



# Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

## Lecture 7: Convolutional Neural Networks (part 3)

# Logistics:

**Assignment 2** was due **yesterday**

**Assignment 3** will be posted **tonight** ...

# Logistics:

**Assignment 2** was due **yesterday**

**Assignment 3** will be posted **tonight** ...

**Final Projects** ... poll results

... 3 more students voted for individual projects than for survey

we will start process of forming groups



# Computer **Vision Problems** (no language for now)

## Categorization

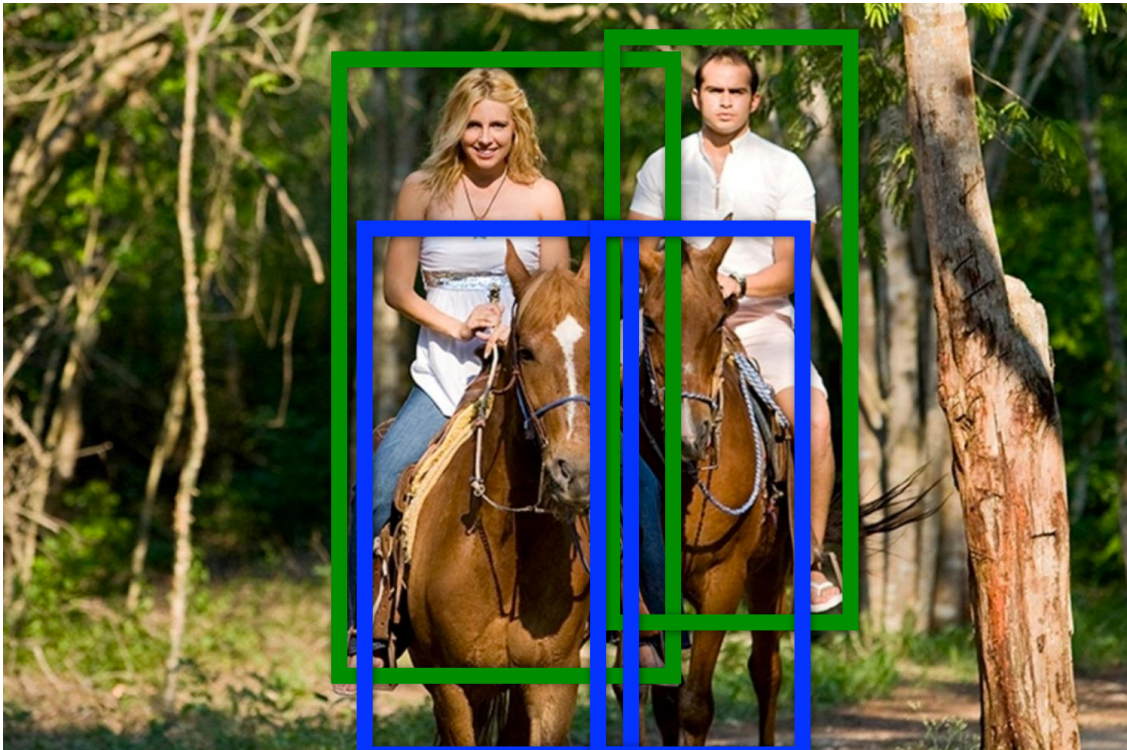


Multi-**class**:  
Horse  
Church  
Toothbrush  
**Person**



Multi-**label**:  
**Horse**  
Church  
Toothbrush  
**Person**

## Detection



Horse (x, y, w, h)  
Horse (x, y, w, h)  
Person (x, y, w, h)  
Person (x, y, w, h)



## Segmentation



Horse  
Person



## Instance Segmentation

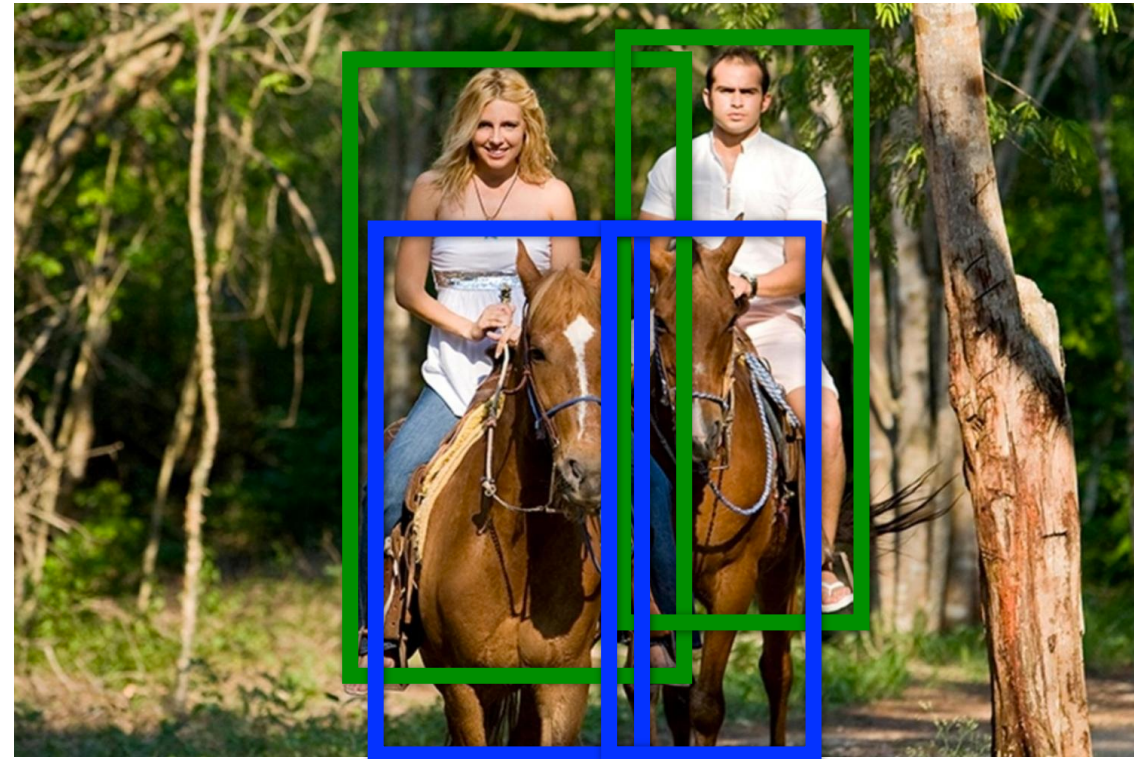


Horse1  
Horse2  
Person1  
Person2

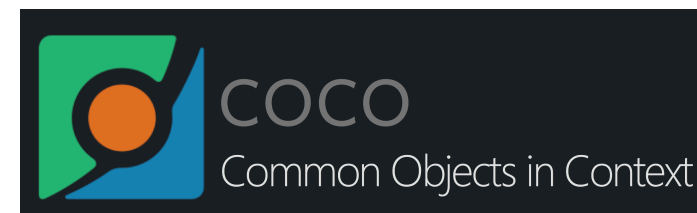


# Computer **Vision Problems** (no language for now)

## Detection



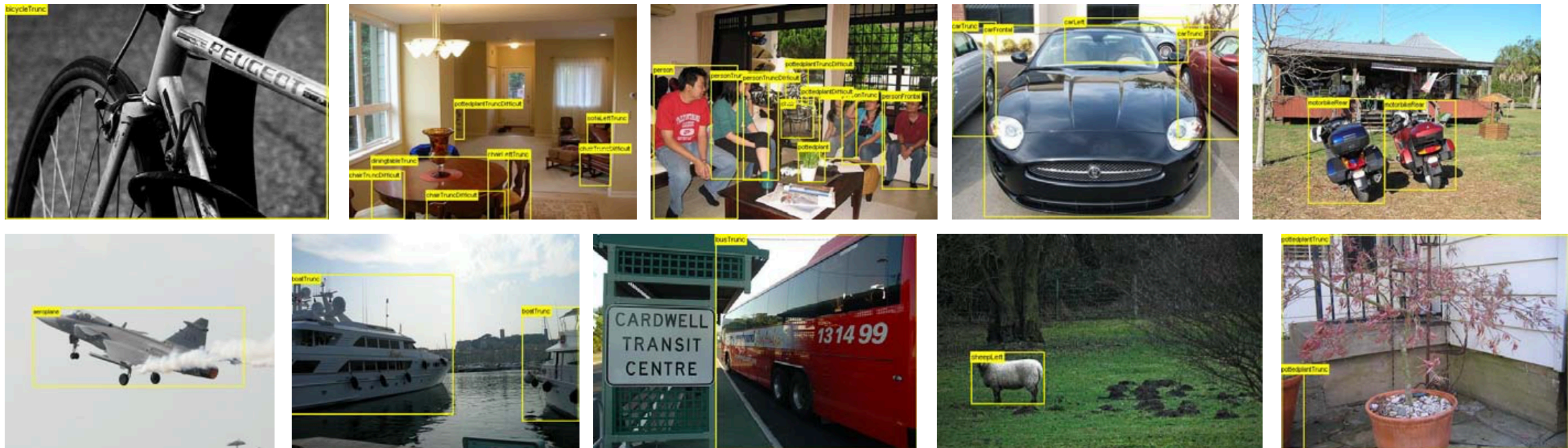
Horse (x, y, w, h)  
Horse (x, y, w, h)  
Person (x, y, w, h)  
Person (x, y, w, h)





# Datasets: Pascal VOC

20 classes: aeroplane, bicycle, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, train, TV



Real images downloaded from flickr, not filtered for “quality”



# Datasets: Pascal VOC

20 classes: aeroplane, bicycle, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, train, TV



|                | Training | Testing |
|----------------|----------|---------|
| <b>Images</b>  | 10,103   | 9,637   |
| <b>Objects</b> | 23,374   | 22,992  |

Real images downloaded from flickr, not filtered for “quality”



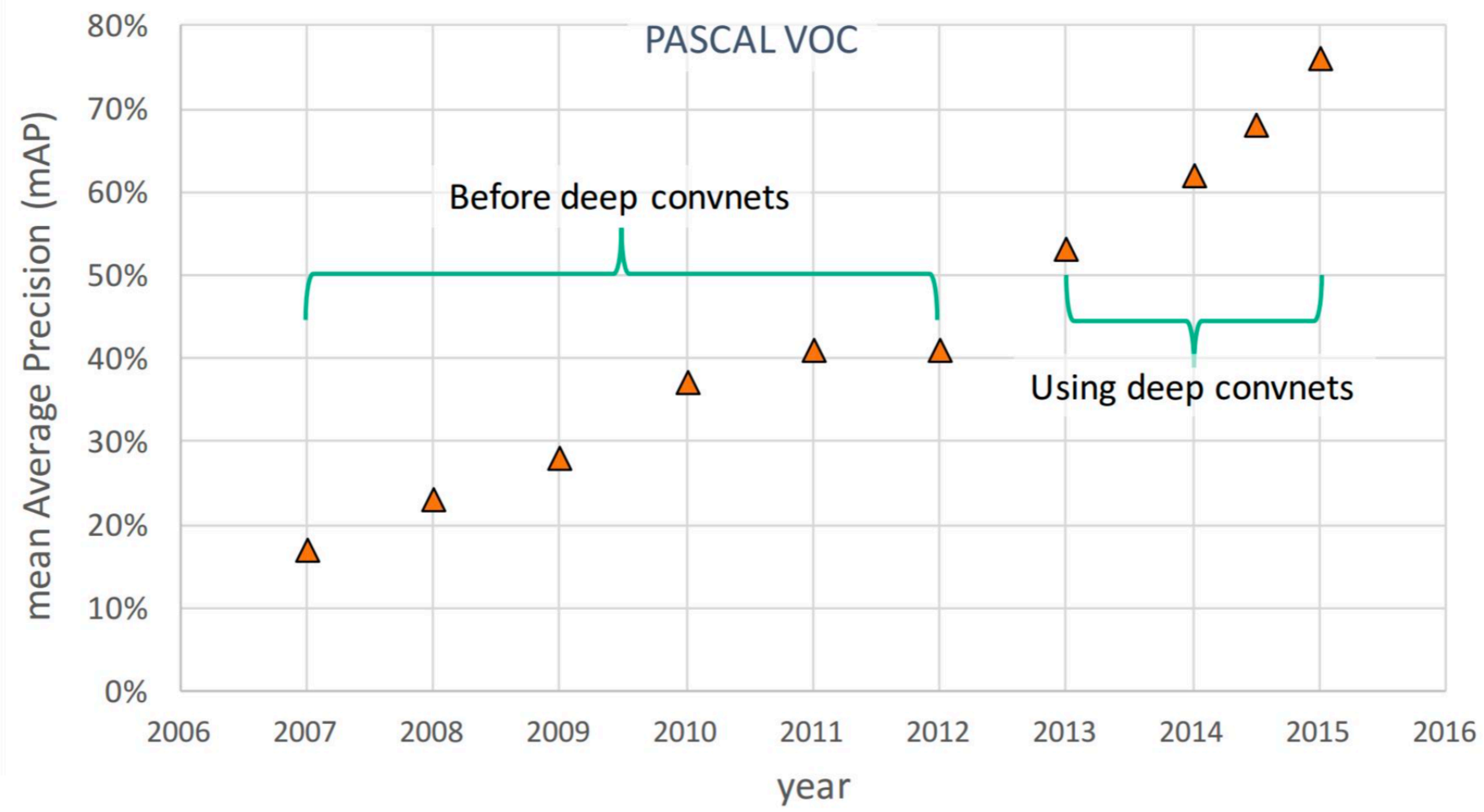
# Datasets: COCO



- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints



# Object Detection



\* plot from Ross Girshick, 2015

# Object **Detection** as Regression Problem



CAT (x, y, w, h)



# Object **Detection** as Regression Problem



CAT (x, y, w, h)

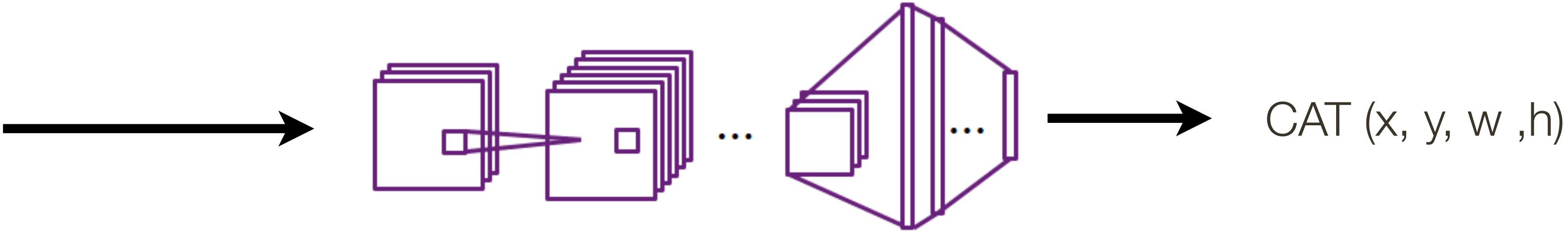


DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
DUCK (x, y, w, h)  
...

