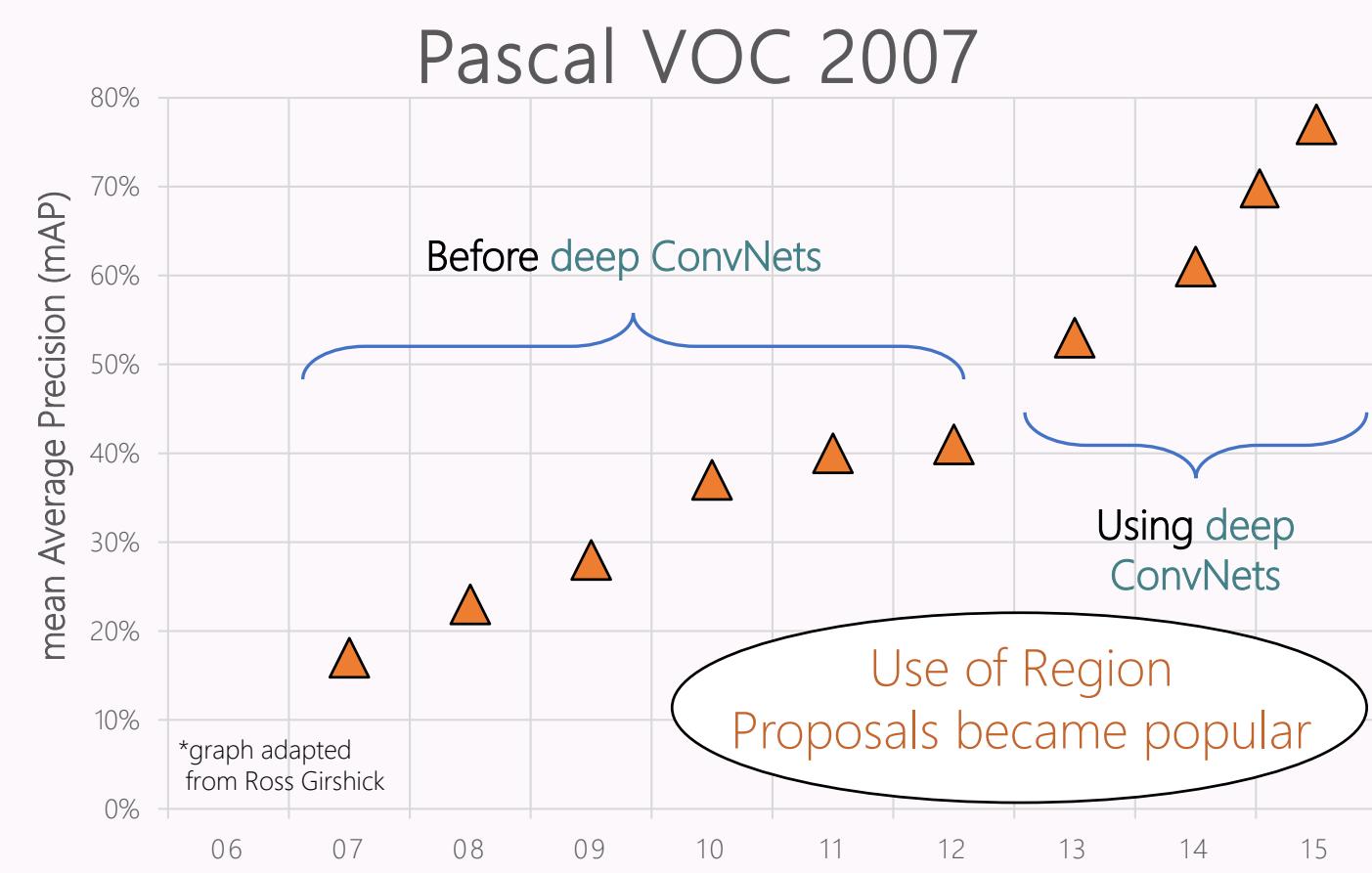
Contribution

Using **Semantic segmentation** for **contextually priming** region proposal & object detection modules, and providing **iterative feedback** to the entire network

Results

Improvement across all three tasks: object detection, semantic segmentation and region proposals.

