

Segmentation and Attention

PRESENTED BY DeepLearningPKU



Segmentation & Attention

01

PART ONE
Semantic Segmentation

02

PART TWO
Instance Segmentation

03

PART THREE
Discrete locations

04

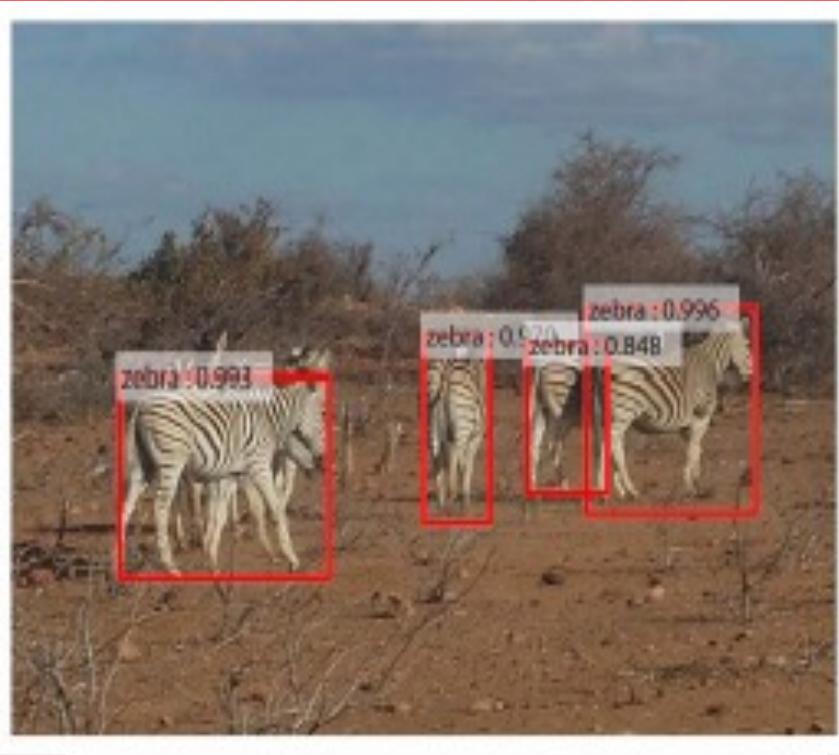
PART FOUR
Continuous locations

Segmentation

Semantic Segmentation

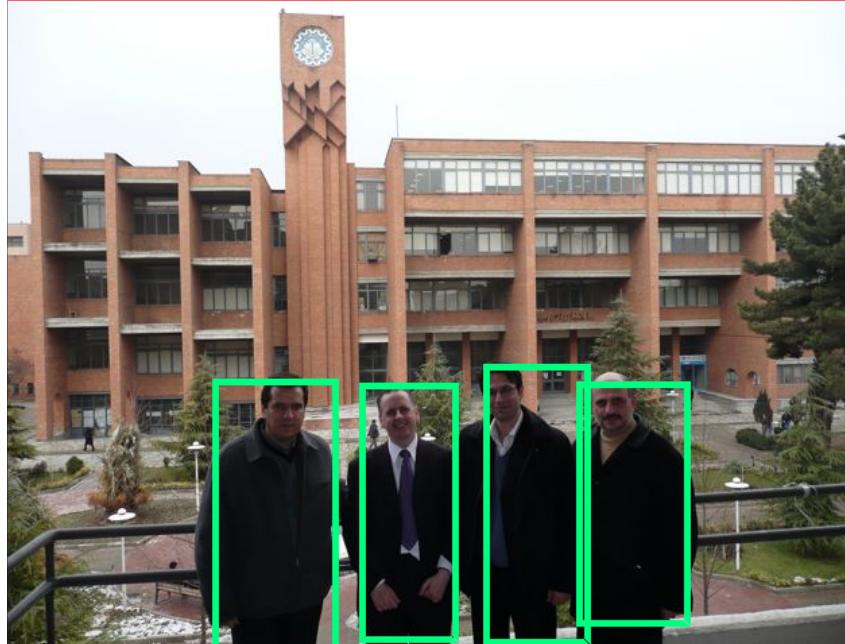
Instance Segmentation

基本概况



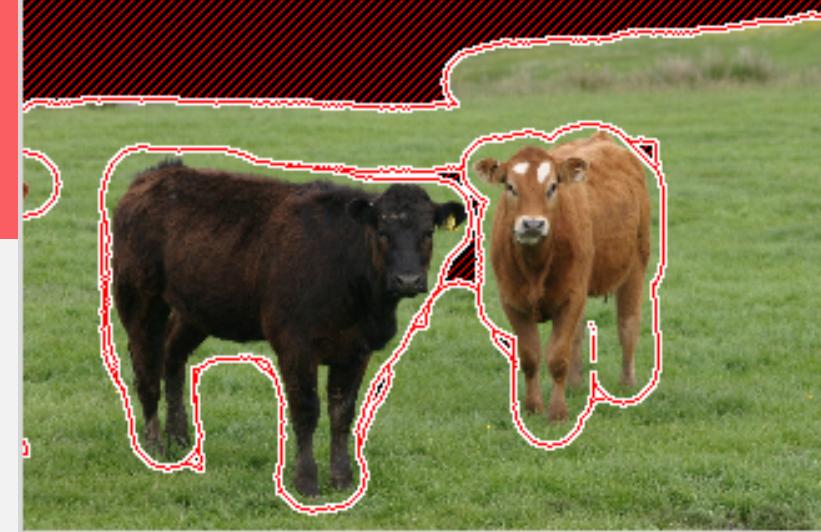
Object Recognition:

Localized them with a bounding box and label that bounding box with a label.



Object Detection:

Object recognition
two class of object classification
which means object bounding boxes
and non-object bounding boxes.



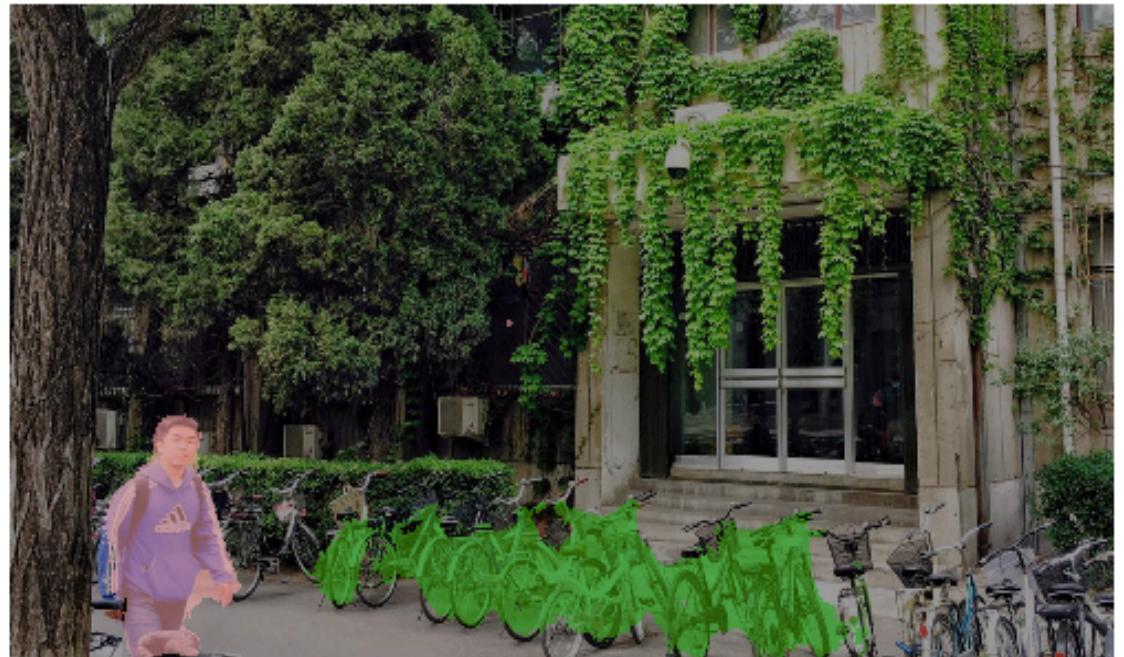
Object Segmentation:

Like object recognition you will
recognize all objects in an image but
your output should show this object
classifying pixels of the image.

Semantic Segmentation



Original image (hover to highlight segmented parts)



Semantic segmentation

Objects appearing in the image:

Bicycle

Person

Objects not appearing in the image:

Aeroplane

Bird

Boat

Bottle

Bus

Car

Cat

Chair

Cow

Dining table

Dog

Horse

Motorbike

Potted plant

Sheep

Sofa

Train

TV/Monitor

Computer Vision Tasks

Classification



CAT

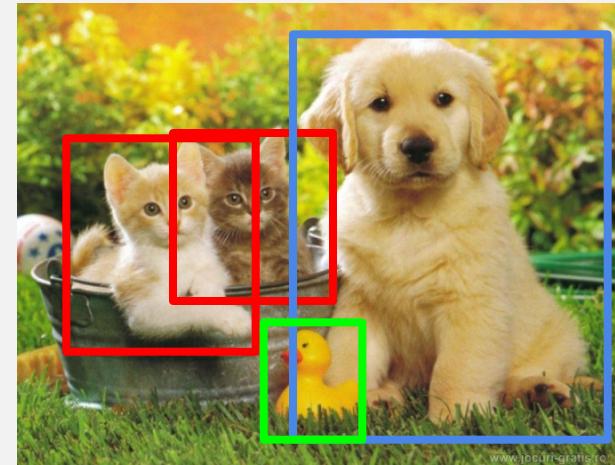


CAT

Single
object

Classification + Localization

Object Detection



CAT, DOG, DUCK

Segmentation

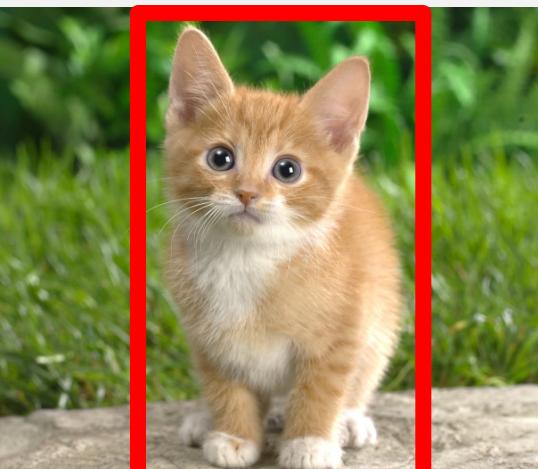


CAT, DOG, DUCK

Multiple
objects

Computer Vision Tasks

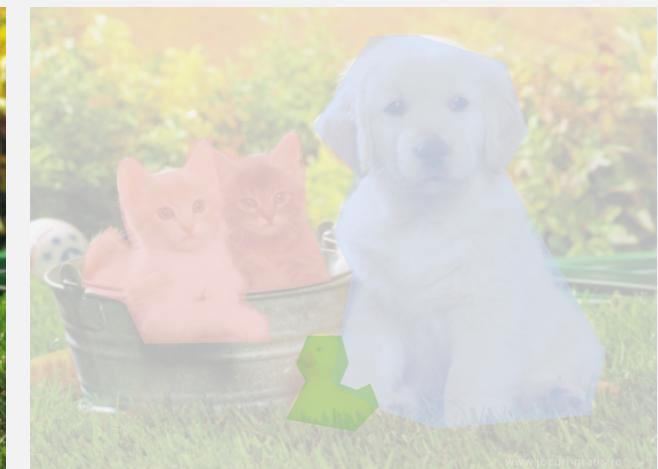
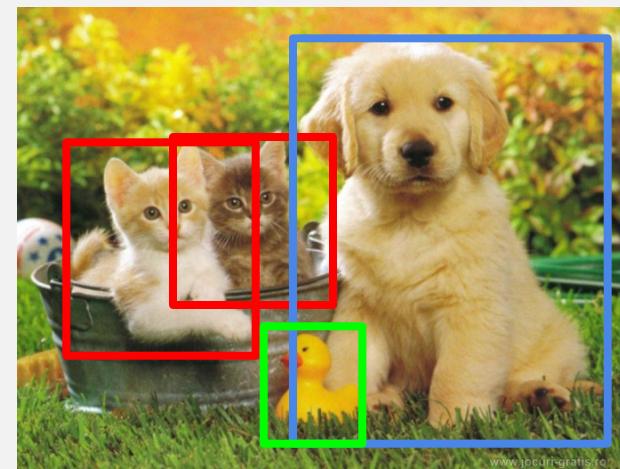
Classification



**Classification
+ Localization**

**Object
Detection**

Segmentation



Lecture 8

Computer Vision Tasks

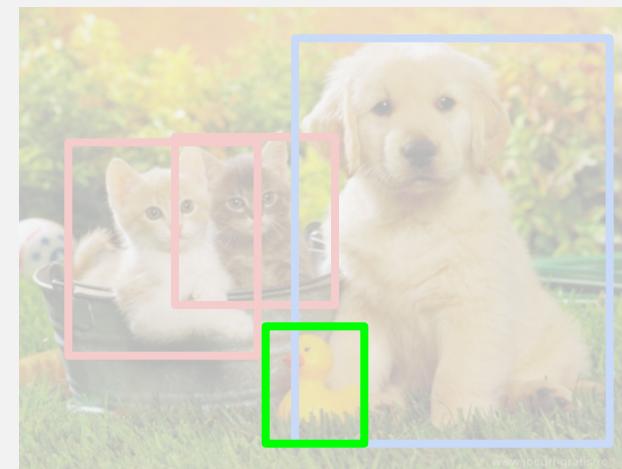
Classification



Classification
+ Localization



Object
Detection



Segmentation



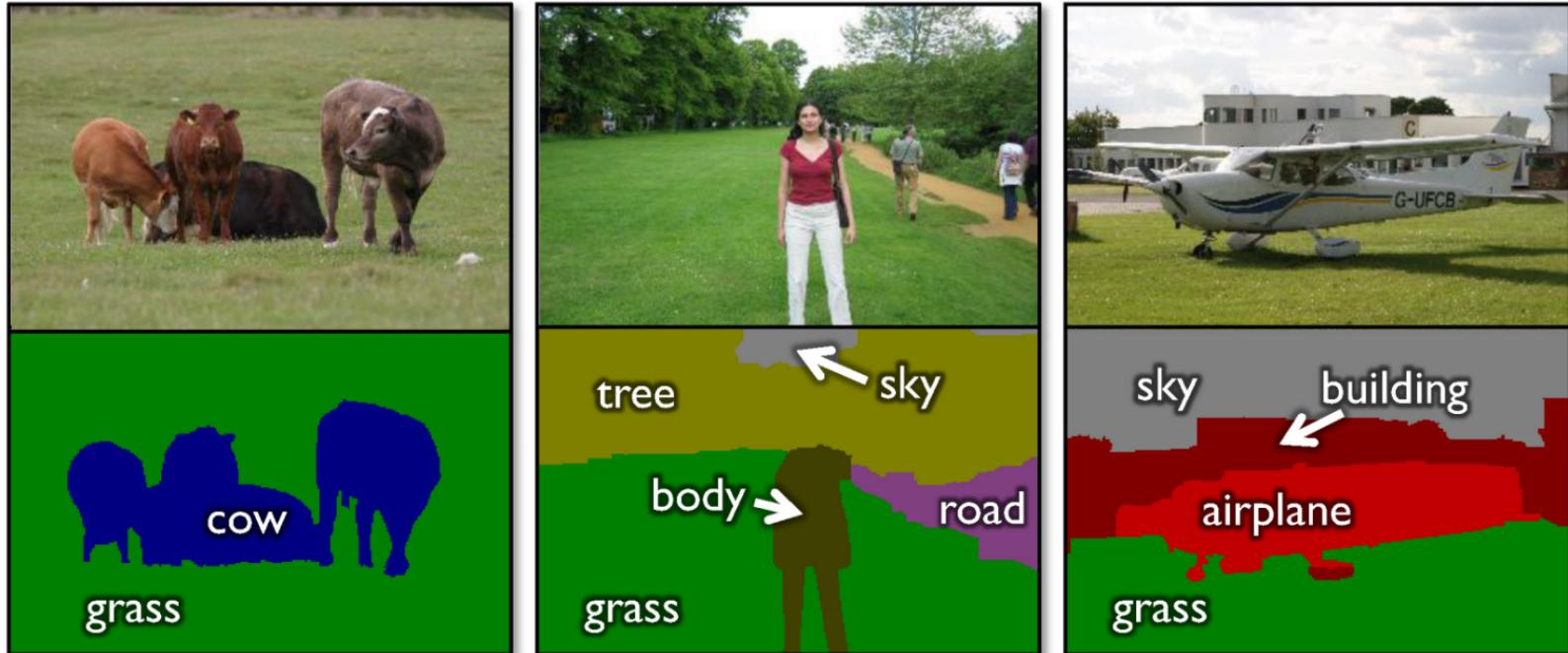
Today

Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)

