

ML Coding Practice  
Lecture 04-1

# Real-time Object Detection

Prof. Jongwon Choi  
Chung-Ang University  
Fall 2022

# Today's Lecture

- **What's Object Detection?**
- **R-CNN / Fast RCNN / Faster RCNN**
- SSD / YOLO

# So far: Image classification



[This image is CC0 public domain](#)

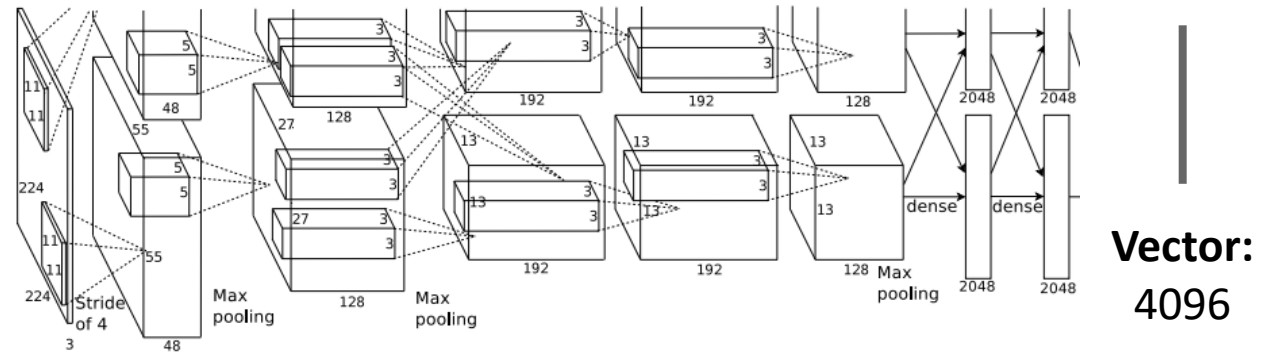
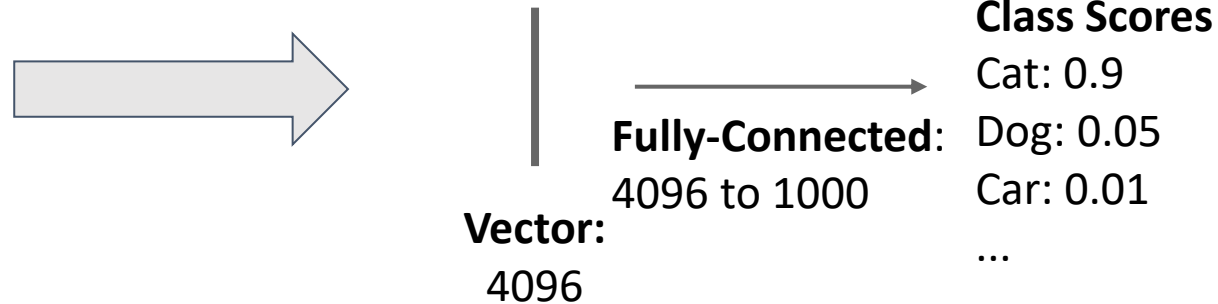


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.



# Other Computer Vision tasks

**Semantic Segmentation**



GRASS, CAT, TREE,  
SKY

No objects, just pixels

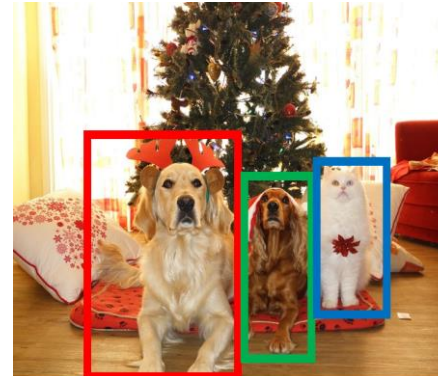
**Classification + Localization**



CAT

Single Object

**Object Detection**



DOG, DOG, CAT

Multiple Object

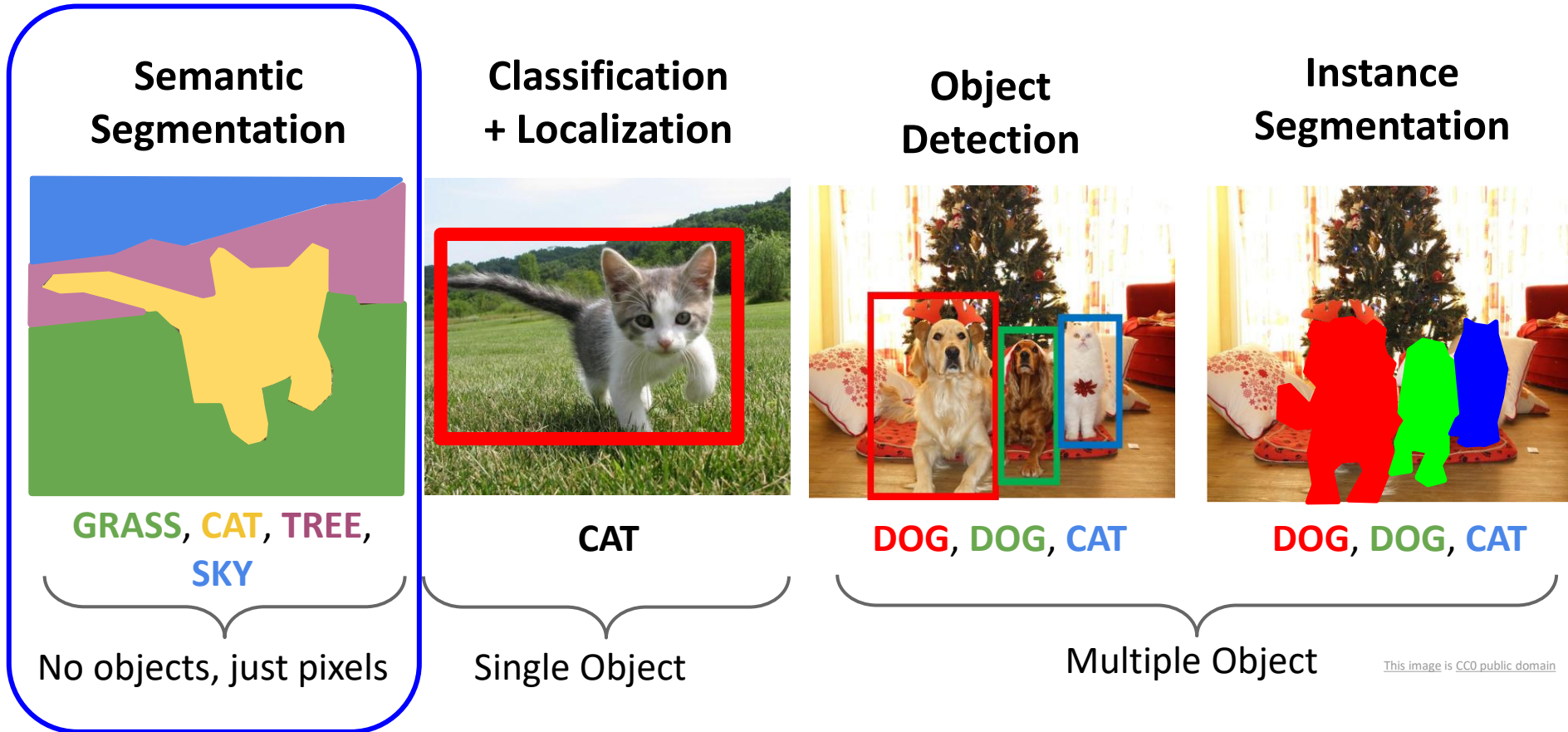
**Instance Segmentation**



DOG, DOG, CAT

This image is CC0 public domain

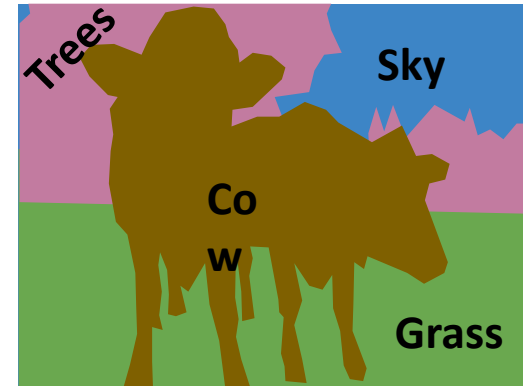
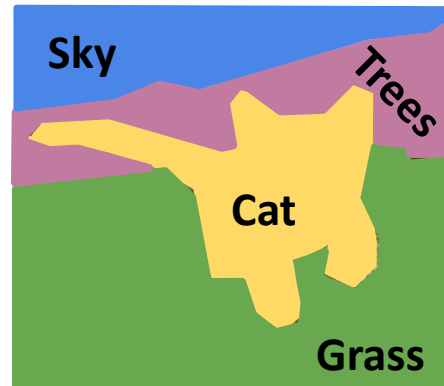
# Other Computer Vision tasks



# Semantic segmentation

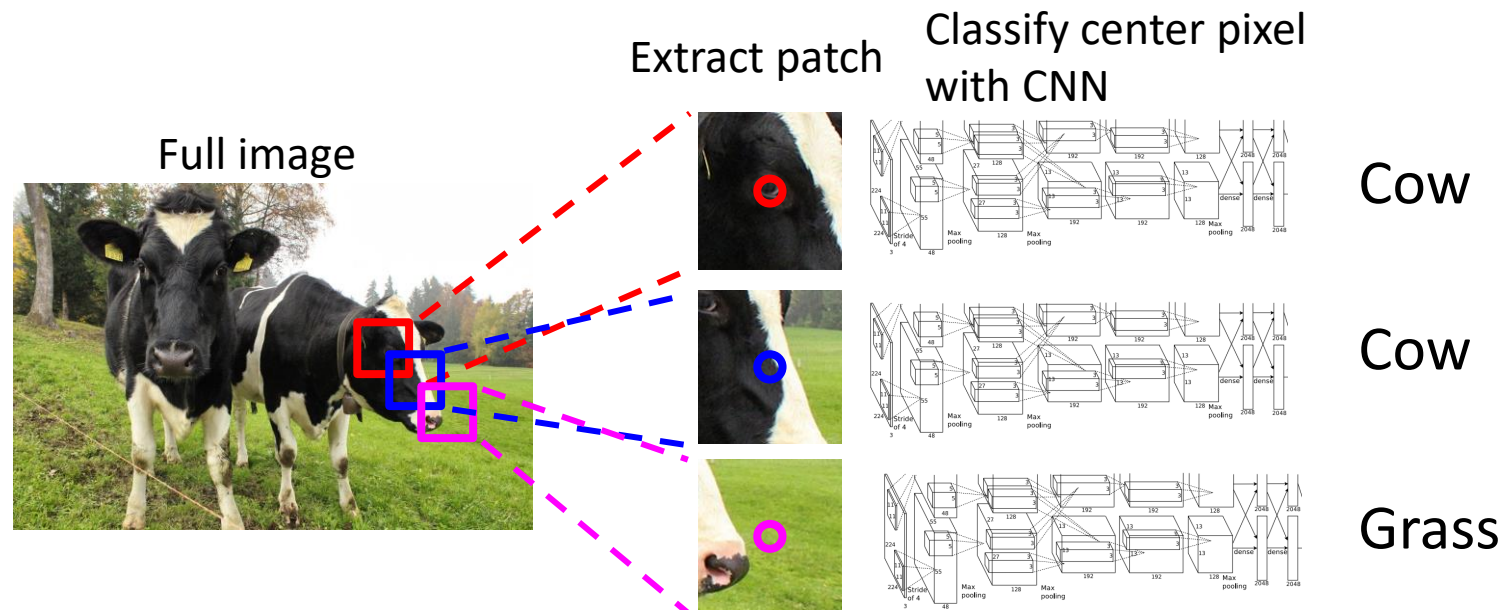
Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



# Semantic segmentation

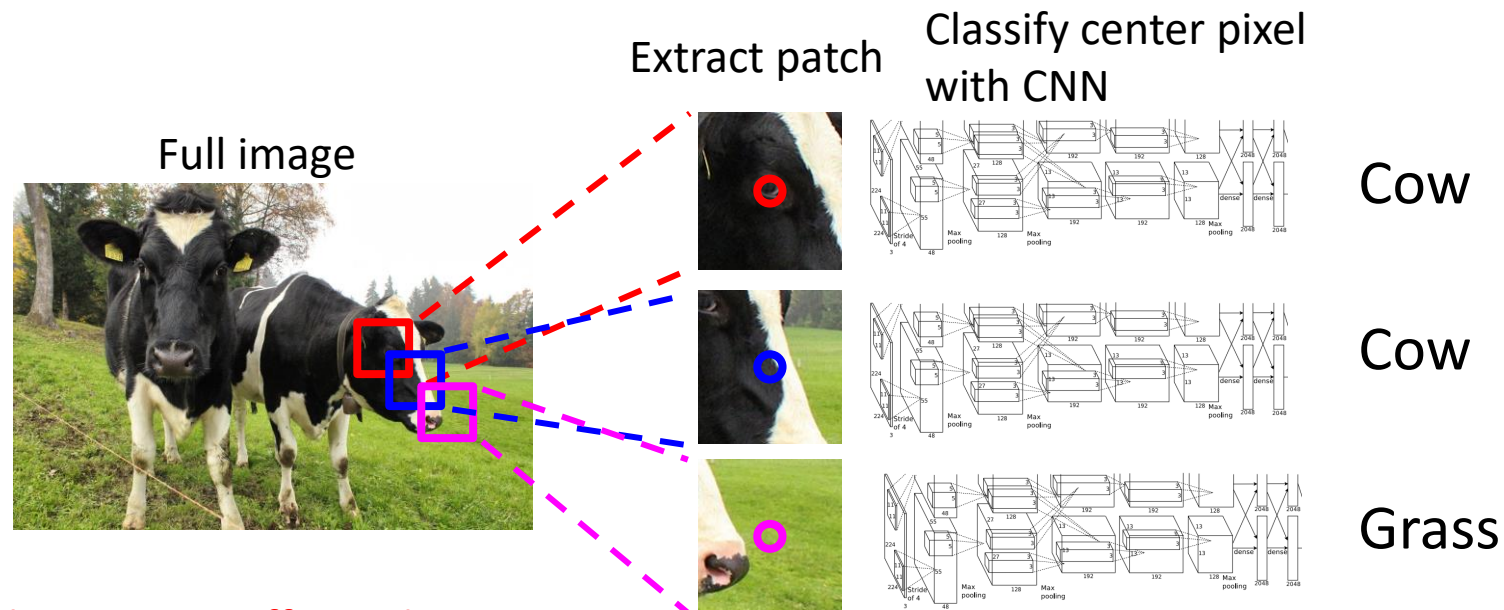
- Sliding window approach





# Semantic segmentation

- Sliding window approach



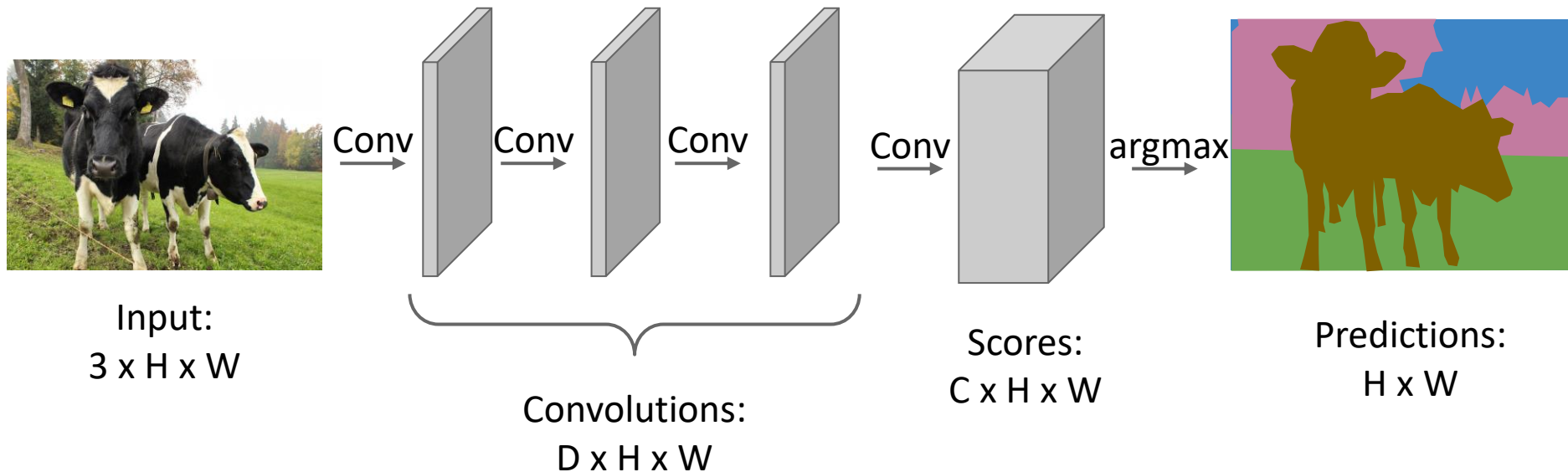
Problem: Very inefficient! Not reusing shared features between overlapping patches