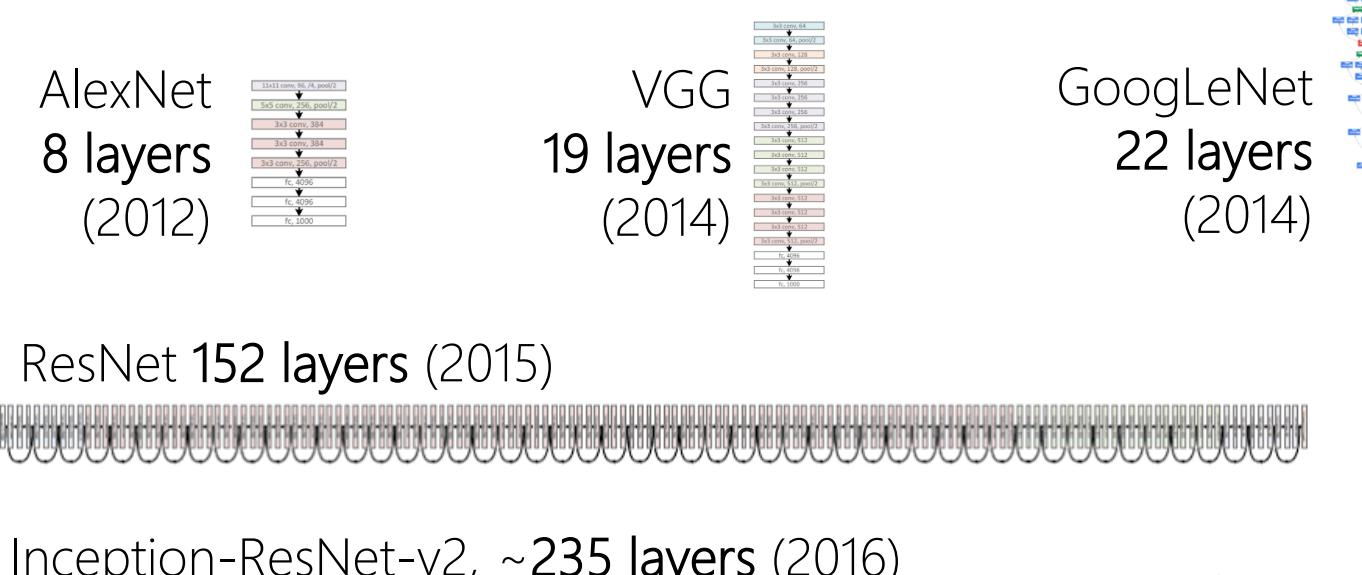
Key Ingredients of a Region-based ConvNet Object Detector [most state-of-the-art in Object Detection systems]



1

Deeper, Feedforward ConvNets

Deeper Network = Better Performance (so far.)



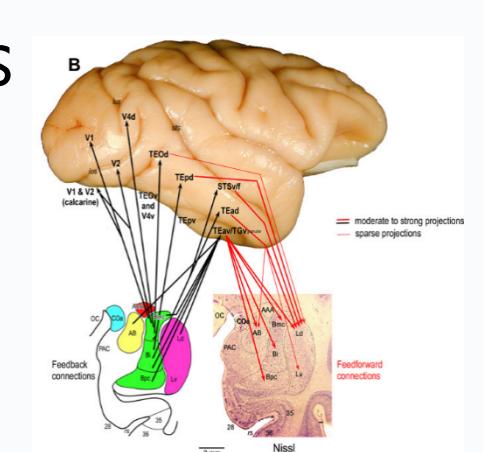
Networks getting deeper, but remain **feedforward**

Are we on the right path?

Human Visual Pathway

Strong evidence of **Feedback connections**

- Outnumber feedforward
- Feedback even to V1



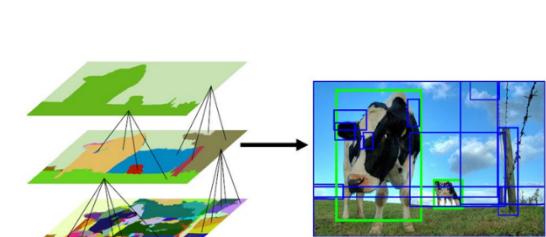
[Hupe et al., 1998], [Kravitz et al., 2013] etc.