Classic computer vision problem



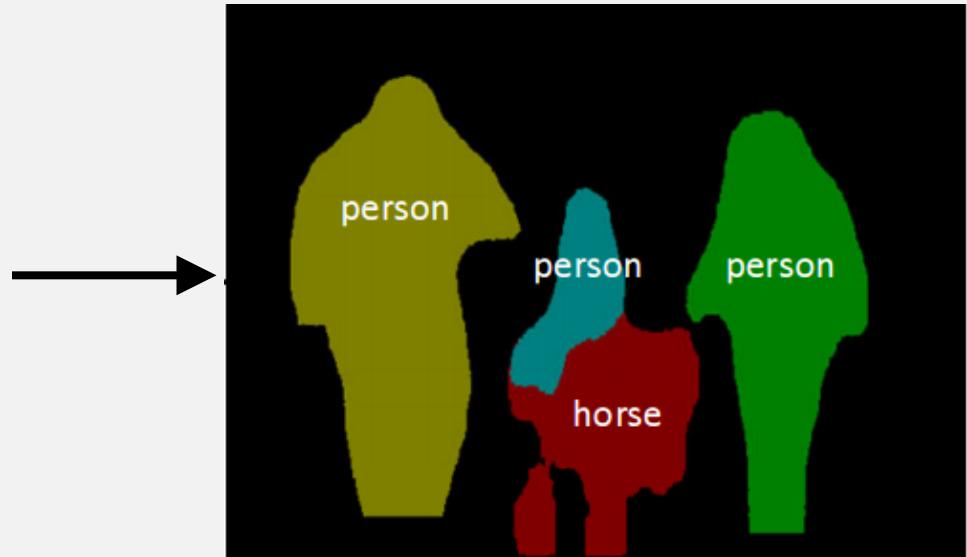
object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

Figure credit: Shotton et al, "TextonBoost for Image Understanding: Multi-Class Object Recognition and Segmentation by Jointly Modeling Texture, Layout, and Context" , IJCV 2007

Instance Segmentation

Detect instances,
give category, label
pixels

“simultaneous
detection and
segmentation”
(SDS)



Lots of recent work
(MS-COCO)

Figure credit: Dai et al, “Instance-aware Semantic Segmentation via Multi-task Network Cascades” , arXiv 2015



PART 1

Semantic Segmentation

Semantic segmentation faces an inherent tension between semantics and location: global information resolves what while local information resolves where.

基本概況



Semantic segmentation

Scene parsing, or semantic segmentation, consists in labeling each pixel in an image with the category of the object it belongs to. It is a challenging task that involves the simultaneous detection, segmentation and recognition of all the objects in the image.

1 Semantic Segmentation

The goal of the scene labeling task is to assign a class label to each pixel in an image.

2 Semantic Segmentation: Multi-Scale

uses purely supervised training from fully-labeled images to learn appropriate low-level and mid-level features

3 Refinement

For a good visual coherence and a high class accuracy, capture long range (pixel) label dependencies in images.

4 Upsampling

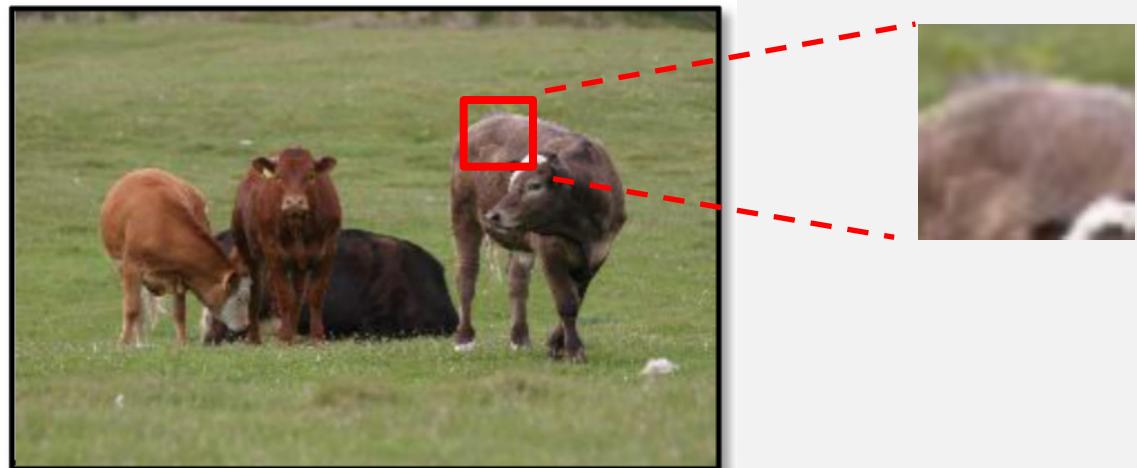
semantic information + appearance information
a deep, coarse layer + a shallow, fine layer
to produce accurate and detailed segmentations.