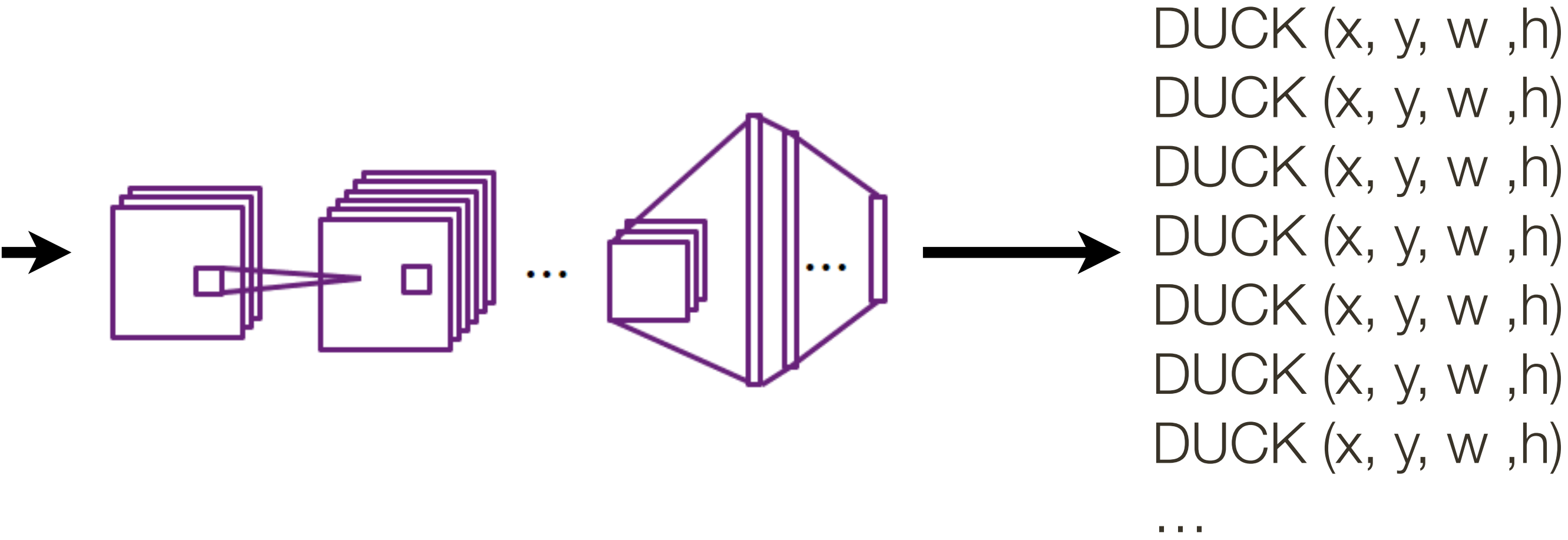
\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**



# Object **Detection** as Regression Problem



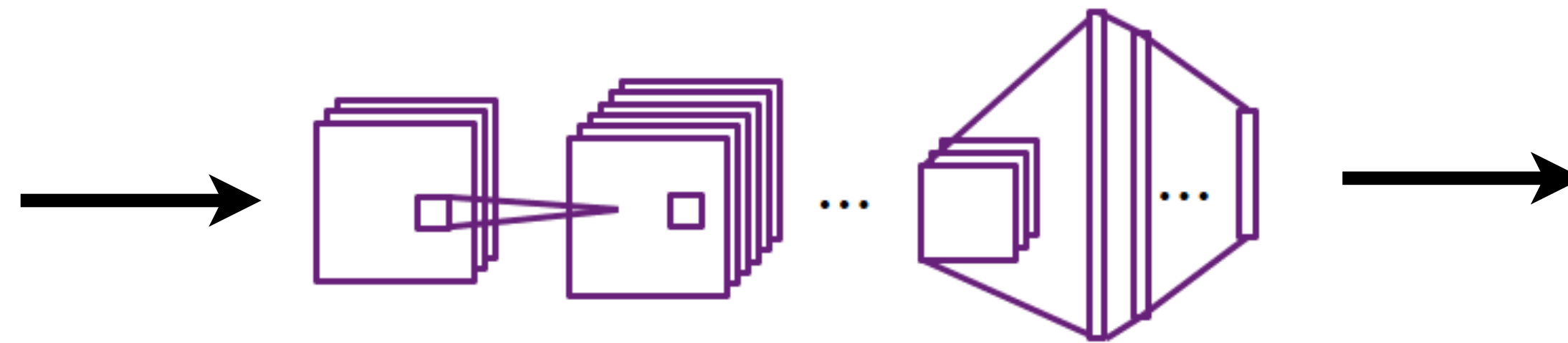
**Problem:** each image needs a different number of outputs



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**



# Object **Detection** as Classification Problem

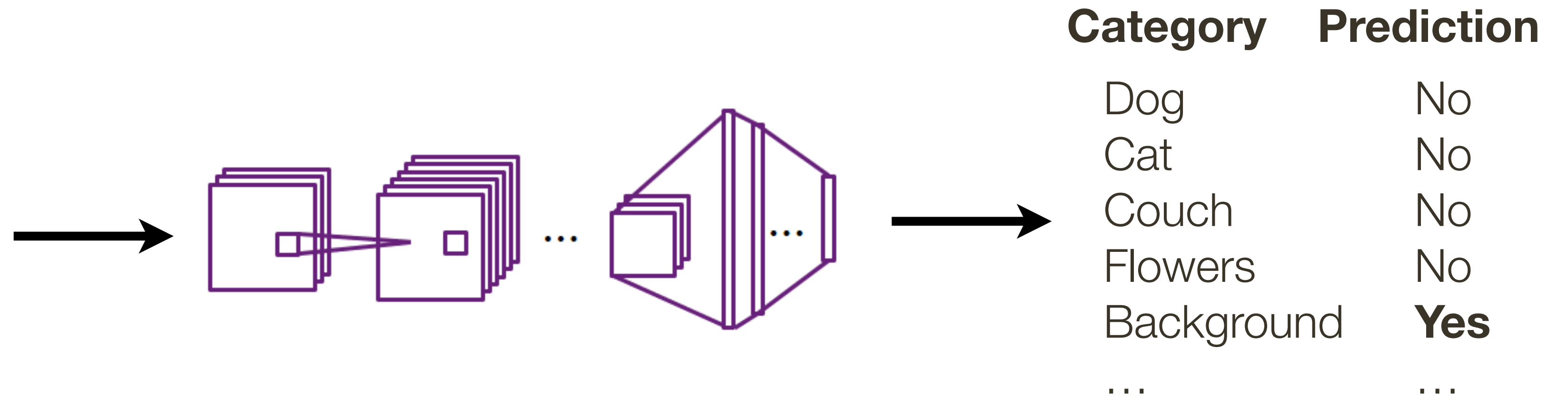
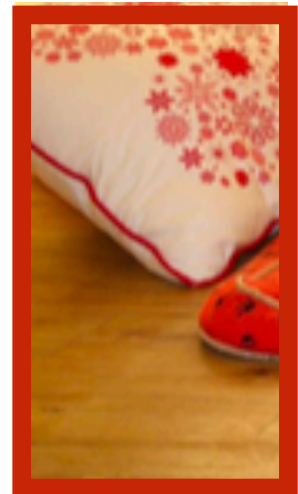


| Category   | Prediction |
|------------|------------|
| Dog        | No         |
| Cat        | No         |
| Couch      | No         |
| Flowers    | No         |
| Background | <b>Yes</b> |
| ...        | ...        |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

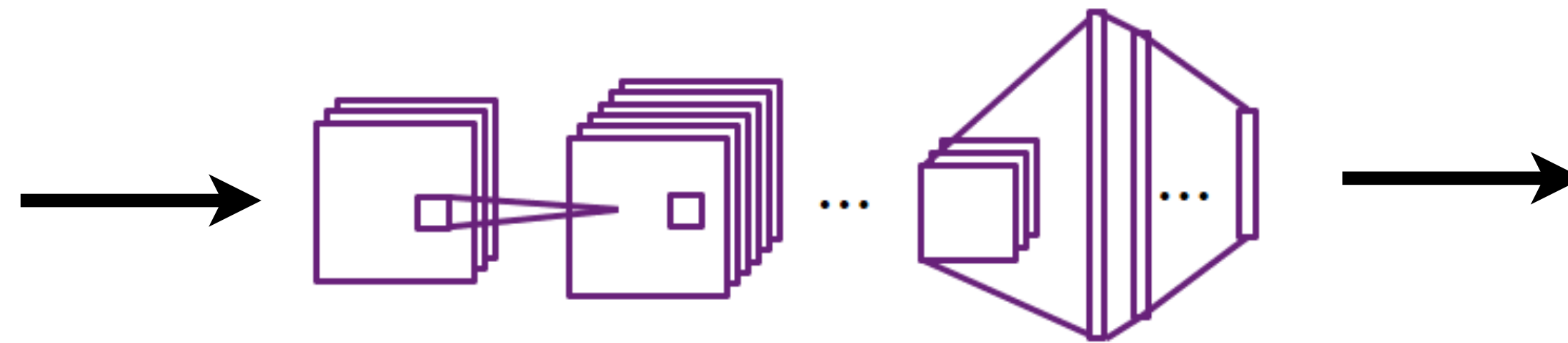


# Object **Detection** as Classification Problem



Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

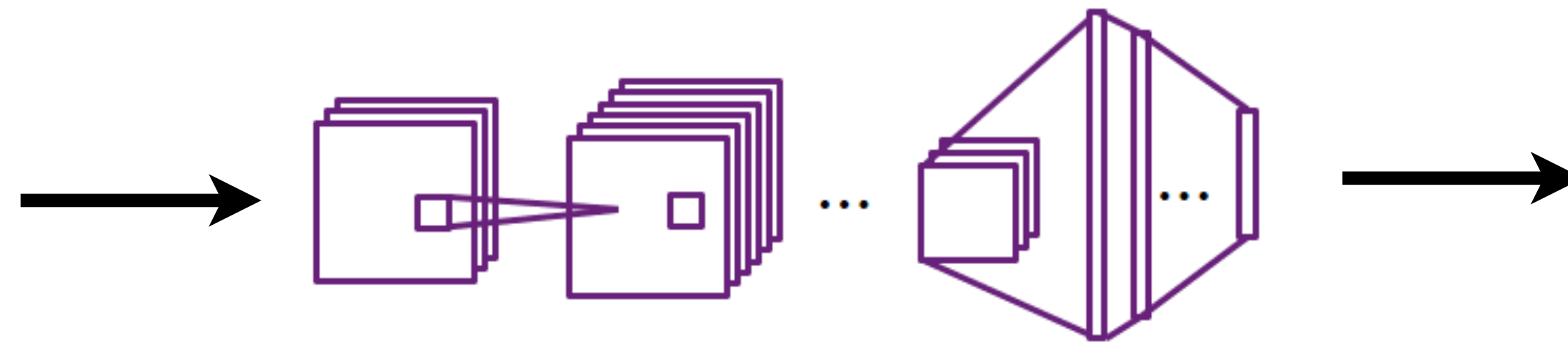
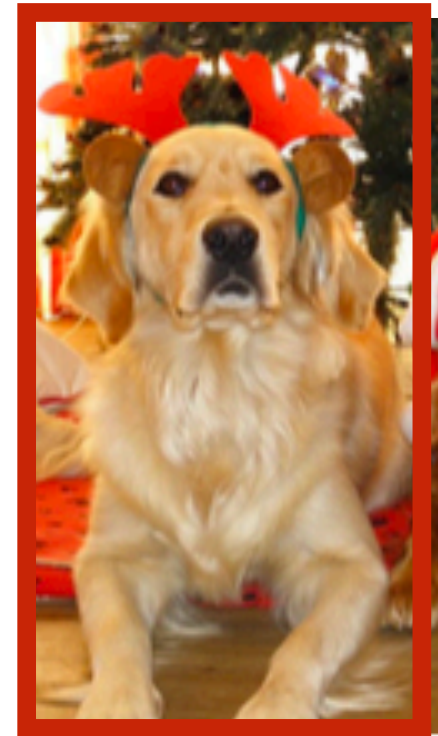
# Object **Detection** as Classification Problem



| Category   | Prediction |
|------------|------------|
| Dog        | <b>Yes</b> |
| Cat        | No         |
| Couch      | No         |
| Flowers    | No         |
| Background | No         |
| ...        | ...        |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem

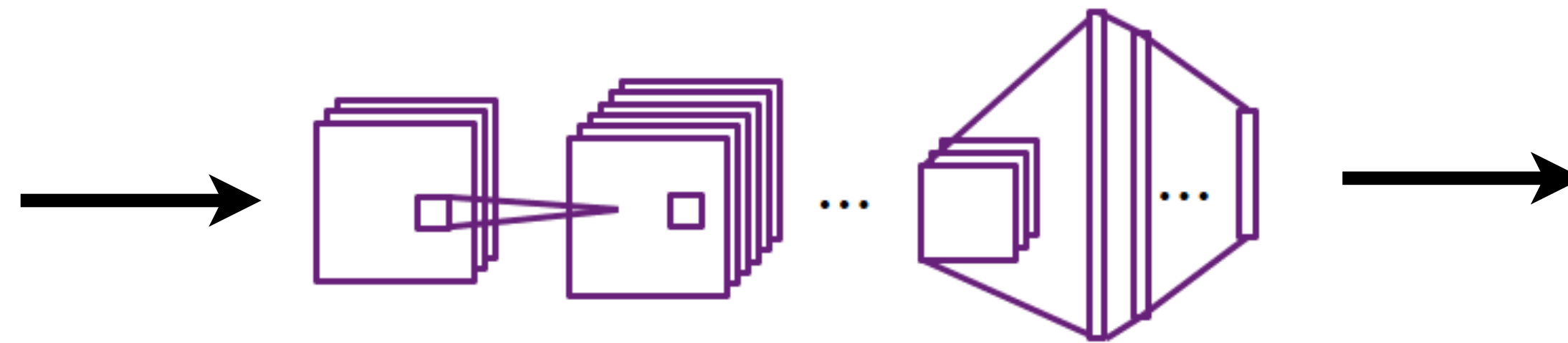


| Category   | Prediction |
|------------|------------|
| Dog        | <b>Yes</b> |
| Cat        | No         |
| Couch      | No         |
| Flowers    | No         |
| Background | No         |
| ...        | ...        |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background



# Object **Detection** as Classification Problem



| Category   | Prediction |
|------------|------------|
| Dog        | No         |
| Cat        | <b>Yes</b> |
| Couch      | No         |
| Flowers    | No         |
| Background | No         |
| ...        | ...        |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background

# Object **Detection** as Classification Problem

**Problem:** Need to apply CNN to **many** patches in each image



| Category   | Prediction |
|------------|------------|
| Dog        | No         |
| Cat        | <b>Yes</b> |
| Couch      | No         |
| Flowers    | No         |
| Background | No         |
| ...        | ...        |

Apply CNN to many different crops in the image and (classification) CNN classifies each patch as object or background



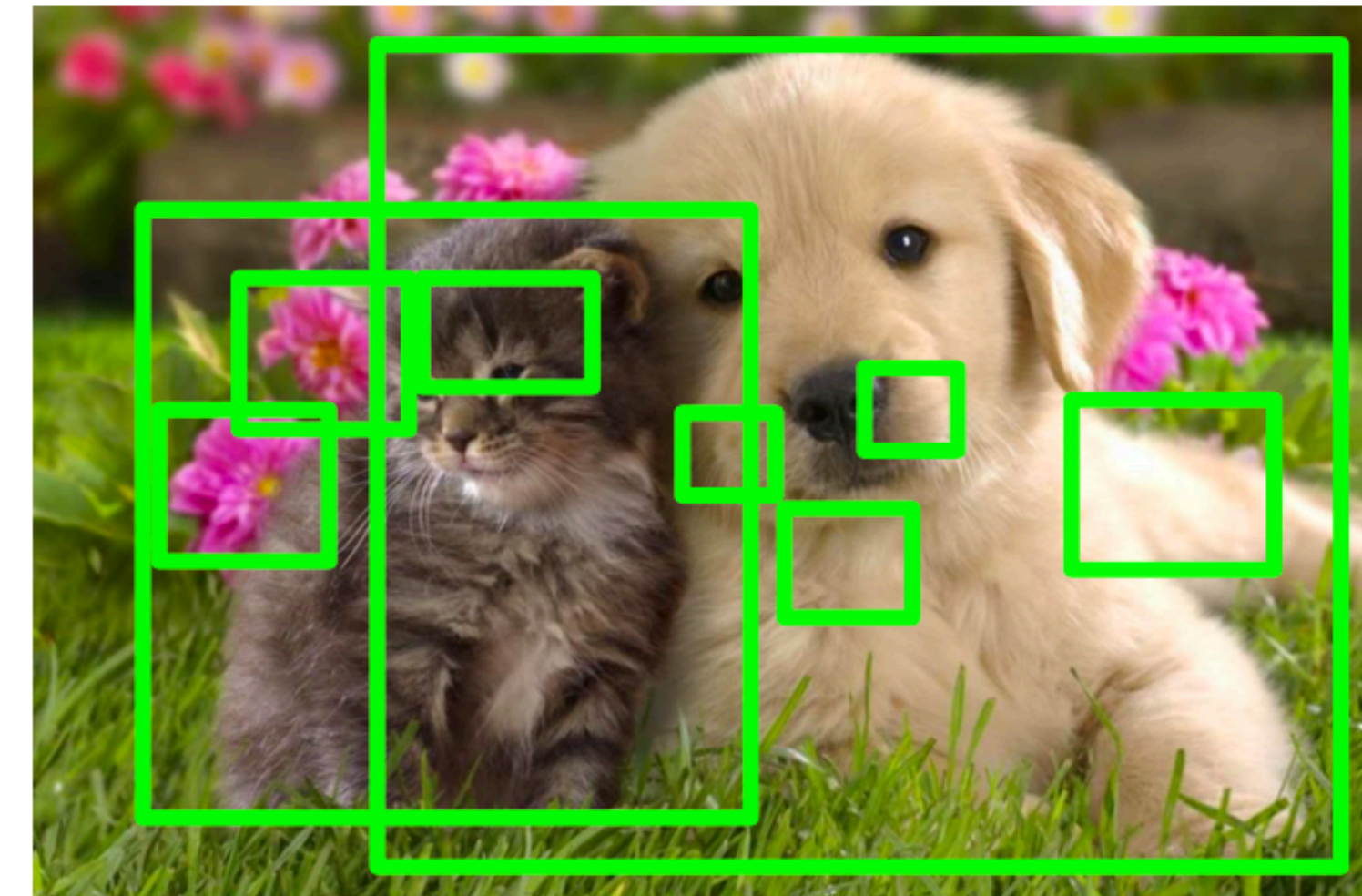
# Region Proposals (older idea in vision)

[ Alexe et al, TPAMI 2012 ]  
[ Ujkings et al, IJCV 2013 ]  
[ Cheng et al, CVPR 2014 ]  
[ Zitnick and Dollar, ECCV 2014 ]

Find image **regions that are likely contain objects** (any object at all)

- typically works by looking at histogram distributions, region aspect ratio, closed contours, coherent color

Relatively **fast to run** (Selective Search gives 1000 region proposals in a few seconds on a CPU)



