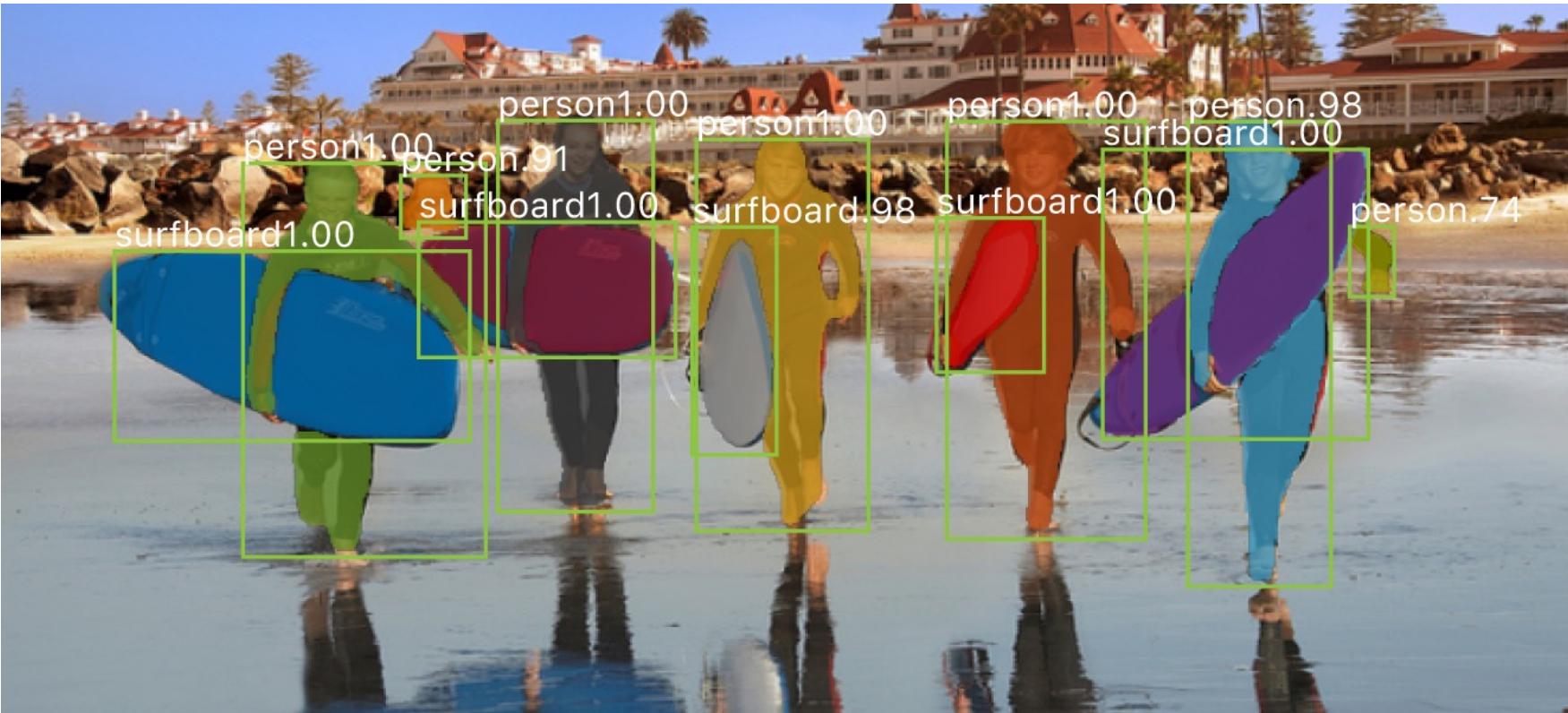


Lecture 13

Object Detection and Image Segmentation



Overview

- Visual Recognition Tasks
- Methodology: ConvNets, Fully Conv Nets
- Semantic Segmentation
- Object Detection
- Instance Segmentation

Visual Recognition Tasks

Classification:

output: class label



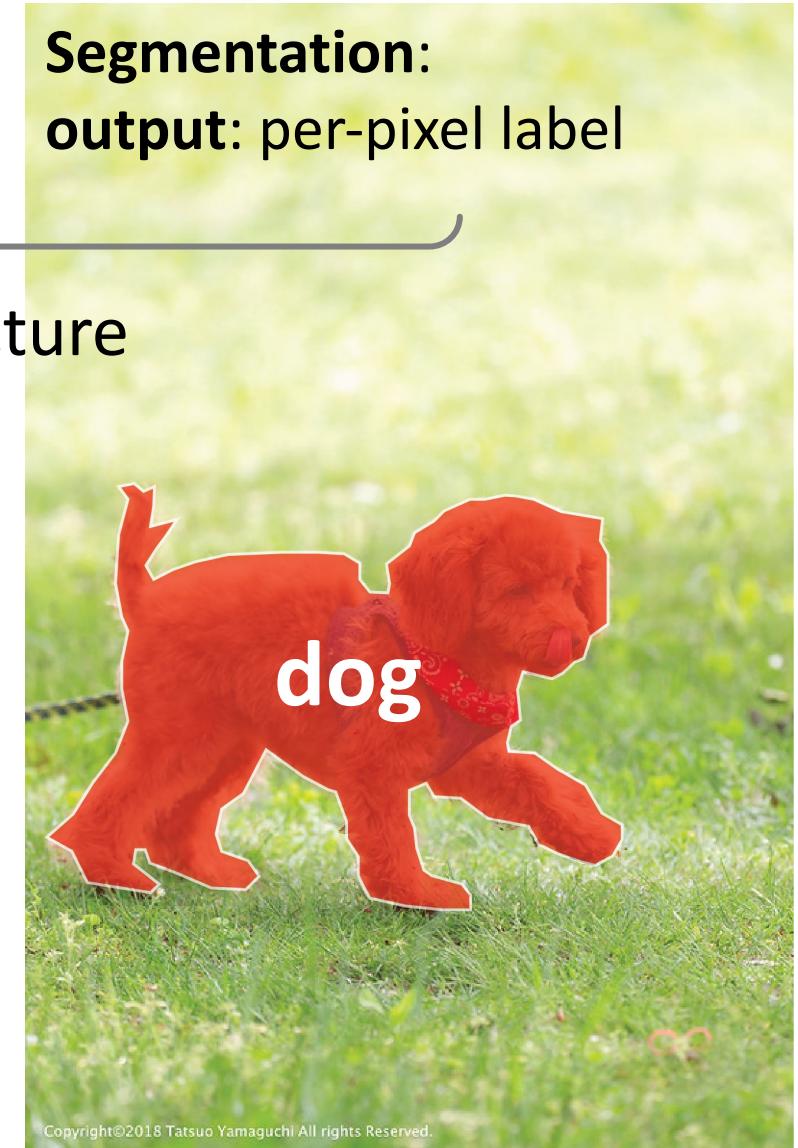
Detection:

output: class label + box



Segmentation:

output: per-pixel label

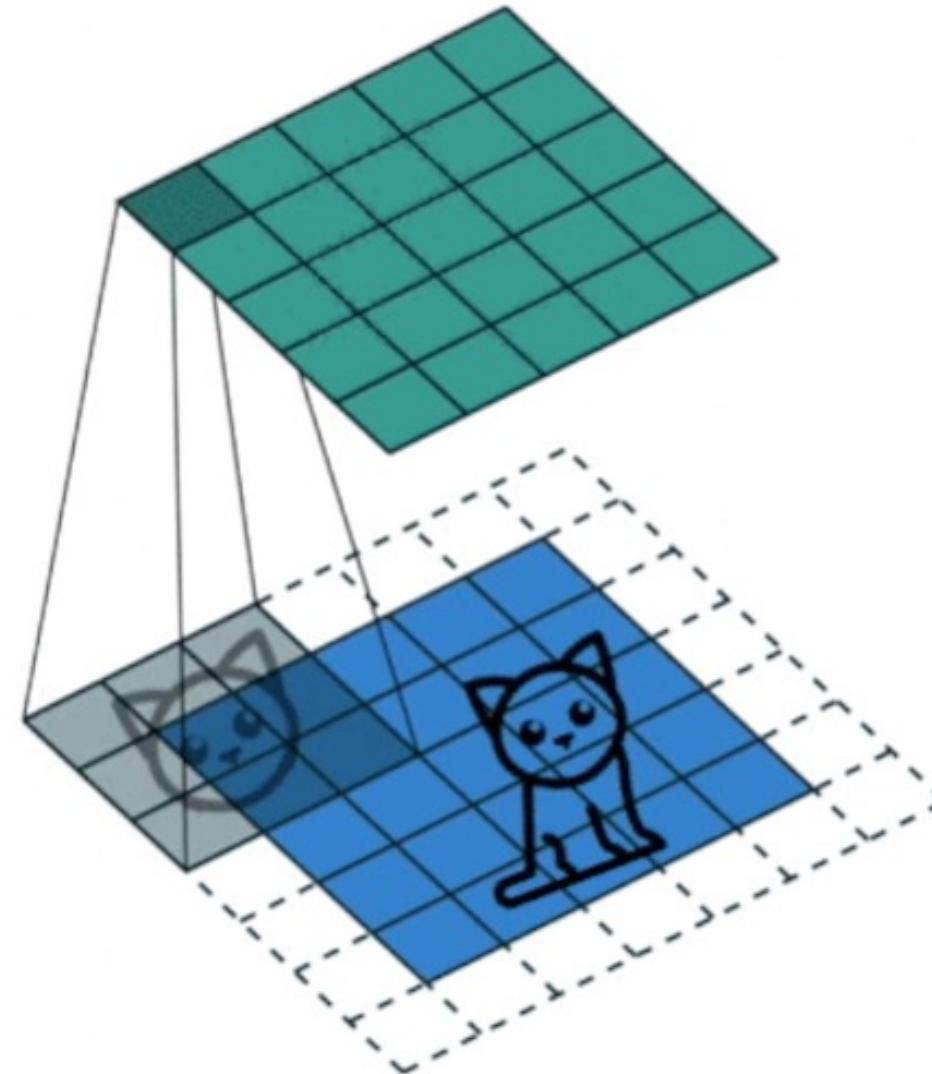


this lecture

Methodology: ConvNets, Fully Conv Nets

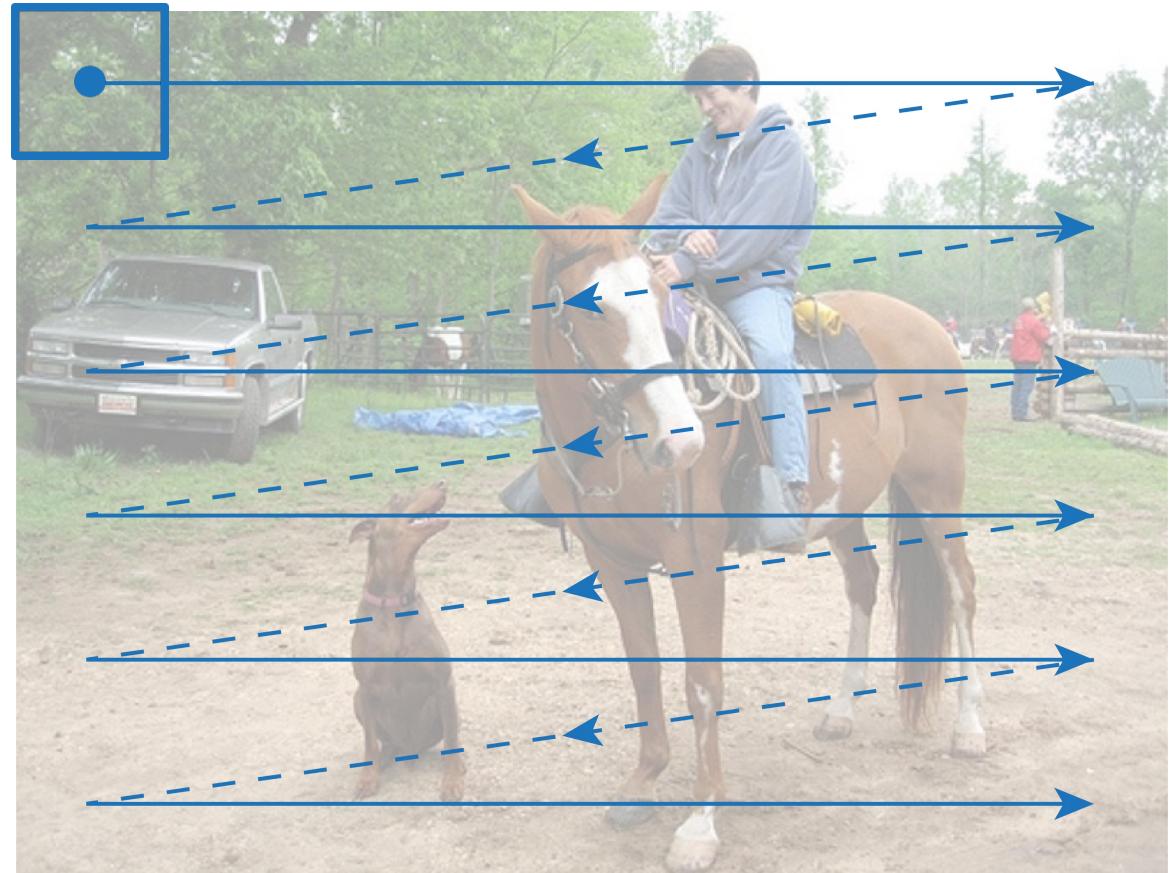
Convolution is template matching

- sliding window
- abstract templates



Convolution is template matching

- sliding window
- abstract templates



ConvNets Are Fully Convolutional

Screenshot from: "What's Wrong With Deep Learning?" Yann LeCun, 2015

Multiple Character Recognition [Matan et al 1992] Y LeCun

- SDNN: Space Displacement Neural Net
- Also known as “replicated convolutional net”, or just ConvNet
 - (are we going to call this “fully convolutional net” now?)
- There is no such thing as a “fully connected layer”
- they are actually convolutional layers with 1×1 convolution kernels.

The diagram illustrates the architecture of a Space Displacement Neural Network (SDNN) for multiple character recognition. It shows two main components: a "Single Character Recognizer" and an "SDNN".

The "Single Character Recognizer" on the left processes a single character input, represented by a green "u" shape, through a series of convolutional layers. These layers are shown as colored blocks (blue, red, pink, grey) with arrows indicating the flow of data. The output of these layers is labeled "Single Character Recognizer".

An arrow points from the "Single Character Recognizer" to the "SDNN" on the right. The "SDNN" consists of multiple parallel "Single Character Recognizer" units, each processing a different character input. The inputs are shown as blue, red, pink, and grey "cursive" shapes. The outputs of these units are combined to recognize the entire word "cursive".

ConvNets Are Fully Convolutional

Screenshot from: "What's Wrong With Deep Learning?" Yann LeCun, 2015

Multiple Character Recognition [Matan et al 1992] Y LeCun

- SDNN: Space Displacement Neural Net
- Also known as "replicated convolutional net", or just ConvNet
 - (are we going to call this "fully convolutional net" now?)
- There is no such thing as a "fully connected layer"
- they are actually convolutional layers with 1×1 convolution kernels.

The diagram illustrates the architecture of a Space Displacement Neural Network (SDNN) for multiple character recognition. It shows two stages: a "Single Character Recognizer" and the "SDNN".

Single Character Recognizer: On the left, a green "us" character is shown at the bottom. Above it, several colored blocks (blue, red, green, pink, grey) represent feature maps. A green line connects the "us" character to these feature maps. A red arrow points from this stage to the next.

SDNN: On the right, the output of the "Single Character Recognizer" is processed by the SDNN. The SDNN consists of multiple layers of feature maps. Red arrows show the flow of data from the input feature maps through the SDNN layers. The final output is a green box containing the word "cursive".

ConvNets Are Fully Convolutional

Screenshot from: "What's Wrong With Deep Learning?" Yann LeCun, 2015

Multiple Character Recognition [Matan et al 1992] Y LeCun

- SDNN: Space Displacement Neural Net
- Also known as "replicated convolutional net", or just ConvNet
 - (are we going to call this "fully convolutional net" now?)
- There is no such thing as a "fully connected layer"
- they are actually convolutional layers with 1x1 convolution kernels.

The diagram illustrates the architecture of a Space Displacement Neural Network (SDNN) for multiple character recognition. It shows two stages: a 'Single Character Recognizer' and an 'SDNN'.

Single Character Recognizer: On the left, a green 'us' character is shown at the bottom. Above it is a stack of colored feature maps (green, pink, grey, blue) representing the output of a single character recognizer. A green arrow points from this stage to the right.

SDNN: On the right, a red arrow points from the 'Single Character Recognizer' stage to a more complex network. This network consists of several stacked layers of colored feature maps. Red arrows indicate the receptive fields of the neurons in the top layer, which are centered on a green 'cursive' character at the bottom. The text 'SDNN' is written below the network diagram.

ConvNets Are Fully Convolutional

Screenshot from: "What's Wrong With Deep Learning?" Yann LeCun, 2015

Multiple Character Recognition [Matan et al 1992] Y LeCun

- SDNN: Space Displacement Neural Net
- Also known as "replicated convolutional net", or just ConvNet
 - (are we going to call this "fully convolutional net" now?)
- There is no such thing as a "fully connected layer"
- they are actually convolutional layers with 1x1 kernels

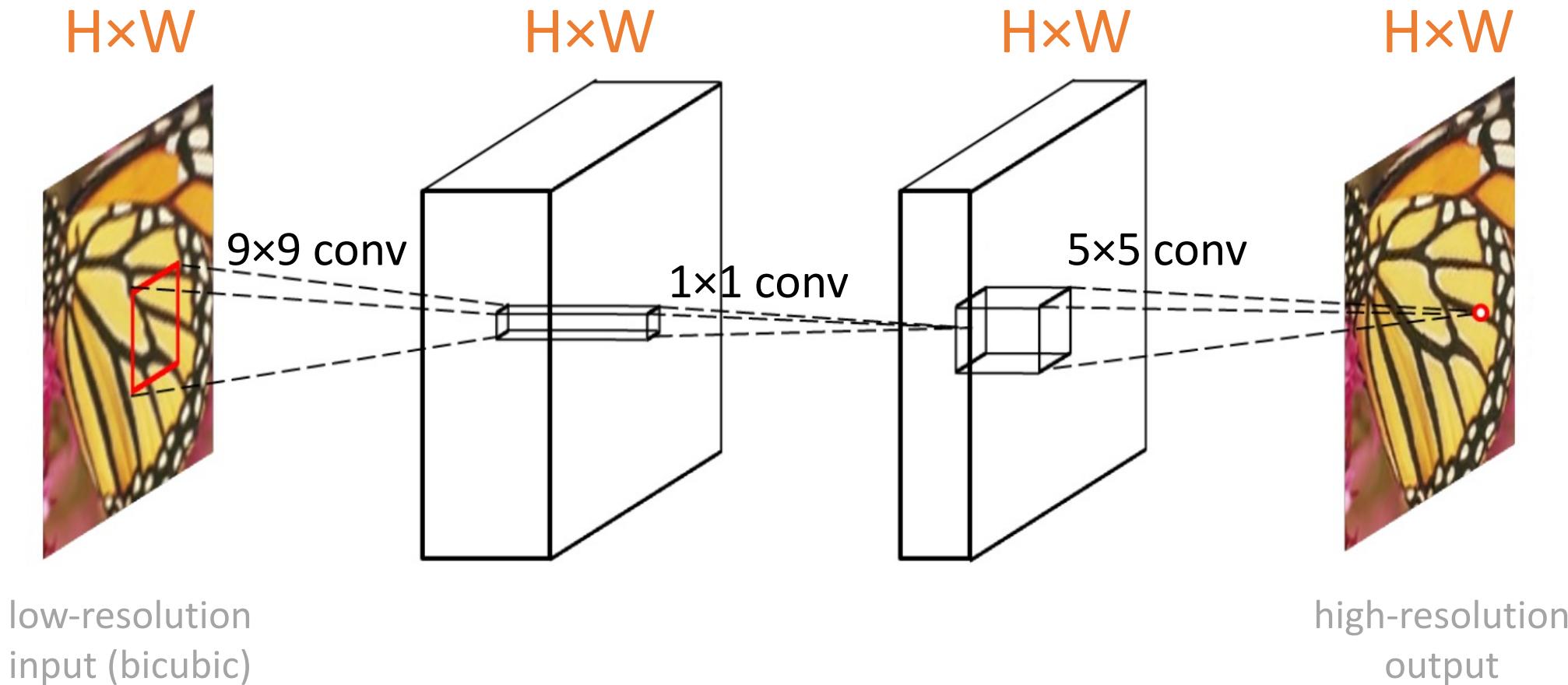
The diagram illustrates the evolution of character recognition architecture. On the left, a 'Single Character Recognizer' is shown as a stack of colored convolutional layers processing a handwritten digit '4'. An arrow points to the right, labeled 'SDNN', where the same digit is processed by a series of stacked convolutional layers. A red box highlights the text '(are we going to call this "fully convolutional net" now?)' in the list above, and another red box encloses the definition of a Fully Convolutional Network (FCN) on the right.

Fully Convolutional Net (FCN):
A network in which all layers are convolutional. [Long et al. CVPR 2015]

ConvNets Are Fully Convolutional

Example: Super-Resolution CNN (SRCNN)

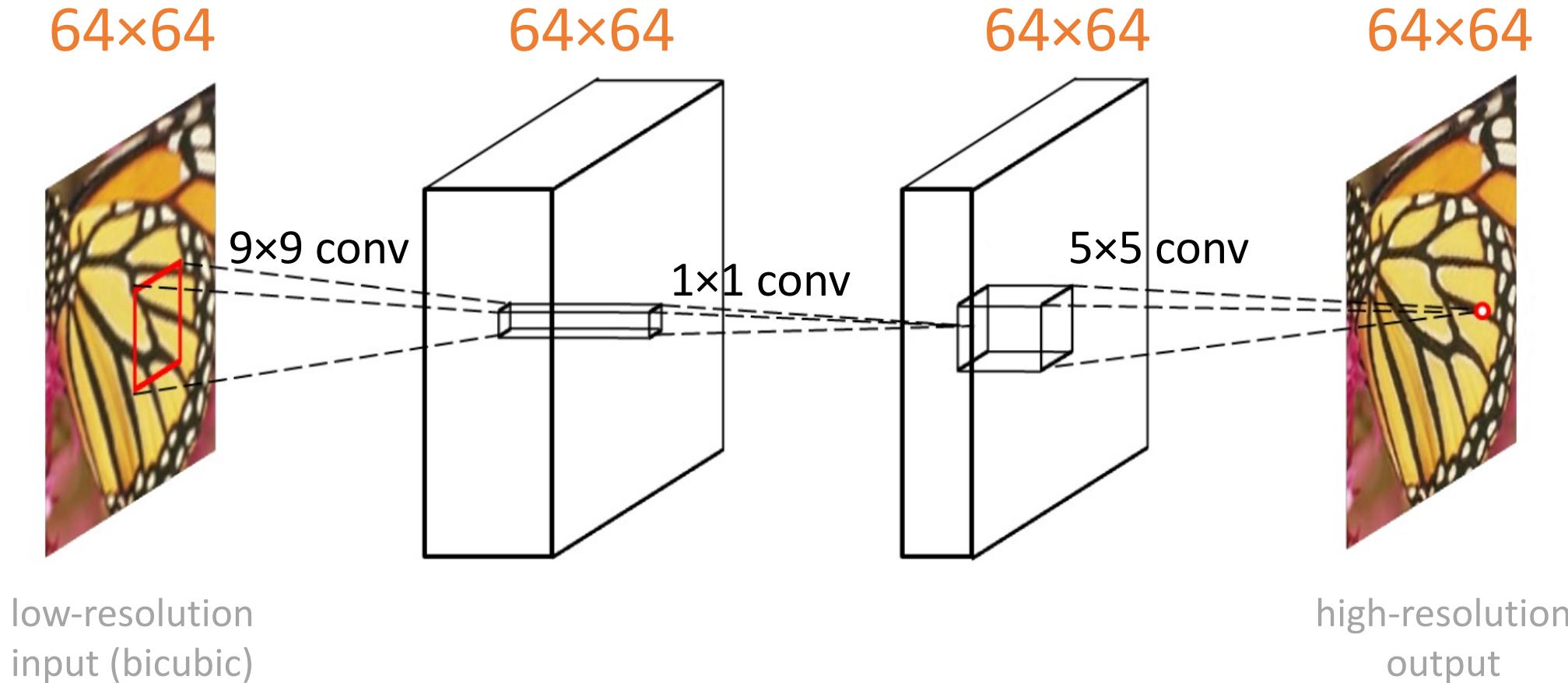
- arbitrary input size
- proportional output size



ConvNets Are Fully Convolutional

Example: Super-Resolution CNN (SRCNN)

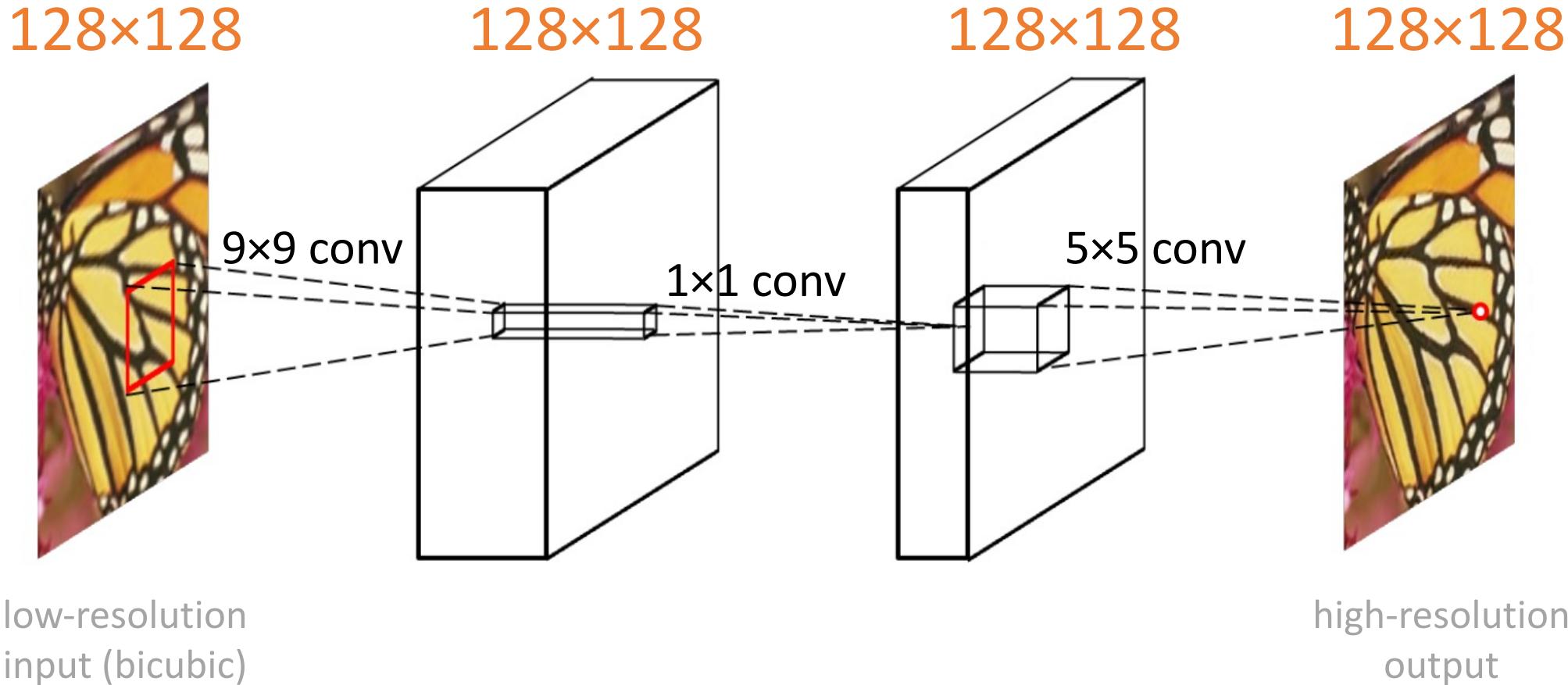
- arbitrary input size
- proportional output size



ConvNets Are Fully Convolutional

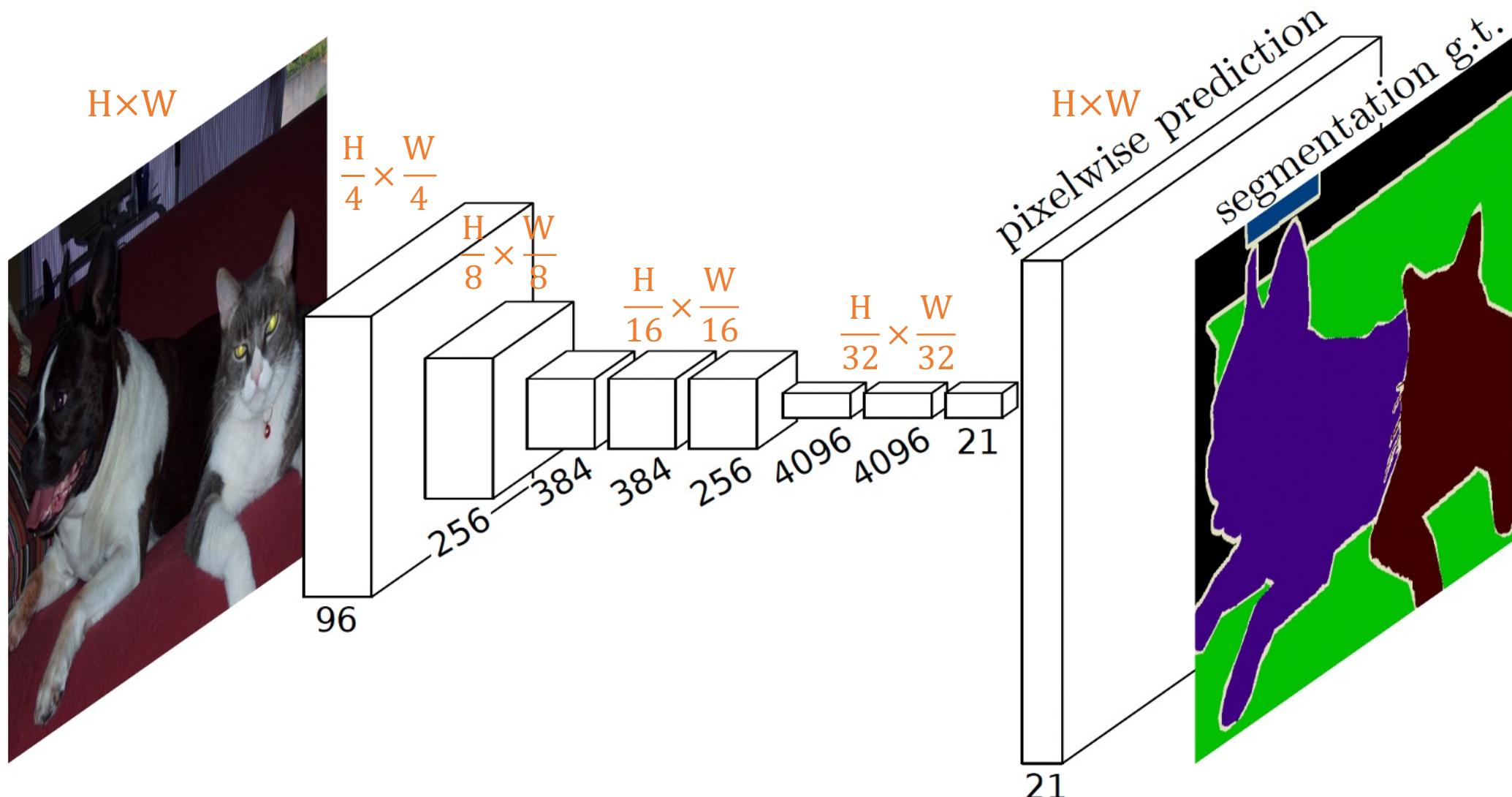
Example: Super-Resolution CNN (SRCNN)

- arbitrary input size
- proportional output size



ConvNets Are Fully Convolutional

Example: Fully Convolutional Network (FCN)

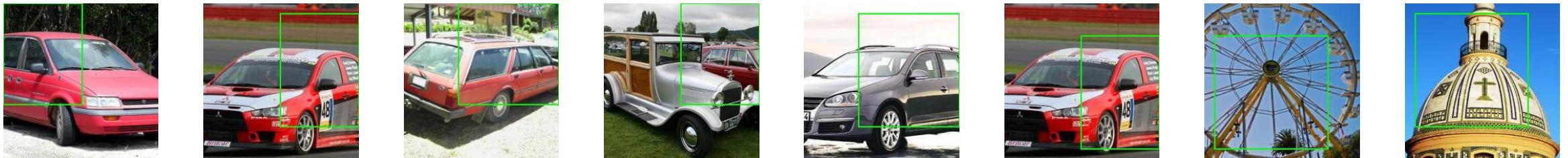
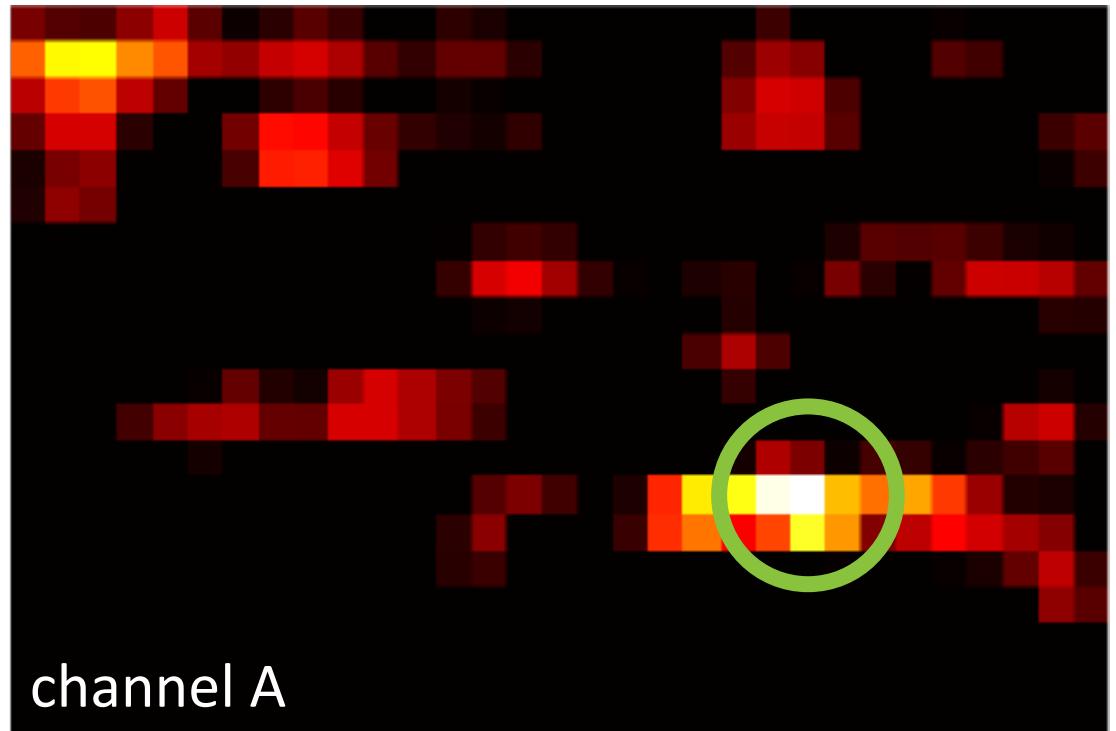


ConvNets Are Fully Convolutional

input $H \times W$



conv5's output $\frac{H}{32} \times \frac{W}{32}$



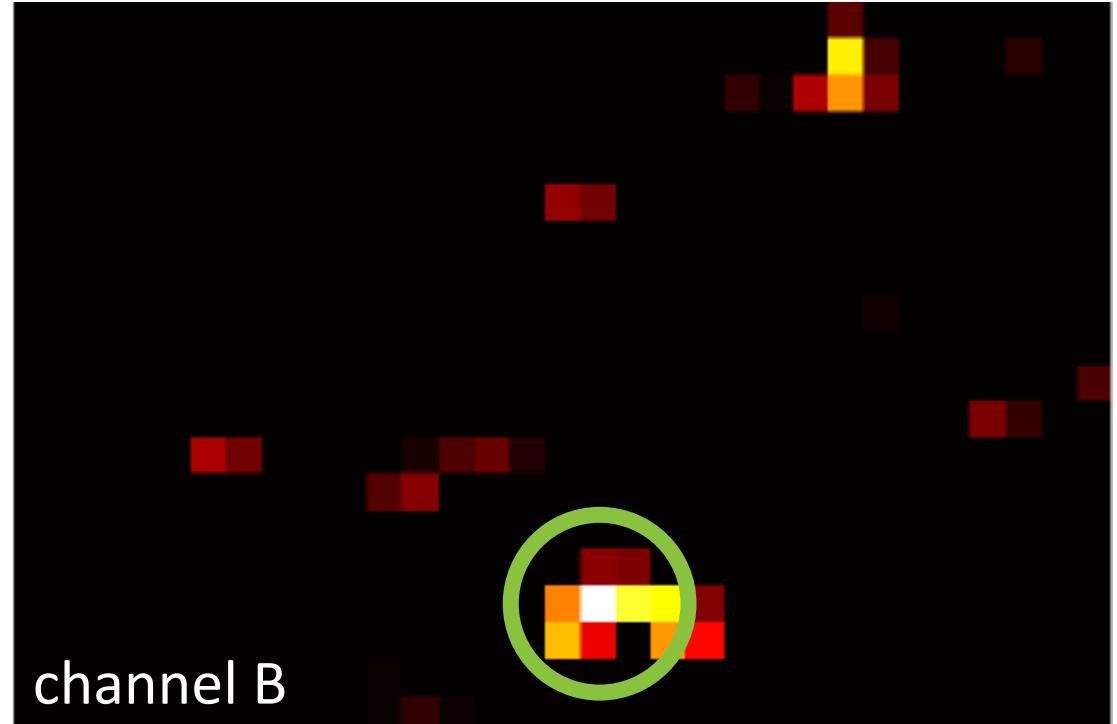
strongest activations of this channel

ConvNets Are Fully Convolutional

input $H \times W$



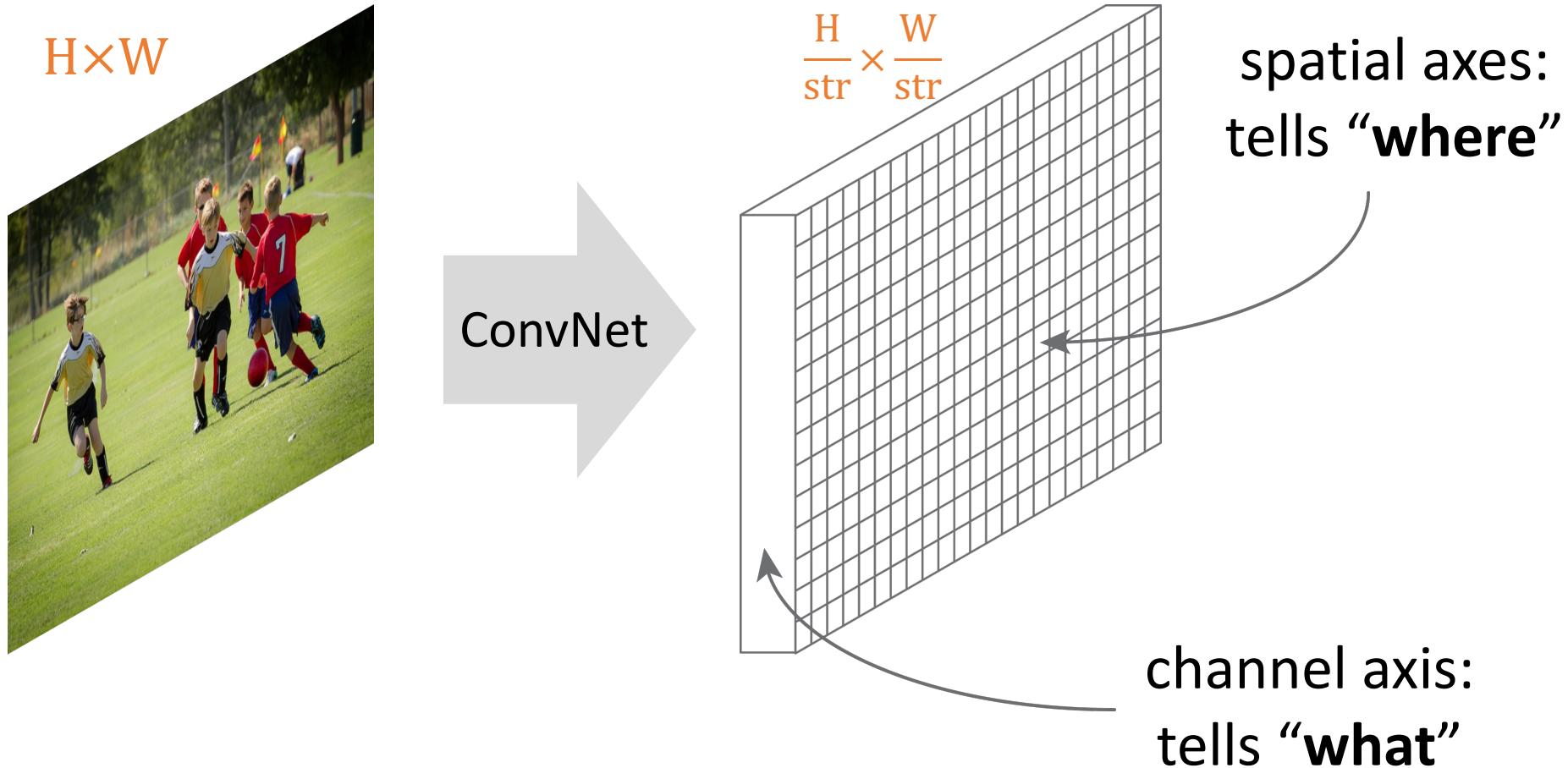
conv5's output $\frac{H}{32} \times \frac{W}{32}$



strongest activations of this channel

ConvNets Are Fully Convolutional

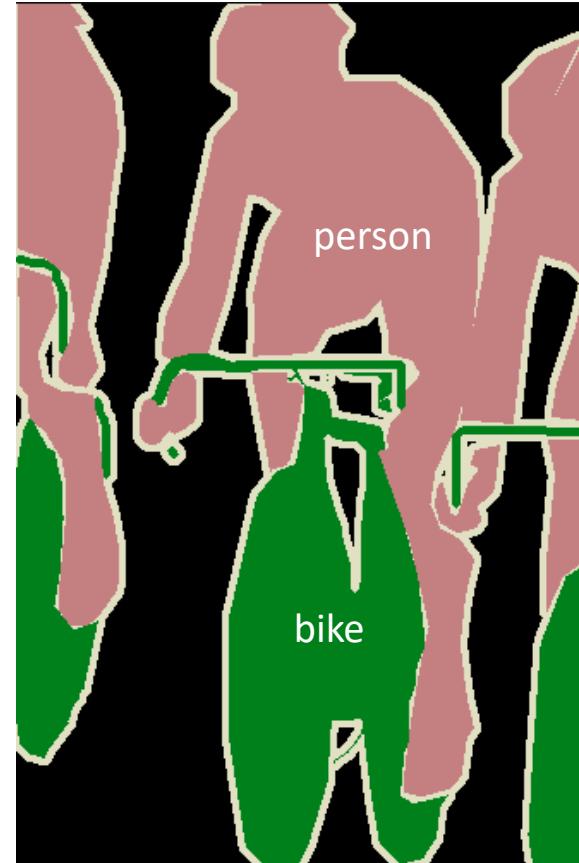
- Conceptually, detection and segmentation are solved in a “fully convolutional” fashion.