Semantic Segmentation

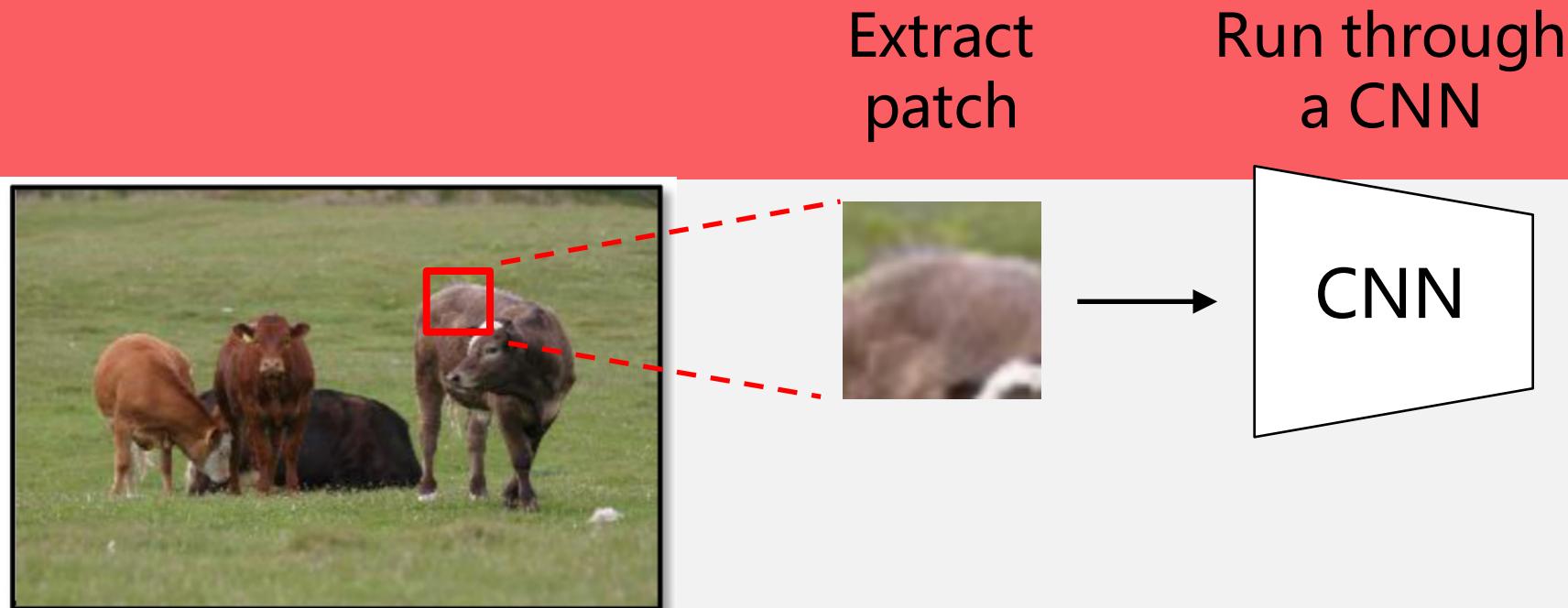


Semantic Segmentation

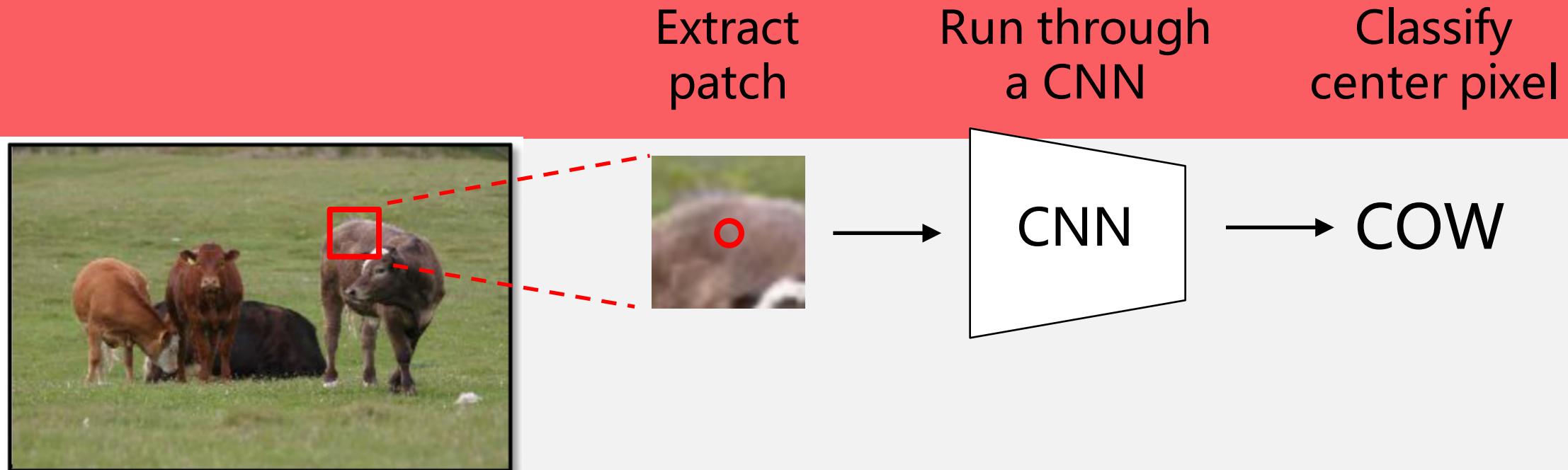
Extract
patch



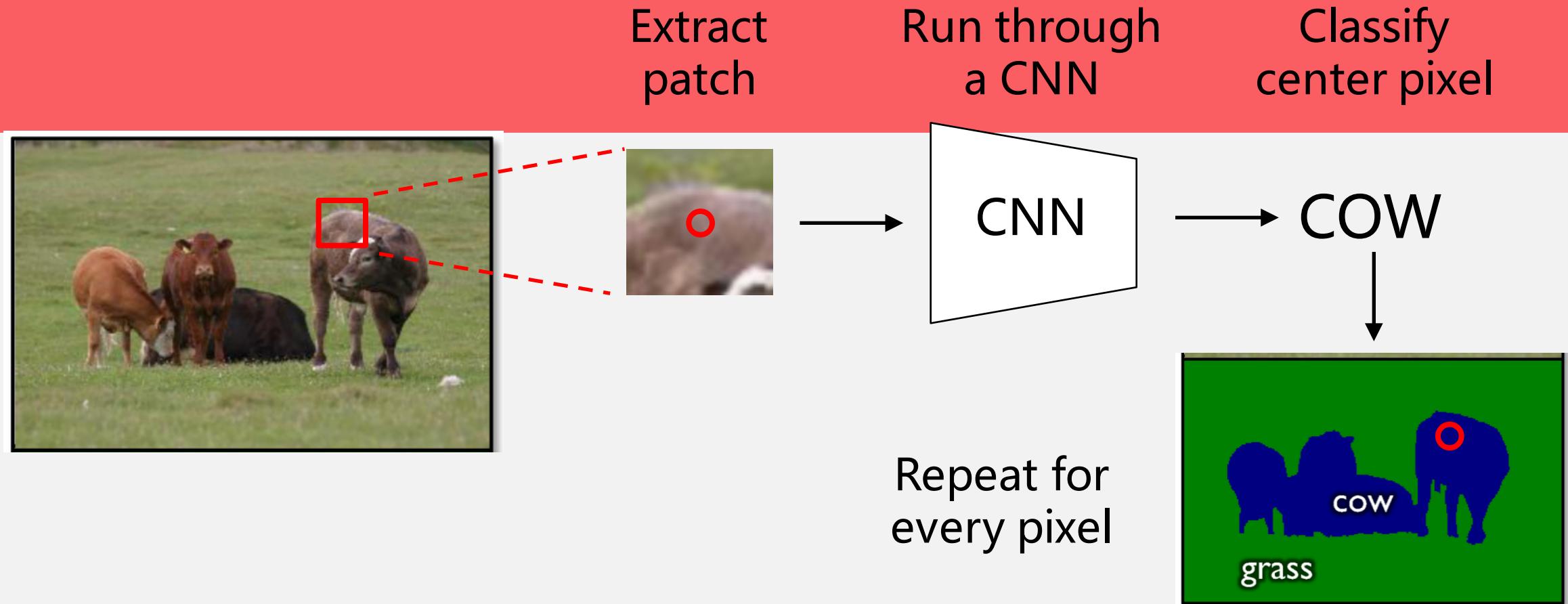
Semantic Segmentation



Semantic Segmentation

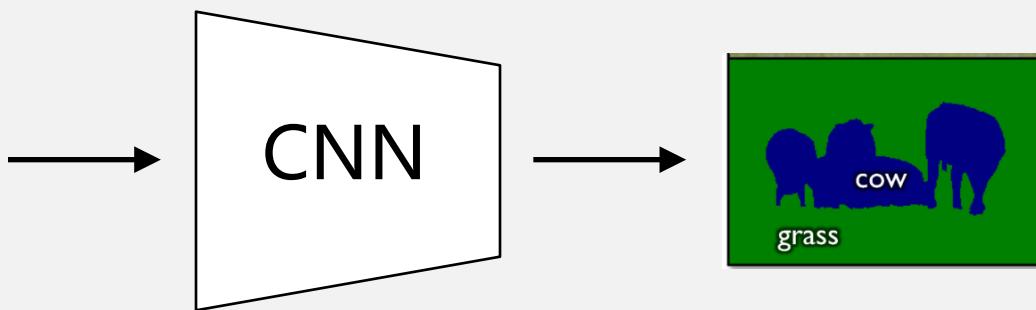


Semantic Segmentation



Semantic Segmentation

Run “fully convolutional” network to get all pixels at once



Smaller output
due to pooling

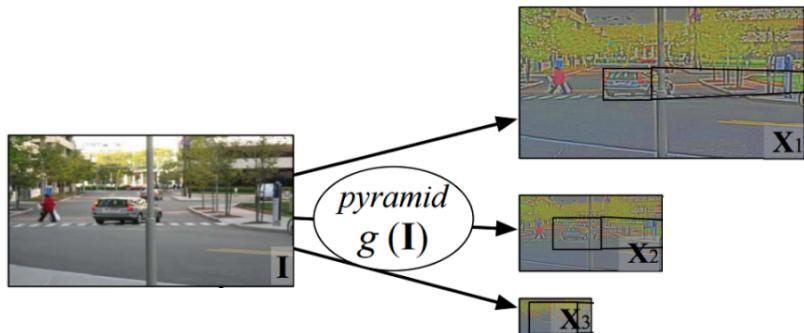
Multi-Scale

Semantic Segmentation: Multi-Scale

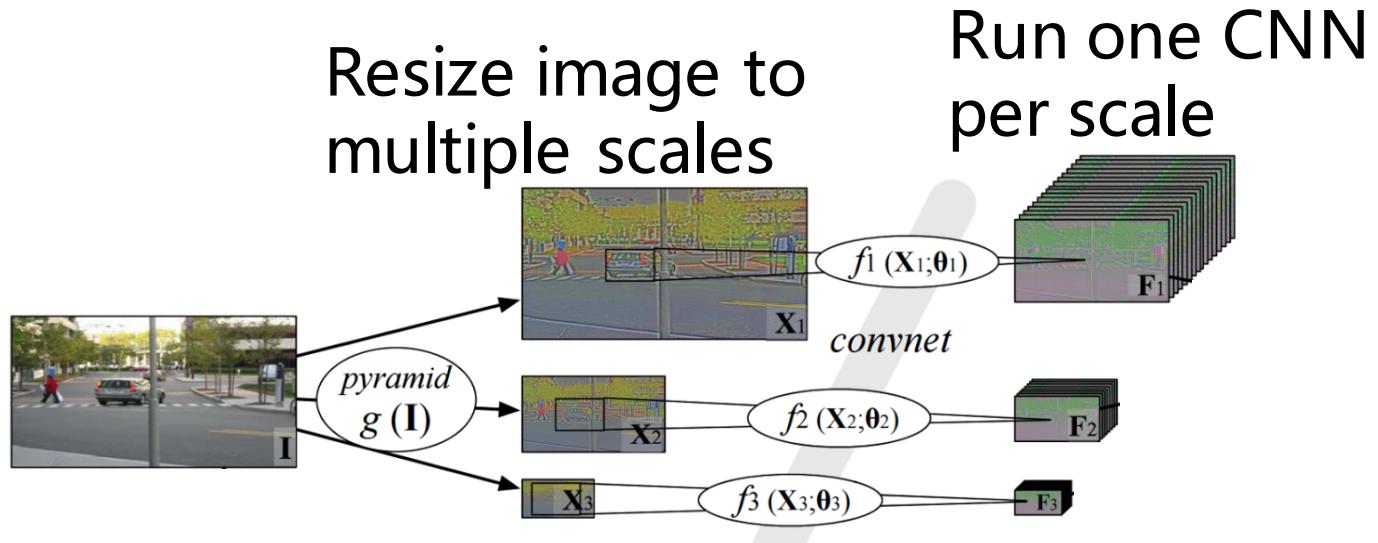


Semantic Segmentation: Multi-Scale

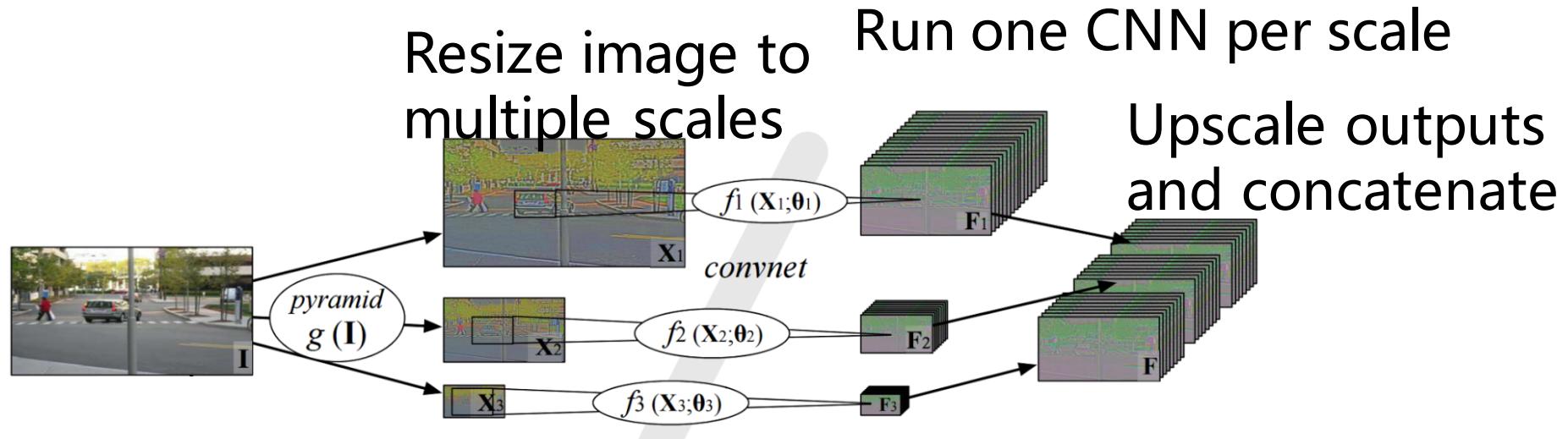
Resize image to
multiple scales



Semantic Segmentation: Multi-Scale

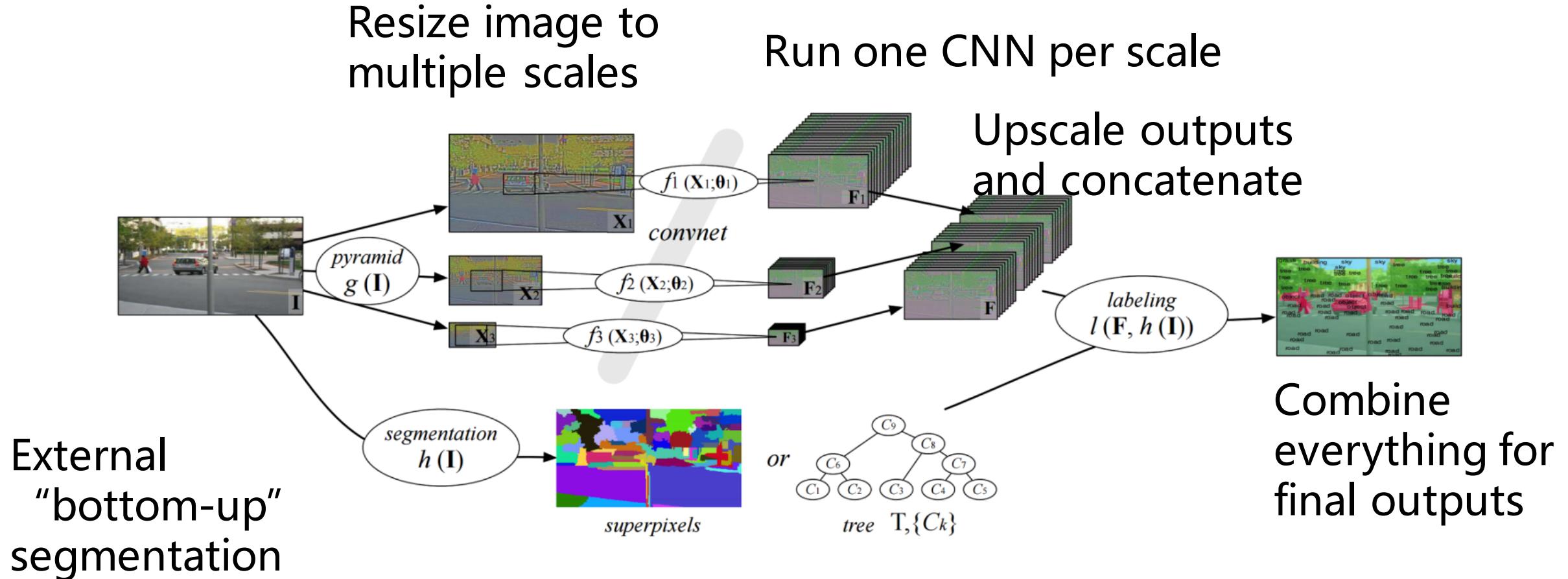


Semantic Segmentation: Multi-Scale

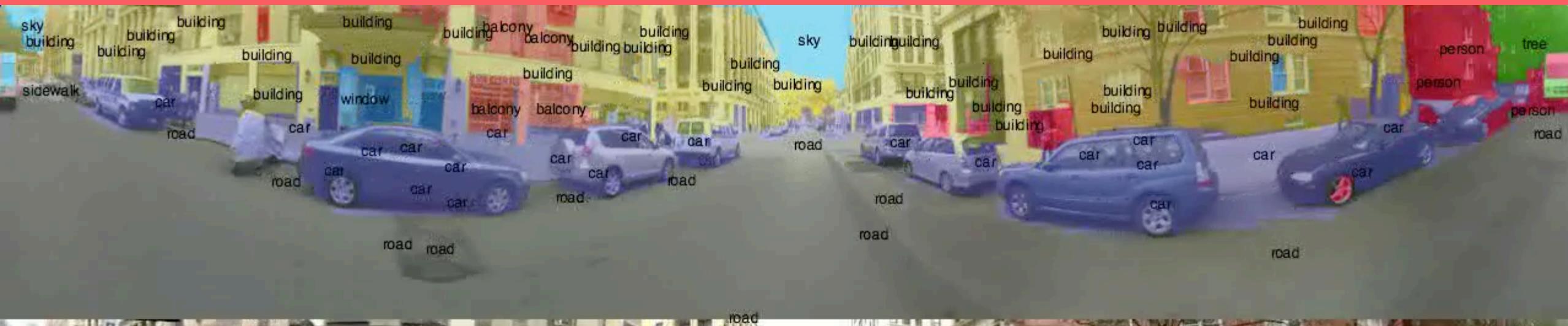


Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Semantic Segmentation: Multi-Scale

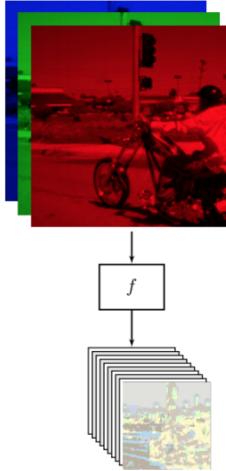


Multi-Scale



Semantic Segmentation: Refinement

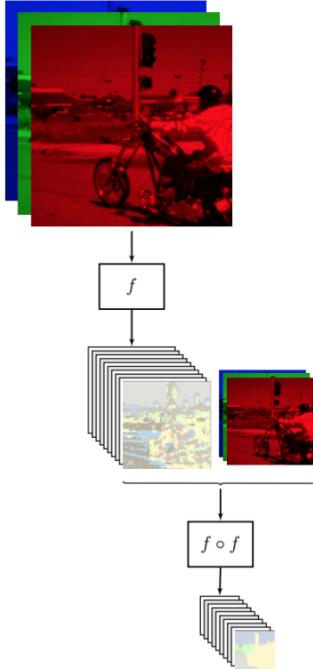
Apply CNN
once to get
labels



Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling" ,
ICML 2014

Semantic Segmentation: Refinement

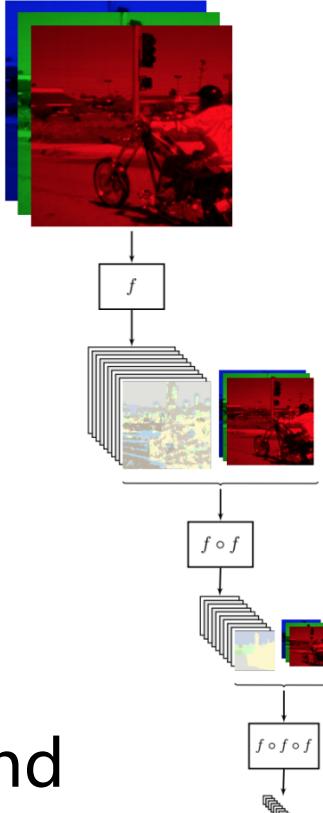
Apply CNN
once to get
labels



Apply
AGAIN to
refine
labels

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling" ,
ICML 2014

Semantic Segmentation: Refinement



Apply CNN
once to get
labels

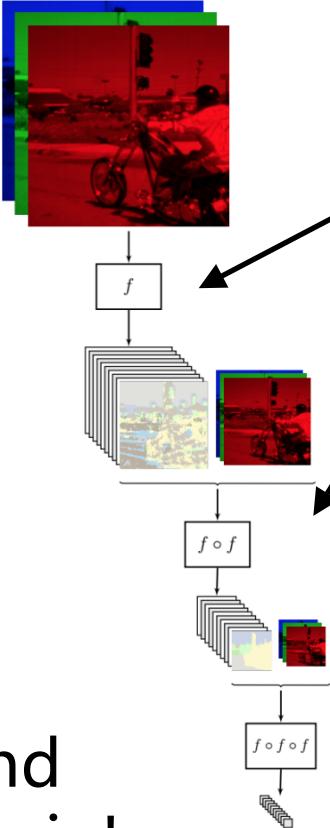
Apply
AGAIN to
refine
labels And
again!

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling" ,
ICML 2014

Semantic Segmentation: Refinement

Apply CNN
once to get
labels

Apply
AGAIN to
refine
labels And
again!



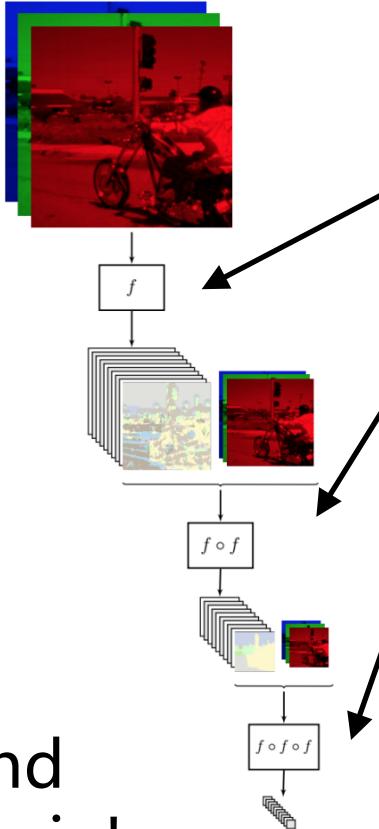
Same CNN weights:
**recurrent convolutional
network**

Semantic Segmentation: Refinement

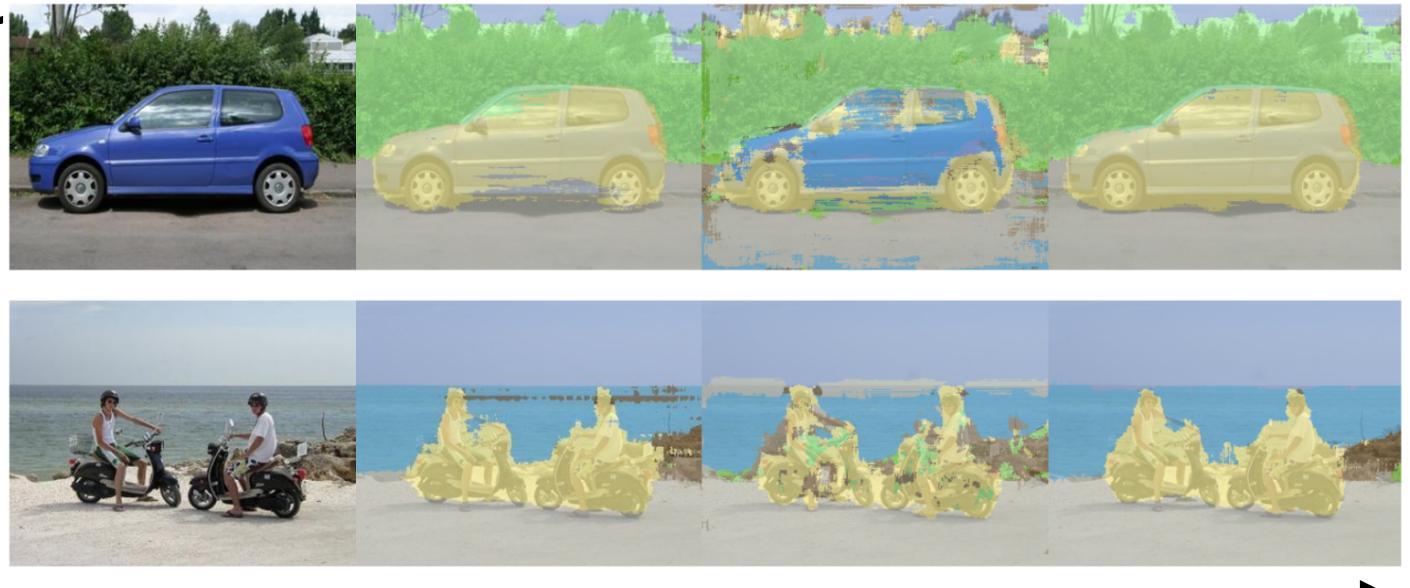
Apply CNN once to get labels

Apply AGAIN to refine labels

And again!