**Goal:** Get “true” object regions to be in as few top K proposals as possible



# R-CNN

[ Girshick et al, CVPR 2014 ]



Input **Image**

\* image from Ross Girshick

# R-CNN

[ Girshick et al, CVPR 2014 ]



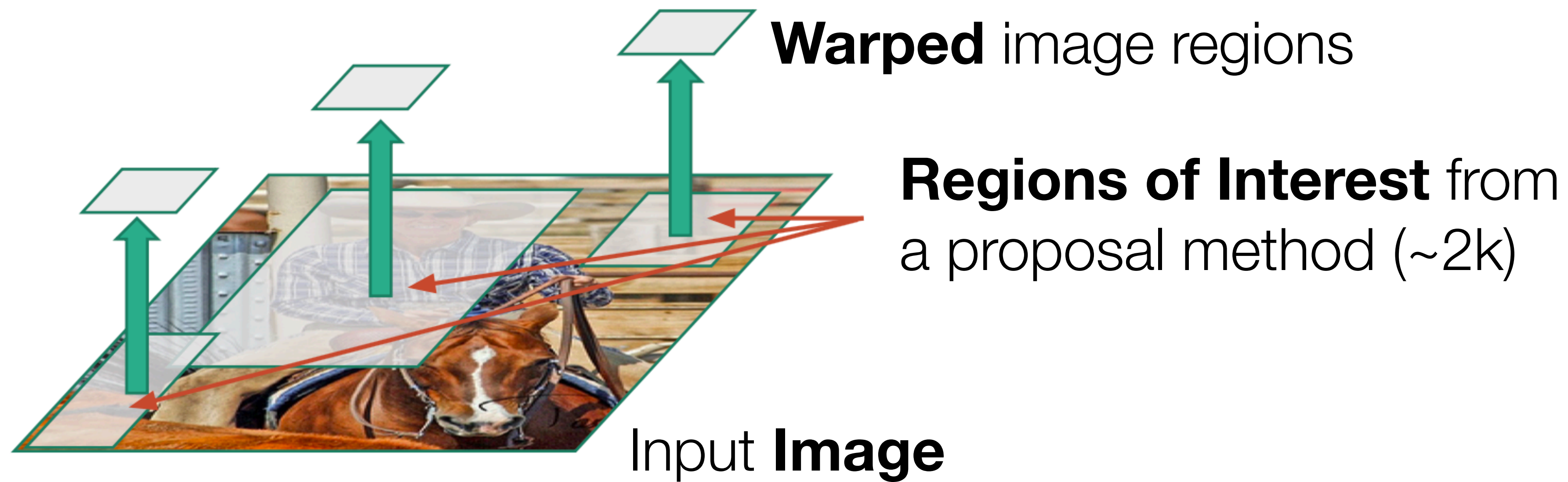
**Regions of Interest** from  
a proposal method (~2k)

Input **Image**



# R-CNN

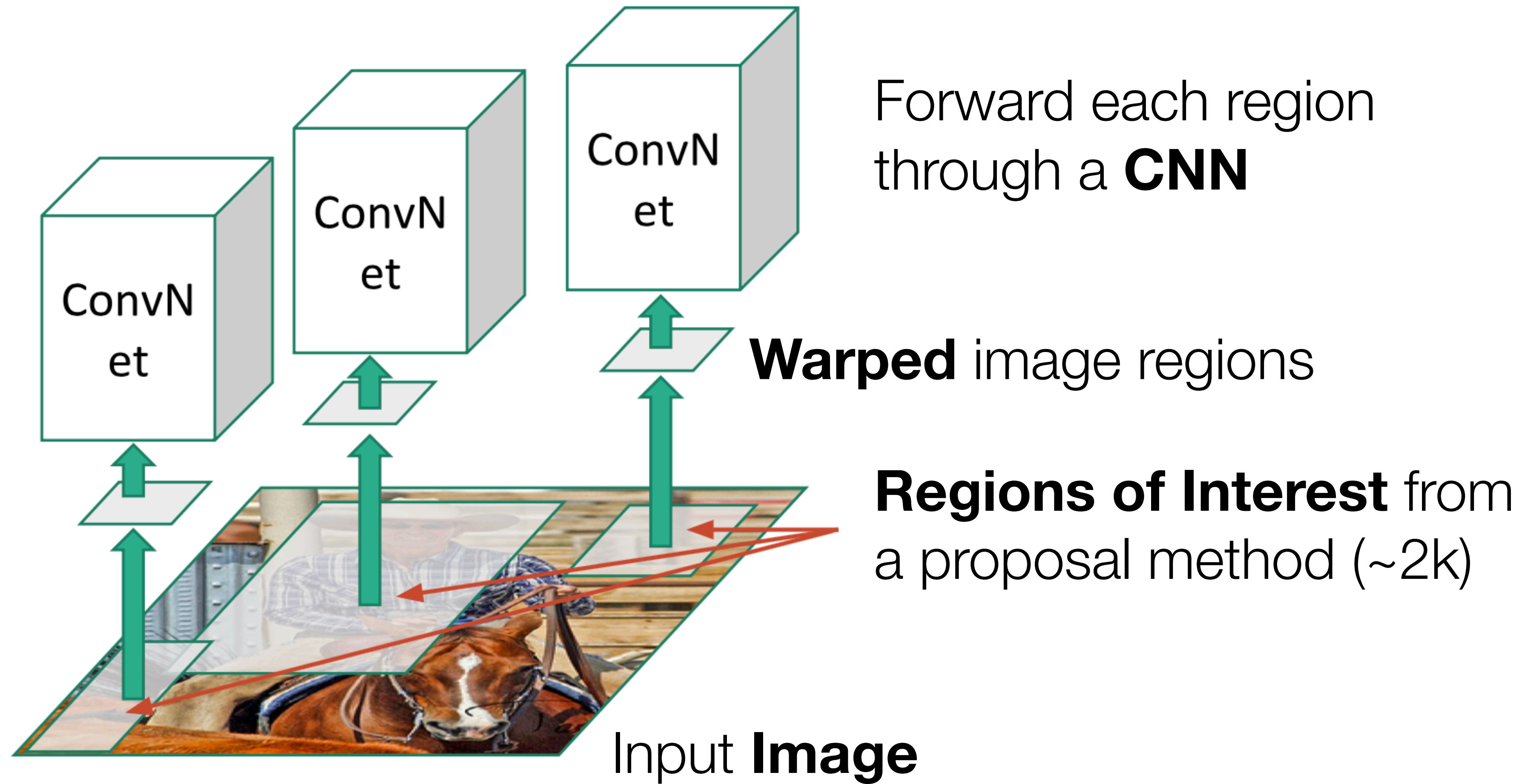
[ Girshick et al, CVPR 2014 ]



\* image from Ross Girshick

# R-CNN

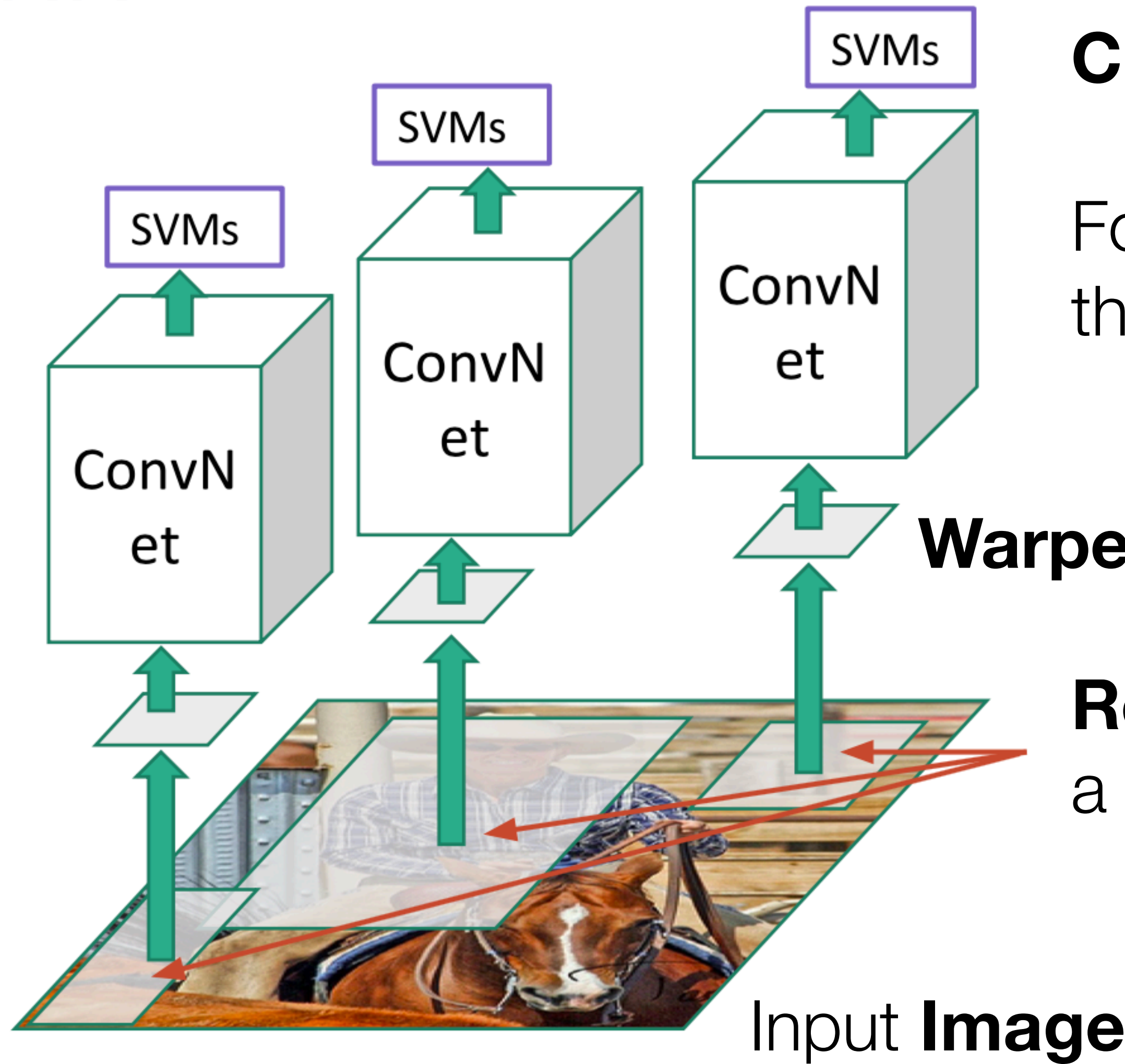
[ Girshick et al, CVPR 2014 ]





# R-CNN

[ Girshick et al, CVPR 2014 ]



**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

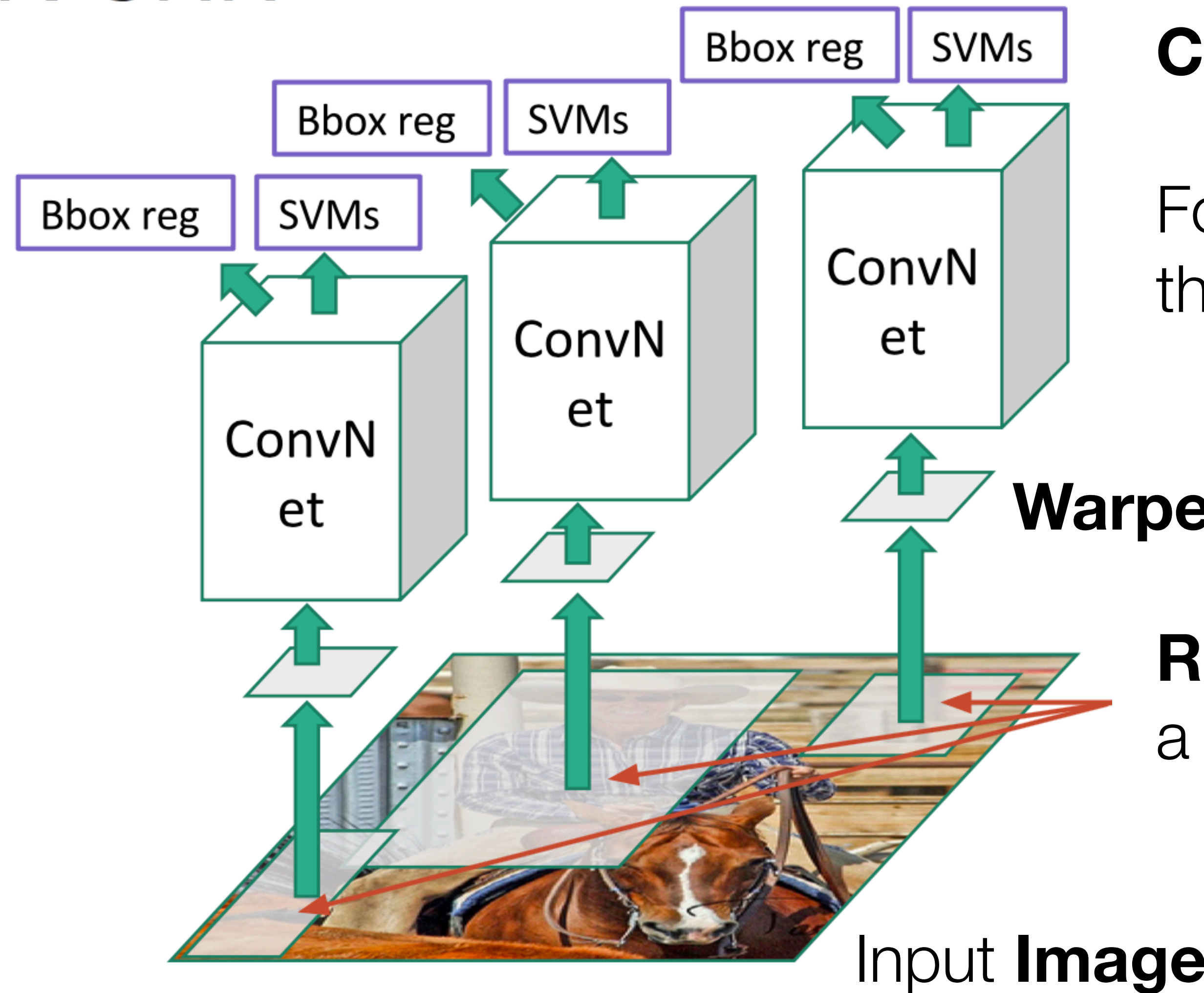
Input **Image**



# R-CNN

**Linear Regression** for bounding box offsets

[ Girshick et al, CVPR 2014 ]



**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

Input **Image**

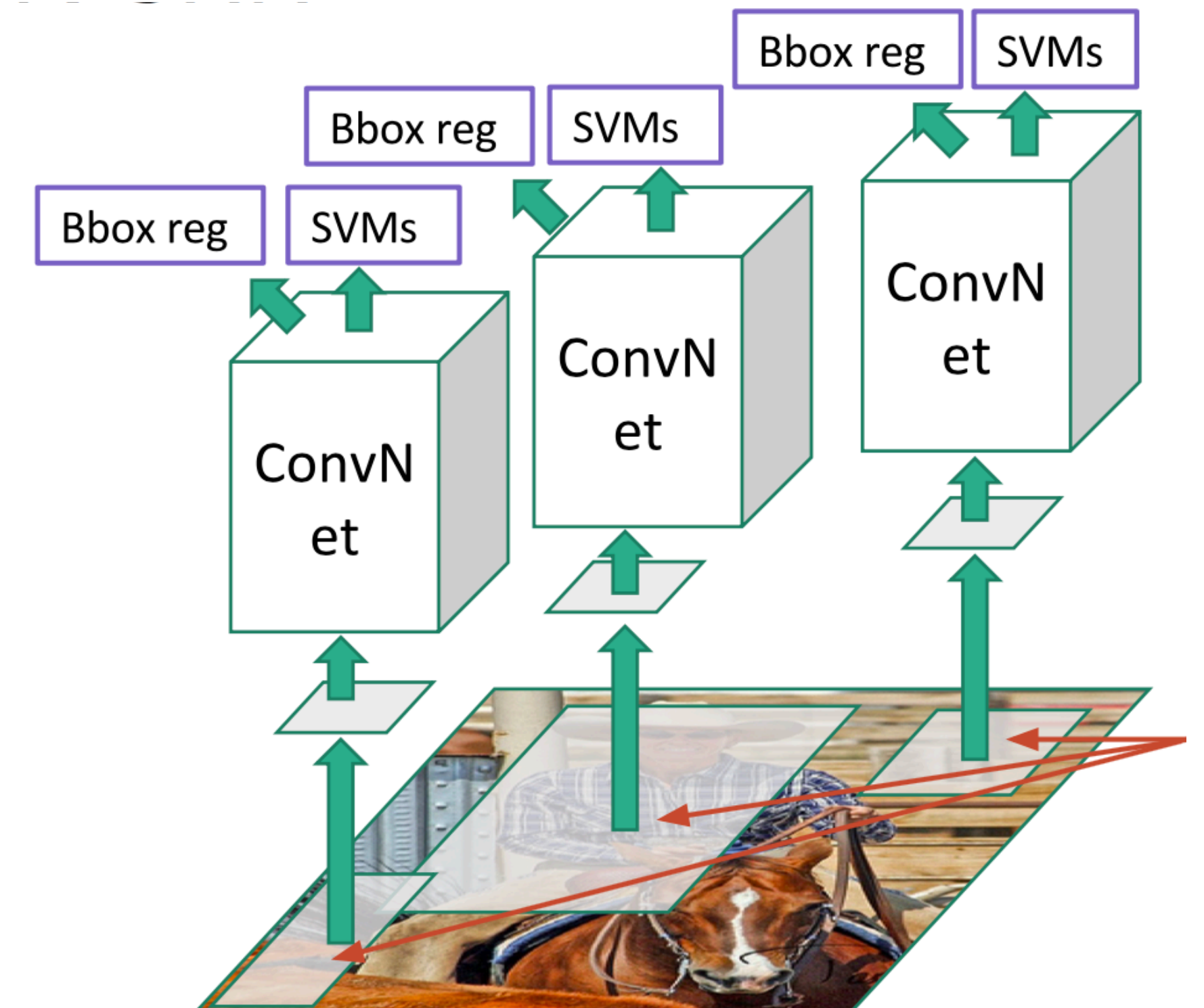


# R-CNN: Training

[ Girshick et al, CVPR 2014 ]