2

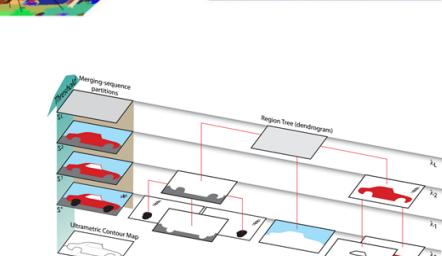
Recognition using Regions

E.g., Selective Search, Randomized Prim's, CPMC, Bing, EdgeBoxes, Rigor, Geodesic, MCG, DeepMask, SharpMask, AttracNet, etc.

Reduces Search Space
Allows use of richer features



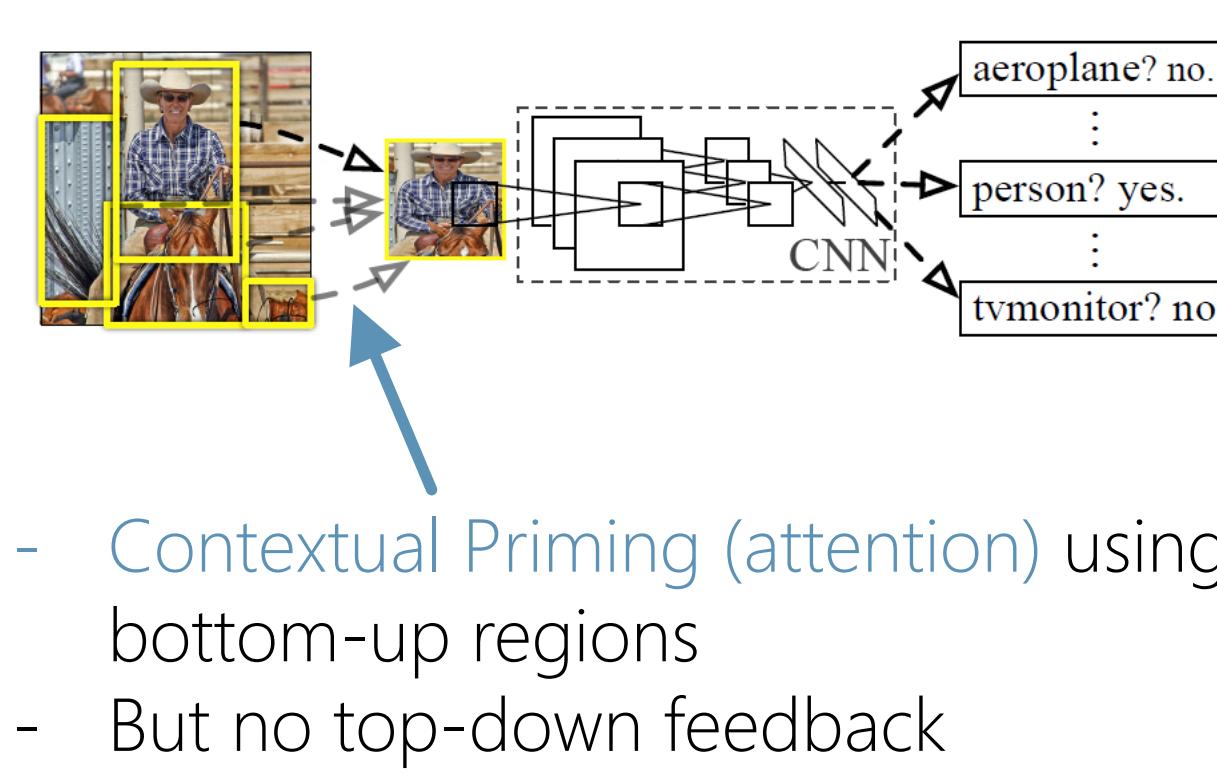
Focuses 'attention' in right areas
Reduces false positives



