

Lecture 8: Image Segmentation

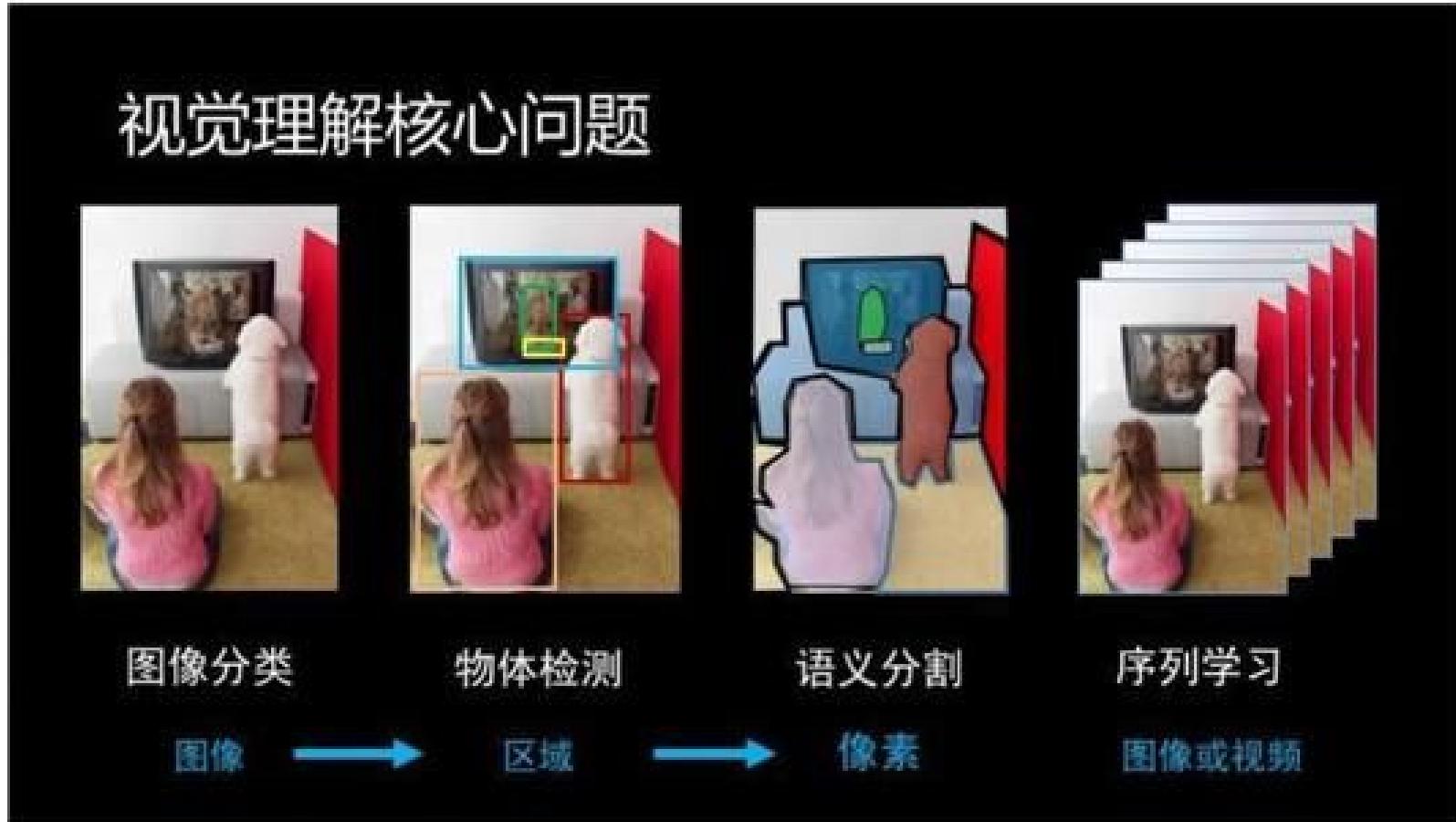
彭 超 Peng Chao

Face++ Researcher

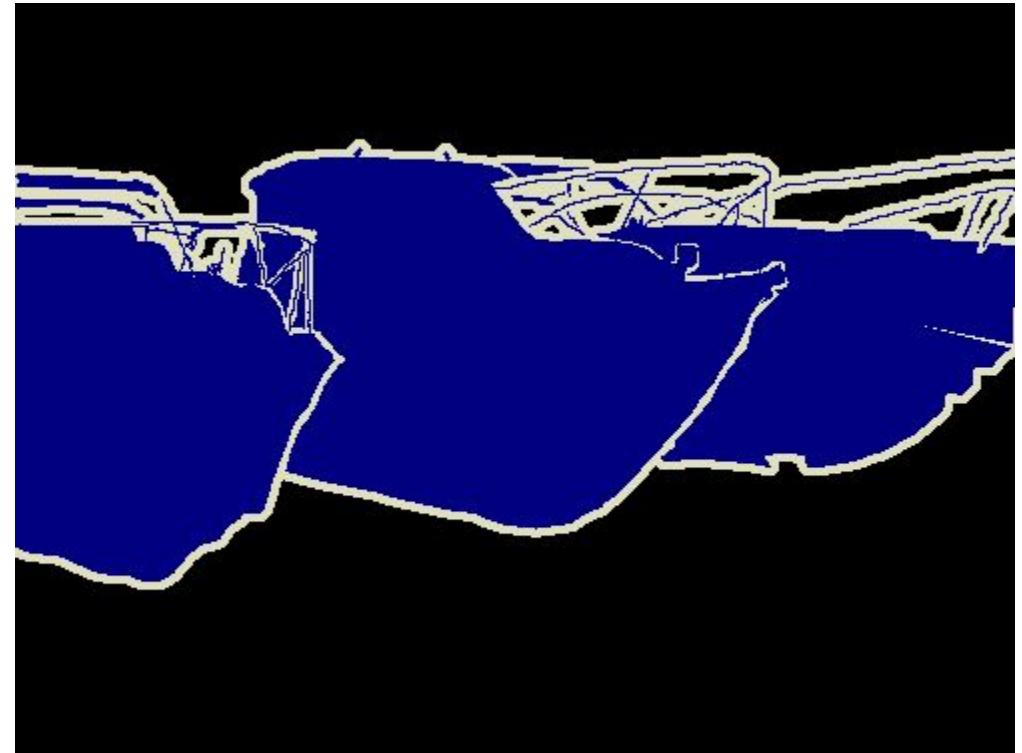
pengchao@megvii.com

Nov. 2017

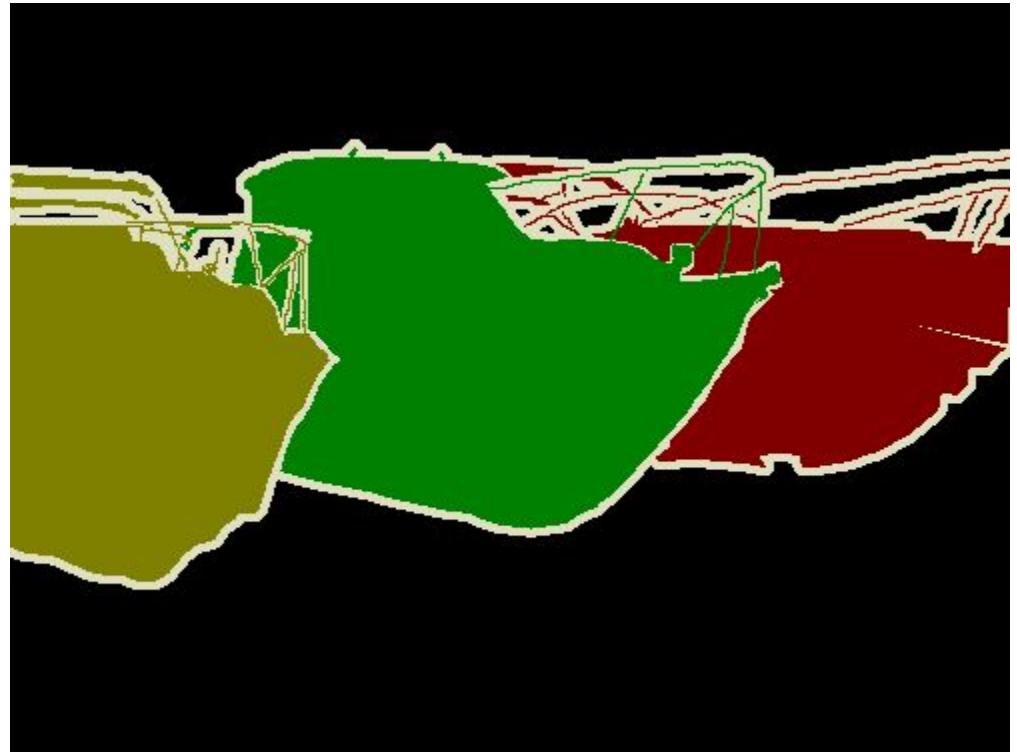
Image Segmentation



Semantic Segmentation



Instance Segmentation



Scene Parsing



Human Parsing

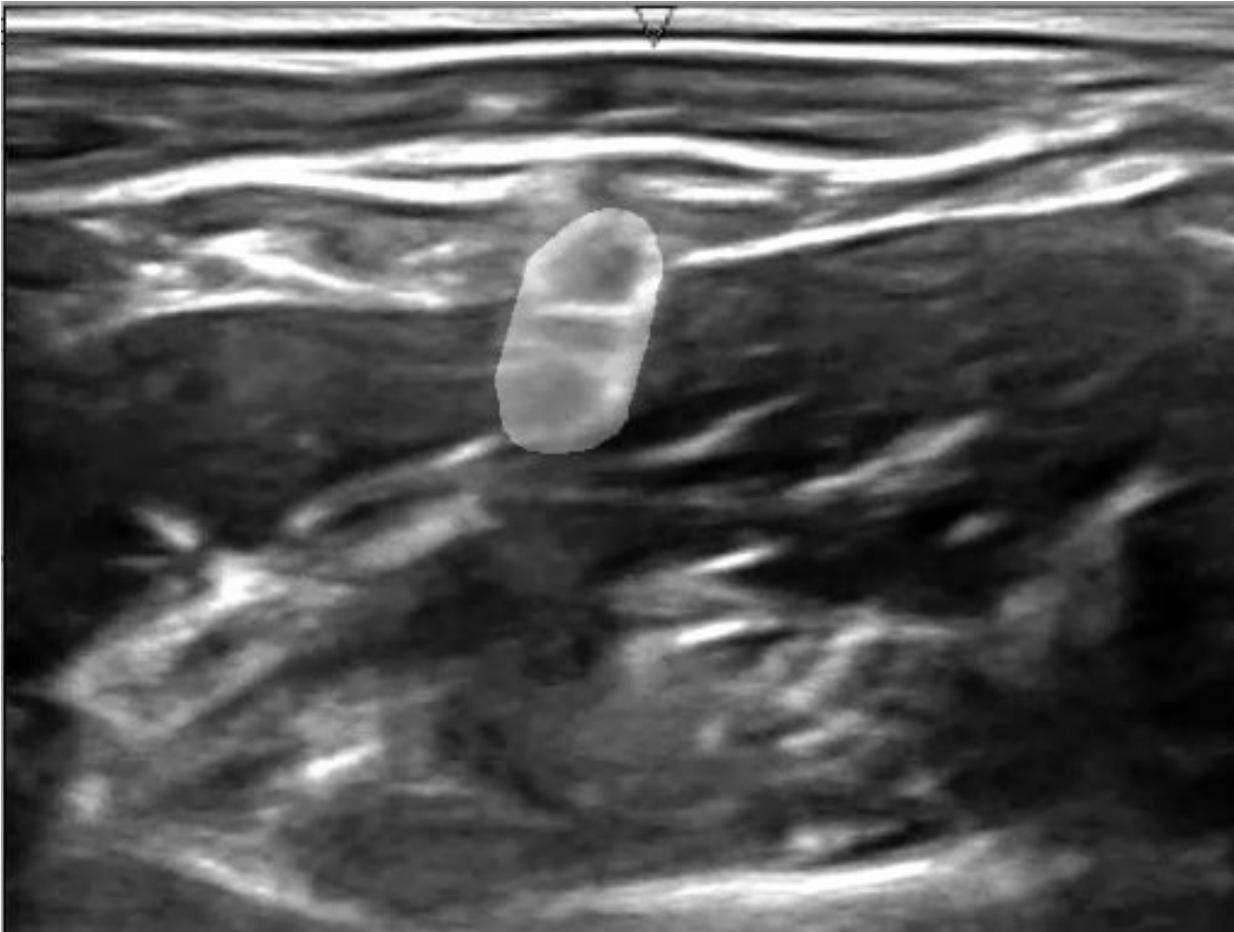


Stuff Segmentation



- New Track in COCO 2017
- Stuff: mountain, grass, wall, sky
-
- Stuff covers about 66% of the pixels in COCO

Ultrasound Segmentation



Selfie Segmentation



Evaluation

$$Accuracy(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{i=0}^n \frac{I[y_i = \hat{y}_i]}{n}$$

$$mean\ IOU(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\sum_c^m \sum_i^n I[y_i = c, \hat{y}_i = c]}{\sum_c^m \sum_i^n I[y_i = c \text{ or } \hat{y}_i = c]}$$

Normally, we use **mean IOU** to judge the results!