# Semantic segmentation

- Fully convolutional

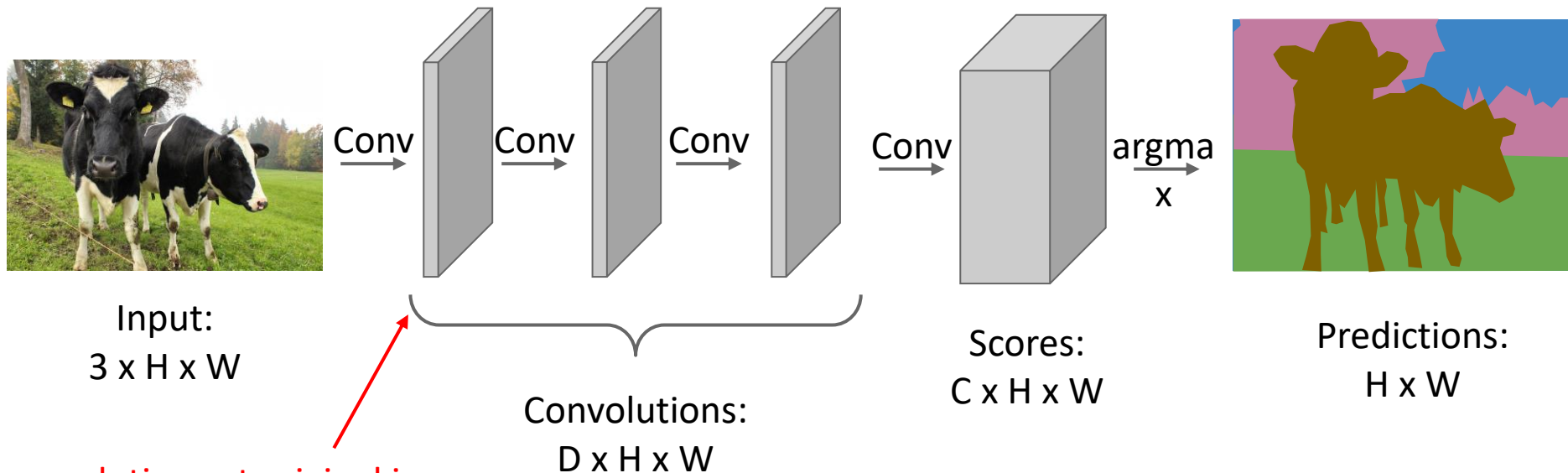
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



# Semantic segmentation

- Fully convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

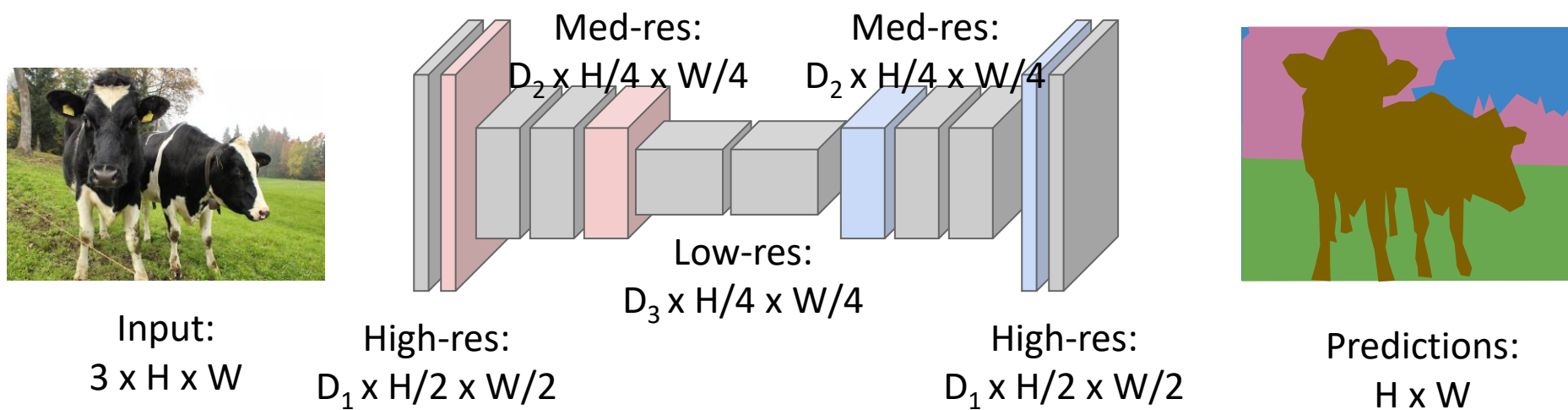


Problem: convolutions at original image resolution will be very expensive ...

# Semantic segmentation

- Fully convolutional

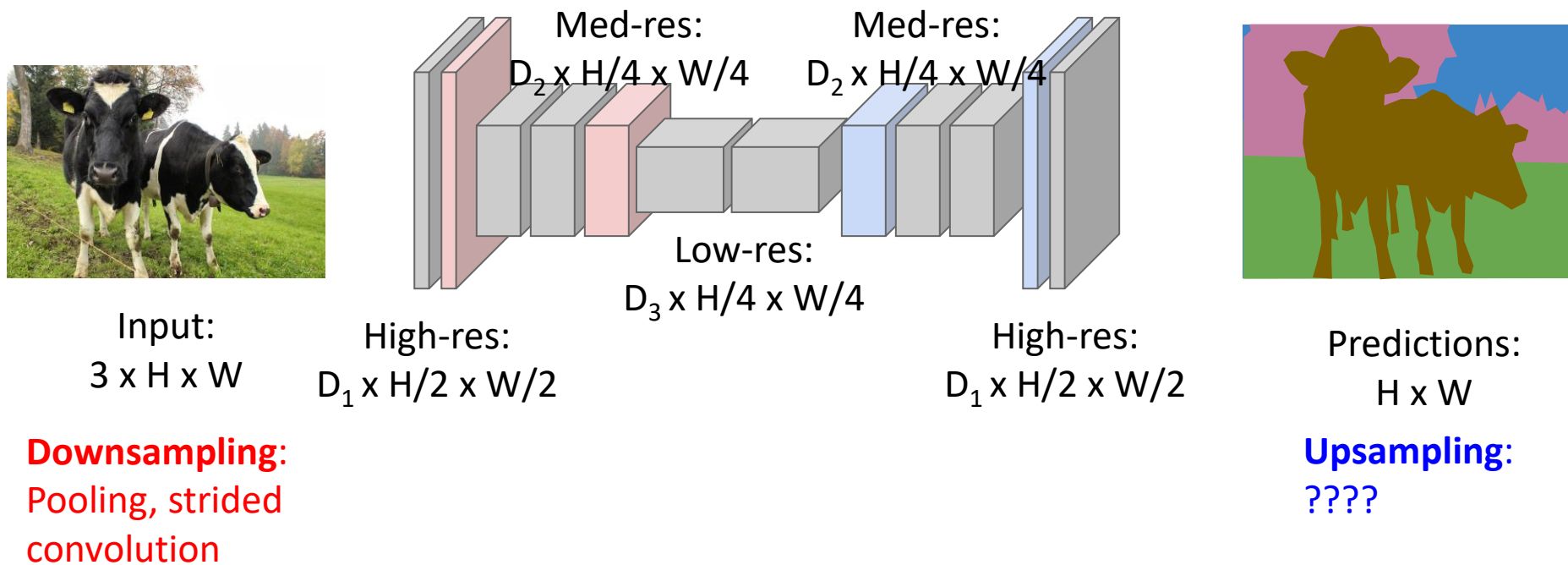
Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