Generally, bottom-up, segmentation driven

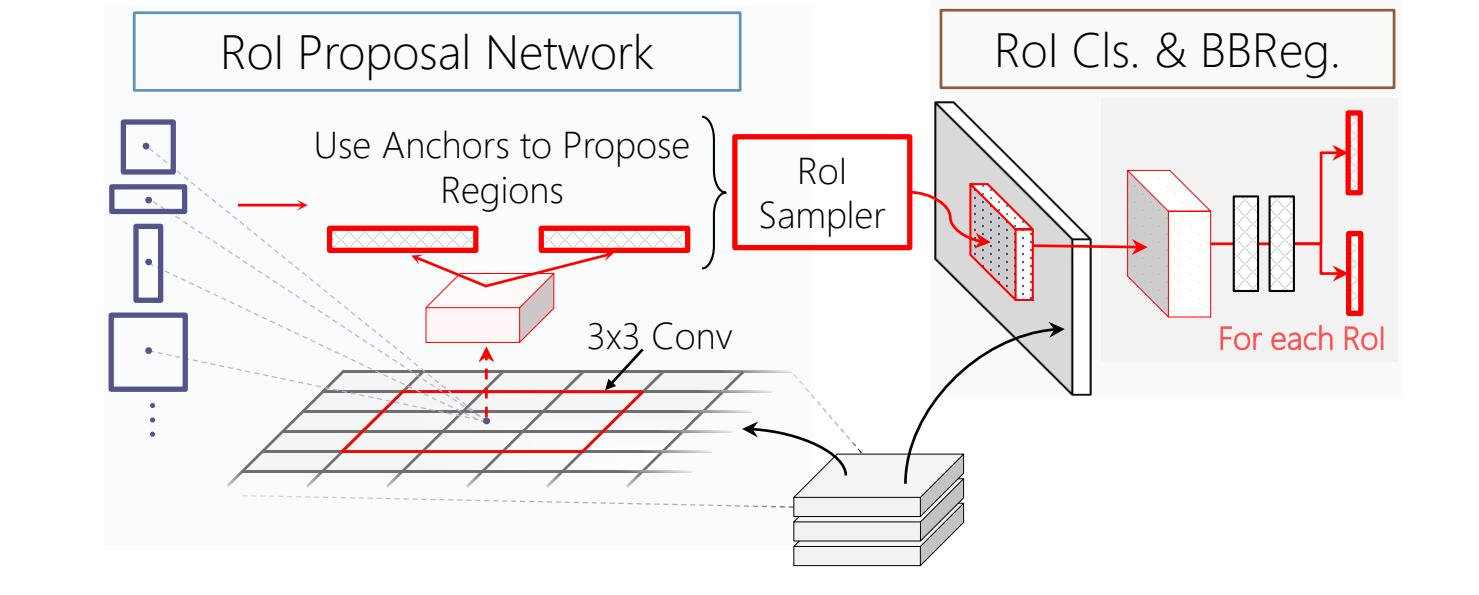
From Fast R-CNN to Faster R-CNN

Fast R-CNN



- Contextual Priming (attention) using bottom-up regions
- But no top-down feedback

Faster R-CNN



- Learns ConvNet to propose regions
- No Segmentation driven or bottom-up regions

Can we bridge this gap between empirical results and theory?

Incorporate top-down information, feedback and/or contextual reasoning in object detection

Contextual Priming and Feedback: Incorporating top-down information Faster R-CNN

Main Contributions:

Semantic segmentation as a top-down signal for:

- **Contextual Priming**
For region proposals & object detection
- **Iterative Feedback**
Top-down feedback to the entire network