Same CNN weights:
recurrent convolutional network

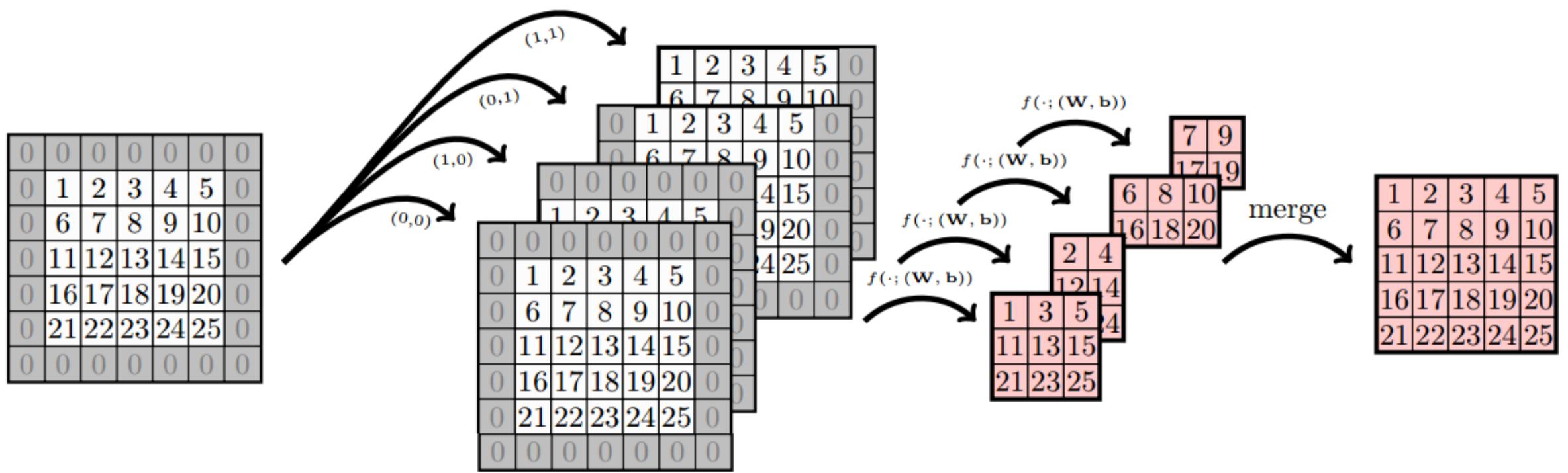


More iterations improve results

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling" ,
ICML 2014

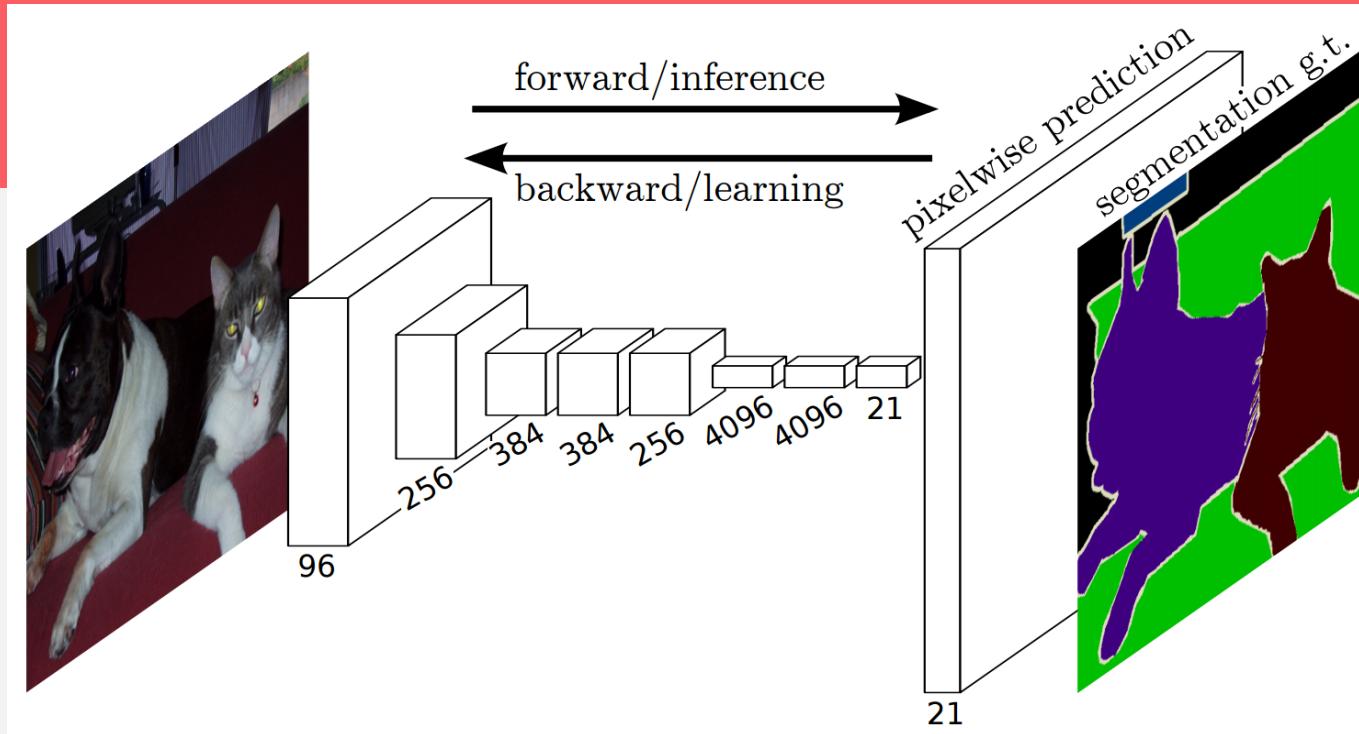
Refinement

- downscaled label planes (compared to the input image) due to pooling layers.
- Downscaled predicted label planes (here in red) are then merged to get back the full resolution label plane in an efficient manner



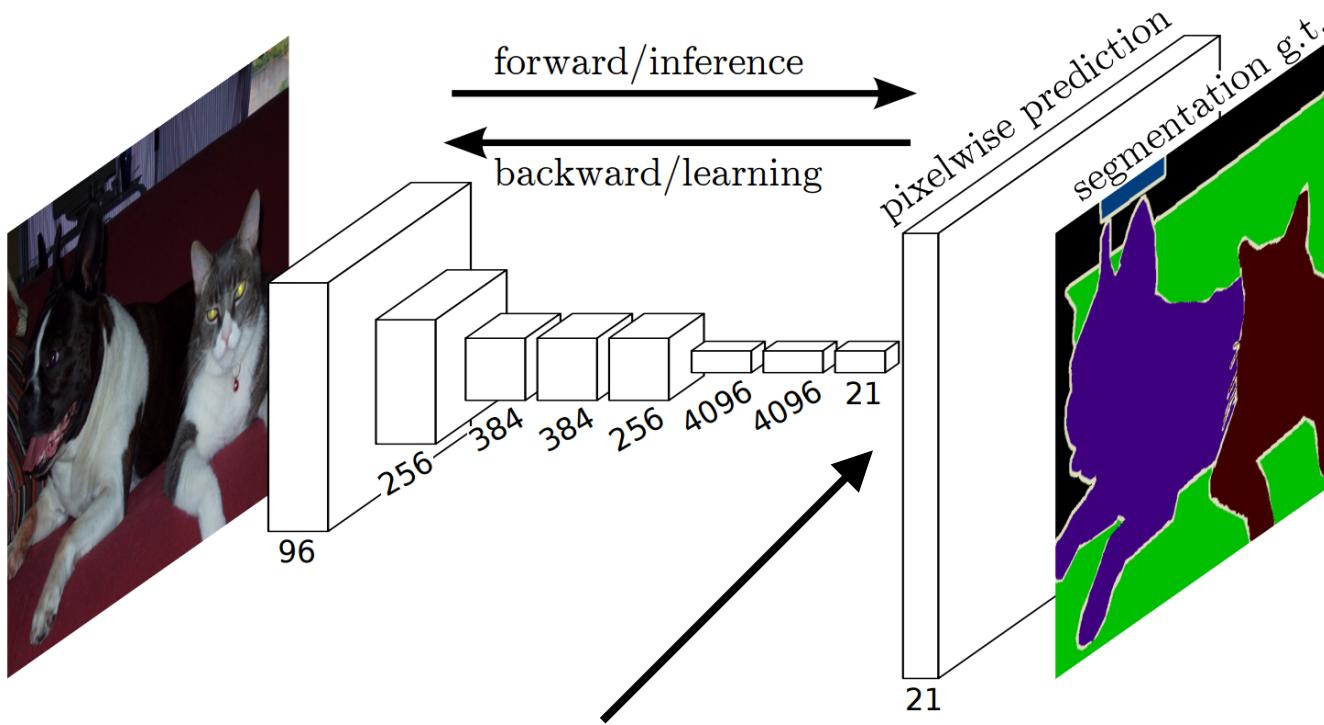
Upsampling

Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation" ,
CVPR 2015

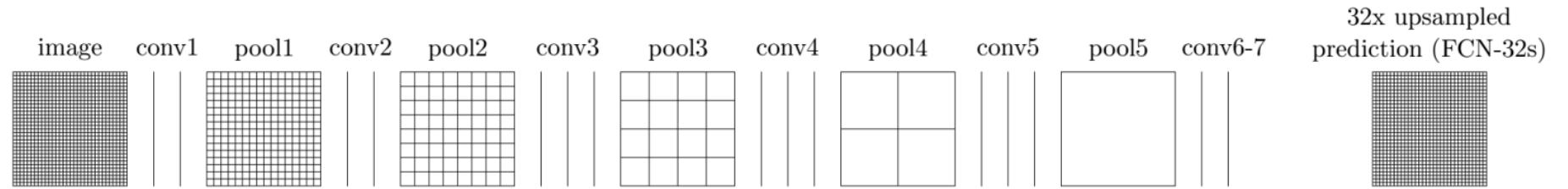
Semantic Segmentation: Upsampling



Learnable
upsampling!

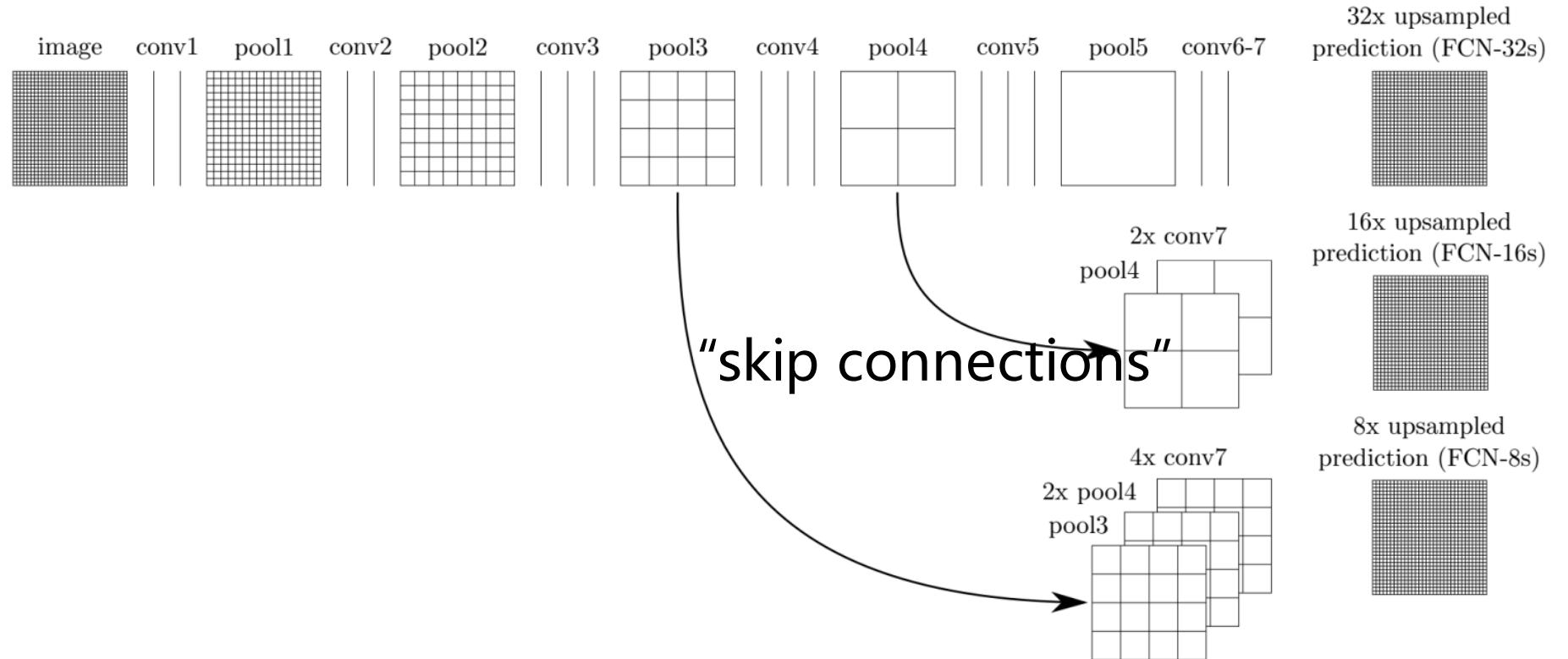
Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation" ,
CVPR 2015

Semantic Segmentation: Upsampling



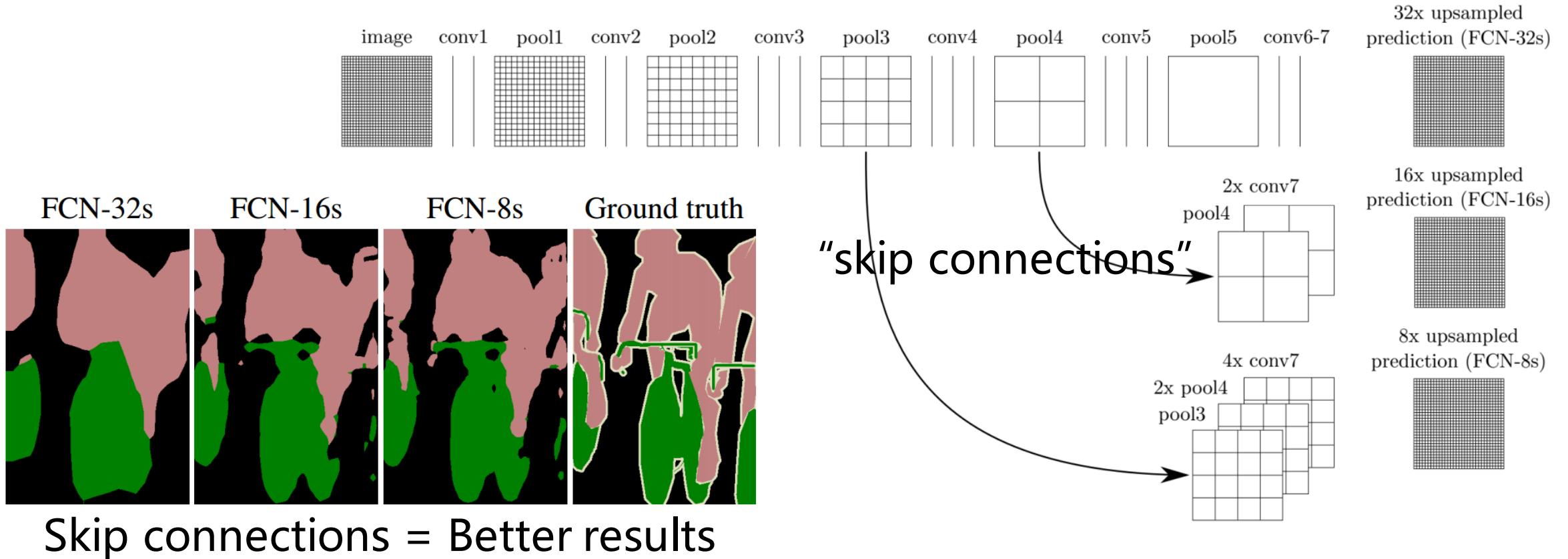
Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation" ,
CVPR 2015