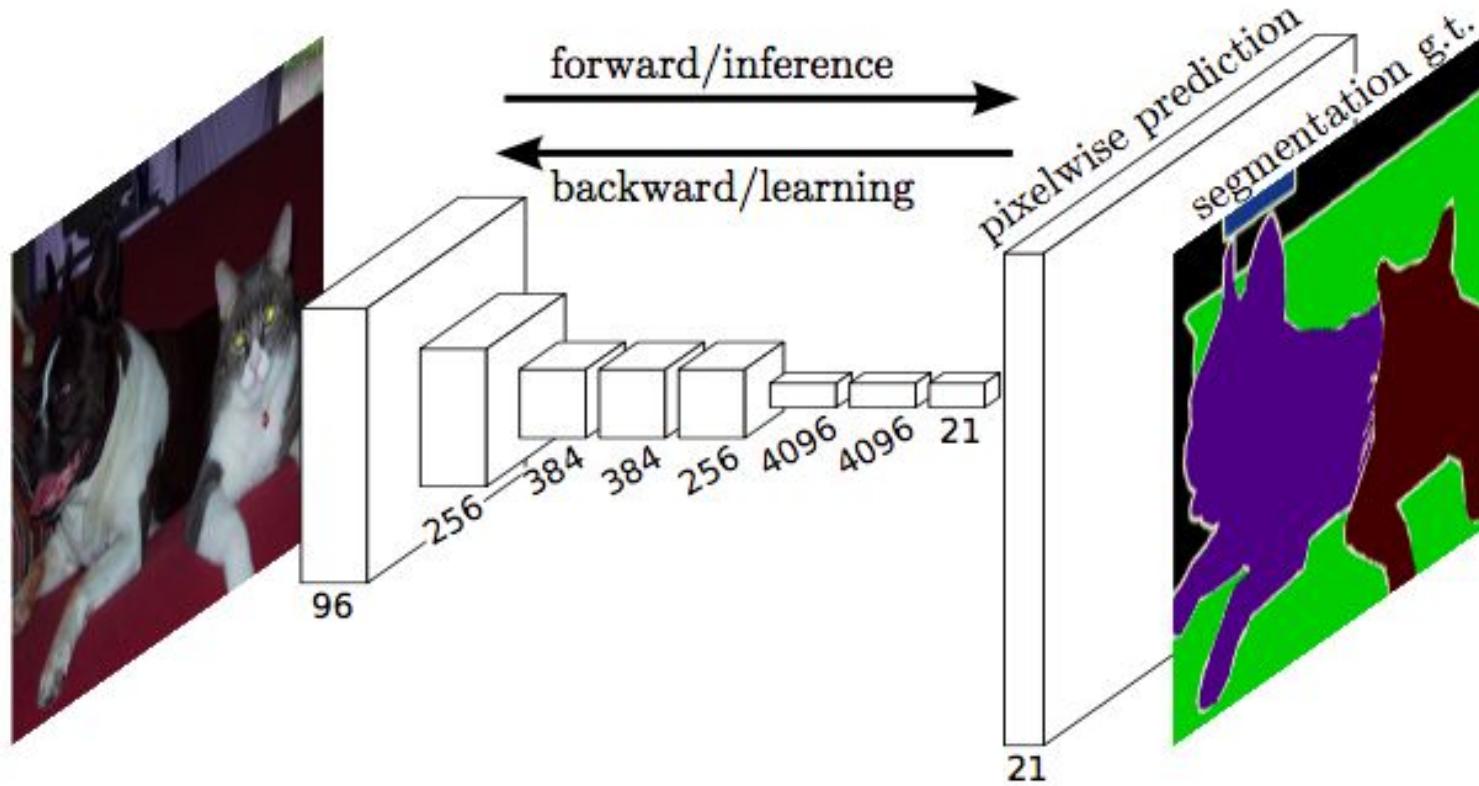
Outline

- Semantic Segmentation
- Instance Segmentation

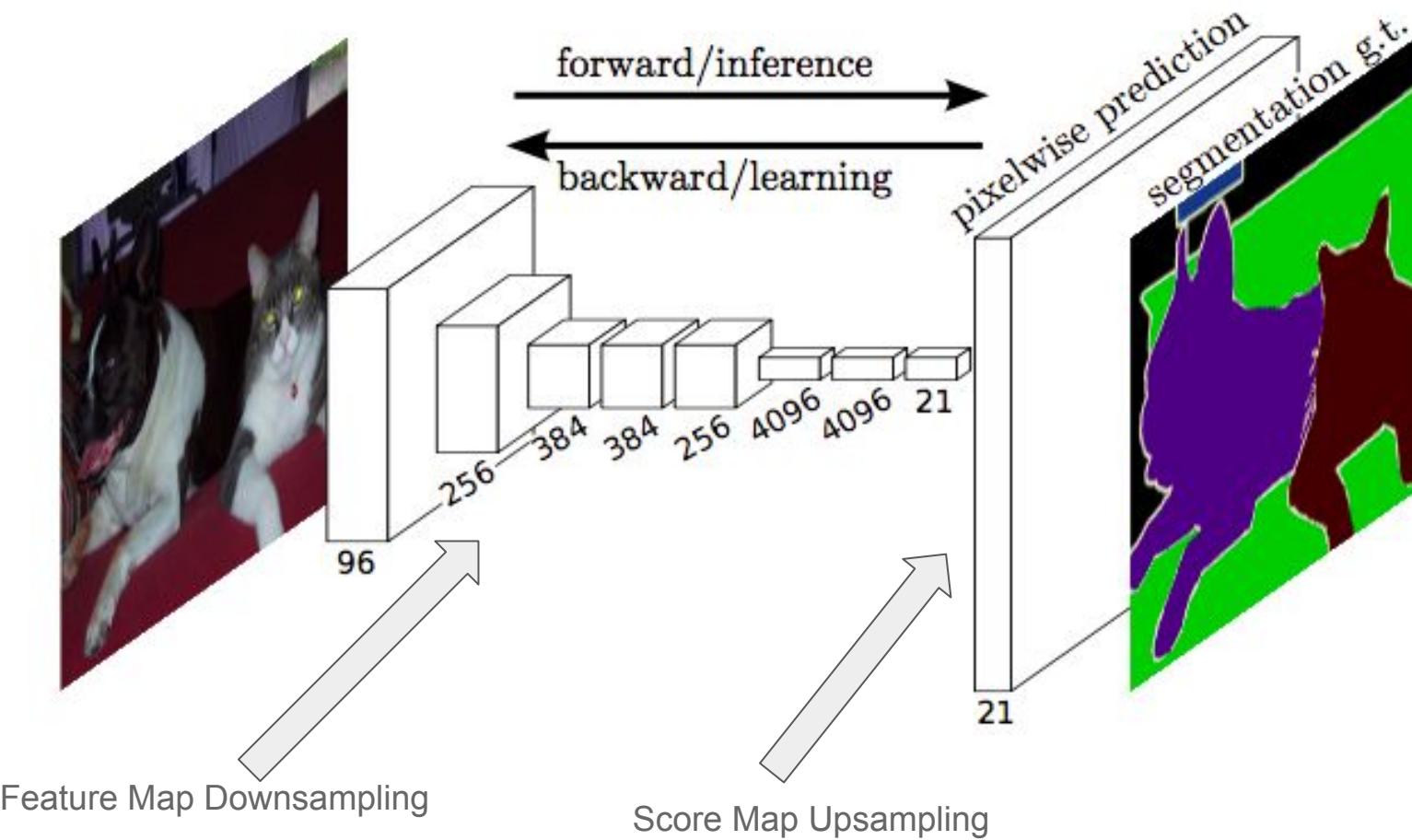
Outline

- Semantic Segmentation
- Instance Segmentation

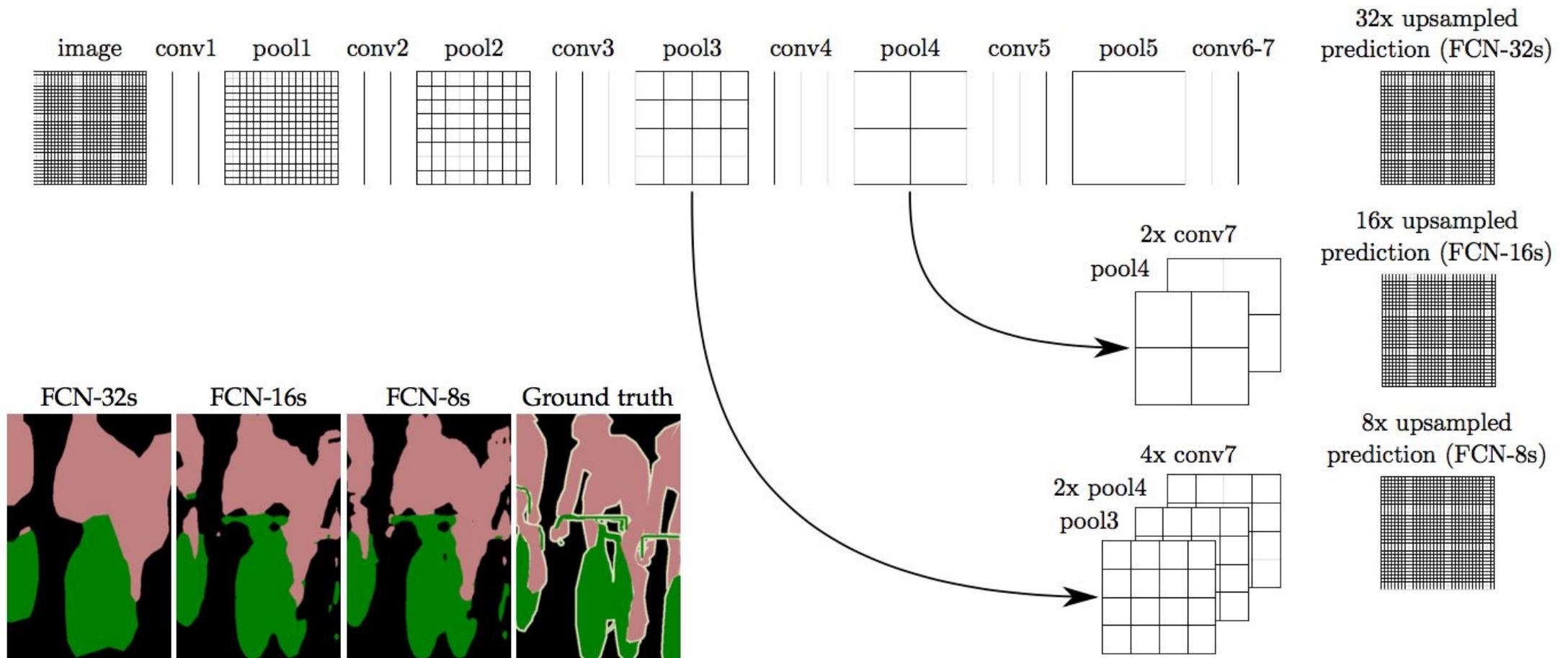
Fully Convolutional Network



Fully Convolutional Network



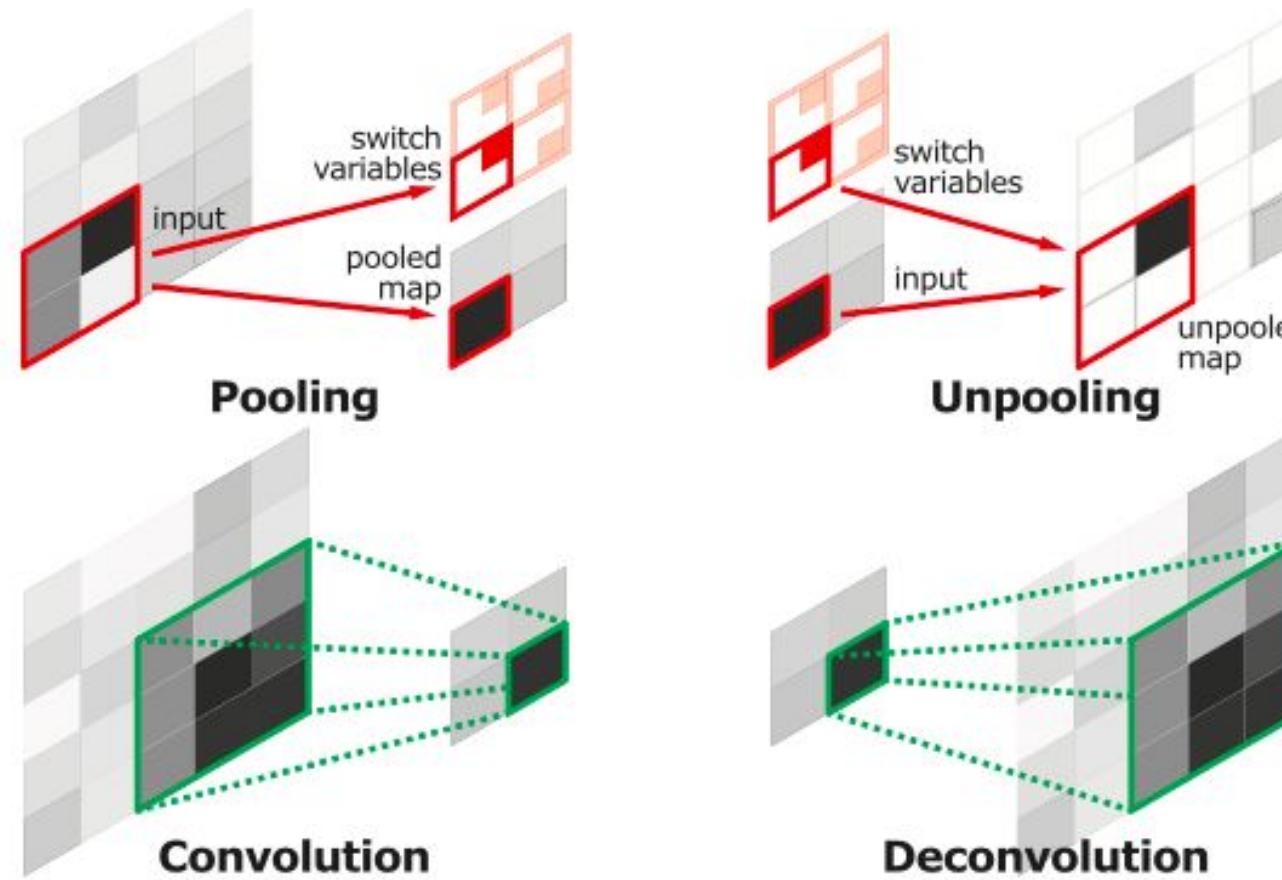
Fully Convolutional Network



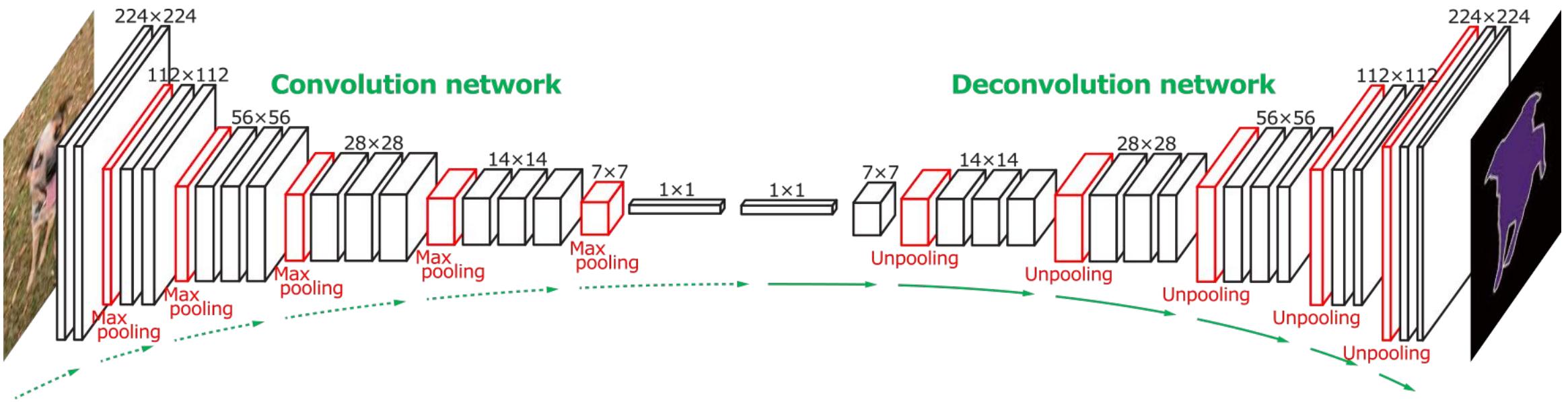
Fully Convolutional Network

- First work using CNN to solve the semantic segmentation
- Introducing skip-net framework
- Large Improvement! (60 vs 30)

Learning Deconvolution Network for Semantic Segmentation



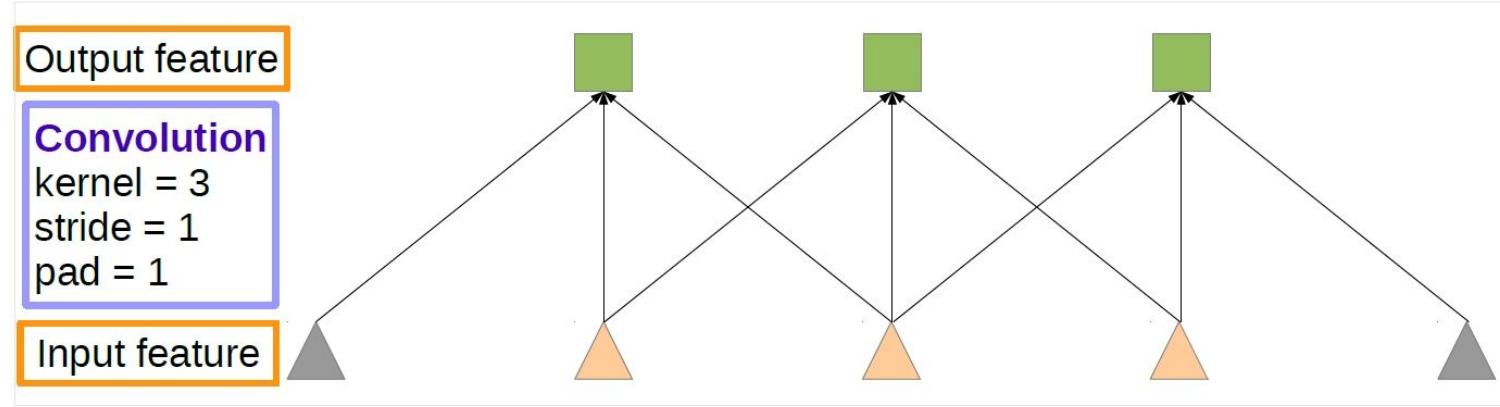
Learning Deconvolution Network for Semantic Segmentation



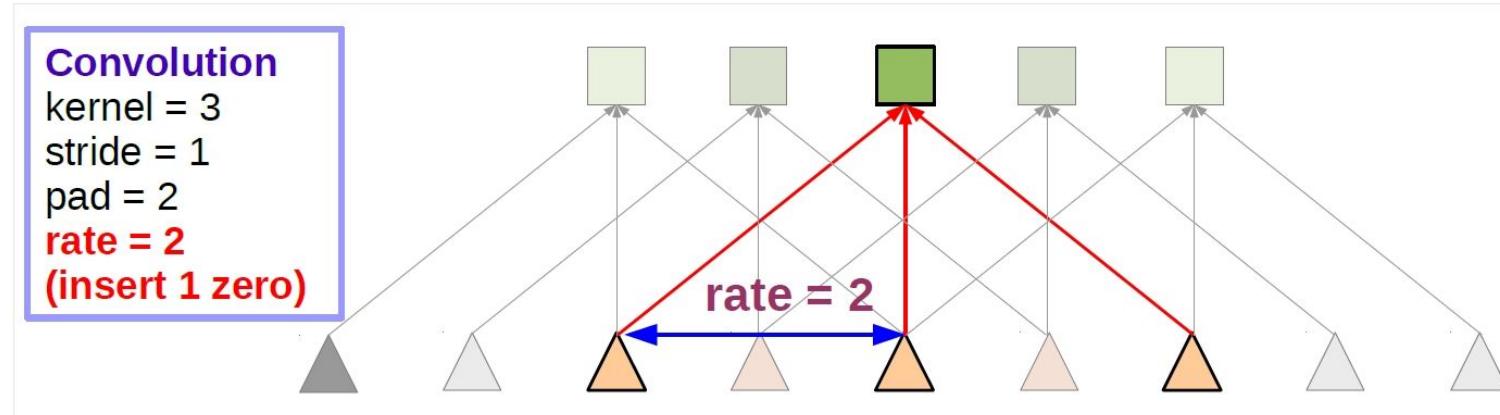
Learning Deconvolution Network for Semantic Segmentation

- Introducing un-pool and de-convolution operations.
- Introducing hourglass-like framework.

DeepLab

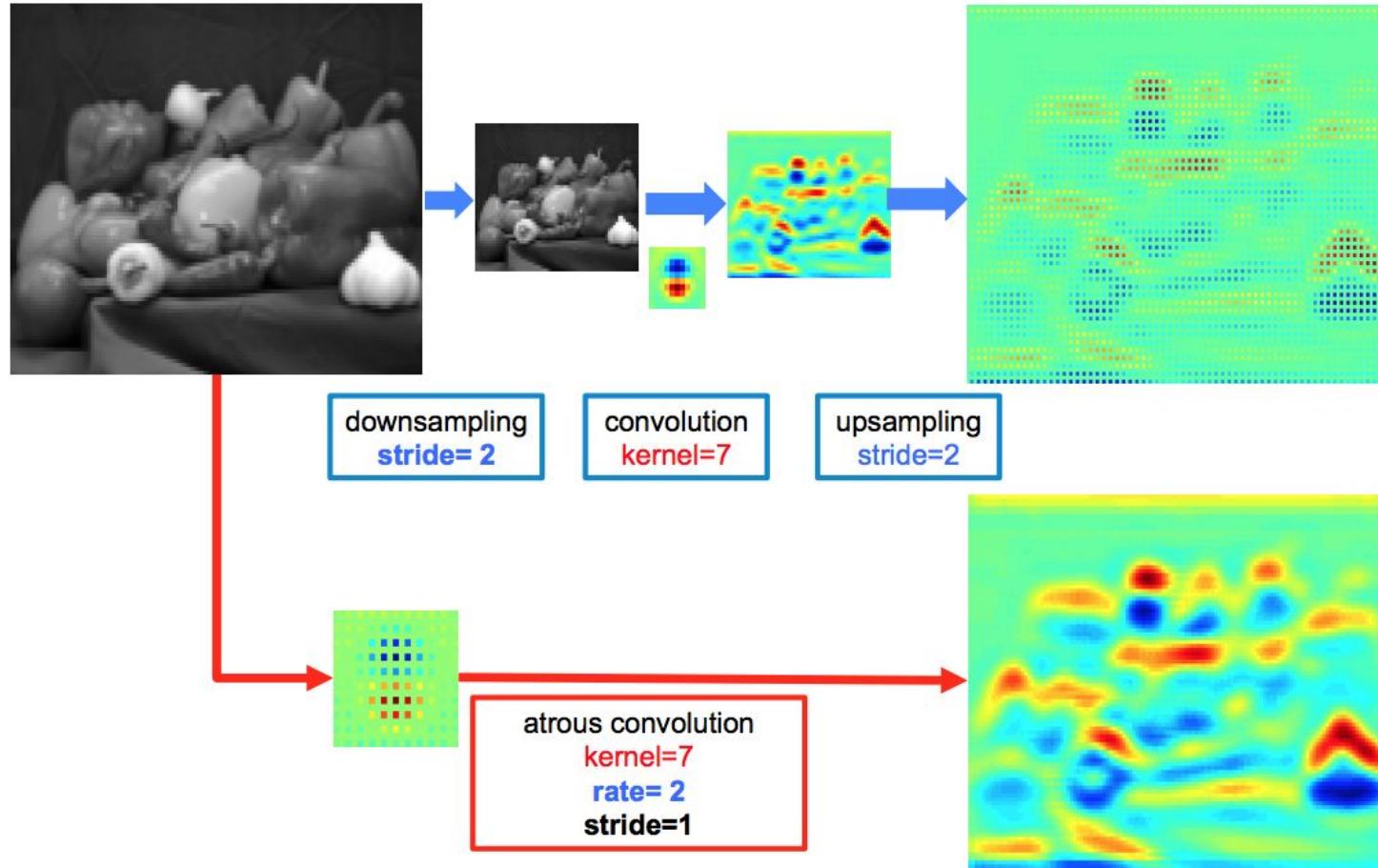


(a) Sparse feature extraction

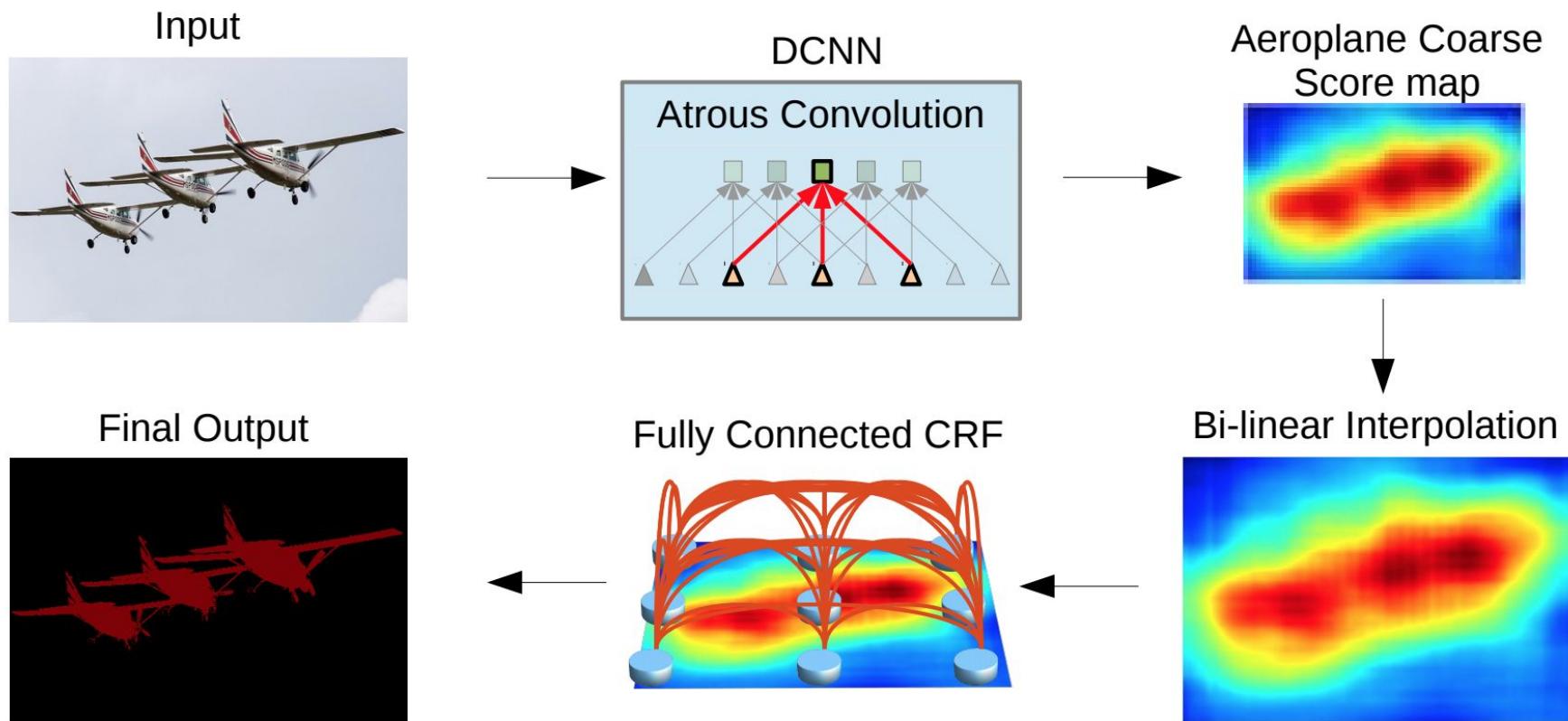


(b) Dense feature extraction

DeepLab



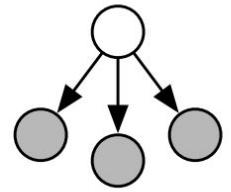
DeepLab



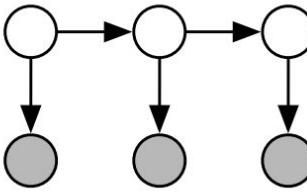
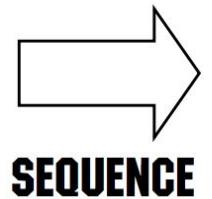
DeepLab

- Introducing dilated-convolution
- Combining traditional method (post processing): **DenseCRF**

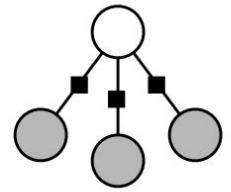
Conditional Random Field



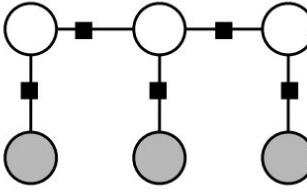
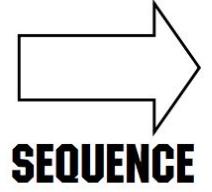
Naive Bayes



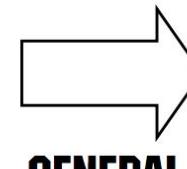
HMMs



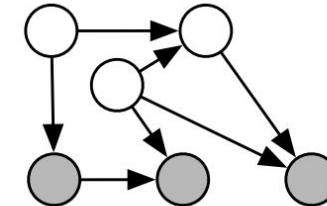
Logistic Regression



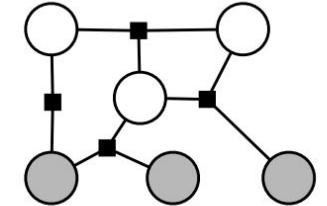
Linear-chain CRFs



GENERAL
GRAPHS



Generative directed models



General CRFs

Conditional Random Field

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp[-E(\mathbf{y}, \mathbf{x})]$$

\mathbf{y} is the label, \mathbf{x} is the image

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \exp[-E(\mathbf{y}, \mathbf{x})]$$

$$E(\mathbf{y}, \mathbf{x}) = \sum_{U \in \mathcal{U}} \sum_{p \in \mathcal{N}_U} U(y_p, \mathbf{x}_p) + \sum_{V \in \mathcal{V}} \sum_{(p,q) \in \mathcal{S}_V} V(y_p, y_q, \mathbf{x}_{pq})$$

U: Unary relation;
V: pairwise relation

CRF Inference

$$\forall p \in \mathcal{N} : P(y_p | \mathbf{x}) = \sum_{\mathbf{y} \setminus y_p} P(\mathbf{y} | \mathbf{x}).$$

- However, for **loopy graph**, the above problem is NP-hard. (The nodes relations are complex, making computing marginal probability harder)
- Approximated methods:
 - MCMC (Gibbs Sampling)
 - Loopy Belief propagation
 - **Mean Field**

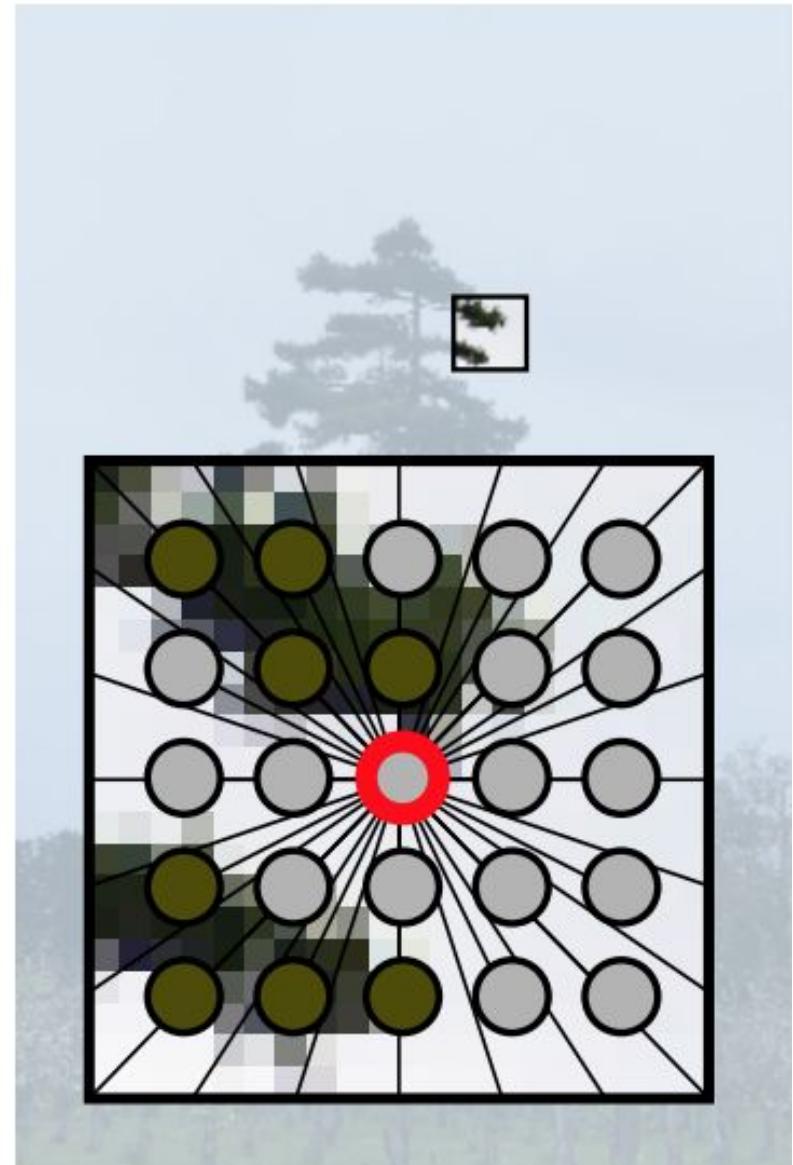
CRF Inference: Mean field

$$\min \sum_{\mathbf{y}} Q(\mathbf{y}) \log\left(\frac{Q(\mathbf{y})}{P(\mathbf{y}|\mathbf{x})}\right)$$