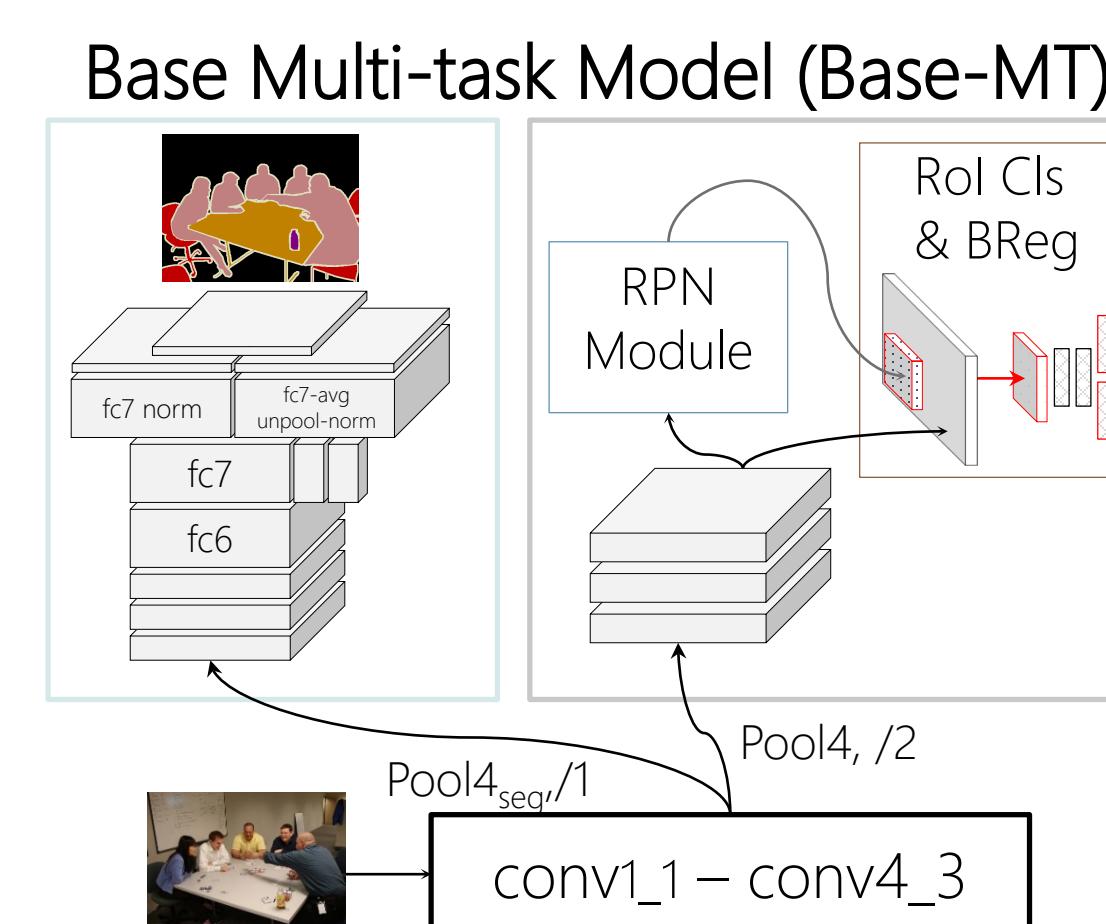
0

Faster R-CNN + Segmentation

Ideal Segmentation Network:

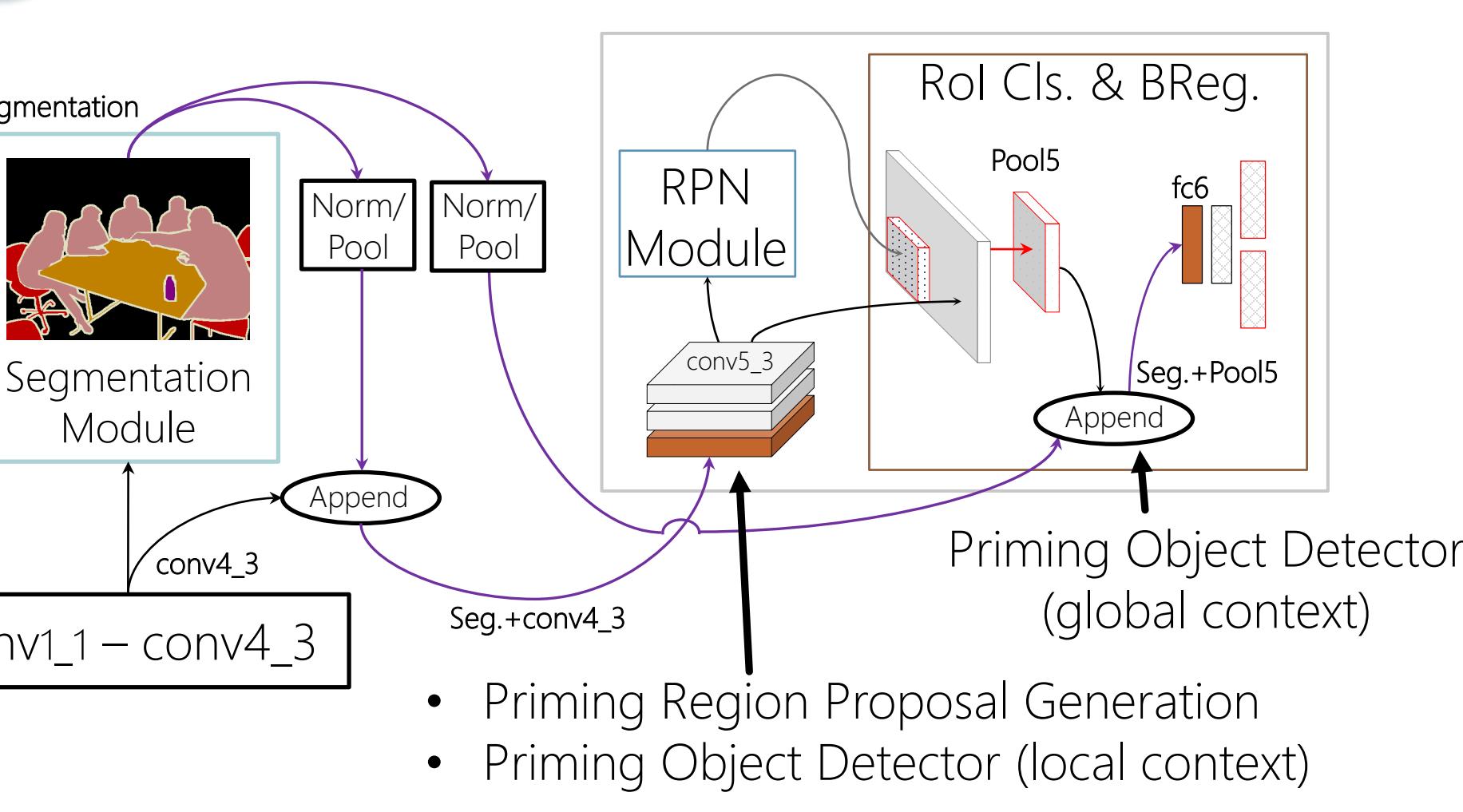
- Should be Fast
- Closely follow Faster R-CNN network (e.g., VGG16)
- No post-processing (e.g., CRFs)
- Helps with end-to-end training

We use ParseNet [Liu 2015].