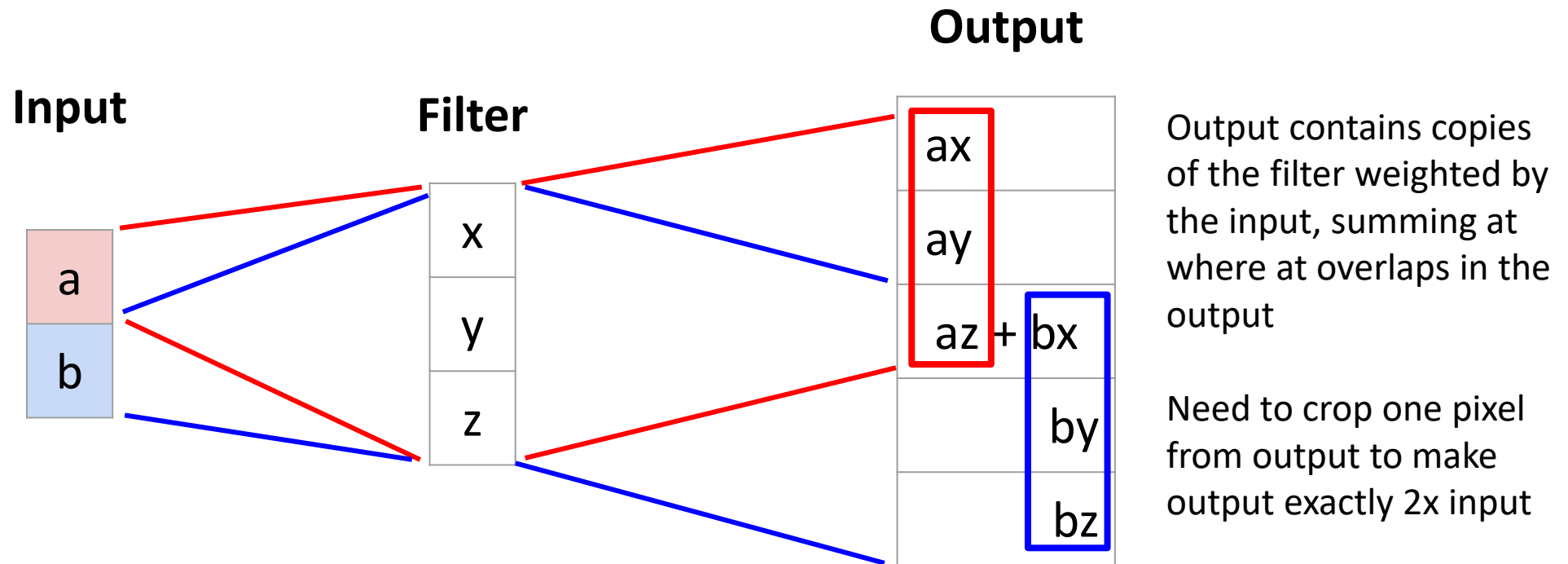
# Semantic segmentation

- Fully convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



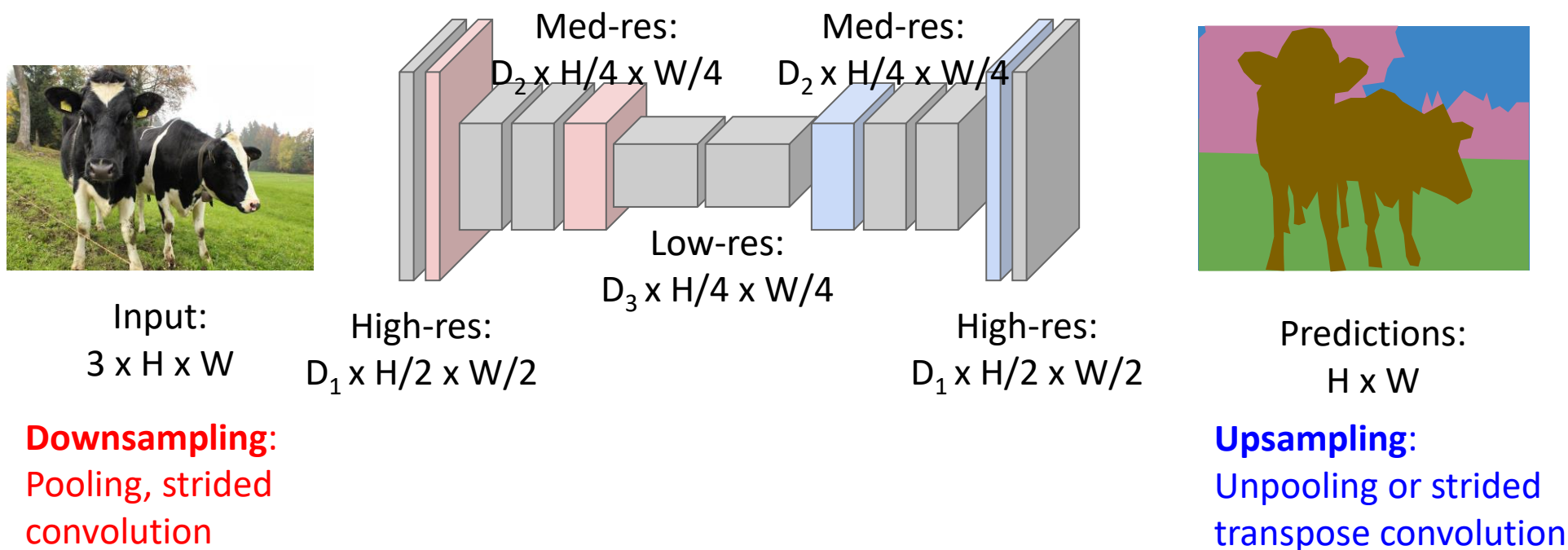
# Recall: Transposed convolution



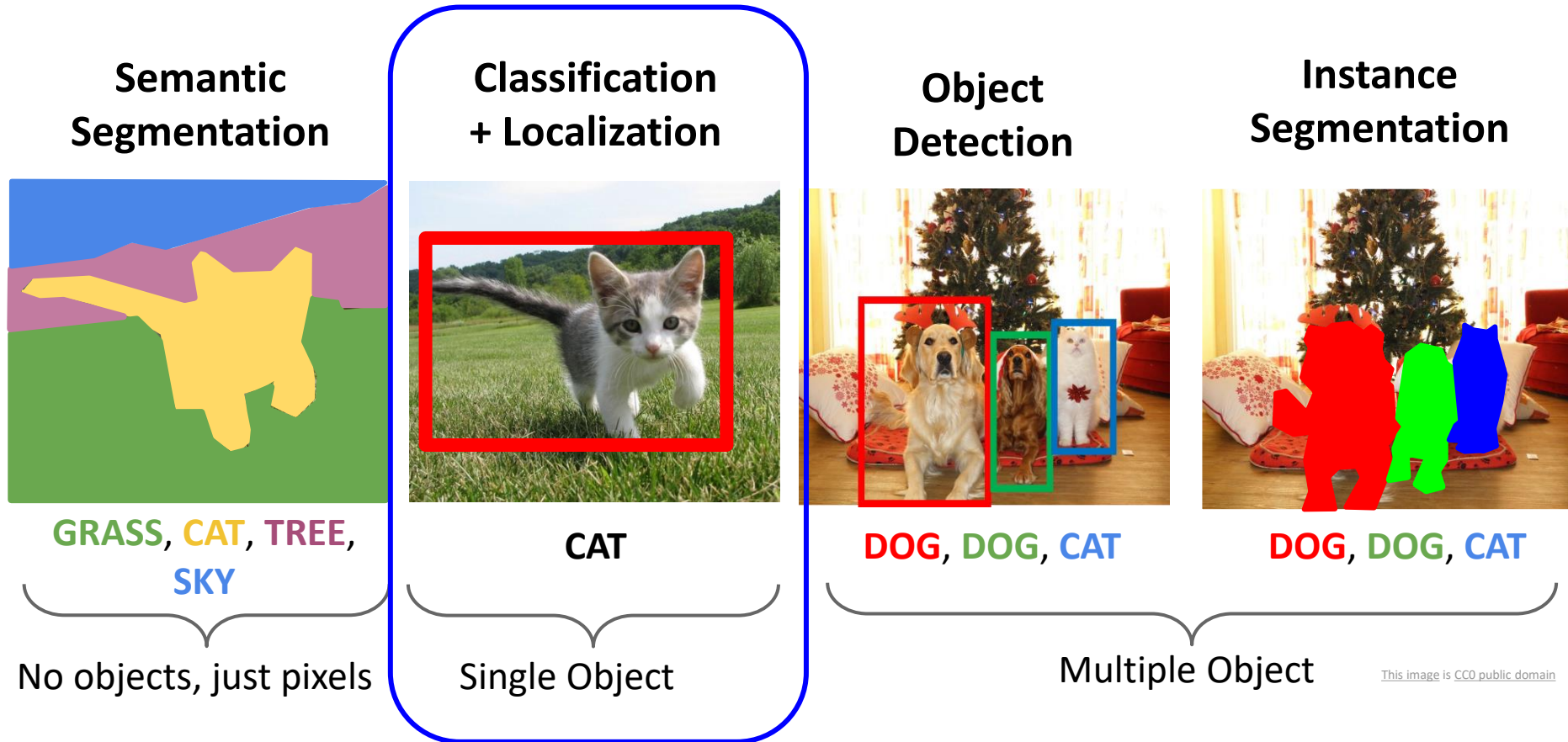
# Semantic segmentation

- Fully convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



# Other Computer Vision tasks

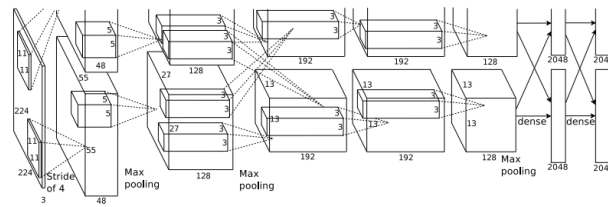




# Classification + Localization



This image is CC0 public domain



**Fully  
Connected:**  
4096 to 1000

## Class Scores

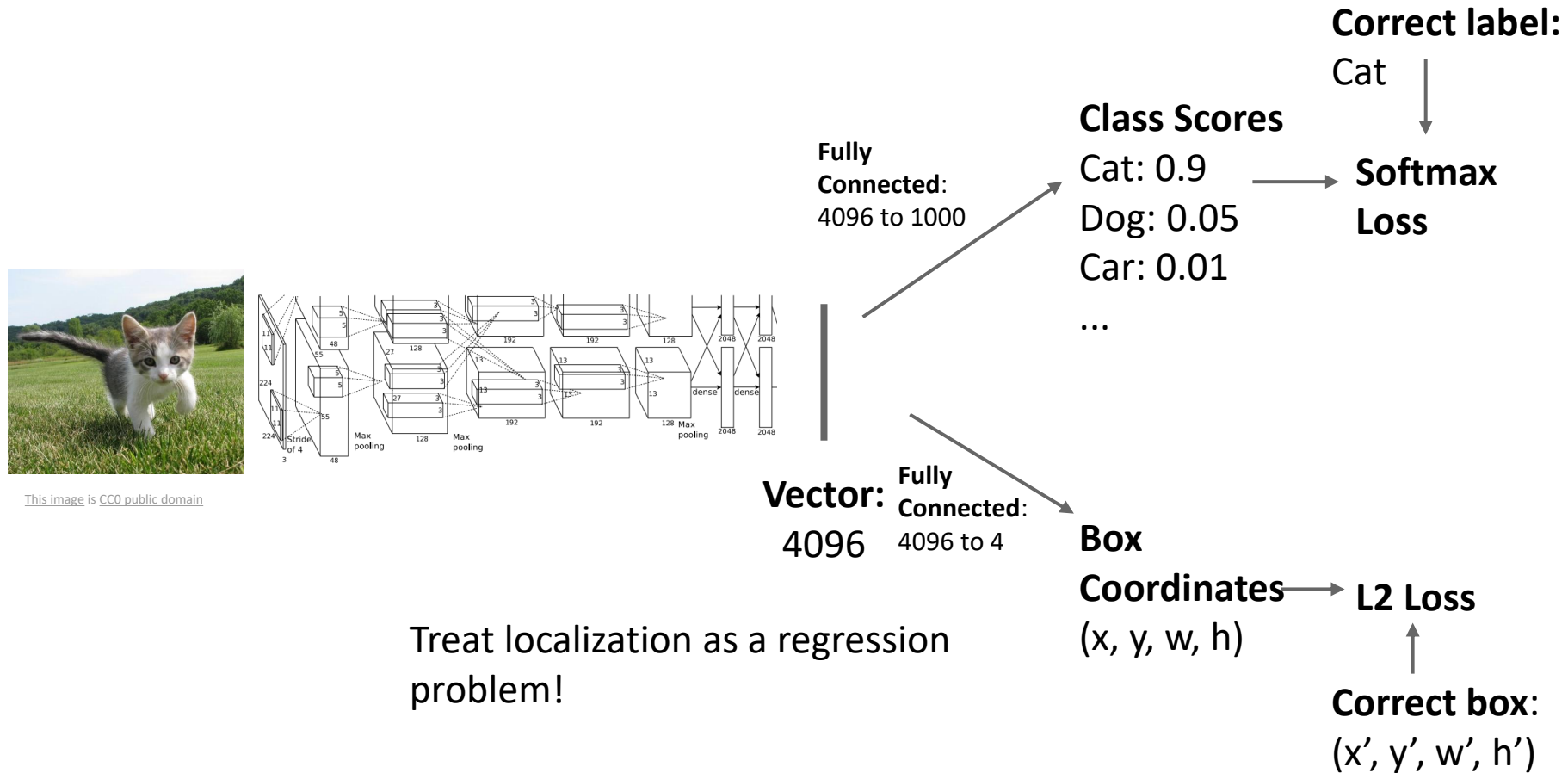
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

**Vector:** 4096  
**Fully  
Connected:**  
4096 to 4

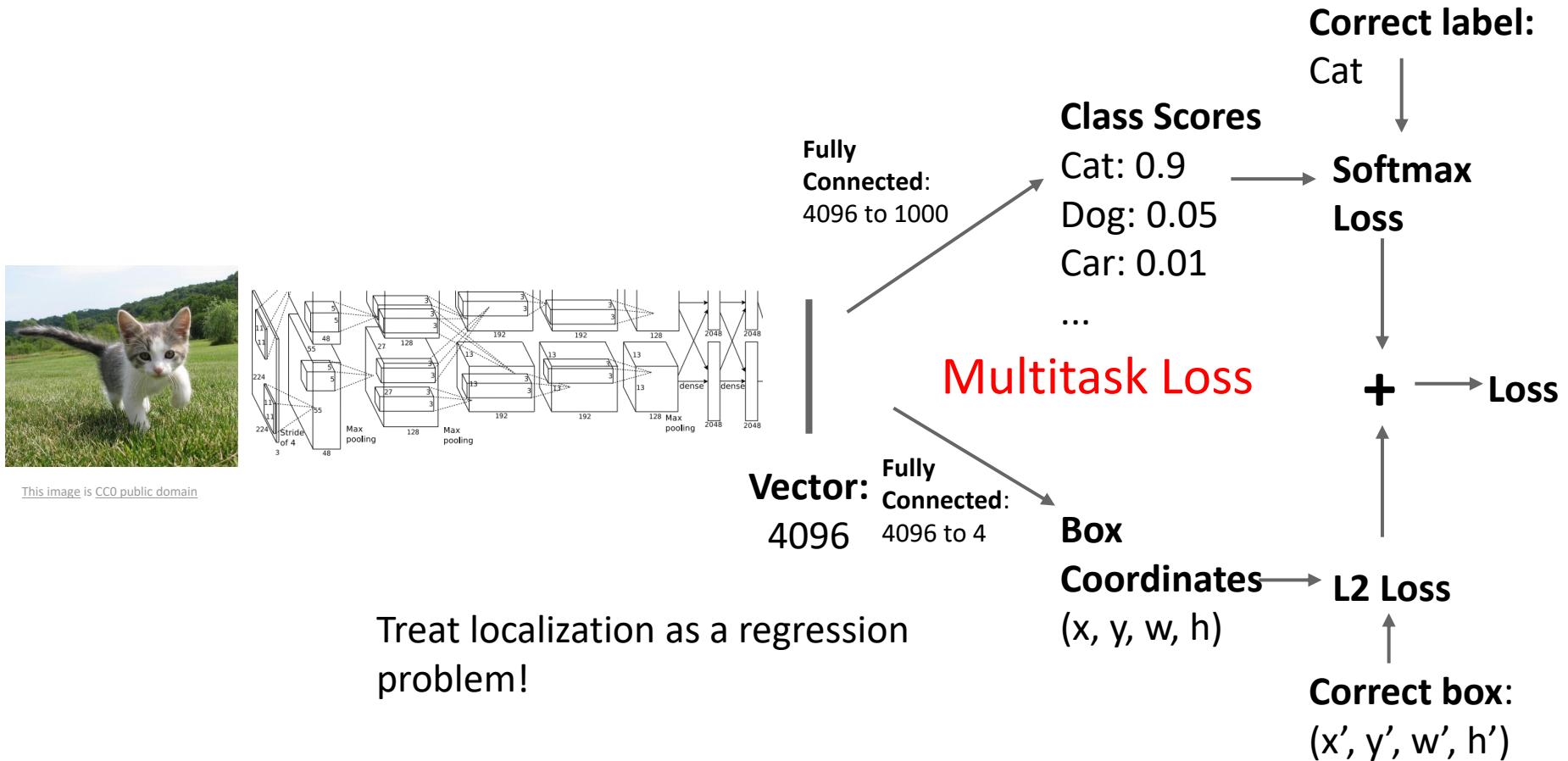
**Box  
Coordinates**  
(x, y, w, h)

Treat localization as a regression problem!

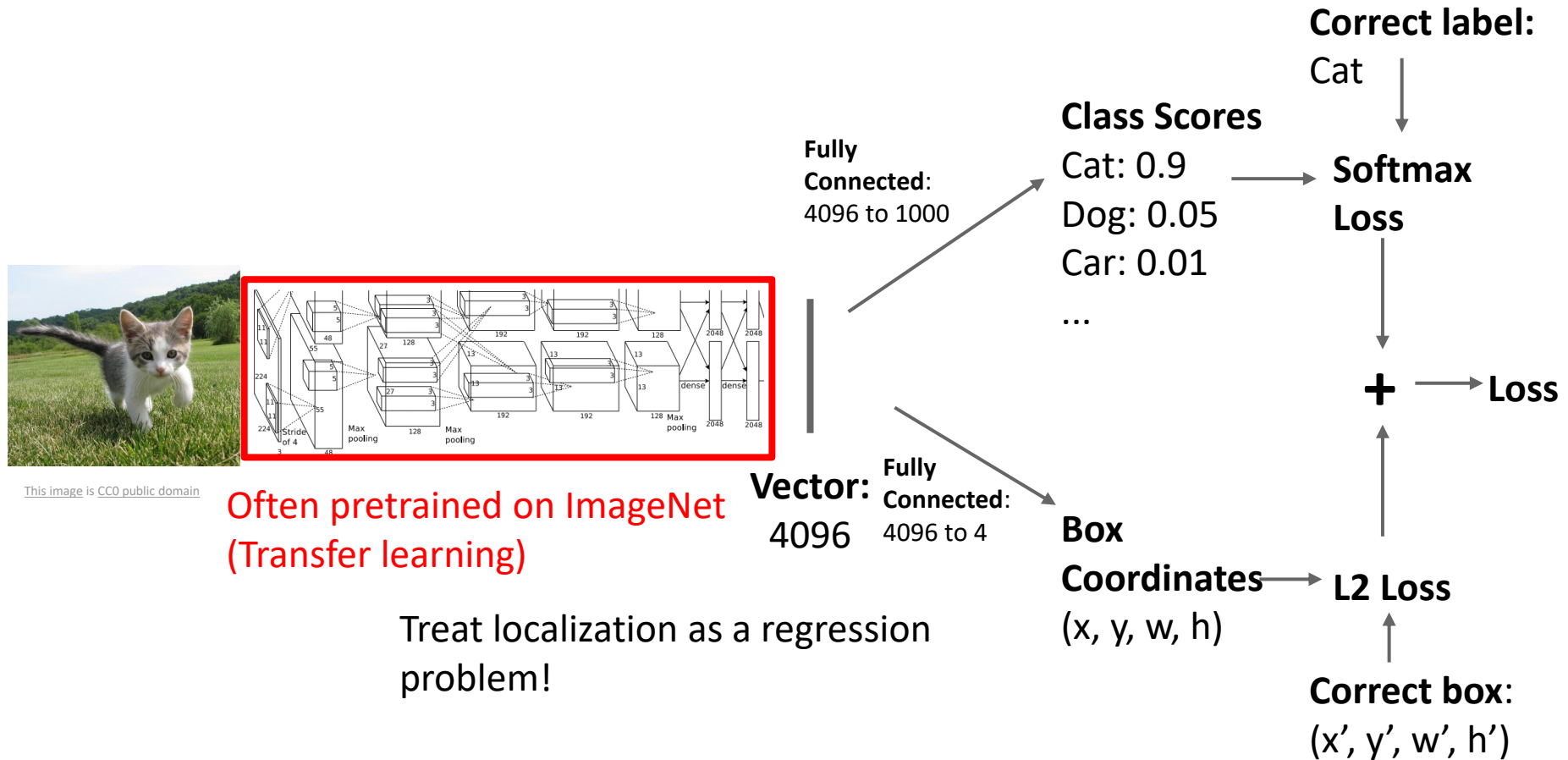
# Classification + Localization



# Classification + Localization



# Classification + Localization



# Aside: Human pose estimation



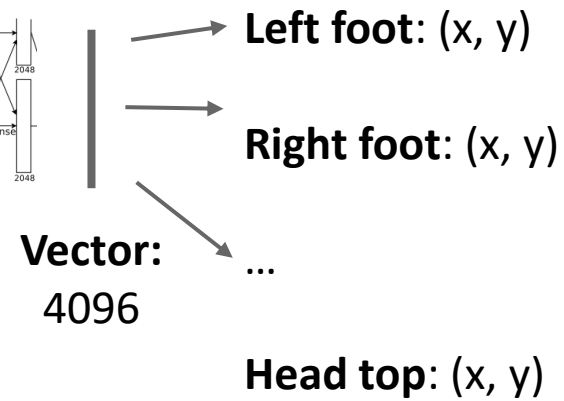
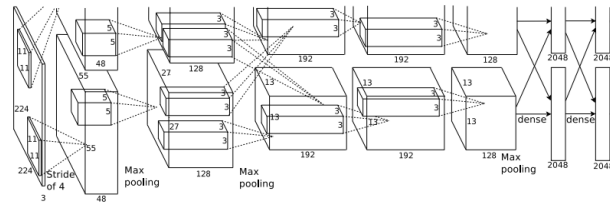
[This image](#) is licensed under [CC-BY 2.0](#).