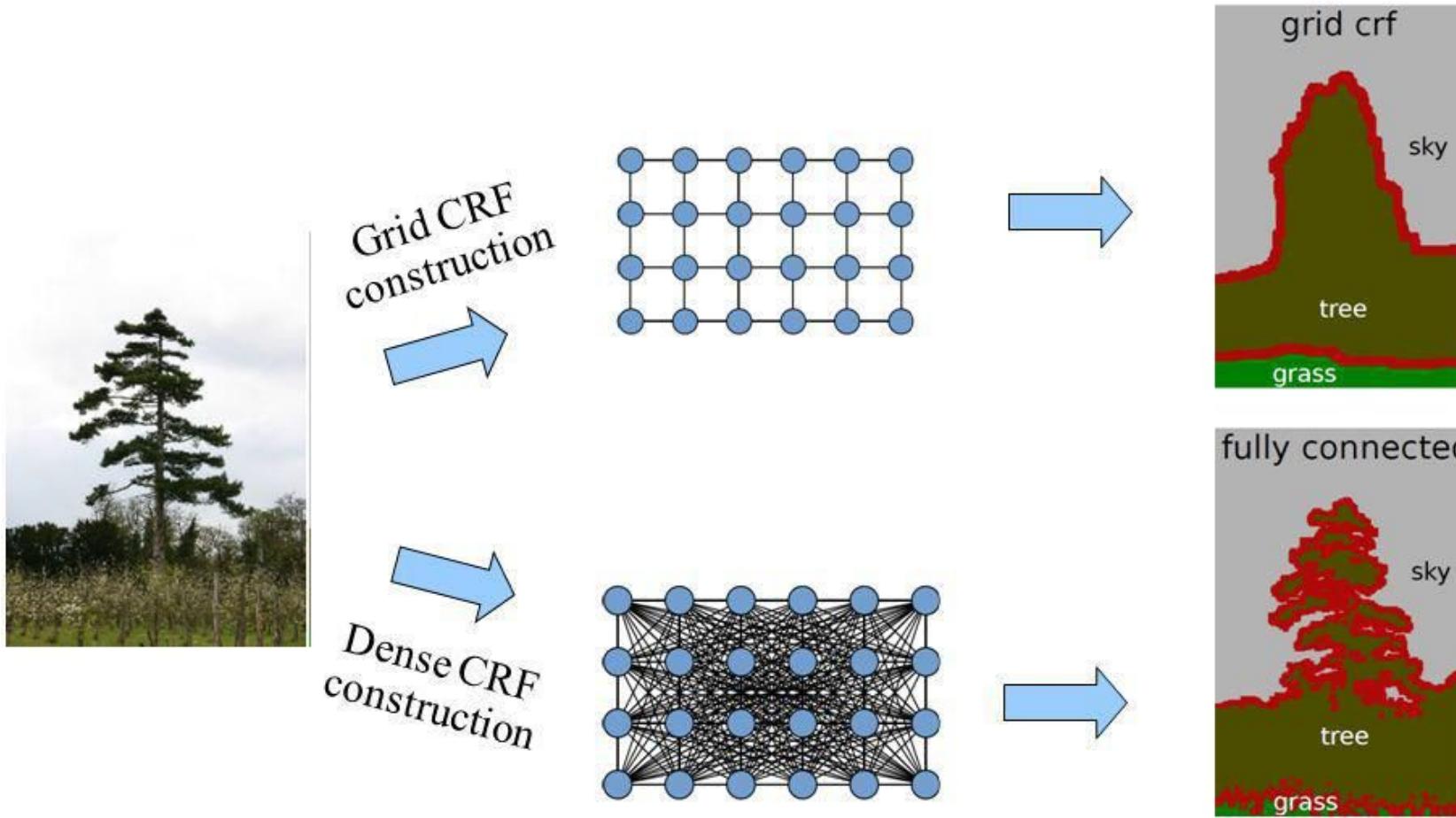
$$s.t. \quad Q(\mathbf{y}) = \prod_i Q(y_i)$$

$$\sum_{\mathbf{y}} Q(\mathbf{y}) = 1$$

DenseCRF



DenseCRF

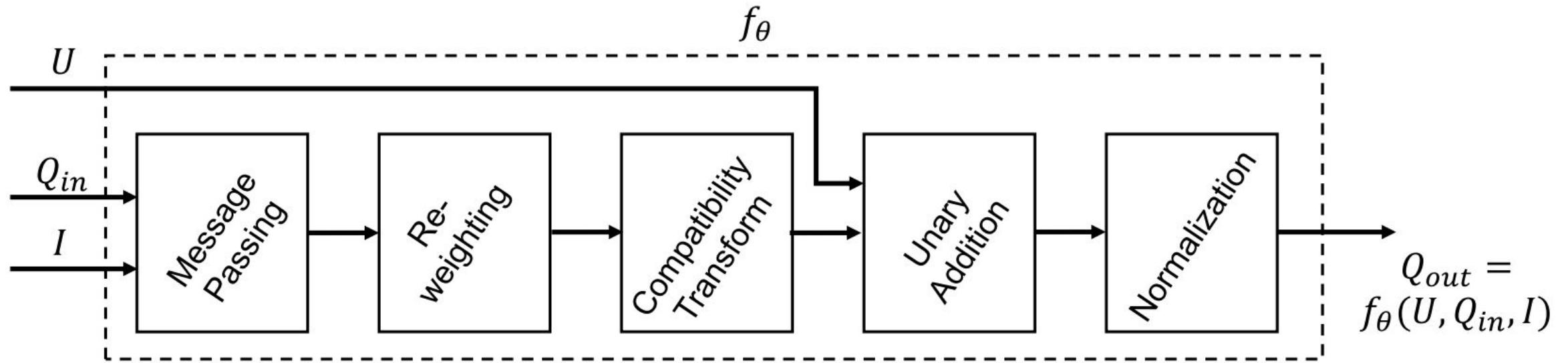


2

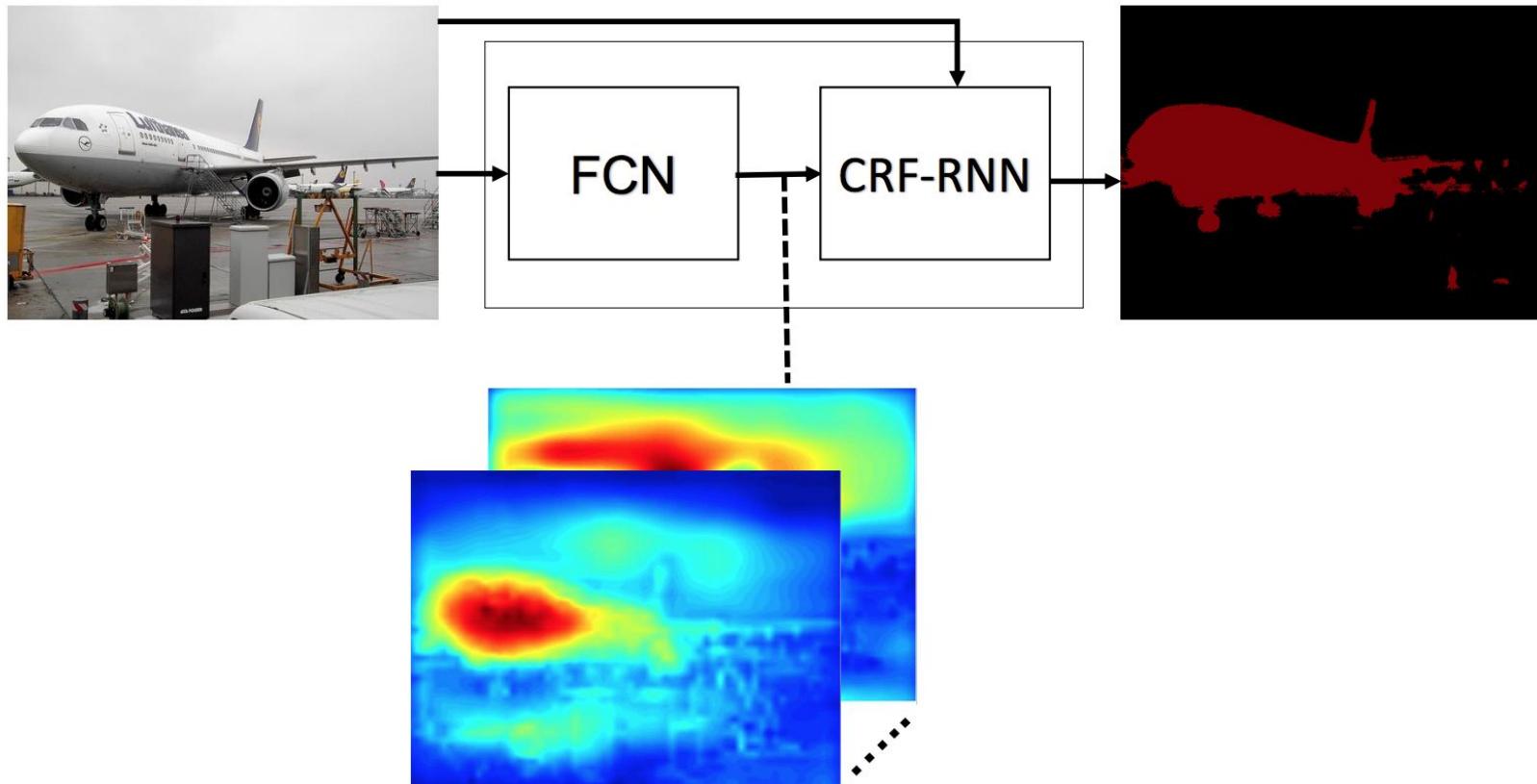
DenseCRF

- Best Traditional Method!
- Poor accuracy on segmentation! (**poor feature**)

CRF AS RNN



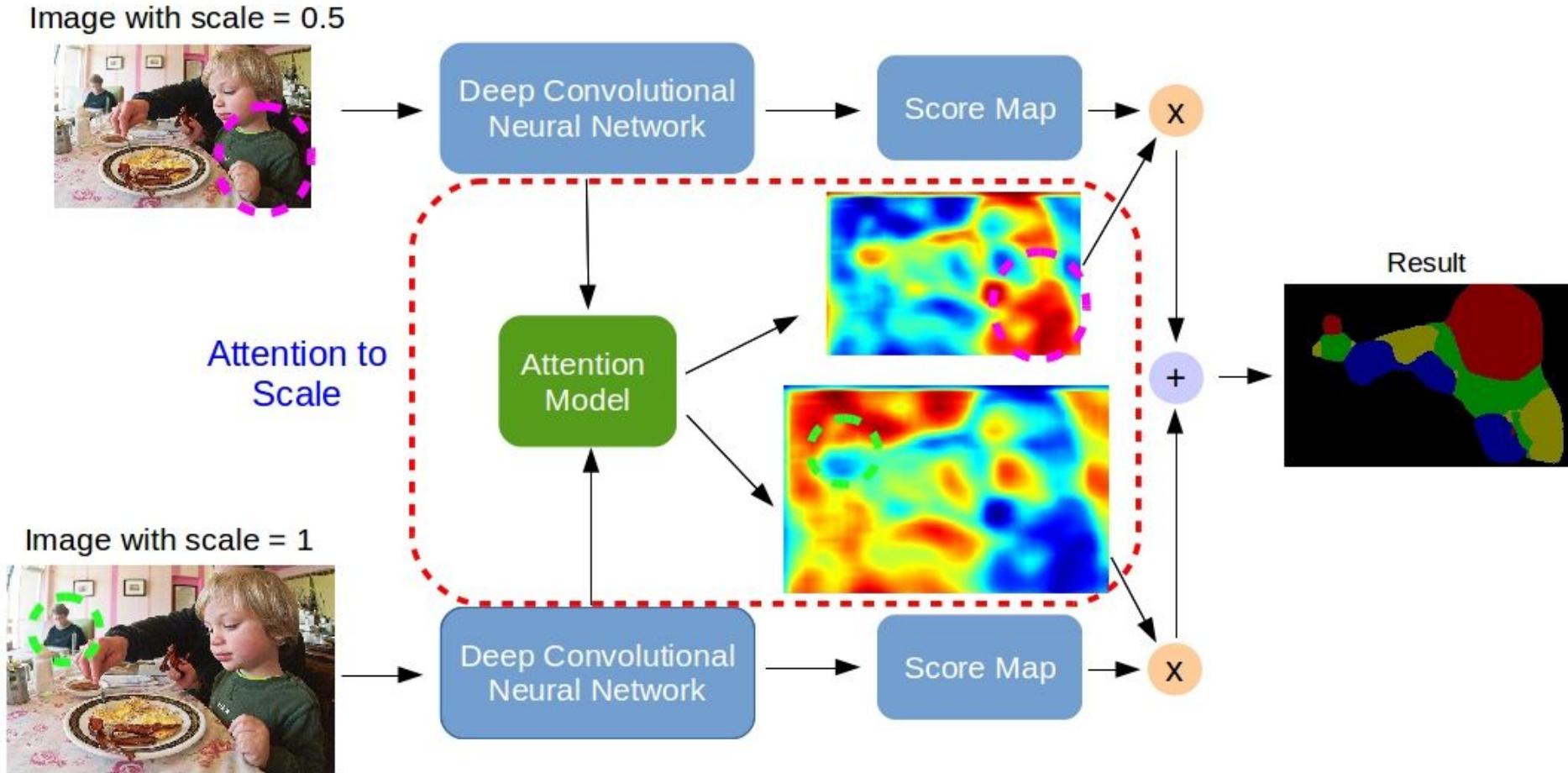
CRF AS RNN



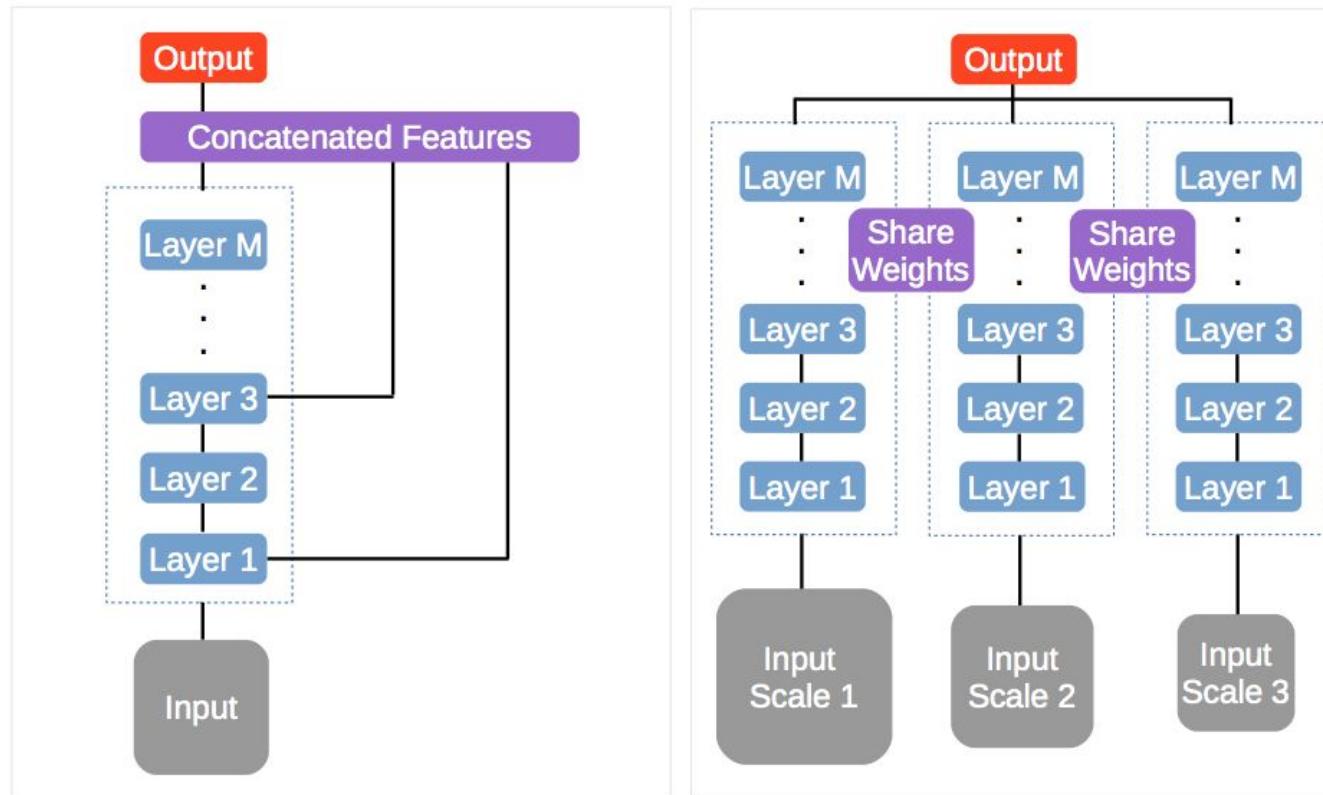
CRF AS RNN

- Encoding the DenseCRF into the CNN framework!
- Better results than DenseCRF in post-processing (Deeplab)

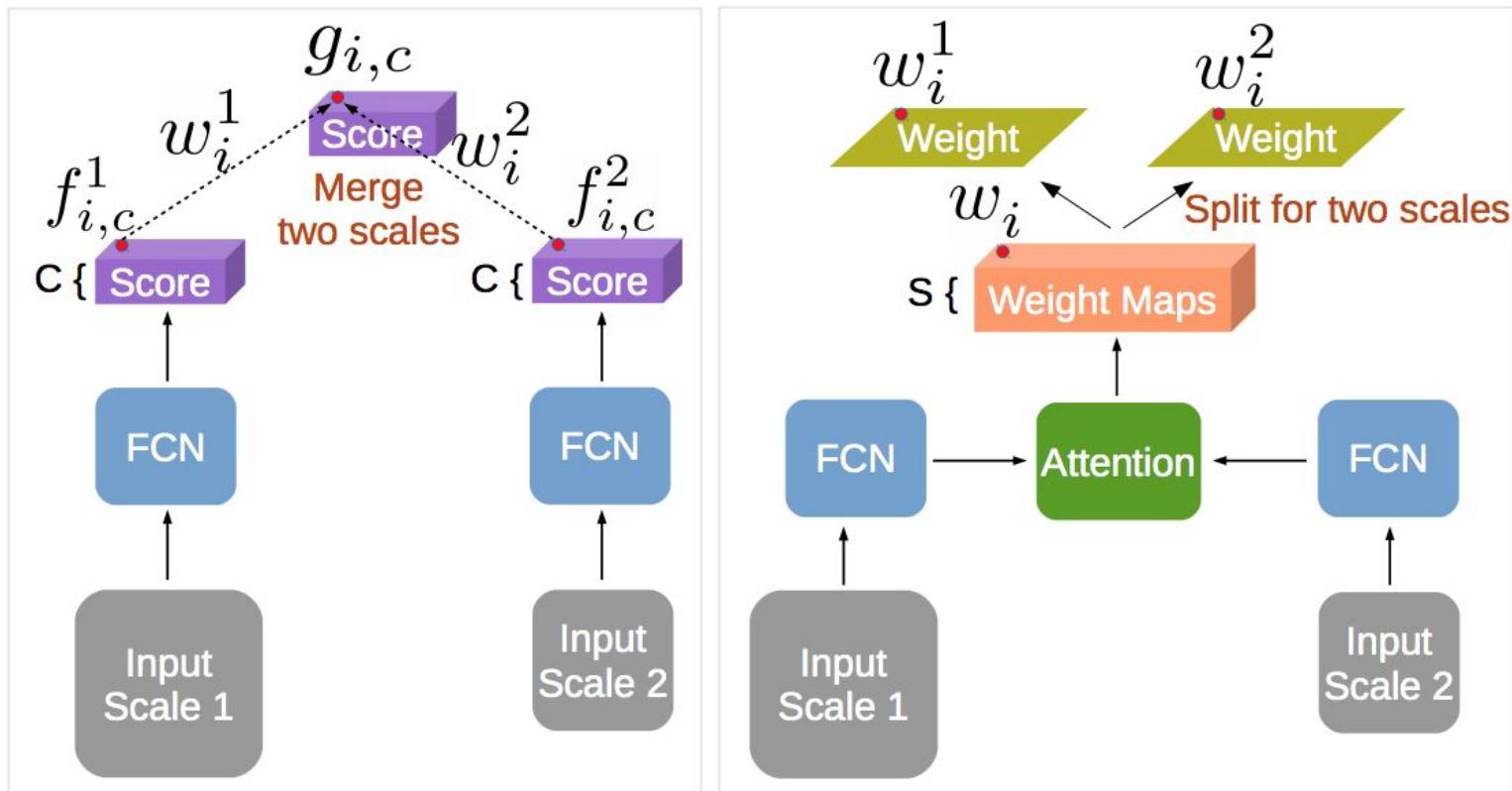
Deeplab Attention



Deeplab Attention



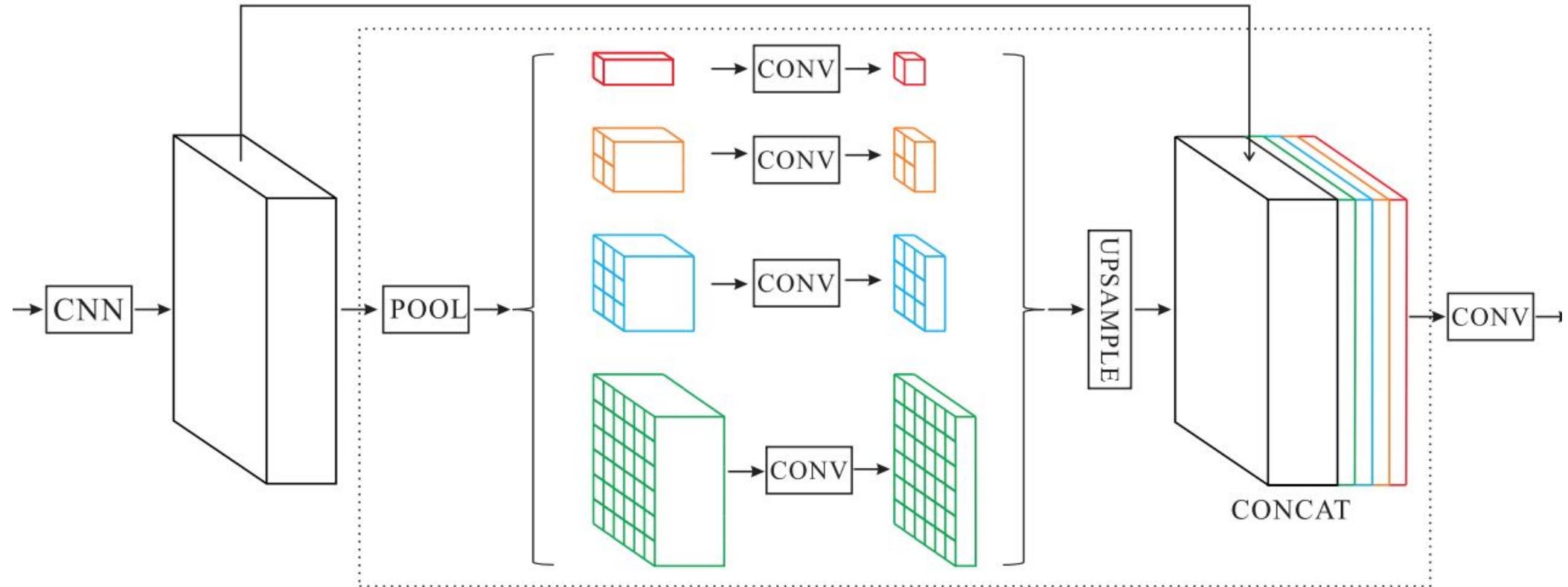
Deeplab Attention



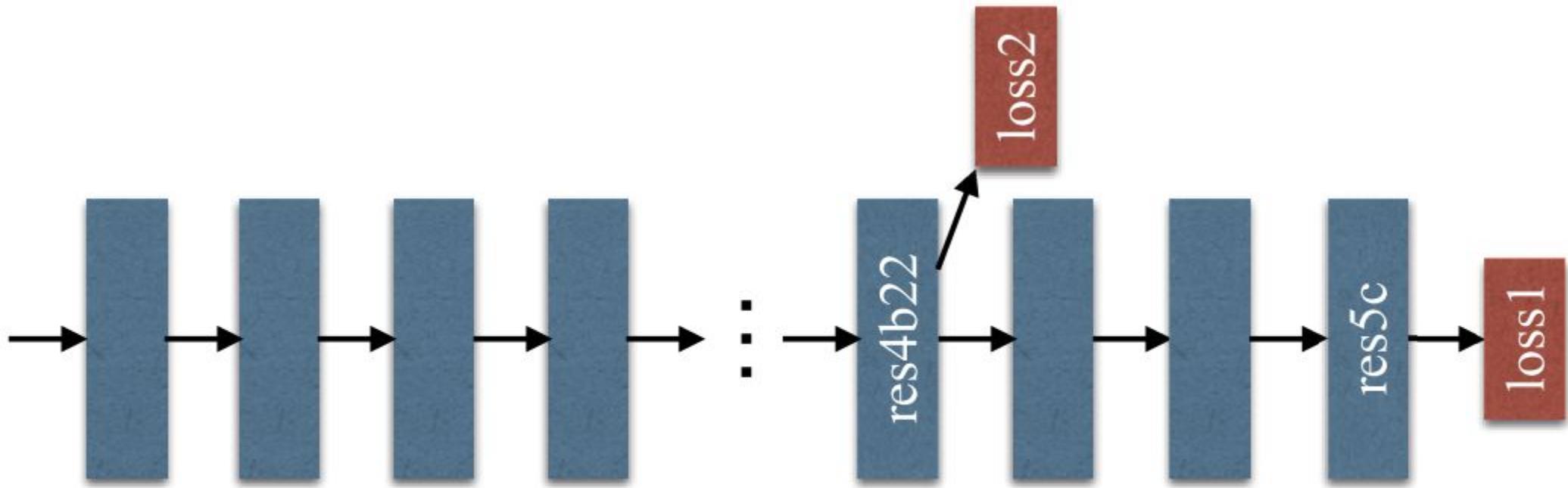
Deeplab Attention

- Fusion framework for multi-scale training and inference!
- Combining Attention model into Segmentation framework!