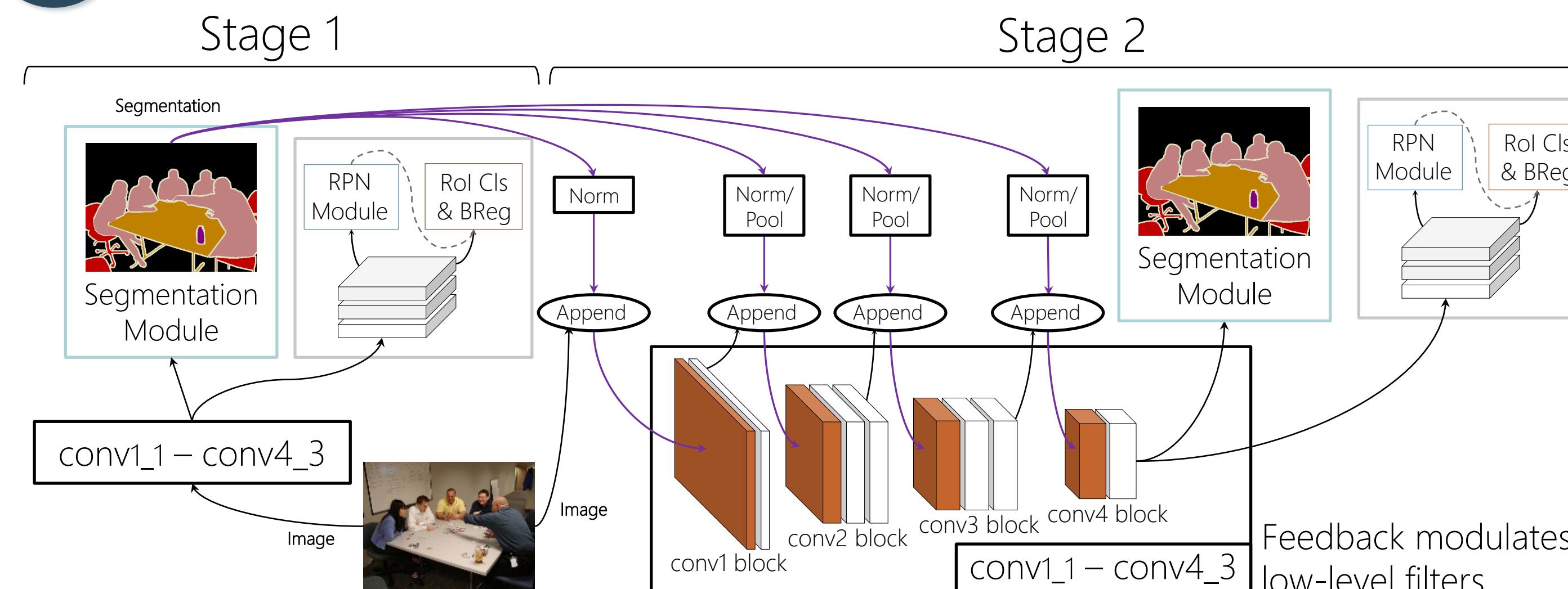
1

Contextual Priming via Segmentation



- Priming Region Proposal Generation
- Priming Object Detector (local context)

2 Iterative Feedback via Segmentation



Experiments to study the impact of Priming & Feedback

Ablation Analysis: Contextual Priming

	mAP	mIOU
Base-MT	75.6	65.8
Priming to conv5_1	77.0	65.8
Priming to conv5_1, each fc6	77.8	65.3

+ Priming to each RoI (which adds global context) helps detection.

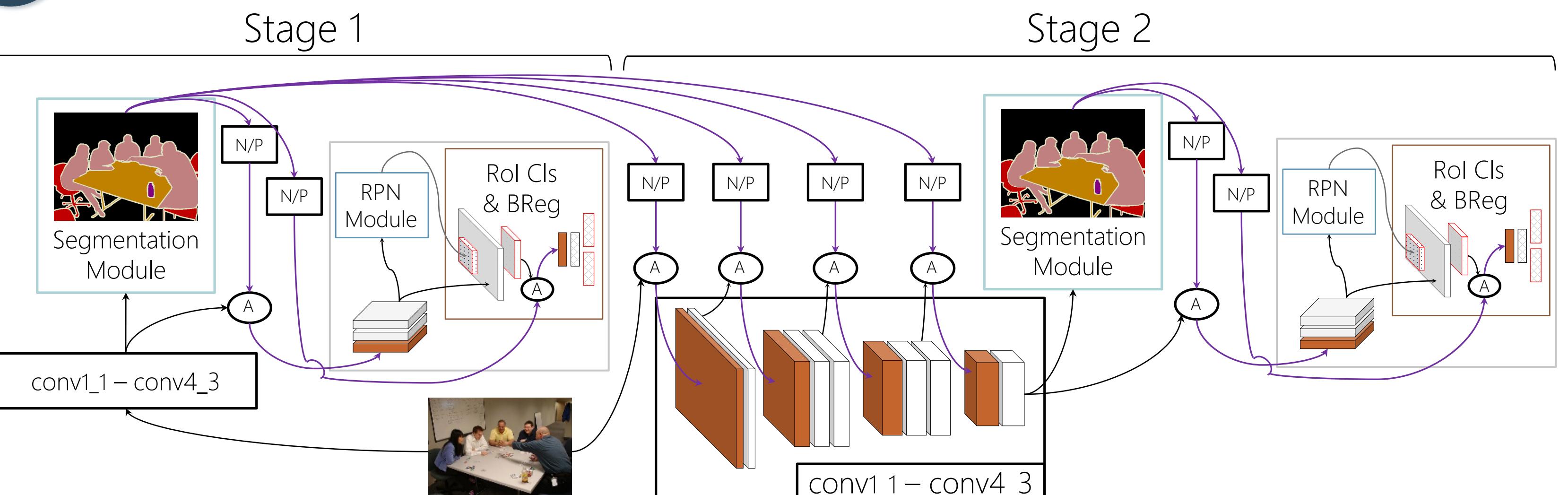
- Gradients from each RoI overpower segmentation network.