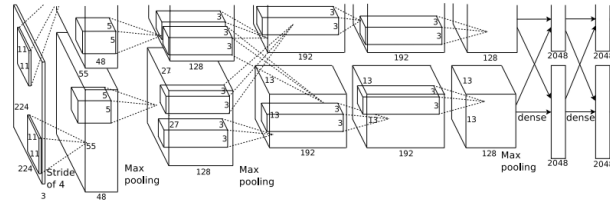
Represent pose as a set of 14 joint positions:

- Left / right foot
- Left / right knee
- Left / right hip
- Left / right shoulder
- Left / right elbow
- Left / right hand
- Neck
- Head top

# Aside: Human pose estimation



# Aside: Human pose estimation



**Vector:**  
4096

**Left foot:**  $(x, y)$

$\rightarrow$  L2 loss

**Right foot:**  $(x, y)$

$\rightarrow$  L2 loss

**Head top:**  $(x, y)$

$\rightarrow$  L2 loss

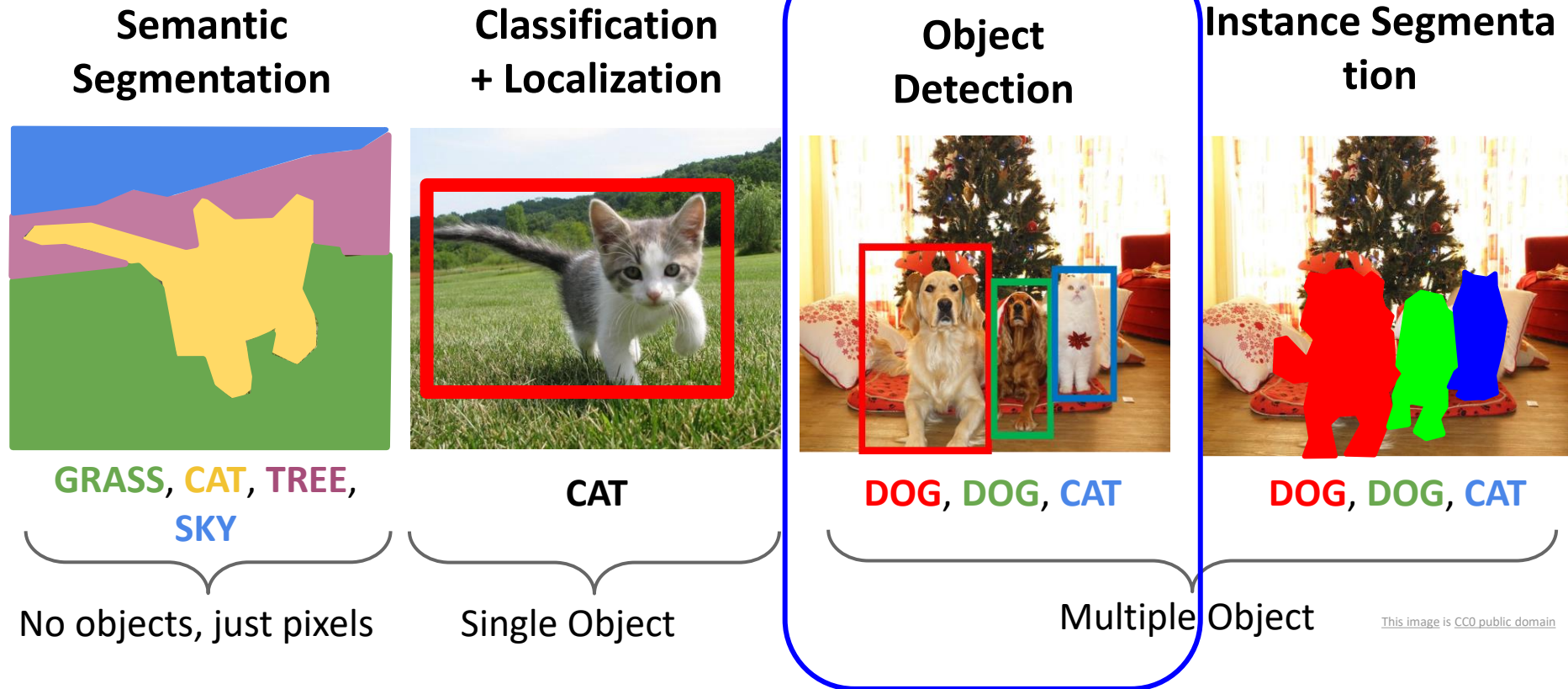
**Correct left foot**  
 $: (x', y')$

**Correct head top:**  
 $(x', y')$

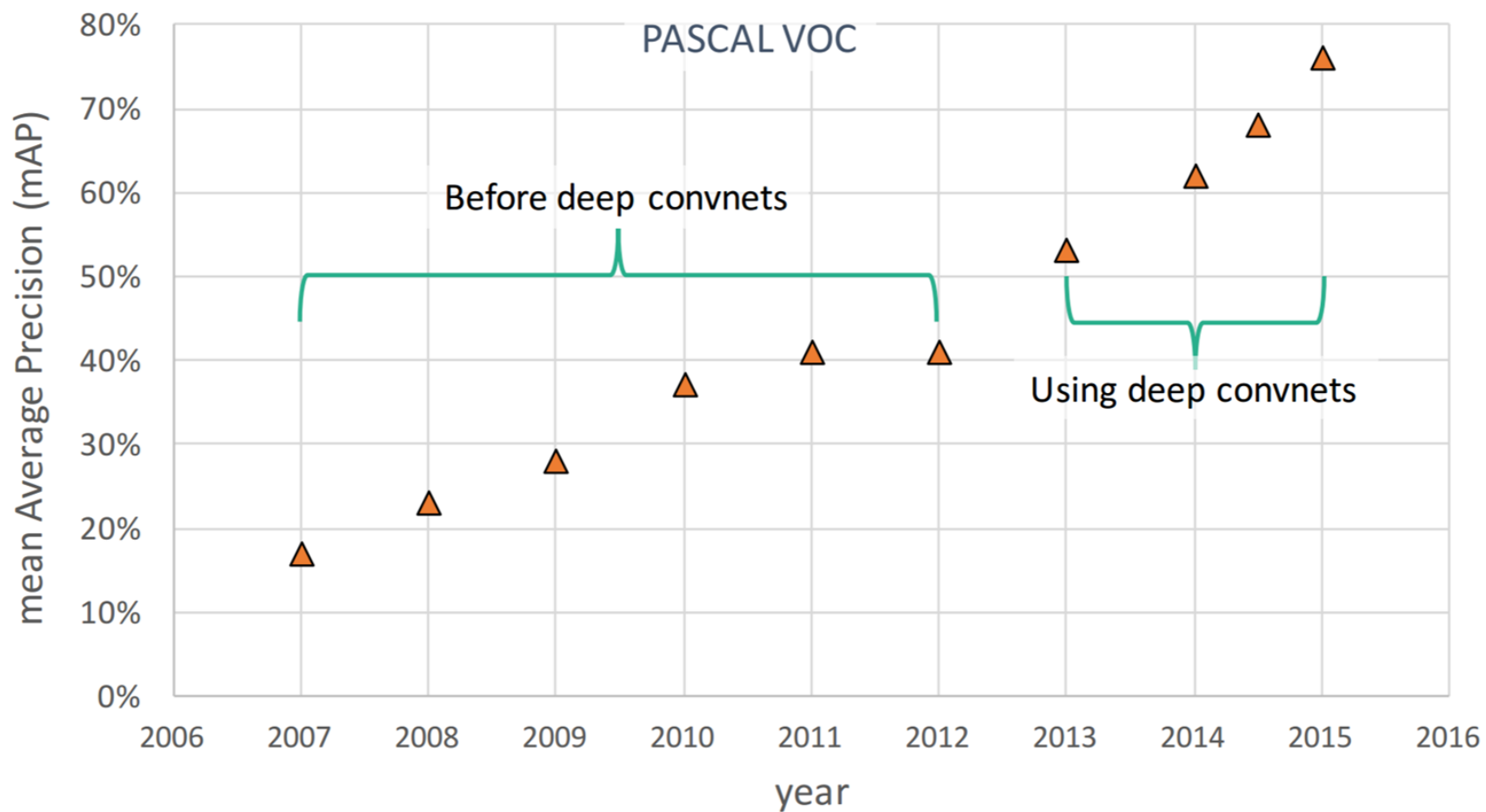
**+**  $\rightarrow$  **Loss**



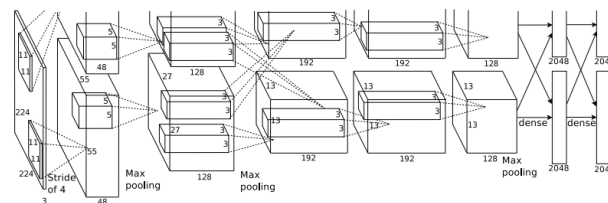
# Other Computer Vision tasks



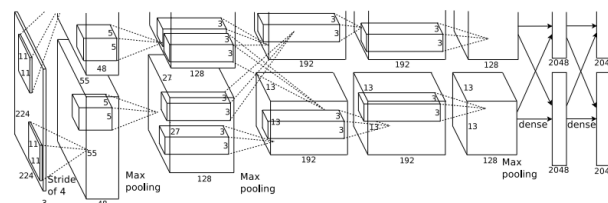
# Object detection: preface



# Object detection: as regression?



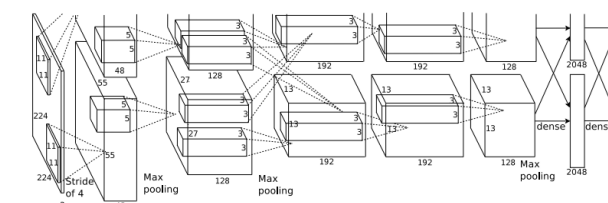
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



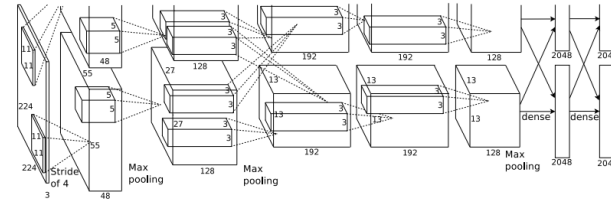
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

....

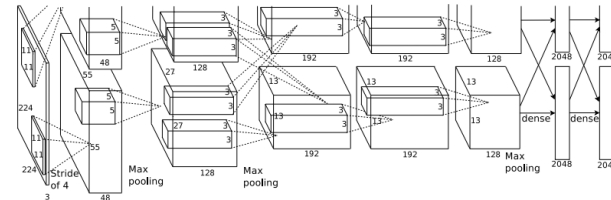
# Object detection: as regression?

Each image needs a different number of outputs!



CAT:  $(x, y, w, h)$

4 numbers

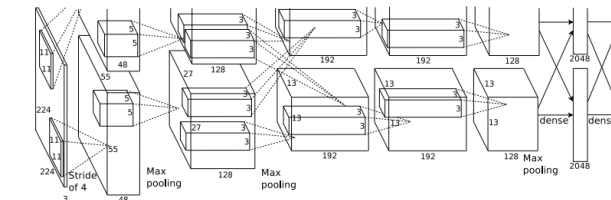


DOG:  $(x, y, w, h)$

DOG:  $(x, y, w, h)$

CAT:  $(x, y, w, h)$

16 numbers



DUCK:  $(x, y, w, h)$

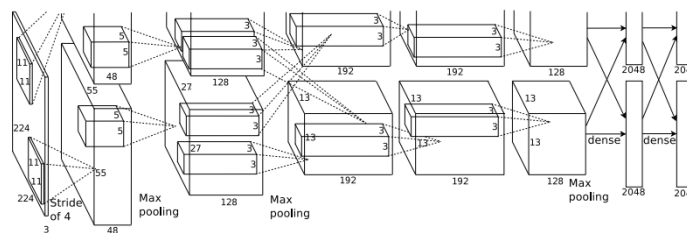
DUCK:  $(x, y, w, h)$

....

Many numbers!

# Object detection: as classification!

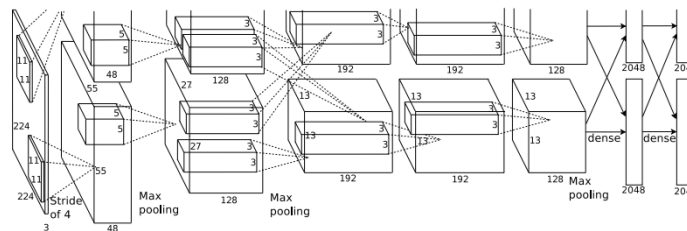
Apply a CNN to many different crops of the image,  
CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object detection: as classification!