PSPNet

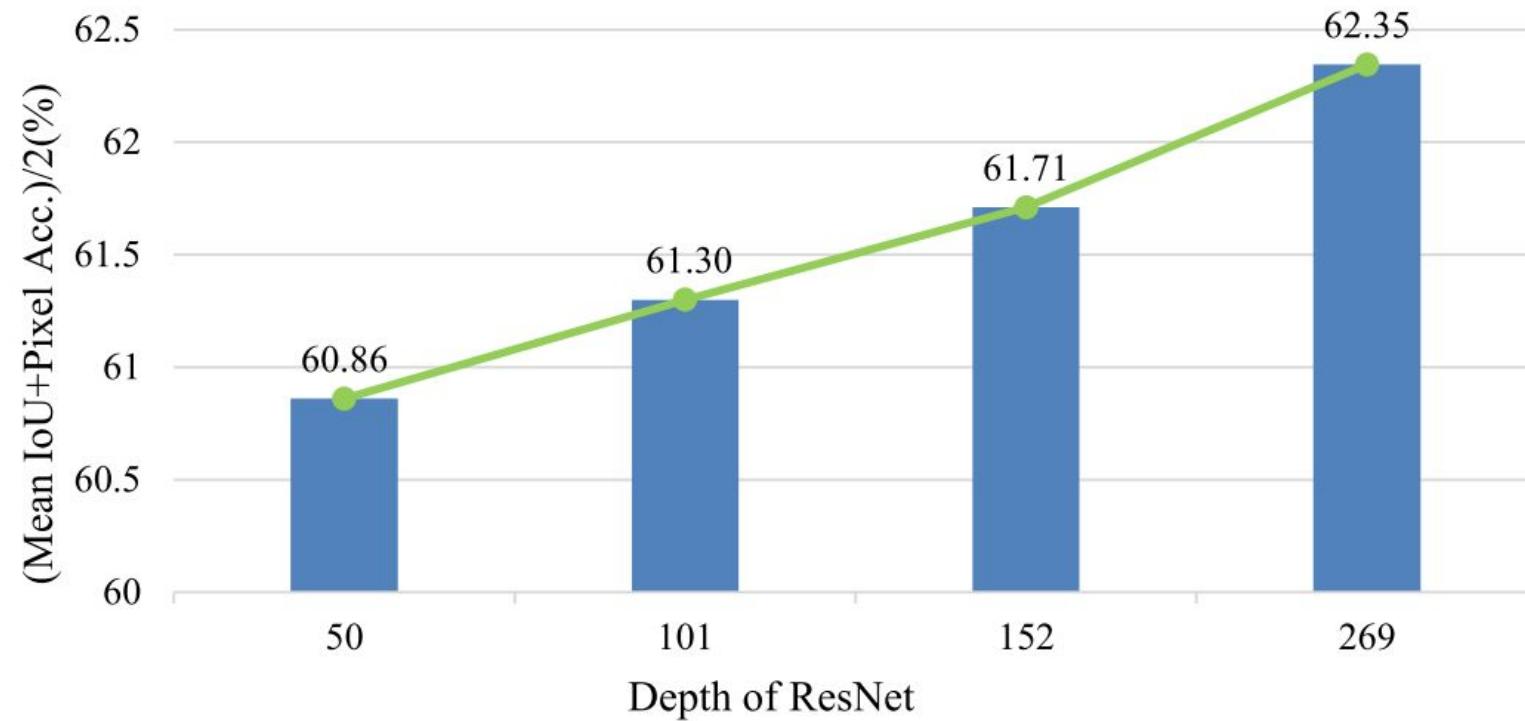


PSPNet



PSPNet

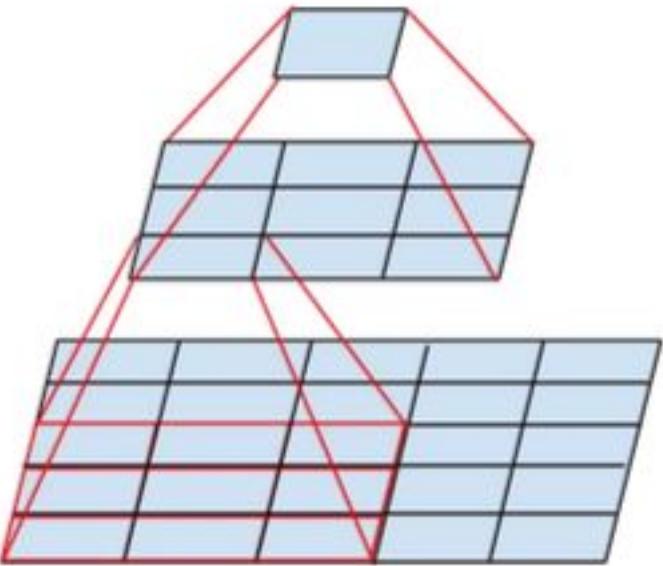
Performance of PSPNet with different pre-trained ResNet on ADE20K validation set



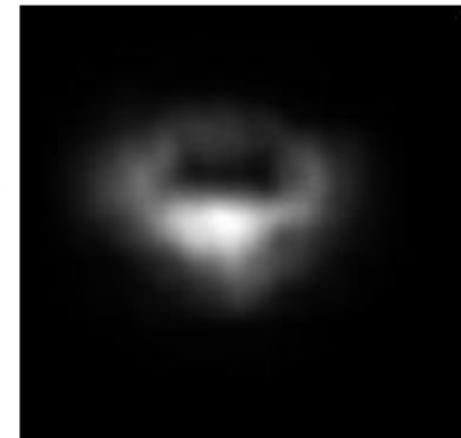
PSPNet

- Propose the Pyramid Pooling Module!
- Hard to reproduce!

Global Convolutional Network



TRF of two 3x3 conv is: 5

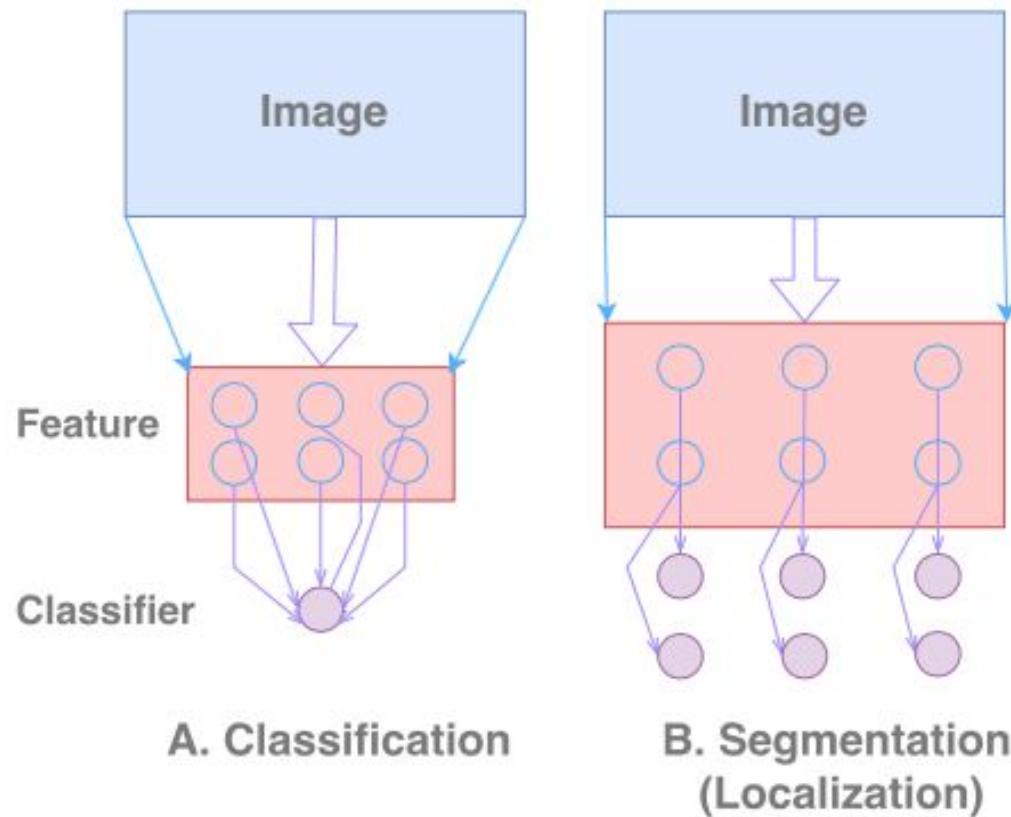


receptive field

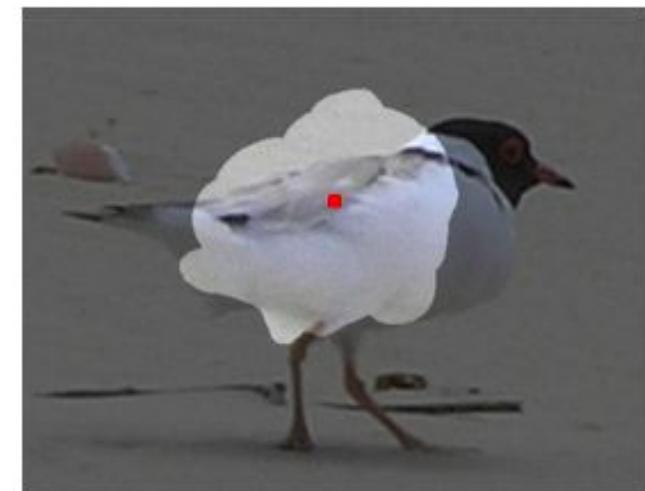
However the VRF maybe different !

Figure credit: Szegedy C, Vanhoucke V, Ioffe S, et al. Rethinking the Inception Architecture for Computer Vision[J]. Computer Science, 2016.
Zhou B, Khosla A, Lapedriza A, et al. Object Detectors Emerge in Deep Scene CNNs[J]. Computer Science, 2015.

Global Convolutional Network



A

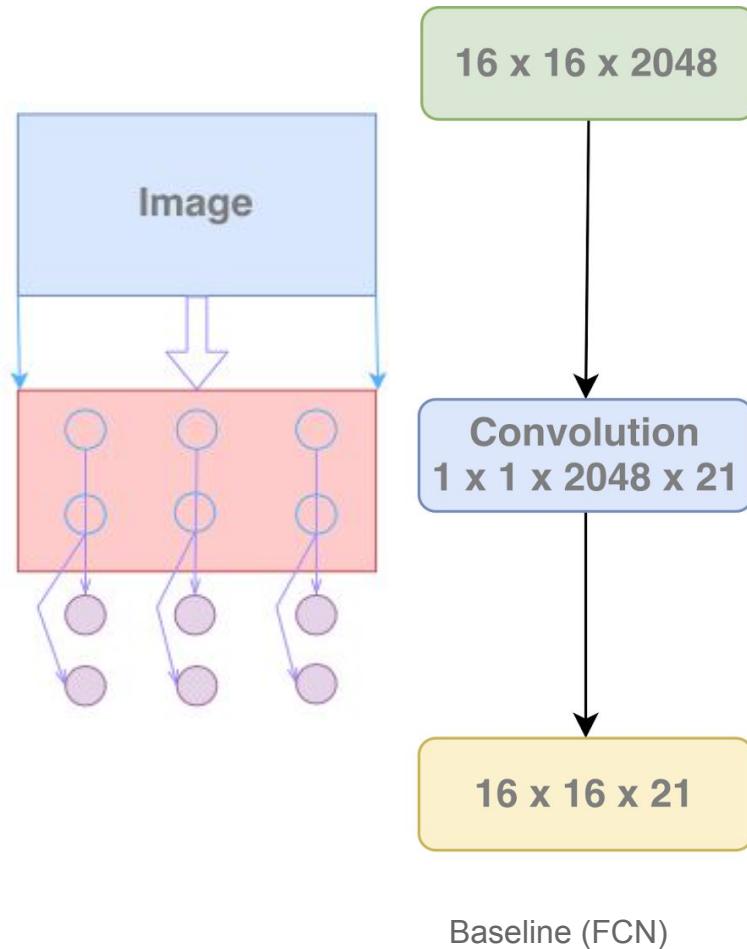


B

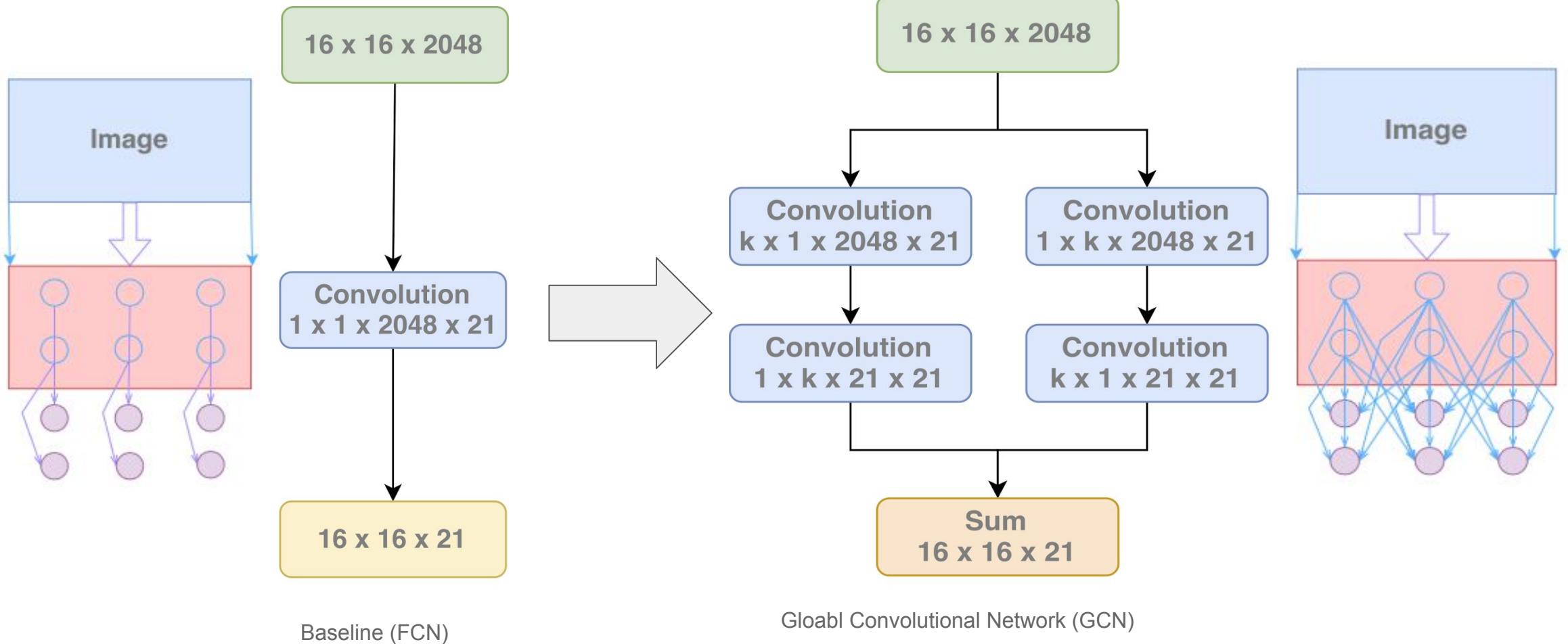
Global Convolutional Network



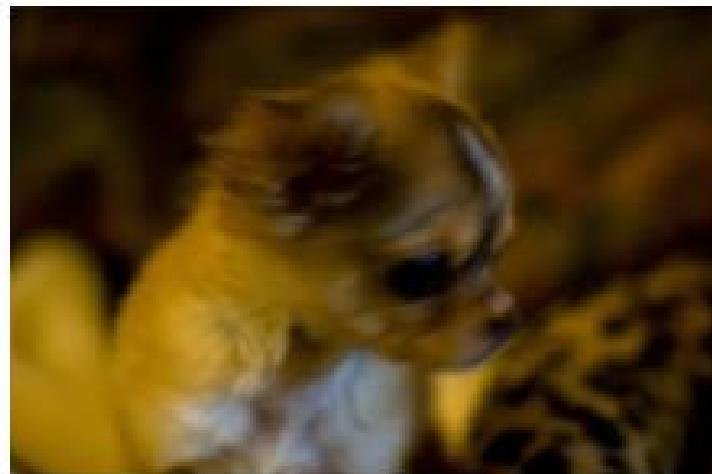
Global Convolutional Network



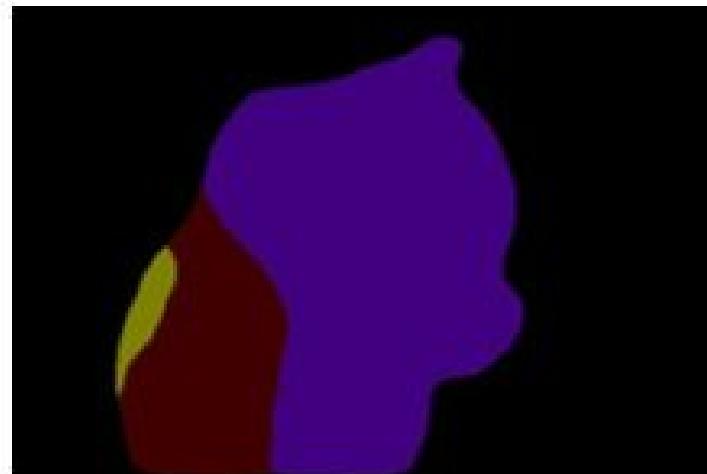
Global Convolutional Network



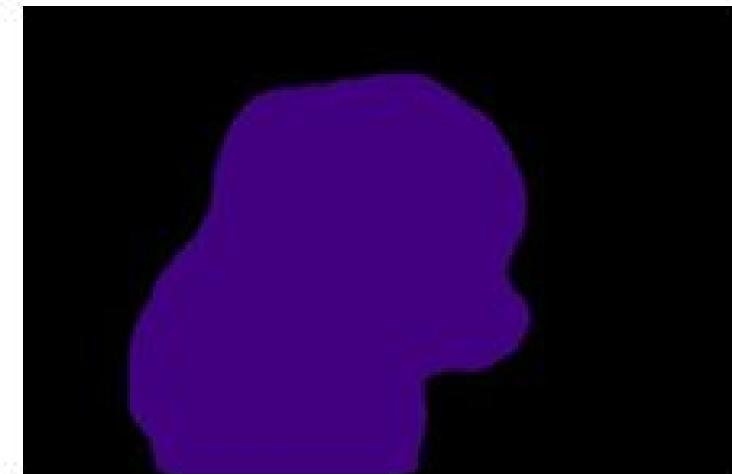
Global Convolutional Network



Image



Baseline (FCN)

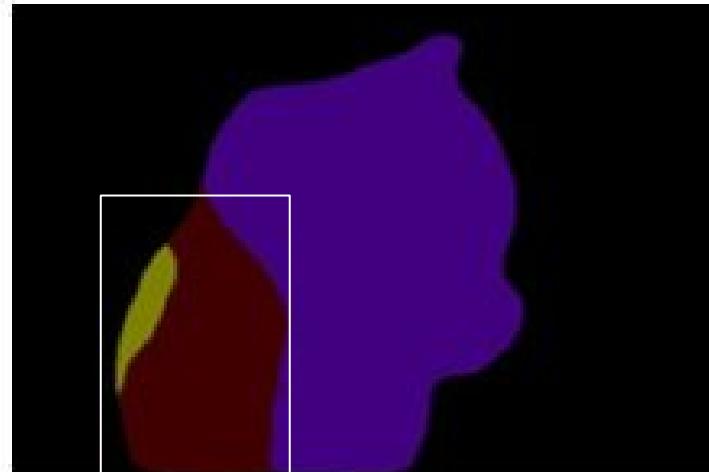


GCN

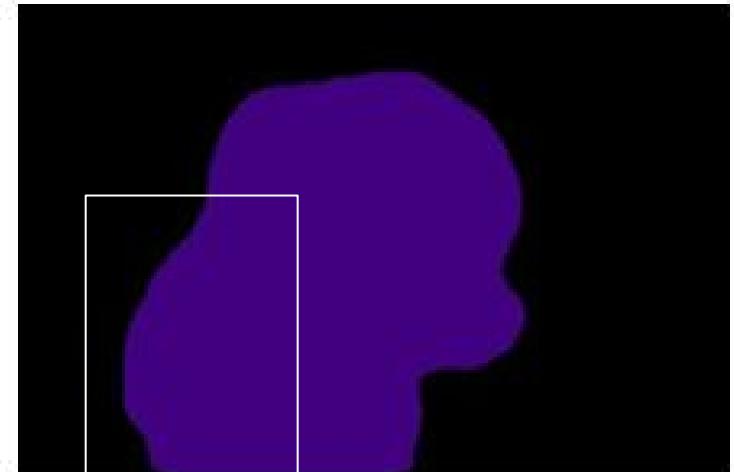
Global Convolutional Network



Image



Baseline (FCN)

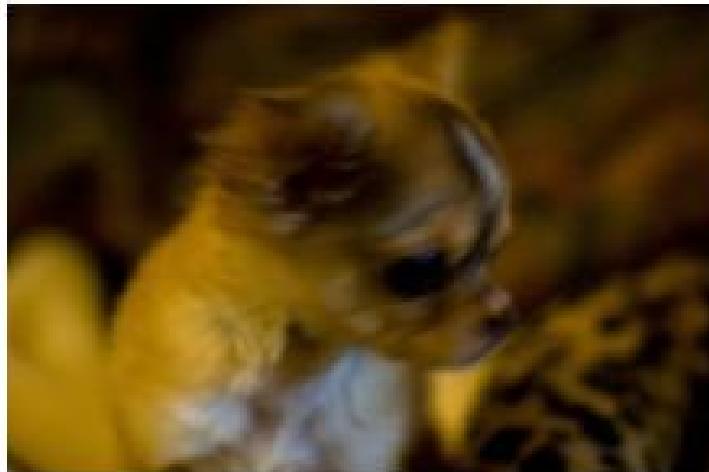


GCN

Region Mis-Classifications are corrected!

Global Convolutional Network

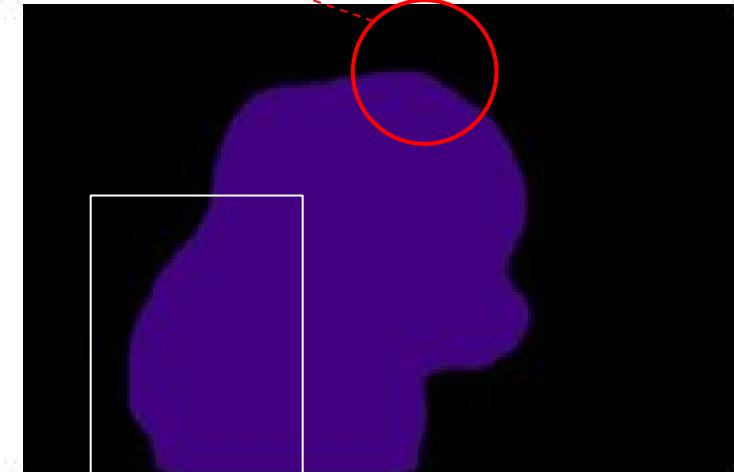
The Details are lost!



Image

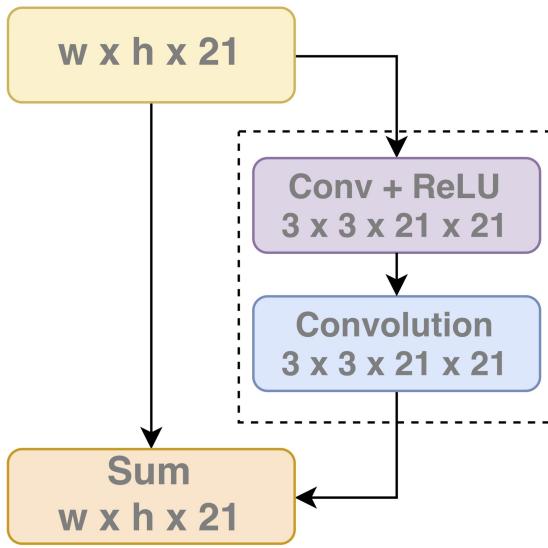


Baseline (FCN)

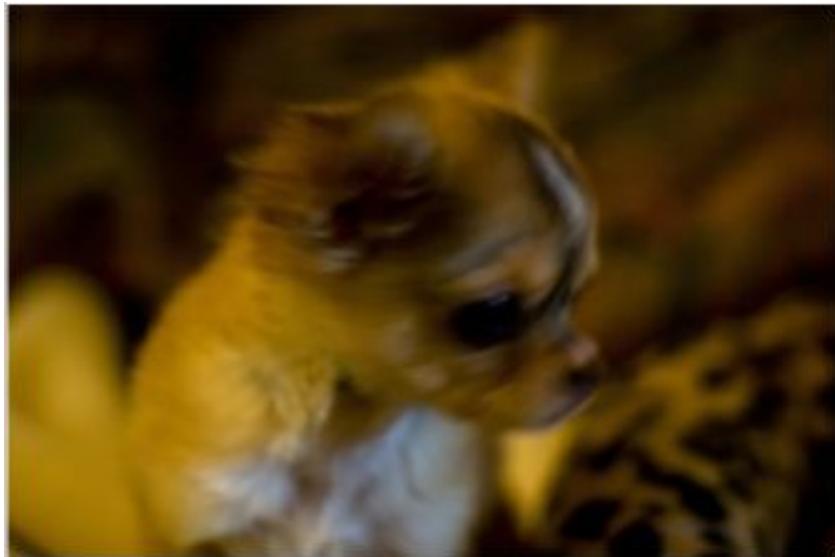


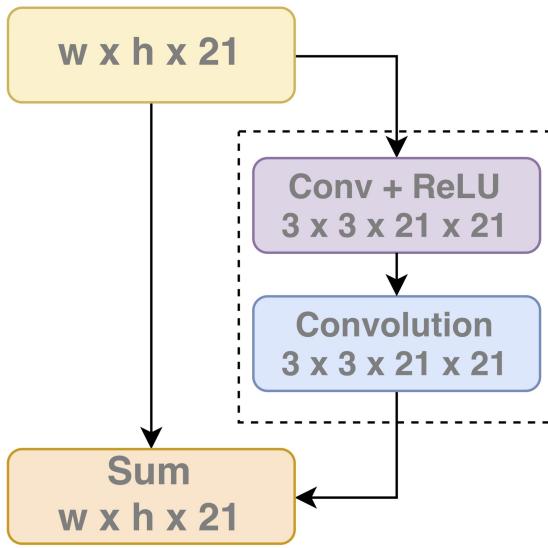
GCN

Region Mis-Classifications are corrected!