Ablation Analysis: Iterative Feedback

	Stage-2 Init.	mAP	mIOU
Base-MT	-	75.6	65.8
Feedback to conv1_1	ImageNet	76.5	69.3
Stage-1	76.3	69.3	
Feedback to conv{1,2,3,4}_1	ImageNet	76.3	69.1
Stage-1	77.3	69.5	

- More feedback helps when initializing with Stage-1 network (cf., unrolled self-feedback)

3 Joint Model: Contextual Priming and Feedback



Main Results on standard dataset splits

Detection results on VOC07 detection test set. All methods are trained on VOC07 trainval and VOC12 trainval

S	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	
Fast R-CNN	70.0	77.0	78.1	69.3	59.4	38.3	81.6	78.6	86.7	42.8	78.8	68.9	84.7	82	76.6	69.9	31.8	70.1	74.8	80.4	70.4	
Faster R-CNN	73.2	76.5	79.0	70.9	65.5	52.1	83.1	84.7	86.4	52	81.9	65.7	84.8	84.6	77.5	76.7	38.8	73.6	73.9	83.0	72.6	
Base-MT	✓	74.7	78.4	79.3	75.9	63.2	56.8	85.9	88.4	54.9	83.9	68.6	84.6	85.6	78.5	78.1	41.3	74.6	74.8	84.0	72.4	
Ours [joint]	✓	76.4	79.3	80.5	76.8	72.0	58.2	85.1	86.5	89.3	60.6	82.2	69.2	87.0	87.2	81.6	78.2	44.6	77.9	76.7	82.4	71.9

+8.8 +5.7 +3.3 +3.3

Detection results on VOC12 detection test set. All methods are trained on VOC07 trainval+test and VOC12 trainval

S	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv	
Fast R-CNN	68.4	82.3	78.4	70.8	52.3	38.7	77.8	71.6	89.3	44.2	73.0	55.0	87.5	80.5	80.8	72	35.1	68.3	65.7	80.4	64.2	
Faster R-CNN	70.4	84.9	79.8	74.3	53.9	49.8	77.5	75.9	88.5	45.6	77.1	55.3	86.9	81.7	80.9	79.6	40.1	72.6	60.9	81.2	61.5	
Base-MT	✓	71.1	84.2	80.9	73.1	55.1	50.6	78.2	75.6	89.0	48.6	76.7	54.8	87.6	82.5	83.0	80.0	41.7	74.2	60.7	81.4	63.1
Ours [joint]	✓	72.6	84.0	81.2	75.9	60.4	51.8	81.2	77.4	90.9	50.2	77.6	58.7	88.4	83.6	82.0	80.4	41.5	75.0	64.2	82.9	65.1

+5.3 +3.0 +3.9 +3.9

+3.5

Segmentation results on VOC12 segmentation test set. All methods are trained on 07 trainval+test and 12 trainval

S	mIOU	bg	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	persn	plant	sheep	sofa	train	tv

<tbl_r cells="22" ix="3" maxcspan="1" maxrspan="1" used