## Fine-tuning ImageNet CNN on object proposal patches

- $> 50\%$  Intersection-over-Union overlap with GT considered “object” others “background”
- batches of 128 (**32 positives, 96 negatives**)



\* image from Ross Girshick

# R-CNN: Issues

[ Girshick et al, CVPR 2014 ]

## Ad-hoc training objectives

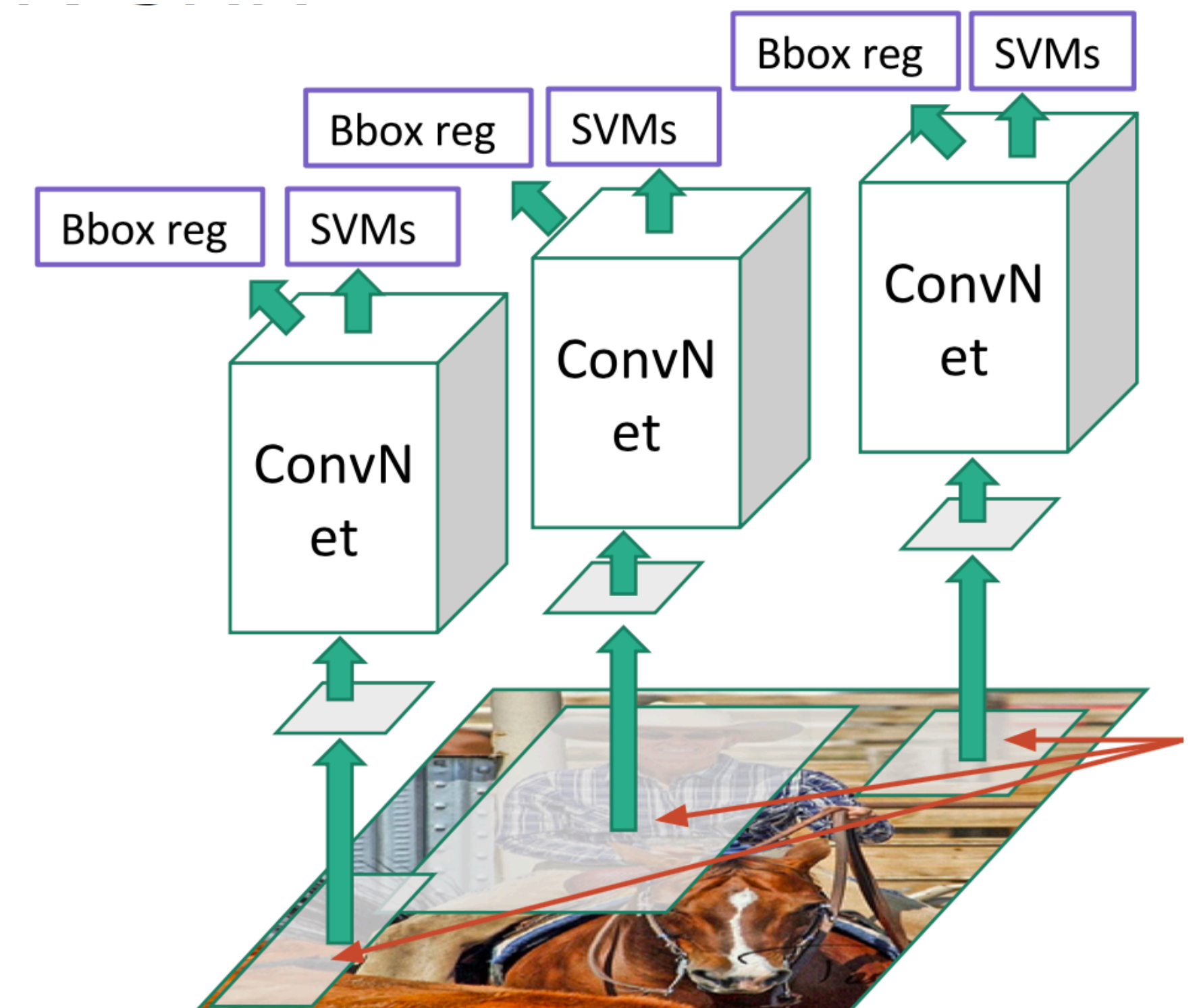
- Fine-tune network with softmax objective (**log** loss)
- Train post-hoc linear SVM (**hinge** loss)
- Train post-hoc bounding-box regression (**least squares**)

## Training is slow

- 84 hours and takes a lot of disk space

## Inference / **Detection is slow**

- 47 sec / image with VGG16 [ Simonyan et al, ICLR 2015 ]

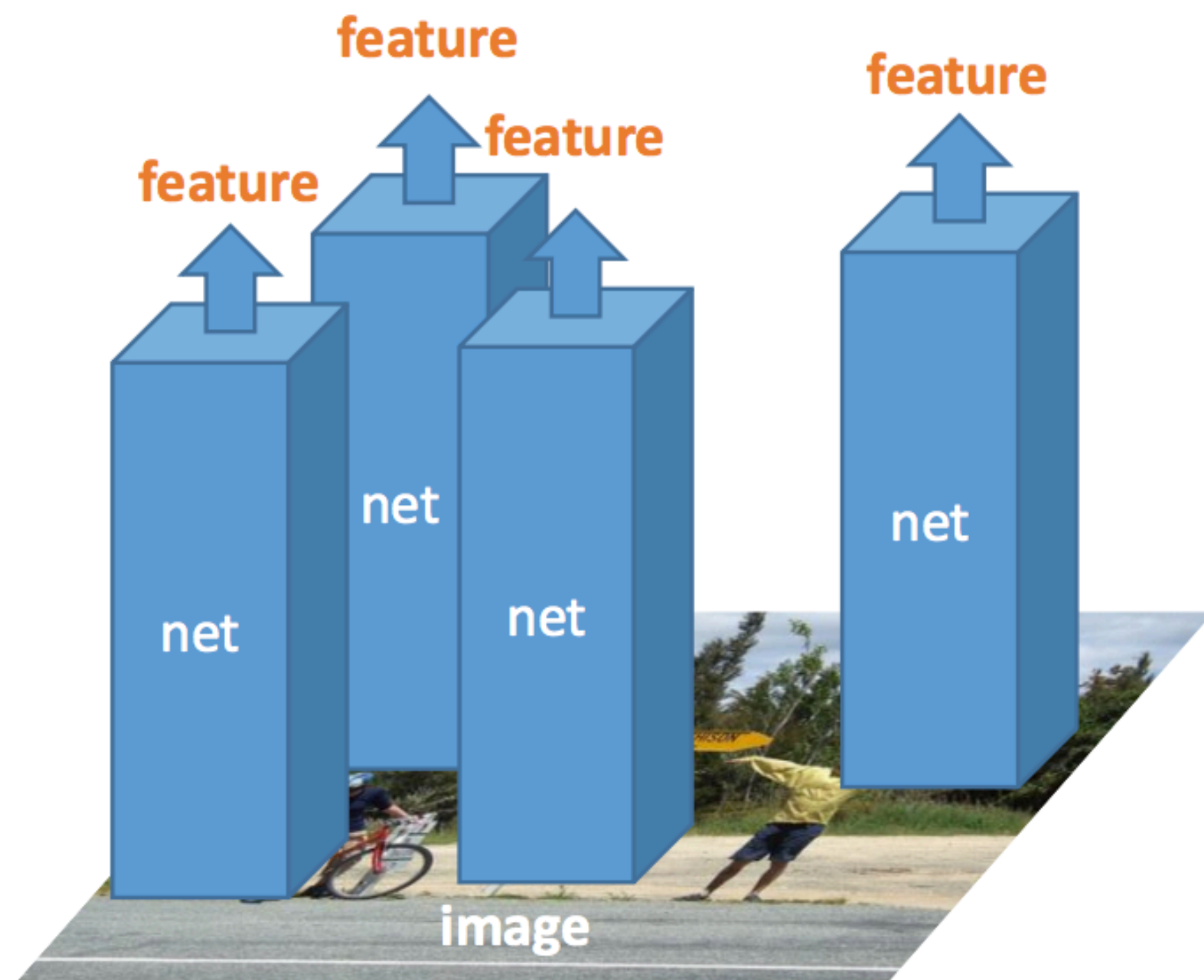


\* image from Ross Girshick



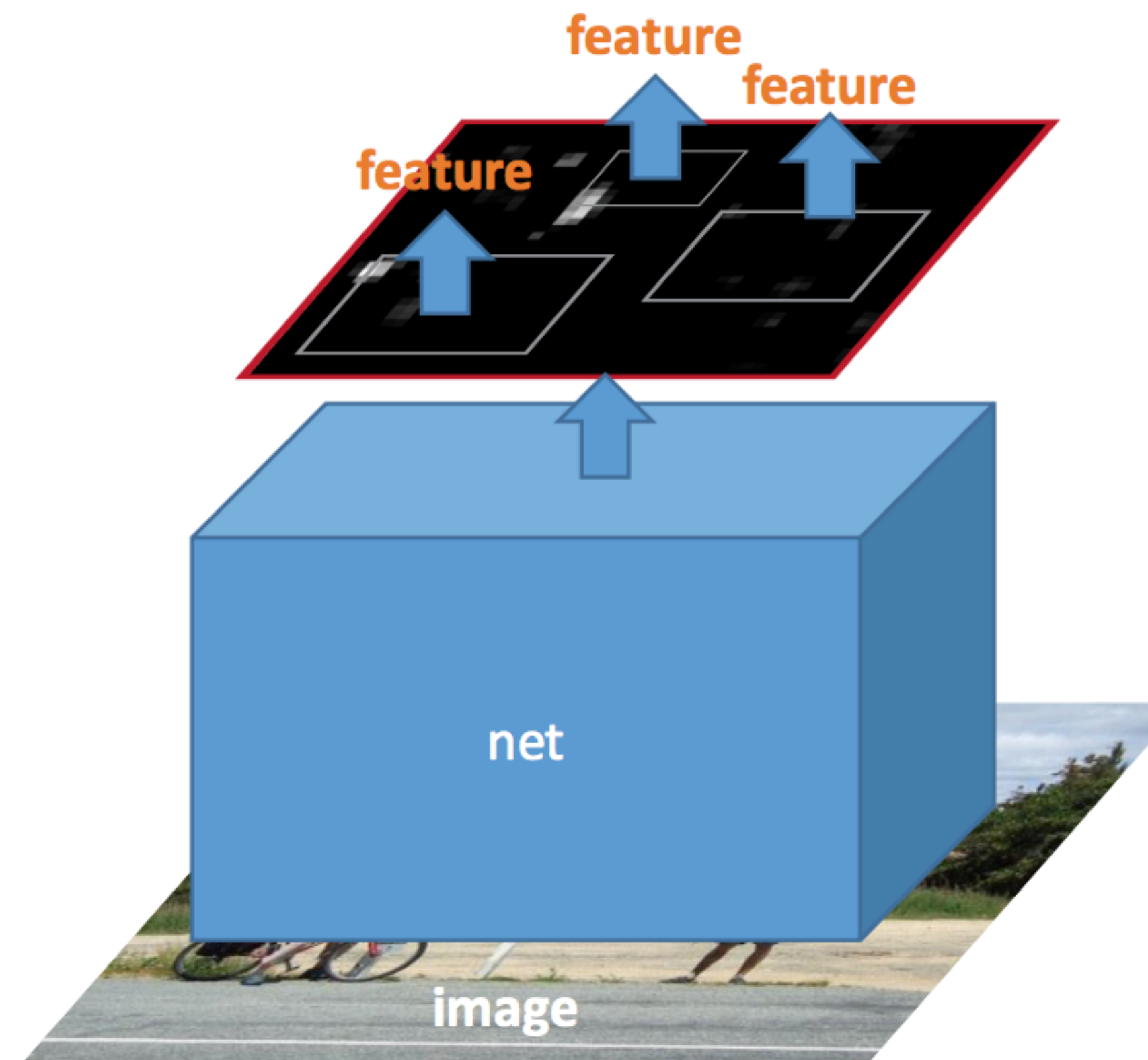
# R-CNN vs. SPP

[ He et al, ECCV 2014 ]



**R-CNN**

2000 nets on image regions



**SPP-net**

**1 net on full image**

# Fast R-CNN

[ Girshick et al, ICCV 2015 ]



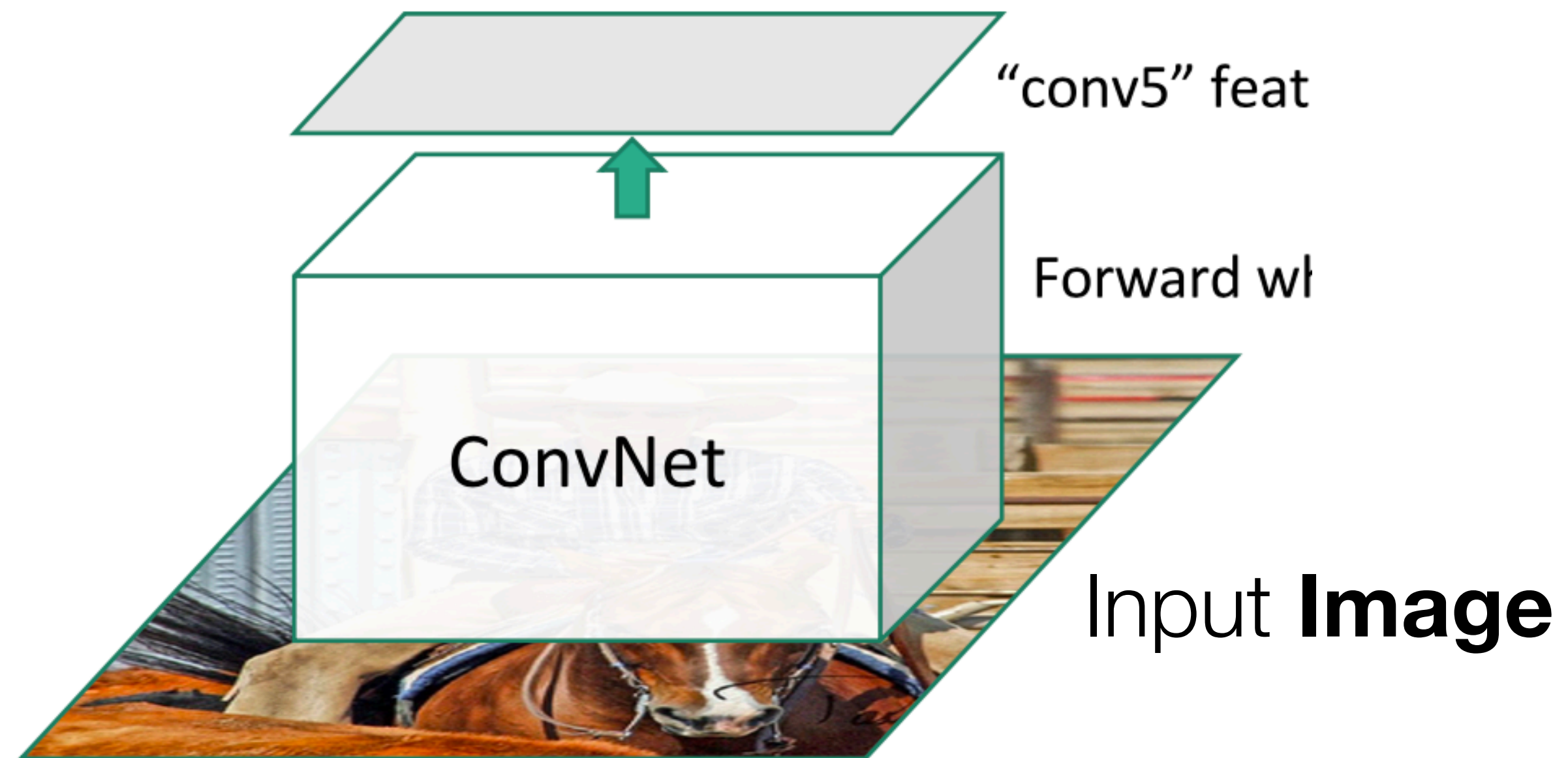
Input **Image**

\* image from Ross Girshick



# Fast R-CNN

[ Girshick et al, ICCV 2015 ]