Semantic Segmentation: Upsampling



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation" ,
CVPR 2015

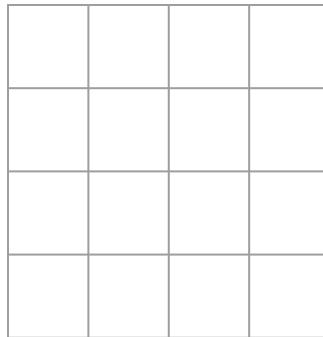
Semantic Segmentation: Upsampling



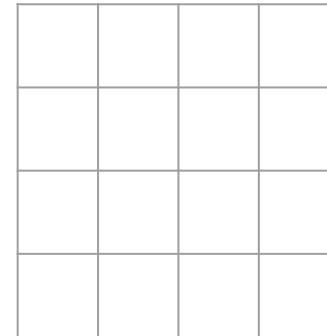
Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation" ,
CVPR 2015

Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1



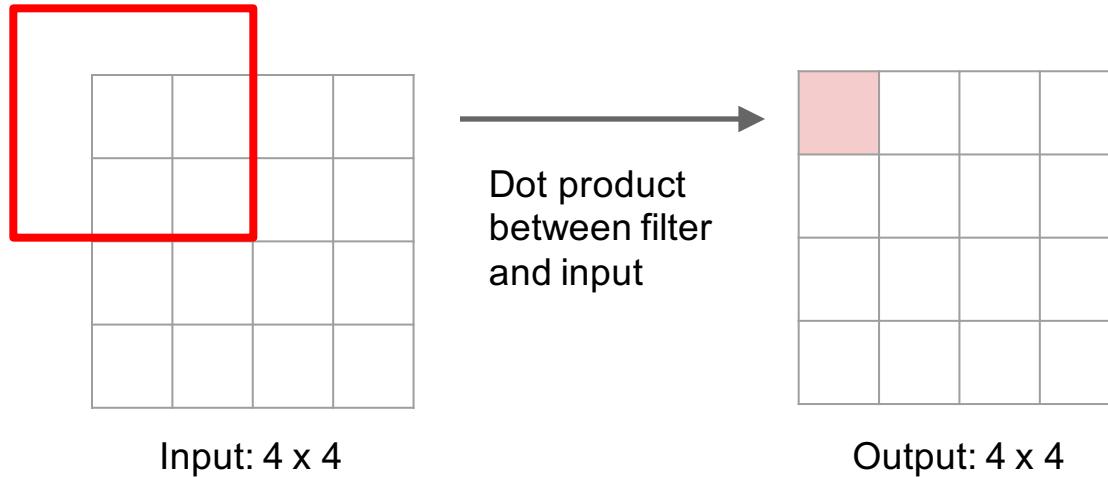
Input: 4×4



Output: 4×4

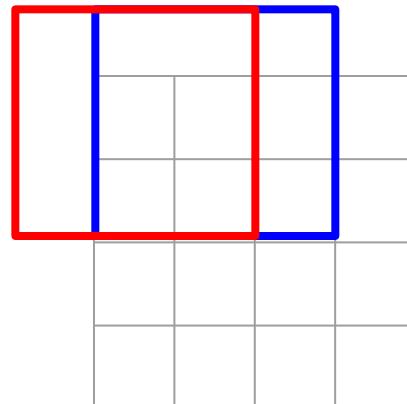
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1



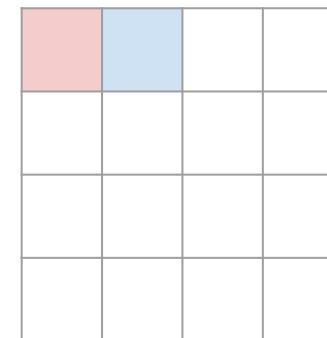
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 1 pad 1



Input: 4×4

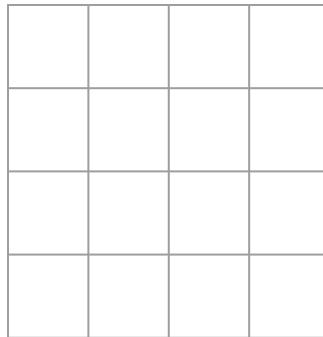
Dot product
between filter
and input



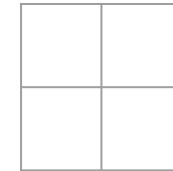
Output: 4×4

Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, **stride 2** pad 1



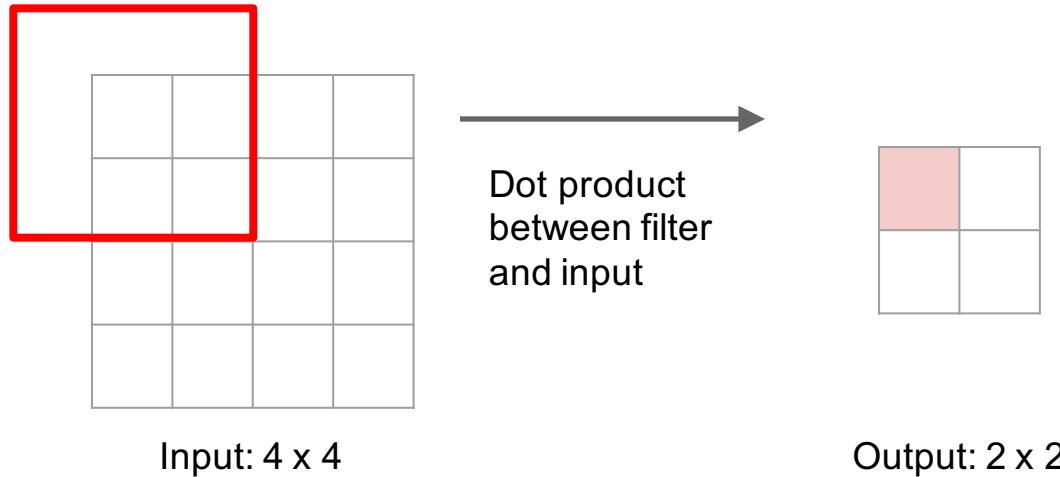
Input: 4×4



Output: 2×2

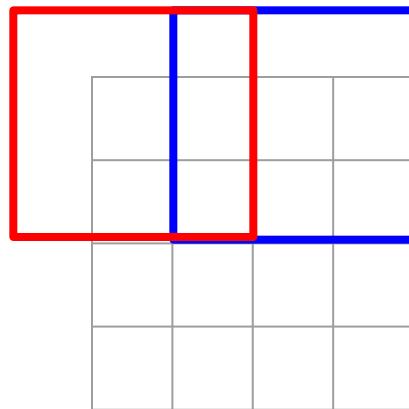
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 2 pad 1



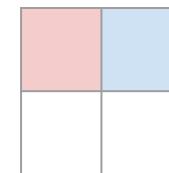
Learnable Upsampling: “Deconvolution”

Typical 3×3 convolution, stride 2 pad 1



Input: 4×4

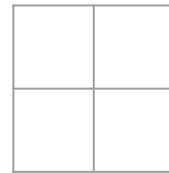
Dot product
between filter
and input



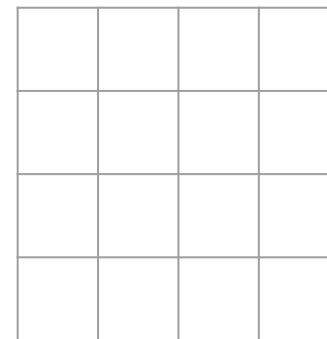
Output: 2×2

Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1



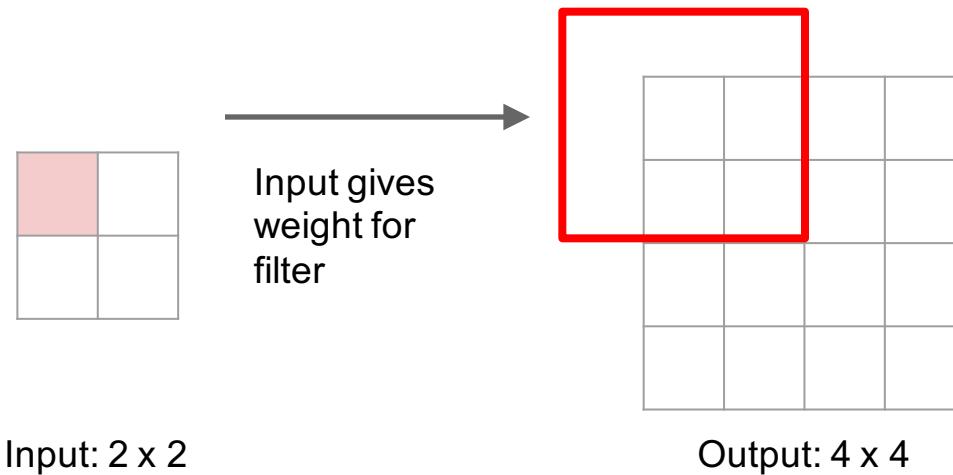
Input: 2 x 2



Output: 4 x 4

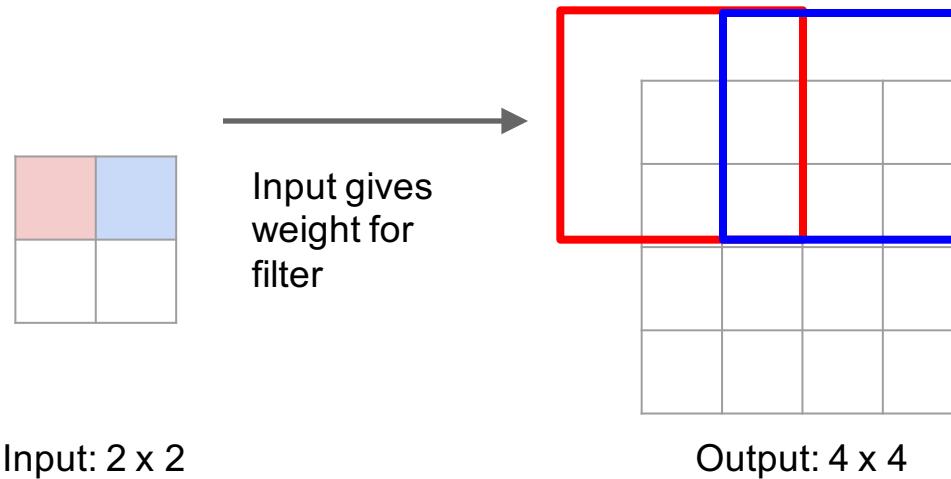
Learnable Upsampling: “Deconvolution”

3 x 3 “deconvolution”, stride 2 pad 1

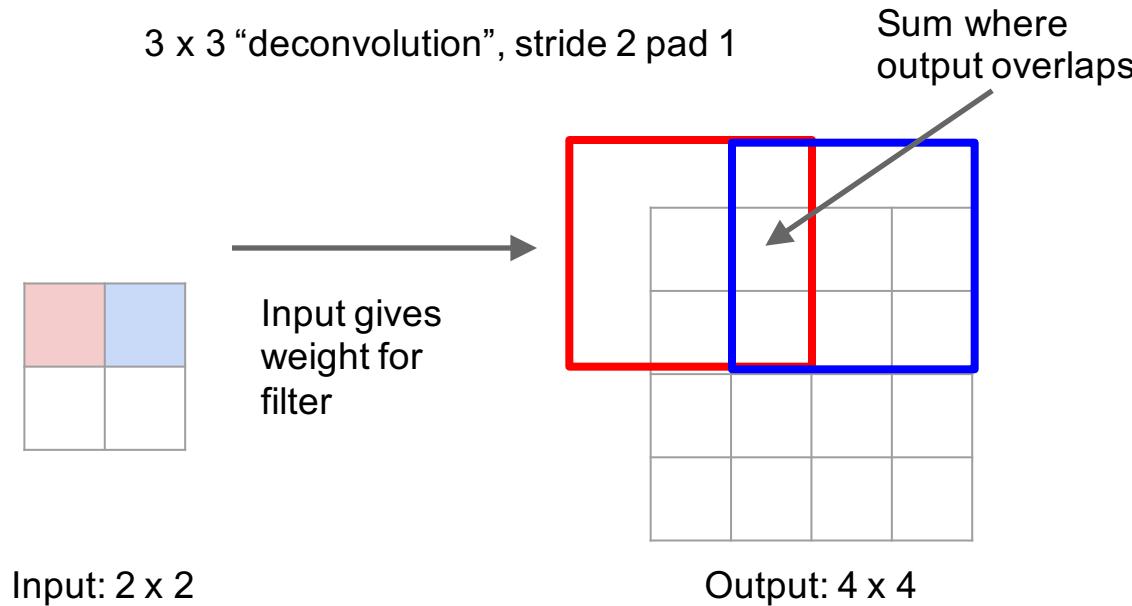


Learnable Upsampling: “Deconvolution”