Semantic Segmentation

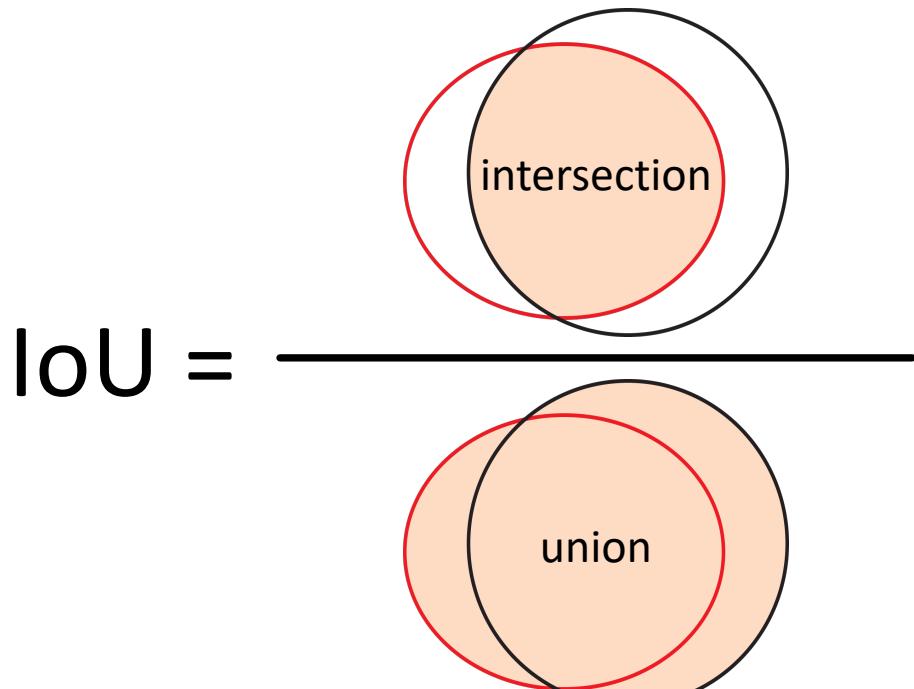
Semantic Segmentation: Problem Definition

- **Task:** Label every pixel
 - Historically: don't differentiate instances

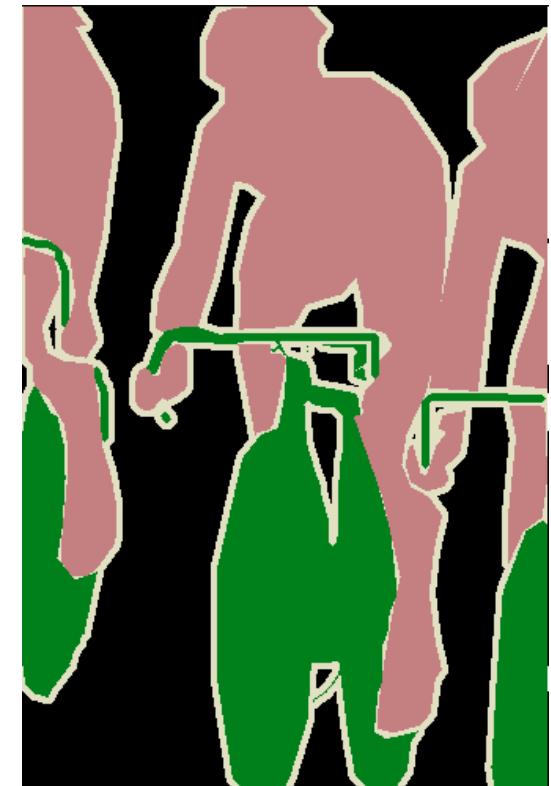


Semantic Segmentation: Evaluation

- Intersection over Union (IoU), averaged over classes



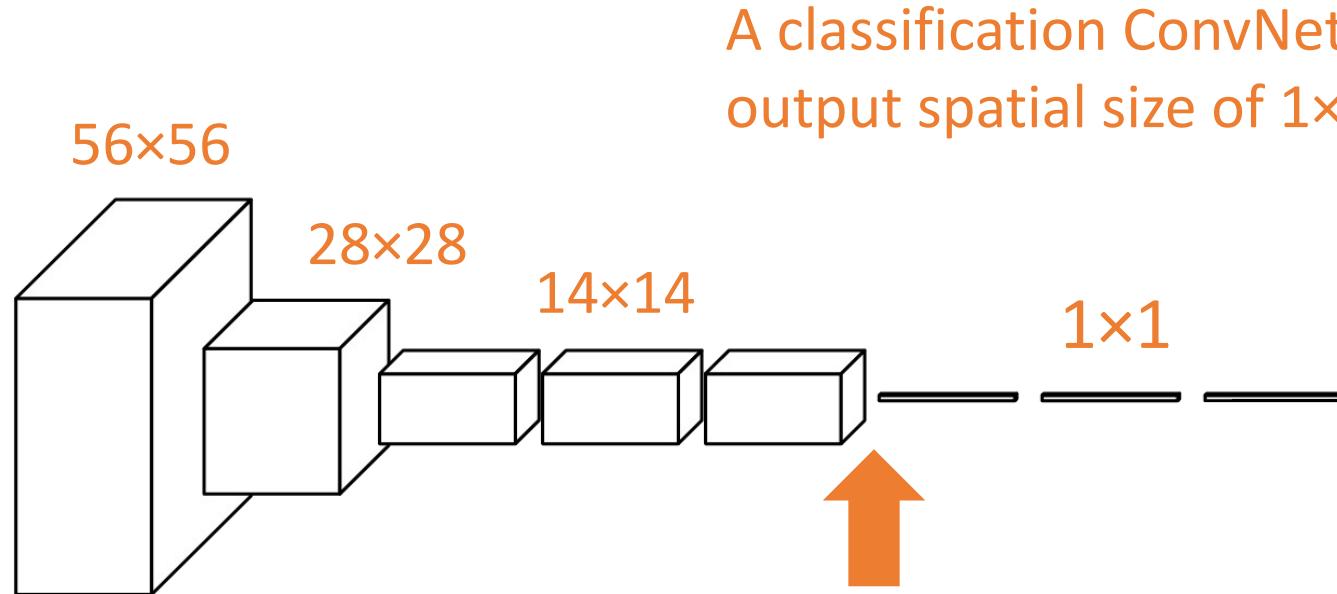
prediction



ground truth

Fully Convolutional Network (FCN)

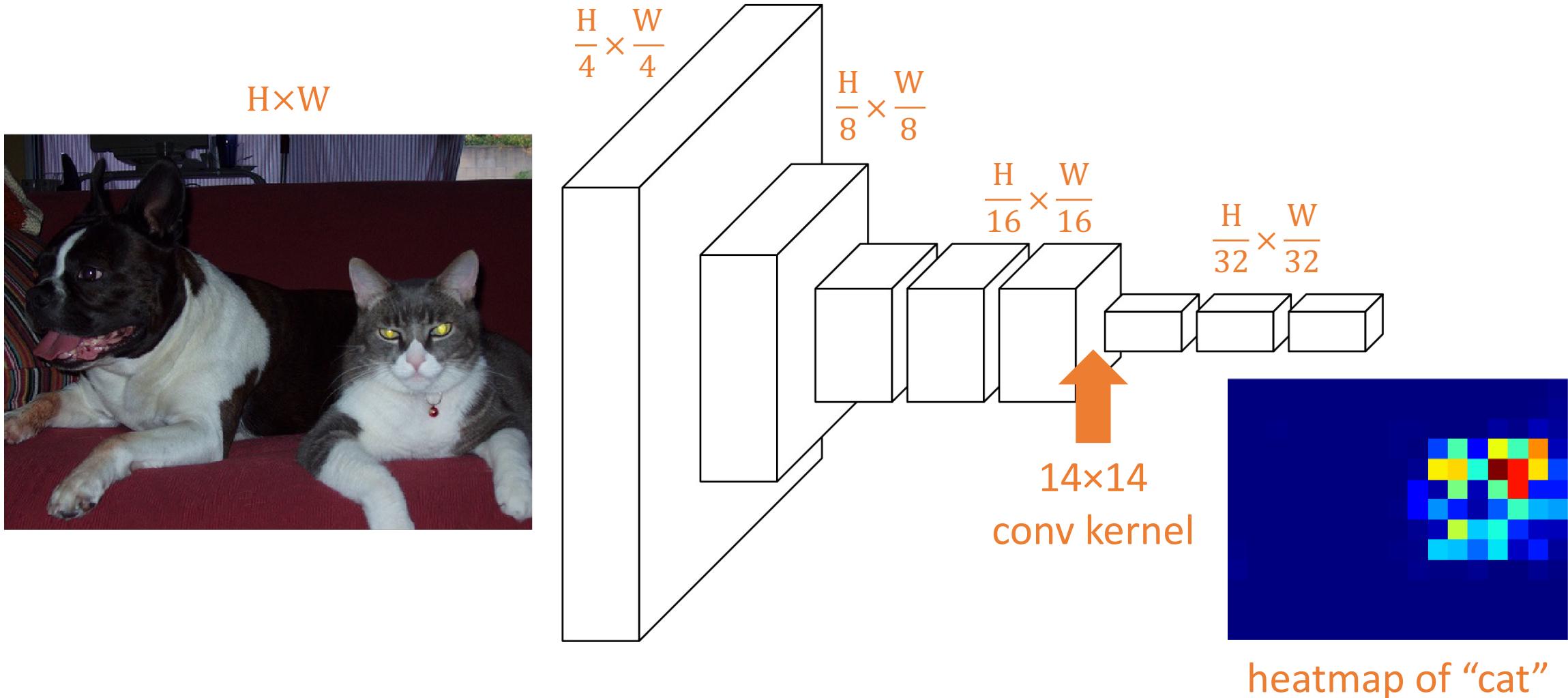
- Repurpose a classification ConvNet for semantic segmentation



This fully-connected (fc) layer can be viewed as a 14×14 conv kernel.

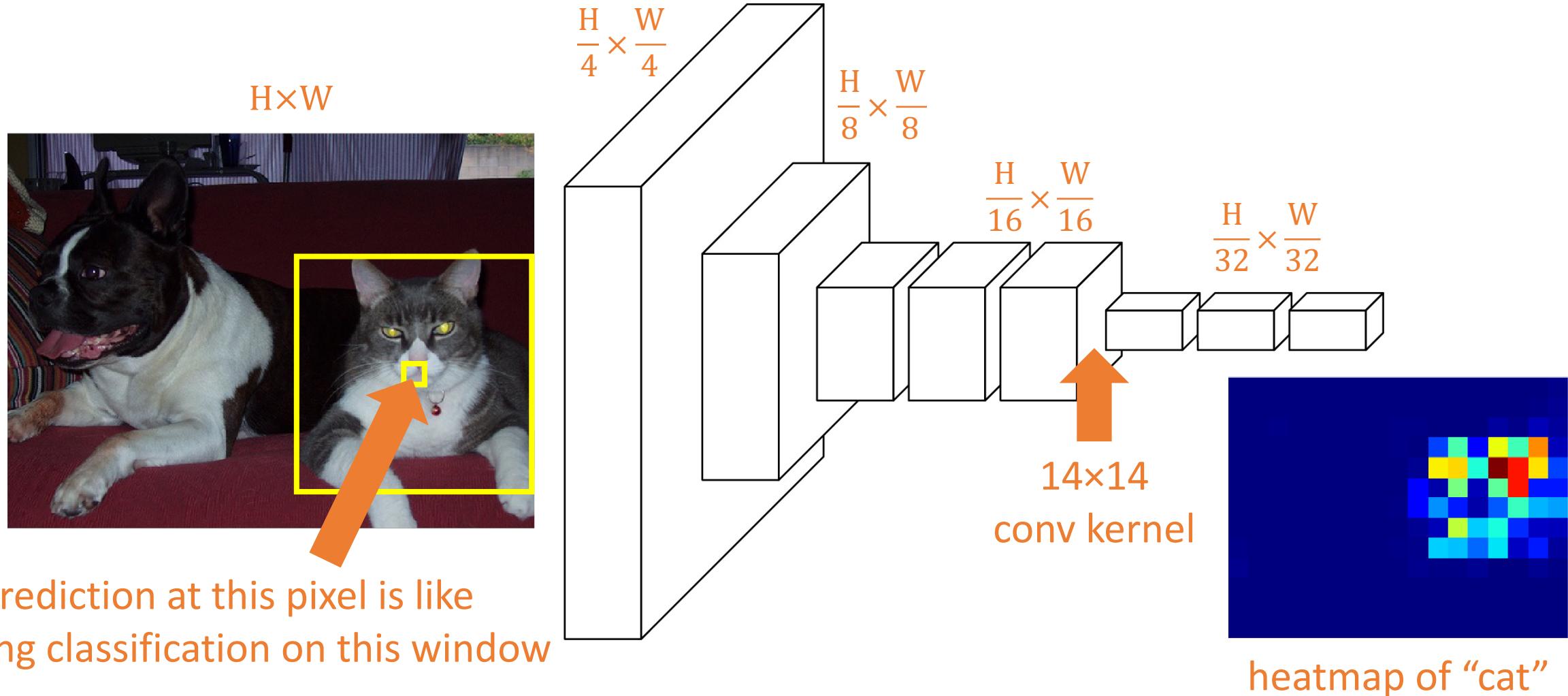
Fully Convolutional Network (FCN)

- Repurpose a classification ConvNet for semantic segmentation



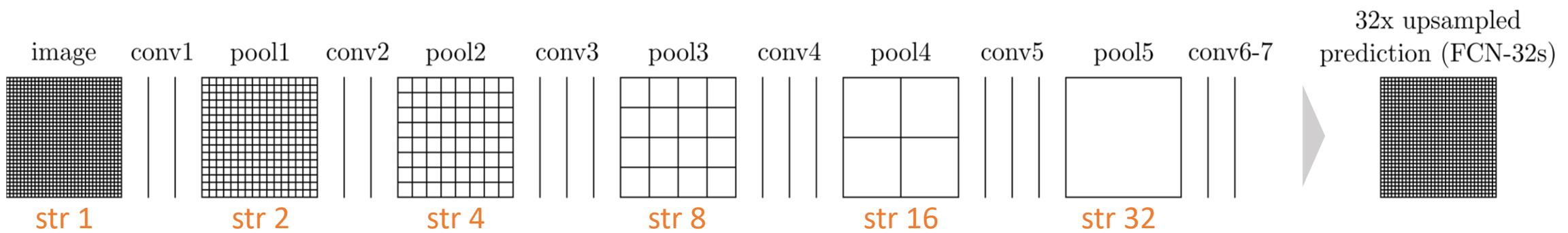
Fully Convolutional Network (FCN)

- Repurpose a classification ConvNet for semantic segmentation



Fully Convolutional Network (FCN)

- Upsampling: combining “where” and “what”

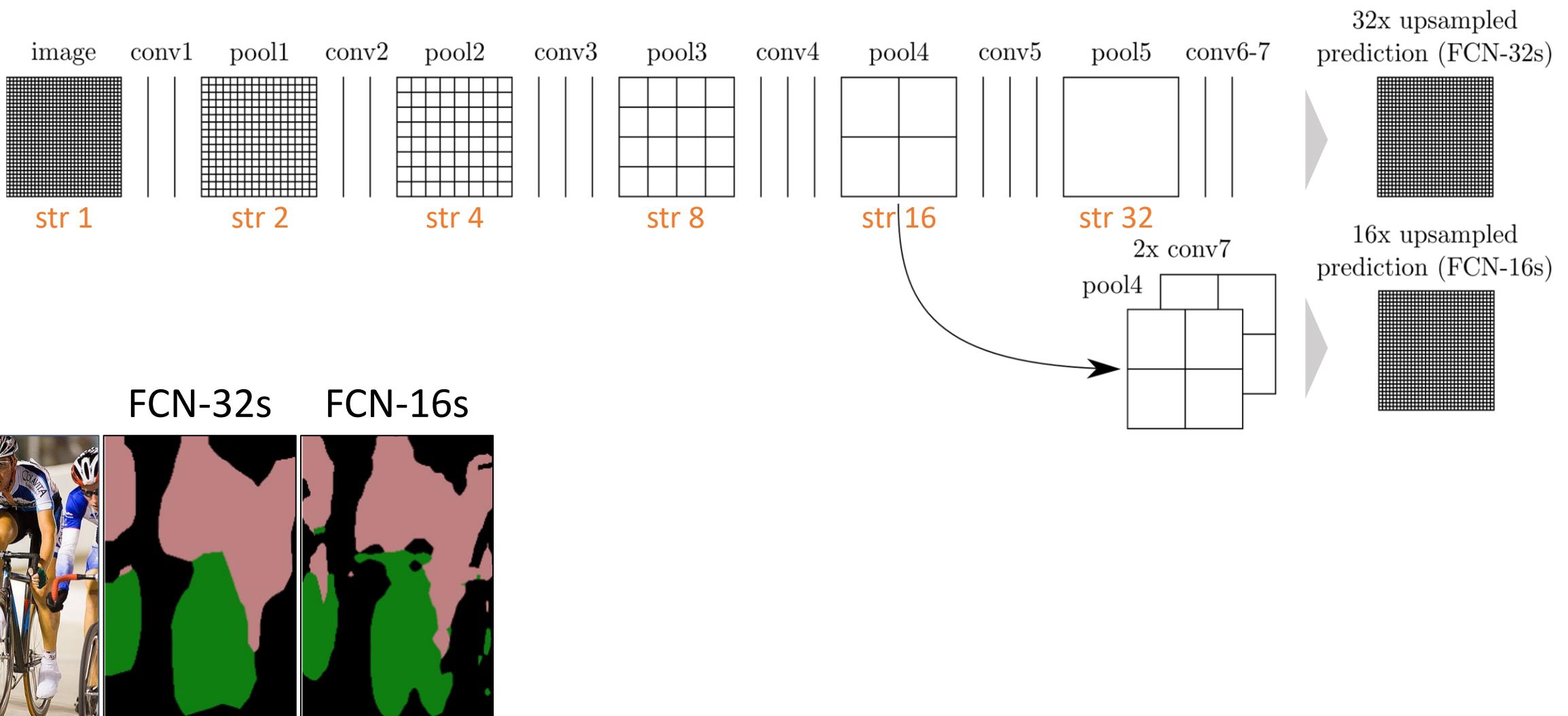


FCN-32s



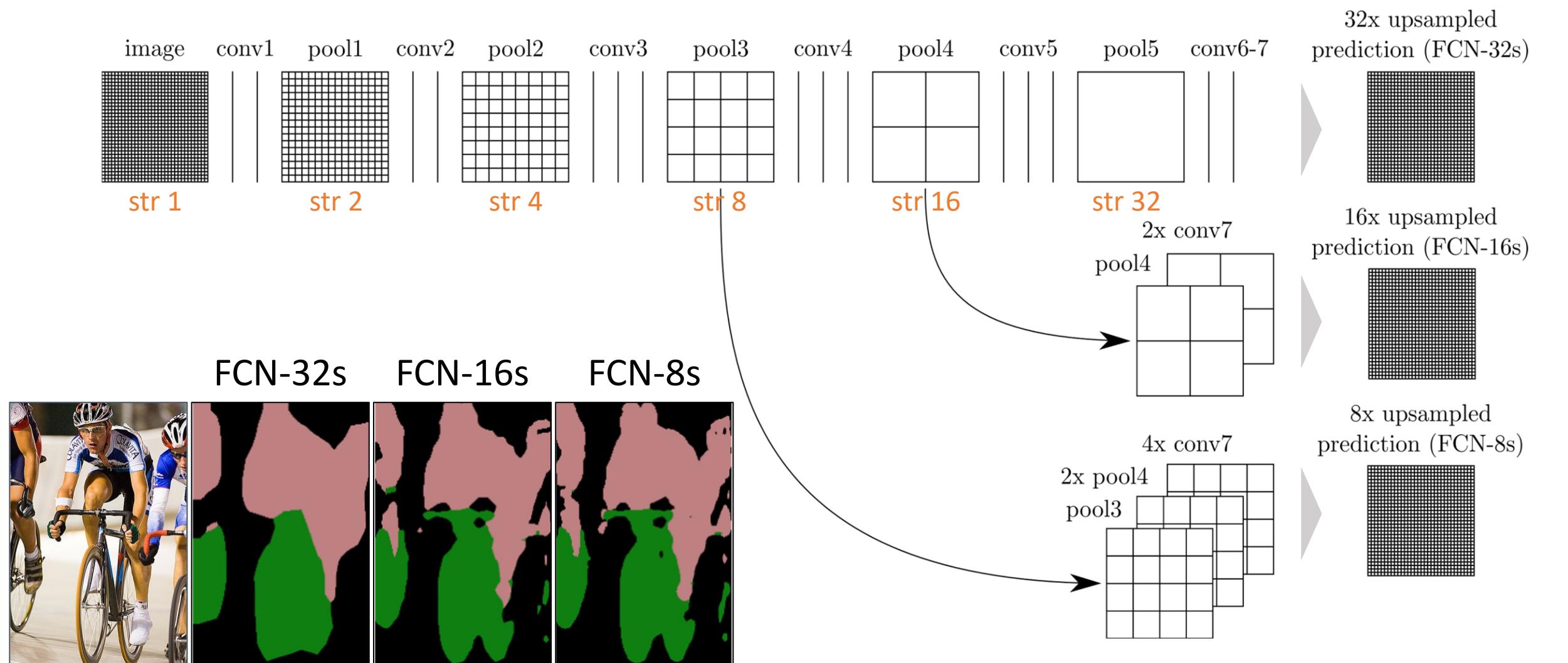
Fully Convolutional Network (FCN)

- Upsampling: combining “where” and “what”

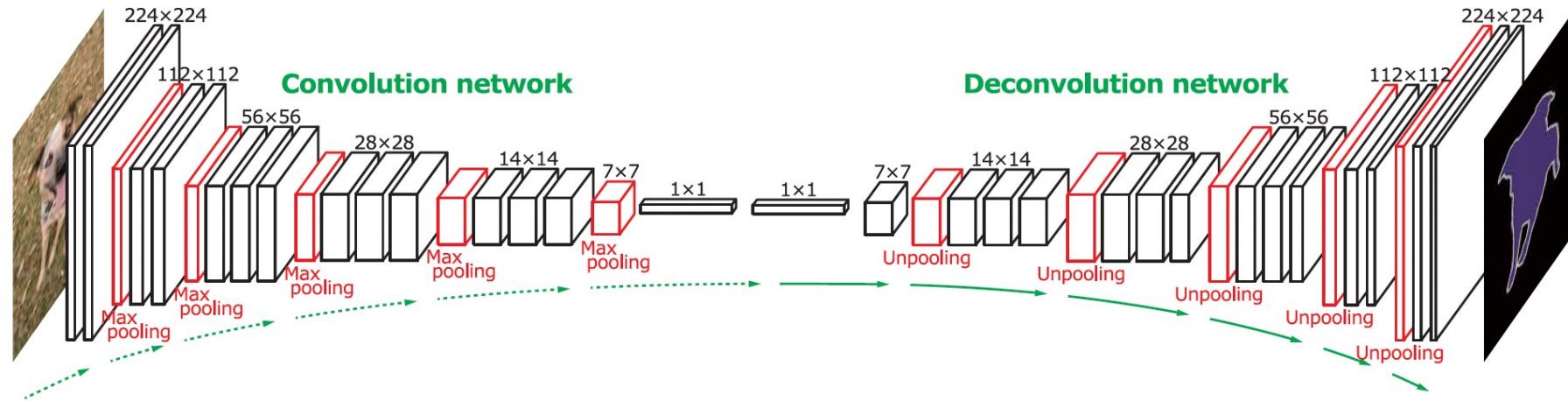
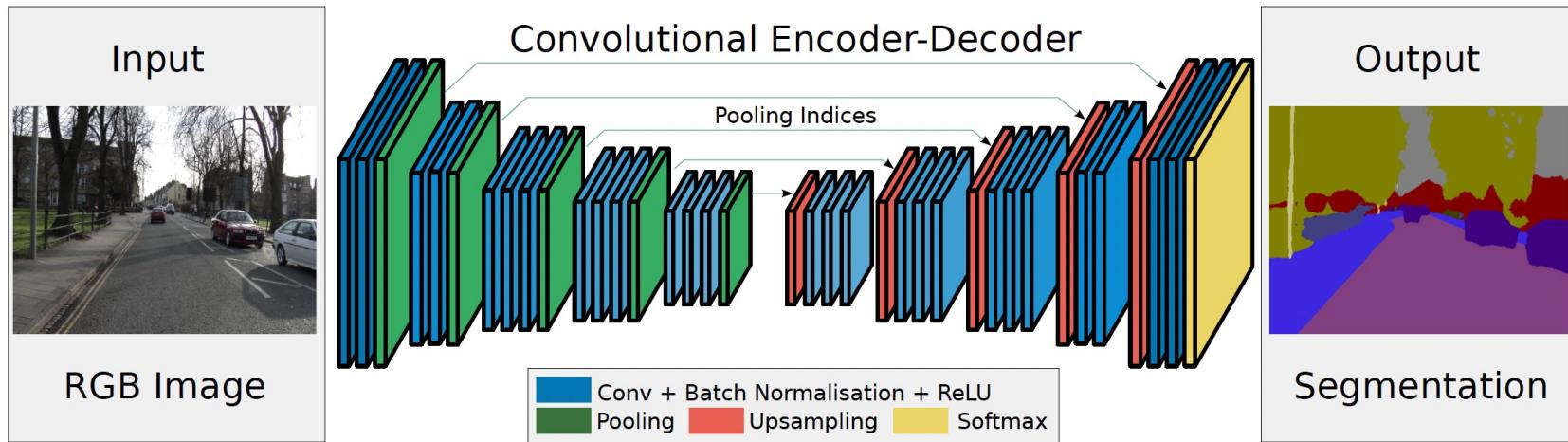


Fully Convolutional Network (FCN)

- Upsampling: combining “where” and “what”



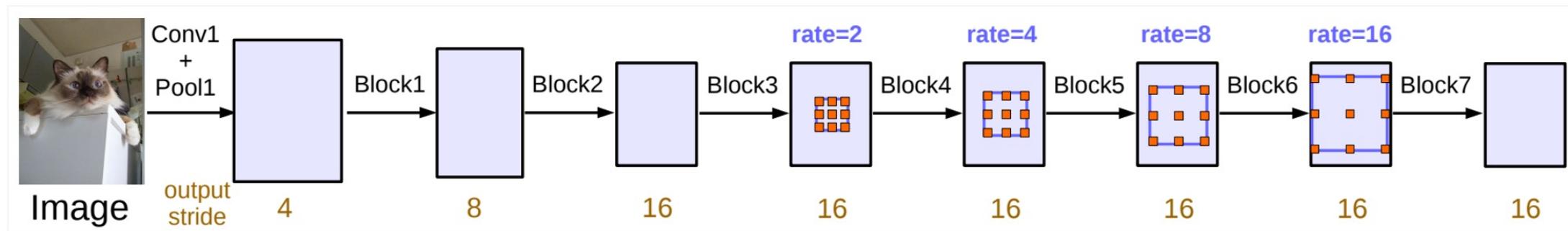
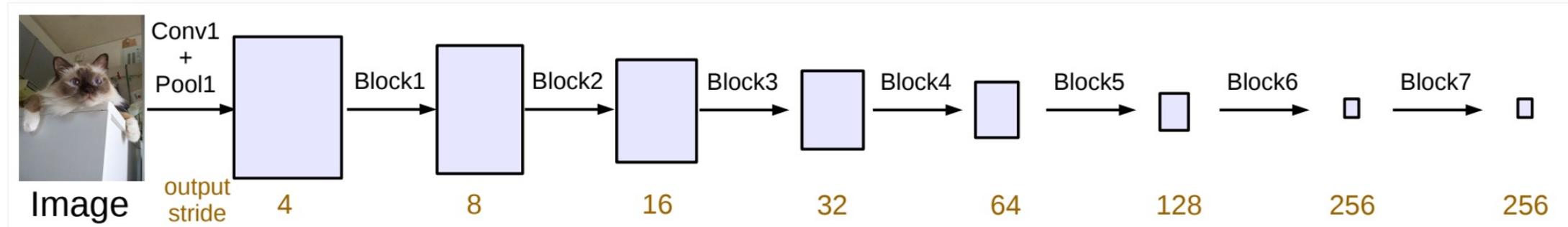
More upsampling strategies



Hyeonwoo Noh, Seunghoon Hong, Bohyung Han, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Vijay Badrinarayanan, Alex Kendall, Roberto Cipolla, "SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation", arXiv 2015, PAMI 2017

More upsampling strategies



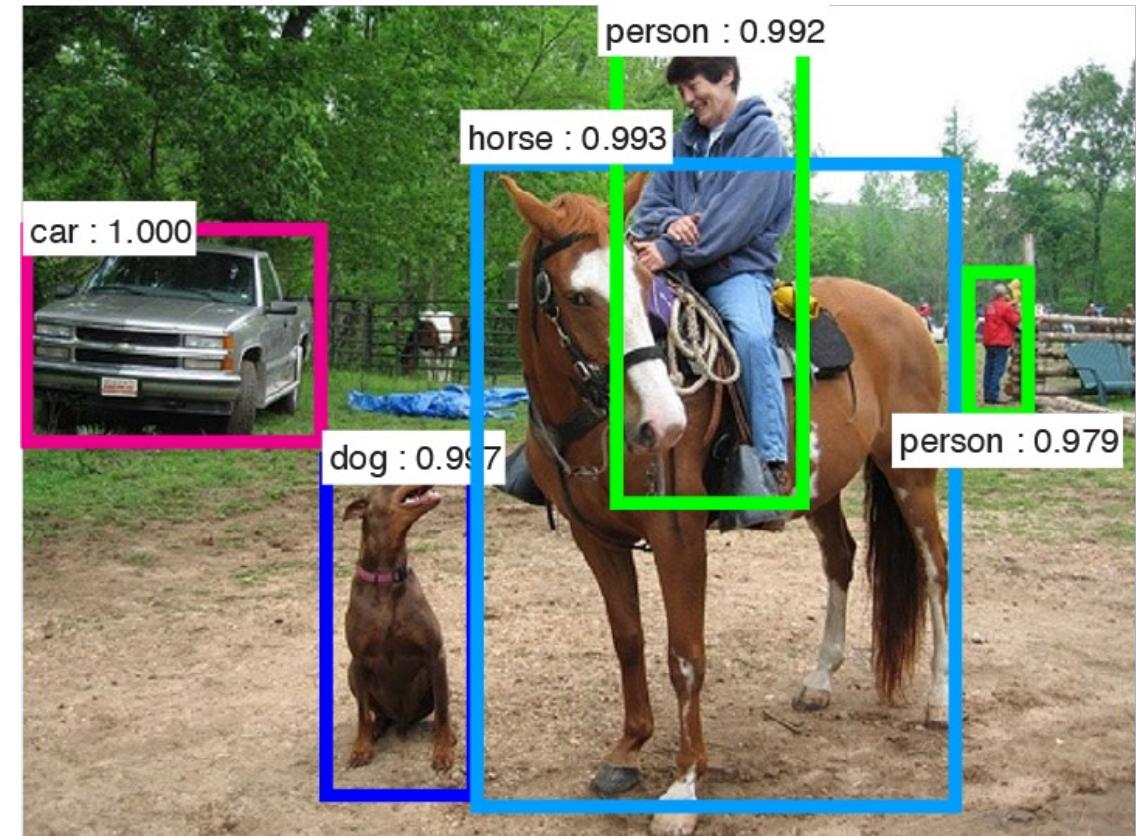
à trous / dilated conv

Don't stride on the feature map; stride on the kernel

Object Detection

Object Detection: Problem Definition

- **Task:** identify, locate, and recognize objects
- **Identify:**
 - Are there objects presented?
- **Locate:**
 - Where is each object?
 - box: (x, y, h, w)
- **Recognize:**
 - What is each object?
 - class label & confidence



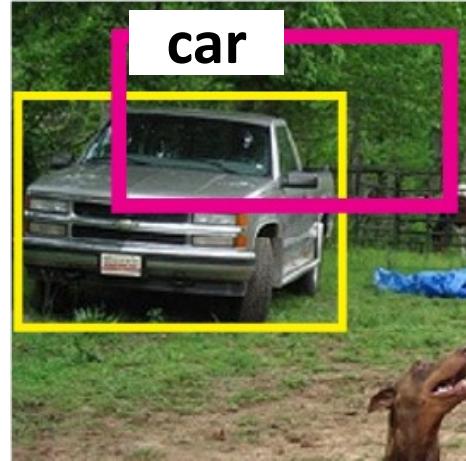
Object Detection: Evaluation



correct
(true positive)



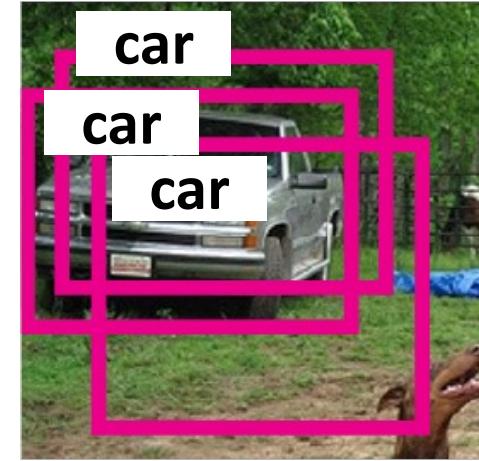
missed
(false negative)



mislocalized
(false positive)

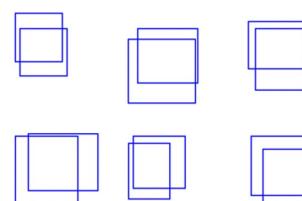


misrecognized
(false positive)

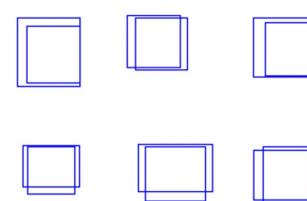


duplicated
(false positive)

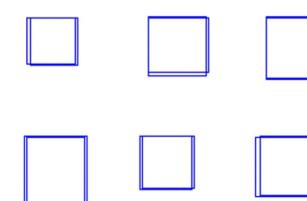
define by an IoU threshold



IoU = 0.5



IoU = 0.7



IoU = 0.9

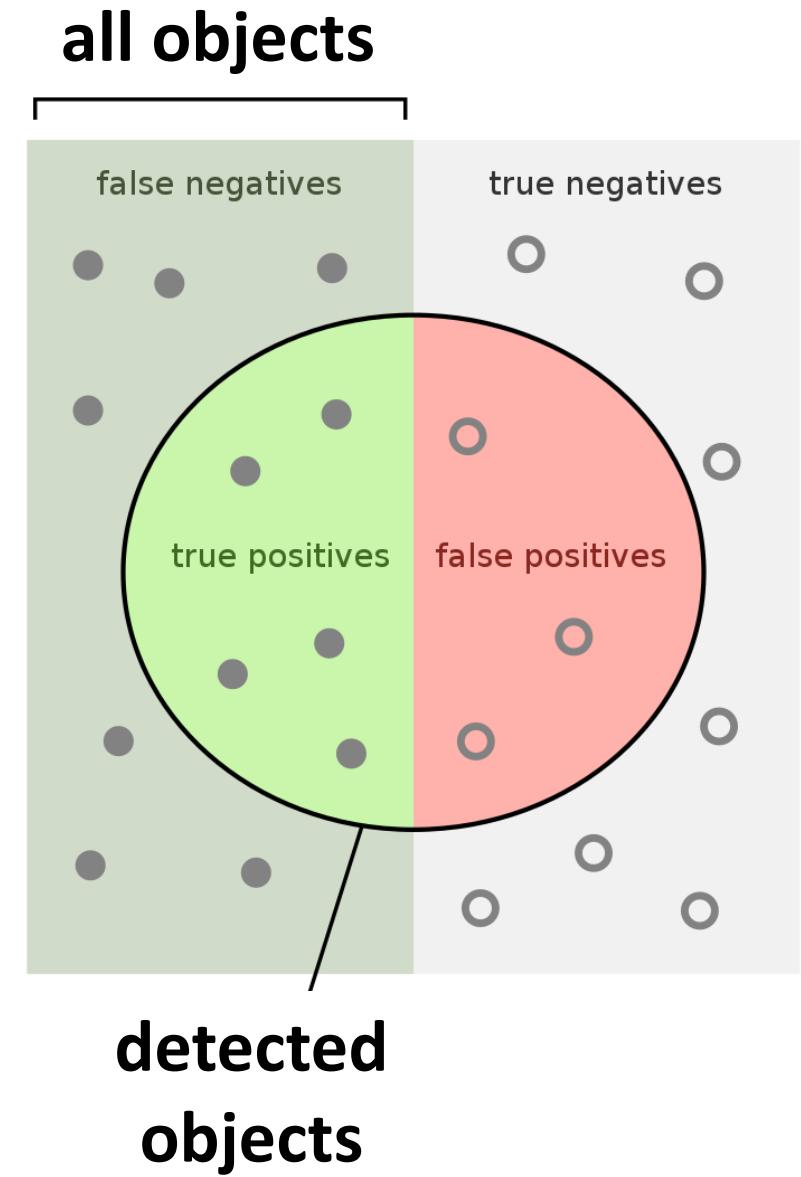
Object Detection: Evaluation

- **Precision:** What % of detected objects are correct?

$$\text{Precision} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$$

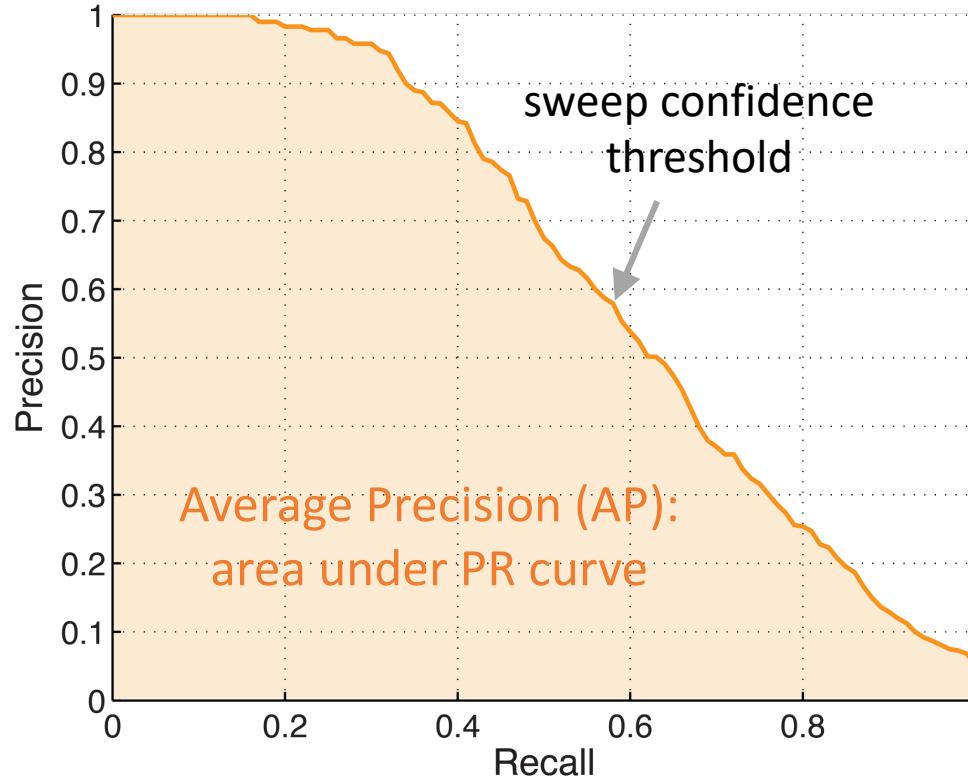
- **Recall:** What % of objects are detected?

$$\text{Recall} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$



Object Detection: Evaluation

- Precision-Recall Tradeoff



Life is a precision-recall tradeoff.