\* image from Ross Girshick

# Fast R-CNN

[ Girshick et al, ICCV 2015 ]



\* image from Ross Girshick



# Fast R-CNN

[ Girshick et al, ICCV 2015 ]

**Regions of Interest**  
from the  
proposal  
method

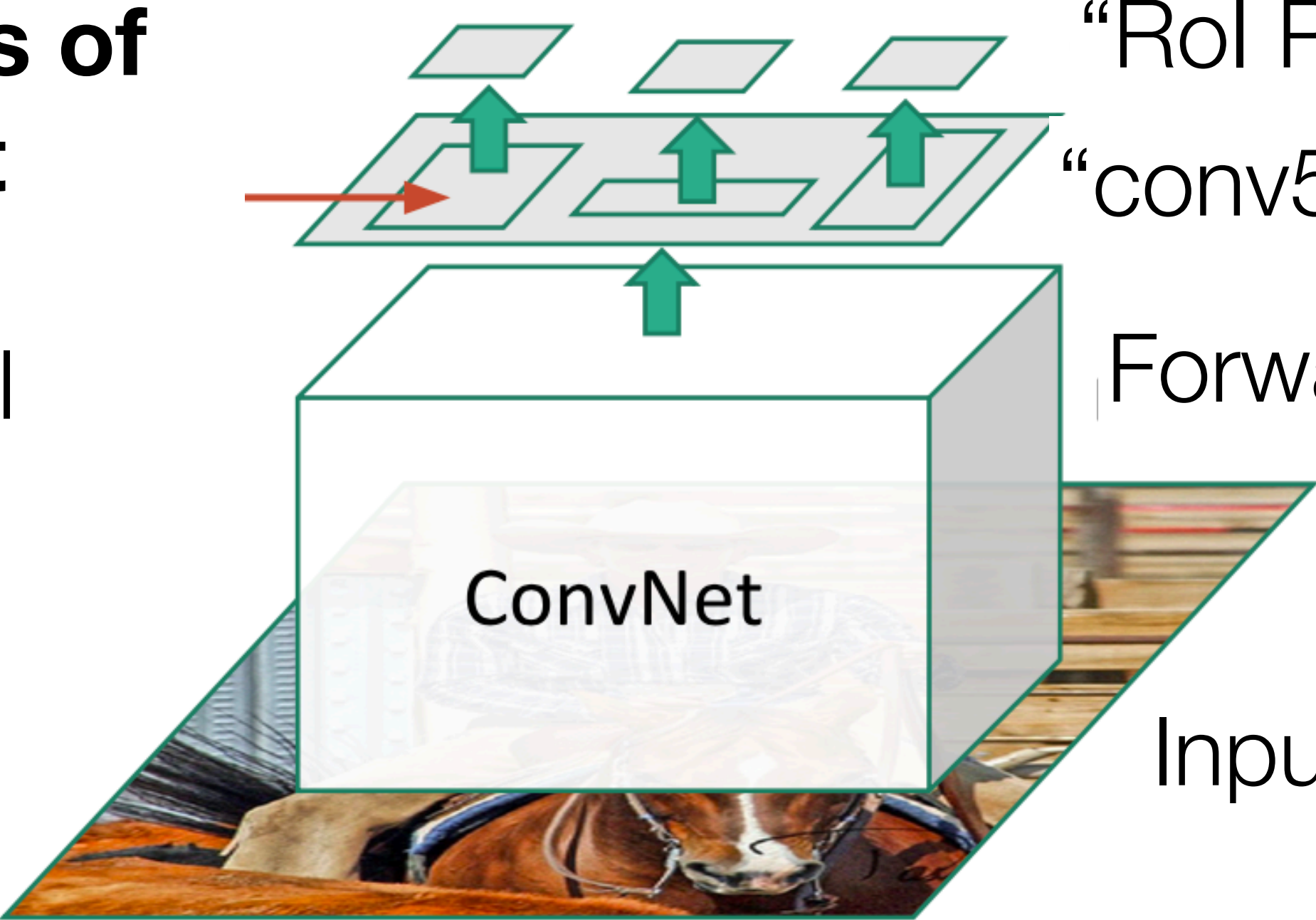


\* image from Ross Girshick

# Fast R-CNN

[ Girshick et al, ICCV 2015 ]

**Regions of Interest**  
from the  
proposal  
method



"RoI Pooling" layer

"conv5" feature map

Forward prop the **whole image** through CNN

Input **Image**

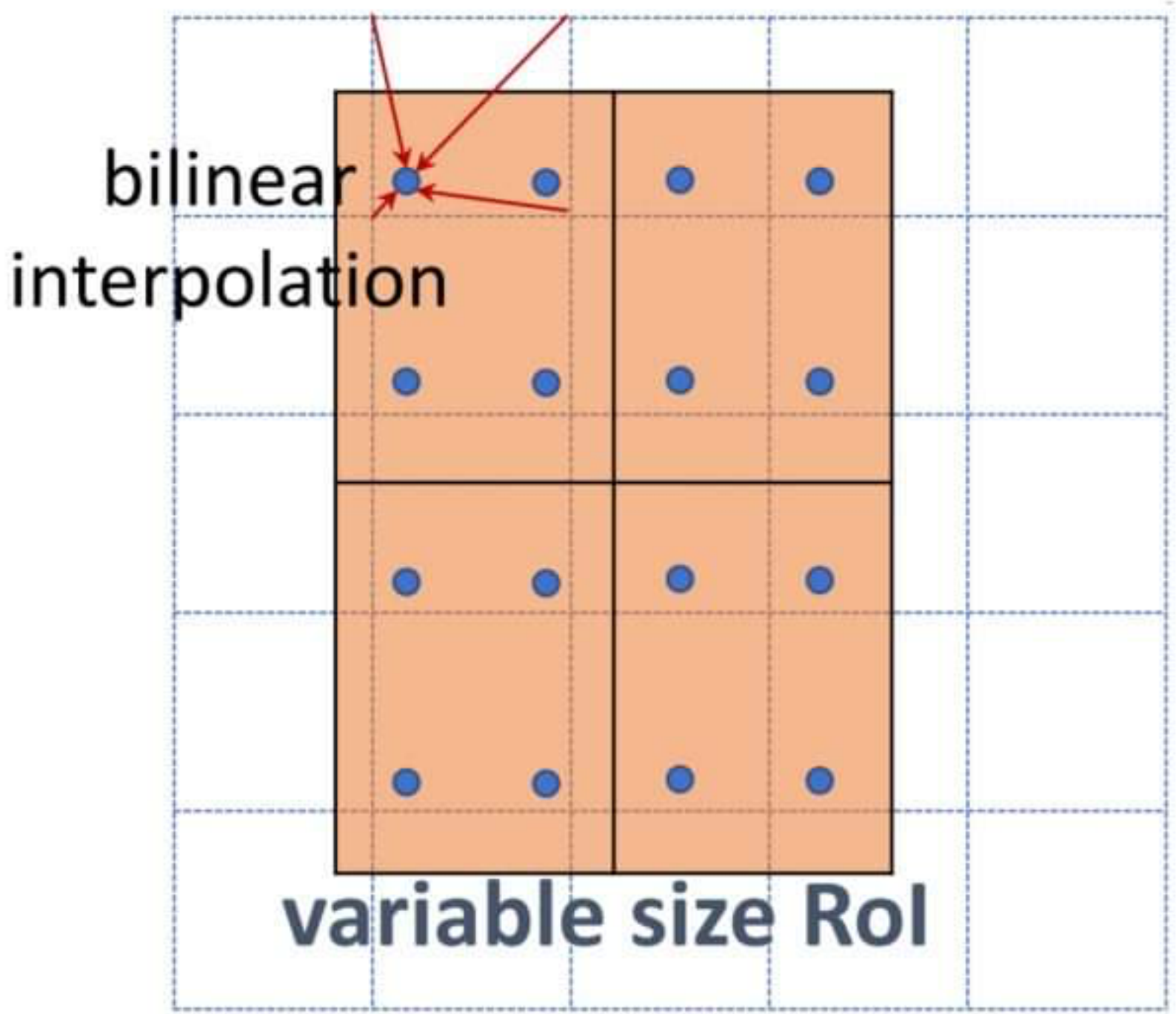
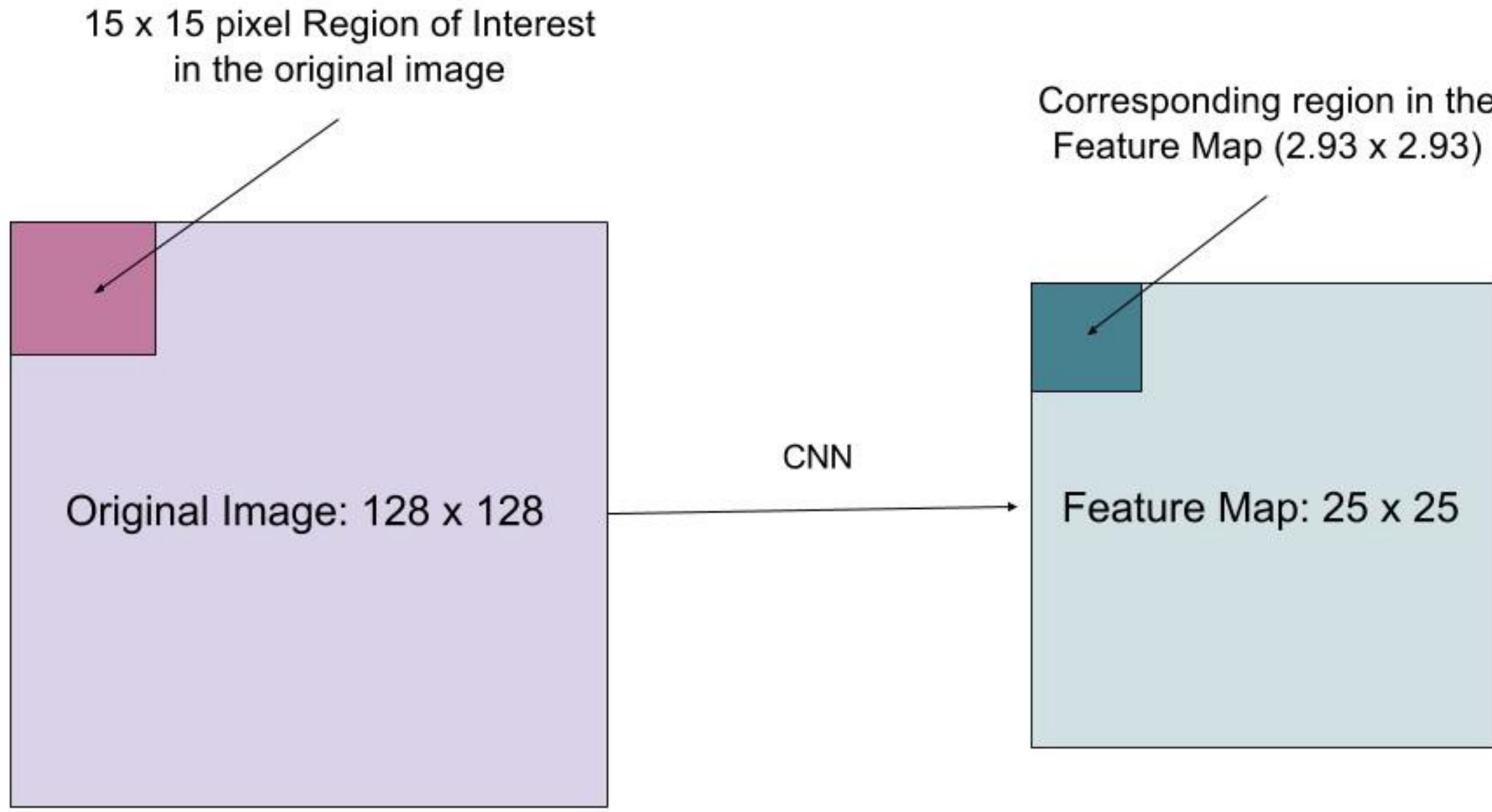
ConvNet

Girshick, "Fast R-C  
Figure copyright R

\* image from Ross Girshick



# RoI Align

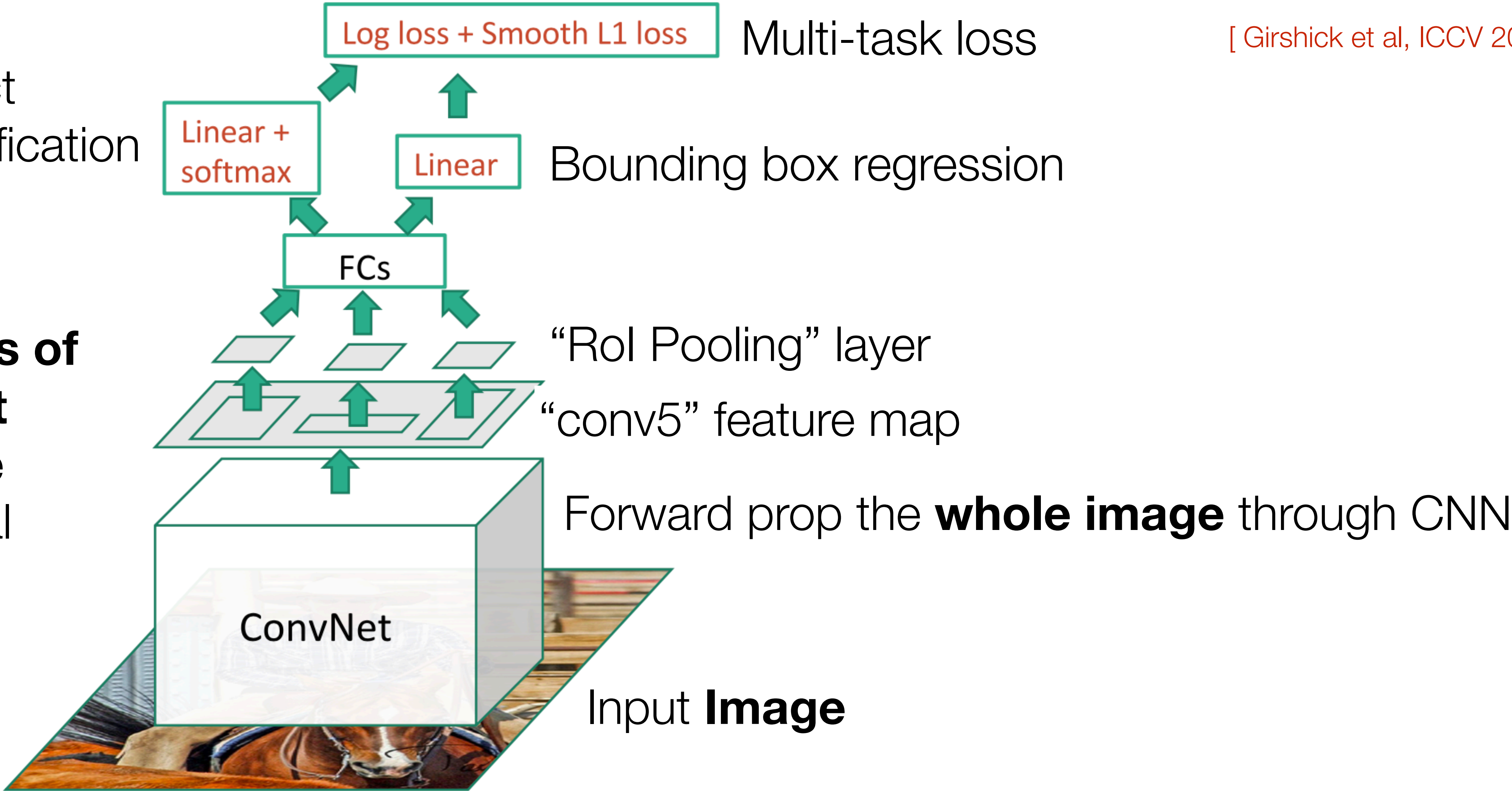


# Fast R-CNN

[ Girshick et al, ICCV 2015 ]