Apply a CNN to many different crops of the image,  
CNN classifies each crop as object or background

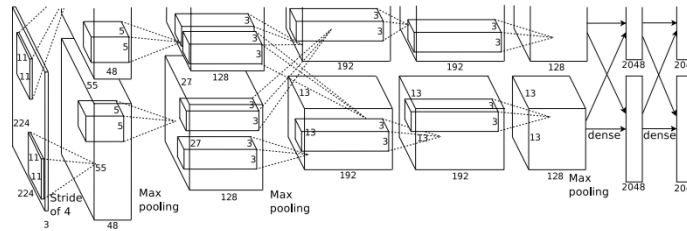


Dog? YES  
Cat? NO  
Background? NO



# Object detection: as classification!

Apply a CNN to many different crops of the image,  
CNN classifies each crop as object or background

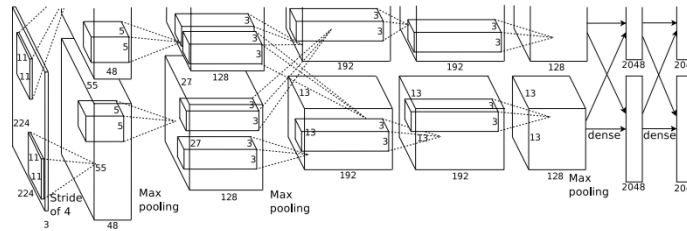


Dog? YES  
Cat? NO  
Background? NO



# Object detection: as classification!

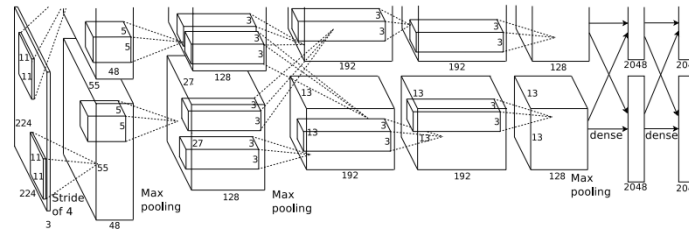
Apply a CNN to many different crops of the image,  
CNN classifies each crop as object or background



Dog? NO  
Cat? YES  
Background? NO

# Object detection: as classification!

Apply a CNN to many different crops of the image,  
CNN classifies each crop as object or background

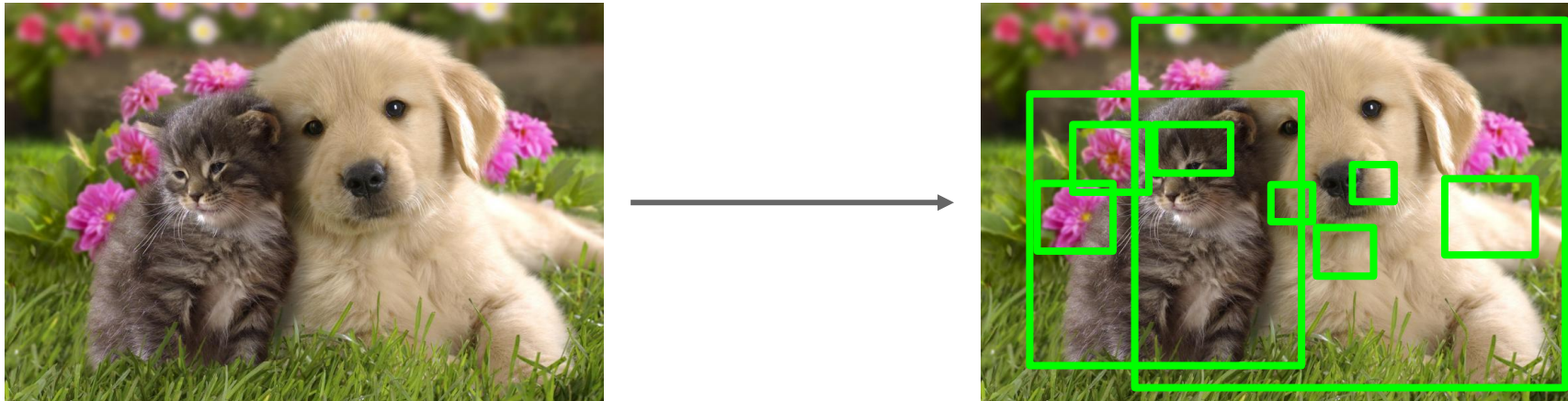


Dog? NO  
Cat? YES  
Background? NO

**Problem: Need to apply CNN to huge number of locations  
and scales, very computationally expensive!**

# Region proposals

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



Alexe et al, “Measuring the objectness of image windows”, TPAMI 2012

Uijlings et al, “Selective Search for Object Recognition”, IJCV 2013

Cheng et al, “BING: Binarized normed gradients for objectness estimation at 300fps”, CVPR 2014

Zitnick and Dollar, “Edge boxes: Locating object proposals from edges”, ECCV 2014

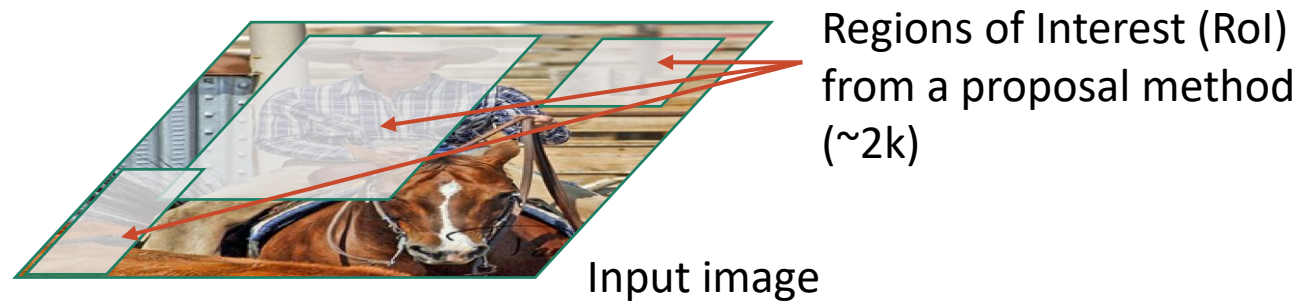
Based on slides for [Stanford cs231n](#) by Li, Jonson, and Young. Modified and reused with permission