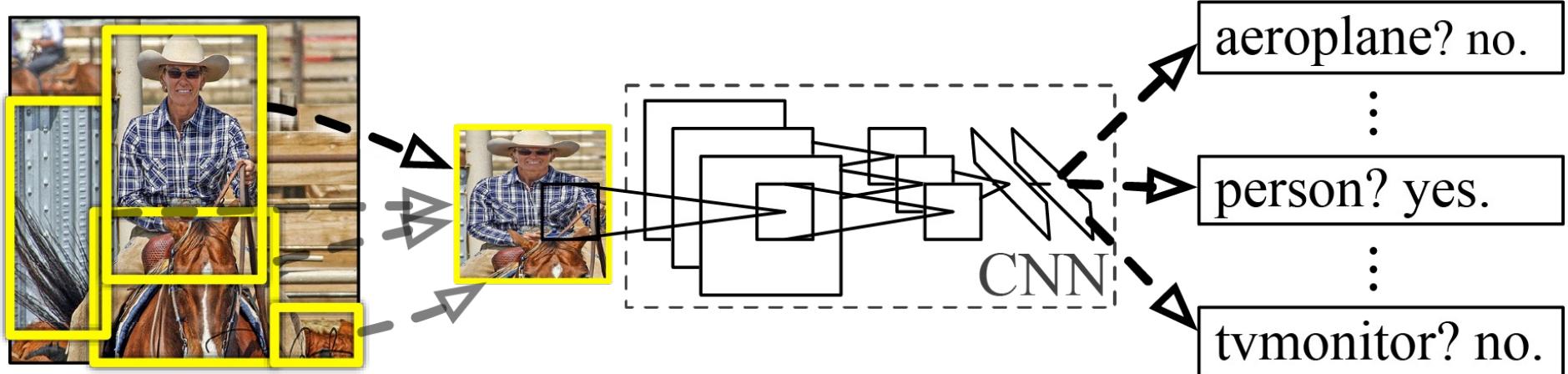
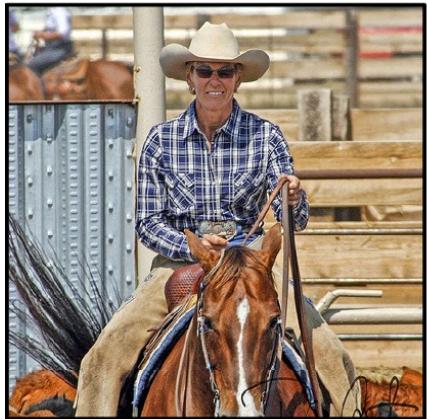


threshold = 0.4: Recall \uparrow , Precision \downarrow



threshold = 0.7: Recall \downarrow , Precision \uparrow

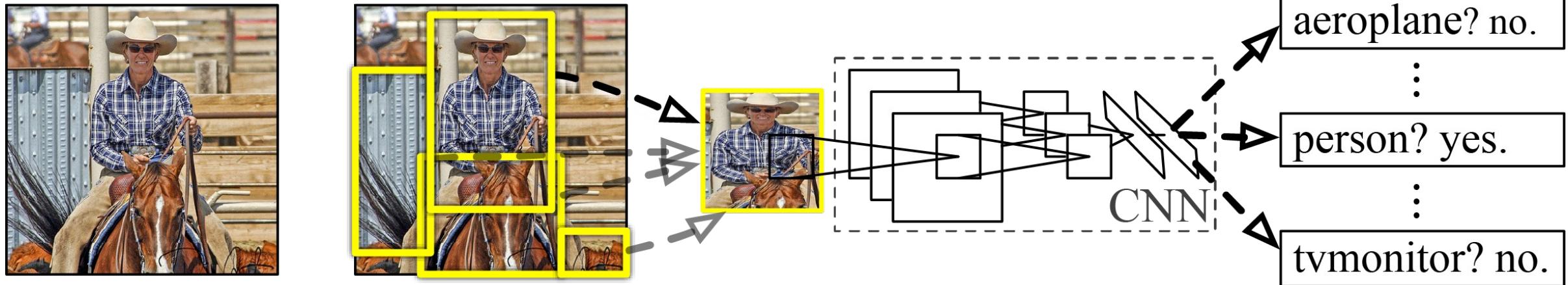
Region-Based CNN (R-CNN)



- The deep learning revolution in the vision community.

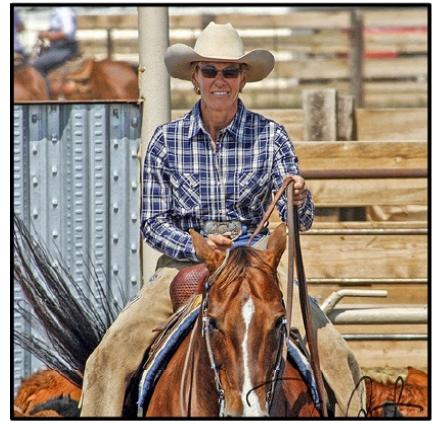
the ILSVRC workshop at ECCV 2012.¹ The central issue of debate can be distilled as: **To what extent do the CNN *classification* results on ImageNet generalize to object *detection* results on the PASCAL VOC Challenge?** In this paper, we study object detection using features computed by a large convolutional neural network, answering this important scientific question.²

Region-Based CNN (R-CNN)

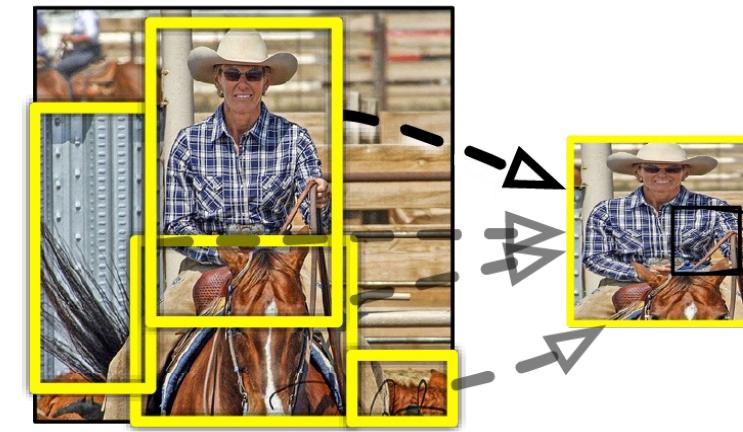


- **The deep learning revolution in the vision community.**
- **R-CNN:** Regions with CNN features
- **40%** relative improvement over classical methods
- Deep Learning made work beyond image classification

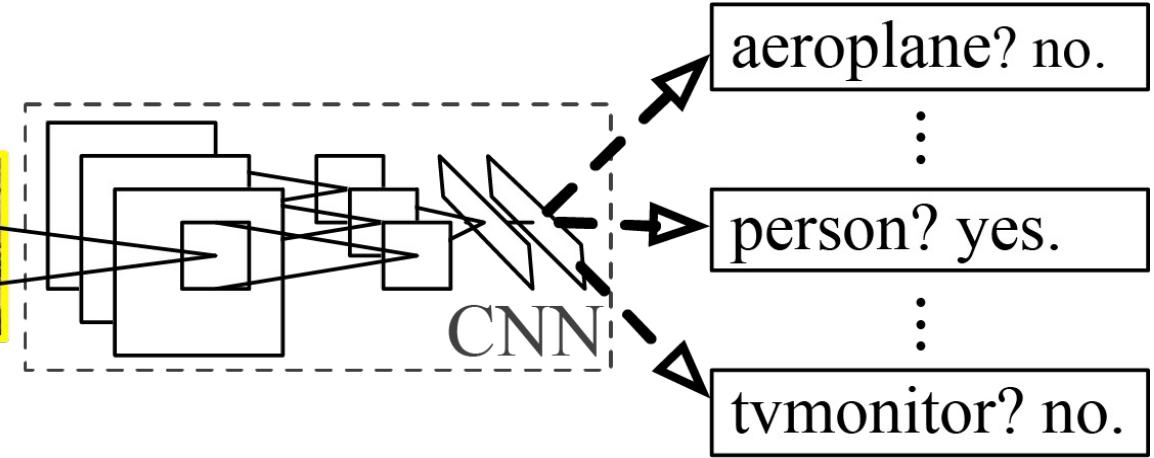
Region-Based CNN (R-CNN)



input



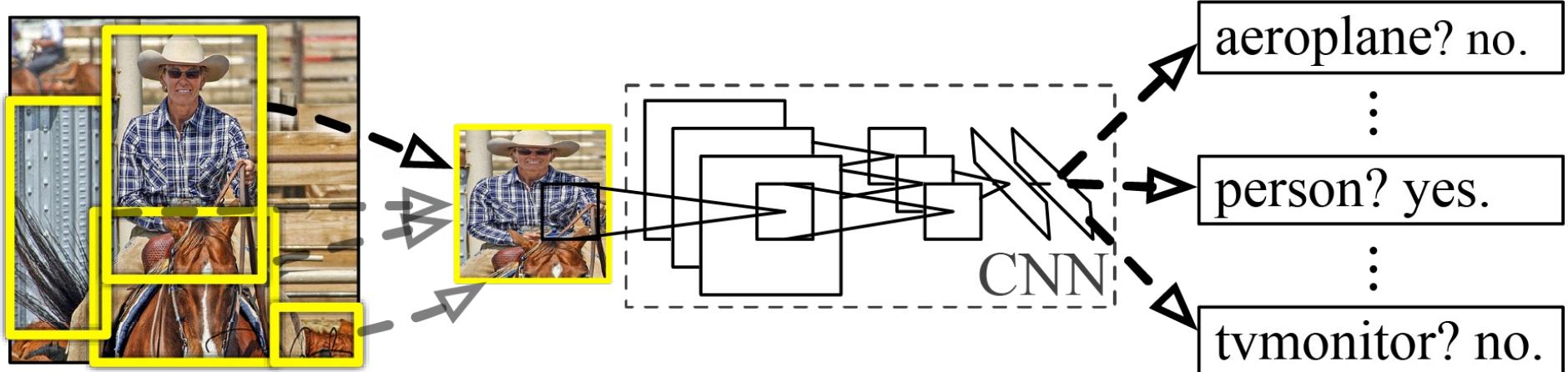
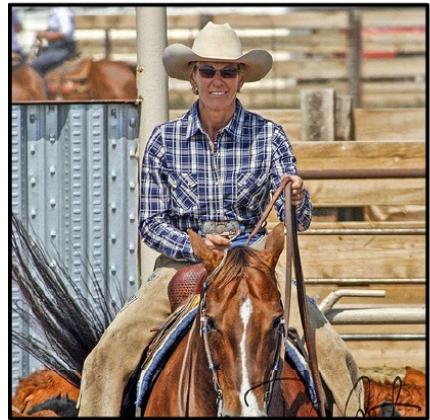
extract
“region proposals”



compute CNN features
for each region

classify regions
& refine boxes

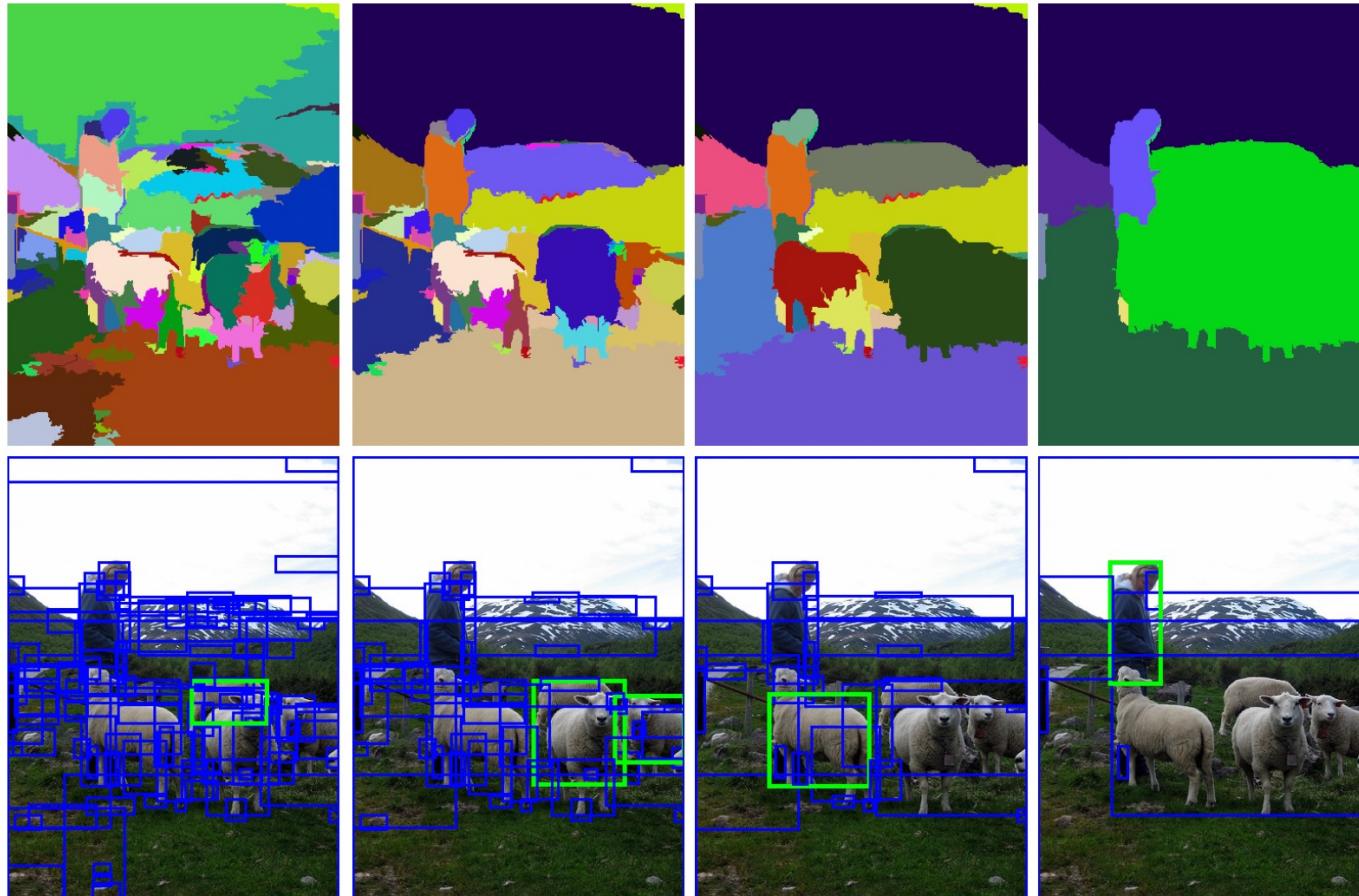
Region-Based CNN (R-CNN)



extract
“region proposals”

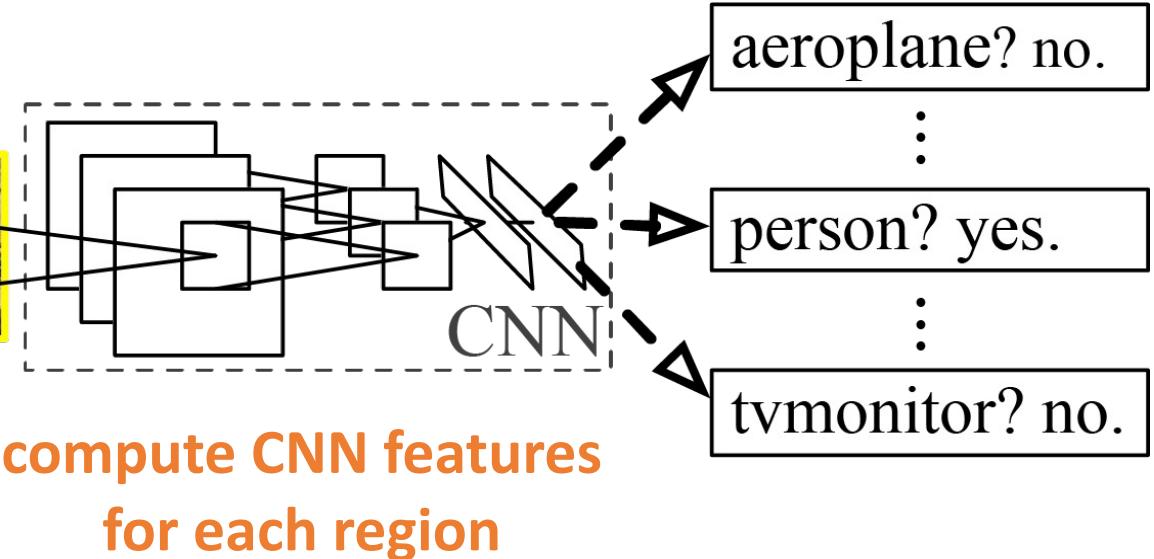
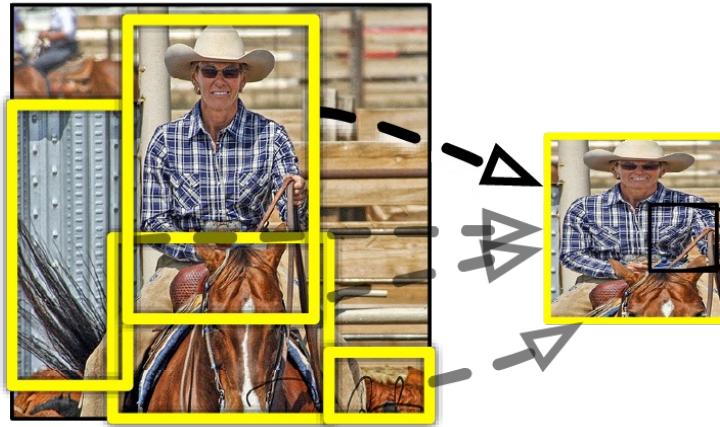
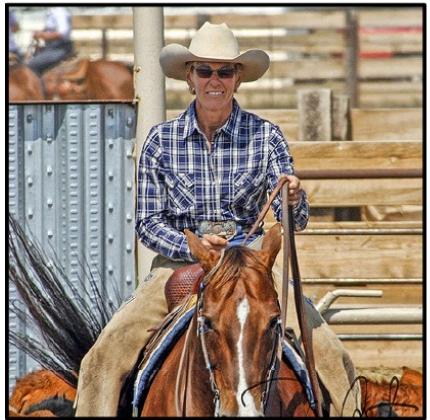
- **Reduce the search space (“selective search”)**
- e.g., ~ 2000 regions / image
- Recall-driven
 - a proposal may not be an object
 - an object is highly likely in the proposals

(e.g.) Selective Search for Region Proposal



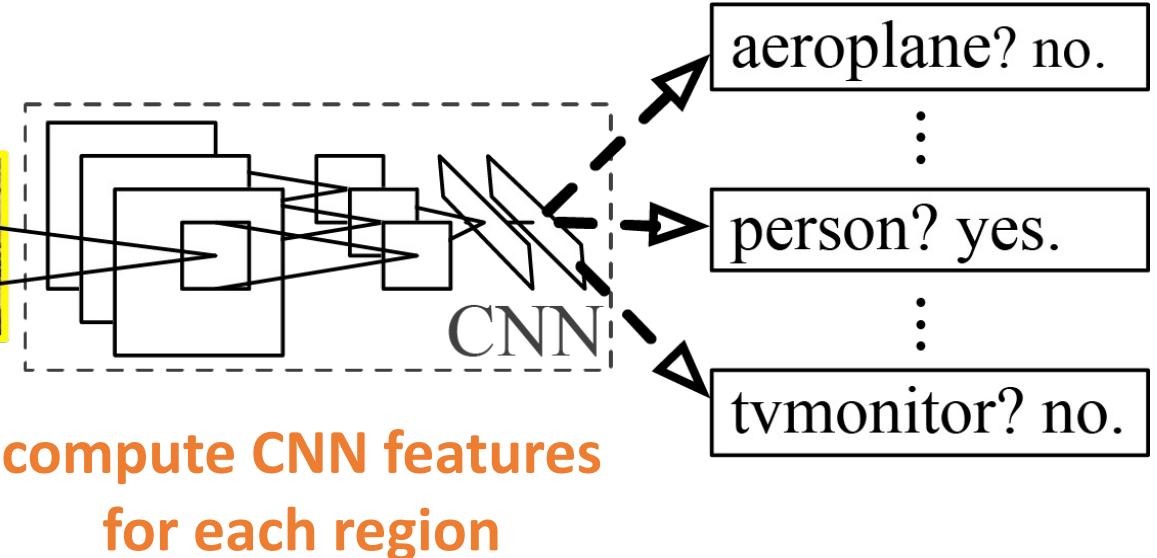
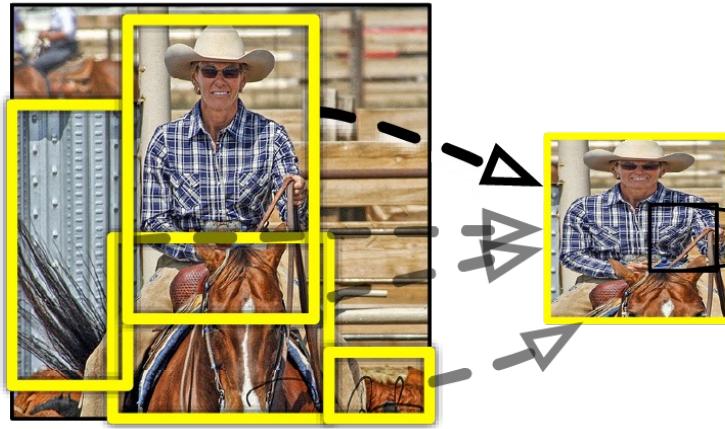
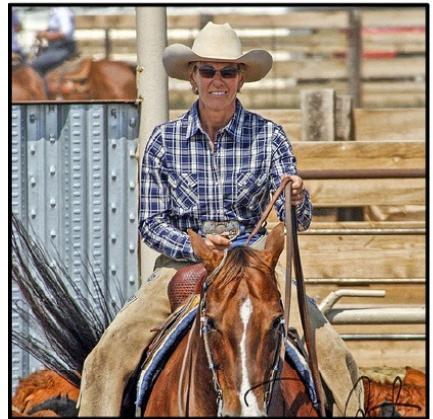
- Multi-scale over-segmentation
- Hierarchically group similar pixels
- “Hand-engineered” features
- Bounding boxes of segments

Region-Based CNN (R-CNN)

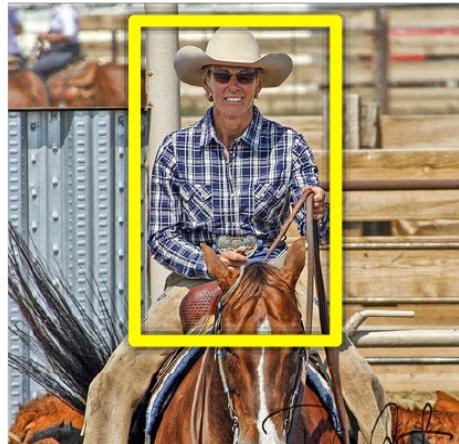


one proposal

Region-Based CNN (R-CNN)



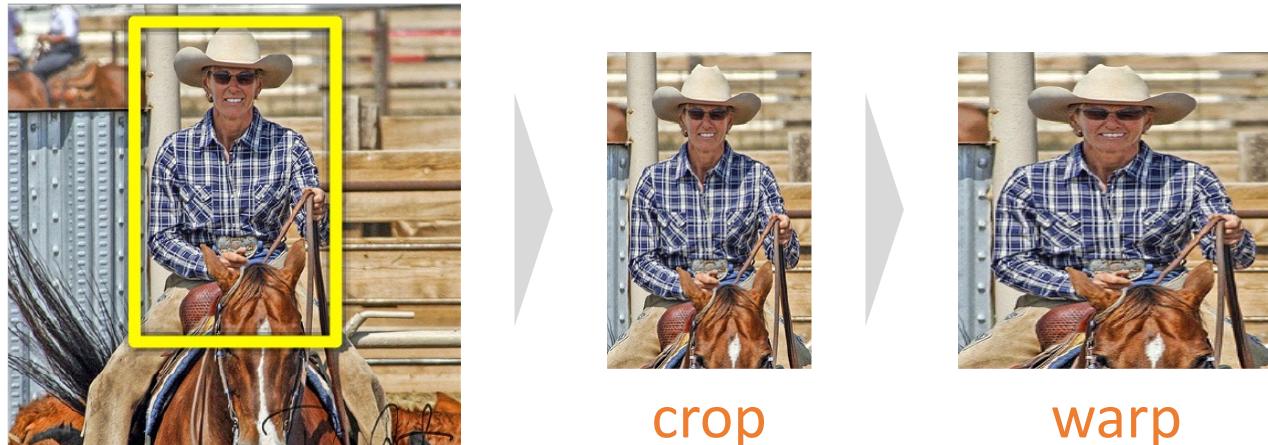
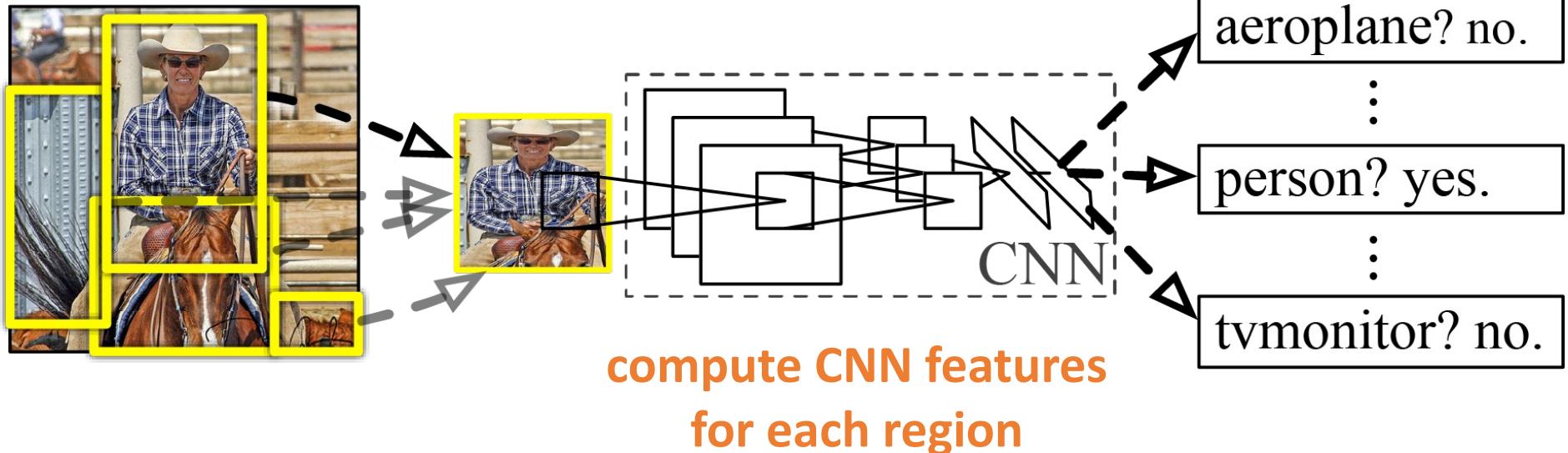
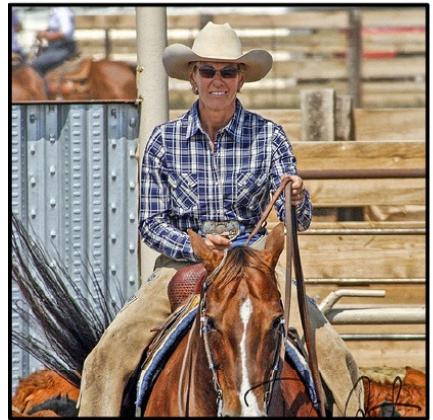
compute CNN features
for each region



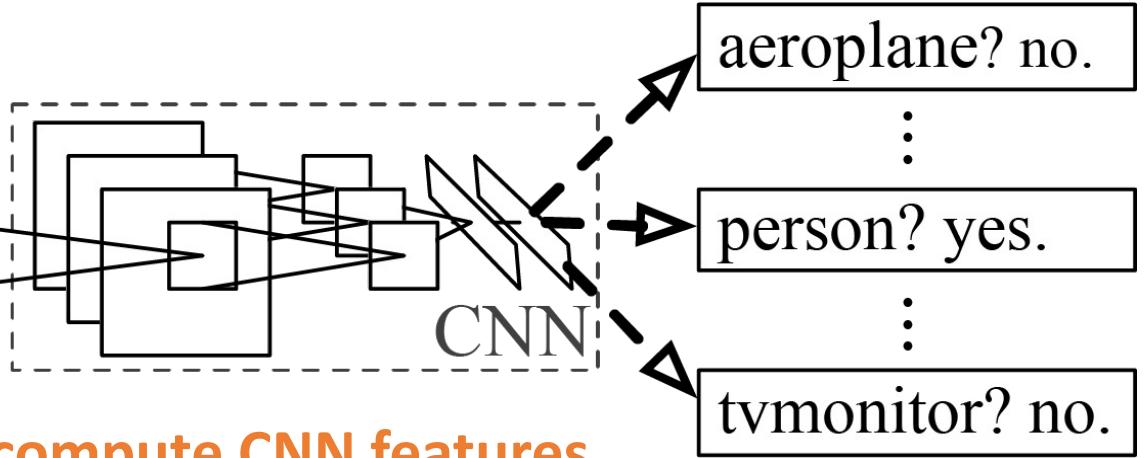
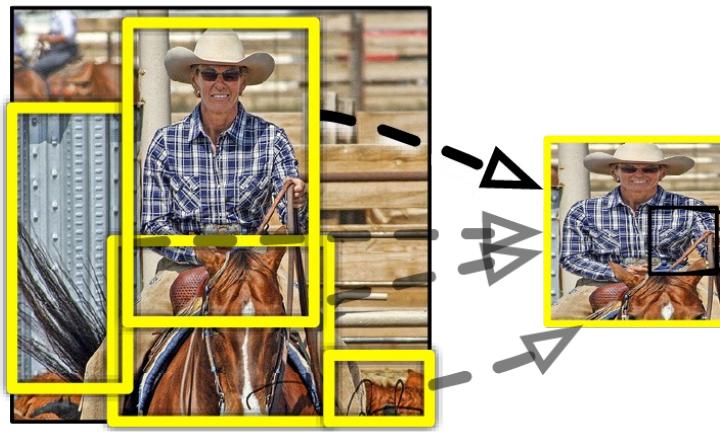
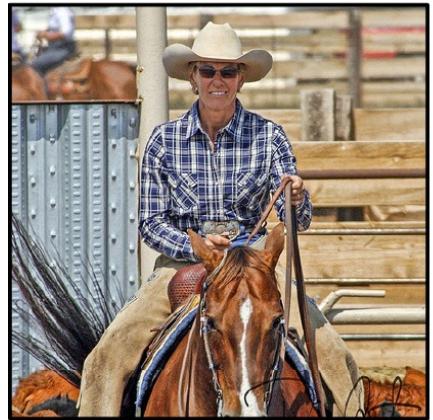
crop

one proposal

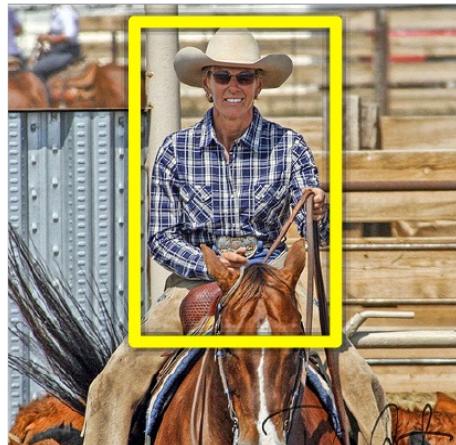
Region-Based CNN (R-CNN)



Region-Based CNN (R-CNN)



compute CNN features
for each region

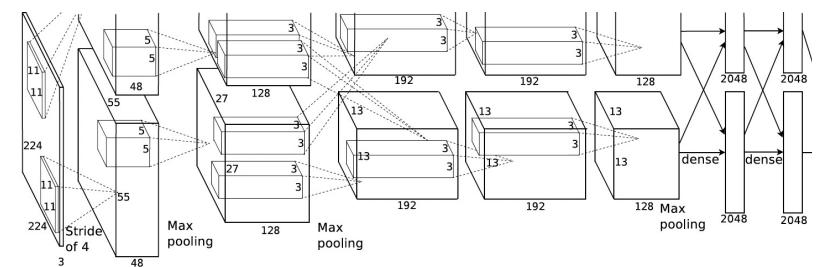


one proposal



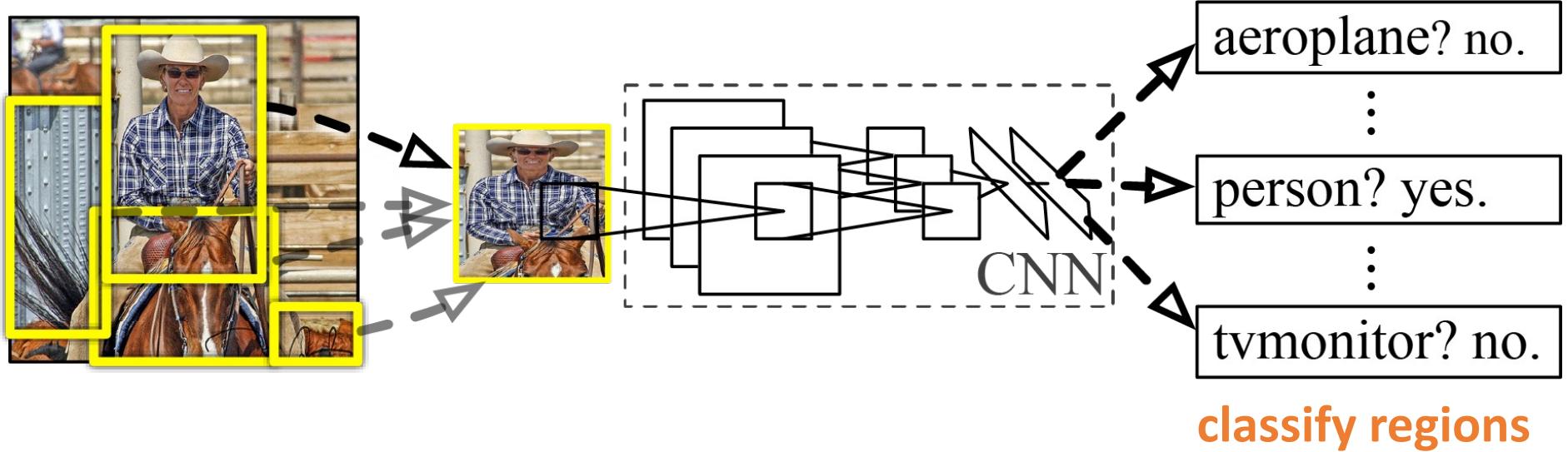
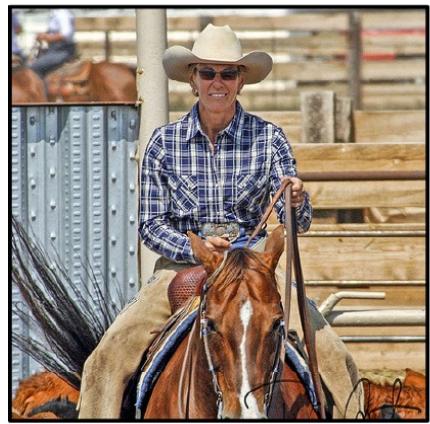
crop

warp
(e.g., 224x224)