Object classification

**Regions of Interest** from the proposal method

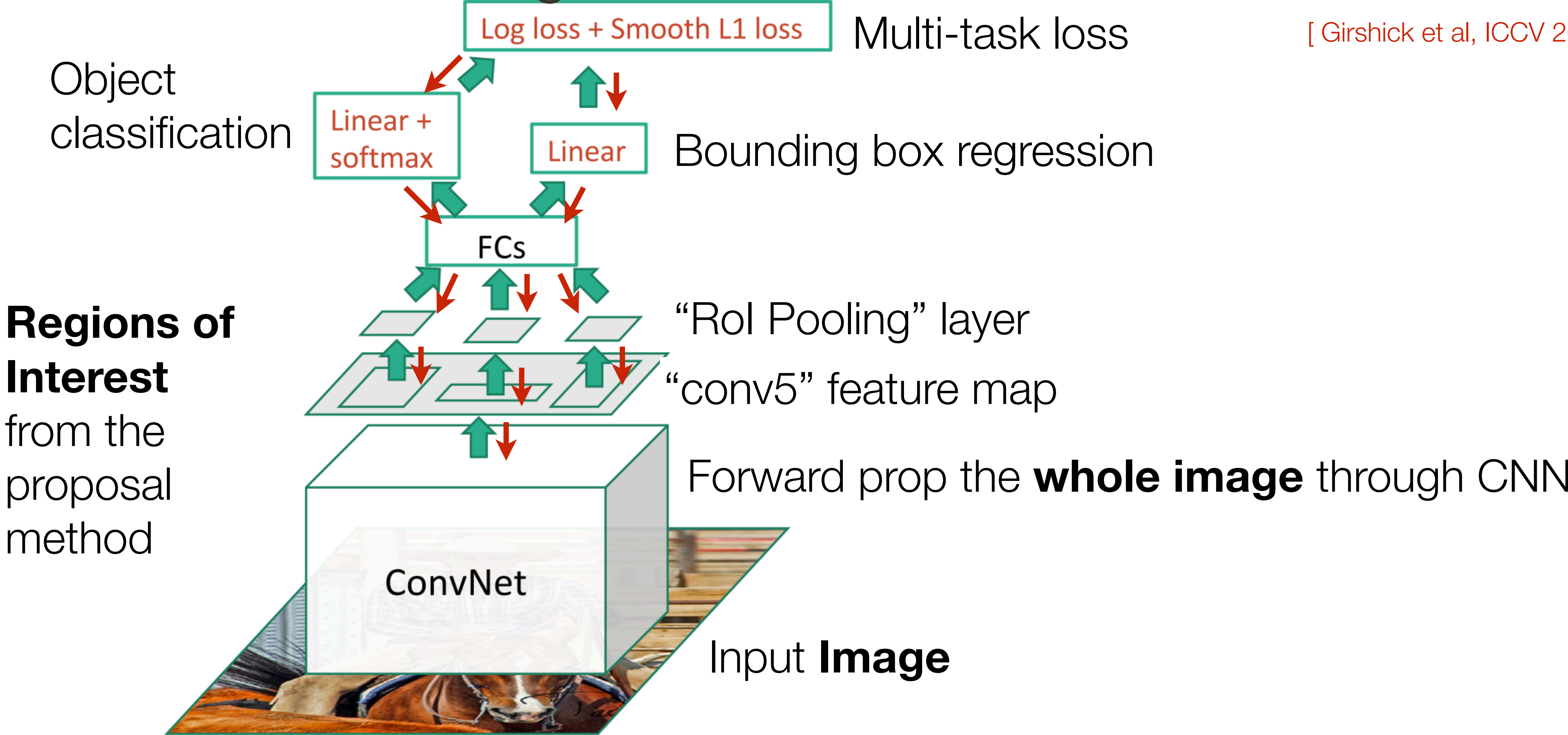


\* image from Ross Girshick



# Fast R-CNN: Training

[ Girshick et al, ICCV 2015 ]



\* image from Ross Girshick

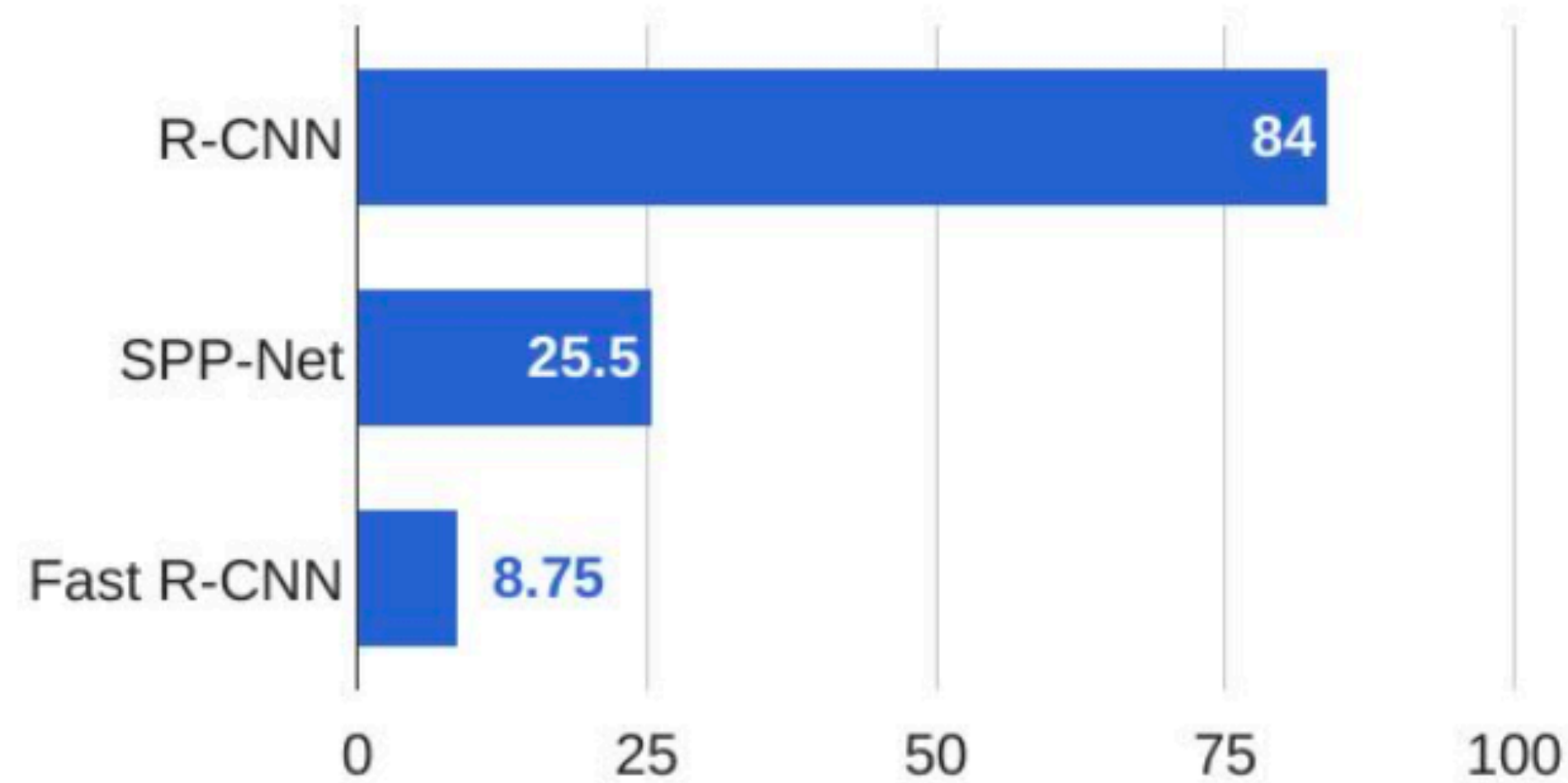
# R-CNN vs. SPP vs. Fast R-CNN

[ Girshick et al, CVPR 2014 ]

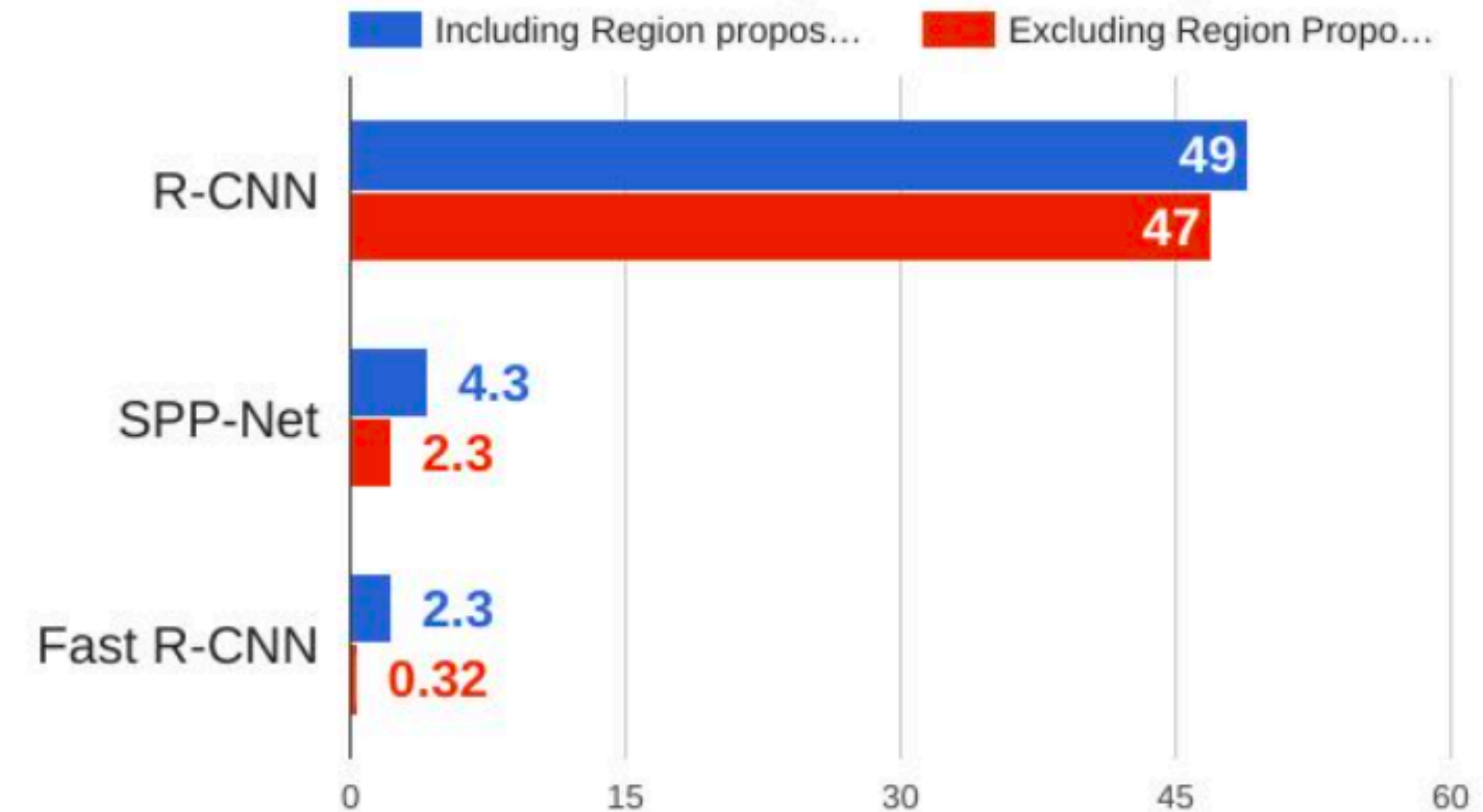
[ Girshick et al, ICCV 2015 ]

[ He et al, ECCV 2014 ]

## Training time (Hours)



## Test time (seconds)





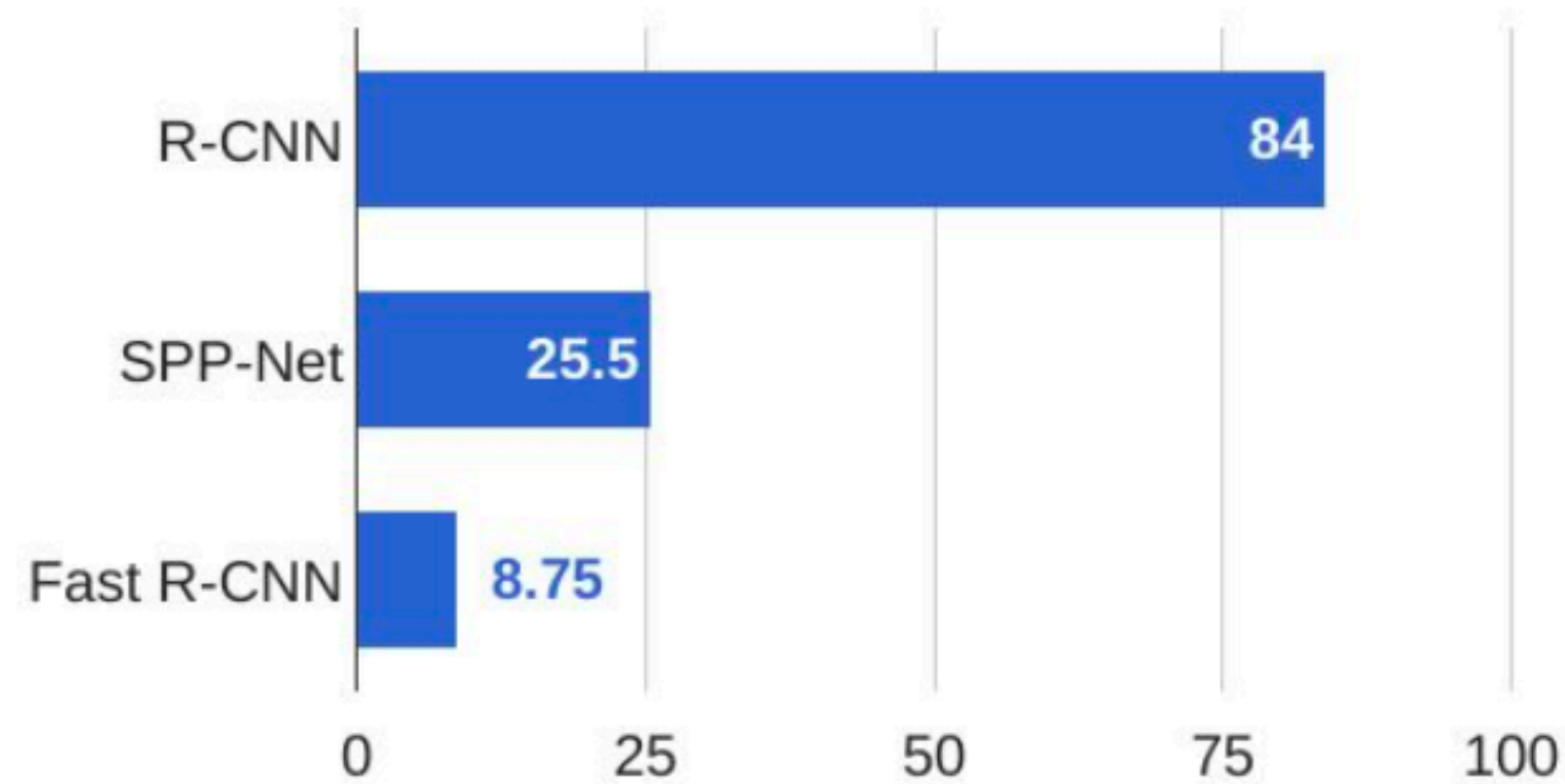
# R-CNN vs. SPP vs. Fast R-CNN

[ Girshick et al, CVPR 2014 ]

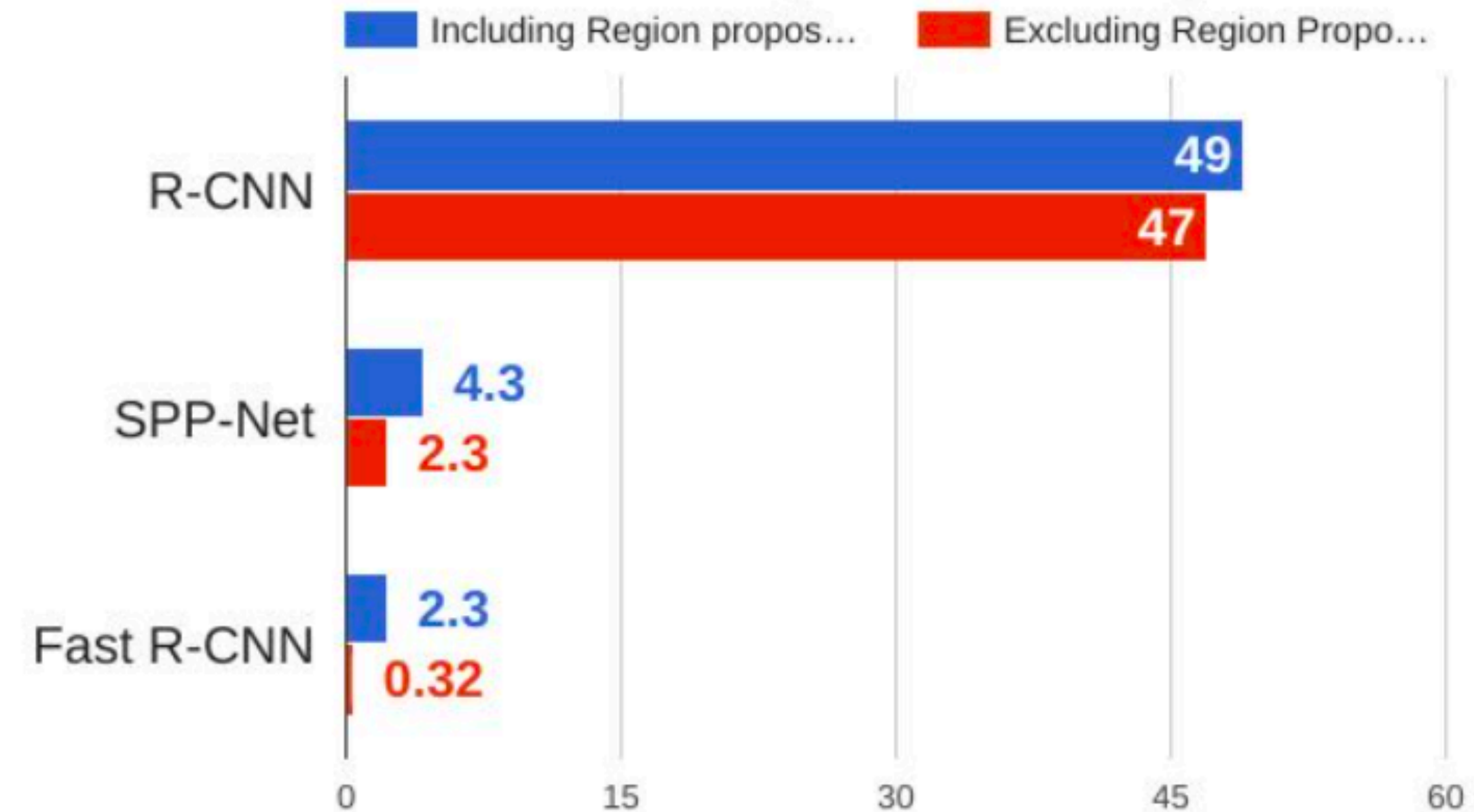
[ Girshick et al, ICCV 2015 ]