# R-CNN

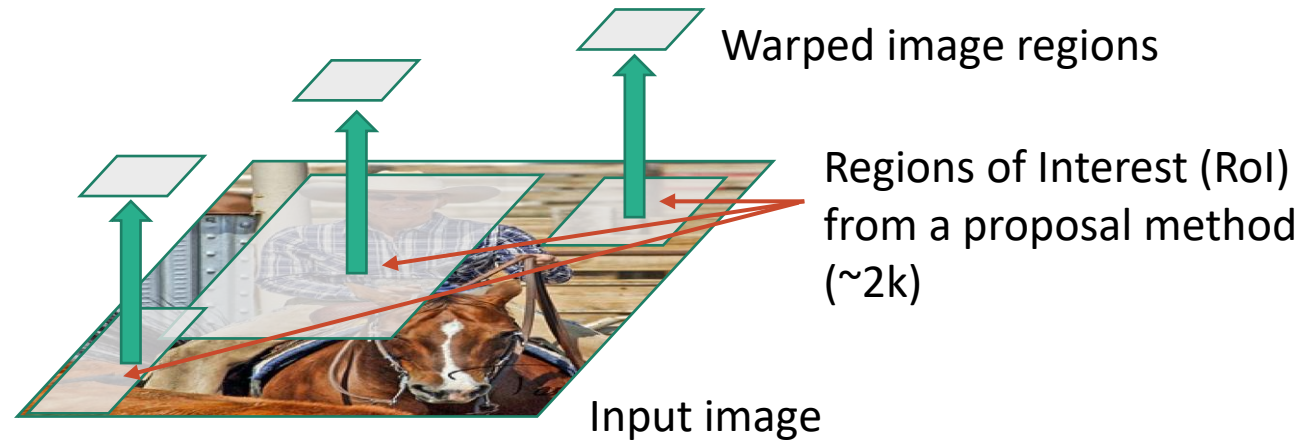


Input image

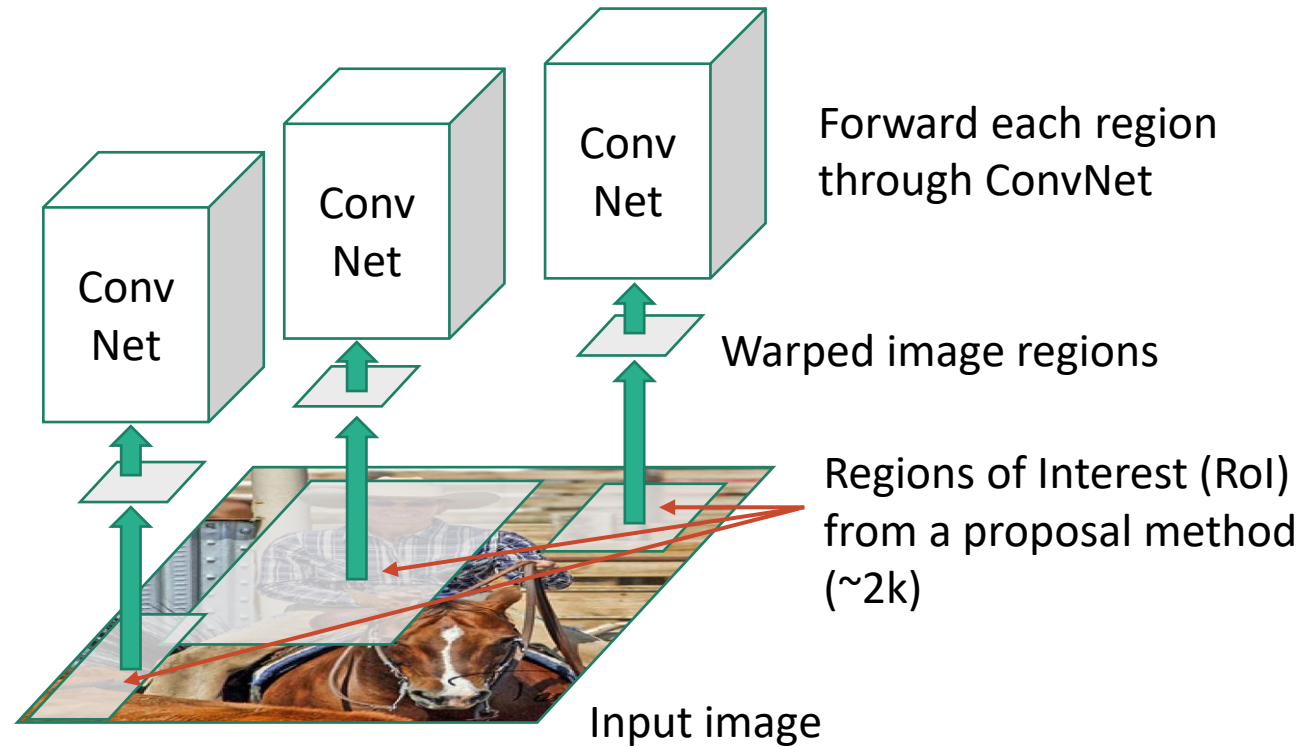
# R-CNN



# R-CNN

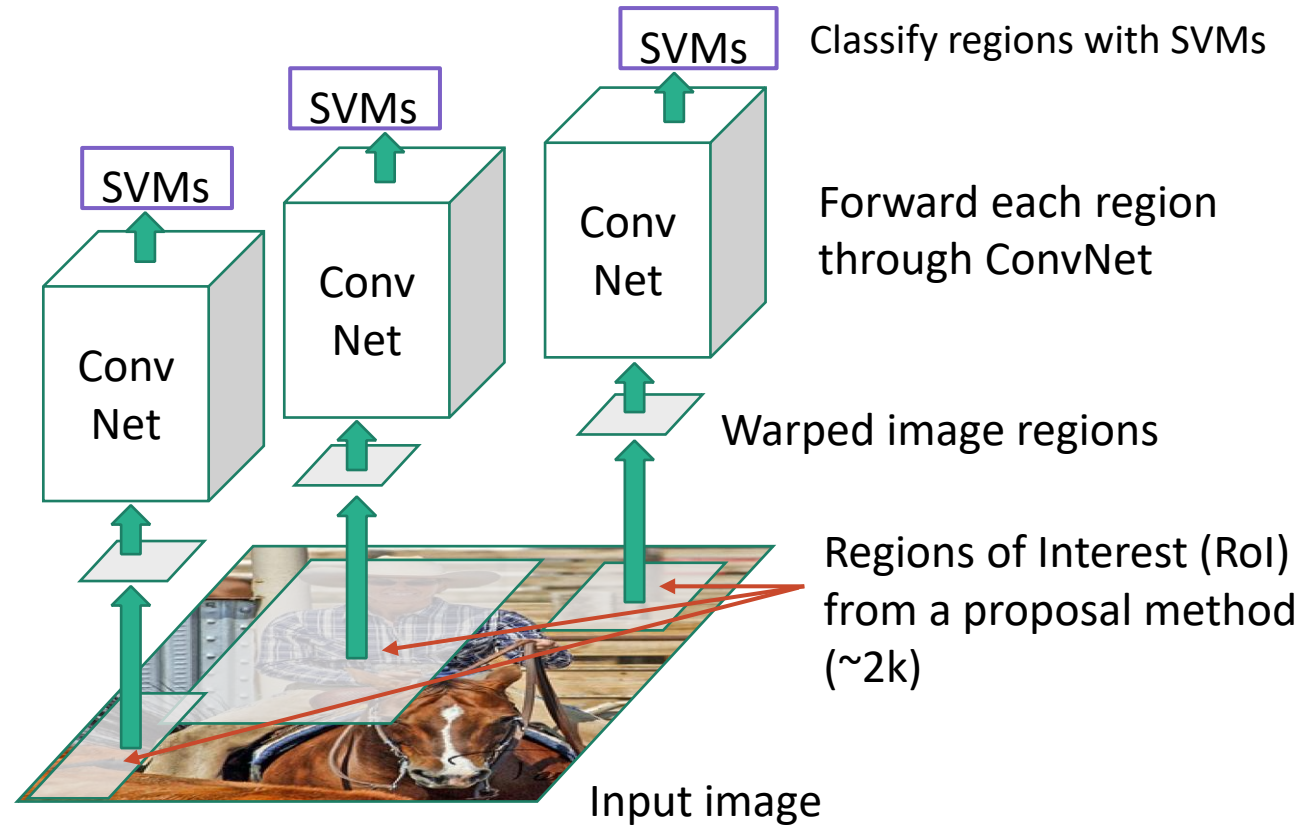


# R-CNN

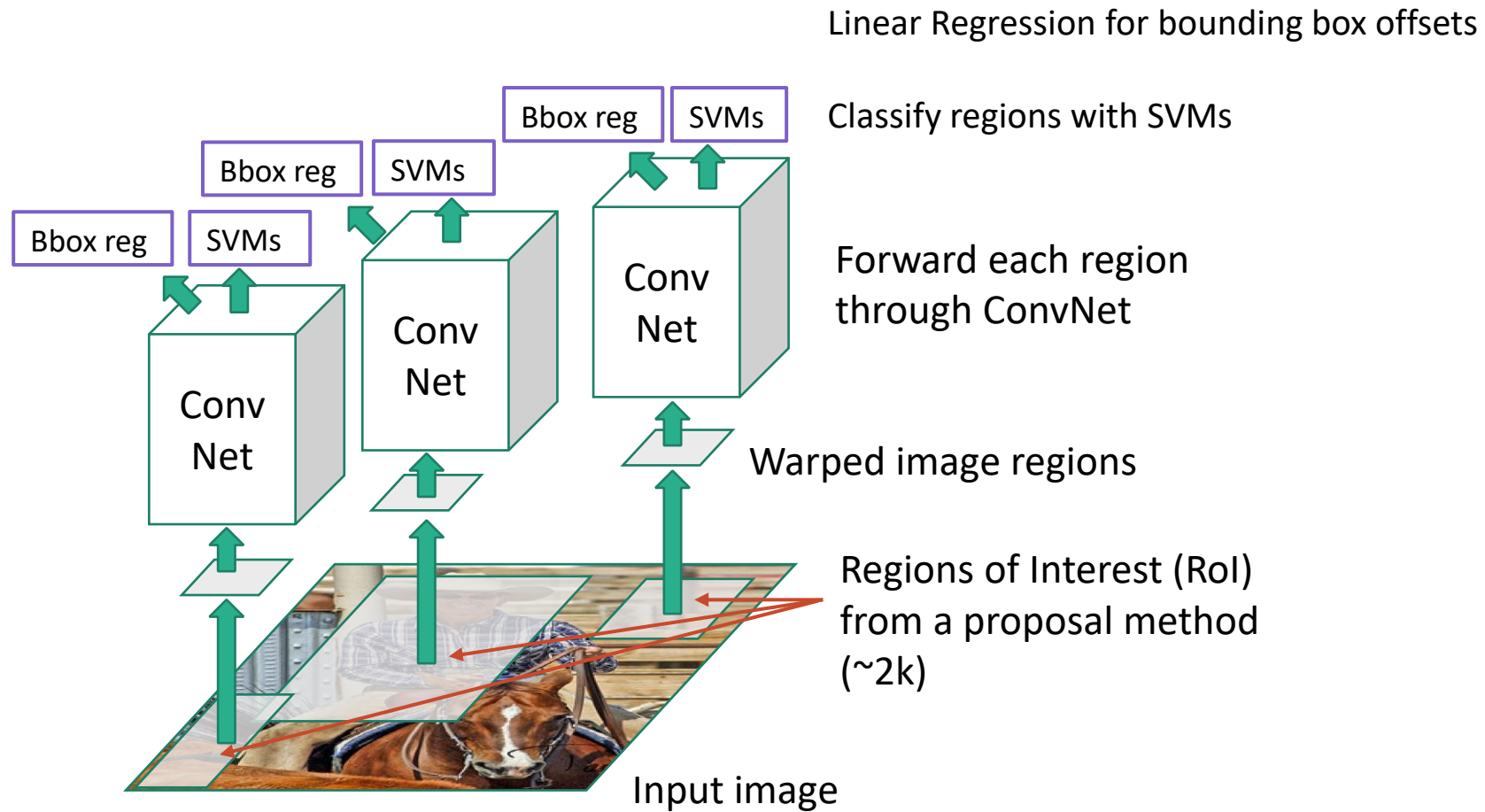




# R-CNN



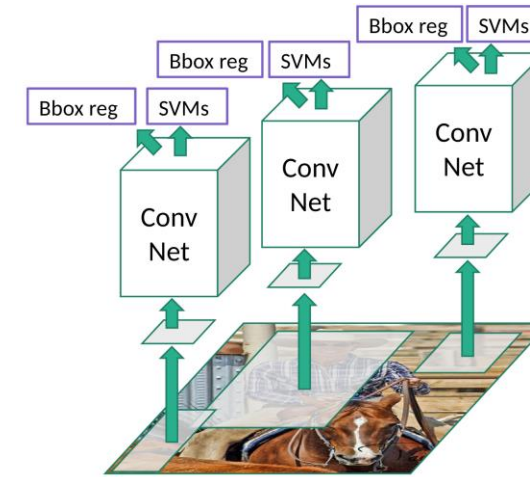
# R-CNN



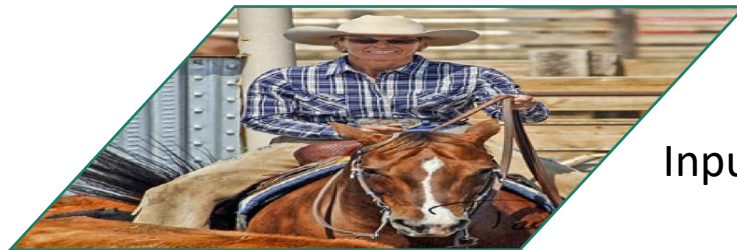
# R-CNN

## ● Problems

- Ad hoc training objectives
  - Fine-tune network with softmax classifier
    - log loss
  - Train post-hoc linear SVMs
    - hinge loss
  - Train post-hoc bounding-box regressions
    - least squares
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
  - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
  - Fixed by SPP-net [He et al. ECCV14]