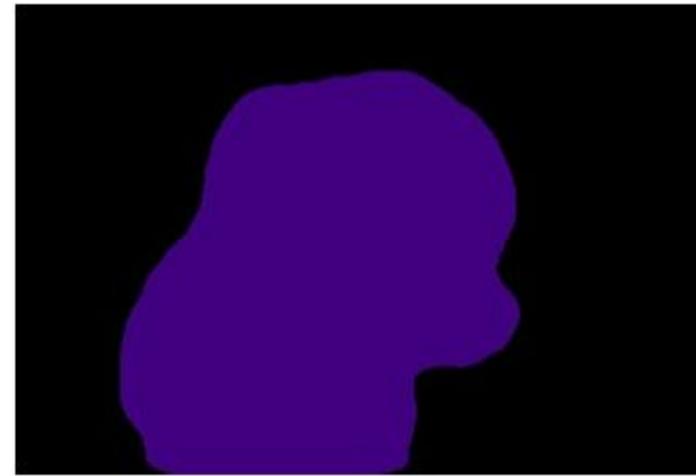
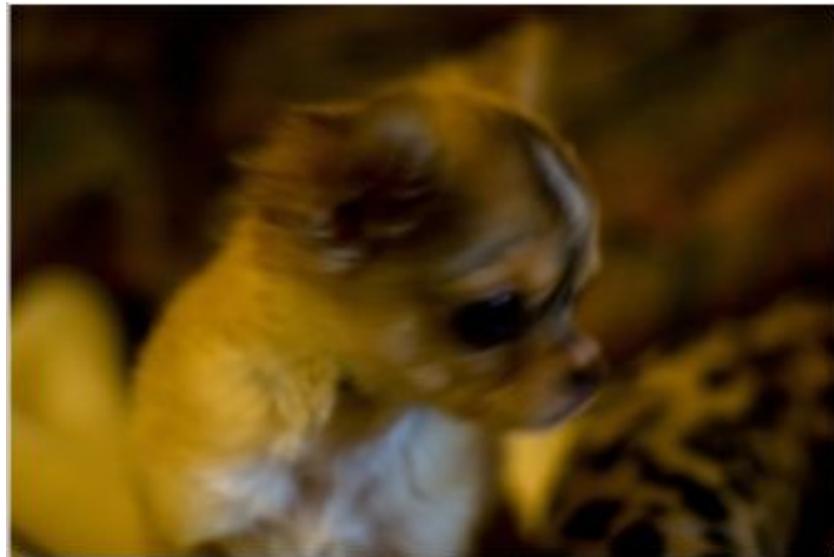


Boundary Refinement (BR)

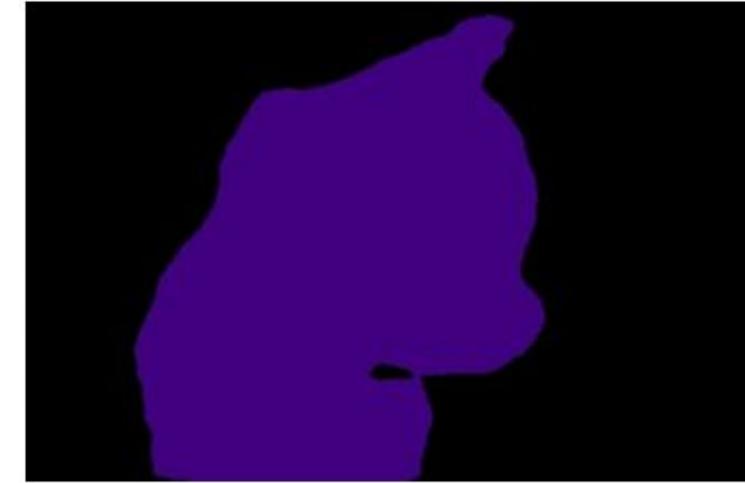




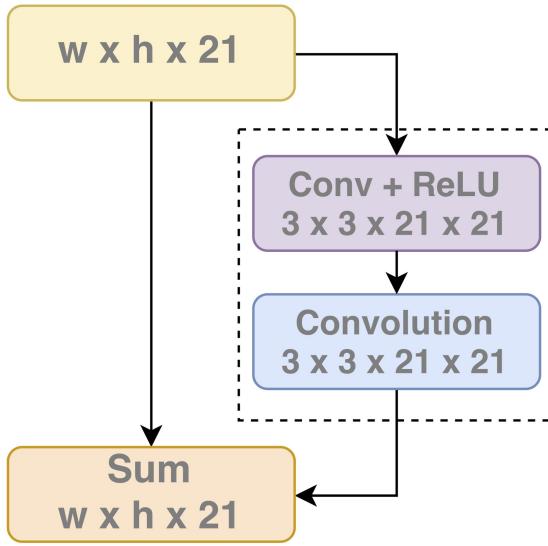
Boundary Refinement (BR)



GCN



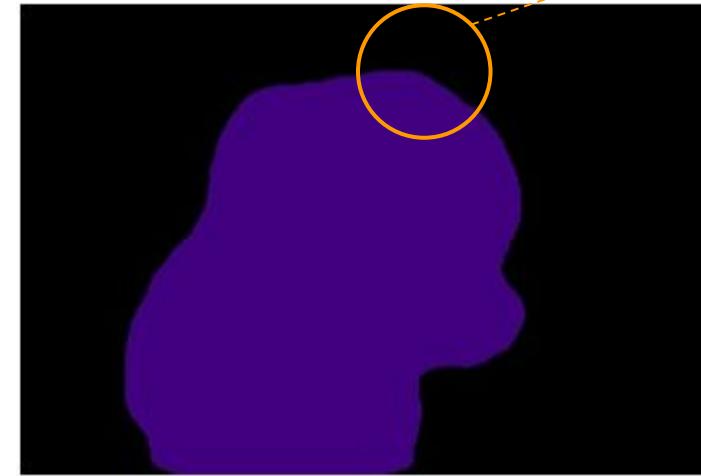
GCN + BR



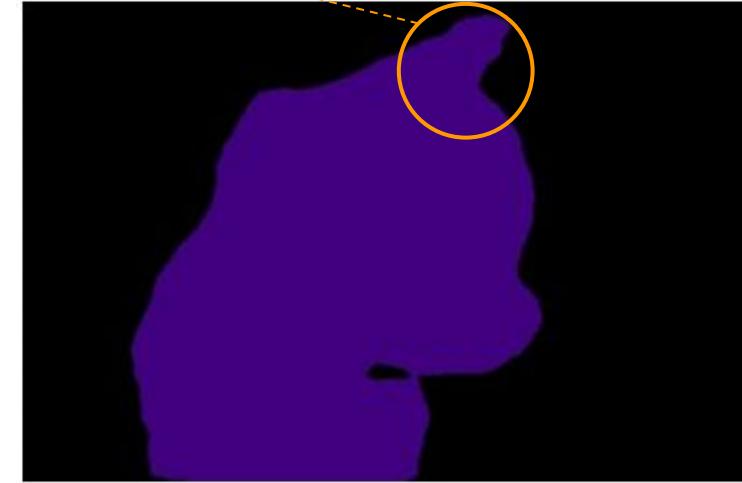
Boundary Refinement (BR)



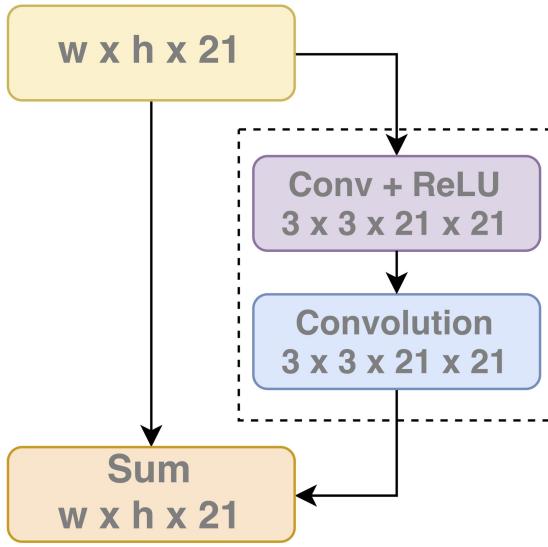
The Details are recovered!



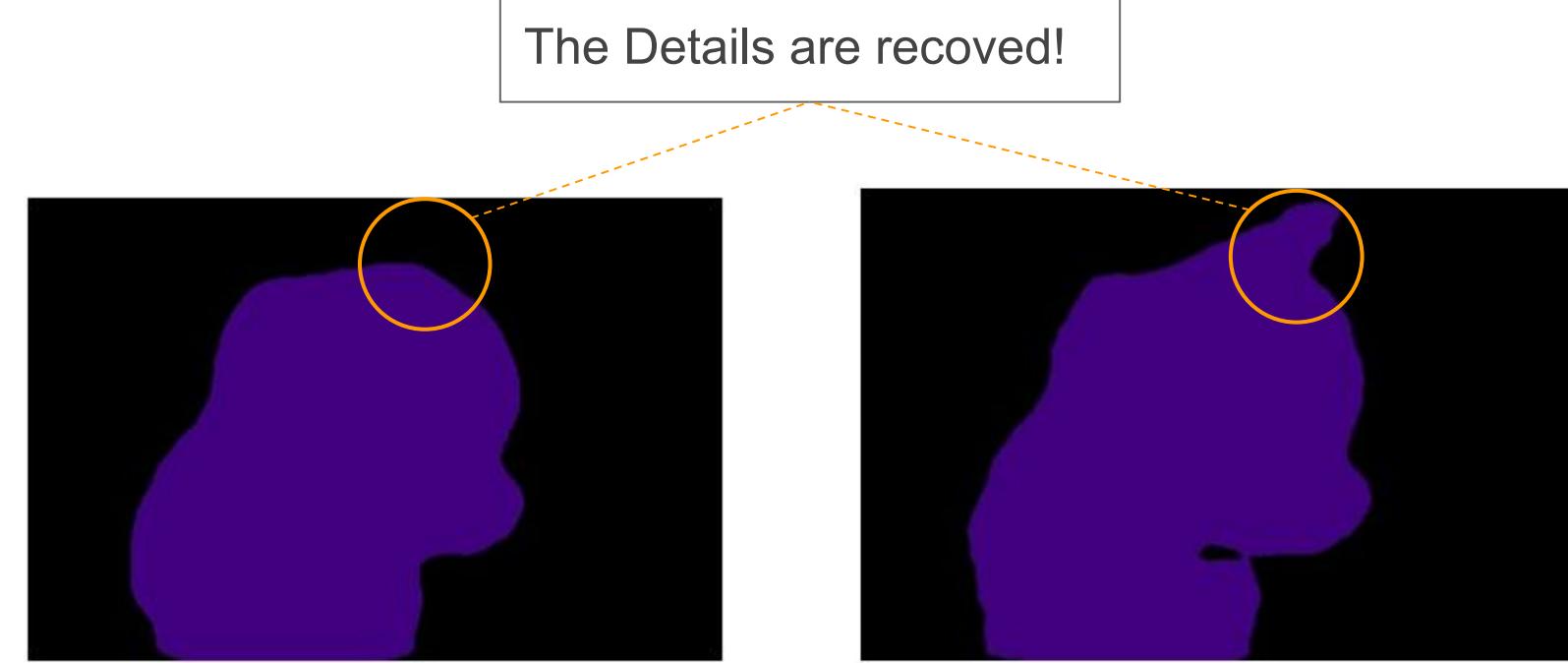
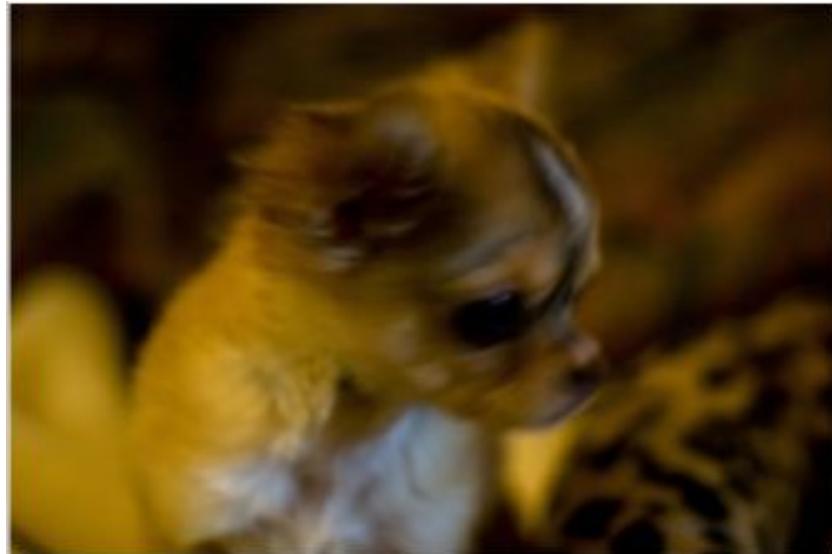
GCN

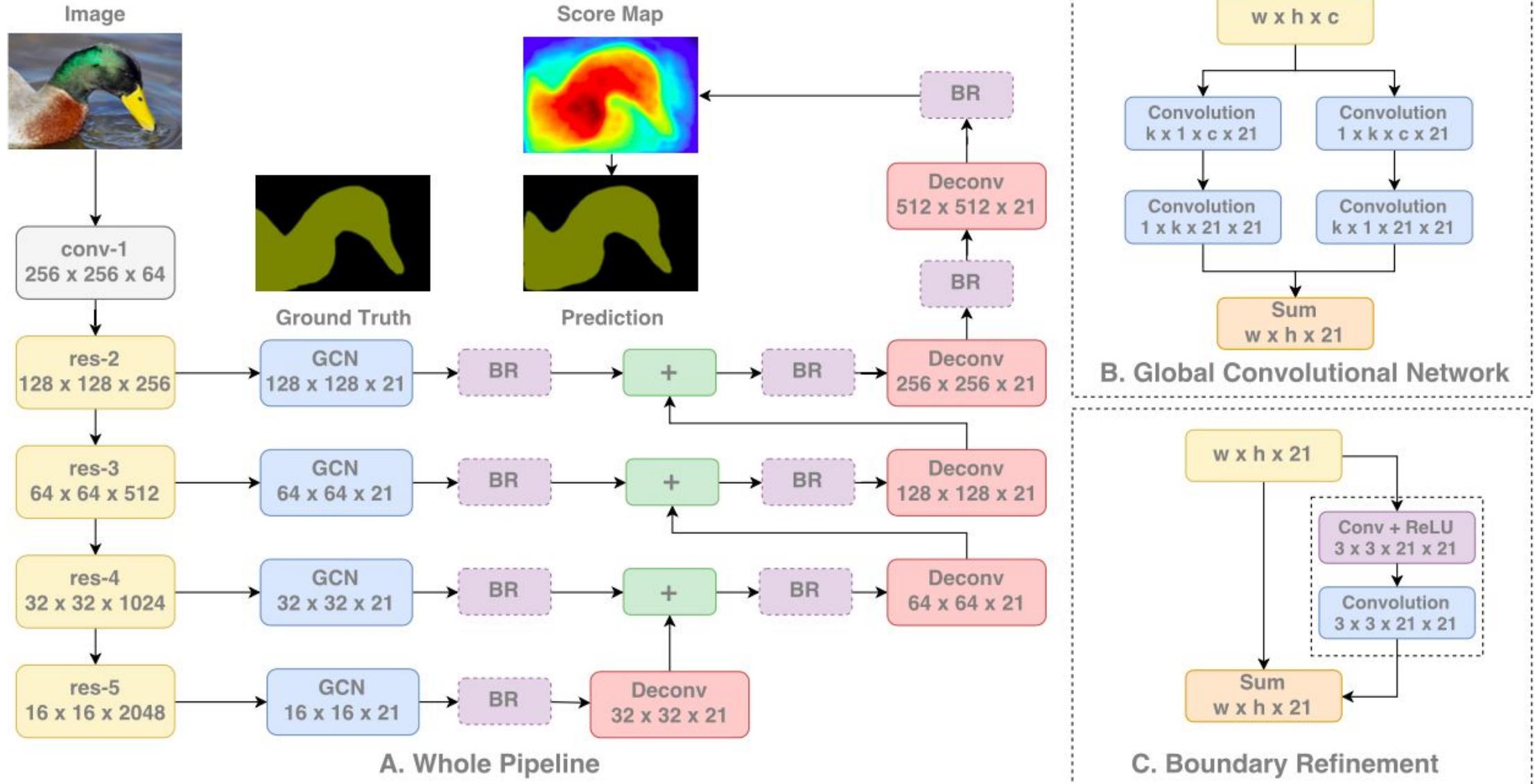


GCN + BR



Boundary Refinement (BR)

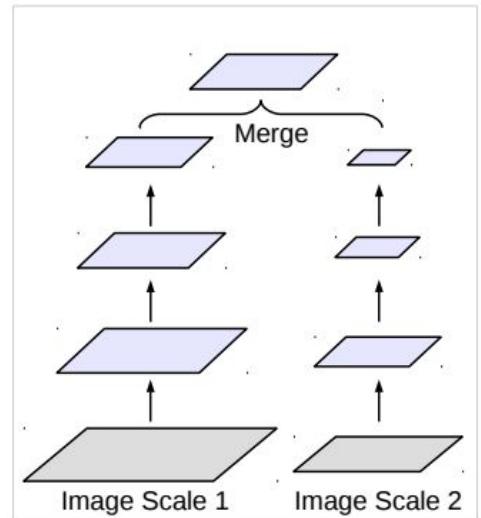




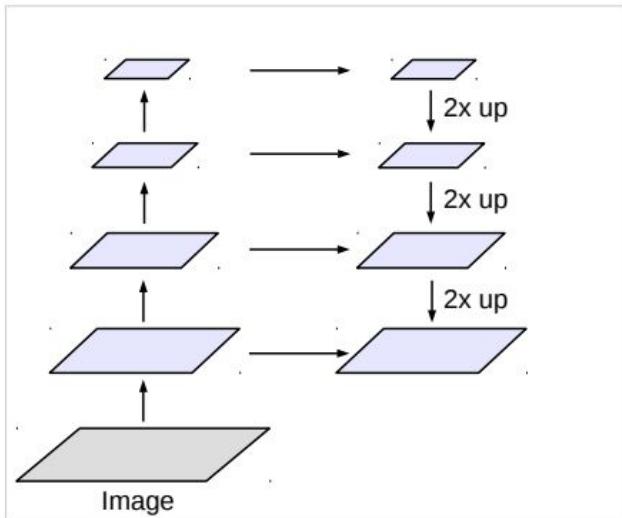
Global Convolutional Network

- Extend the FCN framework!
- Partially Solve the Receptive Field Problem!
- Two key components! (GCN and BRN)

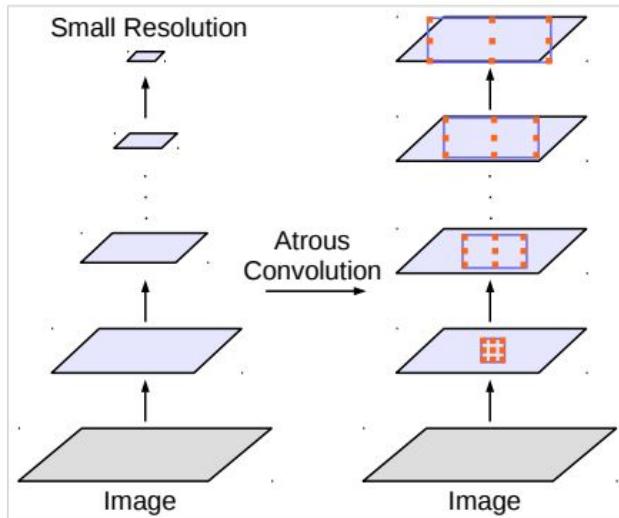
Deeplab V3



(a) Image Pyramid

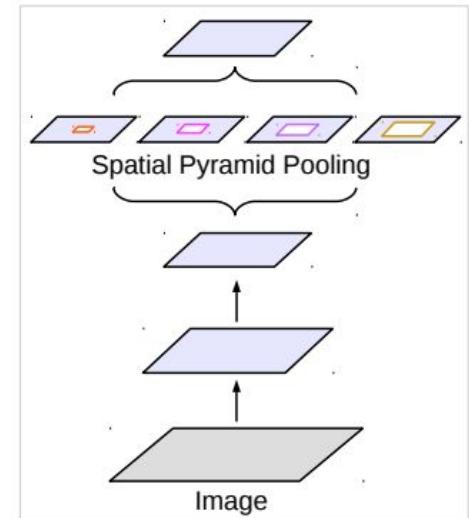


(b) Encoder-Decoder



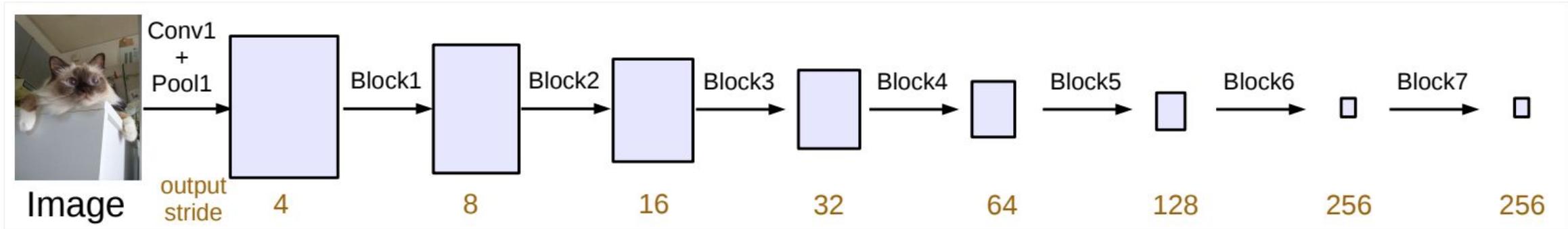
(c) Deeper w. Atrous Convolution

Figure 2. Alternative architectures to capture multi-scale context.