[ He et al, ECCV 2014 ]

## Training time (Hours)



## Test time (seconds)



**Observation:** Performance dominated by the region proposals at this point!

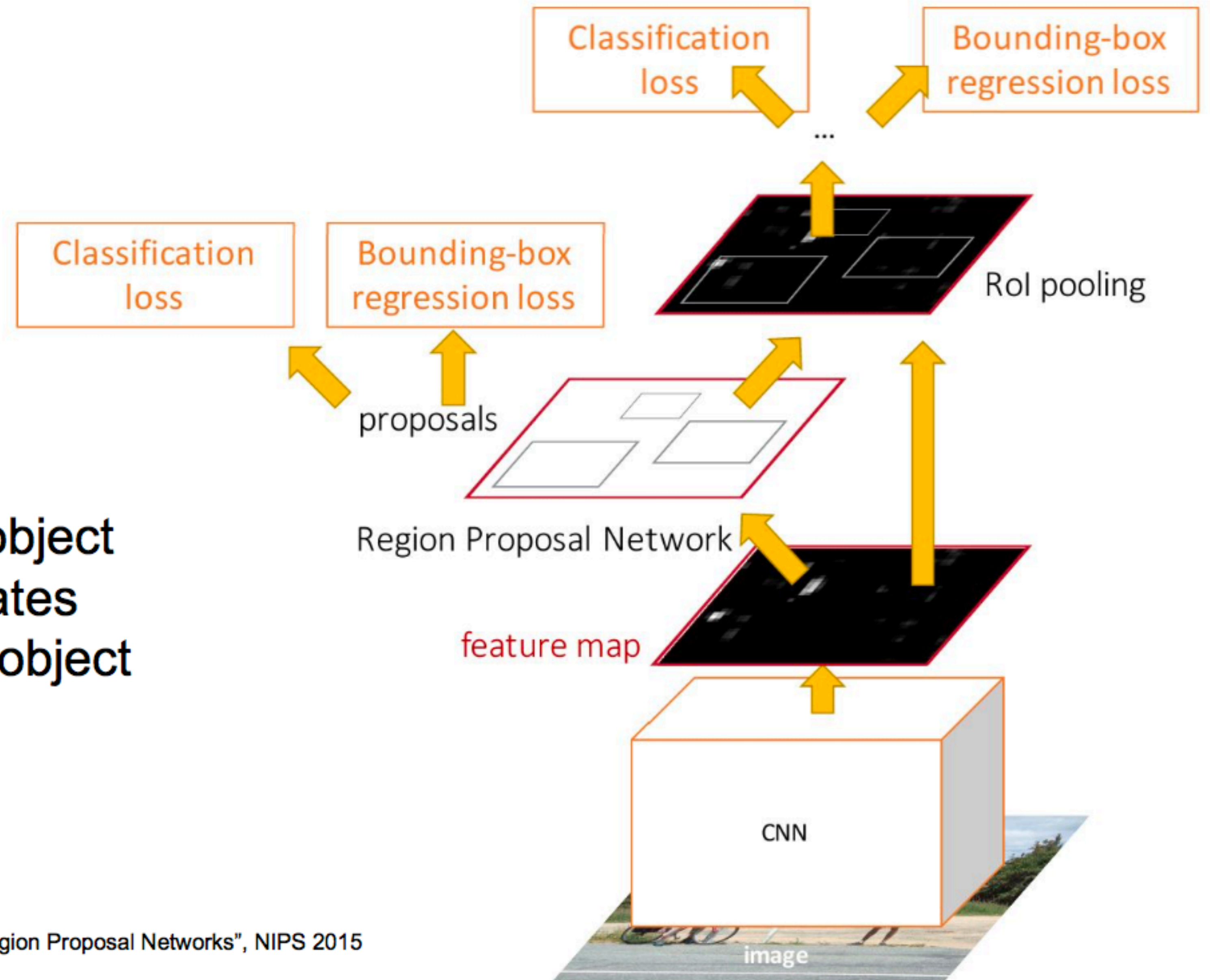
# Faster R-CNN

Make CNN do proposals!

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

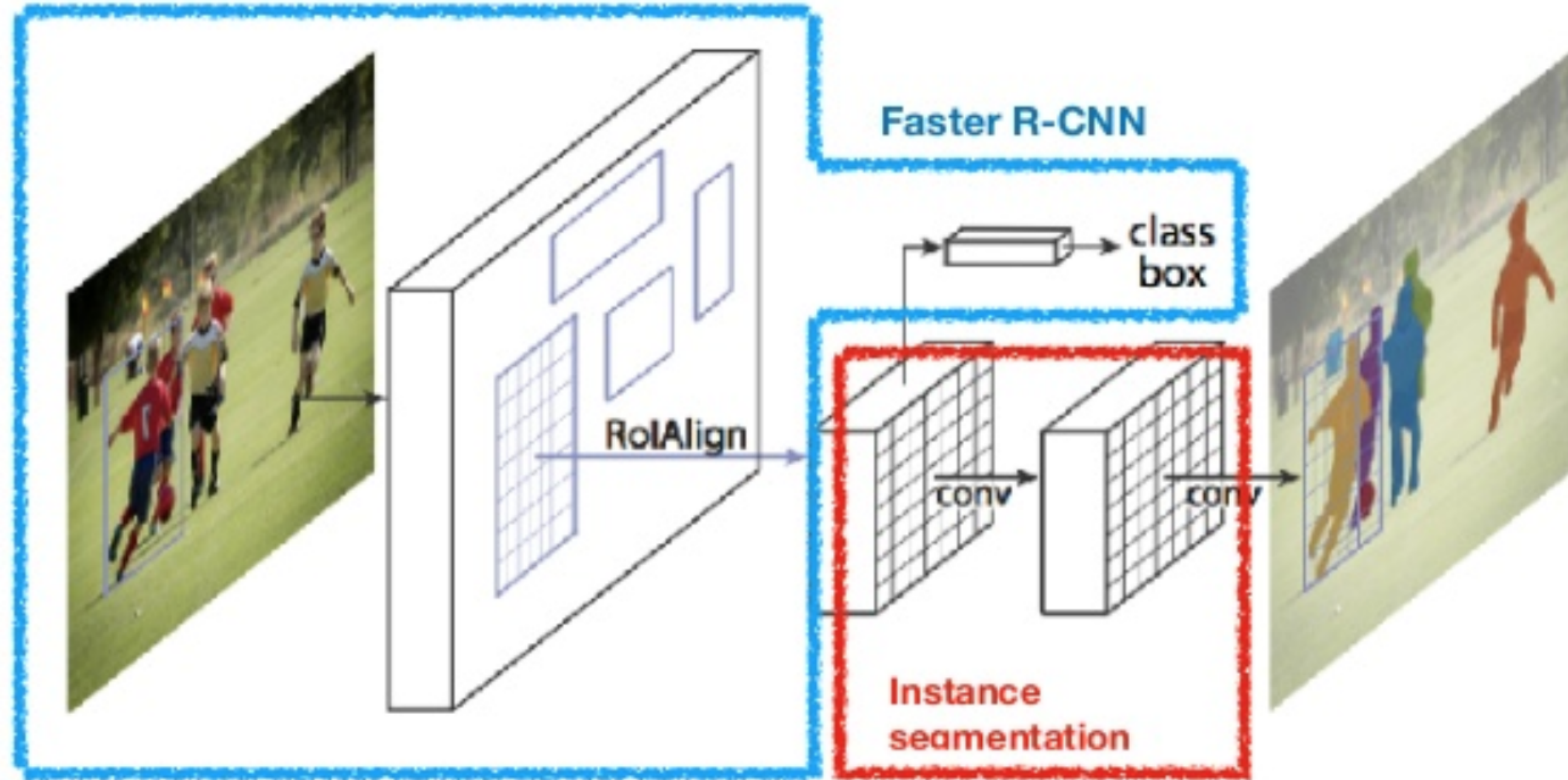
Jointly train with 4 losses:

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



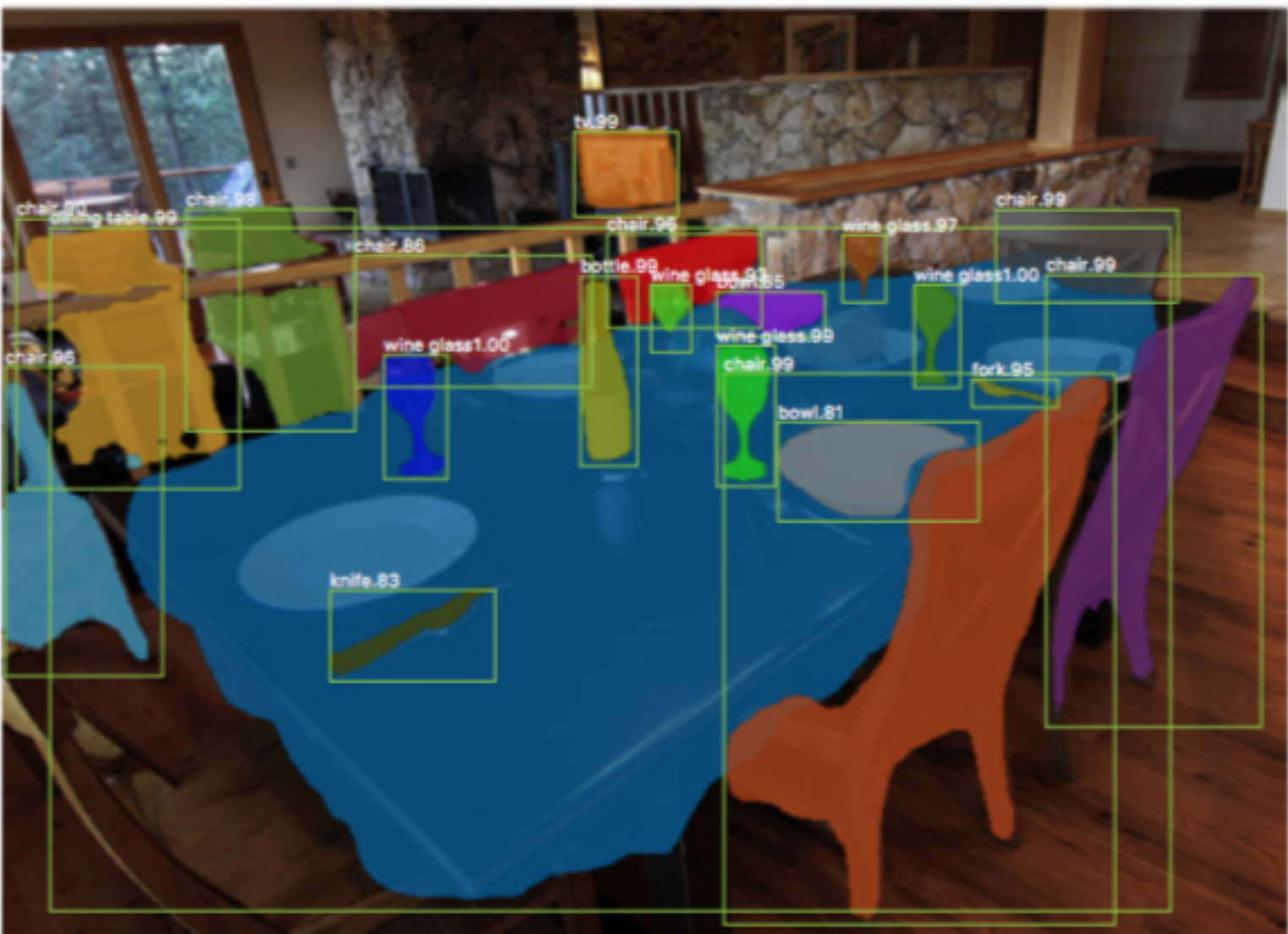
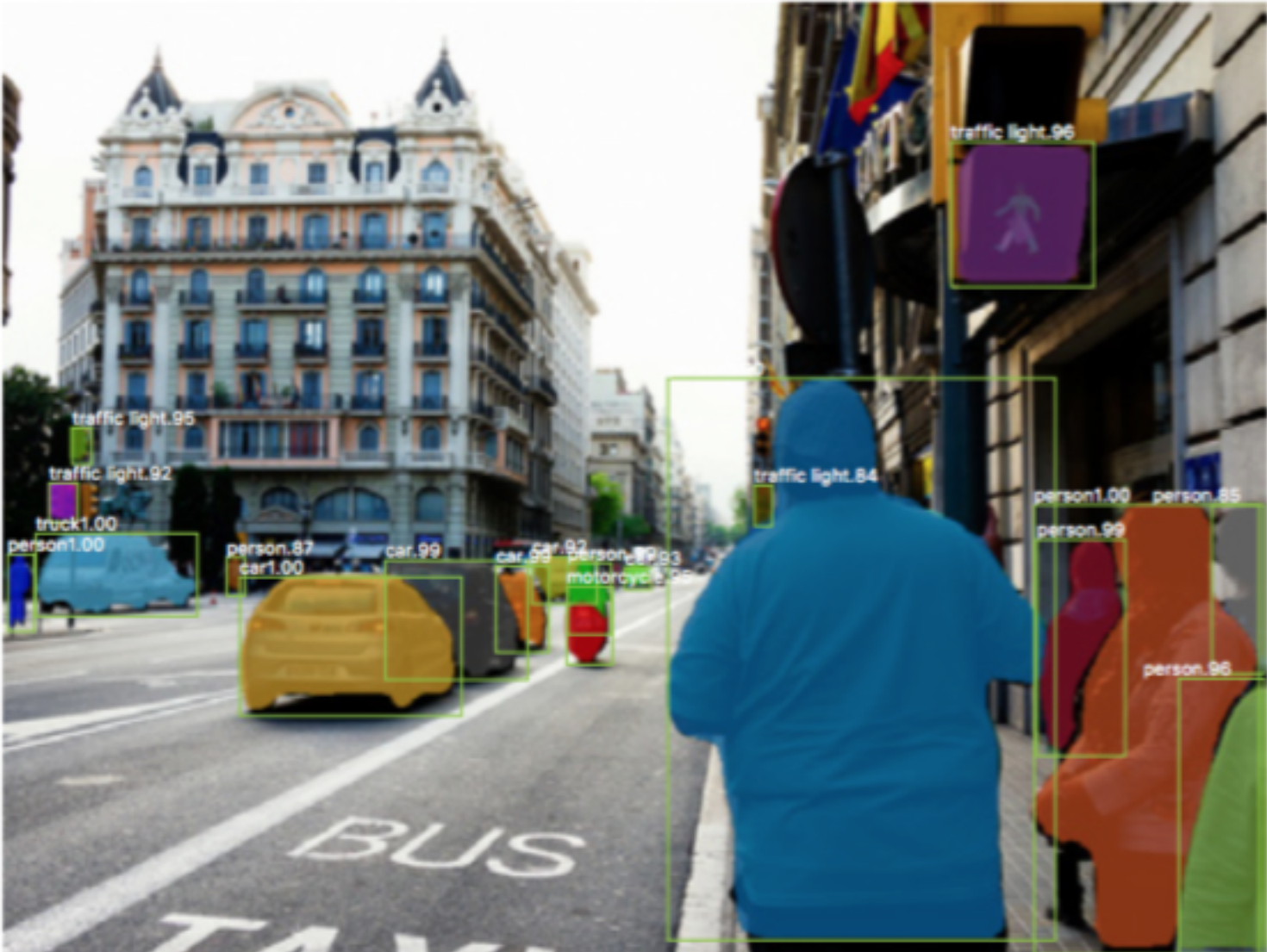
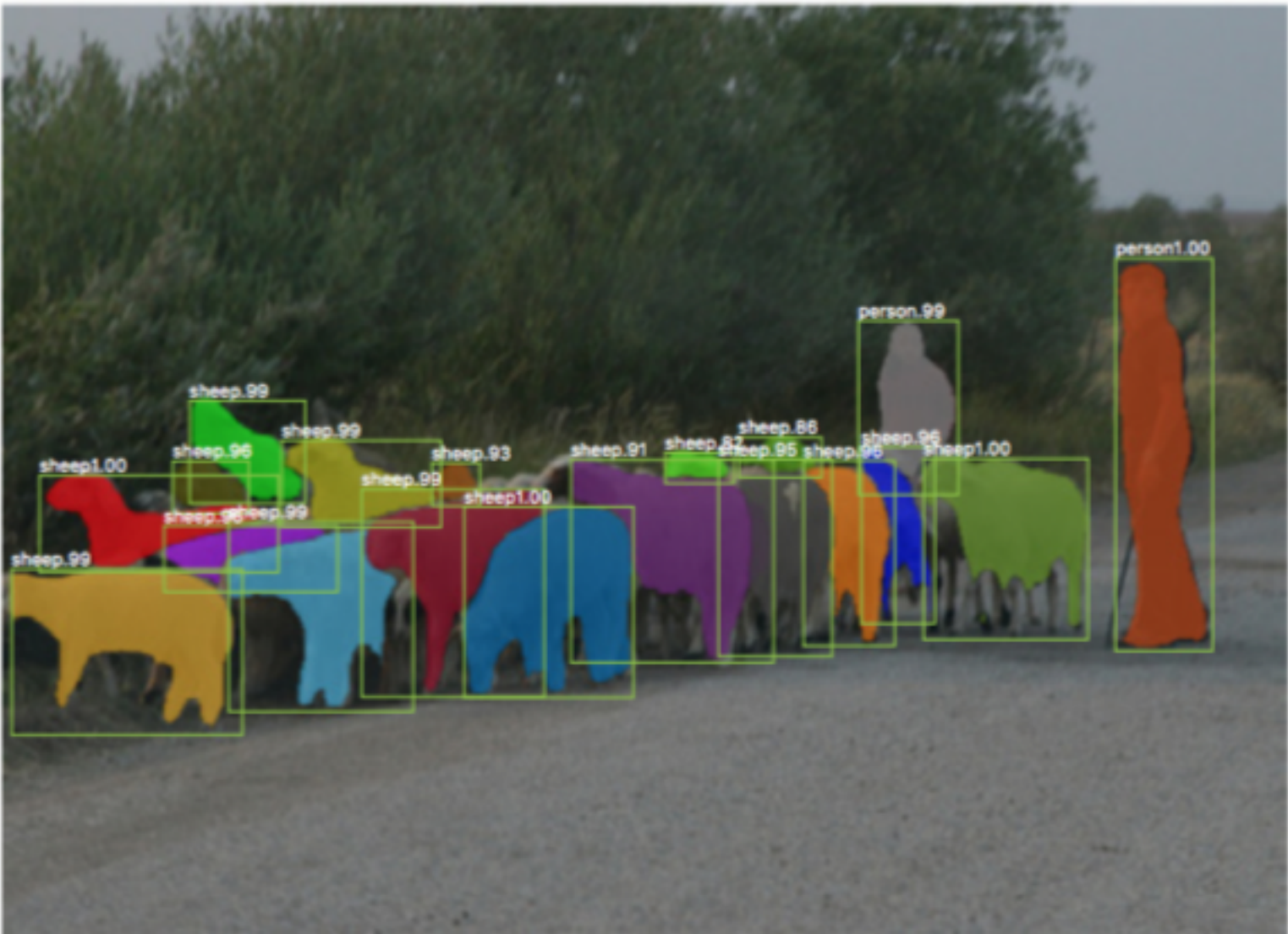