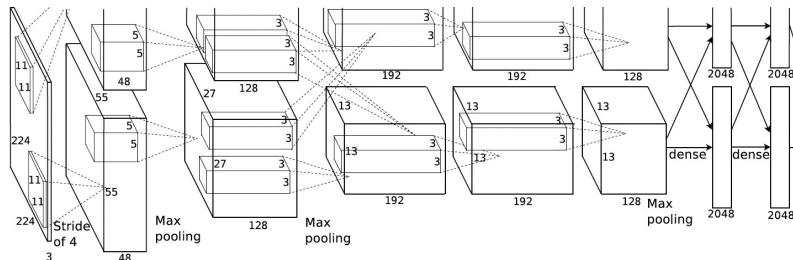
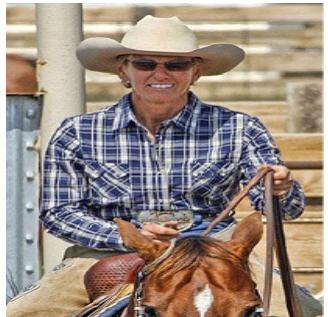
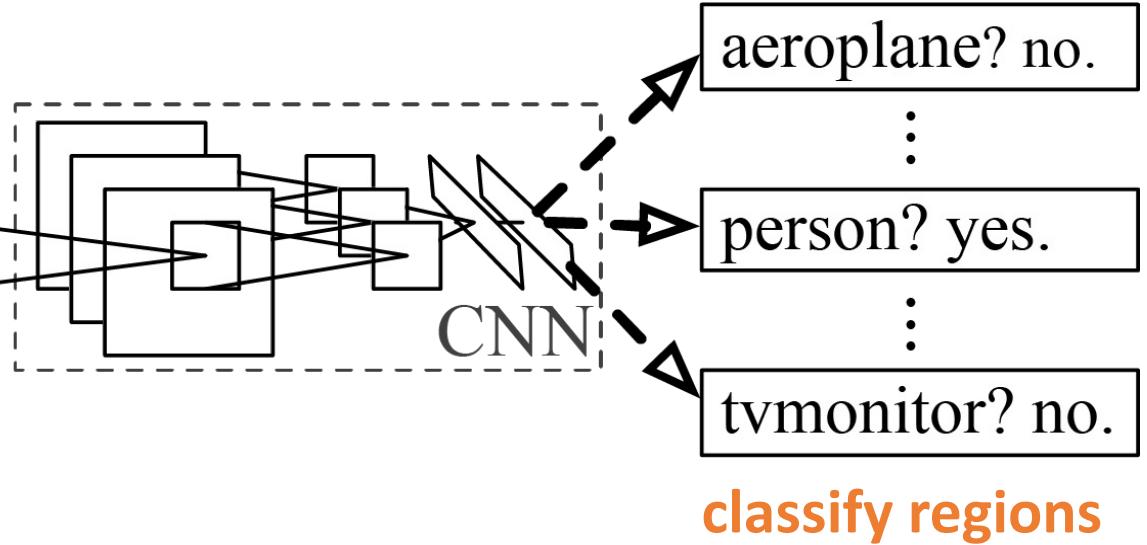
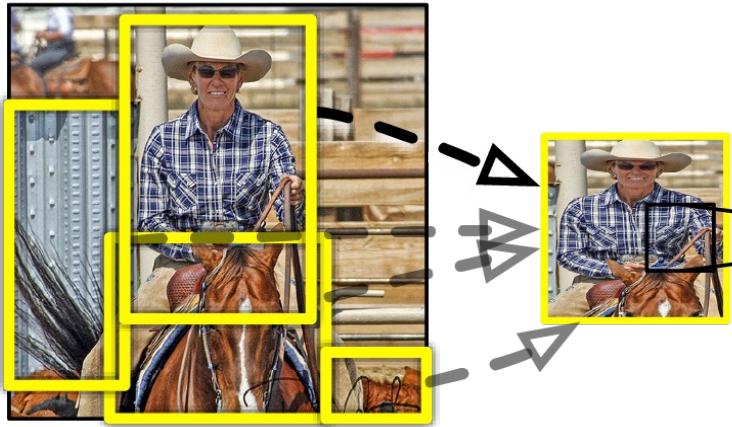
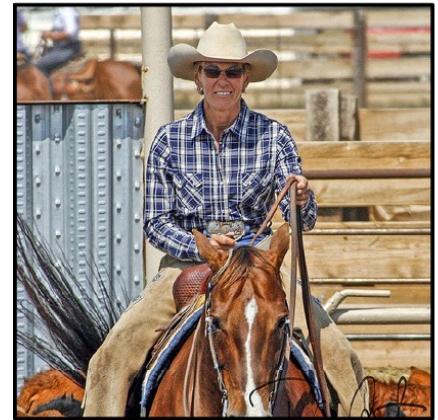


apply a classification ConvNet
(e.g., AlexNet)

Region-Based CNN (R-CNN)



Region-Based CNN (R-CNN)

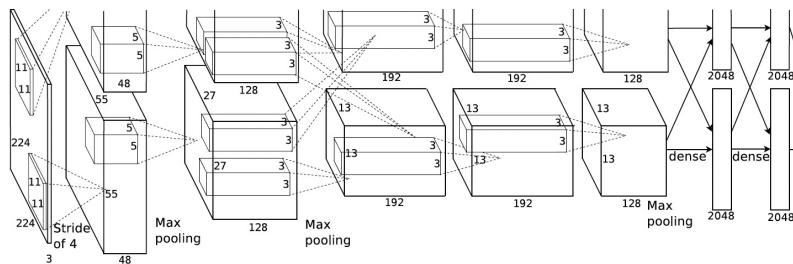
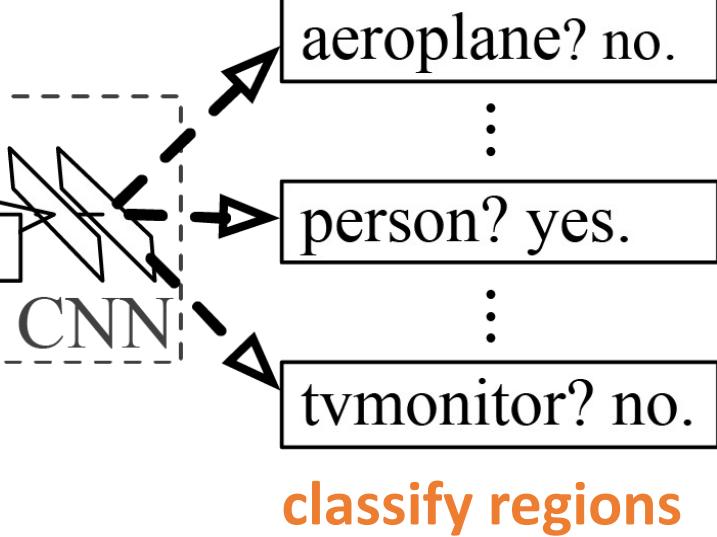
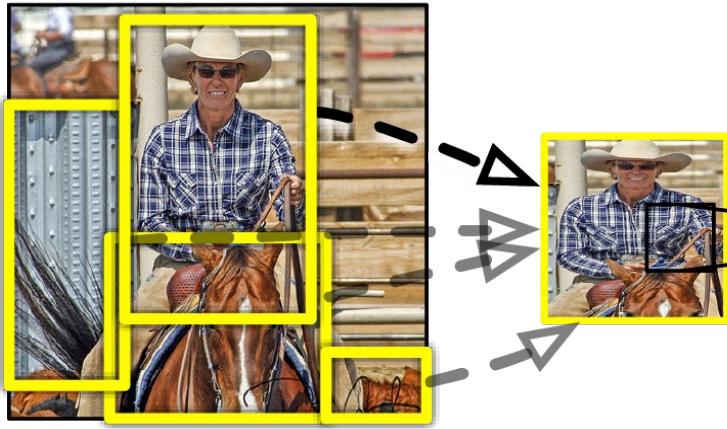
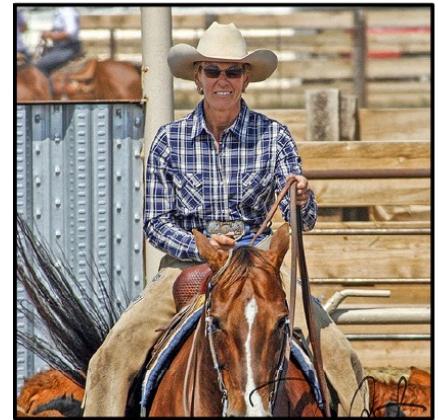


background: 0.0

person: 0.9

horse: 0.1

Region-Based CNN (R-CNN)



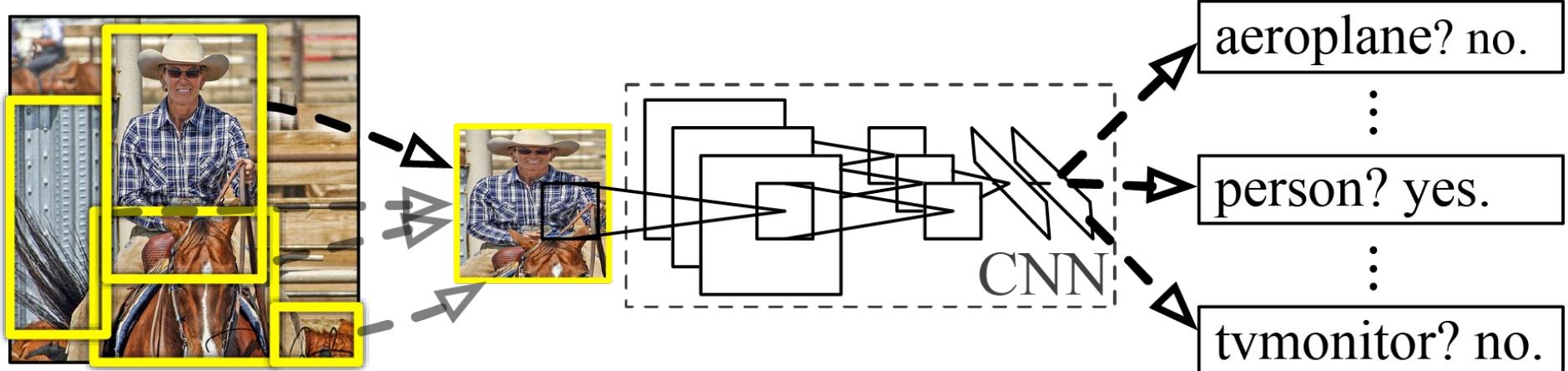
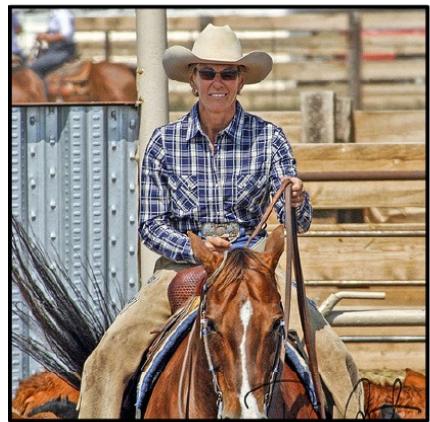
Detection is first of all
a foreground/background
classification problem.

background: 0.9

person: 0.0

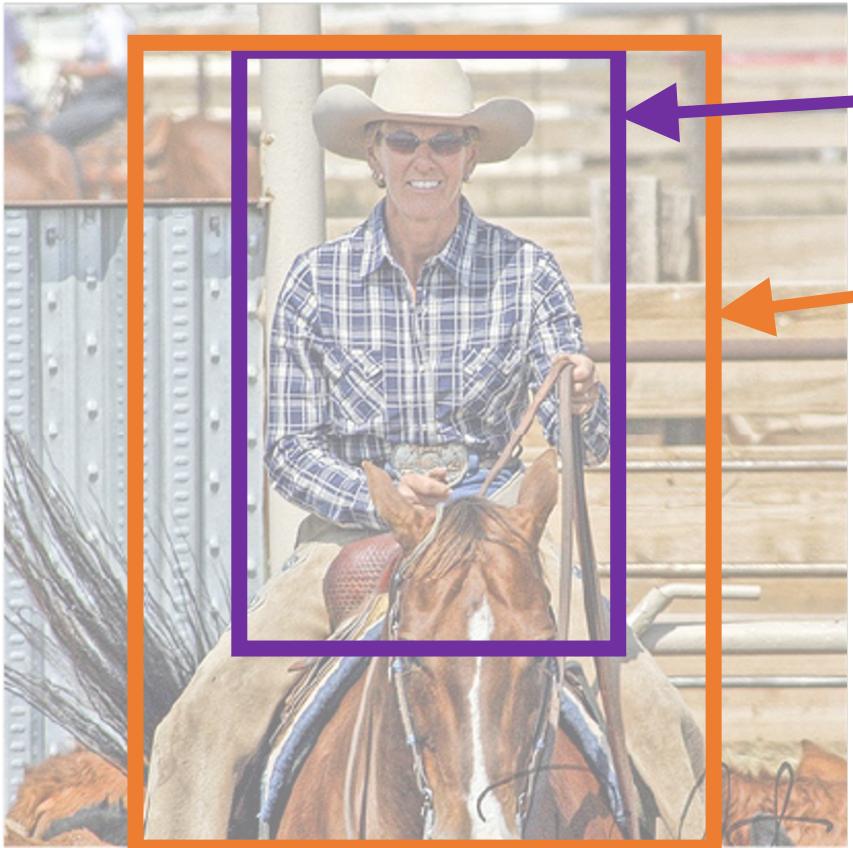
horse: 0.1

Region-Based CNN (R-CNN)



Bounding-box Regression

- Geometric prediction from deep neural networks



reference box: (x, y, w, h)

target box: (x', y', w', h')

$$\Delta x = (x' - x)/w$$

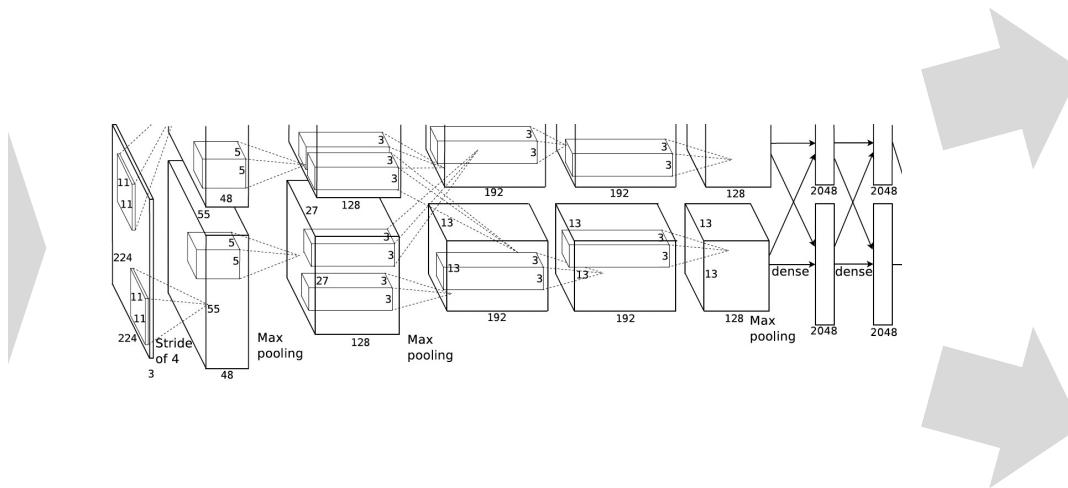
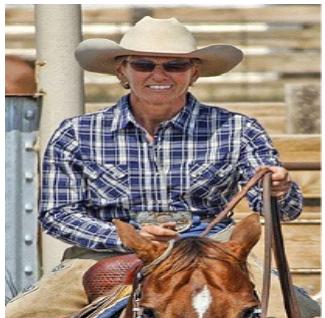
$$\Delta y = (y' - y)/h$$

$$\Delta \log w = \log w' - \log w$$

$$\Delta \log h = \log h' - \log h$$

Bounding-box Regression

- Geometric prediction from deep neural networks
- one network, multiple tasks



classification

bbox regression
 $\Delta x, \Delta y, \Delta \log w, \Delta \log h$

R-CNN



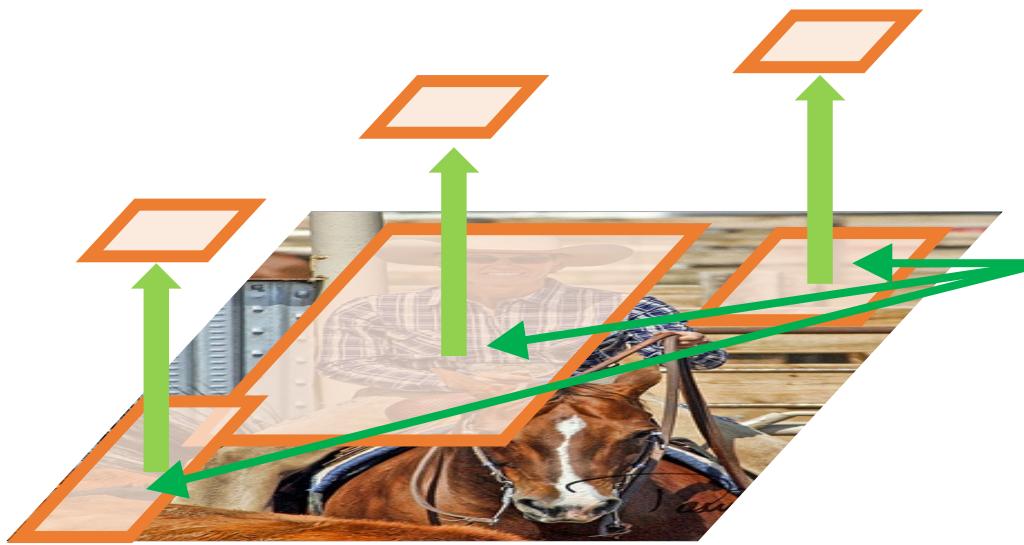
Slides adapted from: Ross Girshick, ICCV 2015
Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation", CVPR 2014

R-CNN



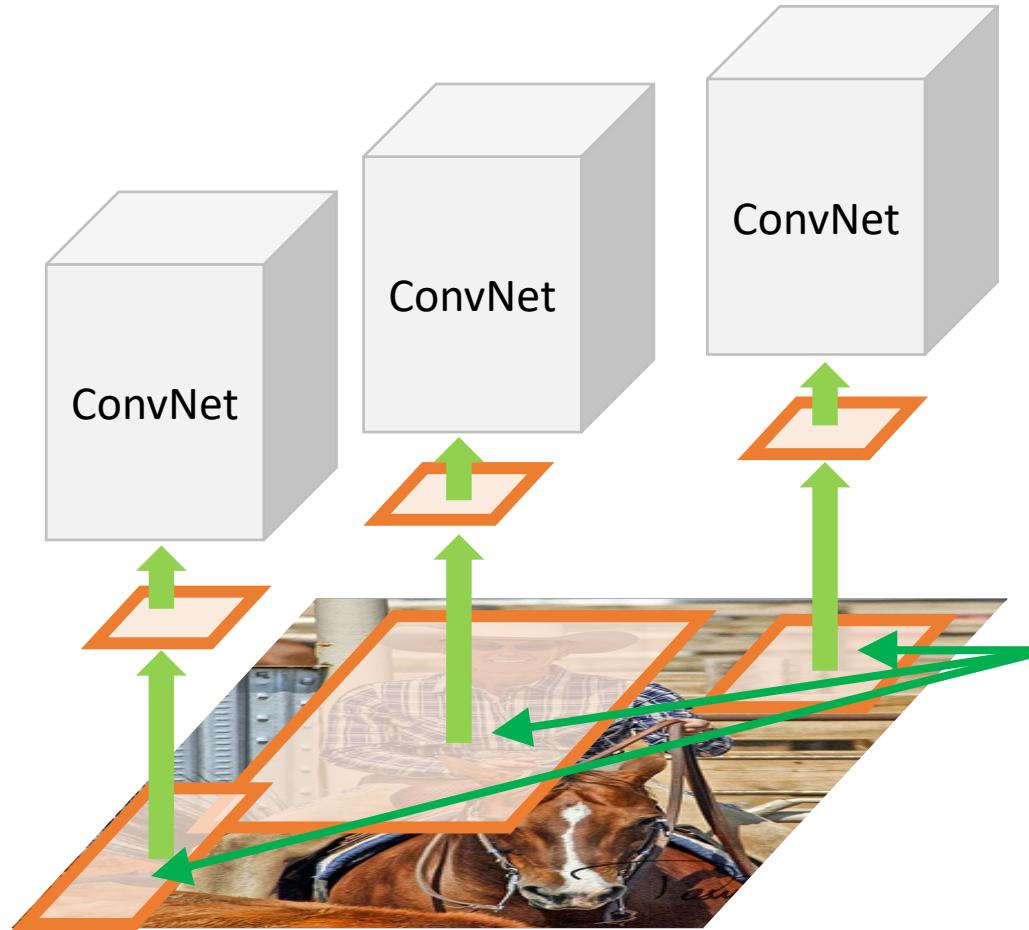
- Propose ~2000 Regions of Interest (RoI)

R-CNN



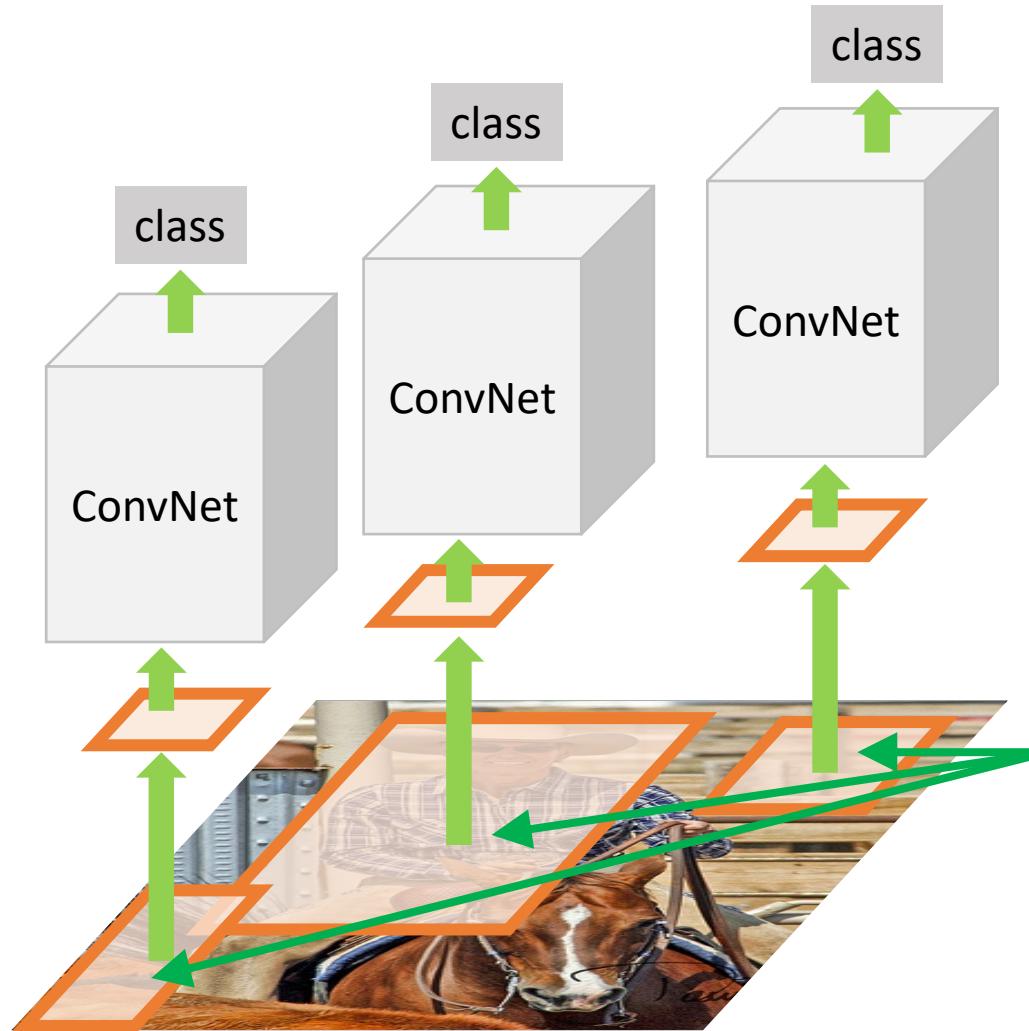
- Warp regions
- Propose ~2000 Regions of Interest (RoI)

R-CNN



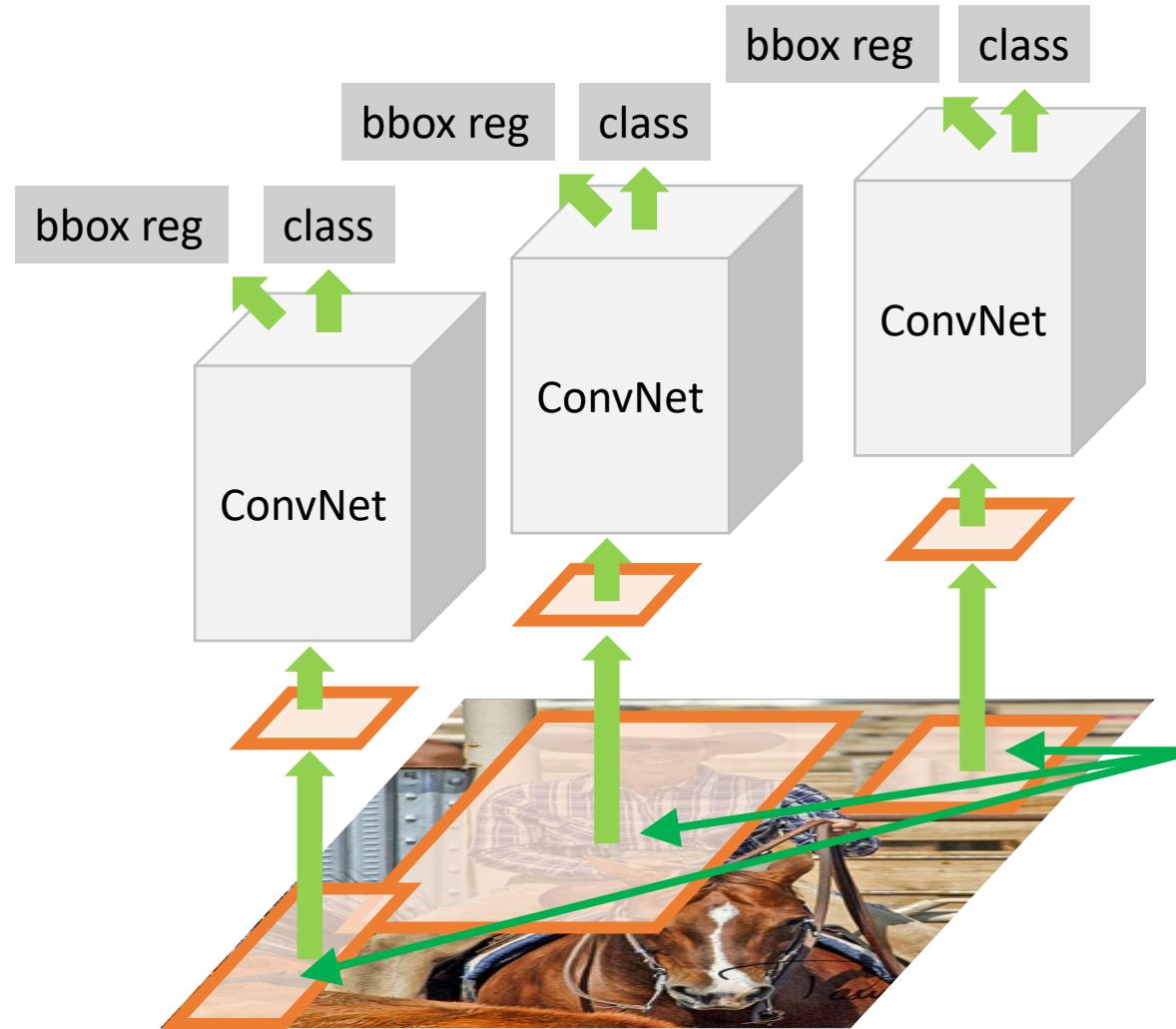
- Apply ConvNets to regions
- Warp regions
- Propose ~2000 Regions of Interest (RoI)

R-CNN



- Classify regions
- Apply ConvNets to regions
- Warp regions
- Propose ~2000 Regions of Interest (RoI)

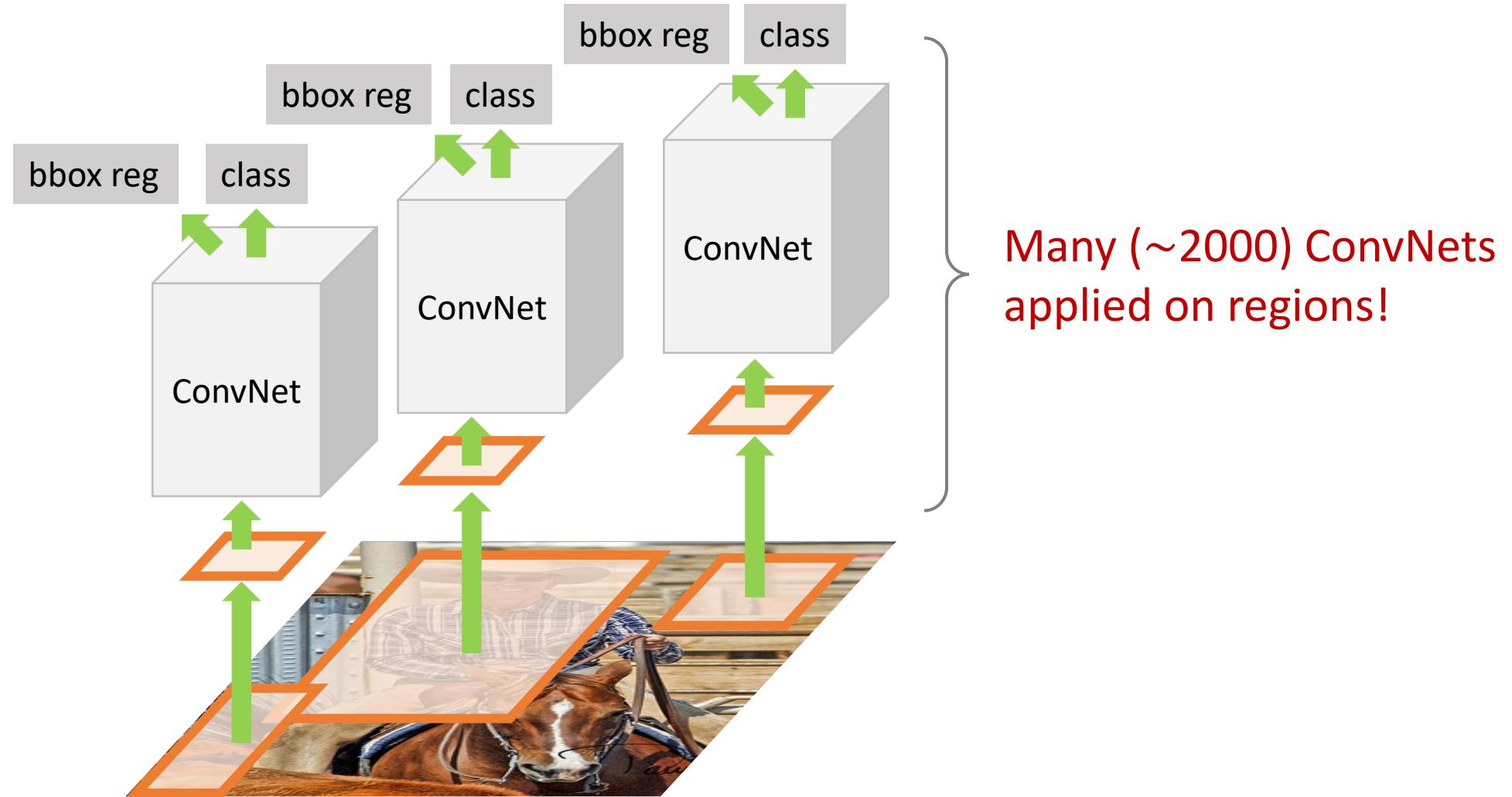
R-CNN



- Regress bounding-boxes
- Classify regions
- Apply ConvNets to regions
- Warp regions
- Propose ~2000 Regions of Interest (RoI)

What's wrong with “Slow” R-CNN?

(quote [Ross Girshick, ICCV 2015])



Slides adapted from: Ross Girshick, ICCV 2015

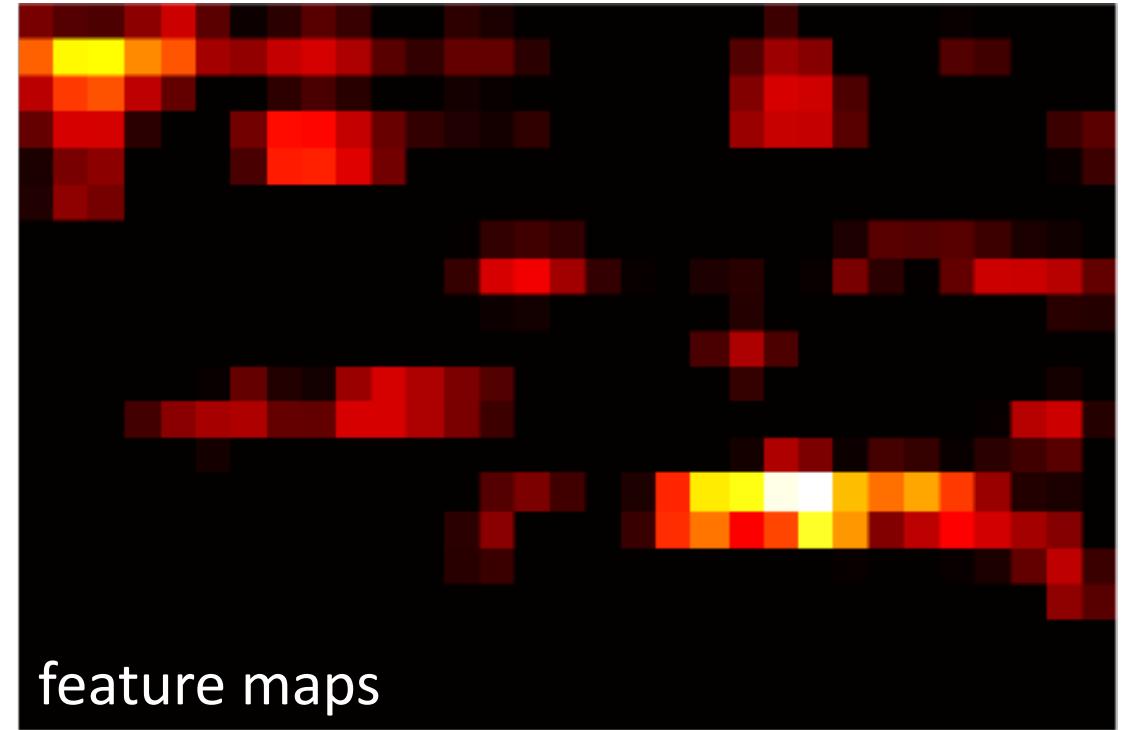
Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik, “Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation”, CVPR 2014

Towards Fast R-CNN (SPPnet)

- ConvNets are Fully Convolutional!



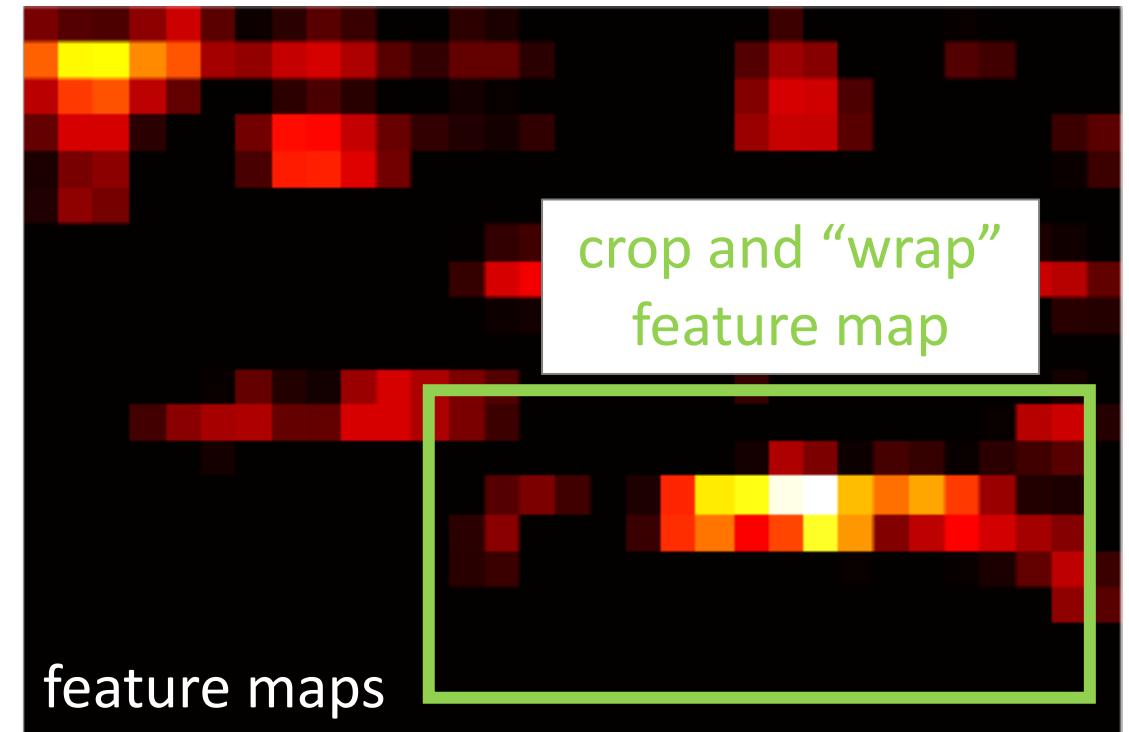
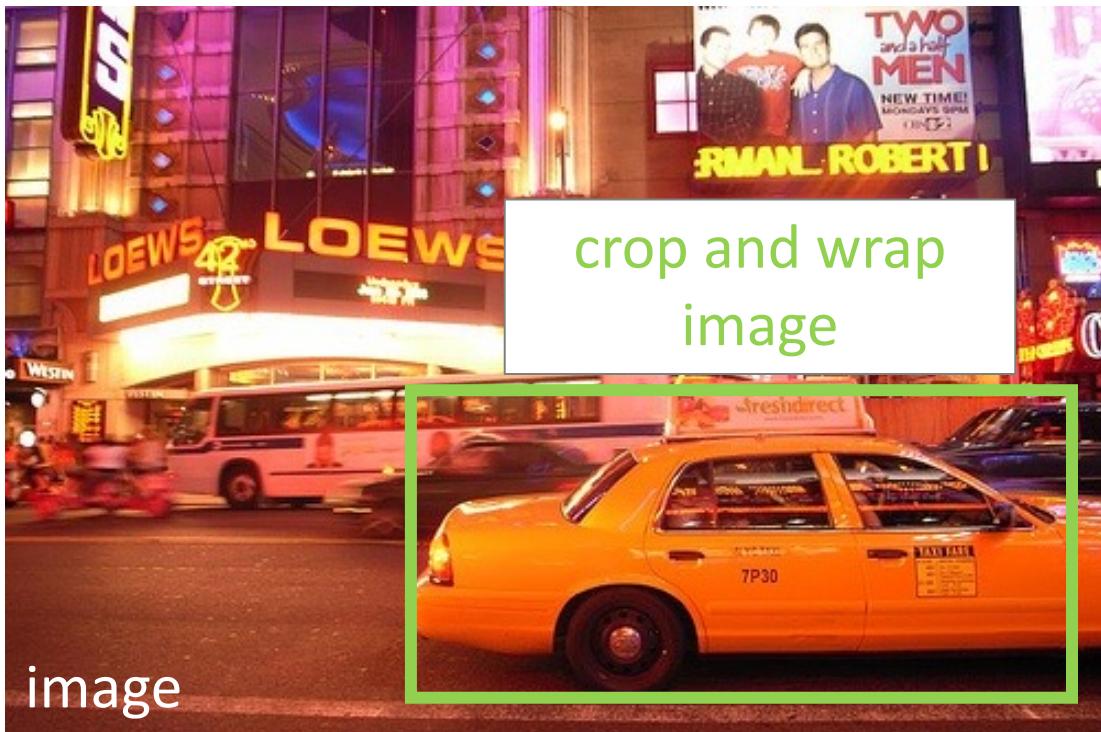
image



feature maps

Towards Fast R-CNN (SPPnet)

- ConvNets are Fully Convolutional!