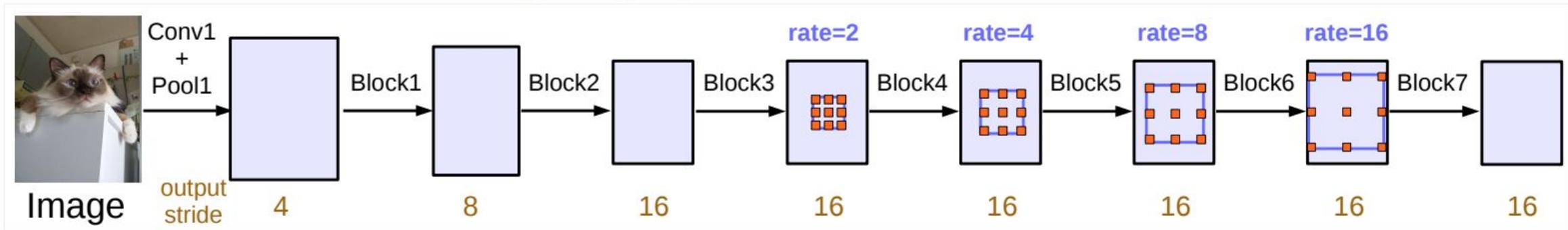


(d) Spatial Pyramid Pooling

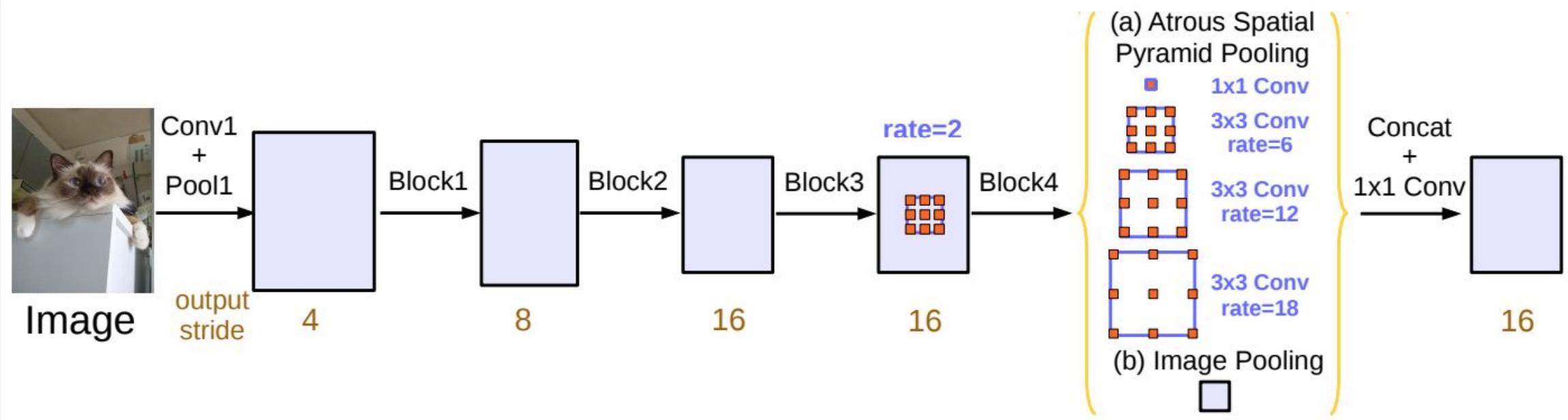
Deeplab V3



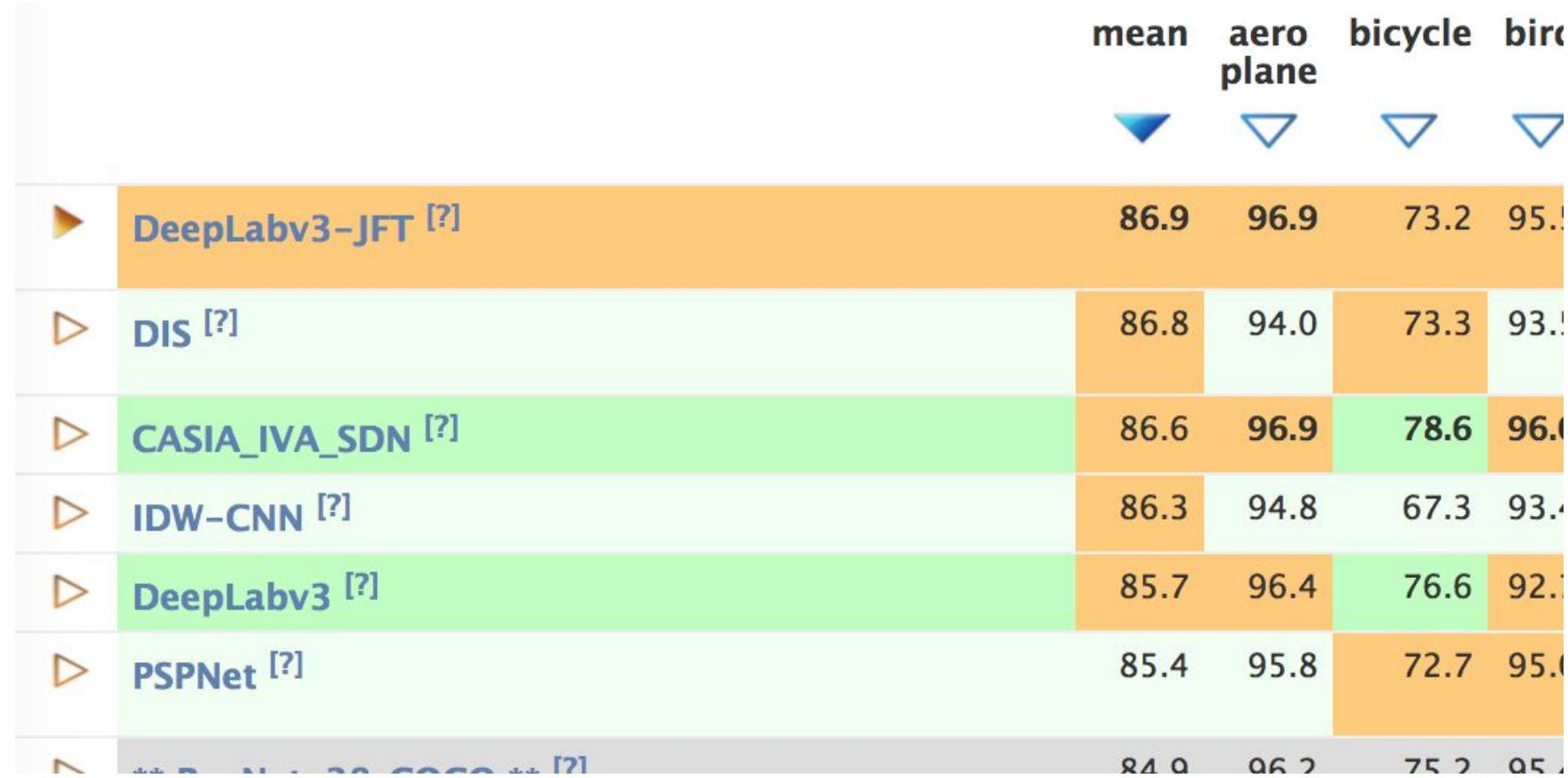
(a) Going deeper without atrous convolution.



Deeplab V3



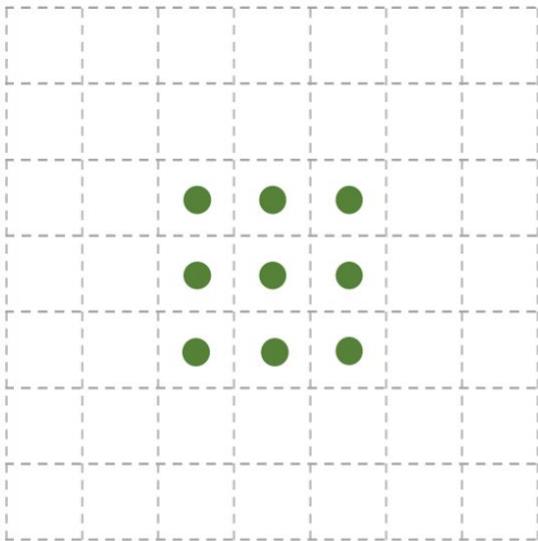
Deeplab V3



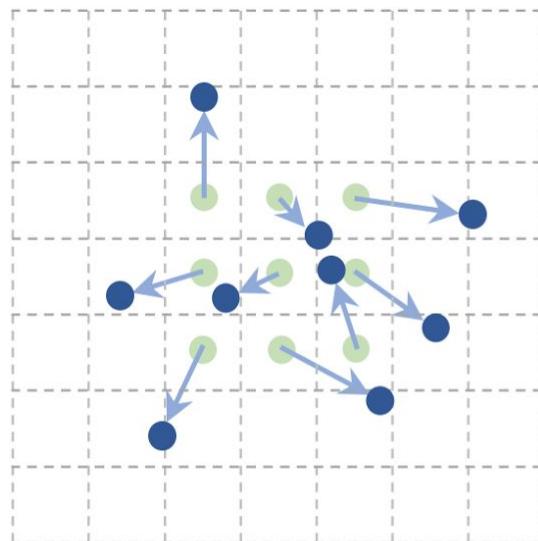
Deeplab V3

- Currently State-Of-Art on PASCAL VOC 2012
- Conclude the dilate-convolution technique on segmentation

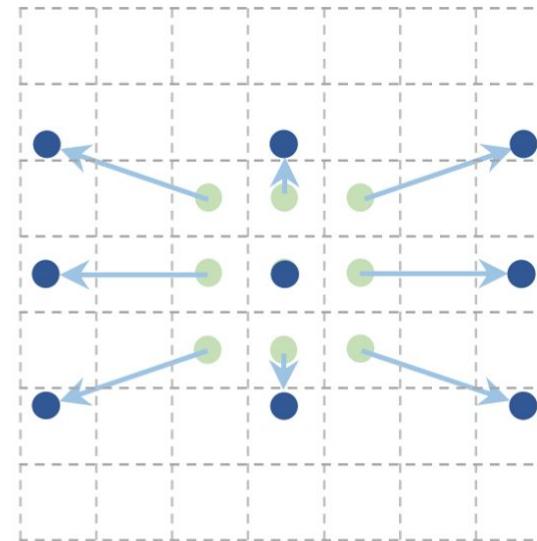
Deformable Convolution



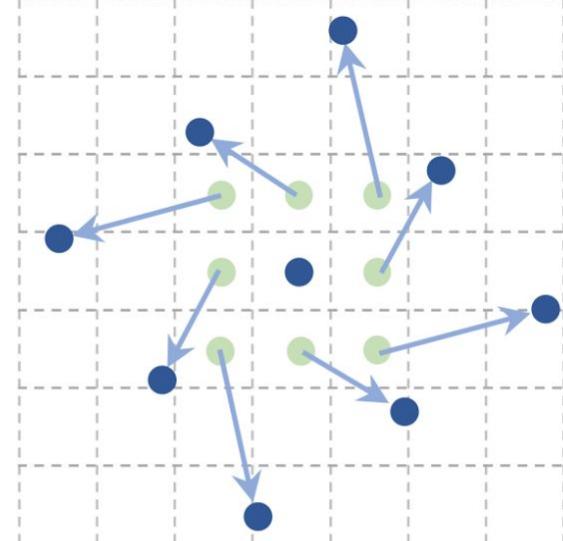
(a)



(b)

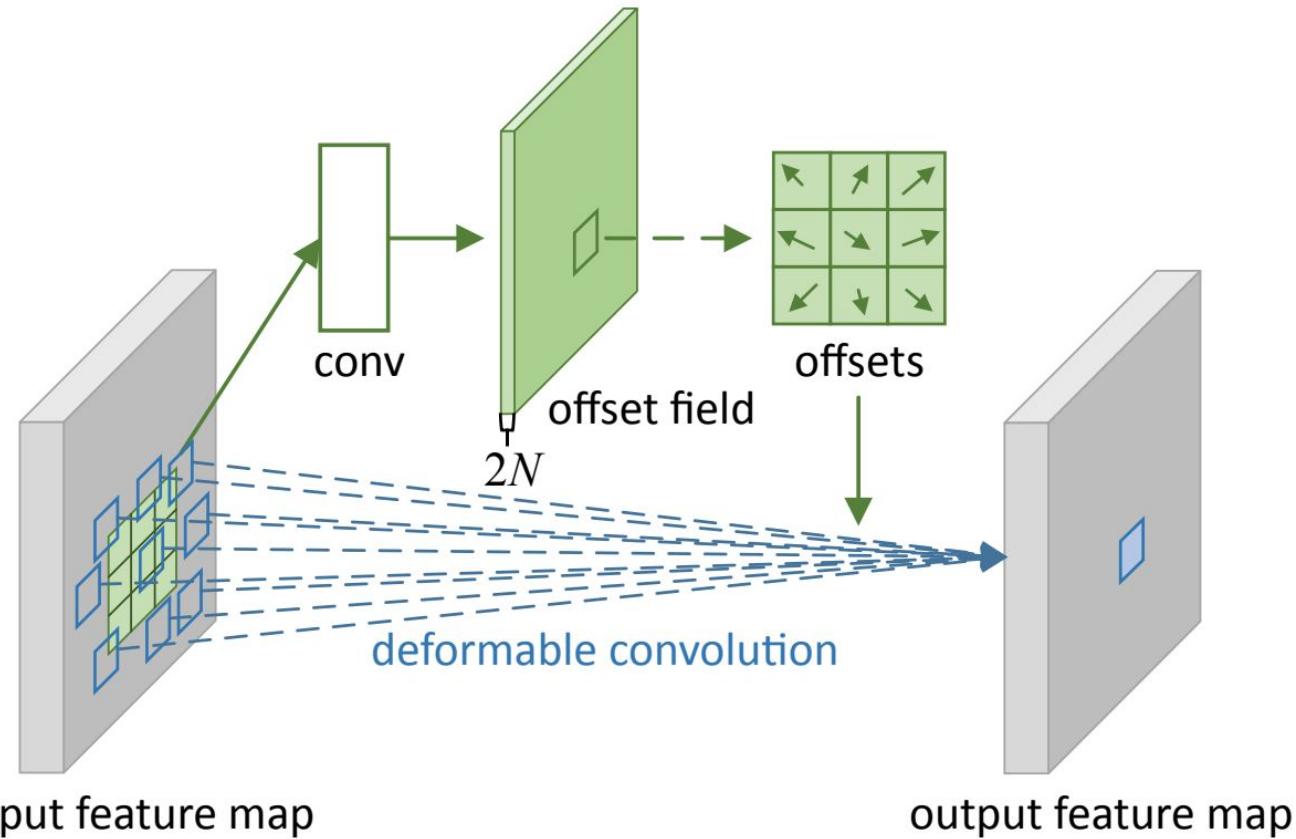


(c)

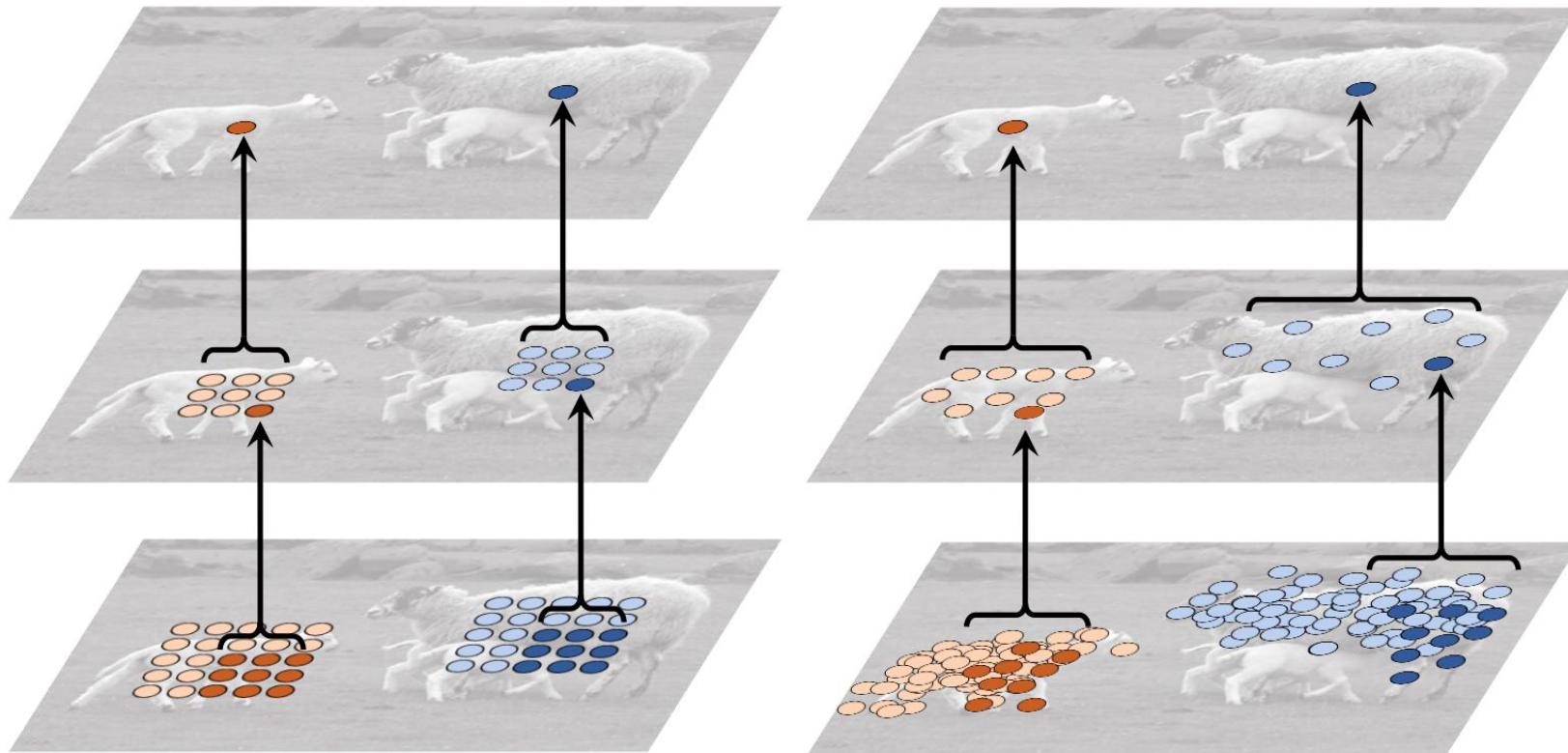


(d)

Deformable Convolution



Deformable Convolution



Deformable Convolution



Deformable Convolution

- Solve the receptive field problem using learned offsets!
- Also valid for detection!

Re-Cap

- Segmentation with CNN: FCN, Deeplab, GCN ...
- Segmentation with CRF: DenseCRF, CRFAsRNN, ...
- Different Convolutions: Dilated Conv, Global Conv, Deformable, ...

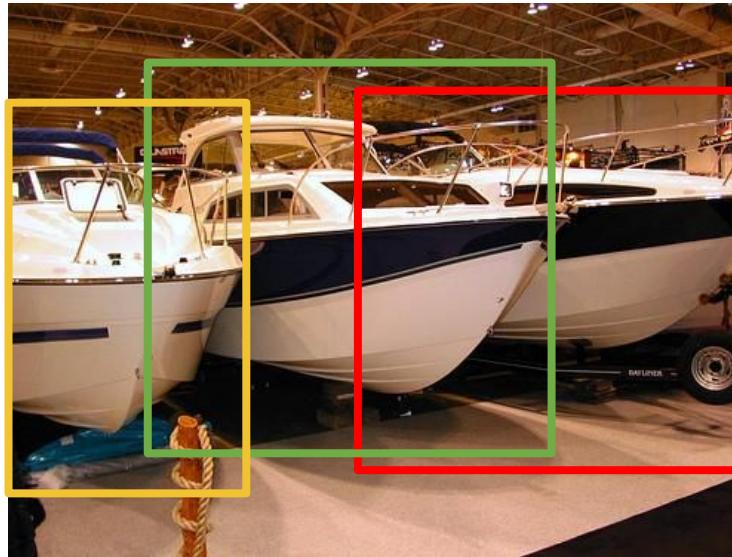
Outline

- Semantic Segmentation
- Instance Segmentation

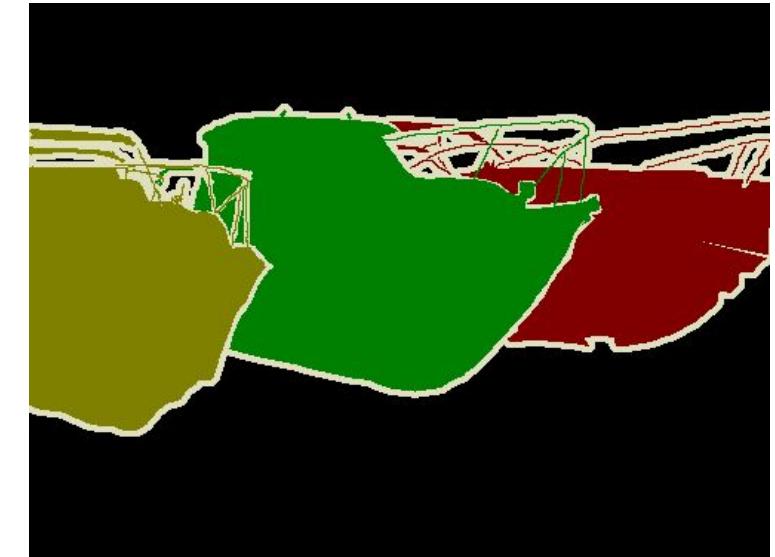
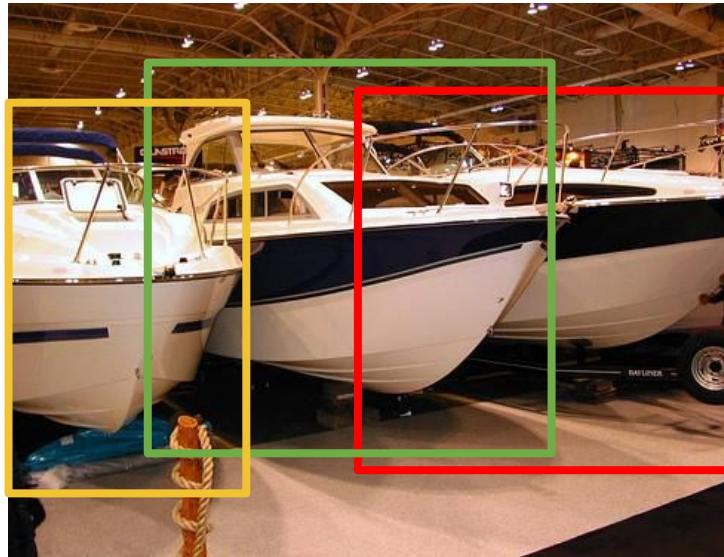
Top-Down Pipeline



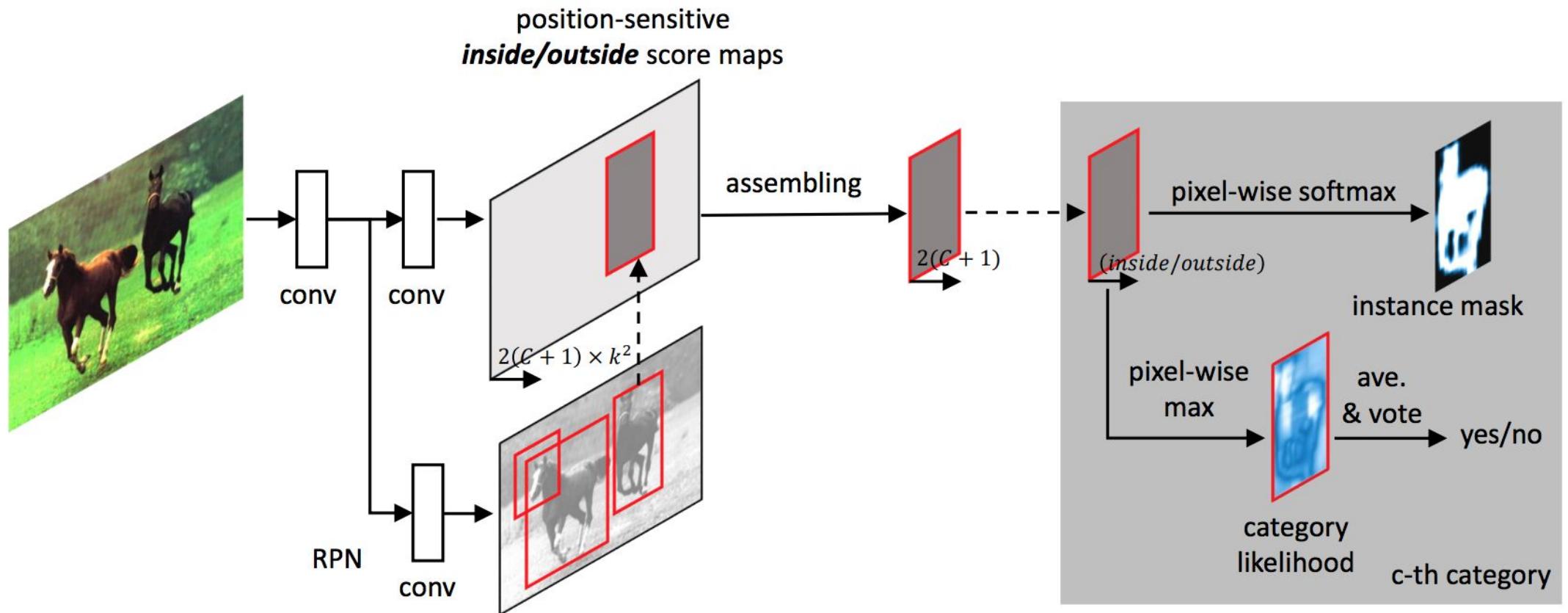
Top-Down Pipeline



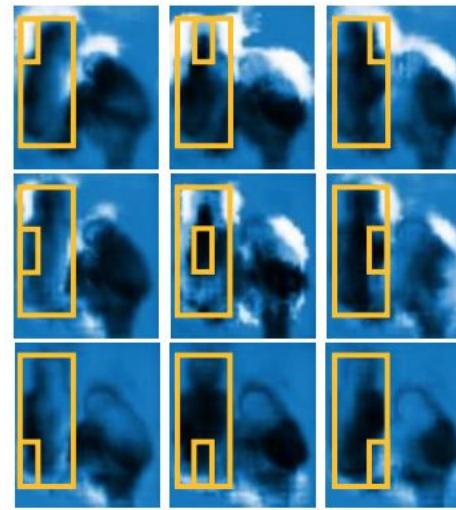
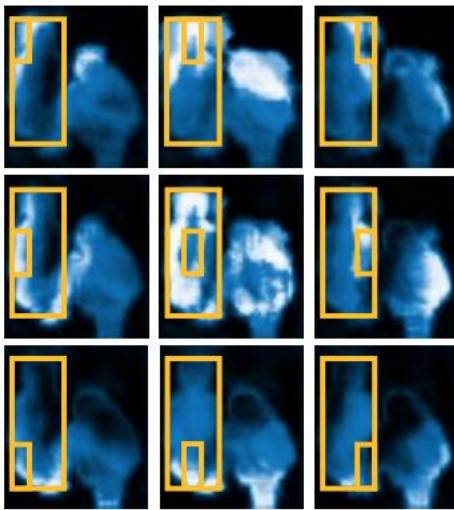
Top-Down Pipeline