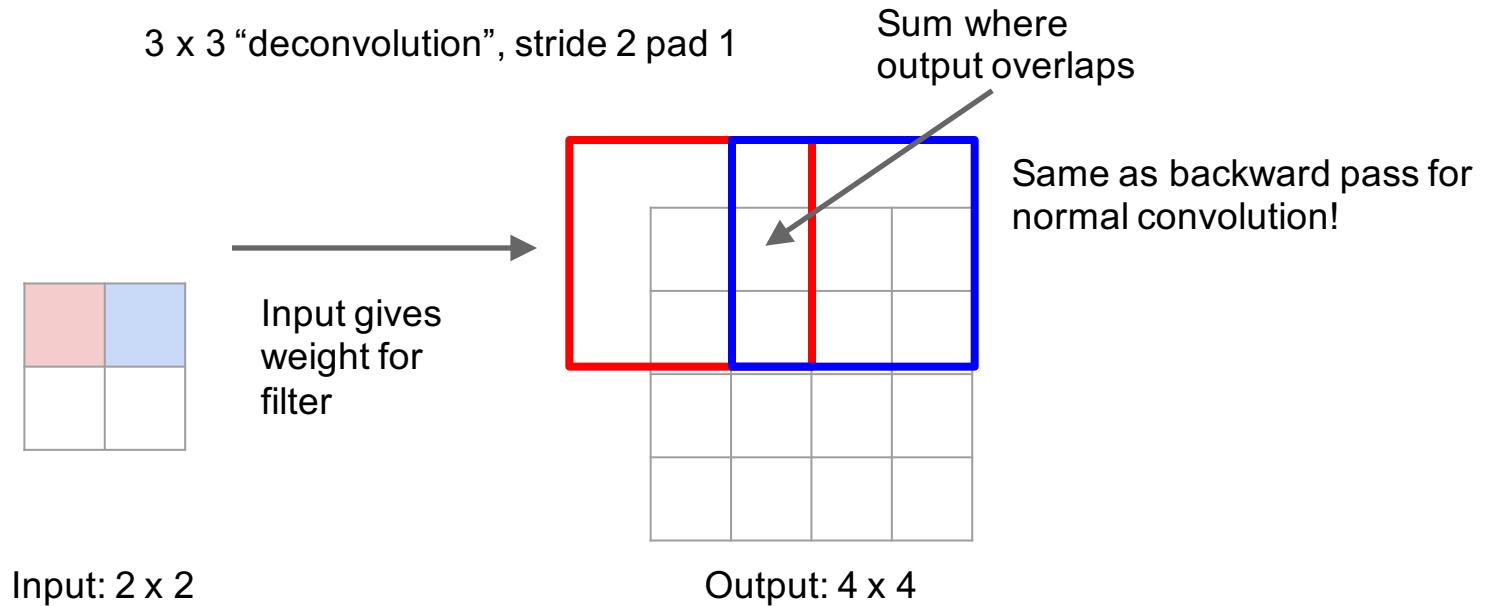
3 x 3 “deconvolution”, stride 2 pad 1



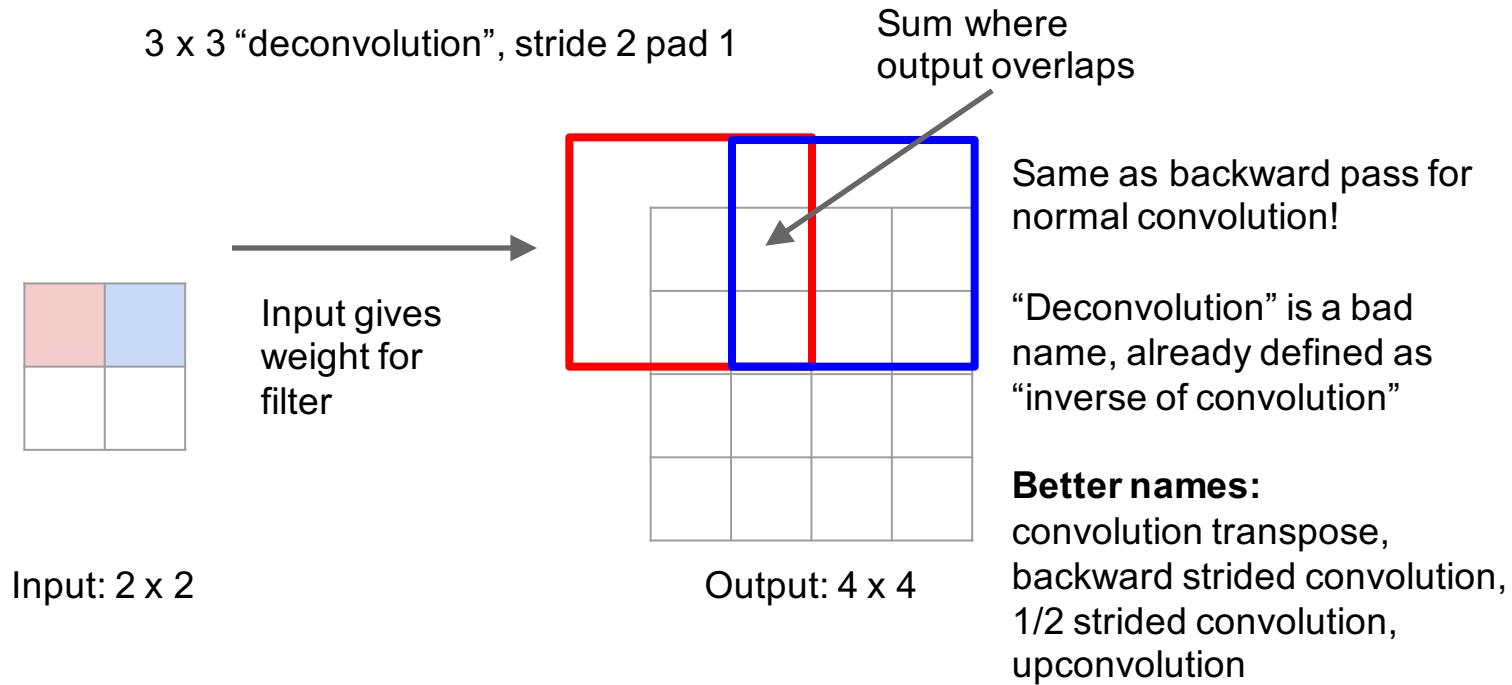
Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”



Learnable Upsampling: “Deconvolution”

¹It is more proper to say “convolutional transpose operation” rather than “deconvolutional” operation. Hence, we will be using the term “convolutional transpose” from now.

Im et al, “Generating images with recurrent adversarial networks” , arXiv 2016

A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions)

Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks” , ICLR 2016

“Deconvolution” is a bad name, already defined as “inverse of convolution”

Better names:
convolution transpose,
backward strided
convolution,
1/2 strided convolution,
upconvolution

Learnable Upsampling: “Deconvolution”

¹It is more proper to say “convolutional transpose operation” rather than “deconvolutional” operation. Hence, we will be using the term “convolutional transpose” from now.

Im et al, “Generating images with recurrent adversarial networks”, arXiv 2016

A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions)

Radford et al, “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”, ICLR 2016

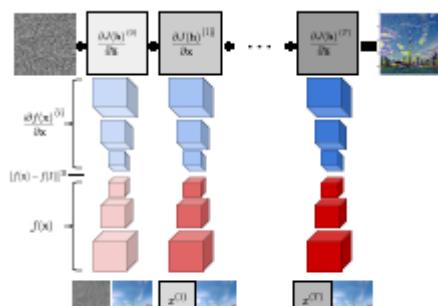


Figure 21. The gradient of convolution at index k .

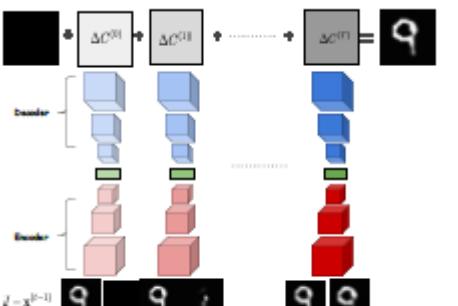
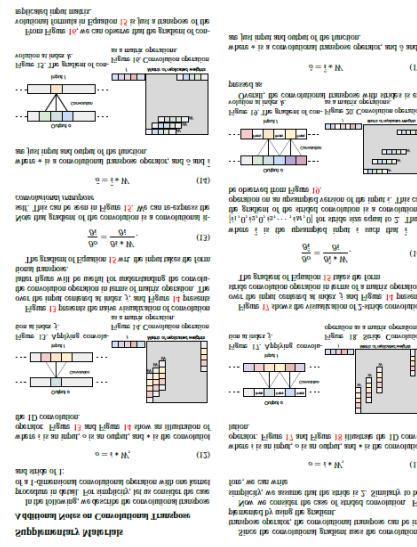


Figure 22. The abstract view of DRAW architecture is delineated.

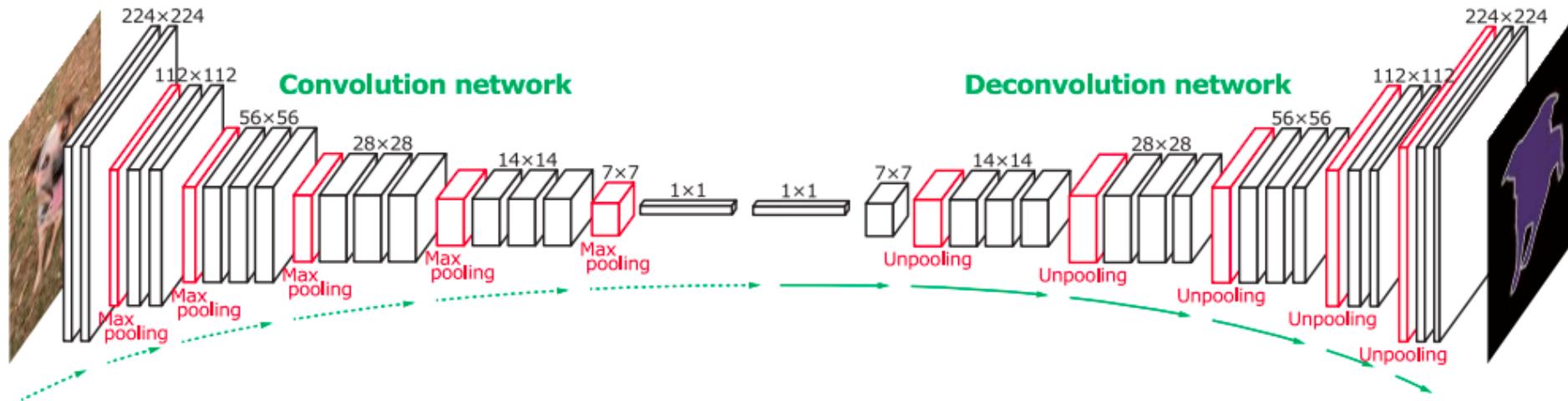
Great explanation in appendix



“Deconvolution” is a bad name, already defined as “inverse of convolution”

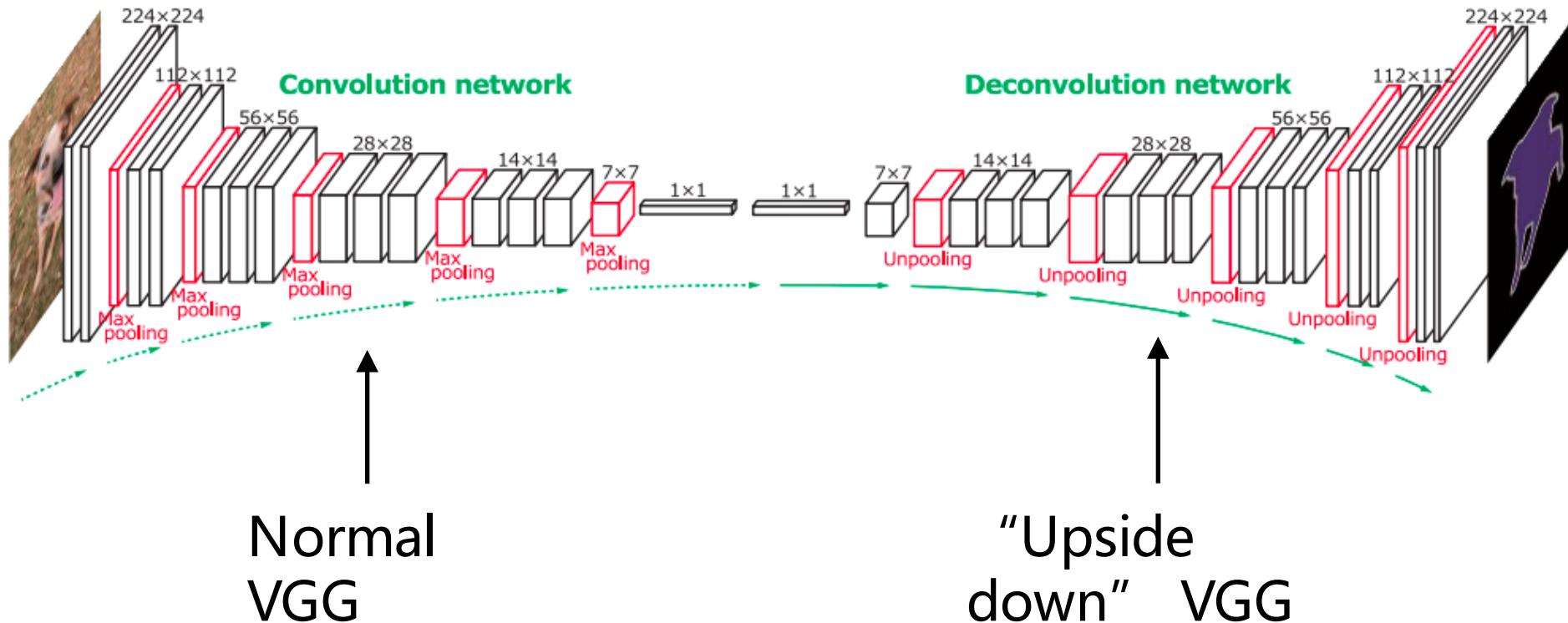
Better names:
convolution transpose,
backward strided
convolution,
1/2 strided convolution,
upconvolution

Semantic Segmentation: Upsampling



Noh et al, "Learning Deconvolution Network for Semantic Segmentation" , ICCV 2015

Semantic Segmentation: Upsampling



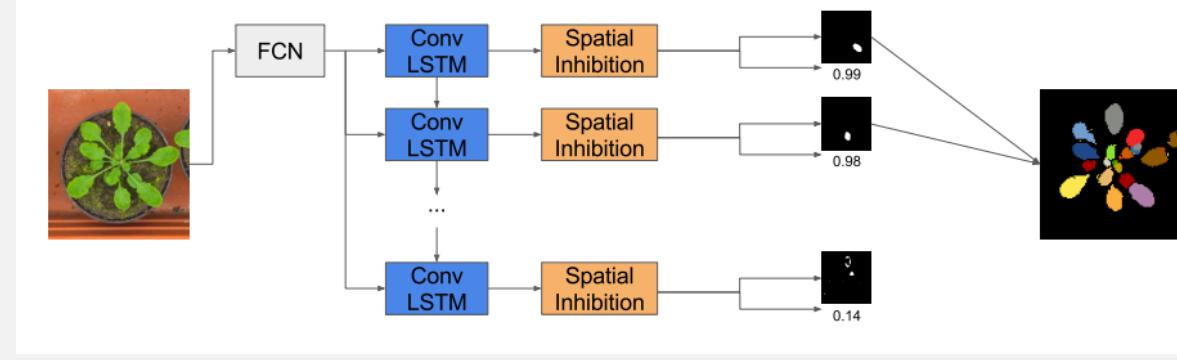
Noh et al, "Learning Deconvolution Network for Semantic Segmentation" , ICCV 2015

6 days of training on
Titan X...
51



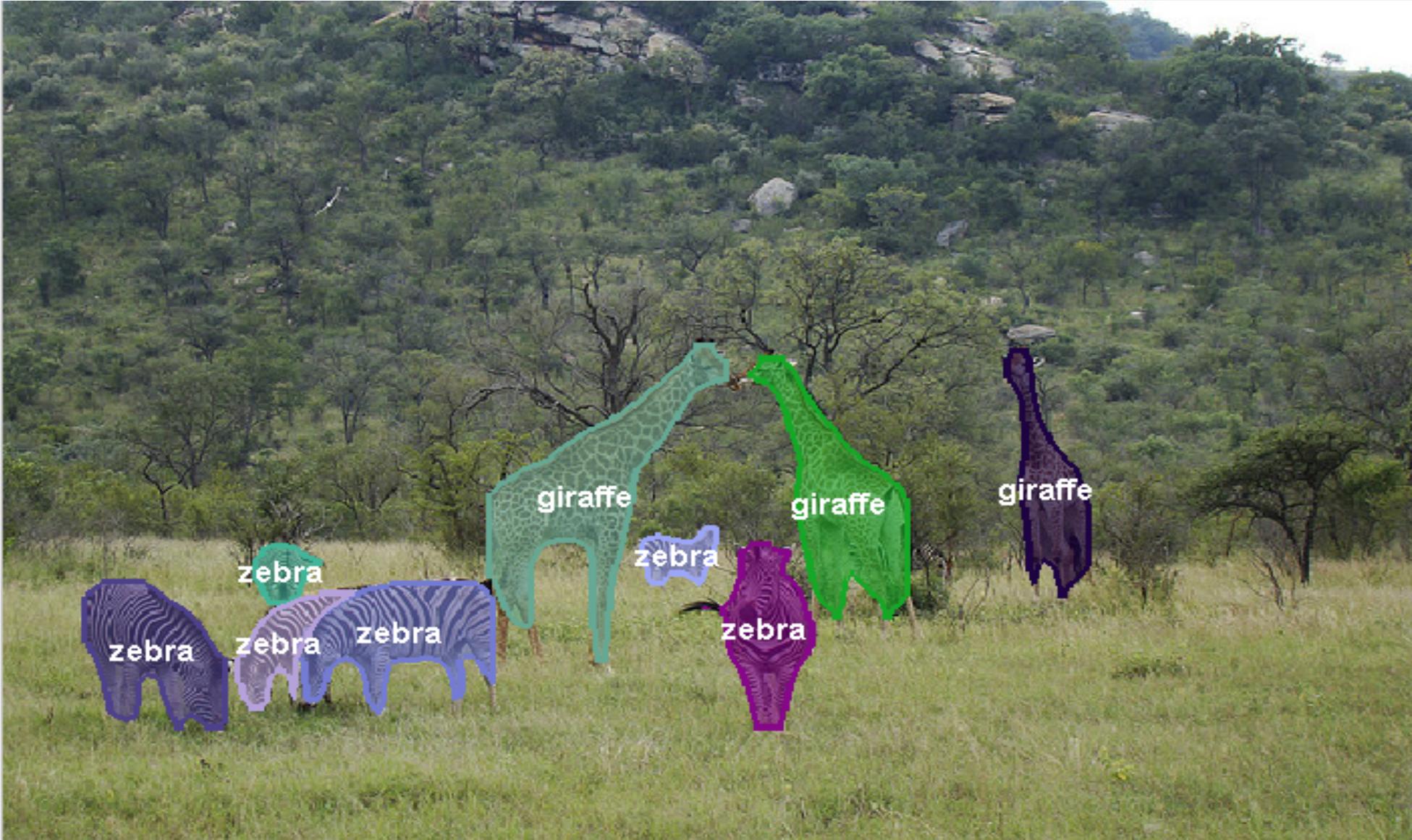
PART 2

Instance Segmentation



- Semantic segmentation research has recently witnessed rapid progress, but many leading methods are unable to identify object instances.