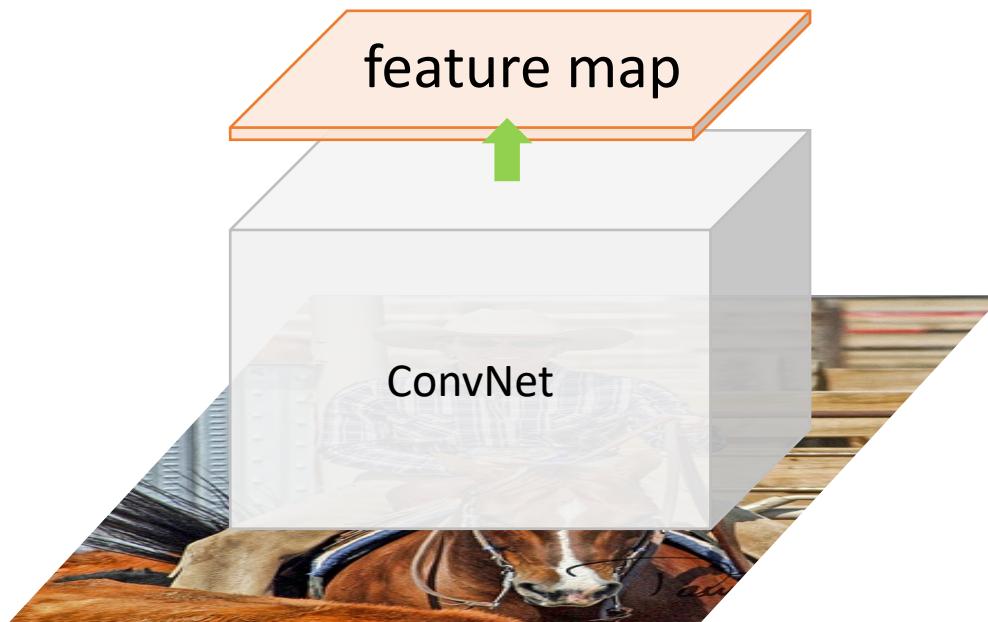
Towards Fast R-CNN (SPPnet)



Slides adapted from: Ross Girshick, ICCV 2015

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, "Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition", ECCV 2014

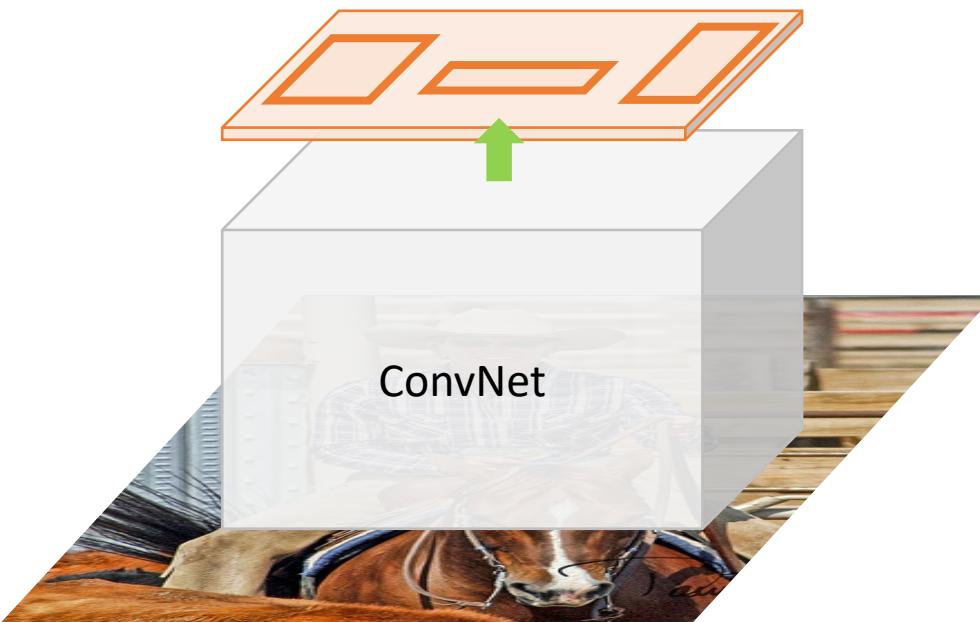
Towards Fast R-CNN (SPPnet)



- Apply ConvNet to **whole image**

Slides adapted from: Ross Girshick, ICCV 2015

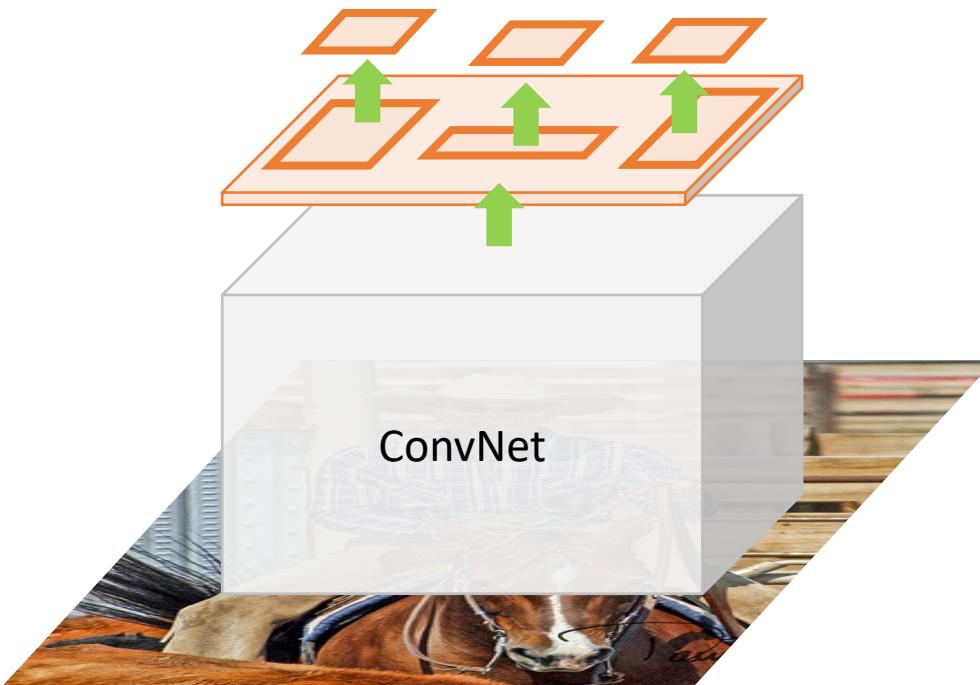
Towards Fast R-CNN (SPPnet)



- Map proposed Rols to feature maps
- Apply ConvNet to **whole image**

Slides adapted from: Ross Girshick, ICCV 2015

Towards Fast R-CNN (SPPnet)



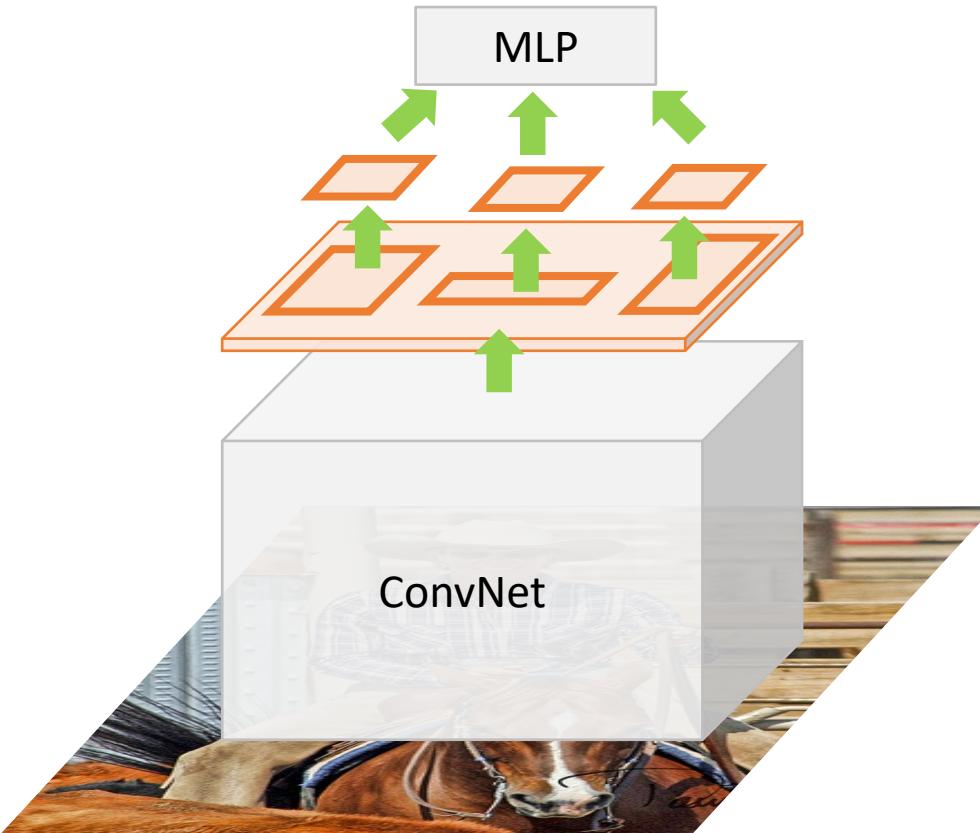
- Crop and **pool** features
- Map proposed Rols to feature maps
- Apply ConvNet to **whole** image

*In original SPPnet, pooling is done with a fancier variant (spatial pyramid pooling)

Slides adapted from: Ross Girshick, ICCV 2015

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, "Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition", ECCV 2014

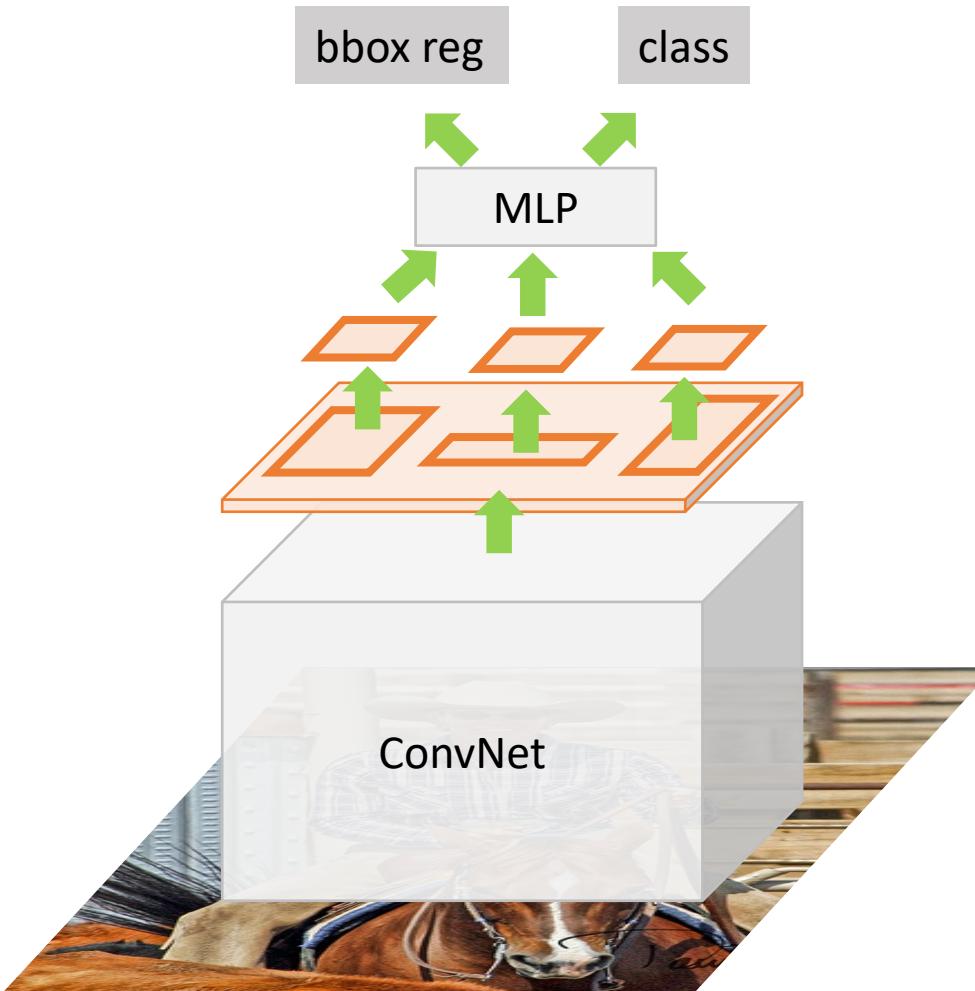
Towards Fast R-CNN (SPPnet)



- Apply MLP on regions
- Crop and **pool** features
- Map proposed Rols to feature maps
- Apply ConvNet to **whole** image

Slides adapted from: Ross Girshick, ICCV 2015

Towards Fast R-CNN (SPPnet)

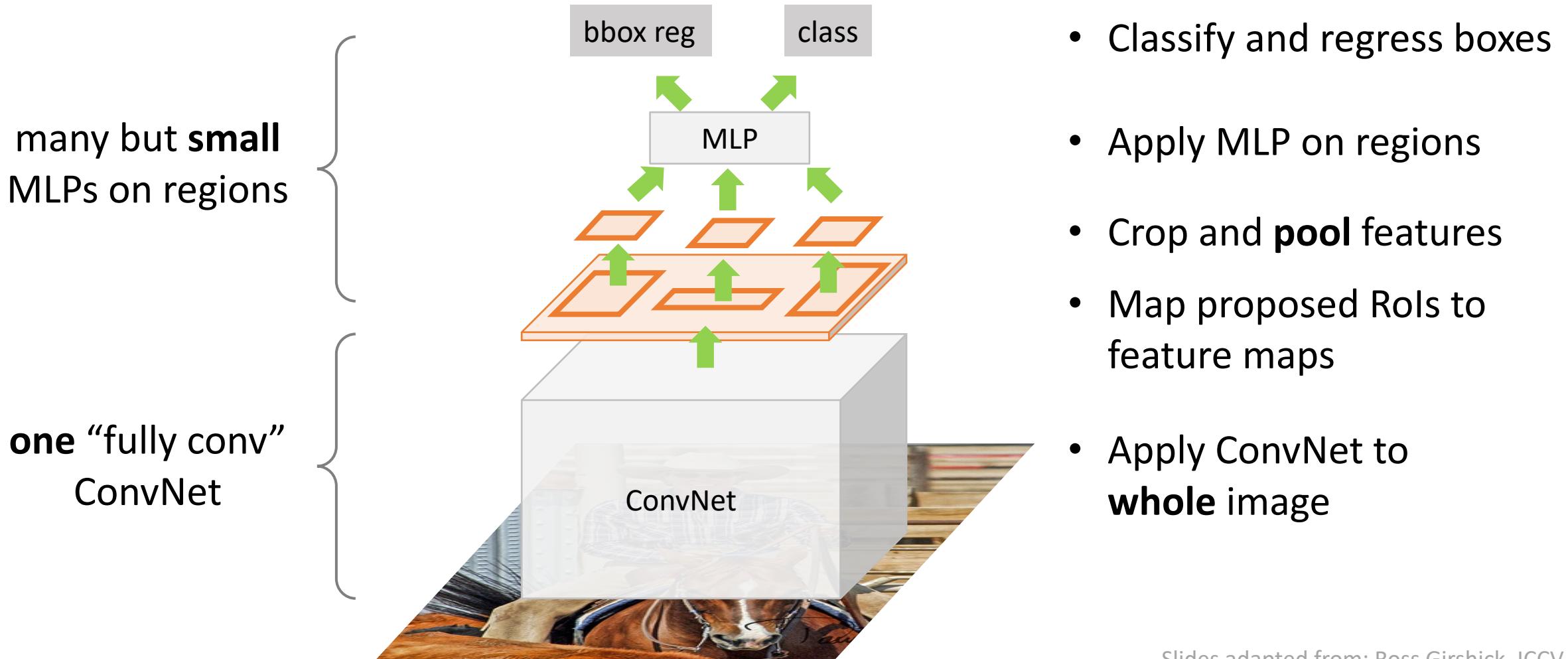


- Classify and regress boxes
- Apply MLP on regions
- Crop and **pool** features
- Map proposed Rols to feature maps
- Apply ConvNet to **whole** image

Slides adapted from: Ross Girshick, ICCV 2015

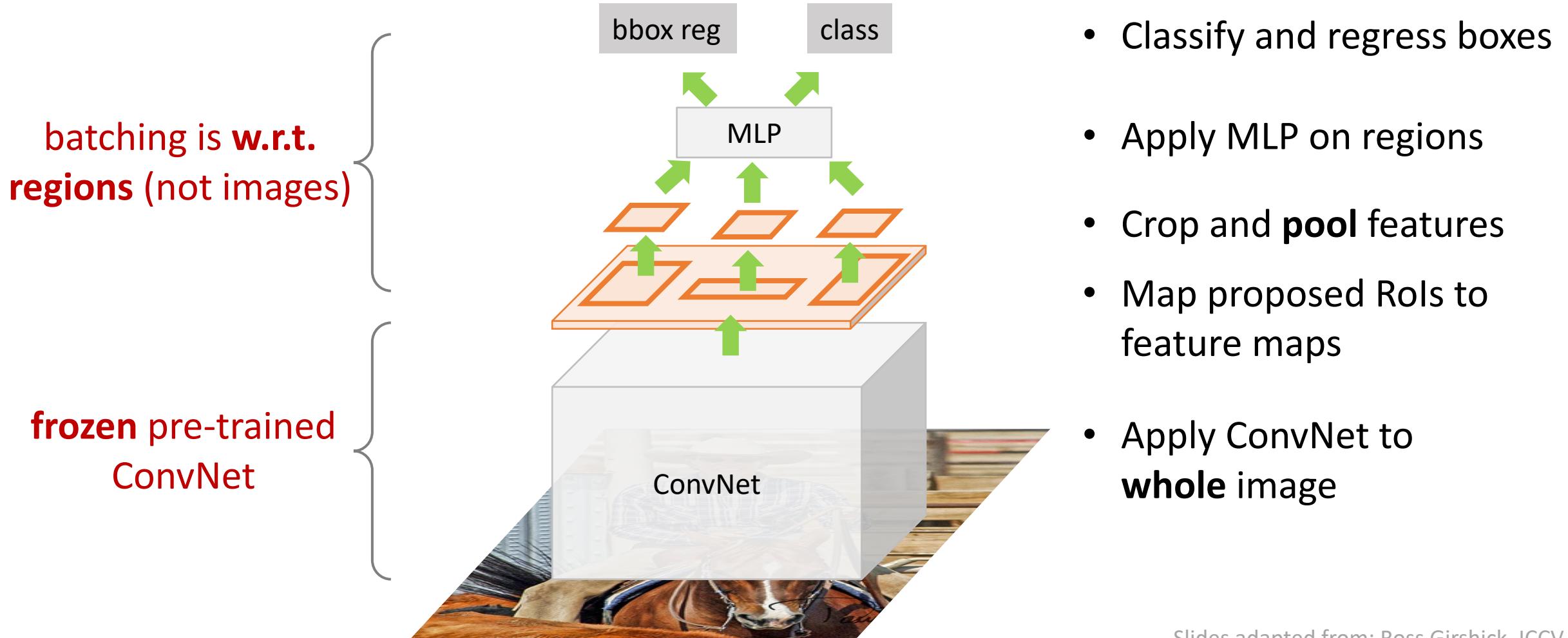
Towards Fast R-CNN (SPPnet)

- What's good with “fully convolutional” ConvNet?



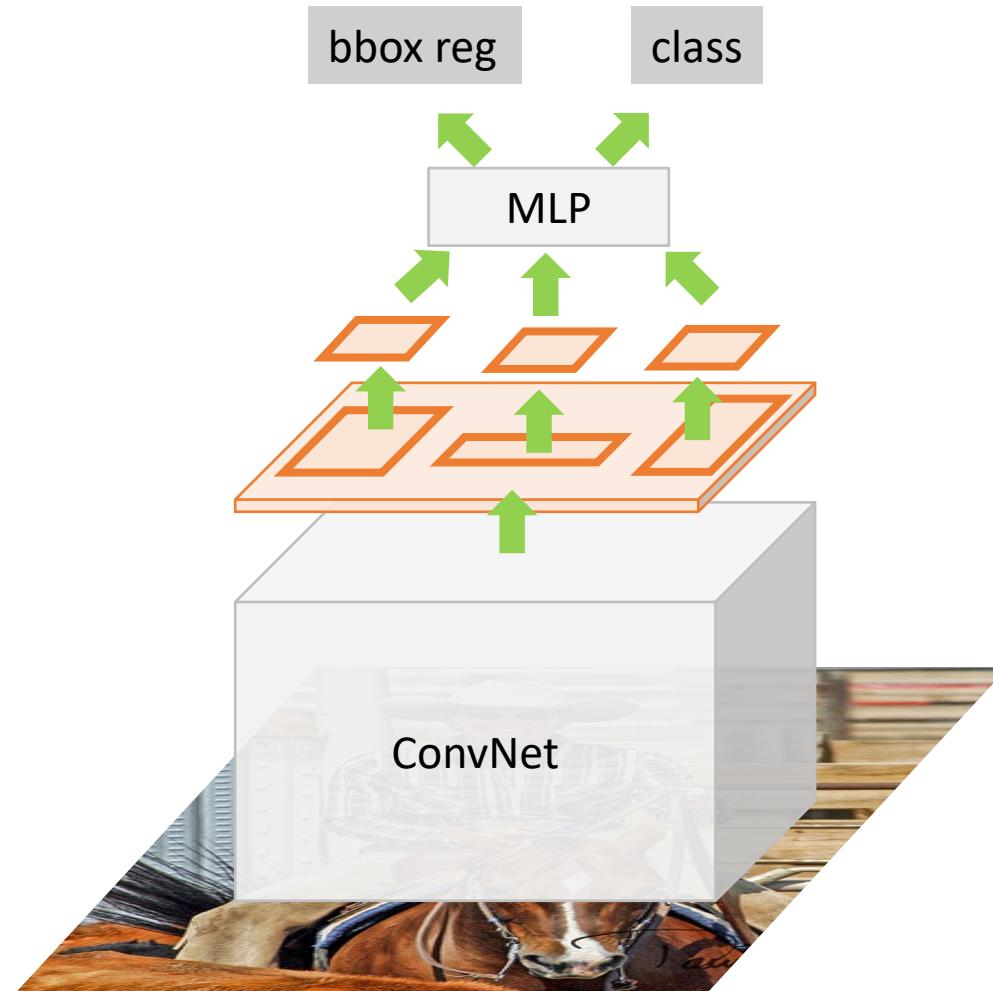
Towards Fast R-CNN (SPPnet)

- What's still wrong? Not yet “Differentiable Programming”



Fast R-CNN

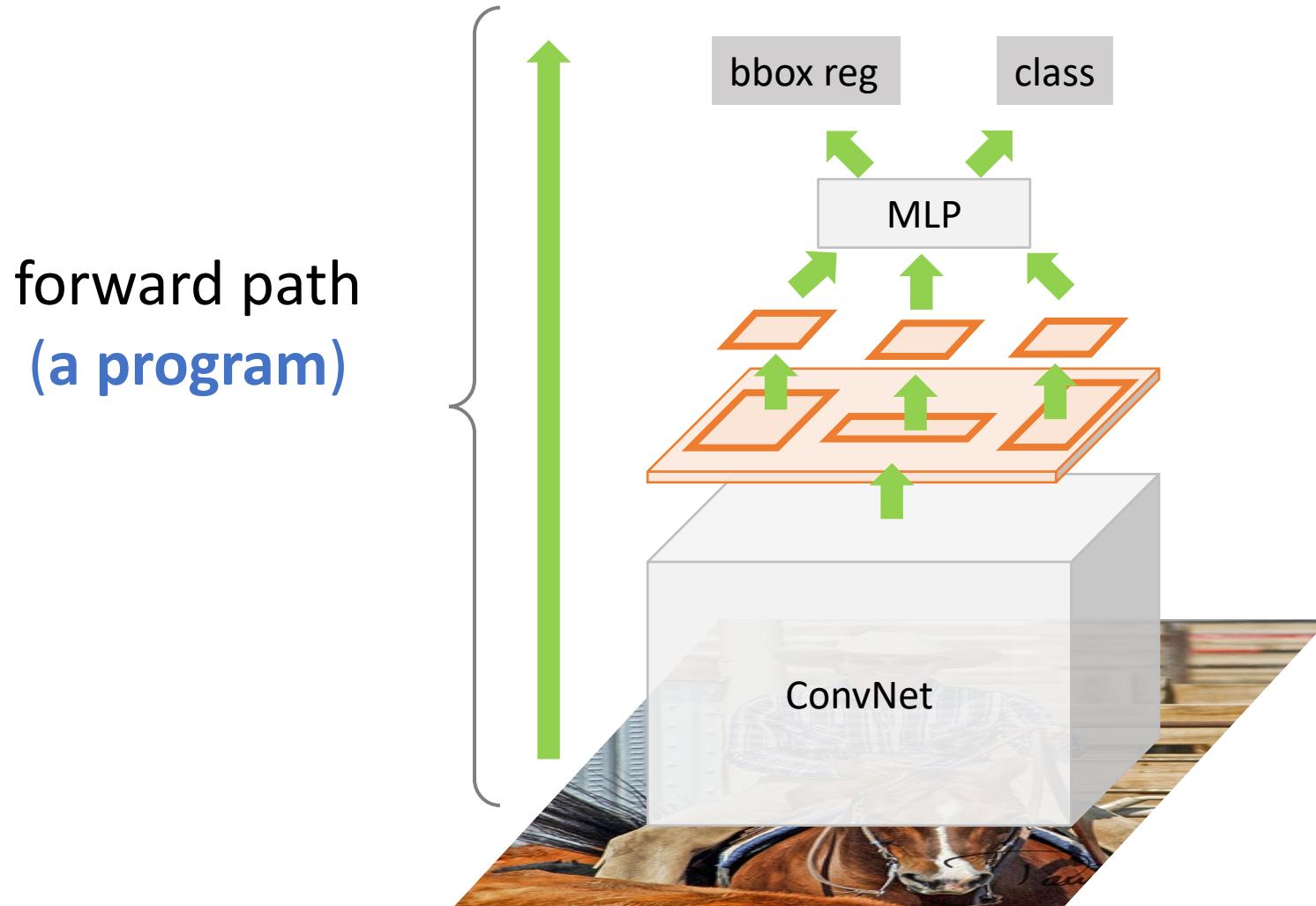
- Differentiable Programming



Slides adapted from: Ross Girshick, ICCV 2015
Ross Girshick, "Fast R-CNN", ICCV 2015