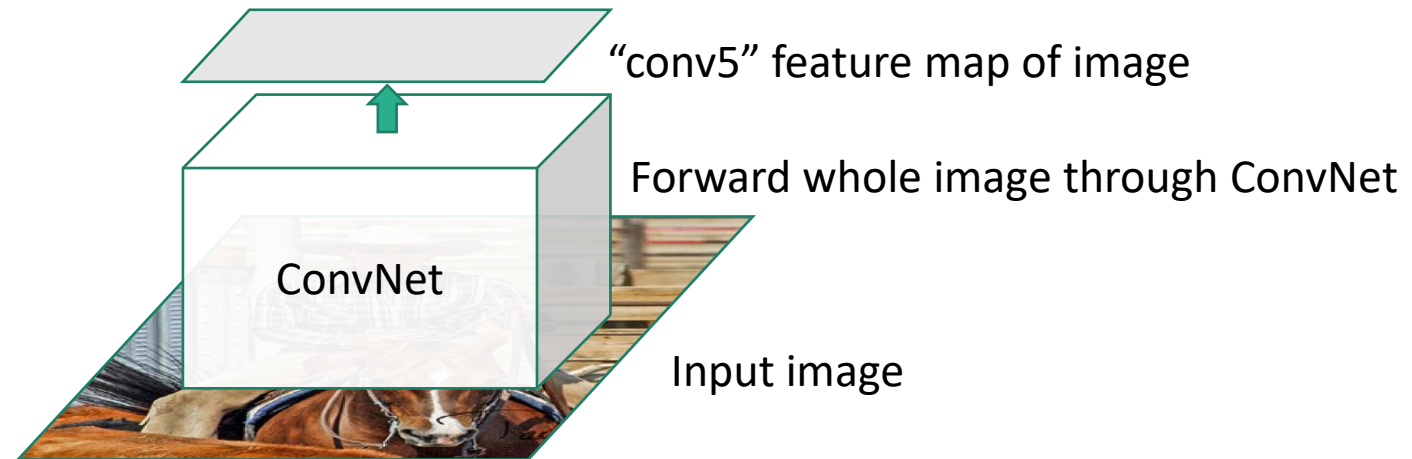


# Fast R-CNN

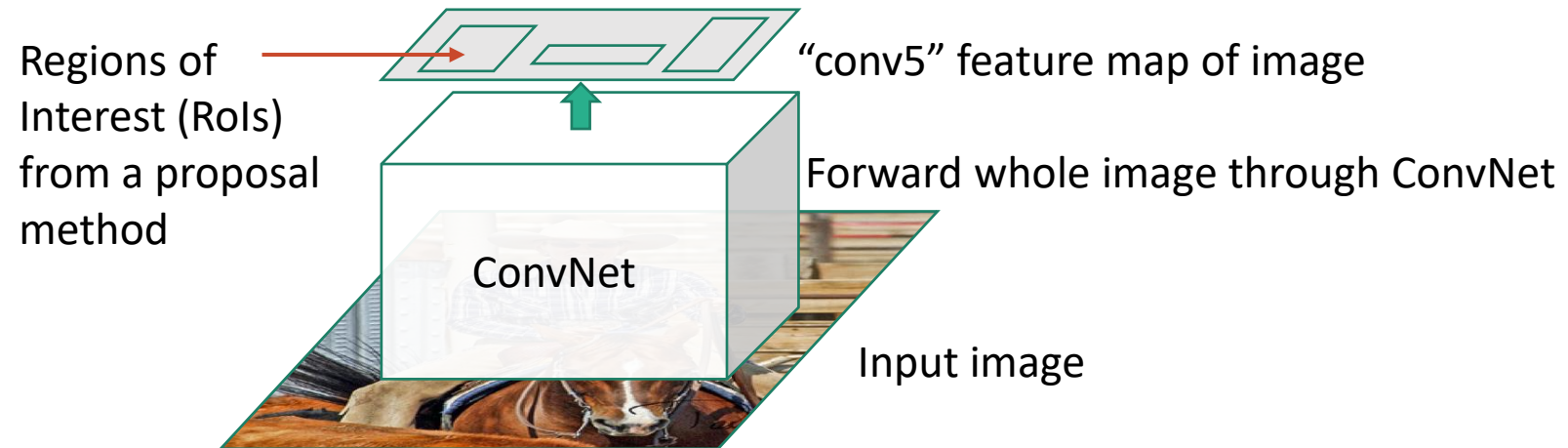


Input image

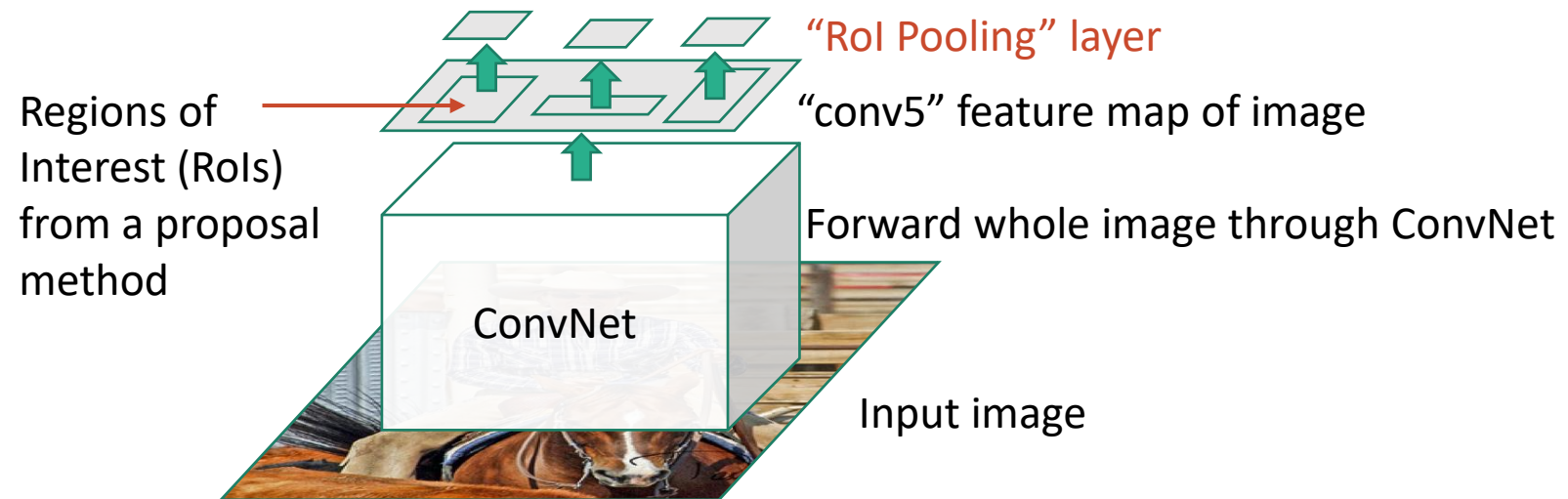
# Fast R-CNN



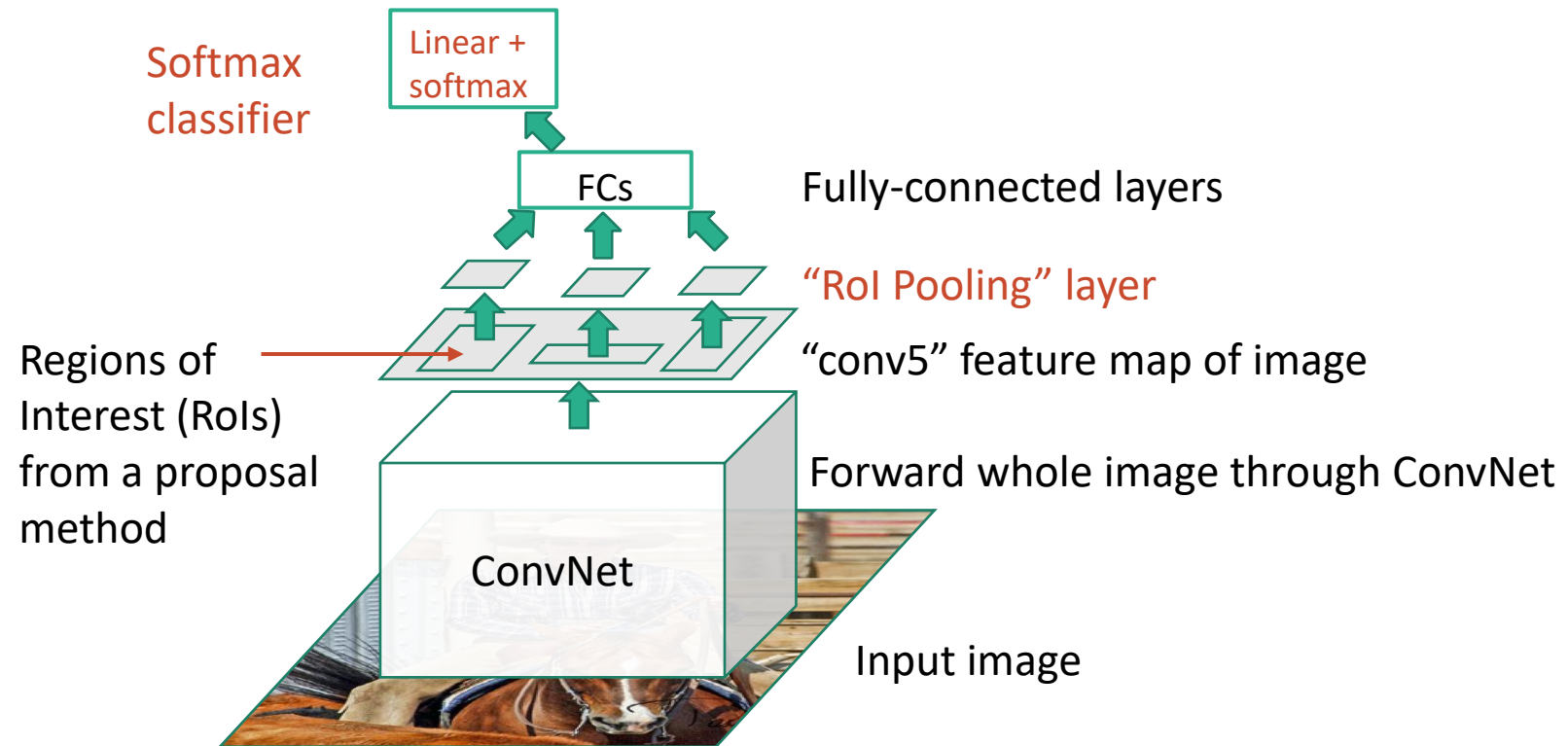
# Fast R-CNN



# Fast R-CNN

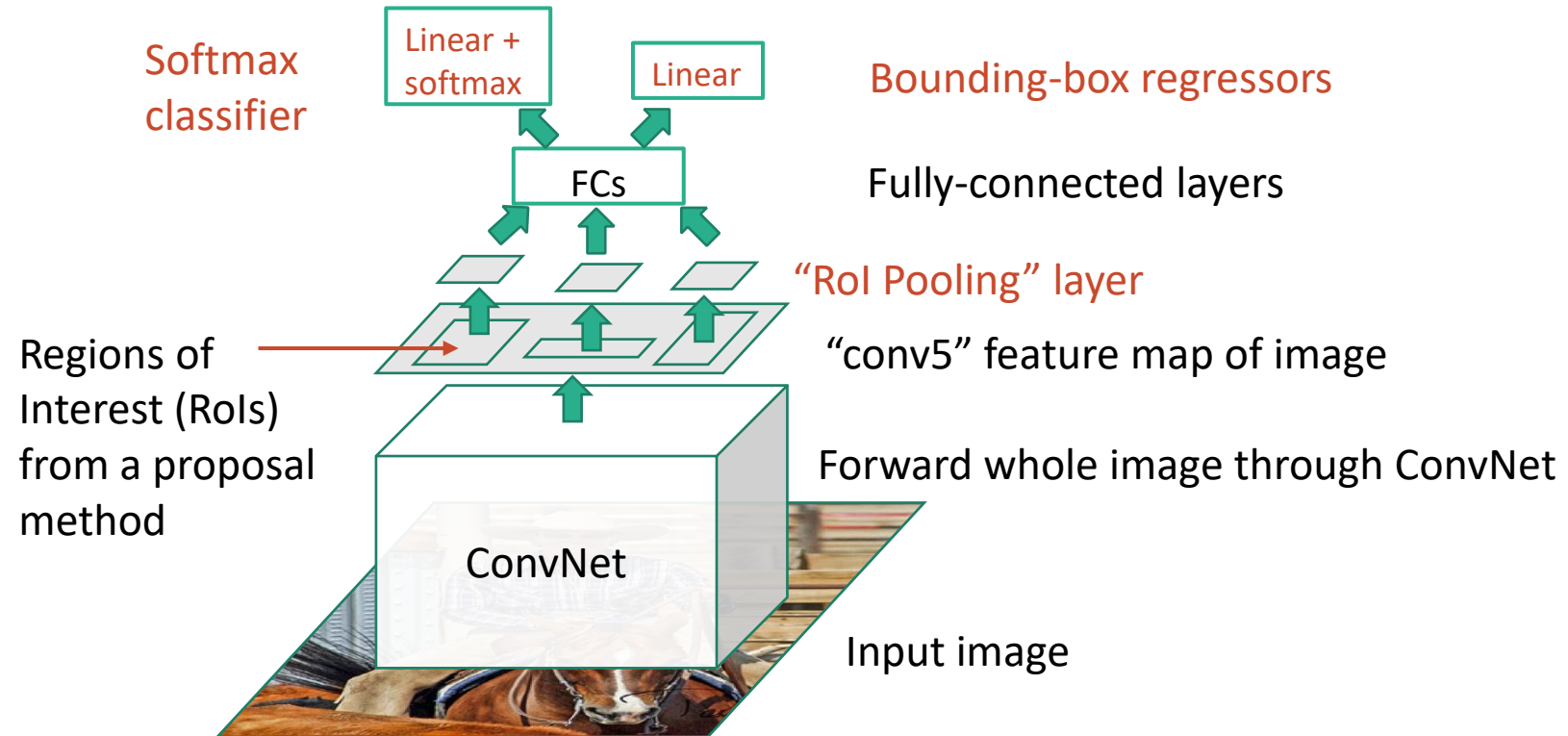


# Fast R-CNN

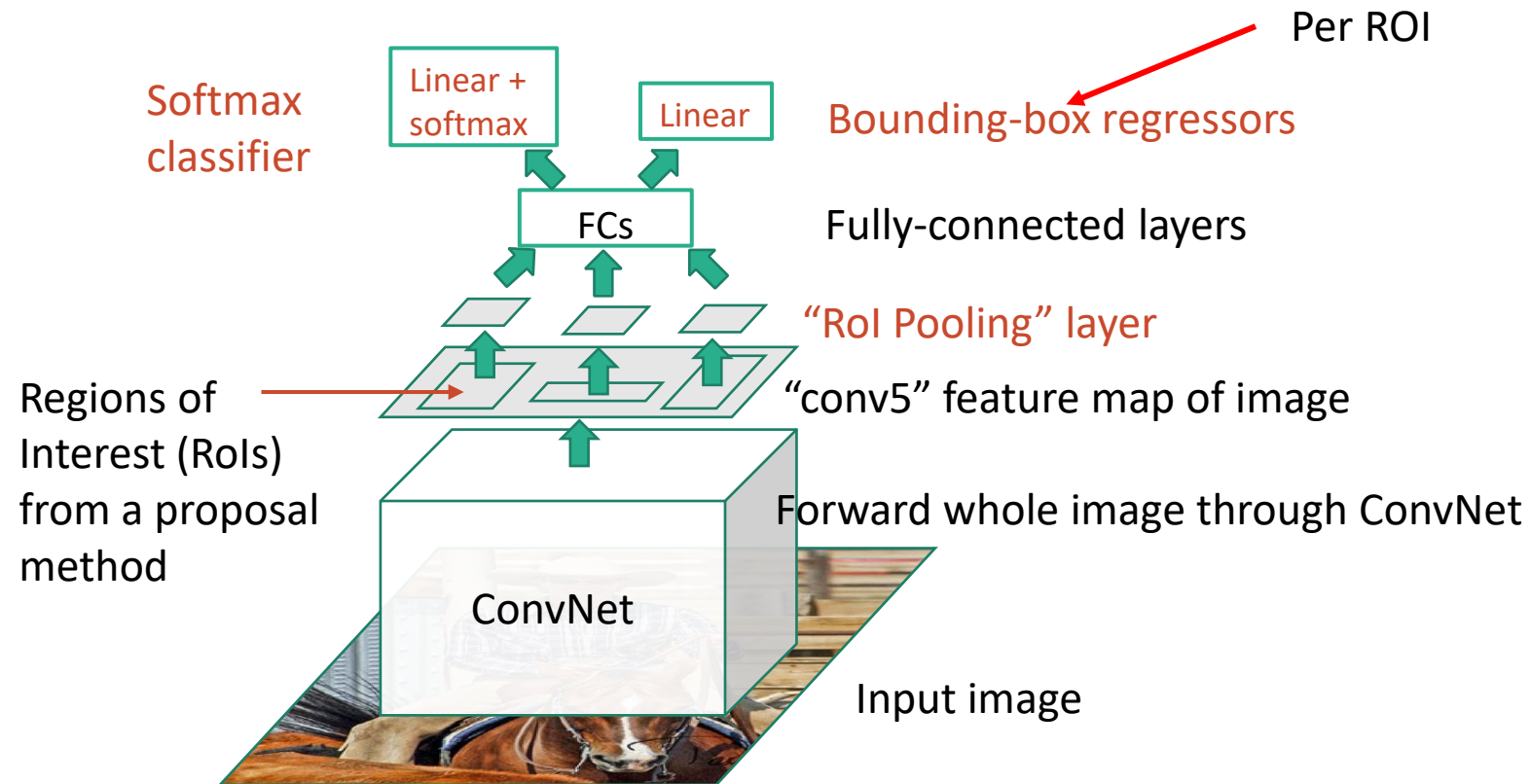




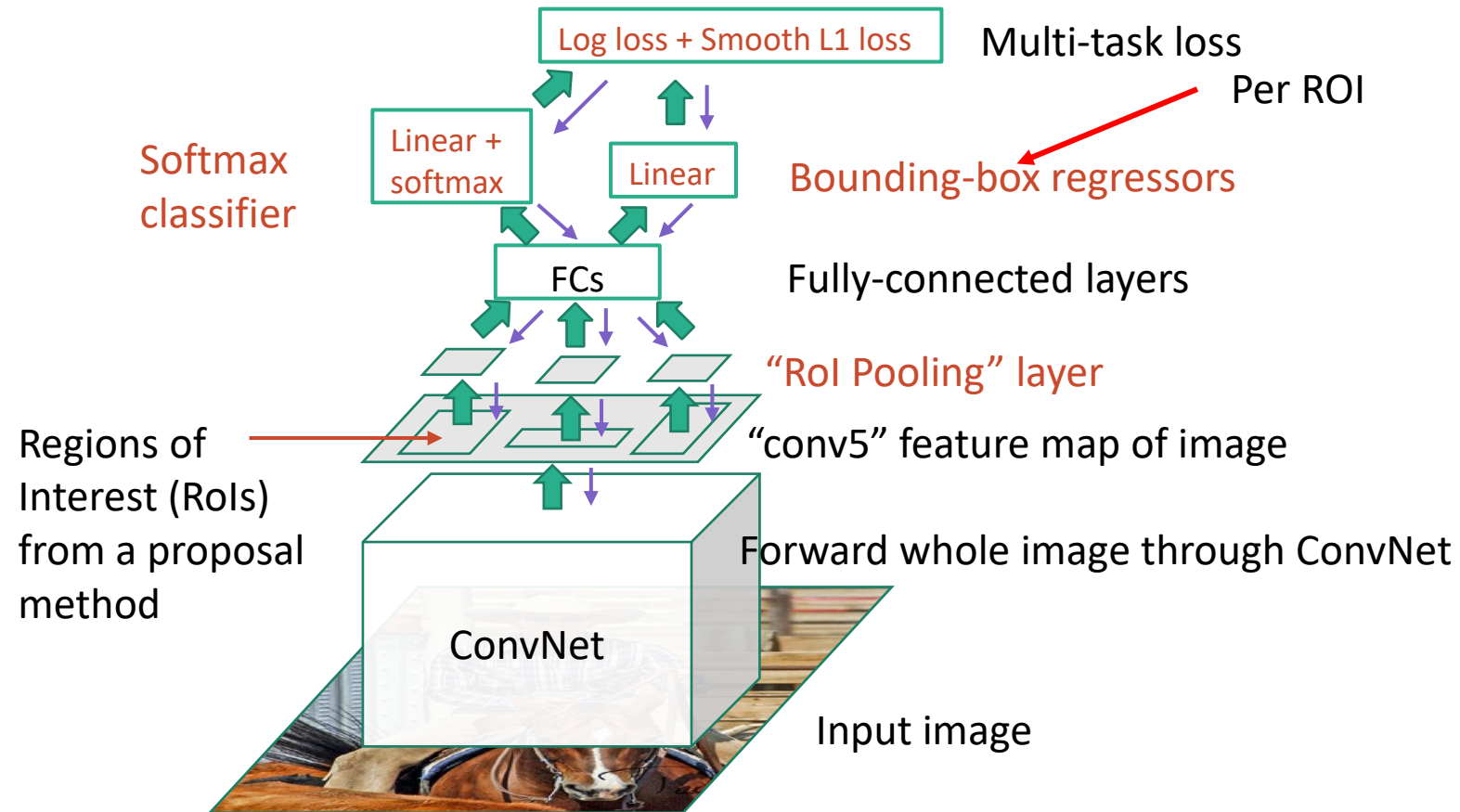
# Fast R-CNN



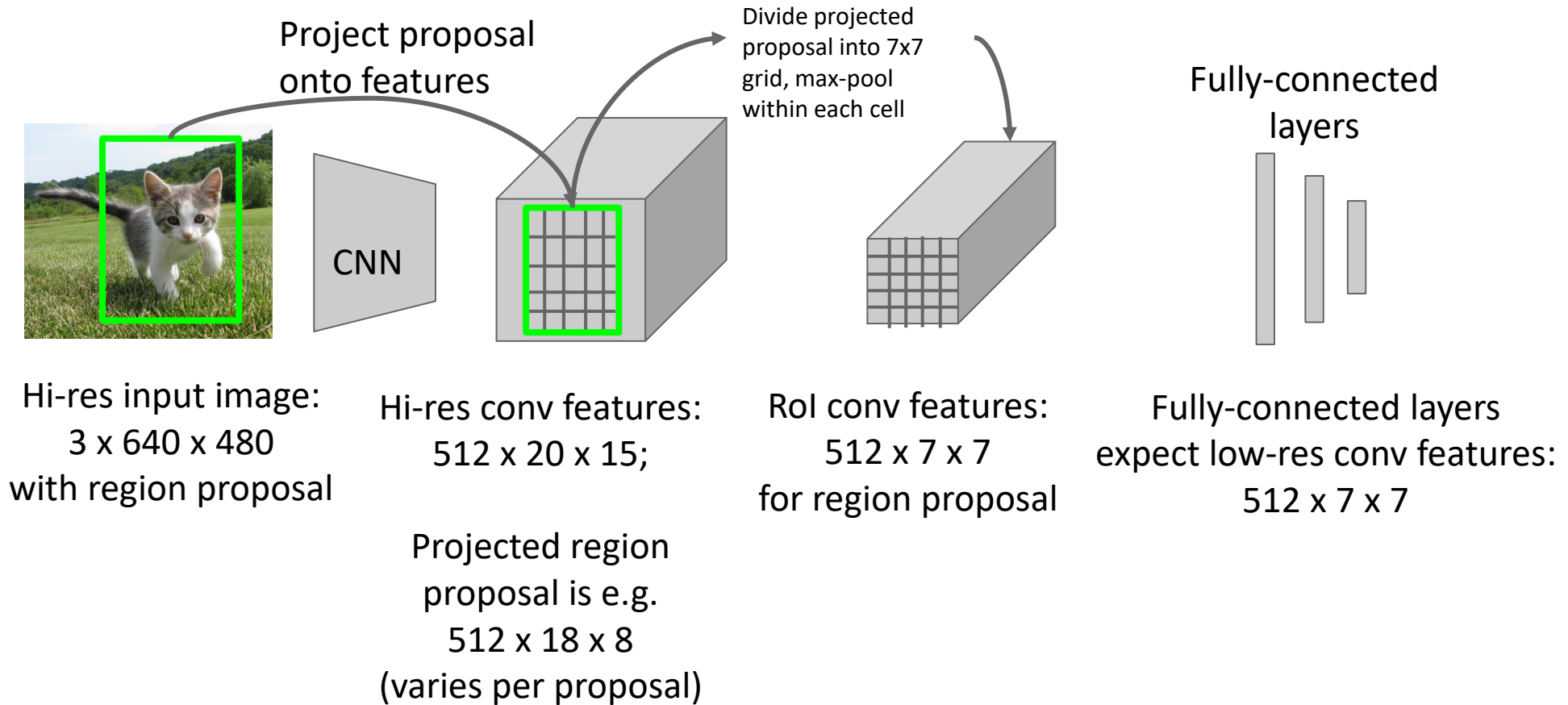
# Fast R-CNN



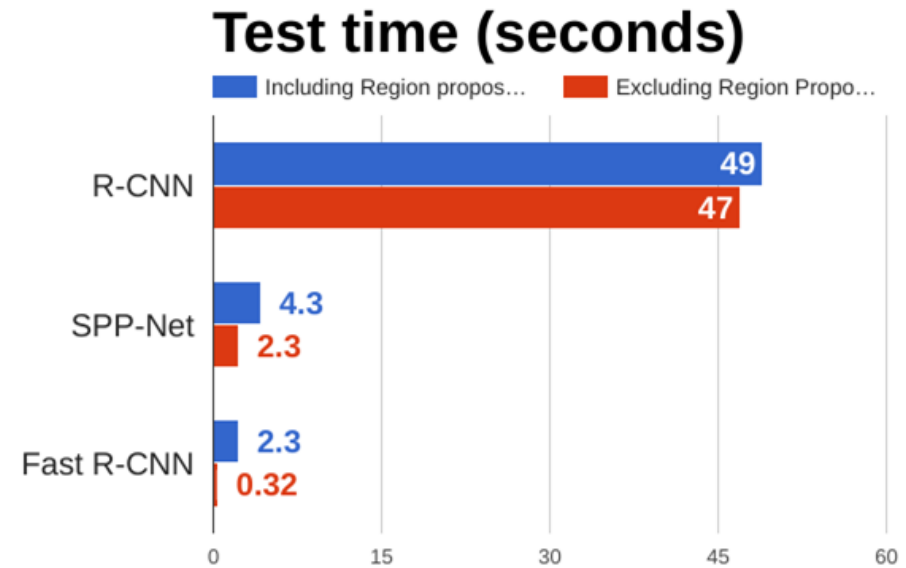
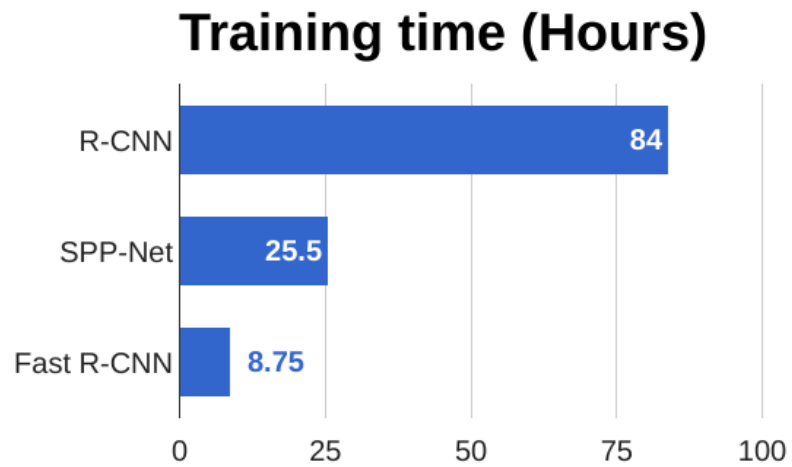
# Fast R-CNN (Training)