FCIS

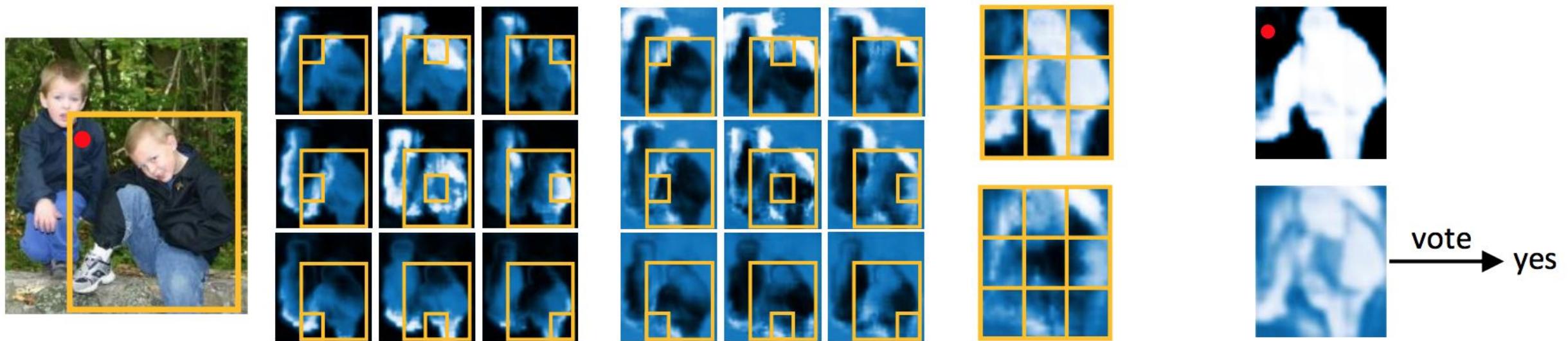


FCIS

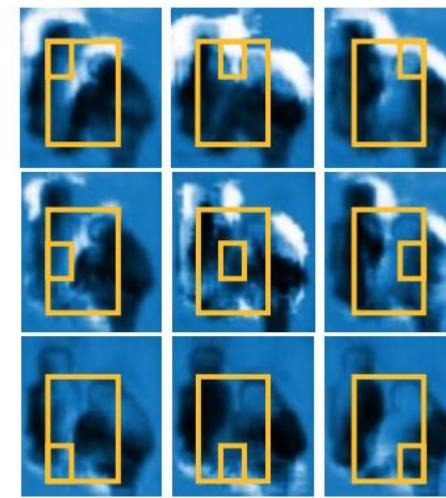
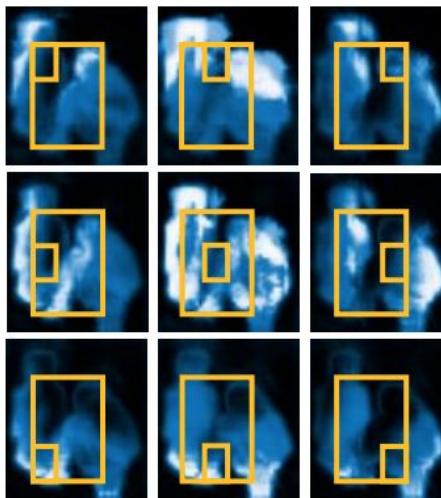
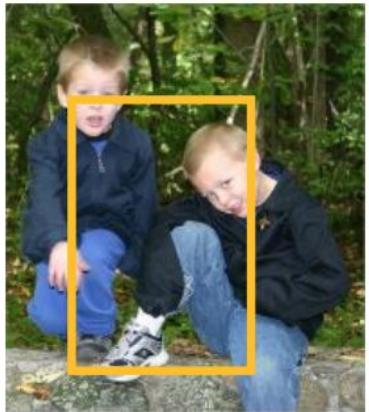


vote → yes

FCIS



FCIS

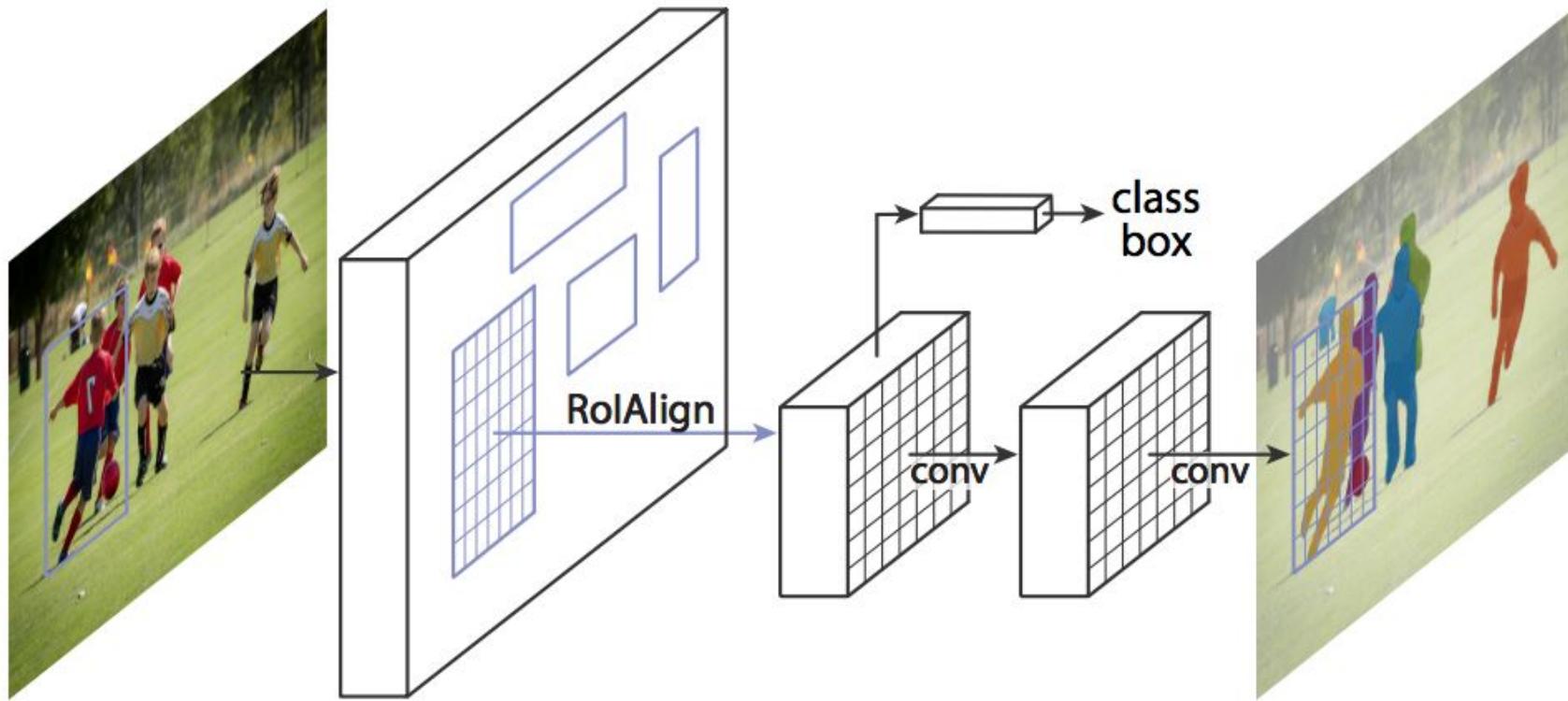


vote → no

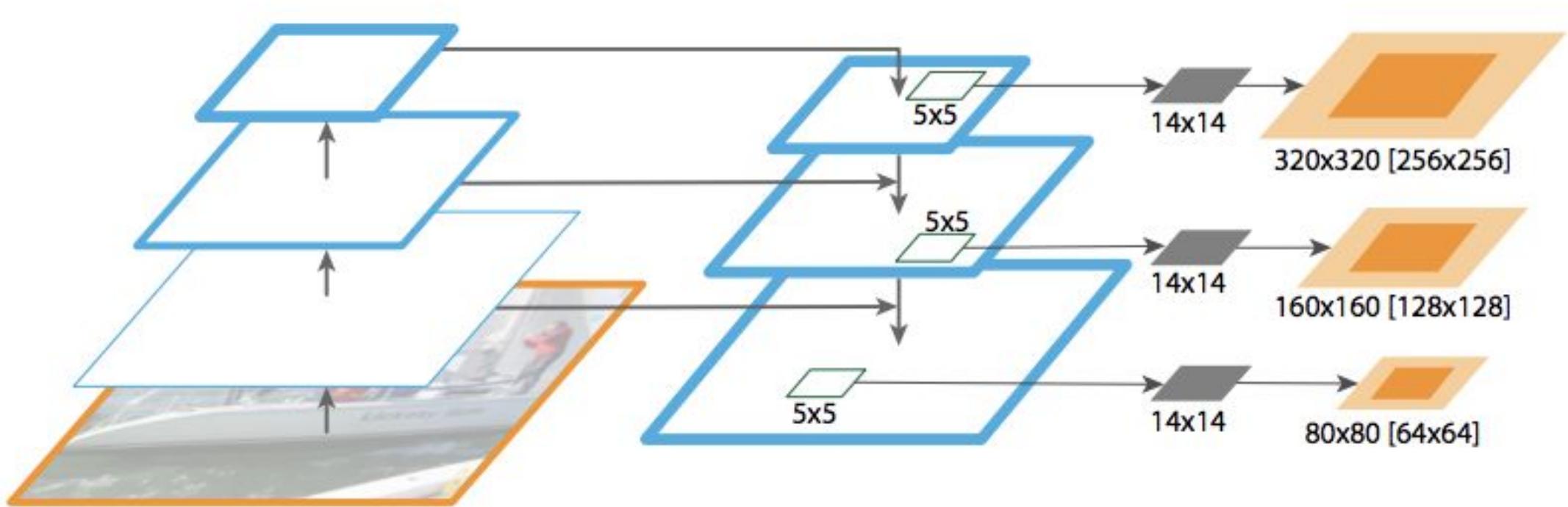
FCIS

- Bases on R-FCN framework
- Propose the inside-outside scoring technique!
- Winner of COCO InstanceSeg Challenge 2016!

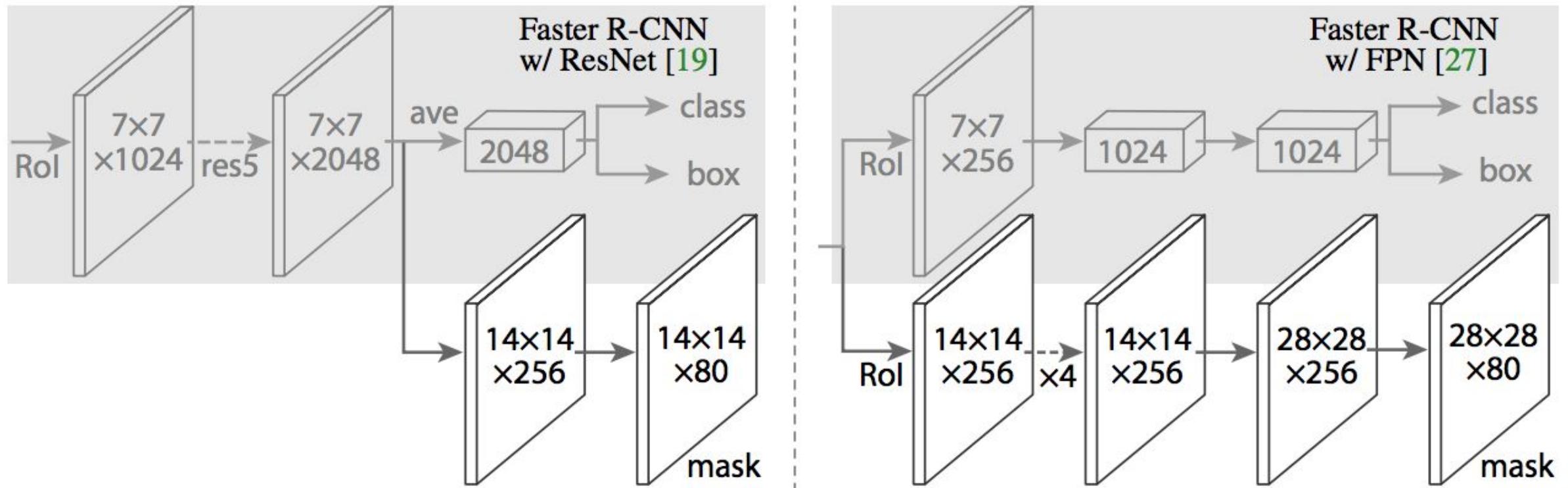
Mask RCNN



Mask RCNN

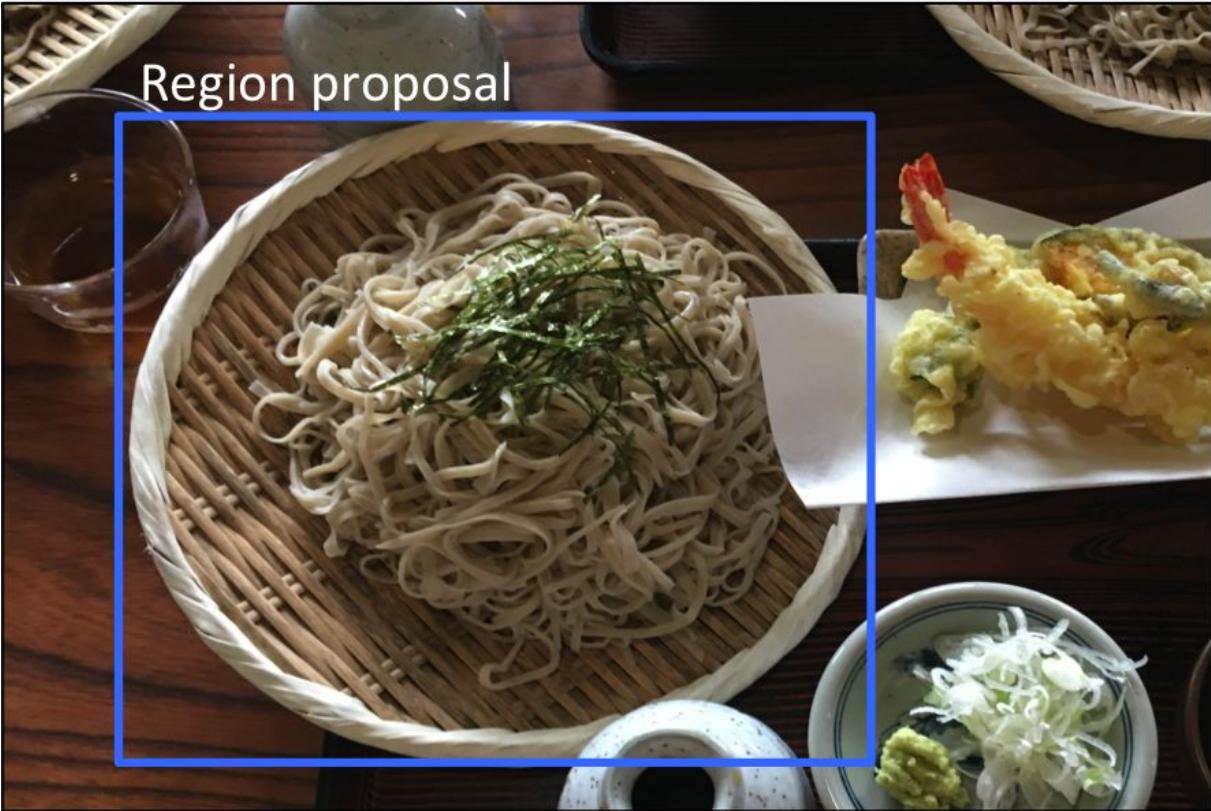


Mask RCNN



Mask RCNN

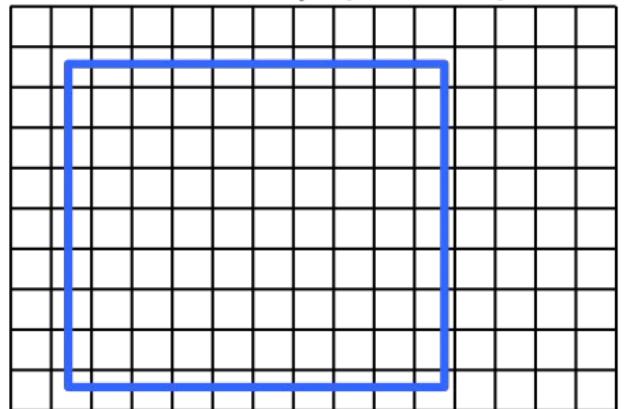
Input image (320x480)



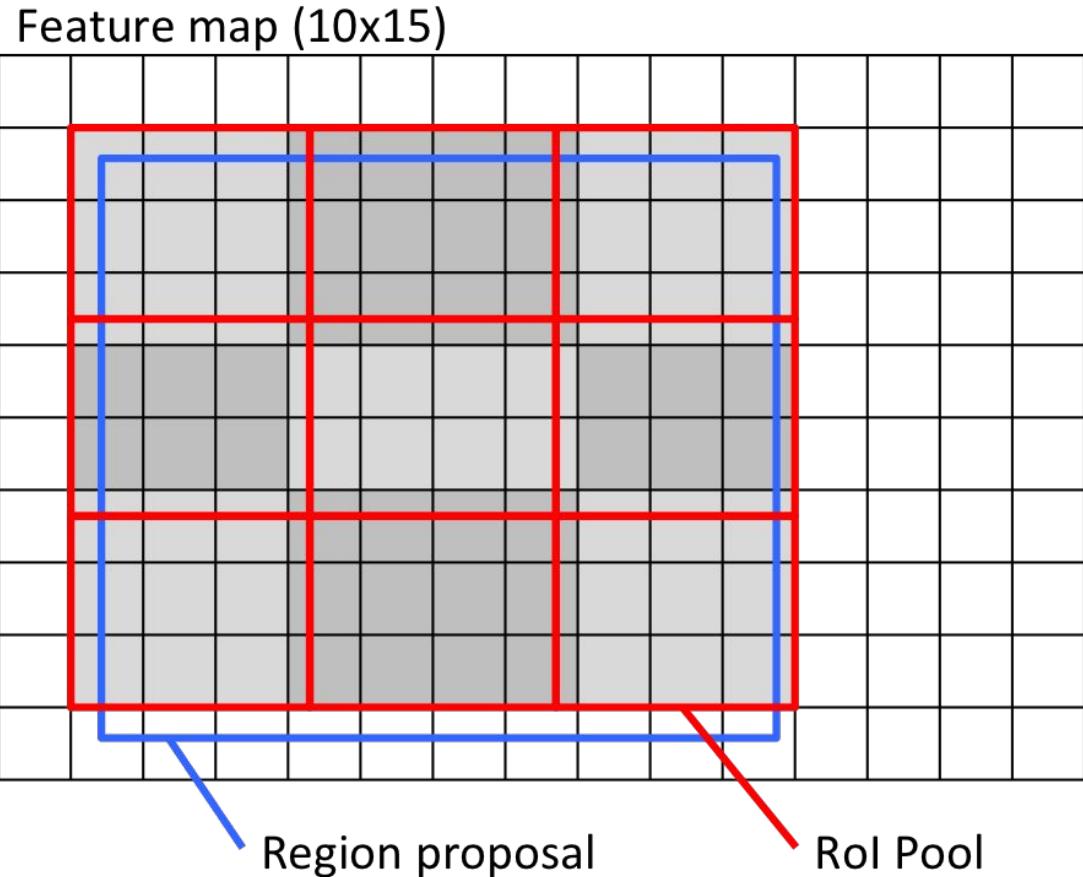
Region proposal

CNN (stride = 32)

Feature map (10x15)

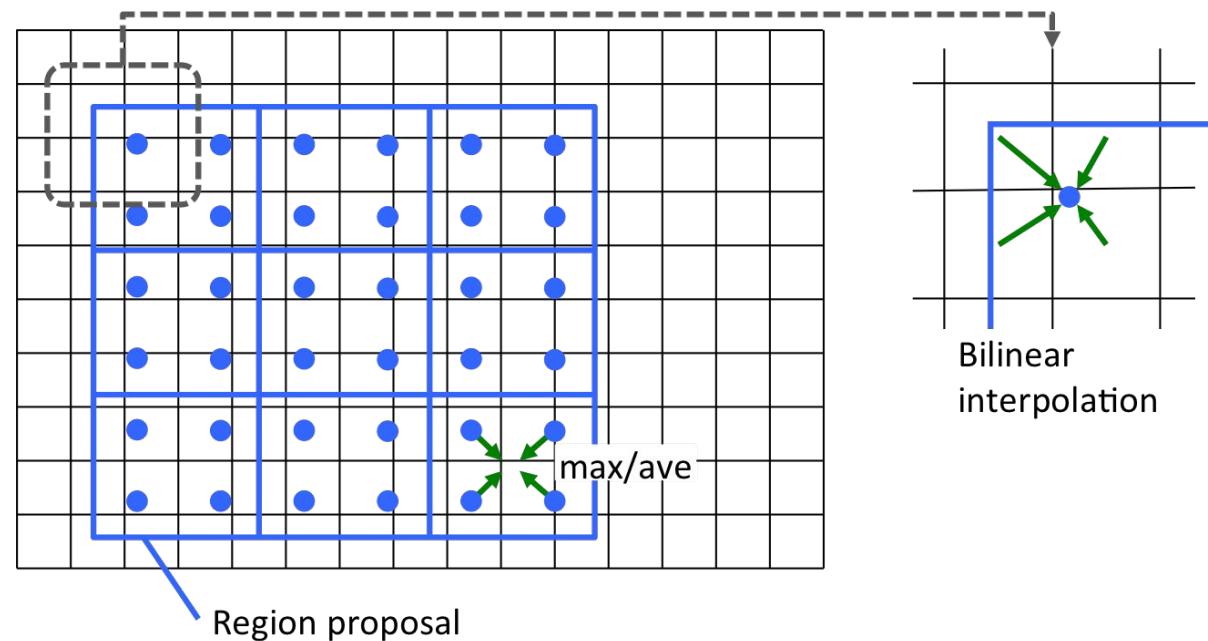
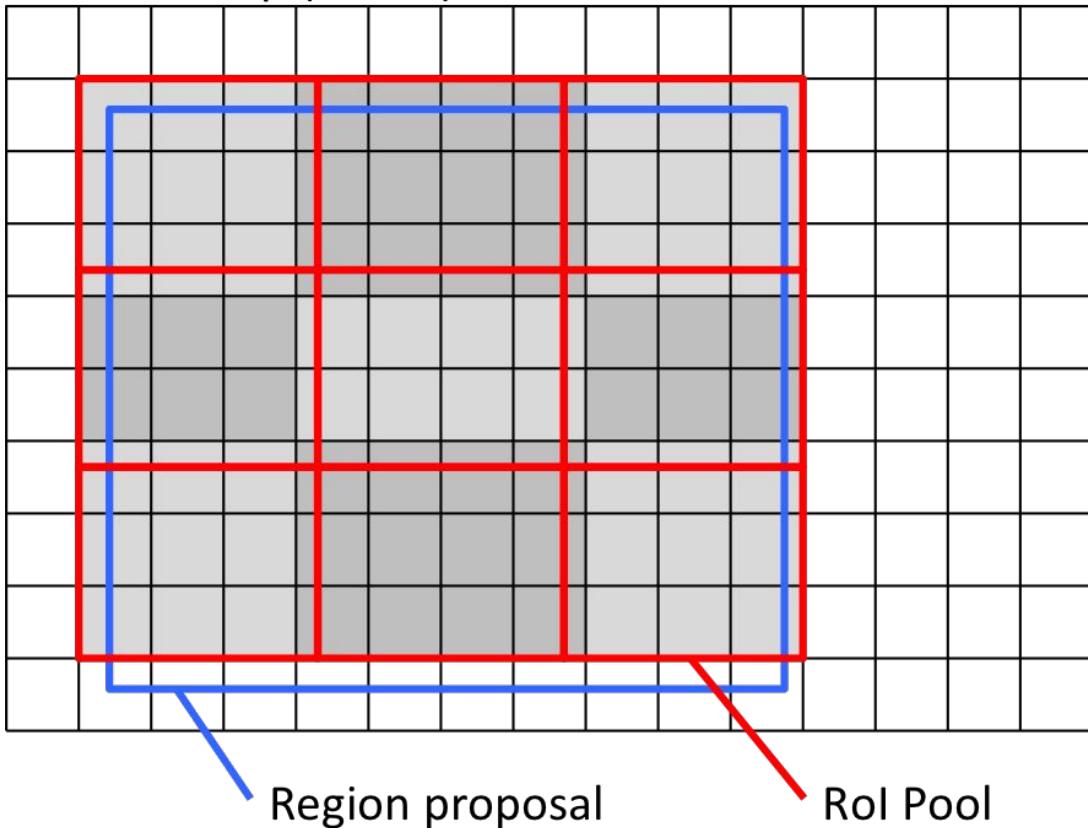


Mask RCNN



Mask RCNN

Feature map (10x15)



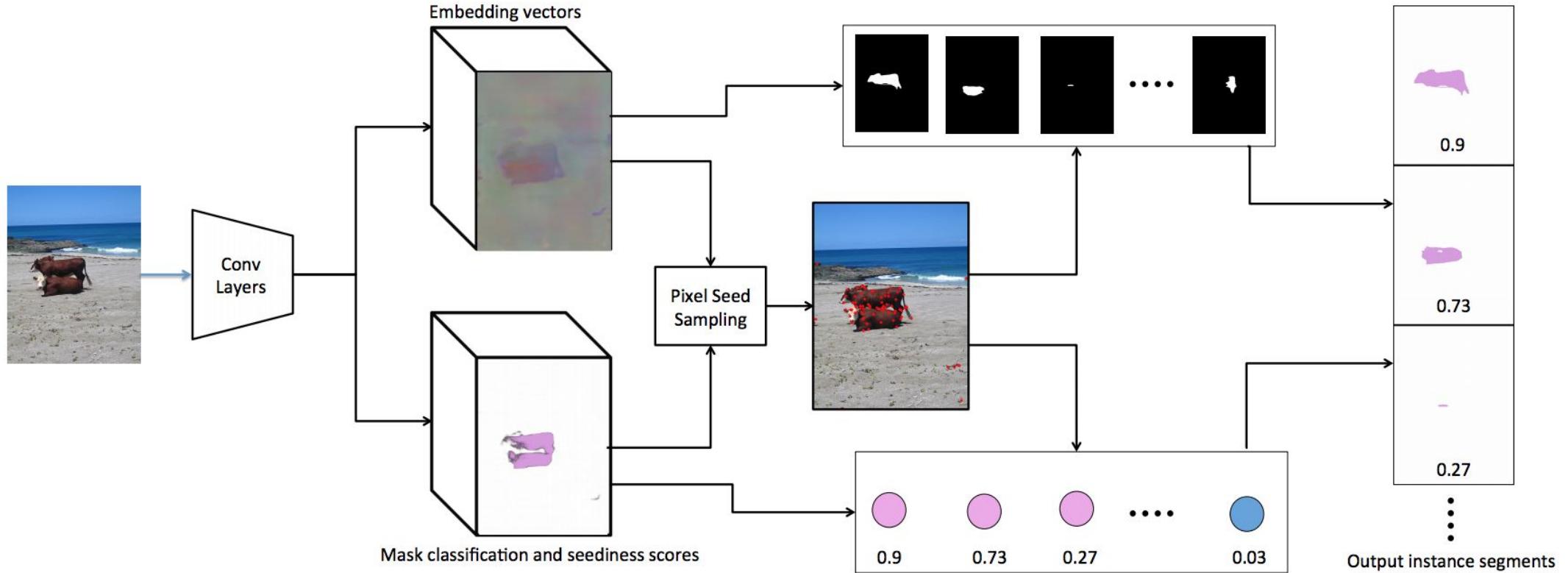
Mask RCNN



Mask RCNN

- State-Of-Art framework on InstanceSeg
- Best paper of ICCV 2017!
- Simple framework!