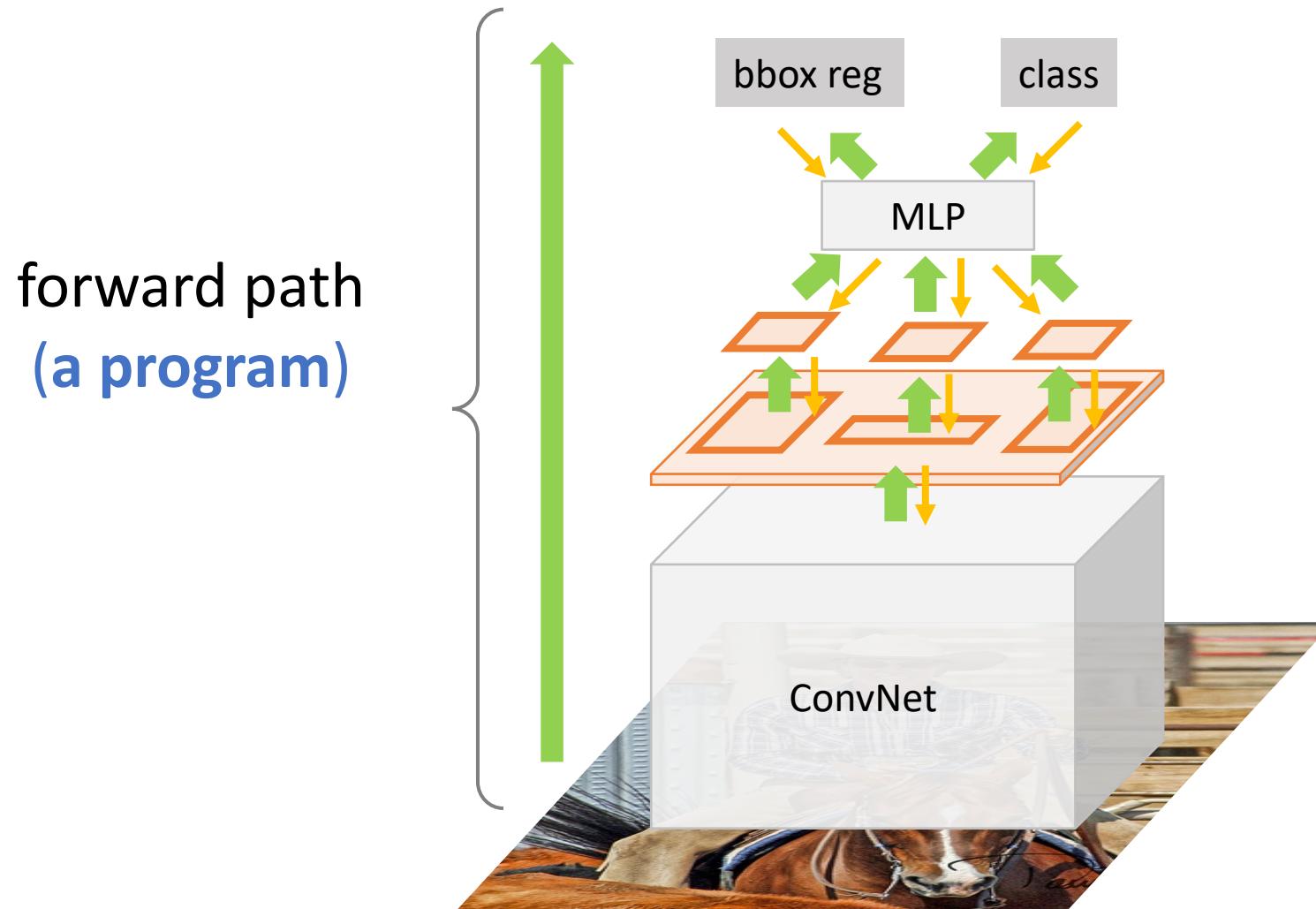
Fast R-CNN

- Differentiable Programming



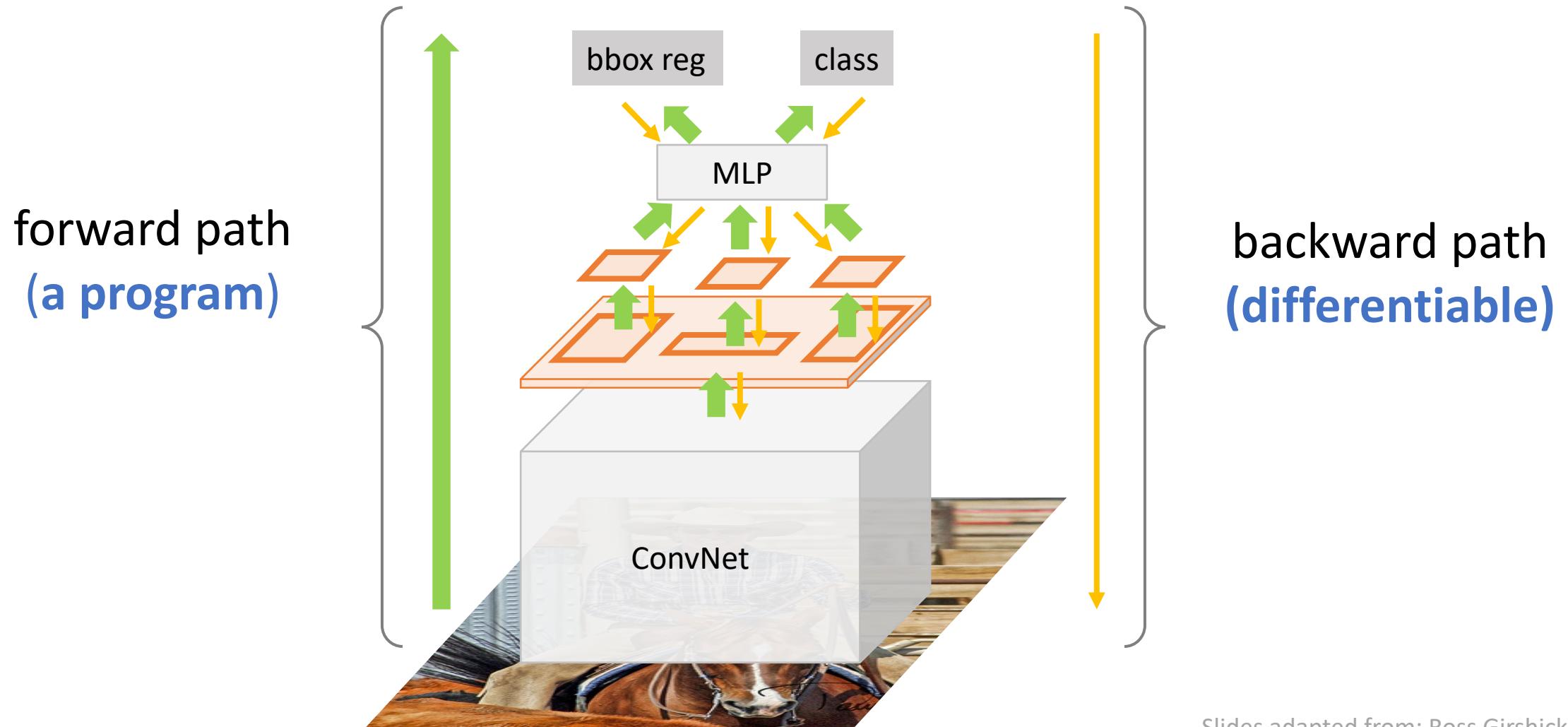
Fast R-CNN

- Differentiable Programming



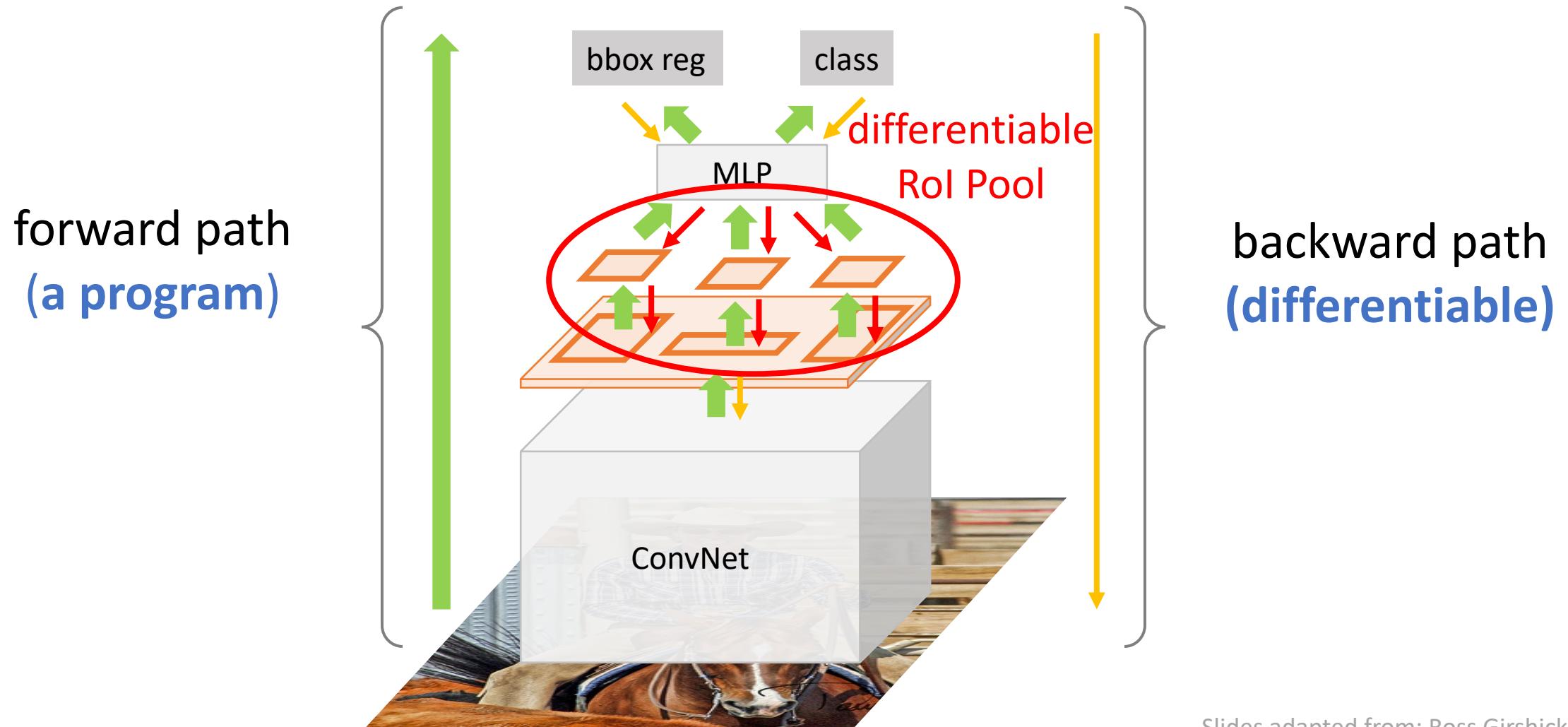
Fast R-CNN

- Differentiable Programming



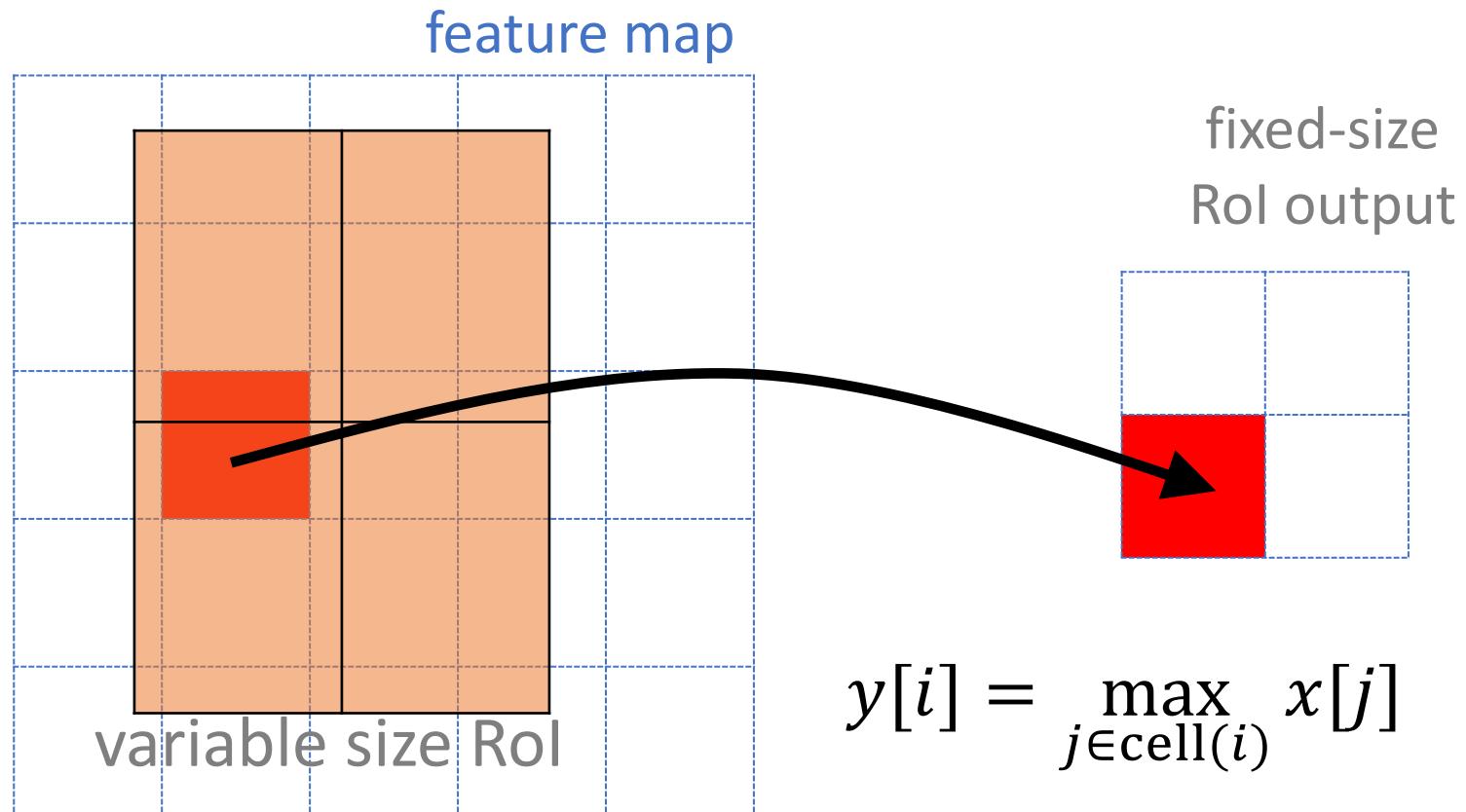
Fast R-CNN

- Differentiable Programming

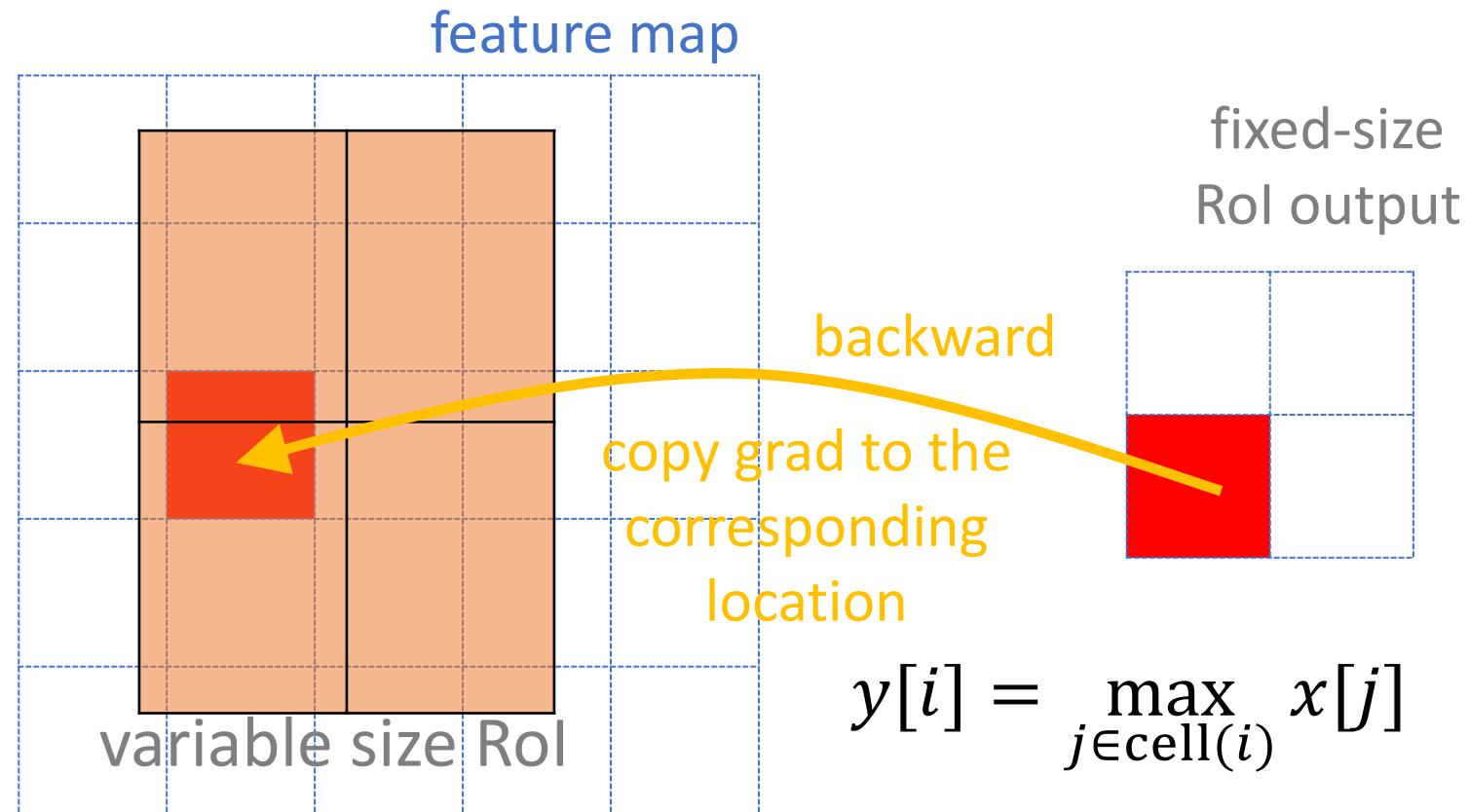


Slides adapted from: Ross Girshick, ICCV 2015
Ross Girshick, "Fast R-CNN", ICCV 2015

Differentiable RoI Pooling

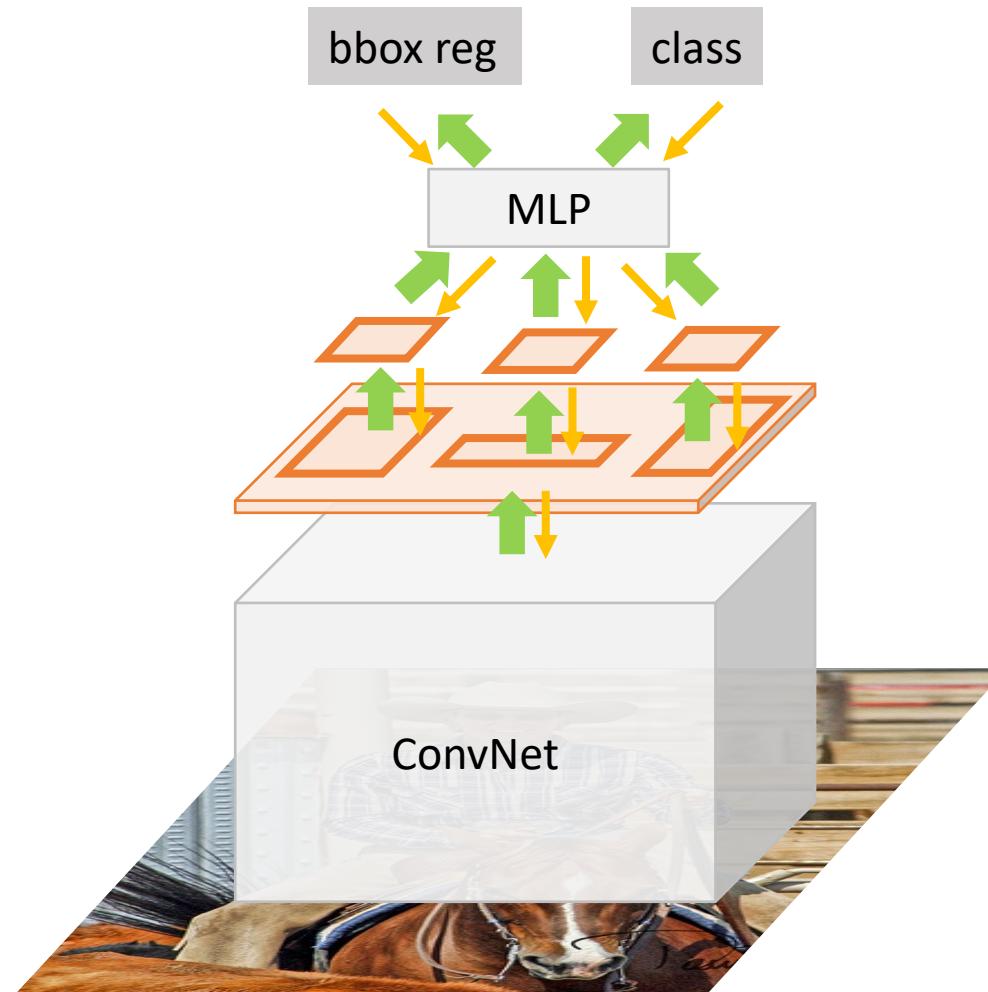


Differentiable RoI Pooling



Fast R-CNN

- **Differentiable Programming**

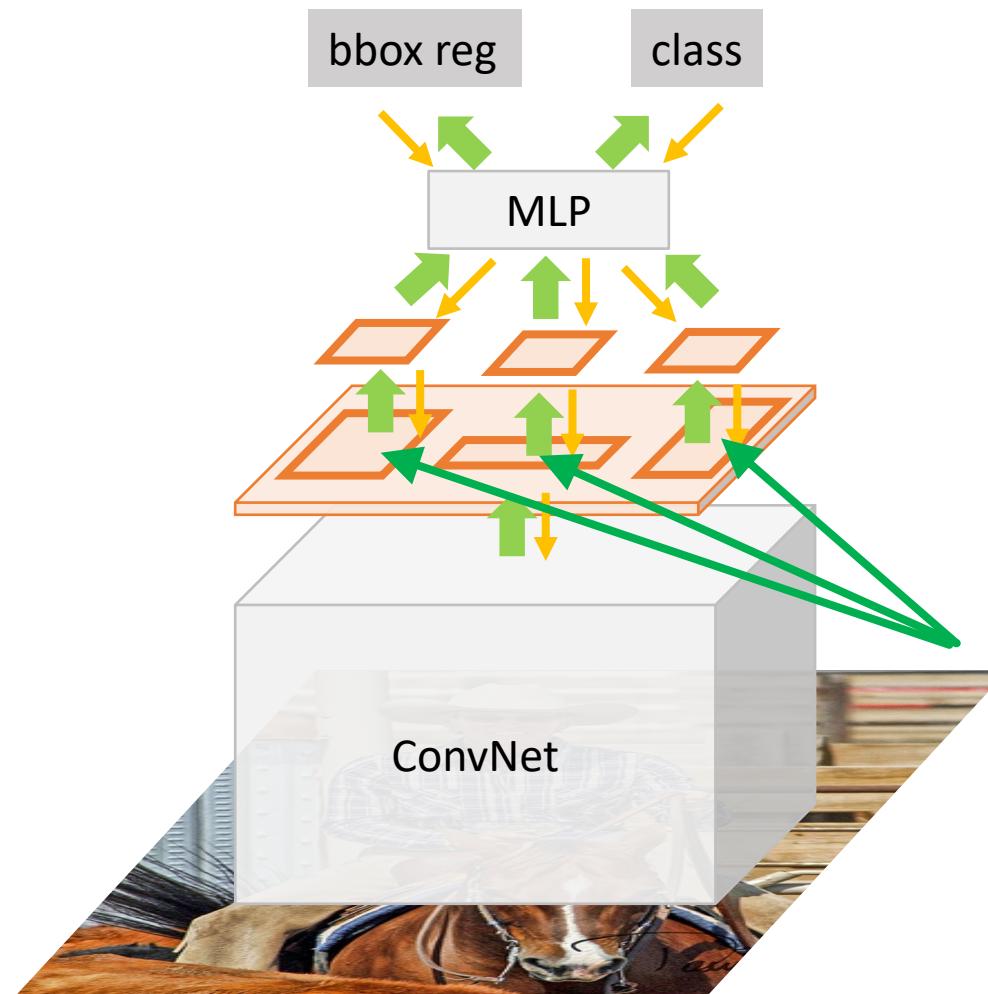


Hierarchical batching

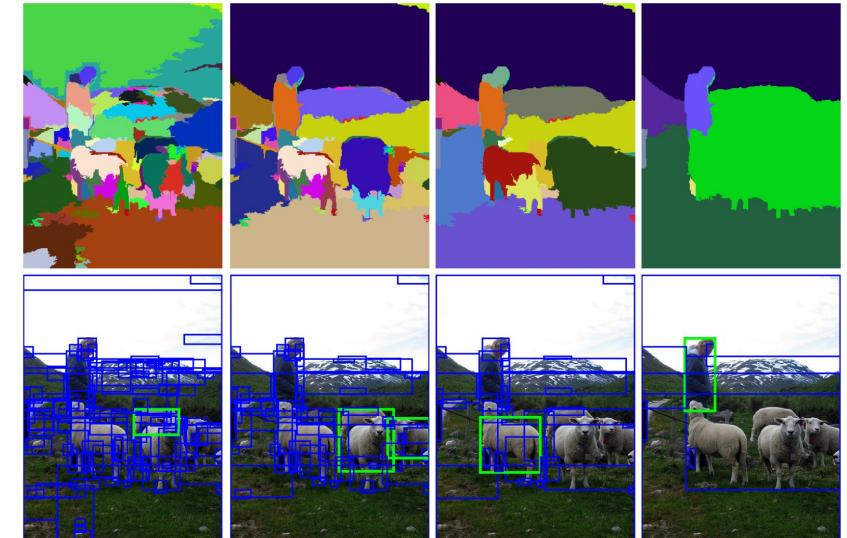
- for each image, sample a batch of Rols
- sample a batch of images

Fast R-CNN

- What's still wrong? Region proposals are not part of the program



recap: Selective Search



- Propose Regions

Faster R-CNN

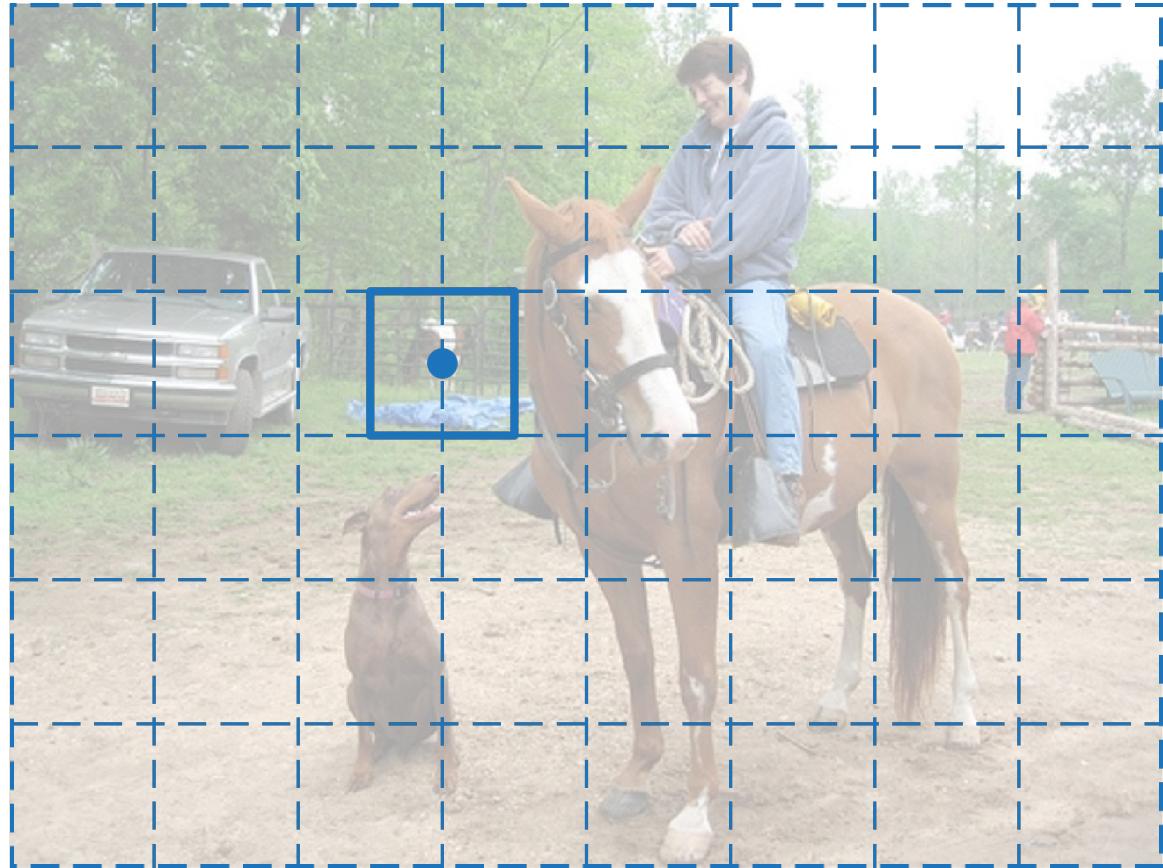
- How can neural nets produce regions?



- ConvNets are Fully Convolutional!
- Search objects w/ sliding windows

Faster R-CNN

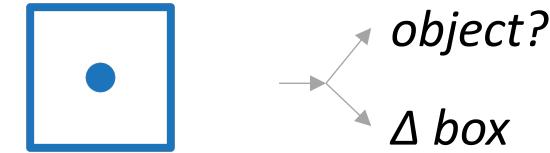
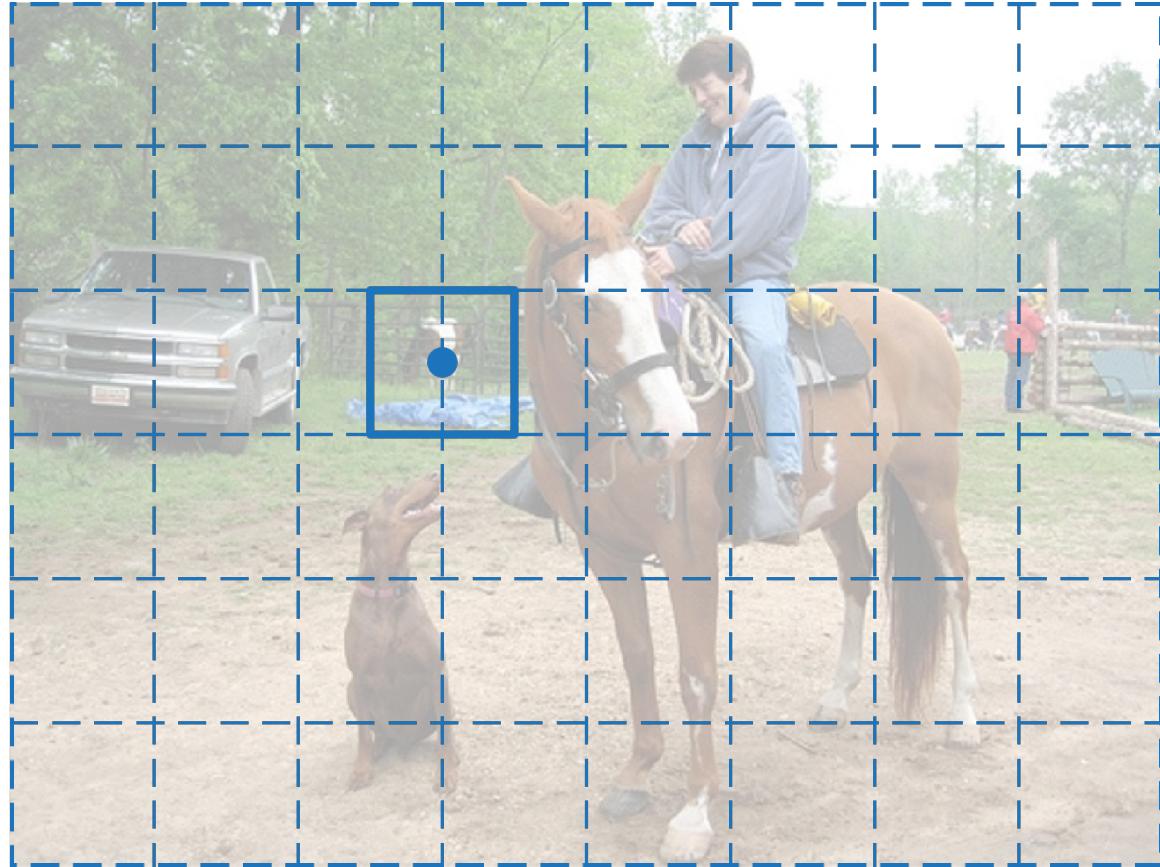
- How can neural nets produce regions?



- ConvNets are Fully Convolutional!
- Search objects w/ sliding windows

Faster R-CNN

- How can neural nets produce regions?

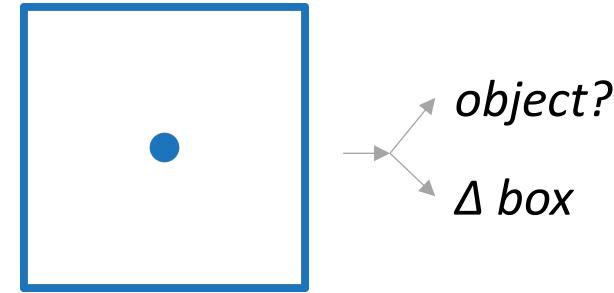
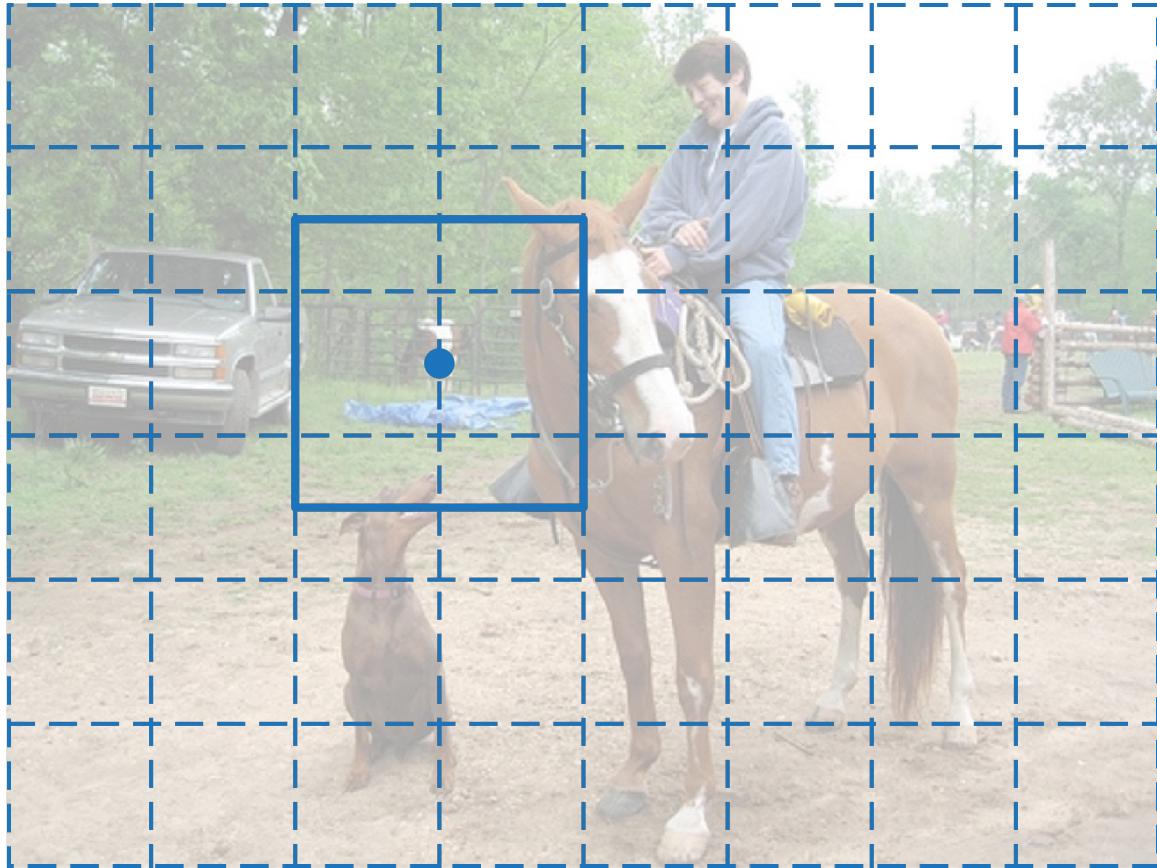


For each window, predict:

- Is it an object?
 - binary classification
- Object location w.r.t. window?
 - bbox regression

Faster R-CNN

- How can neural nets produce regions?

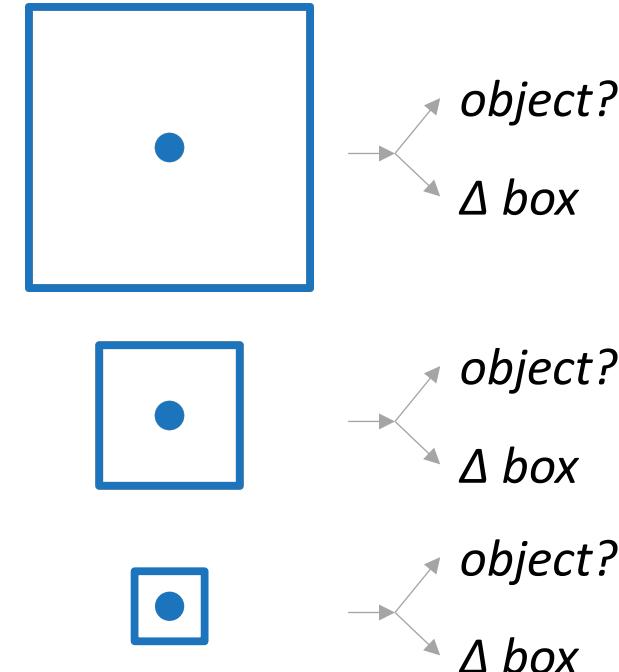
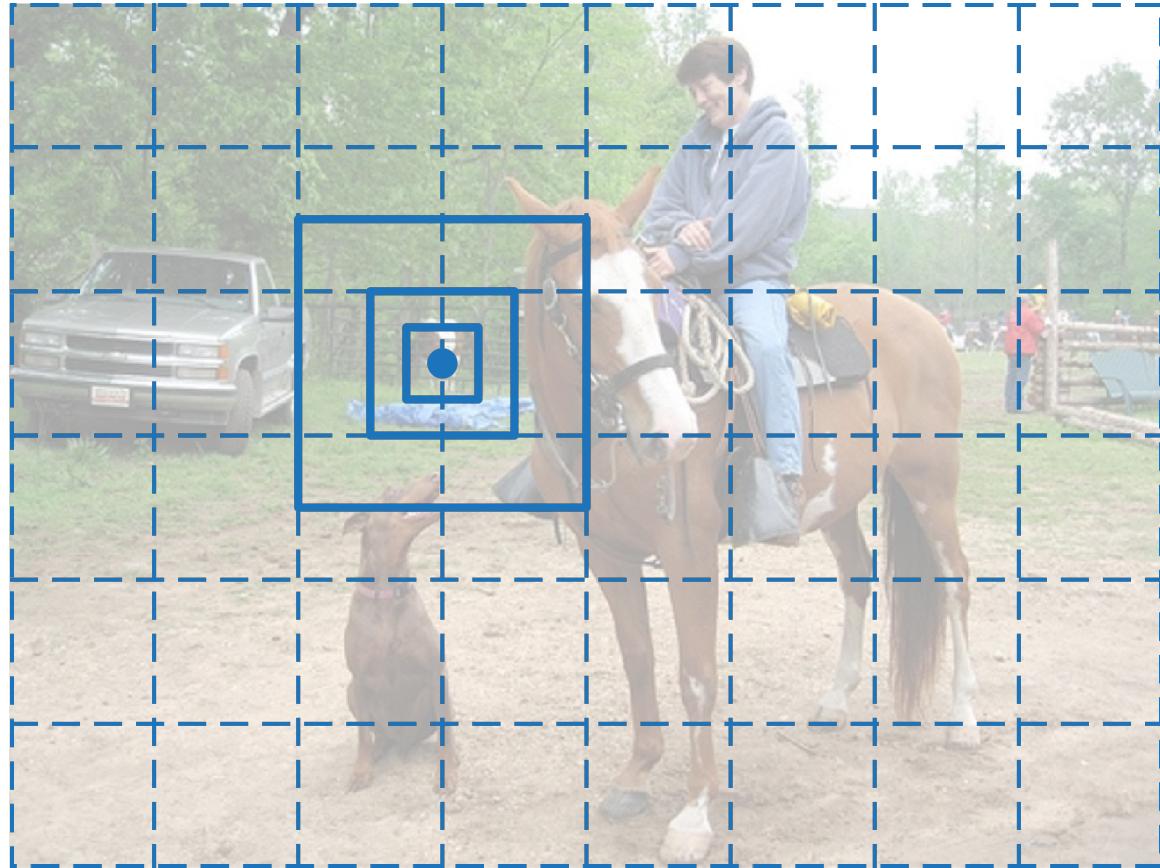


Anchor

- pred w.r.t. a reference window
- reference is called an “anchor”
- be different w/ sliding window

Faster R-CNN

- How can neural nets produce regions?

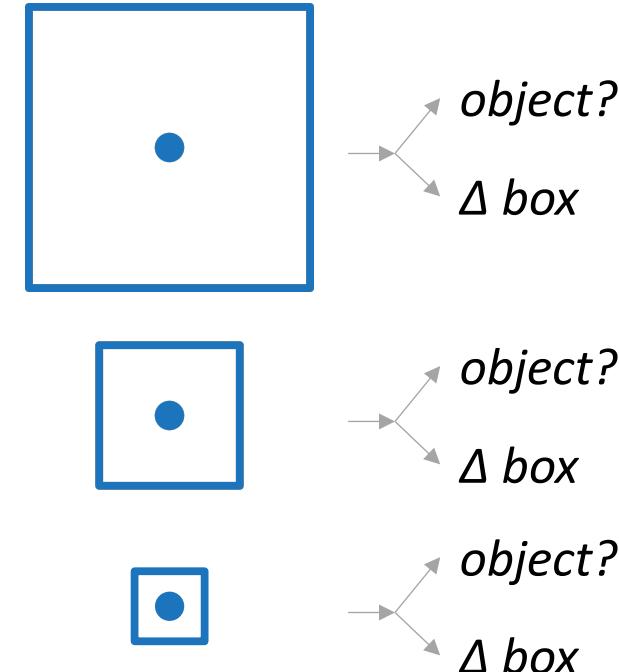
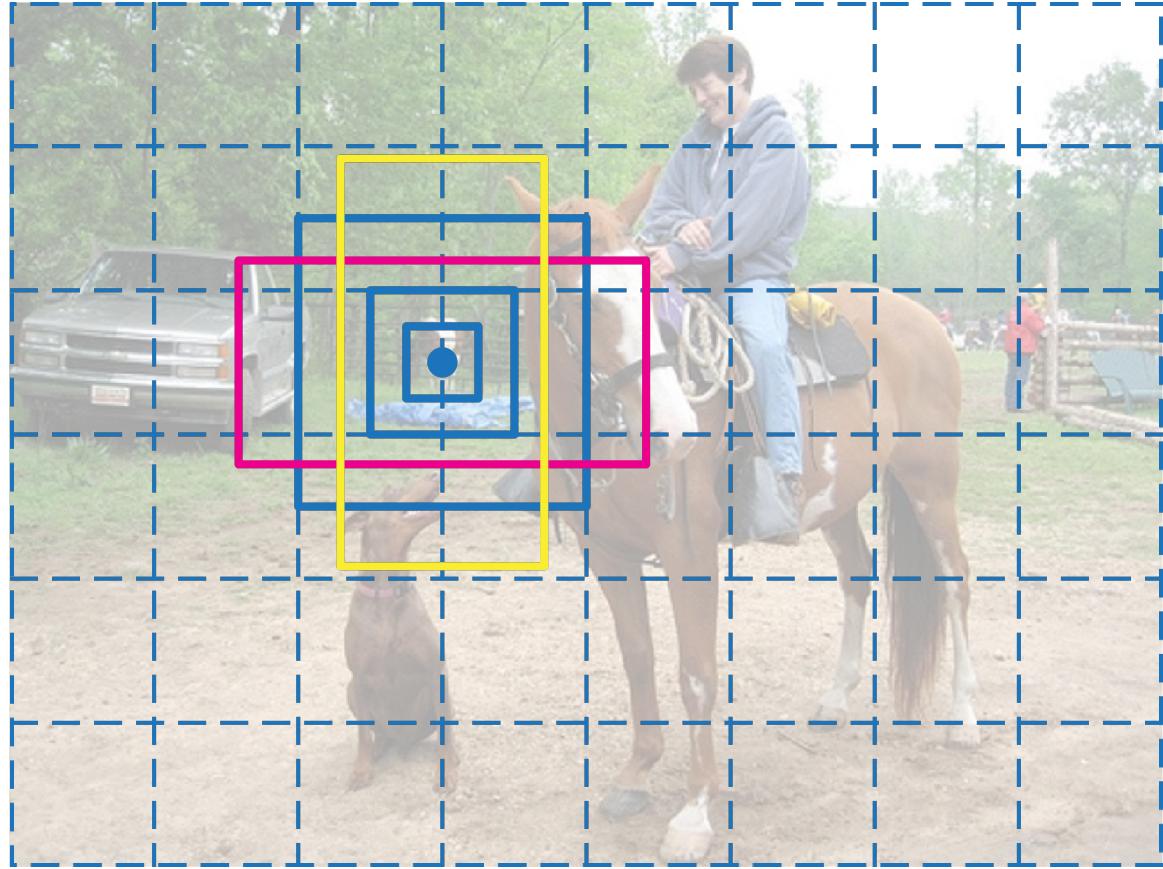


Anchors

- one sliding window, many anchors
- multiple scales

Faster R-CNN

- How can neural nets produce regions?

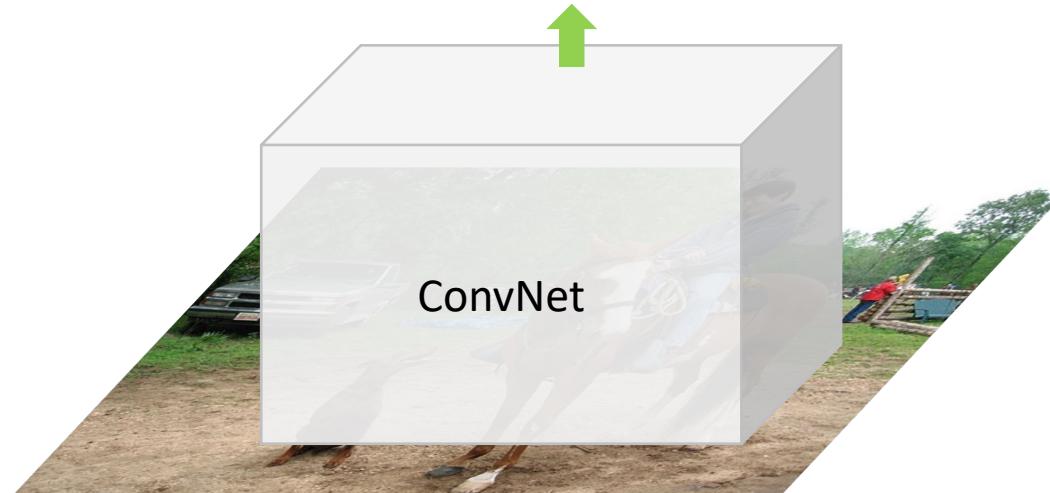


Anchors

- one sliding window, many anchors
- multiple scales
- multiple aspect ratios

Faster R-CNN: Region Proposal Network (RPN)

- How can neural nets produce regions?



- Apply ConvNet to **whole** image

Faster R-CNN: Region Proposal Network (RPN)

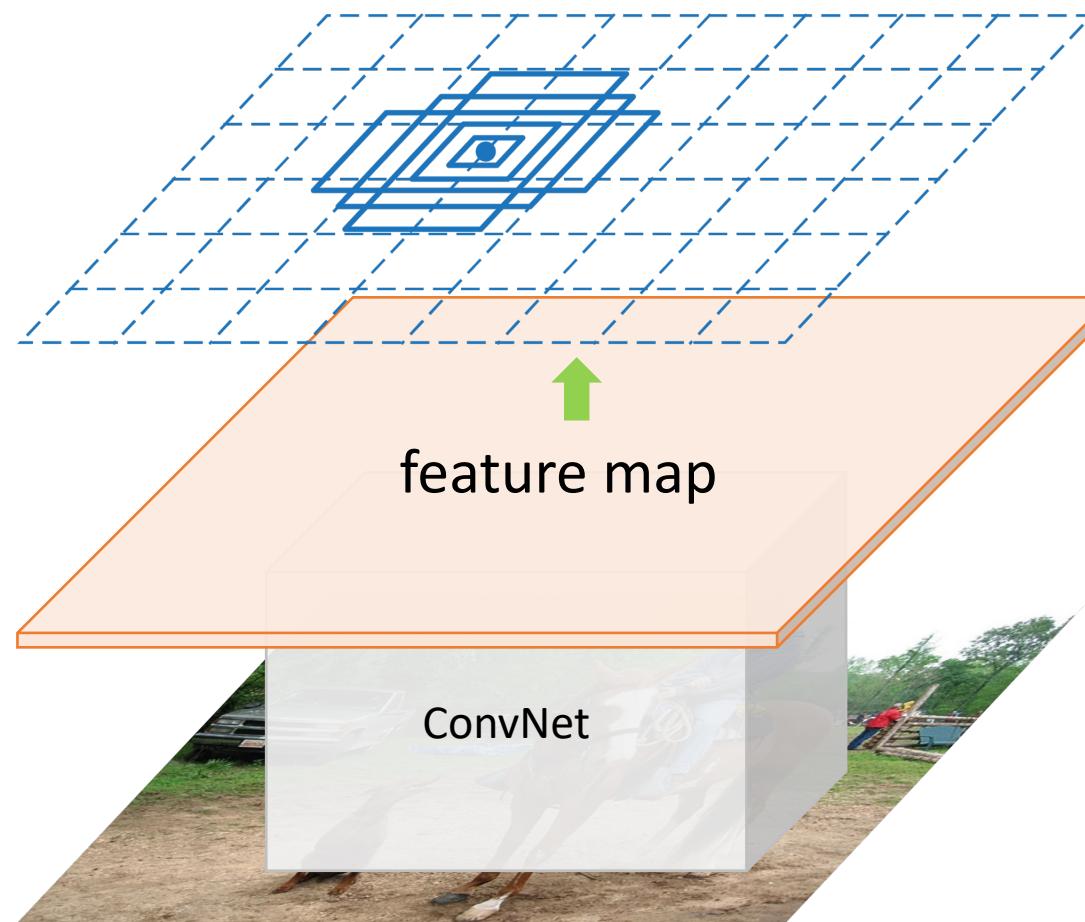
- How can neural nets produce regions?