# Fast R-CNN: RoI Pooling



# R-CNN vs SPP vs Fast R-CNN



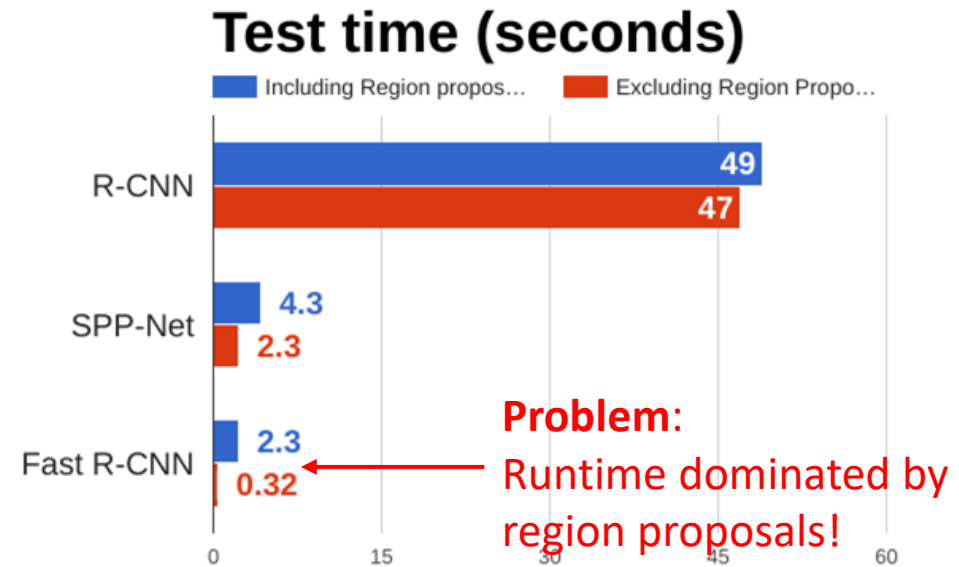
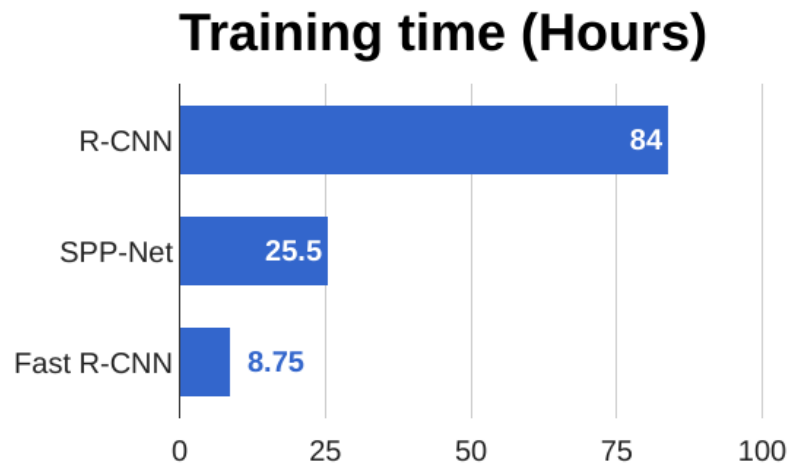
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014

Girshick, "Fast R-CNN", ICCV 2015

Based on slides for [Stanford cs231n](#) by Li, Jonson, and Young. Modified and reused with permission

# R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014

Girshick, "Fast R-CNN", ICCV 2015

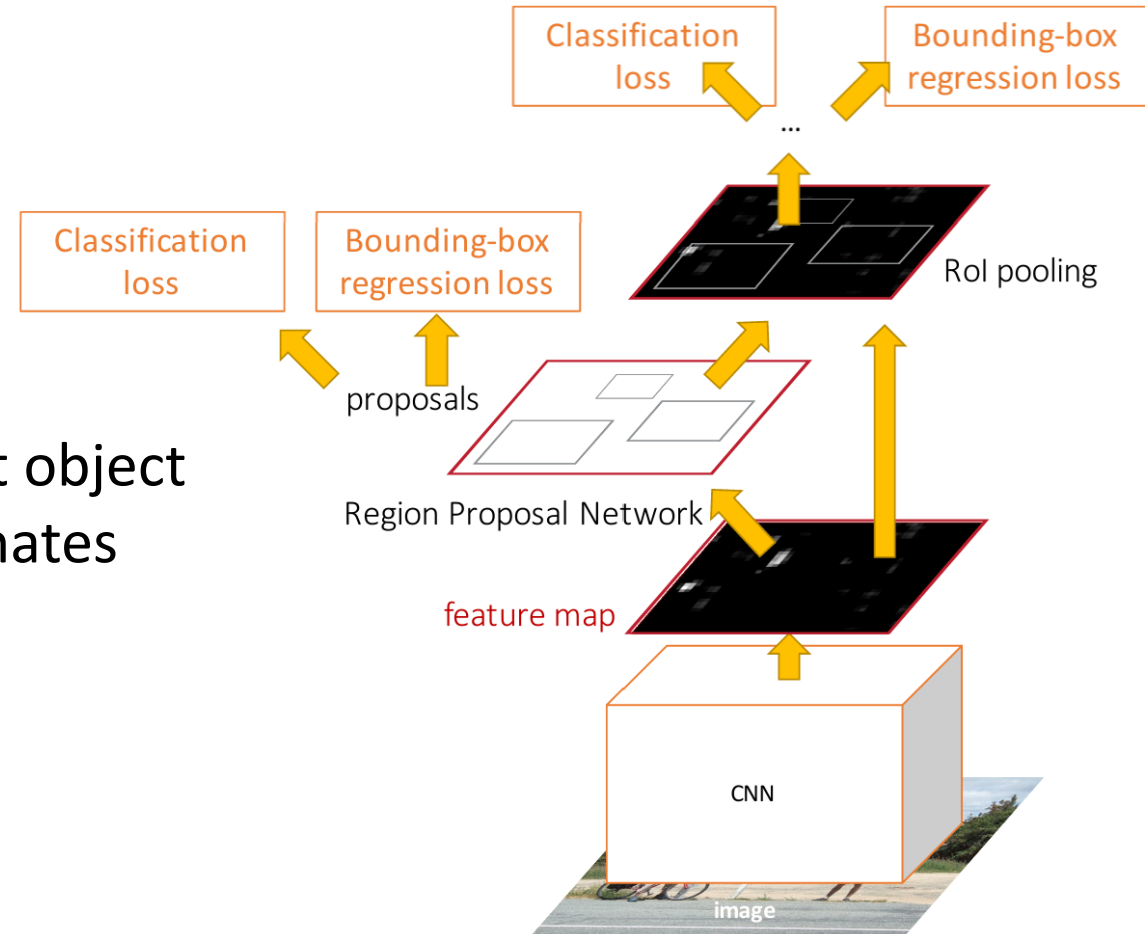
Based on slides for [Stanford cs231n](#) by Li, Jonson, and Young. Modified and reused with permission

# Faster R-CNN

Insert **Region Proposal Network (RPN)** to predict proposals from features

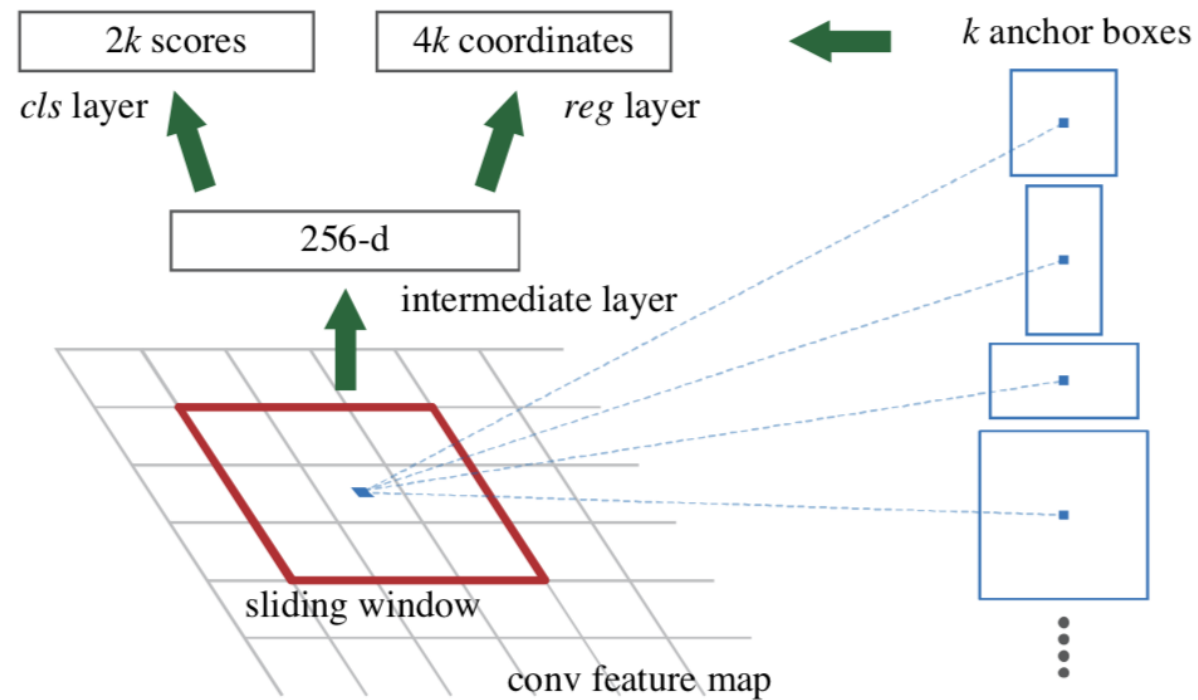
Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



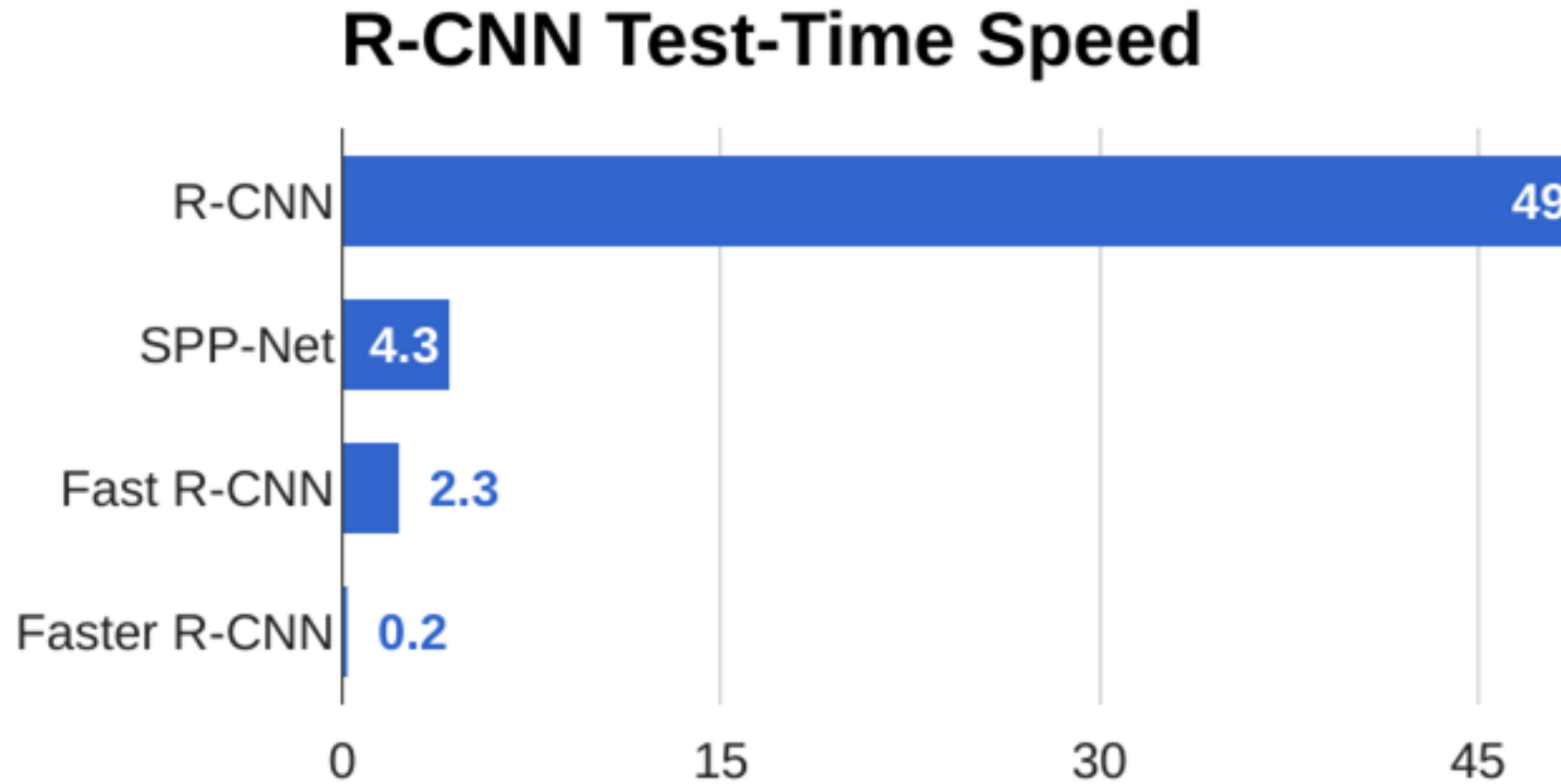
# Faster R-CNN

## Region Proposal Network (RPN)





# Faster R-CNN



# Today's Lecture

- **What's Object Detection?**
- **R-CNN / Fast RCNN / Faster RCNN**
- SSD / YOLO