Difference



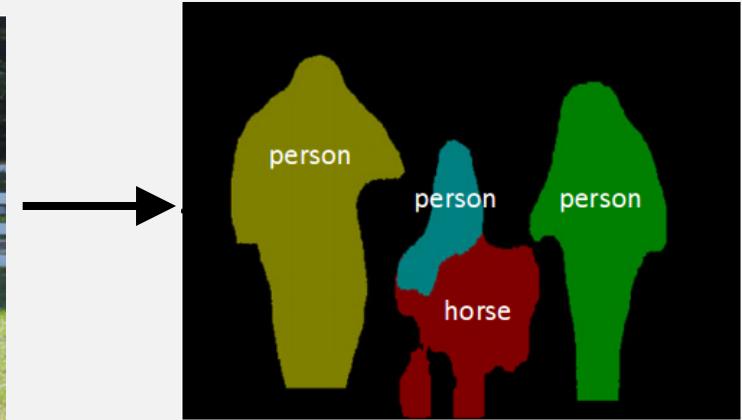
Instance Segmentation

Detect instances,
give category, label
pixels

“simultaneous
detection and
segmentation”
(SDS)

Lots of recent work
(MS-COCO)

detect all instances of a category
in an image and, for each instance,
mark the pixels that belong to it



54

Instance Segmentation

Similar to R-CNN,
but with
segments

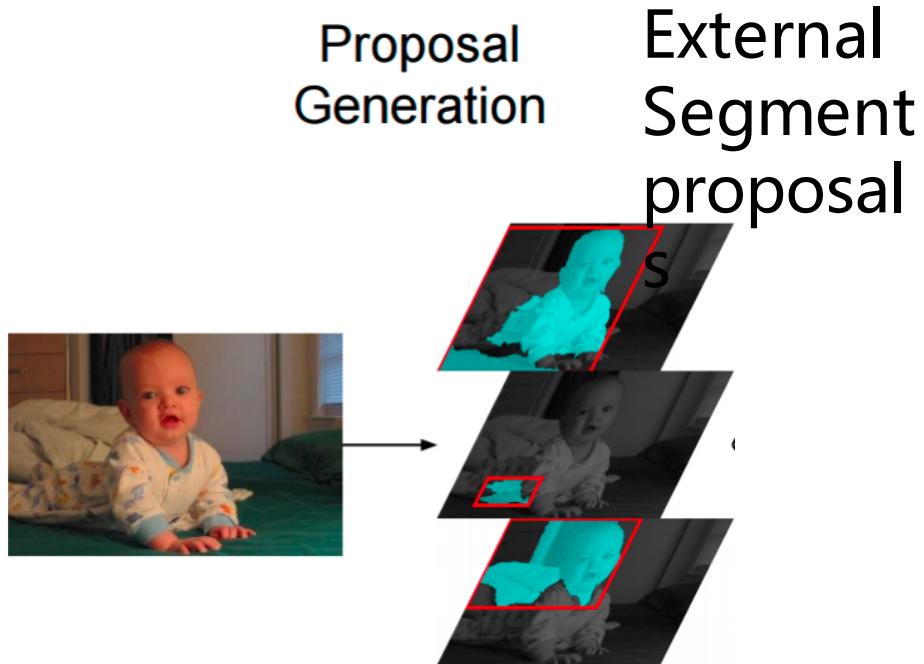


Hariharan et al, "Simultaneous Detection and Segmentation" , ECCV 2014

Instance Segmentation

region proposals: 2000 region candidates

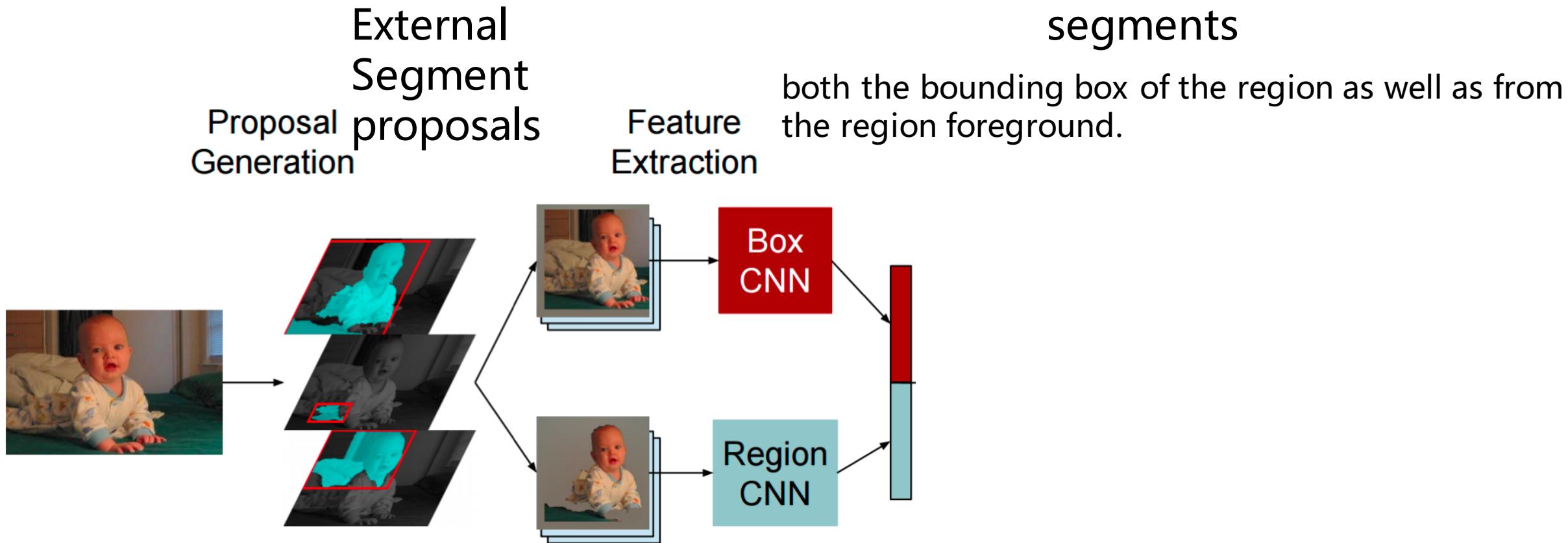
Similar to R-CNN,
but with
segments



Hariharan et al, "Simultaneous Detection and Segmentation" , ECCV 2014

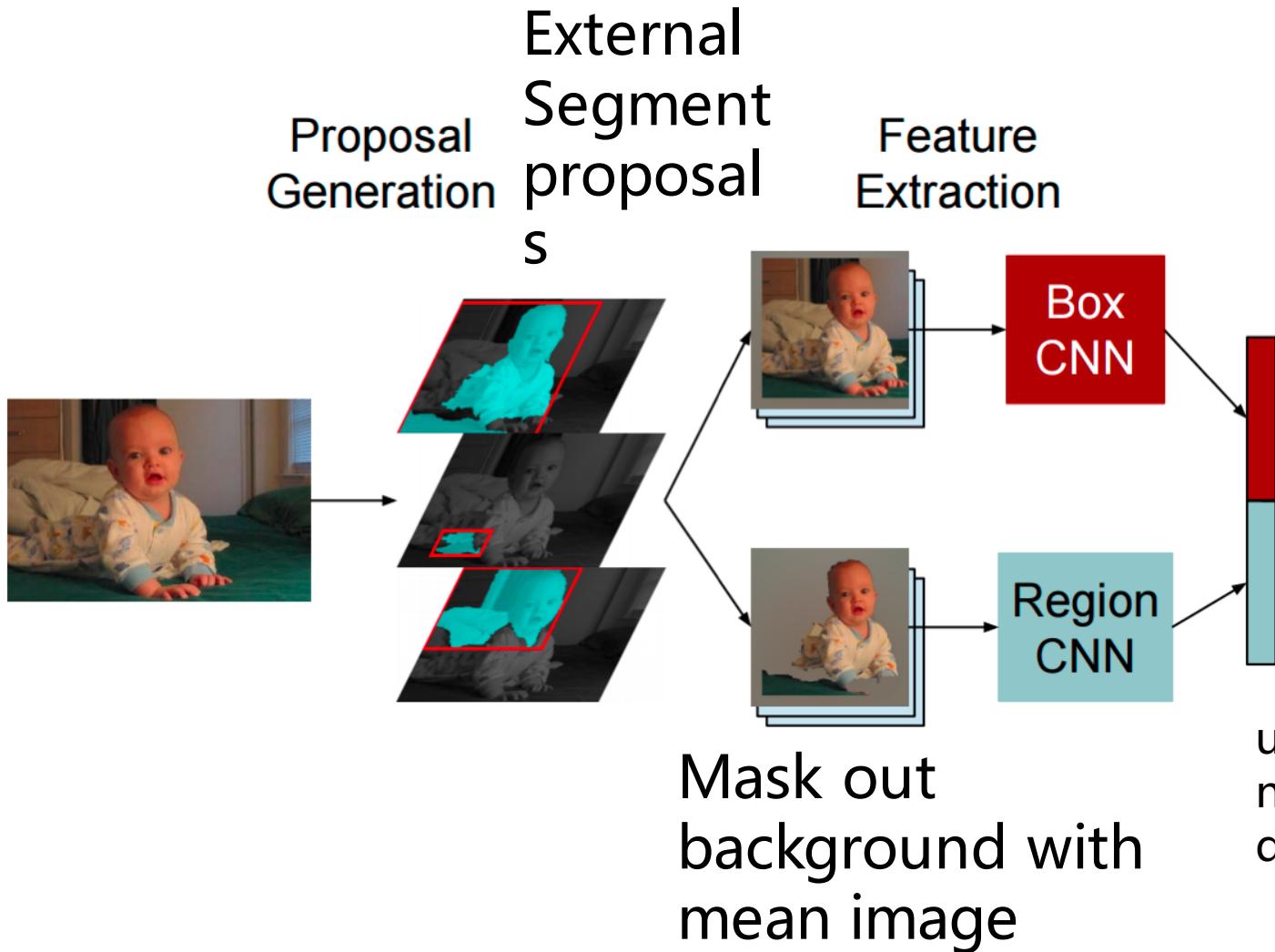
Instance Segmentation

Similar to R-CNN,
but with
segments



Hariharan et al, "Simultaneous Detection and Segmentation" , ECCV 2014

Instance Segmentation

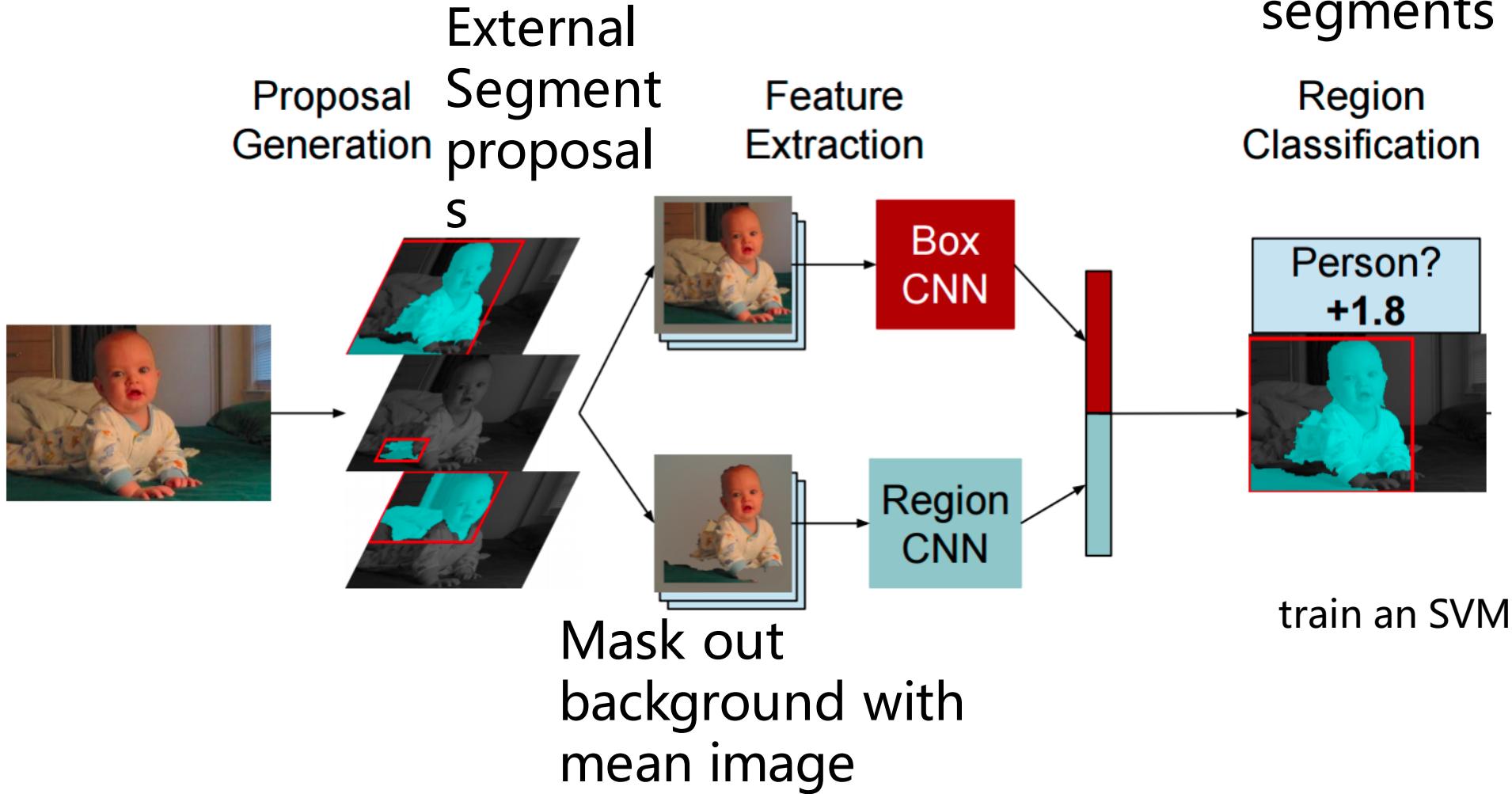


Similar to R-CNN,
but with
segments

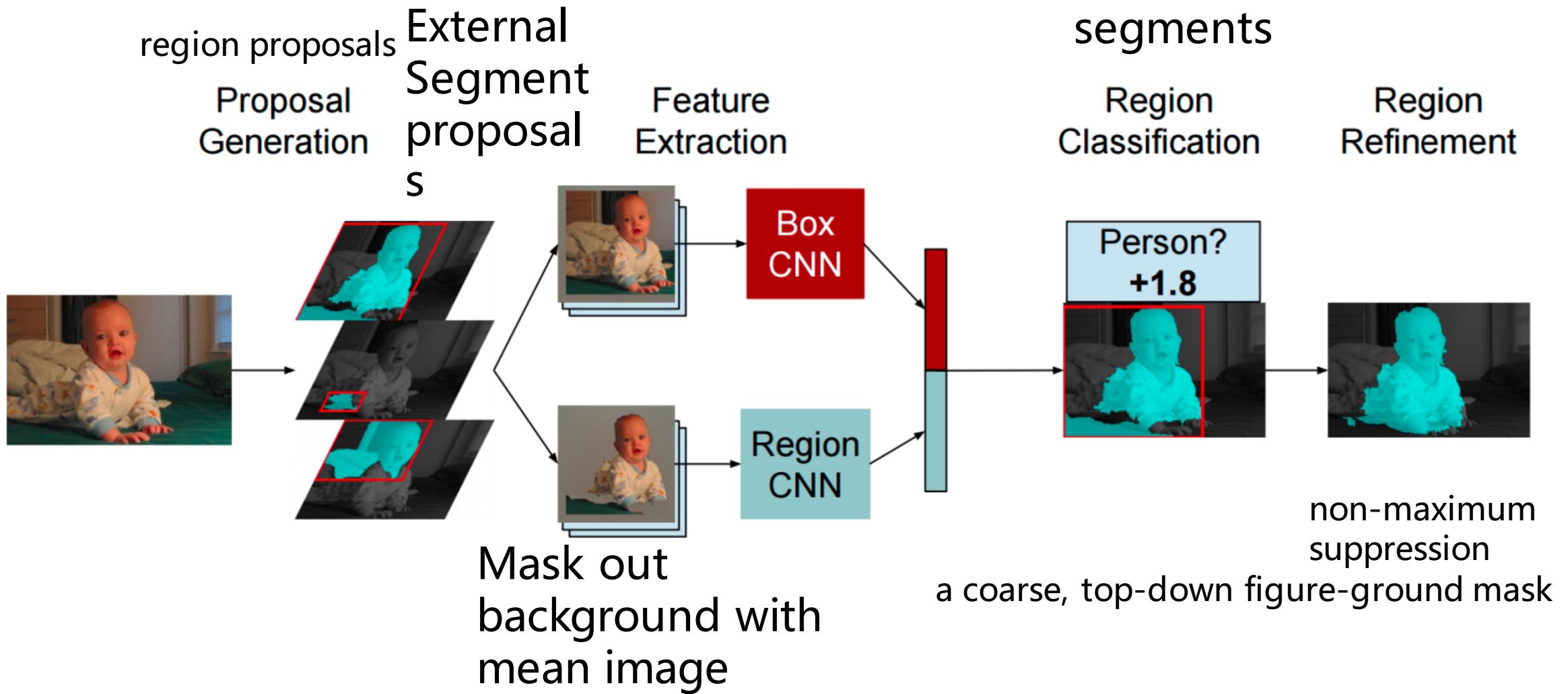
Disjoint but same architecture
different weights but train together

using separate networks where each network is finetuned for its respective role dramatically improves performance.