- Apply ConvNet to **whole image**

Faster R-CNN: Region Proposal Network (RPN)

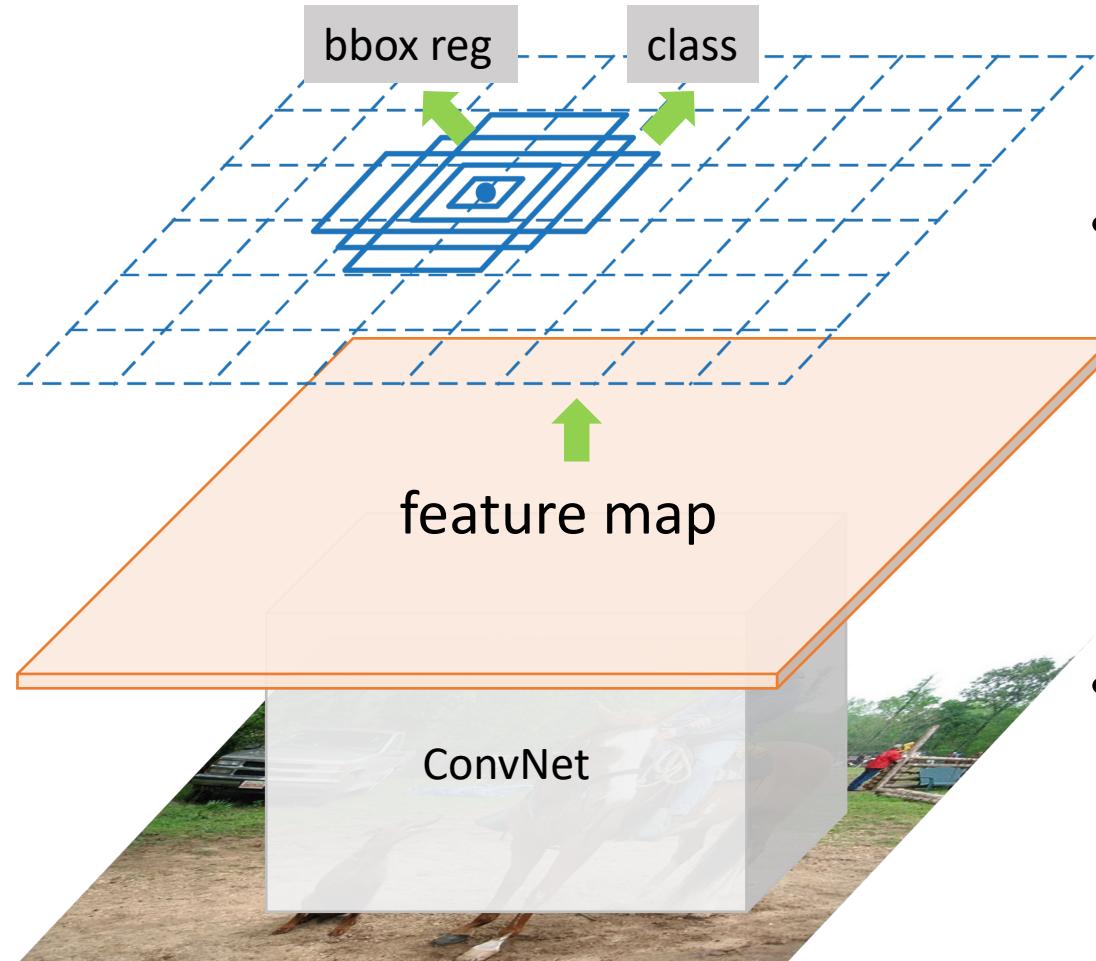
- How can neural nets produce regions?



- Convs for k anchors
 - $2k$ -d: binary classification
 - $4k$ -d: bbox regression
- Apply ConvNet to **whole image**

Faster R-CNN: Region Proposal Network (RPN)

- How can neural nets produce regions?



- Loss: reg + class
 - over all locations
 - over all anchors
- Convs for k anchors
 - $2k$ -d: binary classification
 - $4k$ -d: bbox regression
- Apply ConvNet to **whole** image

Faster R-CNN: Region Proposal Network (RPN)

Example: 1000 regions proposed by RPN



Faster R-CNN: End-to-End Object Detection

- End-to-End Differentiable Programming

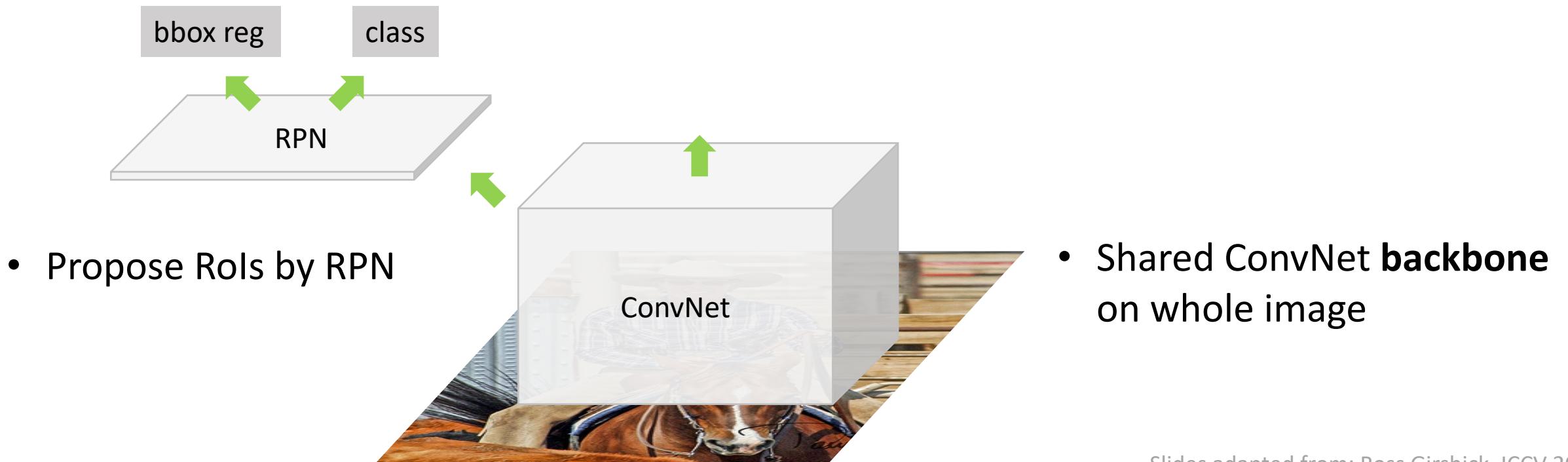


- Shared ConvNet **backbone** on whole image

Slides adapted from: Ross Girshick, ICCV 2015

Faster R-CNN: End-to-End Object Detection

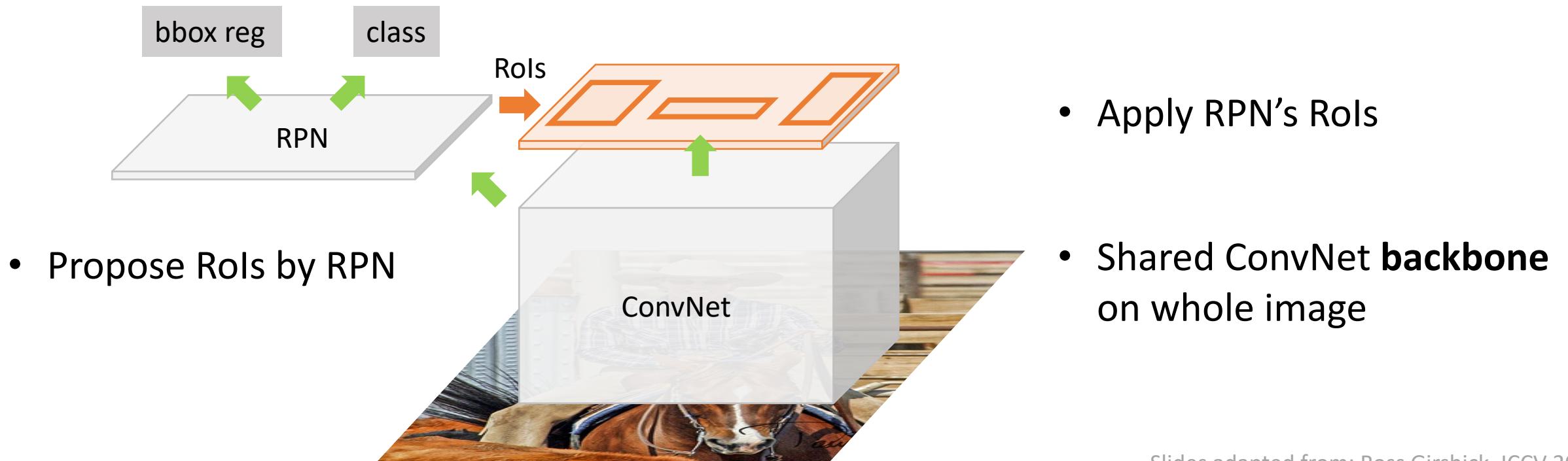
- End-to-End Differentiable Programming



Slides adapted from: Ross Girshick, ICCV 2015

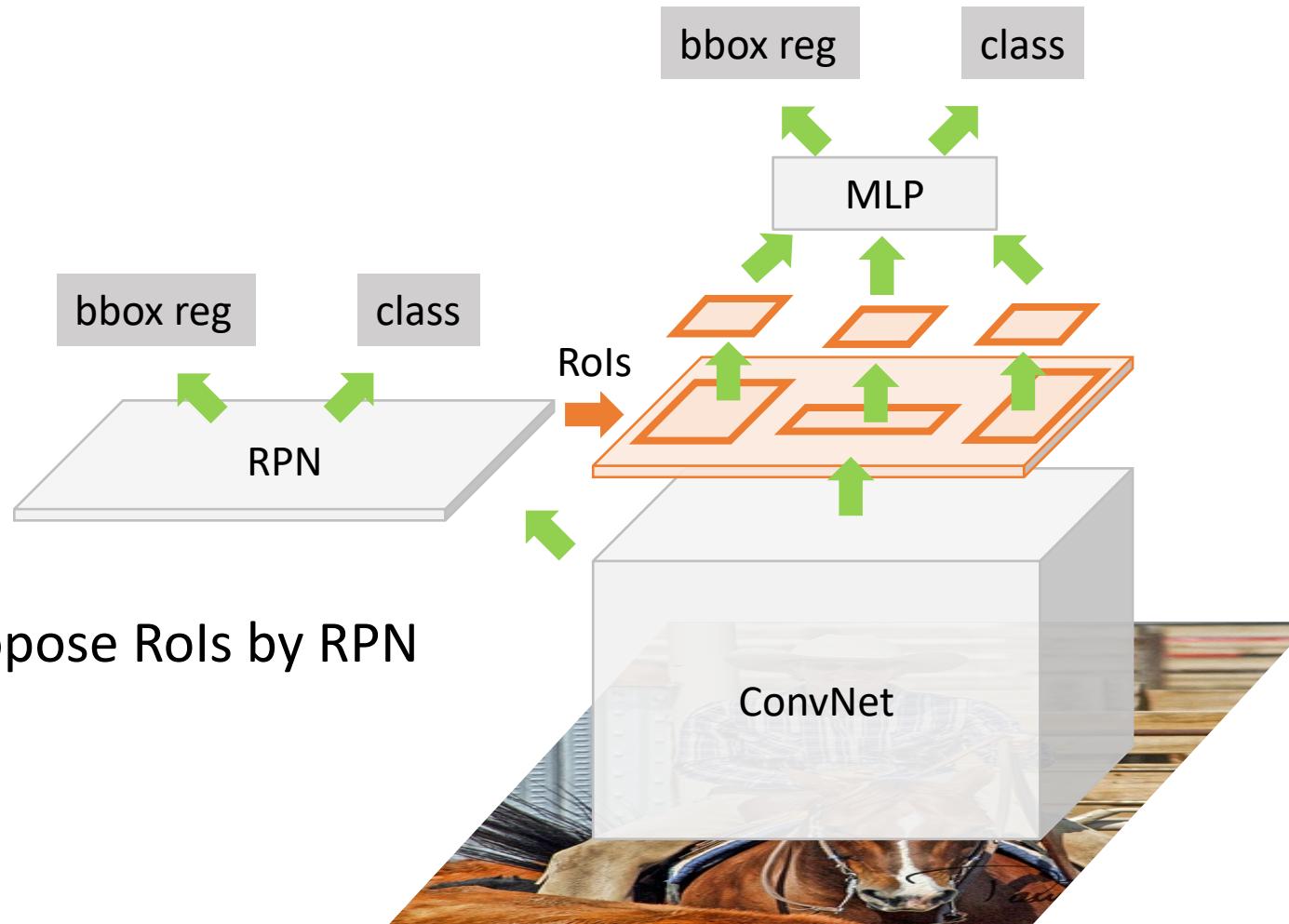
Faster R-CNN: End-to-End Object Detection

- End-to-End Differentiable Programming



Faster R-CNN: End-to-End Object Detection

- End-to-End Differentiable Programming

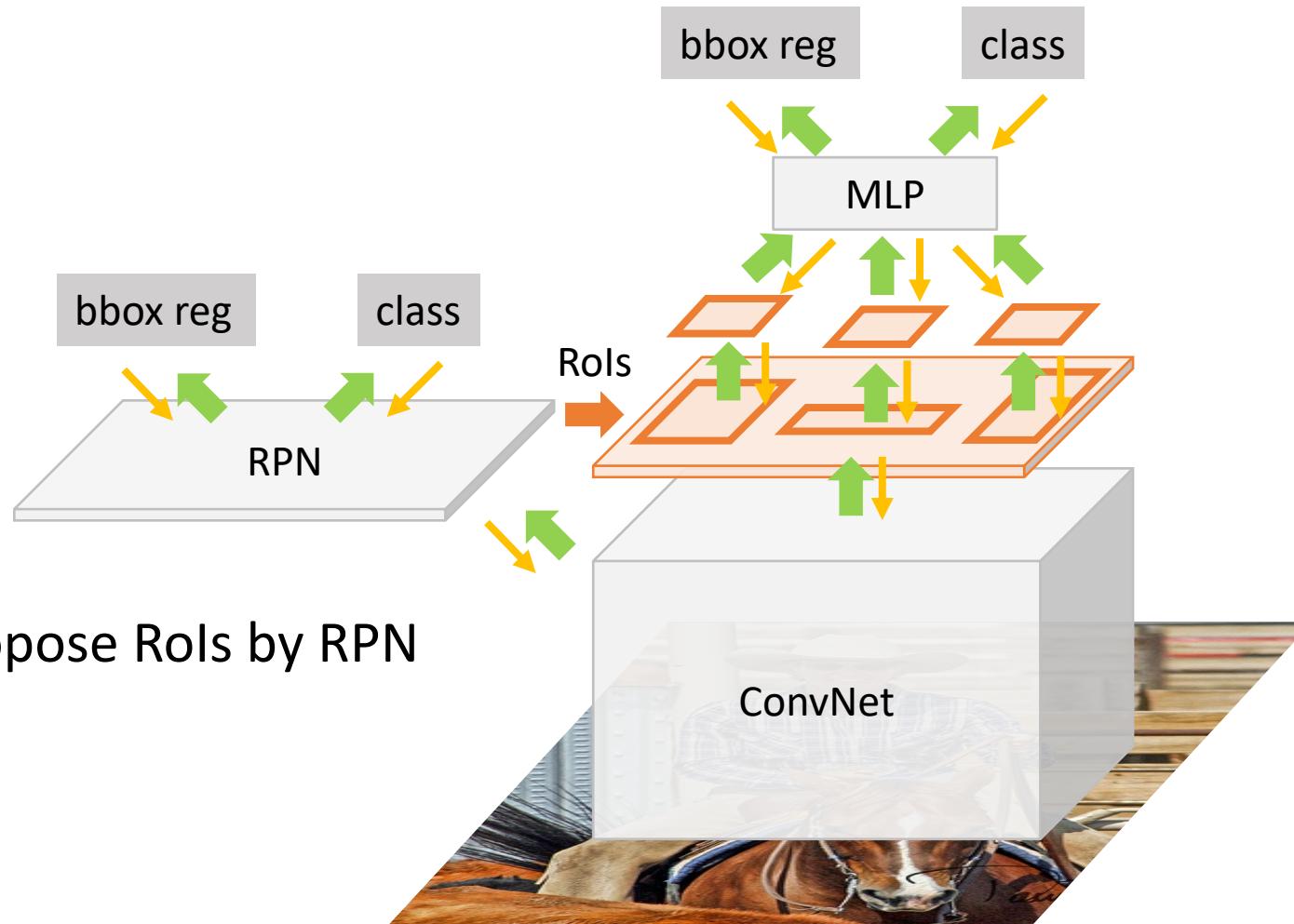


- Propose Rols by RPN

- Same as Fast R-CNN
- Apply RPN's Rols
- Shared ConvNet **backbone** on whole image

Faster R-CNN: End-to-End Object Detection

- End-to-End Differentiable Programming

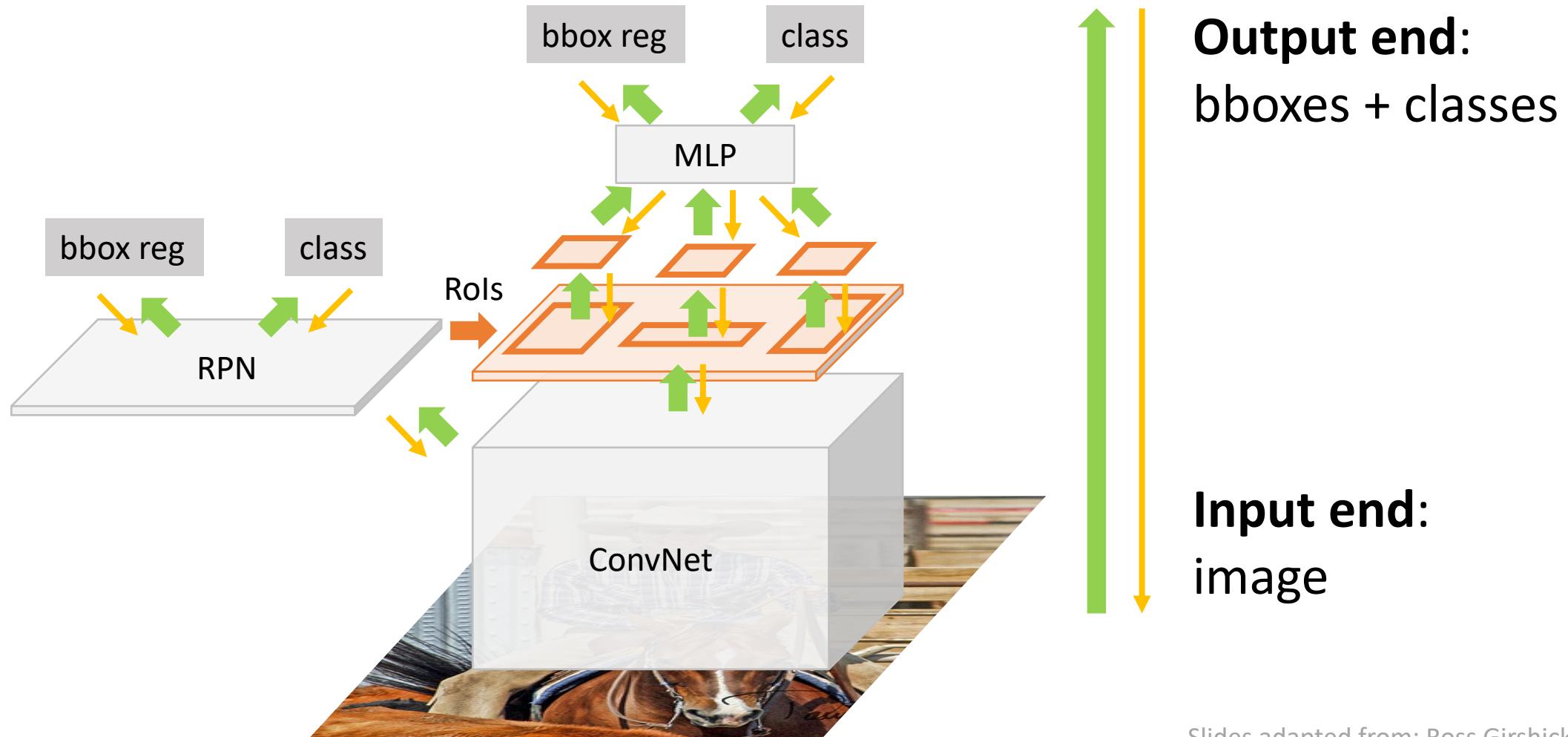


- Propose Rols by RPN

- Same as Fast R-CNN
- Apply RPN's Rols
- Shared ConvNet **backbone** on whole image

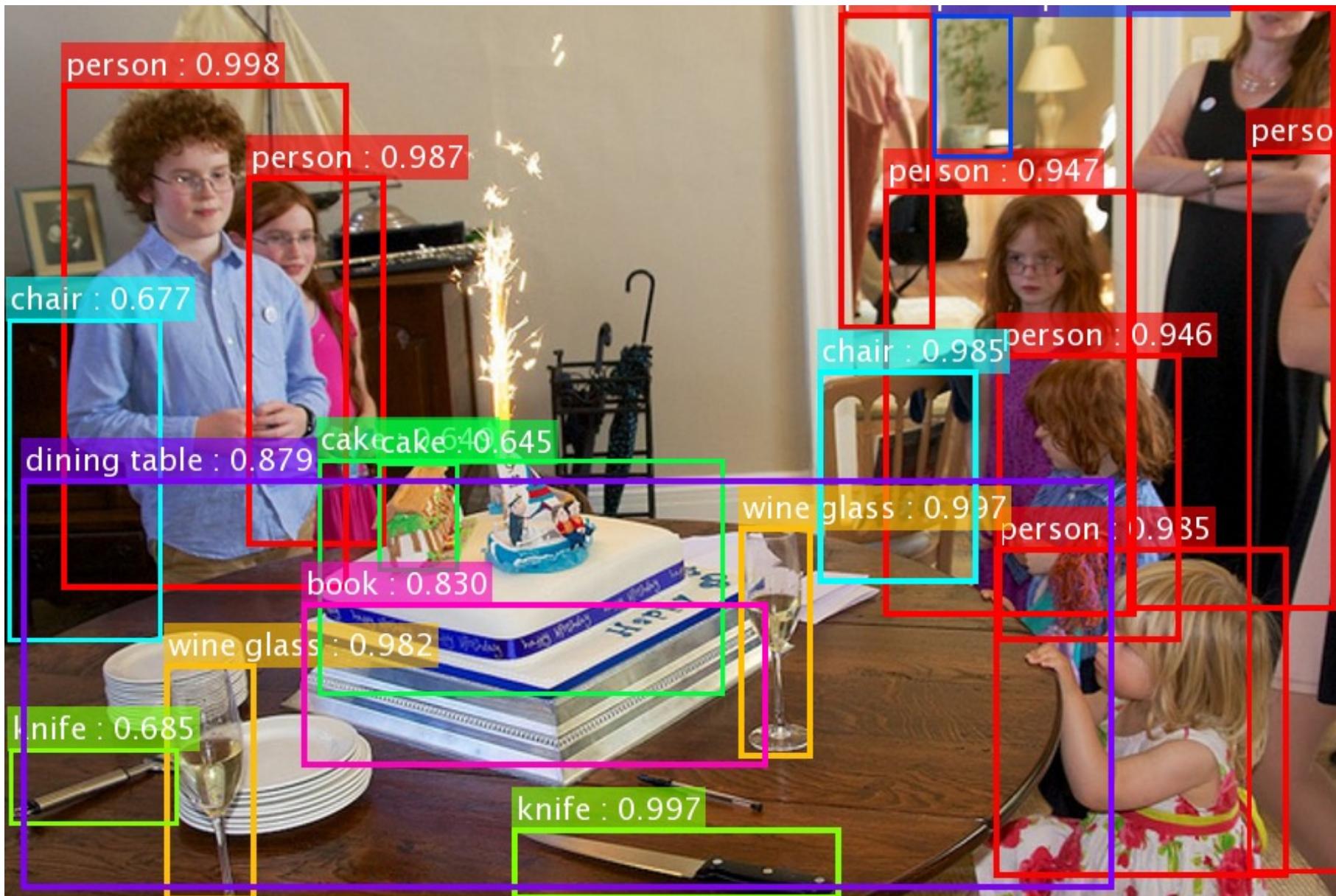
Faster R-CNN: End-to-End Object Detection

- End-to-End Differentiable Programming

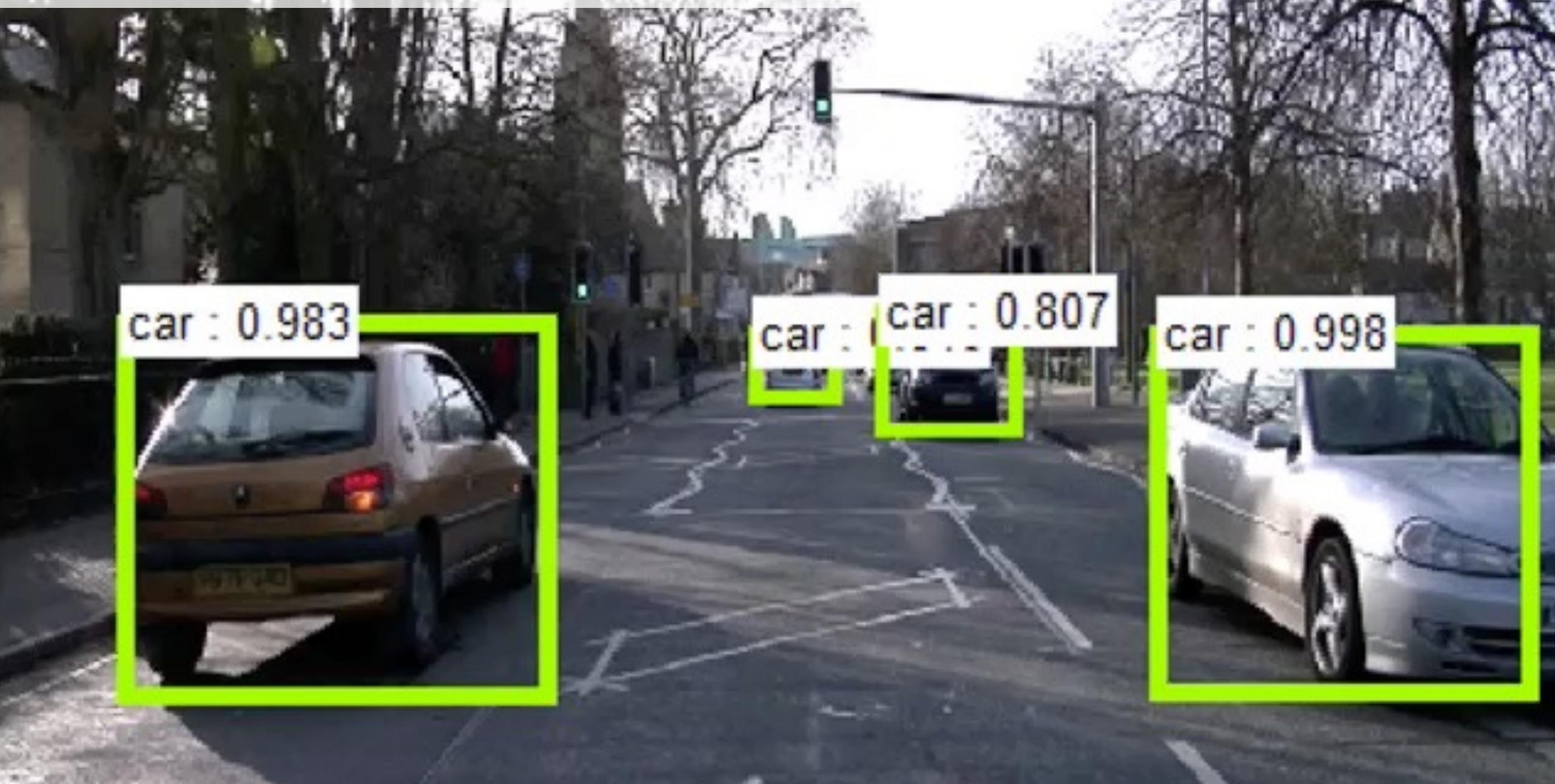


Slides adapted from: Ross Girshick, ICCV 2015

Faster R-CNN result in 2015

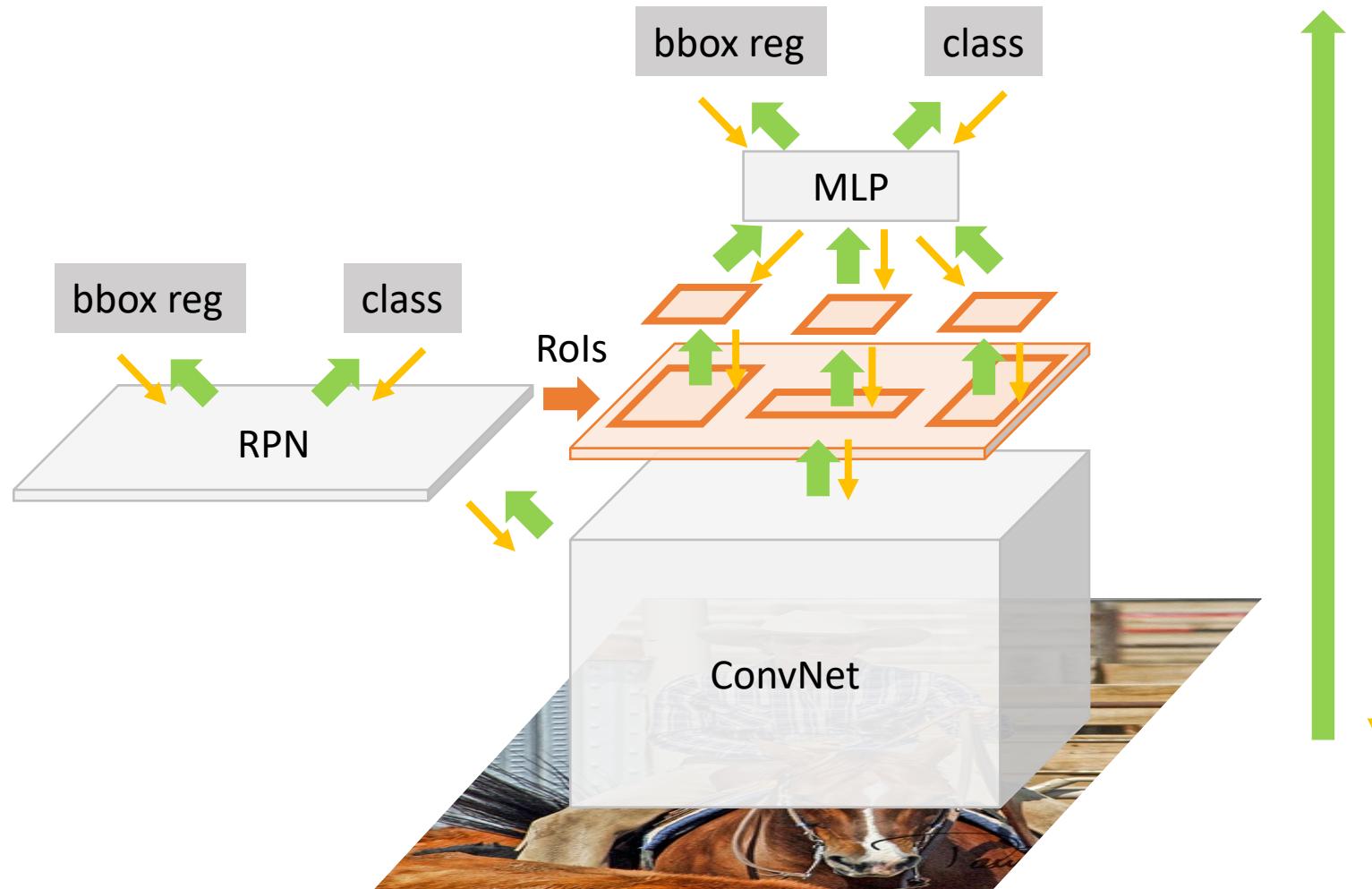


Faster R-CNN result in 2015



Faster R-CNN: End-to-End Object Detection

- What's still wrong? Single-scale feature maps for multi-scale objects



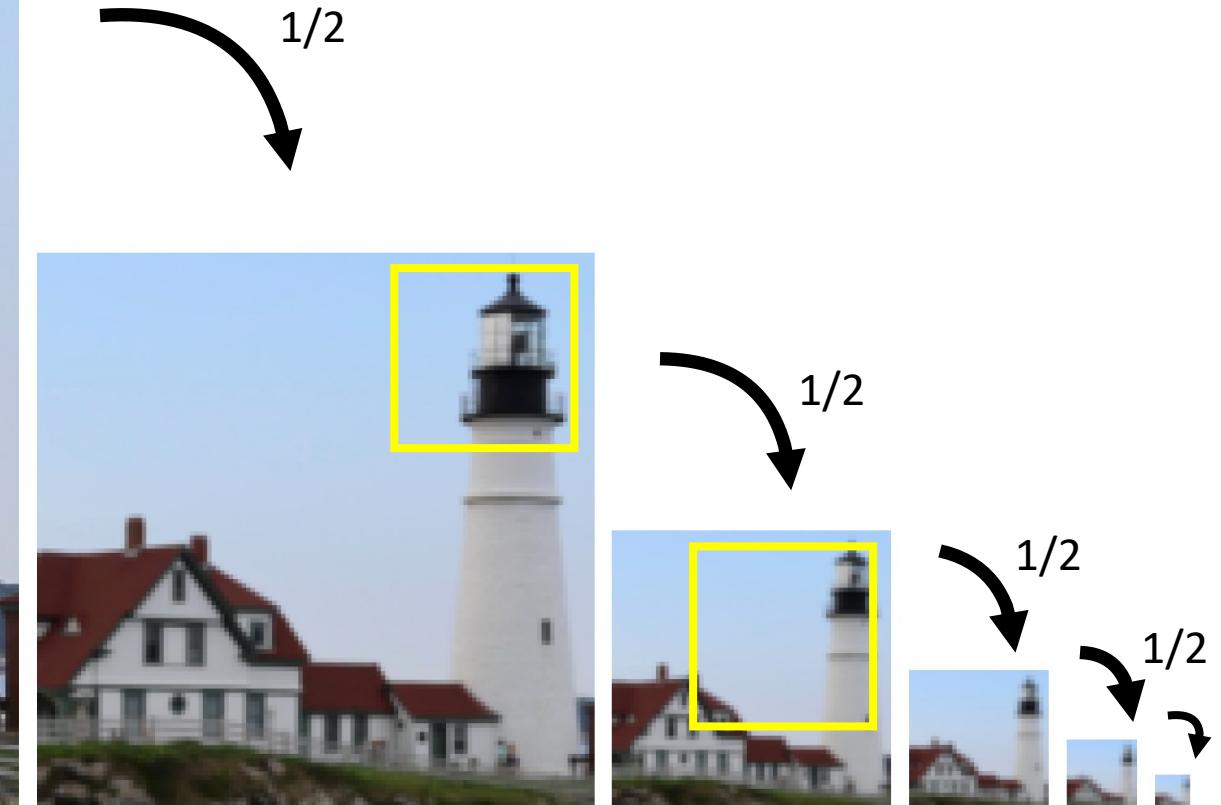
Slides adapted from: Ross Girshick, ICCV 2015