Bottom-up Pipeline

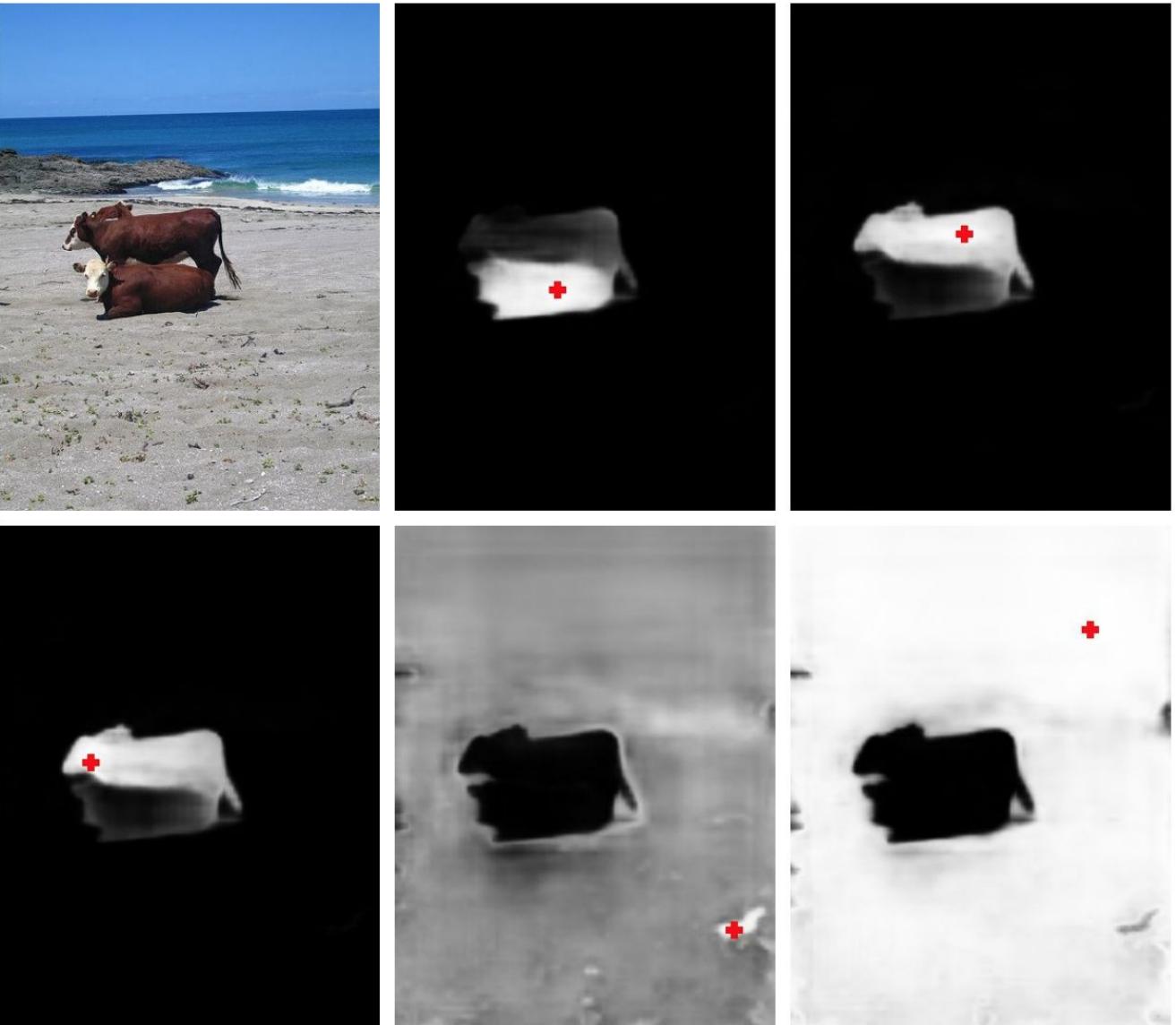


Bottom-up Pipeline

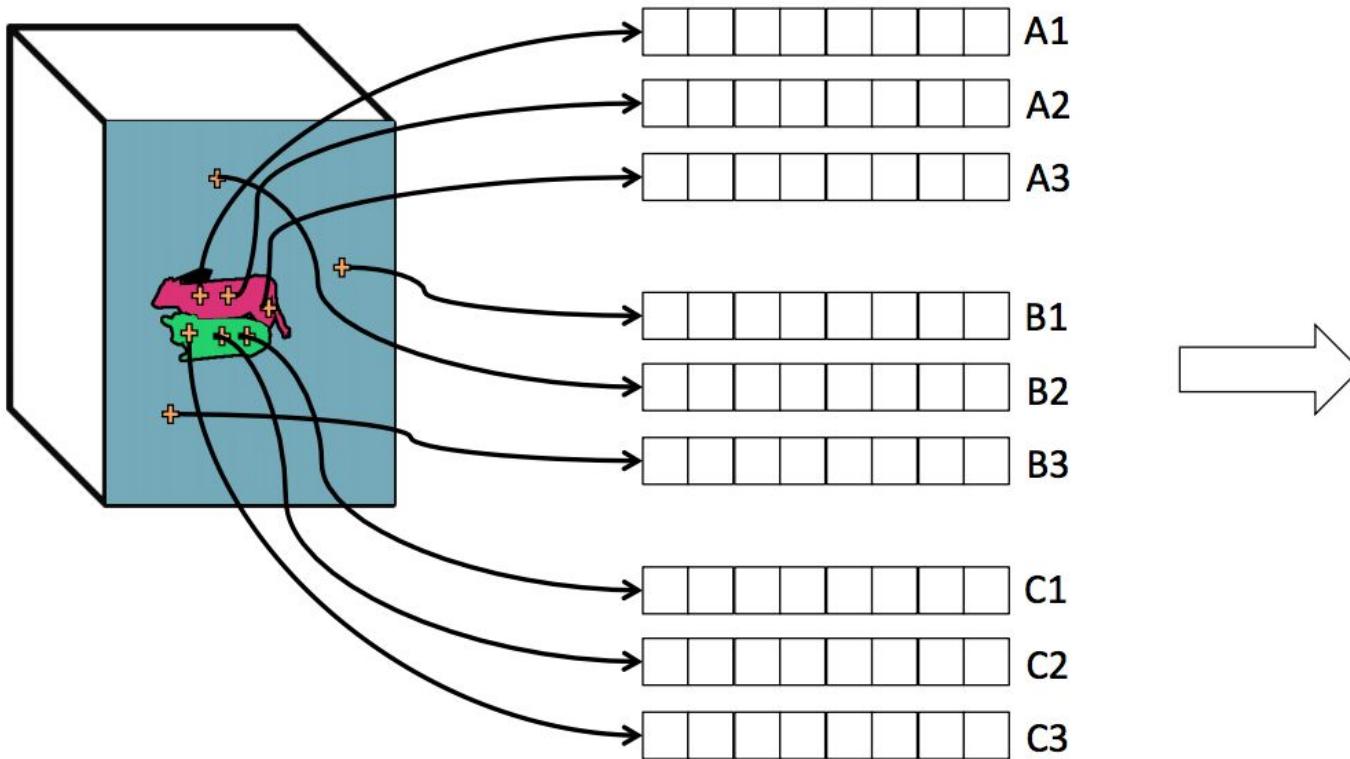
$$p_t = \arg \max_{p \notin p_{1:t-1}} [\log(S_p) + \alpha \log(D(p, p_{1:t-1}))]$$

where

$$D(p, p_{1:t-1}) = \min_{q \in p_{1:t-1}} \|e_p - e_q\|_2^2$$



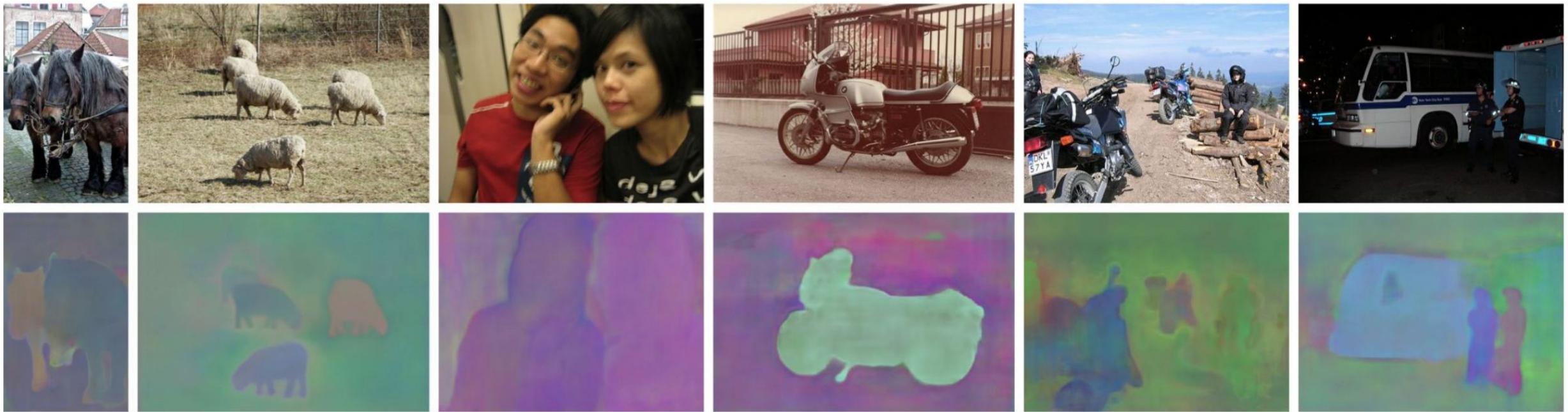
Bottom-up Pipeline



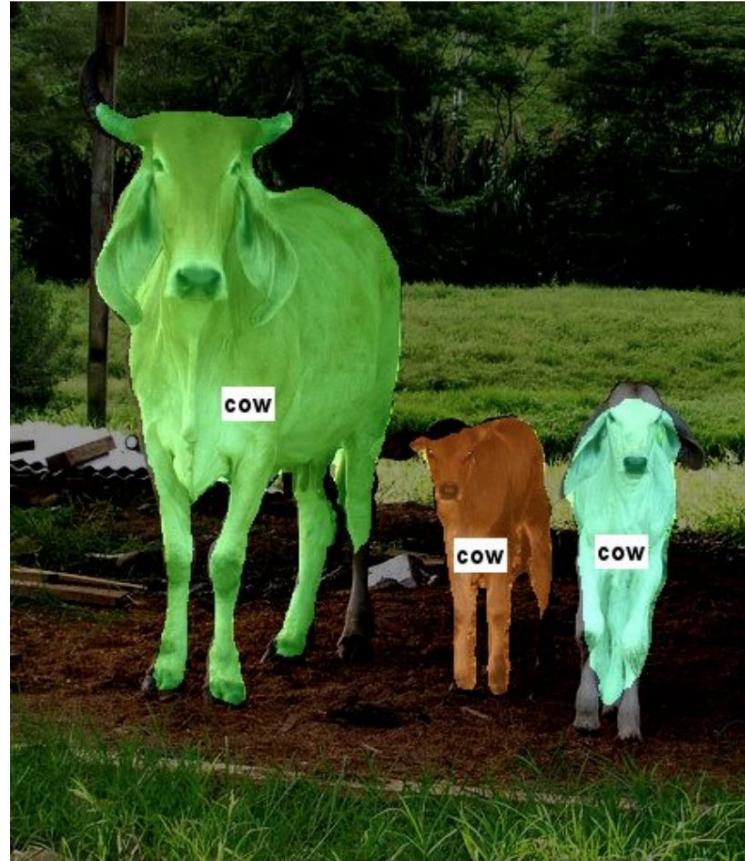
	A1	A2	A3	B1	B2	B3	C1	C2	C3
A1	█	█	█						
A2	█	█	█						
A3	█	█	█						
B1				█	█	█			
B2				█	█	█			
B3				█	█	█			
C1							█	█	█
C2							█	█	█
C3							█	█	█

Sigmoid cross entropy loss on the similarity
between pairs of embedding vectors

Bottom-up Pipeline



Bottom-up Pipeline



Bottom-up Pipeline

- Alternative framework to InstanceSeg
- Tricky to implement
- Incorporating the metric learning

Top-Down

- Instance ---> Segmentation
(for each instance)
- MainStream
- Start-Of-Art
- Easy to implement
- Difficulty: shrink the gap
between det and seg

Bottom-Up

- Segmentation (for image)
---> instance
- Alternative
- Sub-Optimal
- Tricky to implement
- Difficulty: generate better
instance

Re-Cap

- Segmentation with CNN: FCN, Deeplab, GCN ...
- Segmentation with CRF: DenseCRF, CRFAsRNN, ...
- Different Convolutions: Dilated Conv, Global Conv, Deformable, ...
- Top-Down pipeline for Instance Segmentation: FCIS, Mask-RCNN
- Bottom-Up pipeline

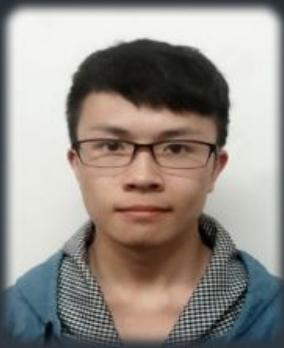
COCO & Places Challenge 2017



Chao PENG*



Tete XIAO*



Zeming LI*



Yuning JIANG



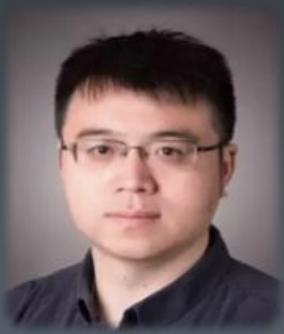
Xiangyu ZHANG



Kai JIA

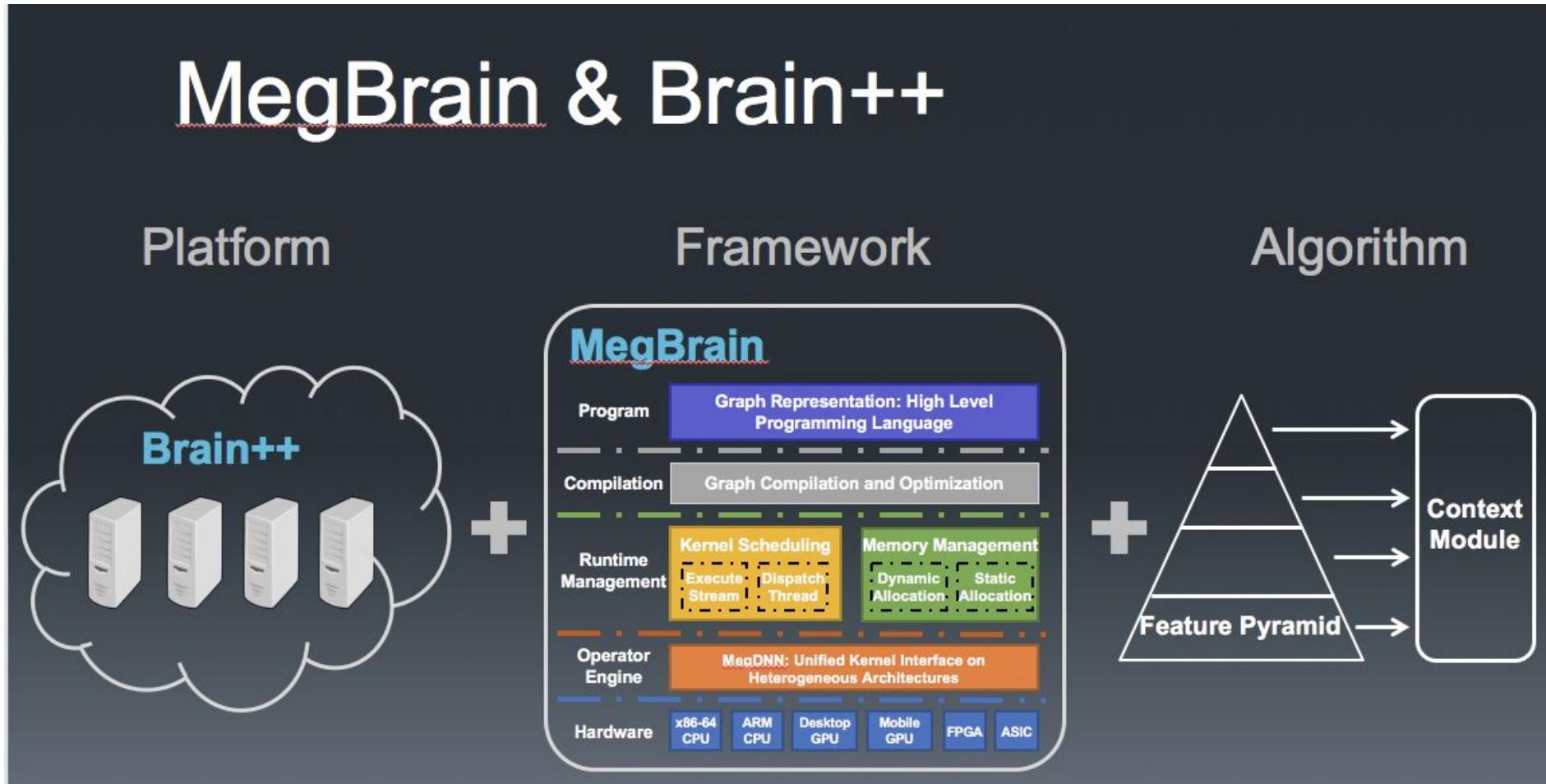


Gang YU



Jian SUN

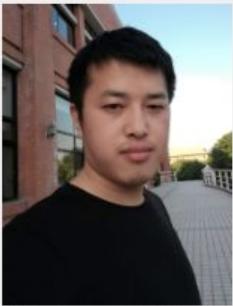
COCO & Places Challenge 2017



COCO & Places Challenge 2017



Yilun Chen*



Zhicheng Wang*



Xiangyu Peng



Zhiqiang Zhang



Gang Yu



Chao Peng



Tete Xiao



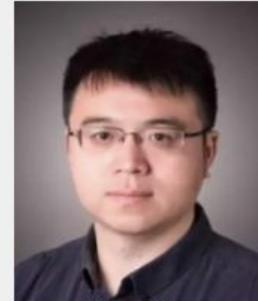
Zeming Li



Xiangyu Zhang

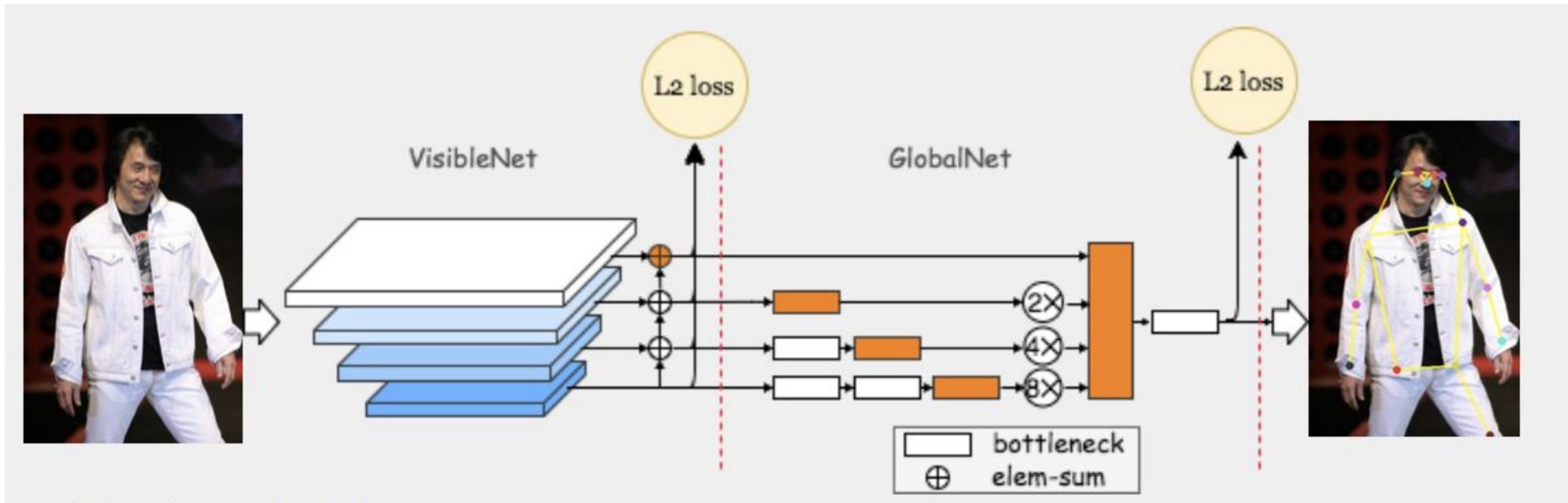


Yuning Jiang



Jian Sun

COCO & Places Challenge 2017



COCO & Places Challenge 2017

Track	Rank	Ensemble	Single
COCO BBox Detection	1 st	52.8	50.5
Places InstanceSeg	1 st	30.7	28.7
COCO Keypoint	1 st	72.6	70.9
COCO InstanceSeg	2 nd	46.4	45.0

COCO Challenge 2017 BBOX

	Megvii (Face++)	0.526	0.730	0.585	0.343	0.556
	UCenter	0.510	0.705	0.558	0.326	0.539
	MSRA	0.507	0.717	0.566	0.343	0.529
	FAIR Mask R-CNN	0.503	0.720	0.558	0.328	0.537

COCO Challenge 2017 BBOX

		AP	AR@10	AR@50	AR@100	AR@300
	Meggvii (Face++)	0.526	0.730	0.585	0.343	0.556
	UCenter	0.510	0.705	0.558	0.326	0.539
	MSRA	0.507	0.717	0.566	0.343	0.529
	FAIR Mask R-CNN	0.503	0.720	0.558	0.328	0.537



Our Single Model is Here: 50.5.

Places Challenge 2017 InstanceSeg

Team Name	mean AP
Megvii (Face++)	0.2977
G-RMI	0.2415
BlueSky	0.1551
...	...
baseline	0.200

Megvii (Face++)

*Tete Xiao, *Ruixuan Luo, *Borui Jiang,
Shuai Shao, Yuning Jiang, Yadong Mu,
Jieqi Shi, Chi Zhang, Jian Sun

Megvii Research and Peking University

G-RMI

Alireza Fathi, Nori Kanazawa, Kai Yang,
Kevin Murphy
Google Research Machine Intelligence

COCO Challenge 2017 Keypoint

	Megvii (Face++)	0.721	0.905	0.789	0.679
 oks		0.714	0.894	0.781	0.659
 bangbangren		0.706	0.880	0.765	0.656
 G-RMI		0.691	0.859	0.752	0.660
 FAIR Mask R-CNN		0.689	0.892	0.752	0.637

COCO Challenge 2017 Keypoint

SenseTime

 Megvii (Face++)	0.721	0.905	0.789	0.679
---	-------	-------	-------	-------

 oks	0.714	0.894	0.781	0.659
---	-------	-------	-------	-------

Google Research

 bangbangren	0.706	0.880	0.765	0.656
---	-------	-------	-------	-------

 G-RMI	0.691	0.859	0.752	0.660
---	-------	-------	-------	-------

 FAIR Mask R-CNN	0.689	0.892	0.752	0.637
--	-------	-------	-------	-------

COCO Challenge 2017 InstanceSeg

	UCenter	0.463	0.694	0.507	0.258
	Megvii (Face++)	0.460	0.707	0.503	0.263
	FAIR Mask R-CNN	0.435	0.687	0.466	0.234
	MSRA	0.426	0.680	0.460	0.240

Thanks