Instance Segmentation



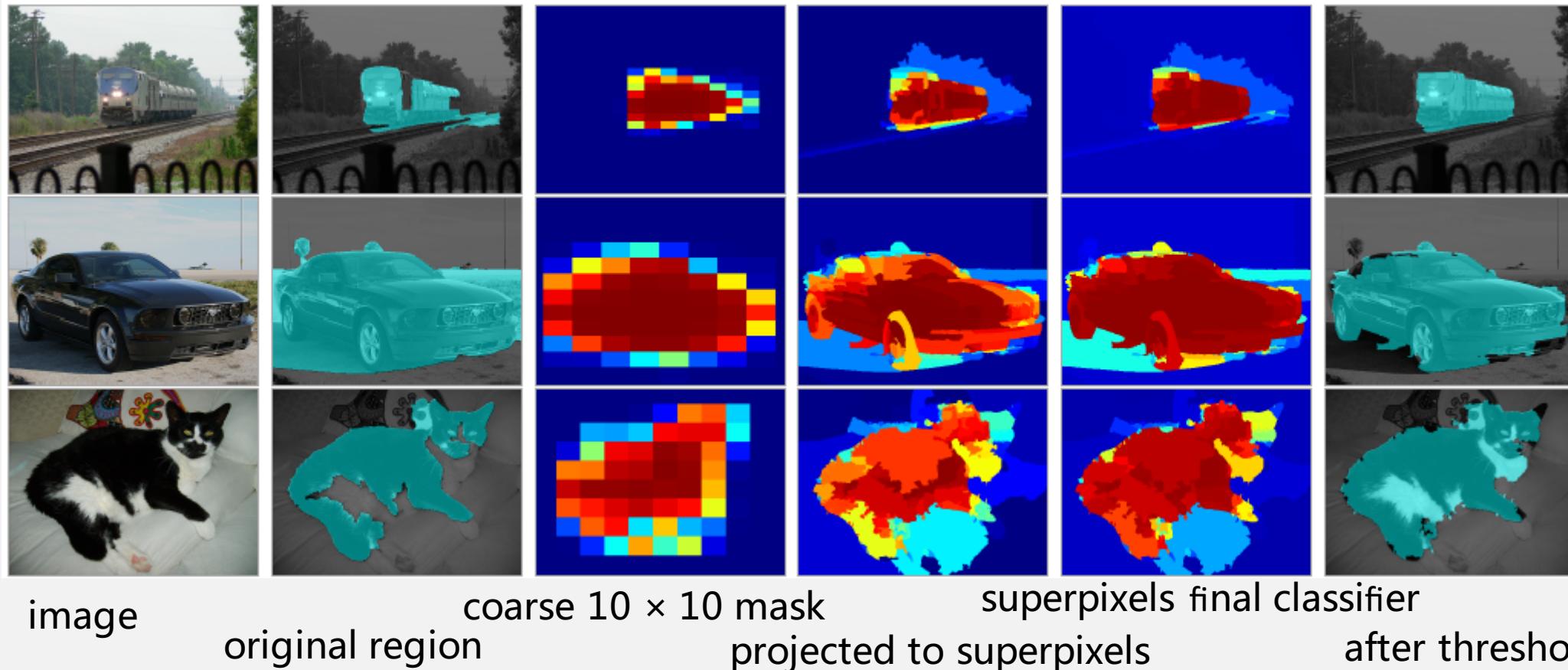
Instance Segmentation



region refinement

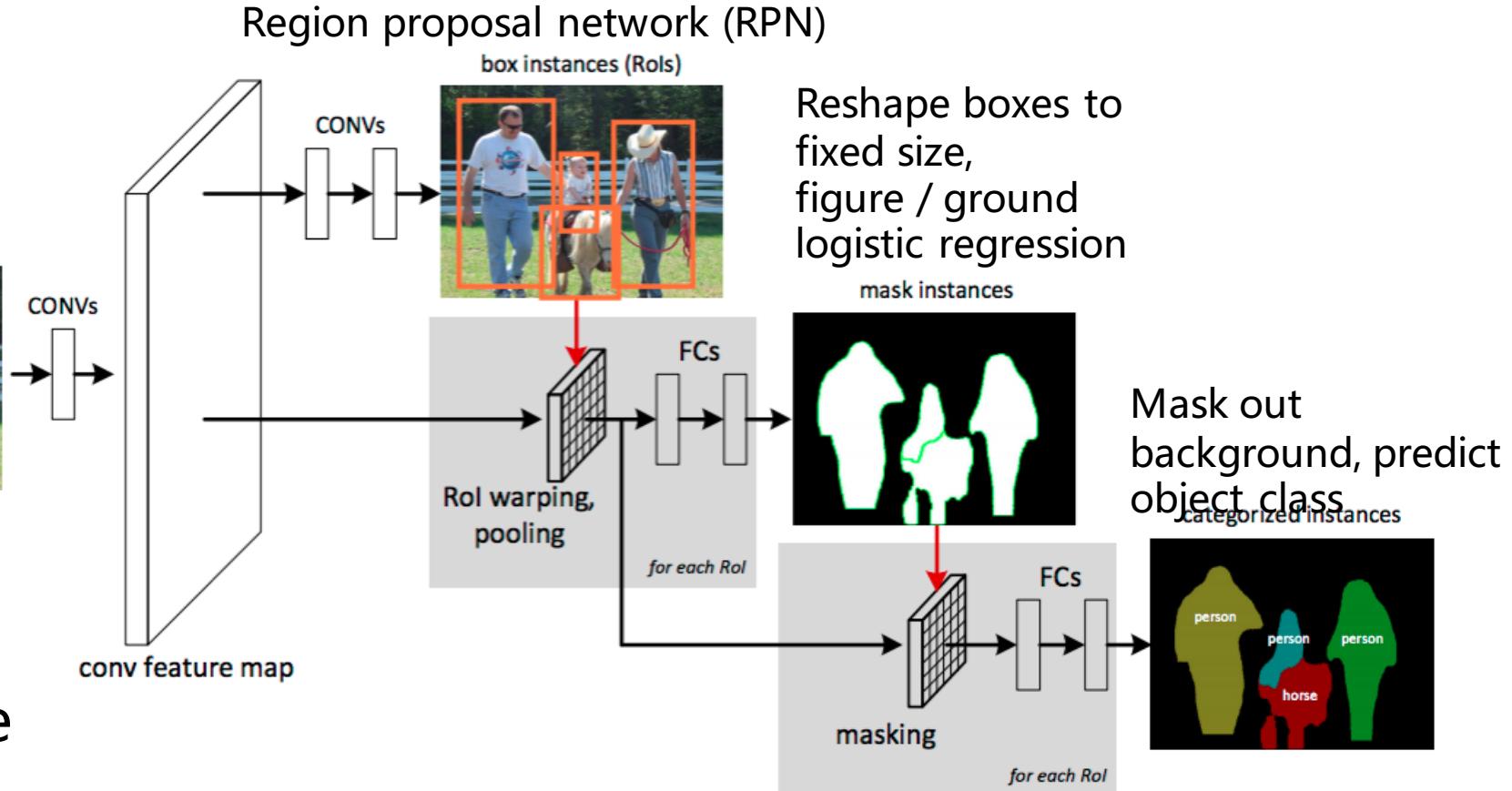
examples

Refinement uses top-down category specific information to fill in the body of the train and the cat and remove the road from the car.



Instance Segmentation: Cascades

Similar to
Faster R-CNN

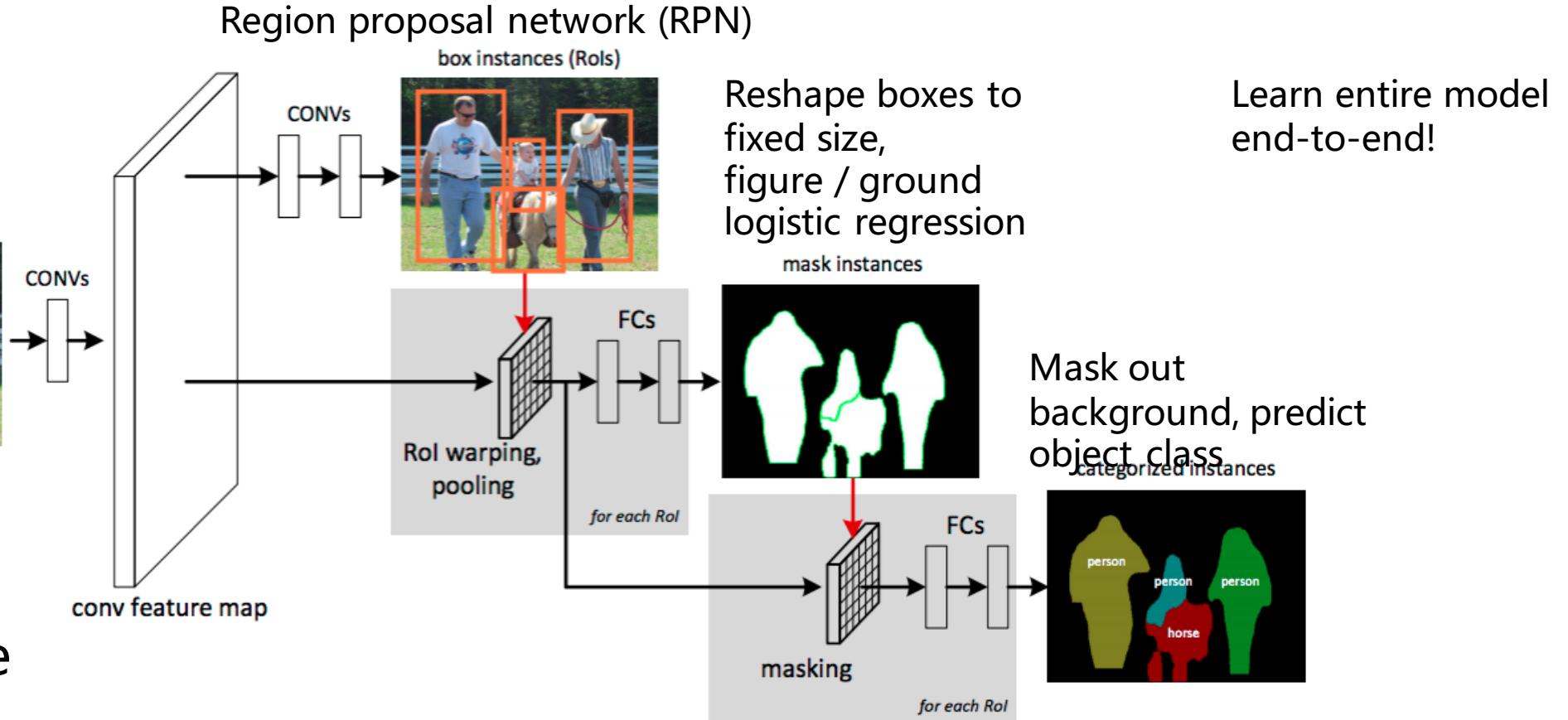


Won COCO
2015 challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades" ,
arXiv 2015

Instance Segmentation: Cascades

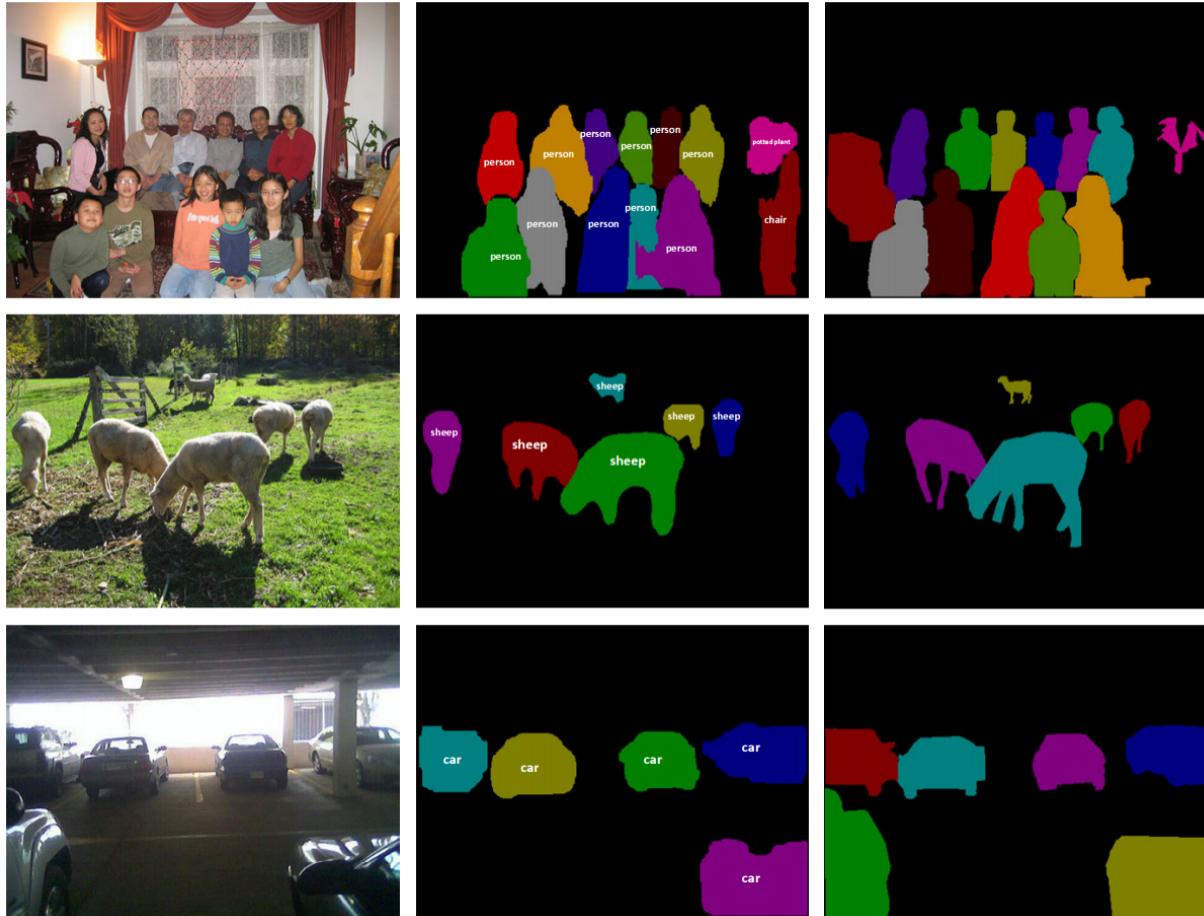
Similar to
Faster R-CNN



Won COCO
2015 challenge
(with ResNet)

Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades" ,
arXiv 2015

Instance Segmentation: Cascades



Dai et al, "Instance-aware Semantic Segmentation via Multi-task Network Cascades", arXiv 2015

**Predicti
ons**

**Ground
truth**

64

Segmentation Overview

Semantic segmentation



Classify all pixels



Fully convolutional models
downsample then upsample

Learnable upsampling
fractionally strided convolution

Skip connections

Detect instance, generate mask



Similar pipelines to object detection

Instance Segmentation



THANK
YOU!

PRESNTED BY DeepLearningPKU