Image Pyramid

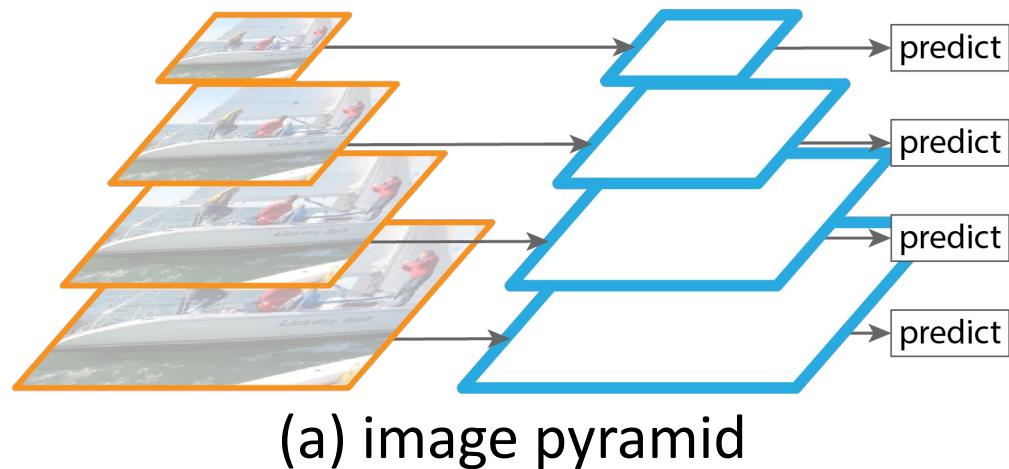


Scale-invariance

- multi-scale input
- scale-independent window

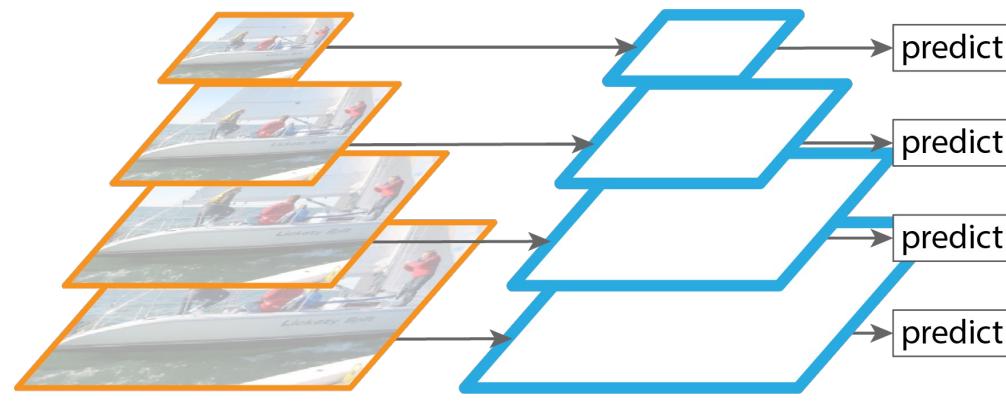


Pyramid and Neural Networks

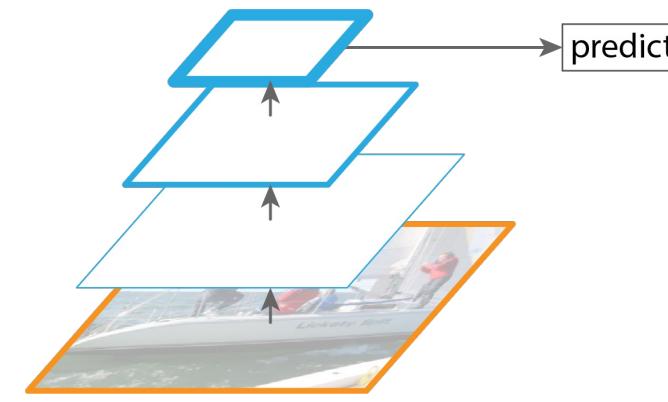


- Apply the same model (e.g., Faster R-CNN) to multi-scale input images

Pyramid and Neural Networks



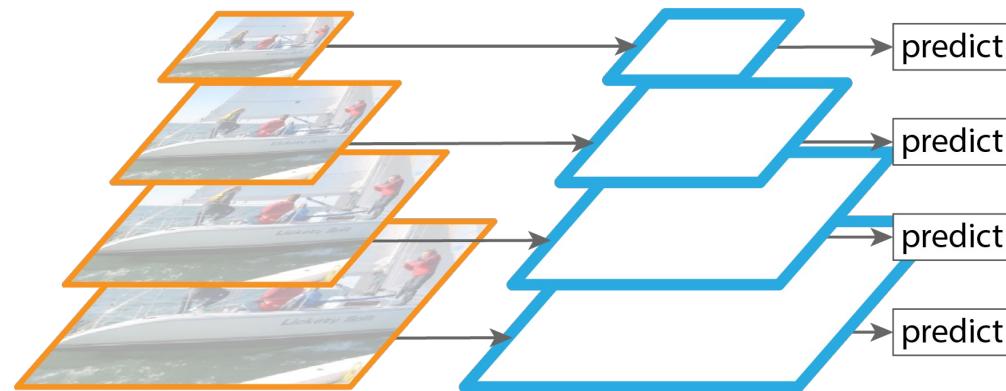
(a) image pyramid



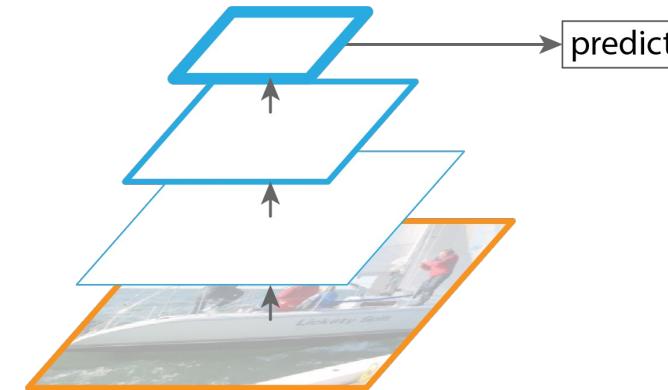
(b) single feature map

- Apply the model to the single, last (lowest-resolution) feature map

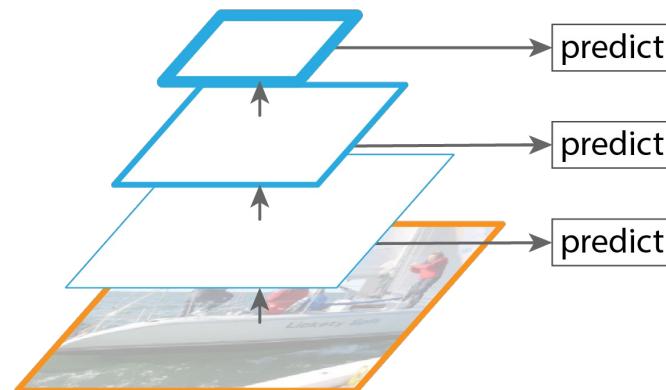
Pyramid and Neural Networks



(a) image pyramid



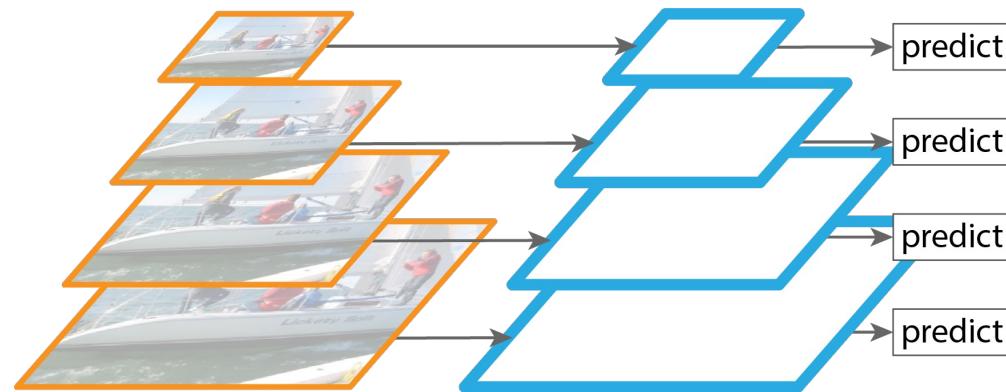
(b) single feature map



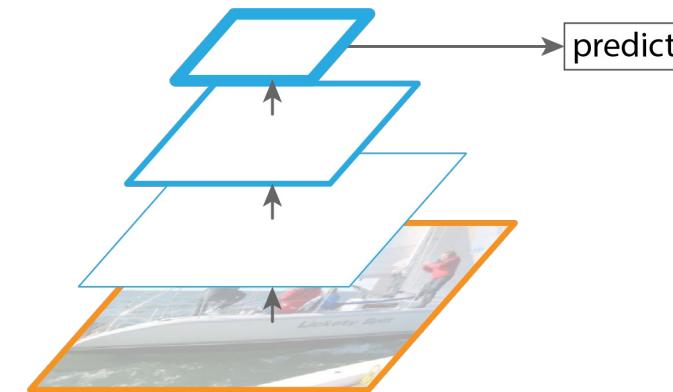
(c) pyramid feature hierarchy

- Reuse the “feature hierarchy” of a ConvNet as a pyramid (e.g., $1/2\times$, $1/4\times$, $1/8\times$)

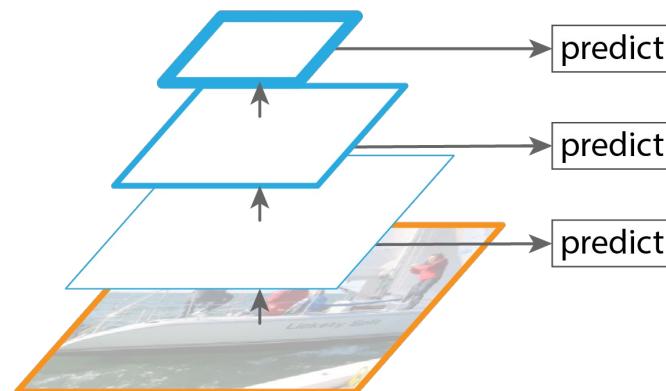
Pyramid and Neural Networks



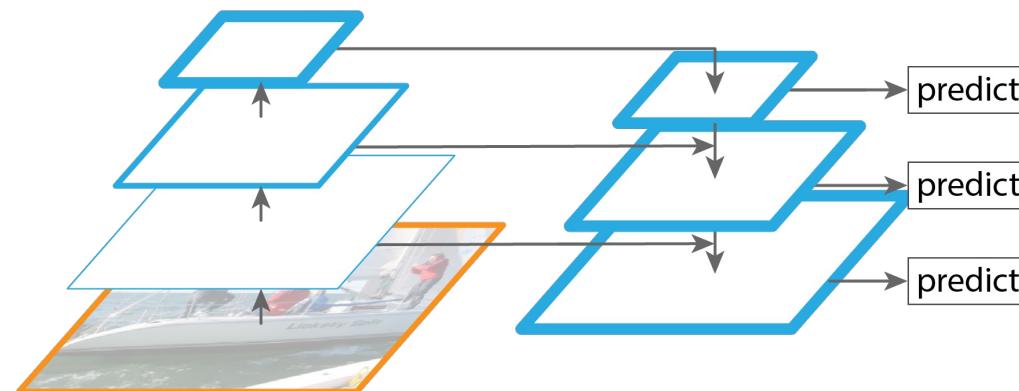
(a) image pyramid



(b) single feature map



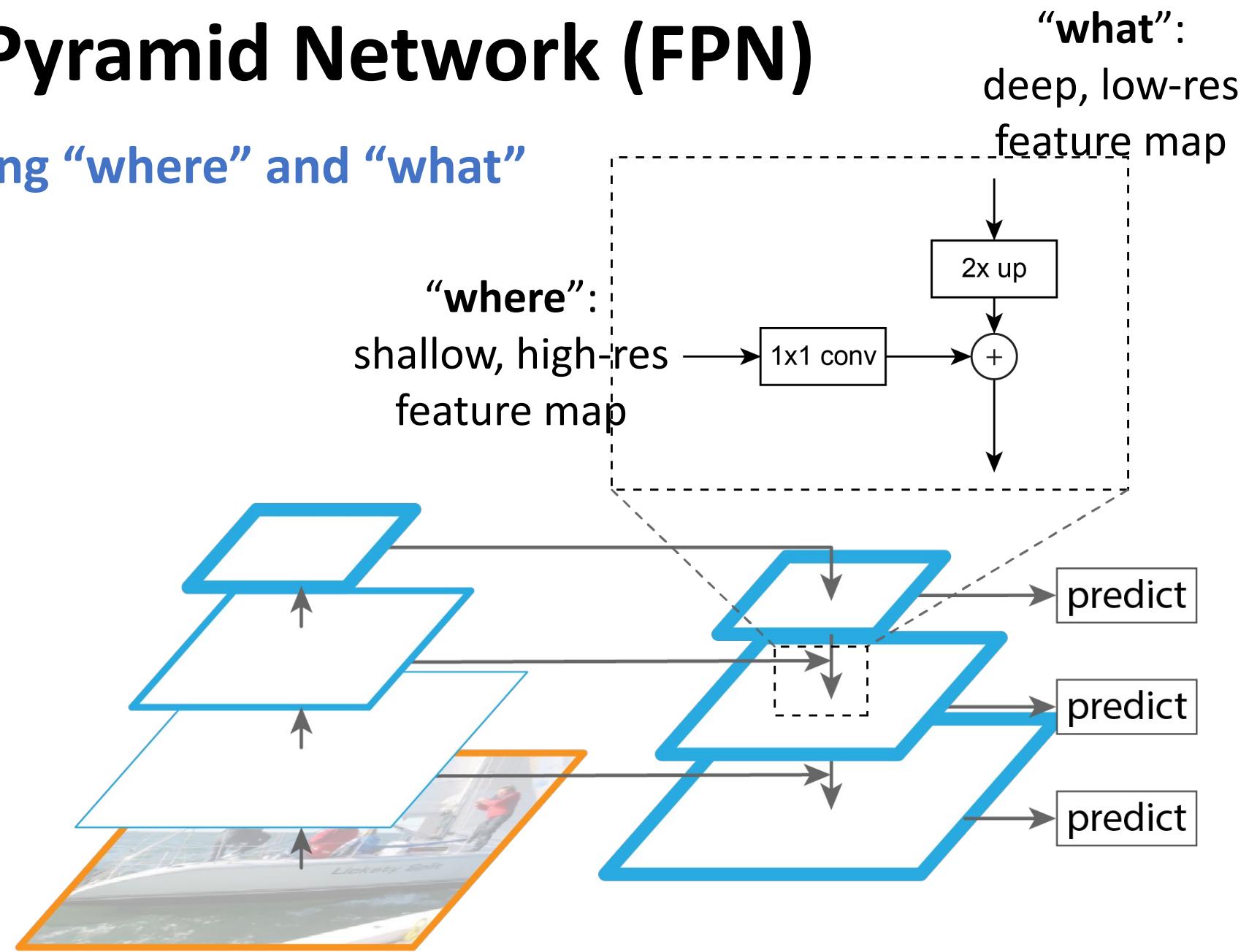
(c) pyramid feature hierarchy



(d) Feature Pyramid Network

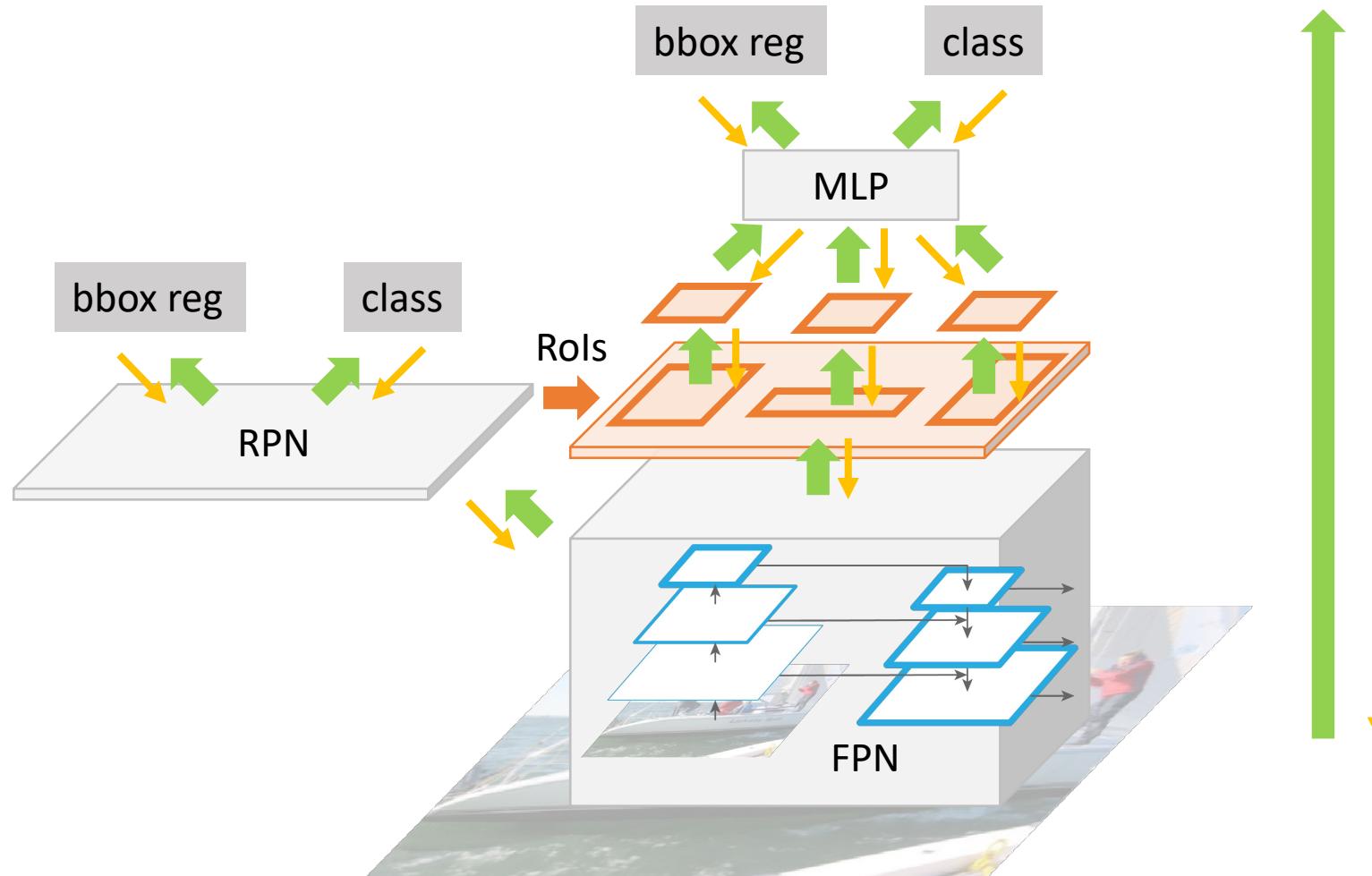
Feature Pyramid Network (FPN)

- Combining “where” and “what”



Feature Pyramid Network for Faster R-CNN

- Apply RPN and Rols to Feature Pyramid

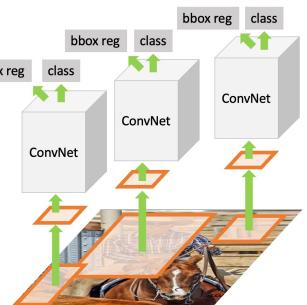


Slides adapted from: Ross Girshick, ICCV 2015

Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, "Feature Pyramid Networks for Object Detection", CVPR 2017

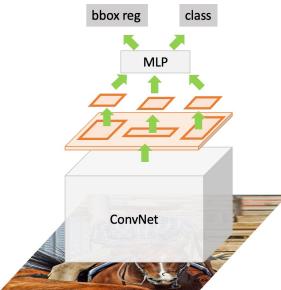
Evolution of R-CNN

R-CNN



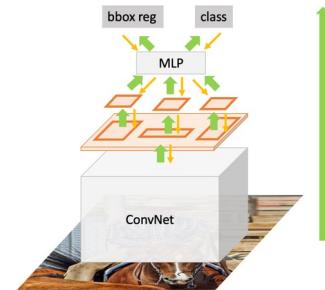
DL for detection
made work!

Fast R-CNN
“v0.1” (SPPnet)



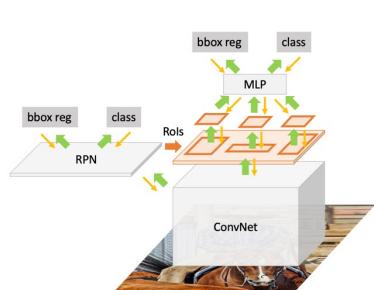
Fully convolutional
feature maps

Fast R-CNN



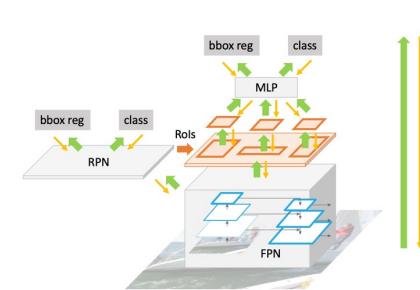
Differentiable
Programming

Faster R-CNN



RPN &
End-to-End detection

Faster R-CNN
w/ FPN

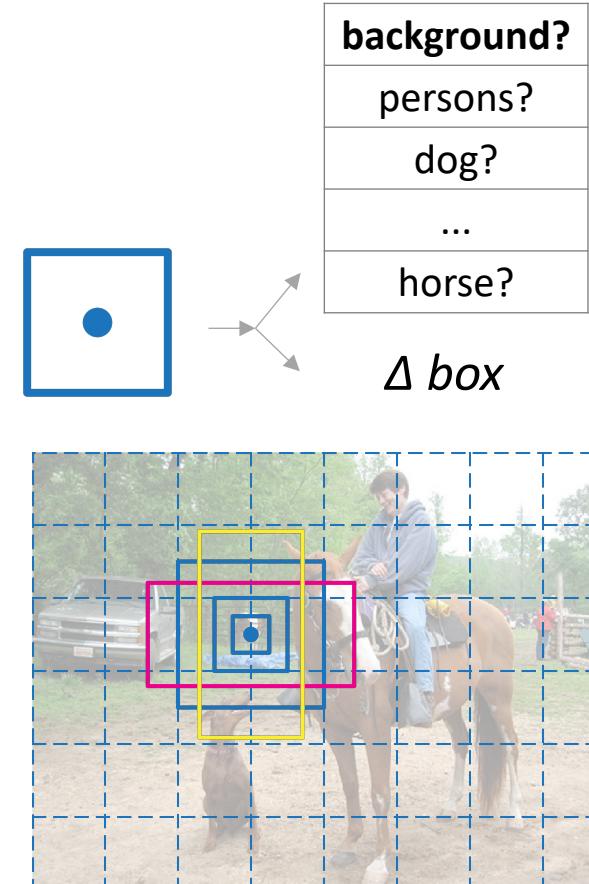


Feature
Pyramid

Single-Stage Object Detection

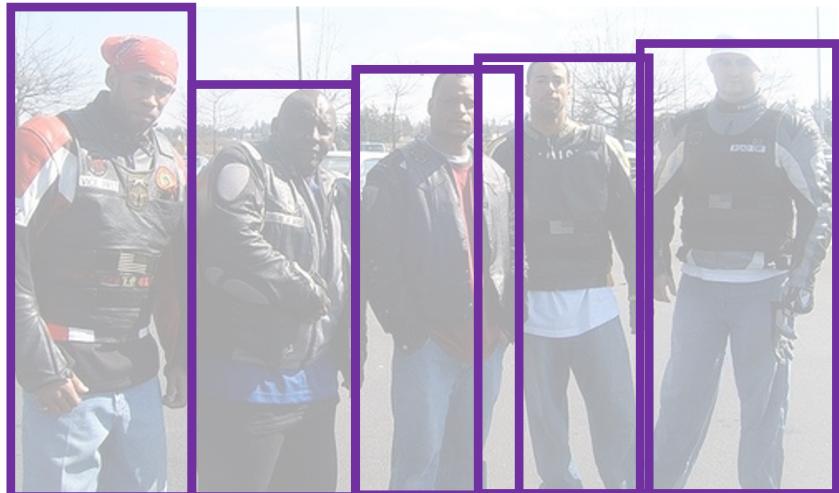
What could still be wrong? Do we need RoI operations?

- “Two-Stage” Detectors: RPN + RoI nets
- “Single-Stage” Detectors
 - YOLO [Redmon+], SSD [Liu+], RetinaNet [Lin+]
- Analogous to “**class-aware** RPN”
- What are still needed?
 - Anchors
 - Feature Pyramid Networks

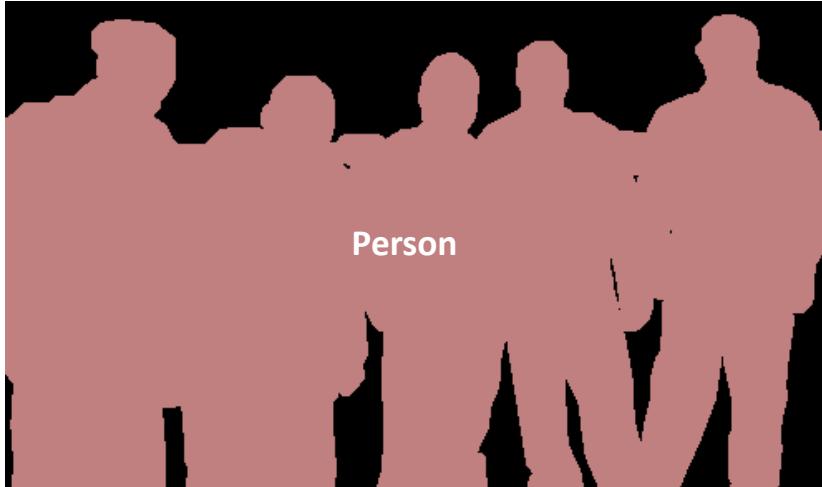


Instance Segmentation

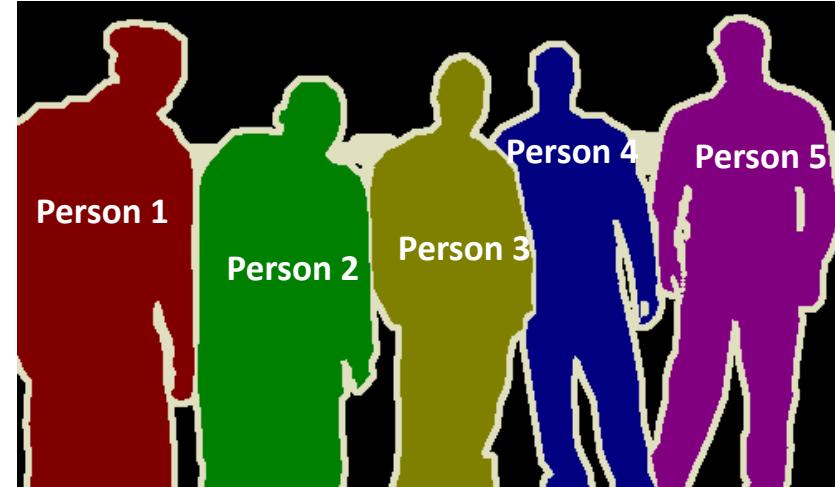
Instance Segmentation: Problem Definition



Object Detection



Semantic Segmentation

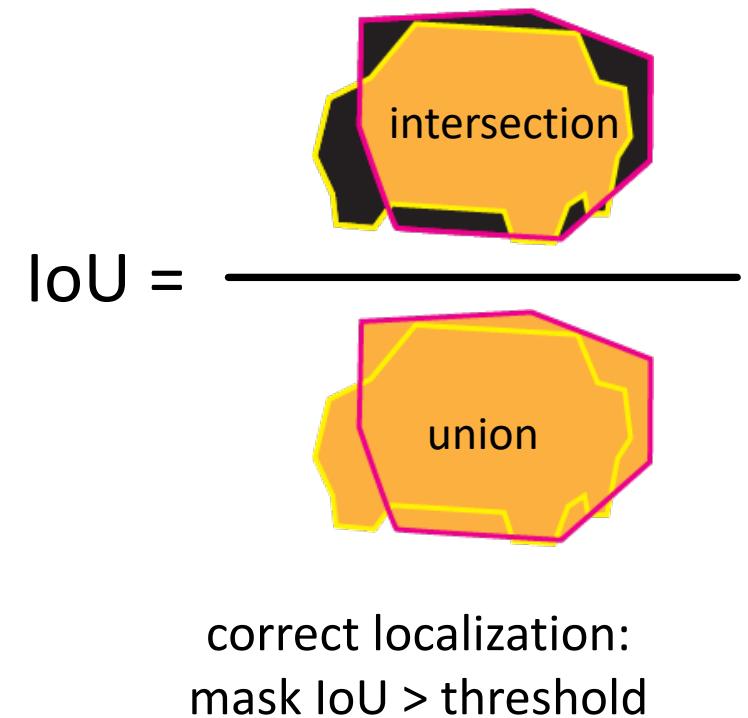
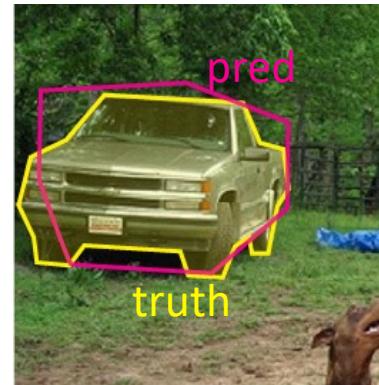
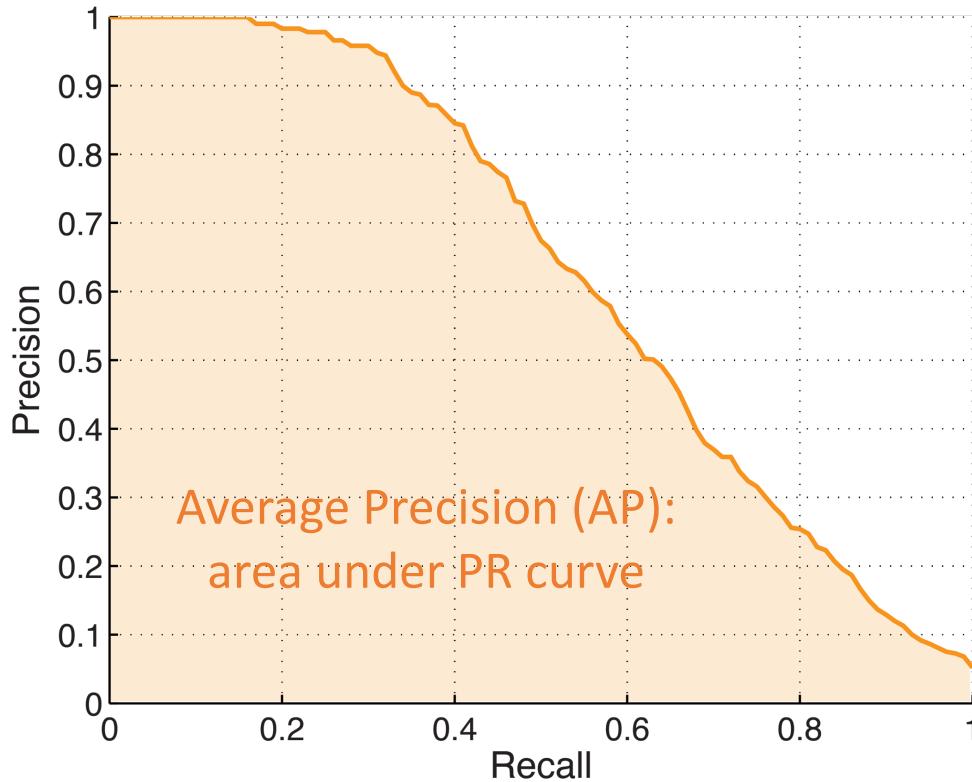


Instance Segmentation

**(Detect and segment
each object)**

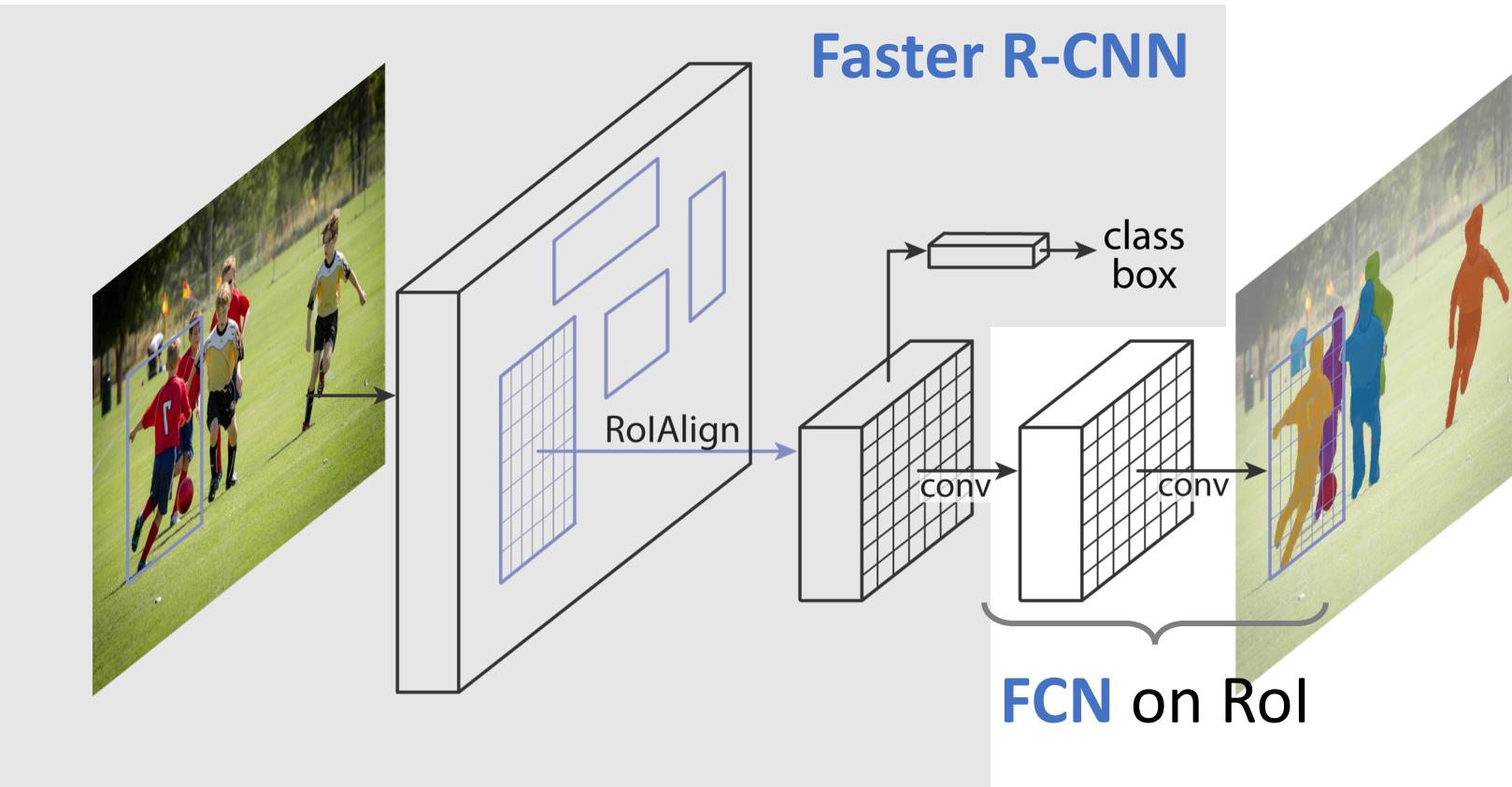
Instance Segmentation: Evaluation Metric

- Similar to Object Detection: Average Precision (AP)
- But localization is determined by mask IoU

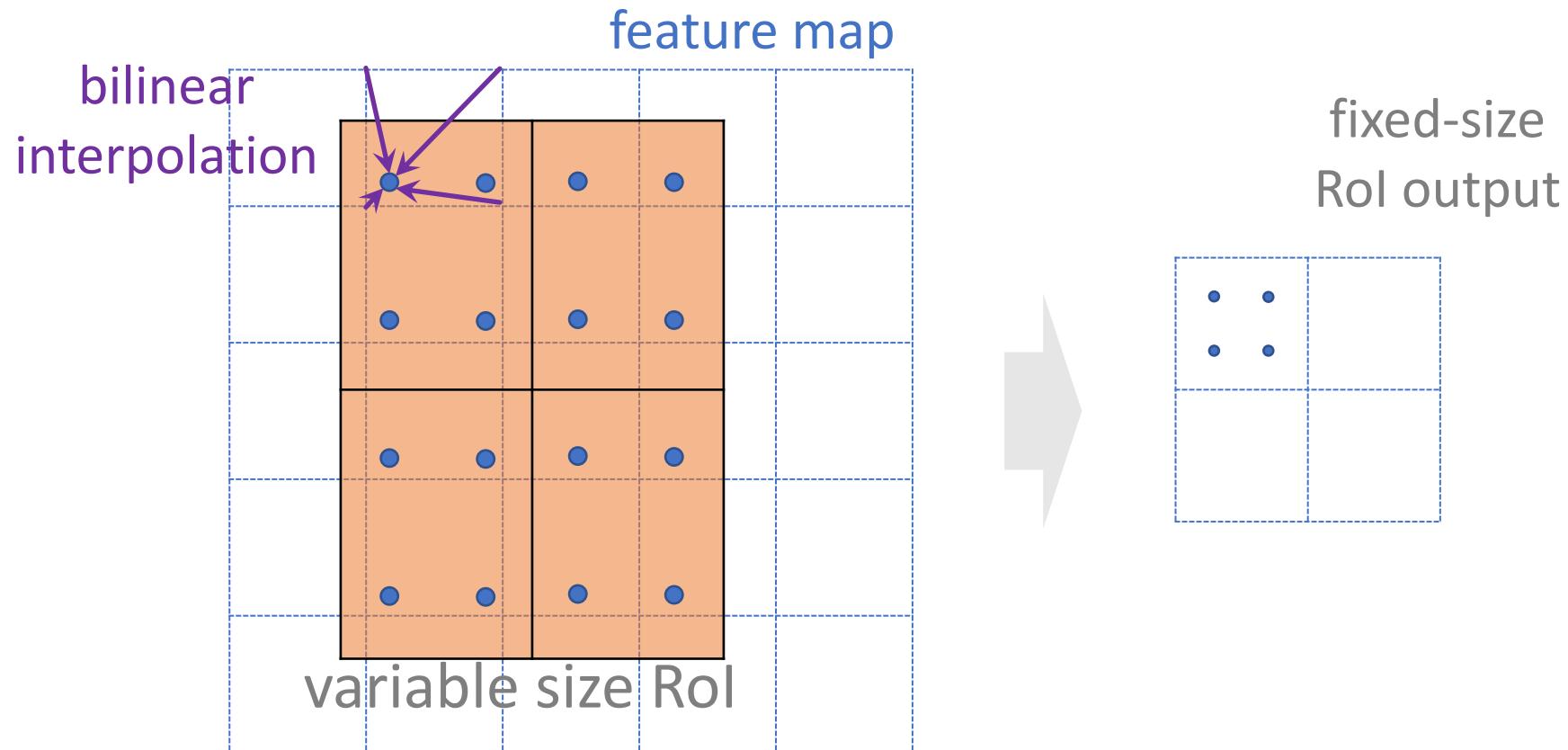


Mask R-CNN

- Difficult problem, simple solution: Faster R-CNN + FCN



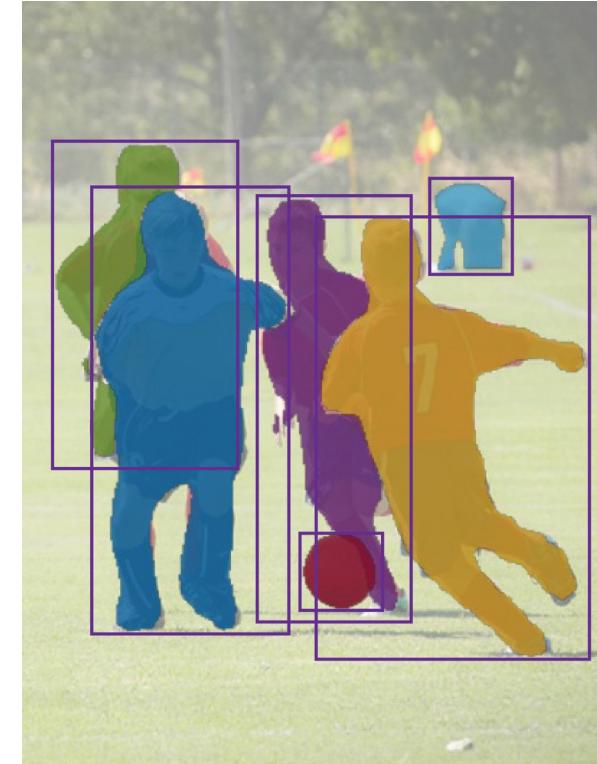
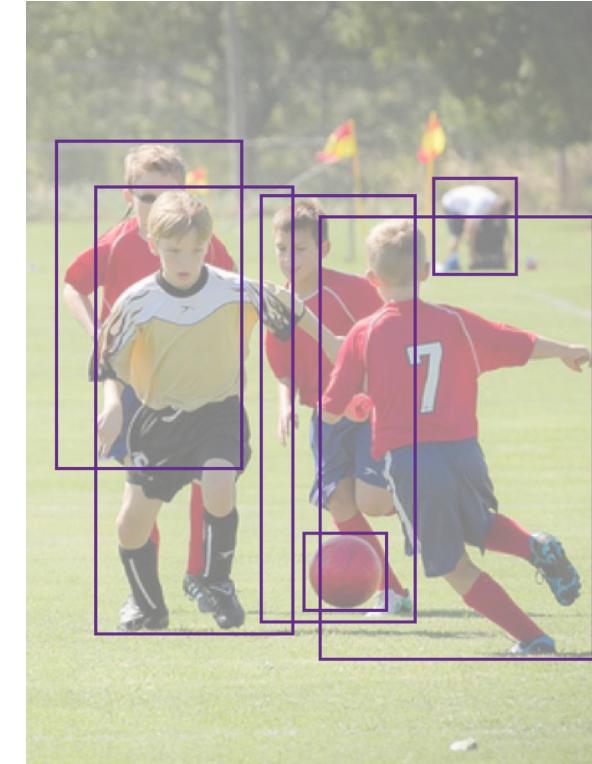
Mask R-CNN: RoI Align



- To apply FCN on each region, we need “aligned” features
- Manipulate feature maps just like images: crop and interpolate

Mask R-CNN

- Difficult problem, simple solution: Faster R-CNN + FCN



Detect individual objects:
Faster R-CNN

For each object, extra an
aligned RoI feature and
apply a tiny **FCN**

Mask R-CNN result in 2017



Mask R-CNN result in 2017



Computer Vision

Algorithms and Applications

Second Edition



Richard Szeliski

Lecture 13: Object Detection and Image Segmentation

- Visual Recognition Tasks
- Methodology: ConvNets, Fully Conv Nets
- Semantic Segmentation
- Object Detection
- Instance Segmentation