# Mask R-CNN





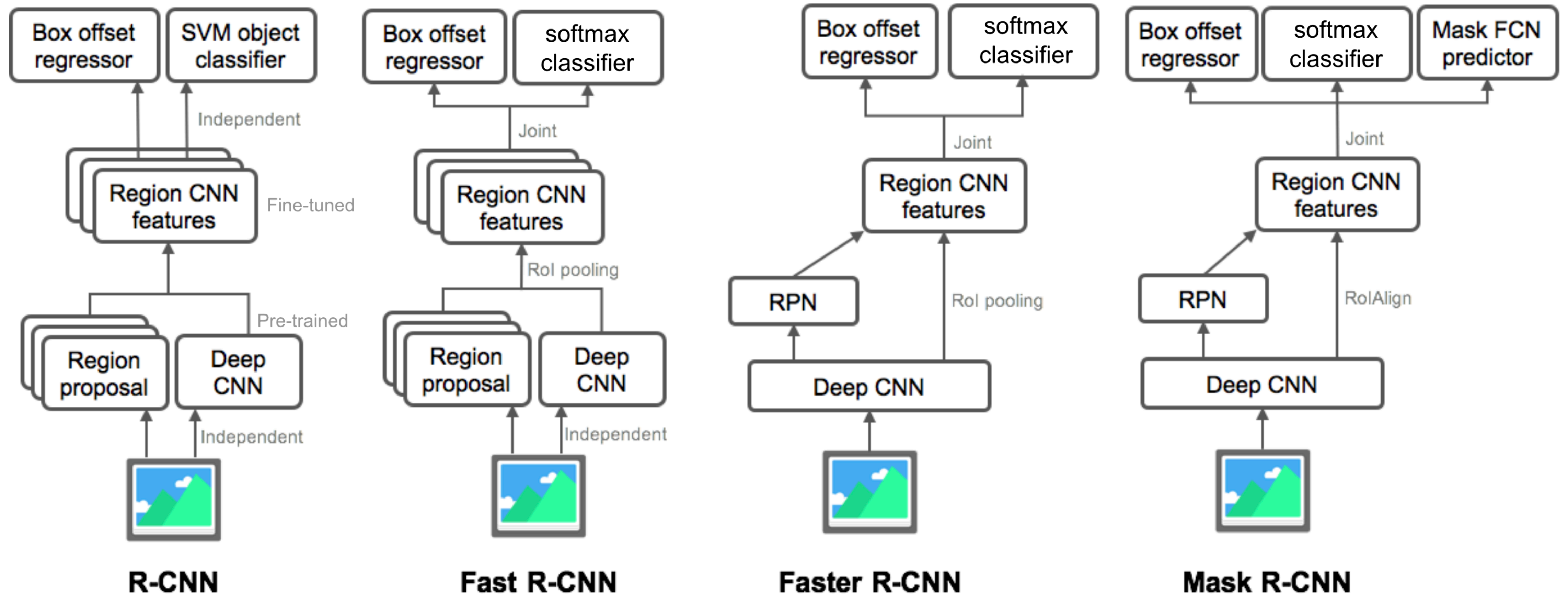
# Mask R-CNN



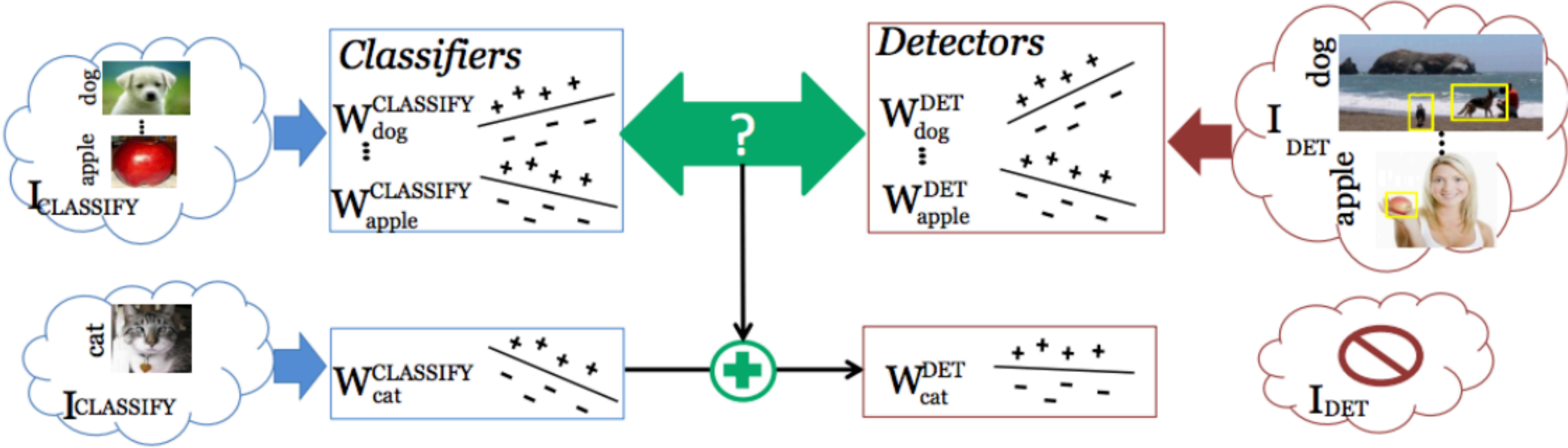
[ He et al, 2017 ]



# Summary of R-CNN Family of Models



# LSDA: Large Scale Detection through Adaptation



$$W_{cat}^{DETECT} = W_{cat}^{CLASSIFY} + \delta W_{cat}$$

[ Hoffman et al, NIPS 2014 ]

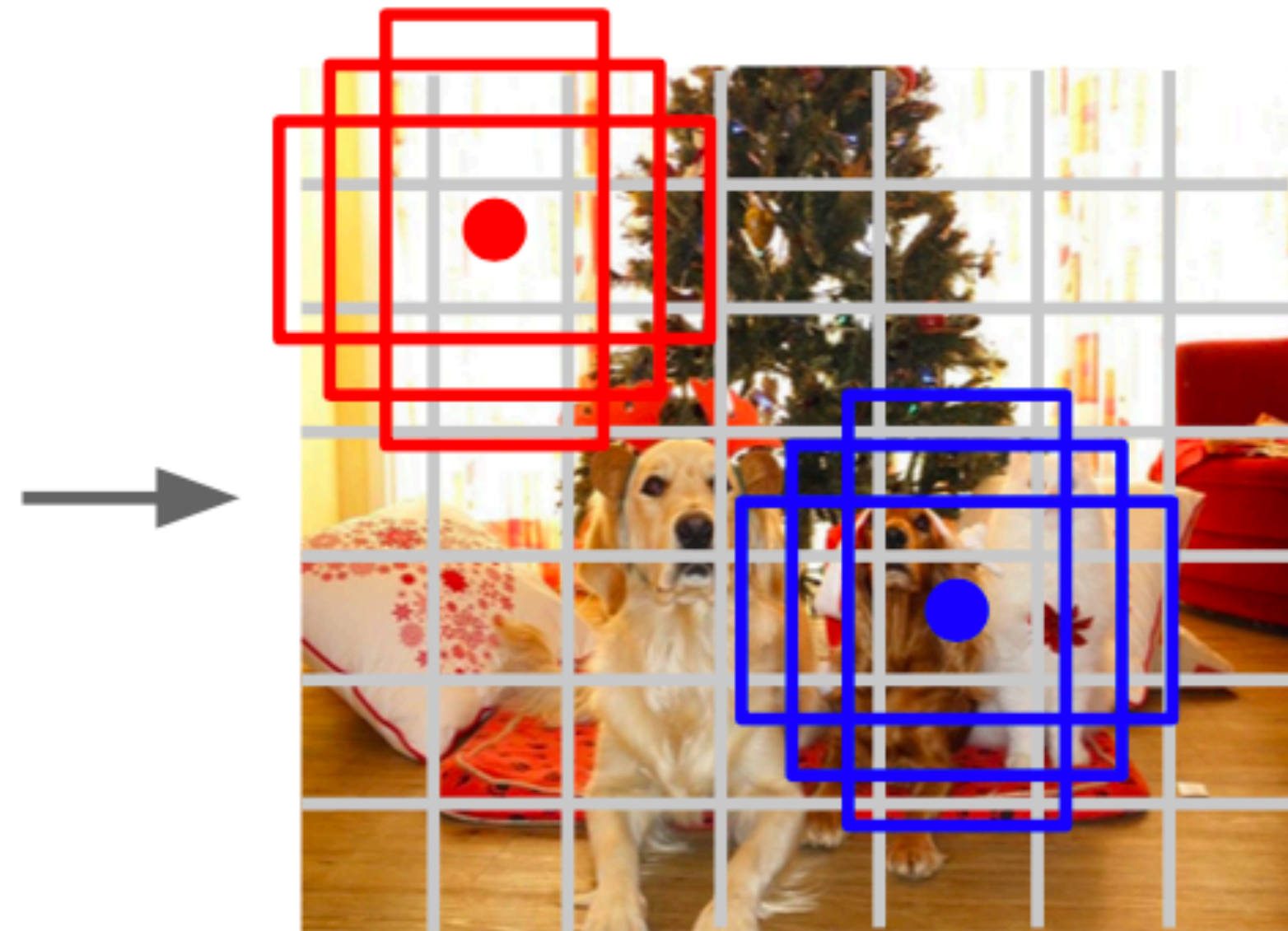


# YOLO: You Only Look Once

[ Redmon et al, CVPR 2016 ]



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
(dx, dy, dh, dw, confidence)
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

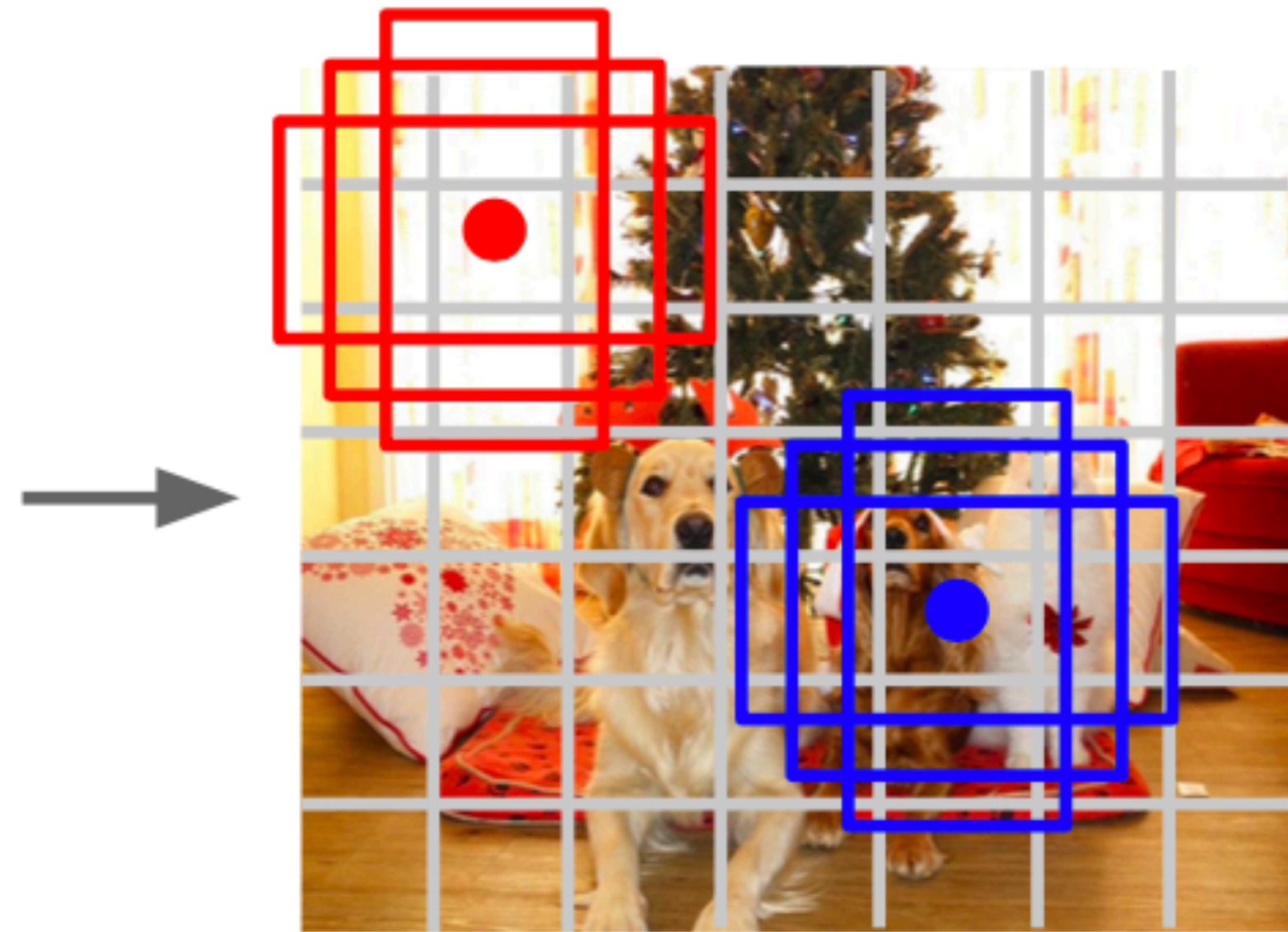


# YOLO: You Only Look Once

[ Redmon et al, CVPR 2016 ]

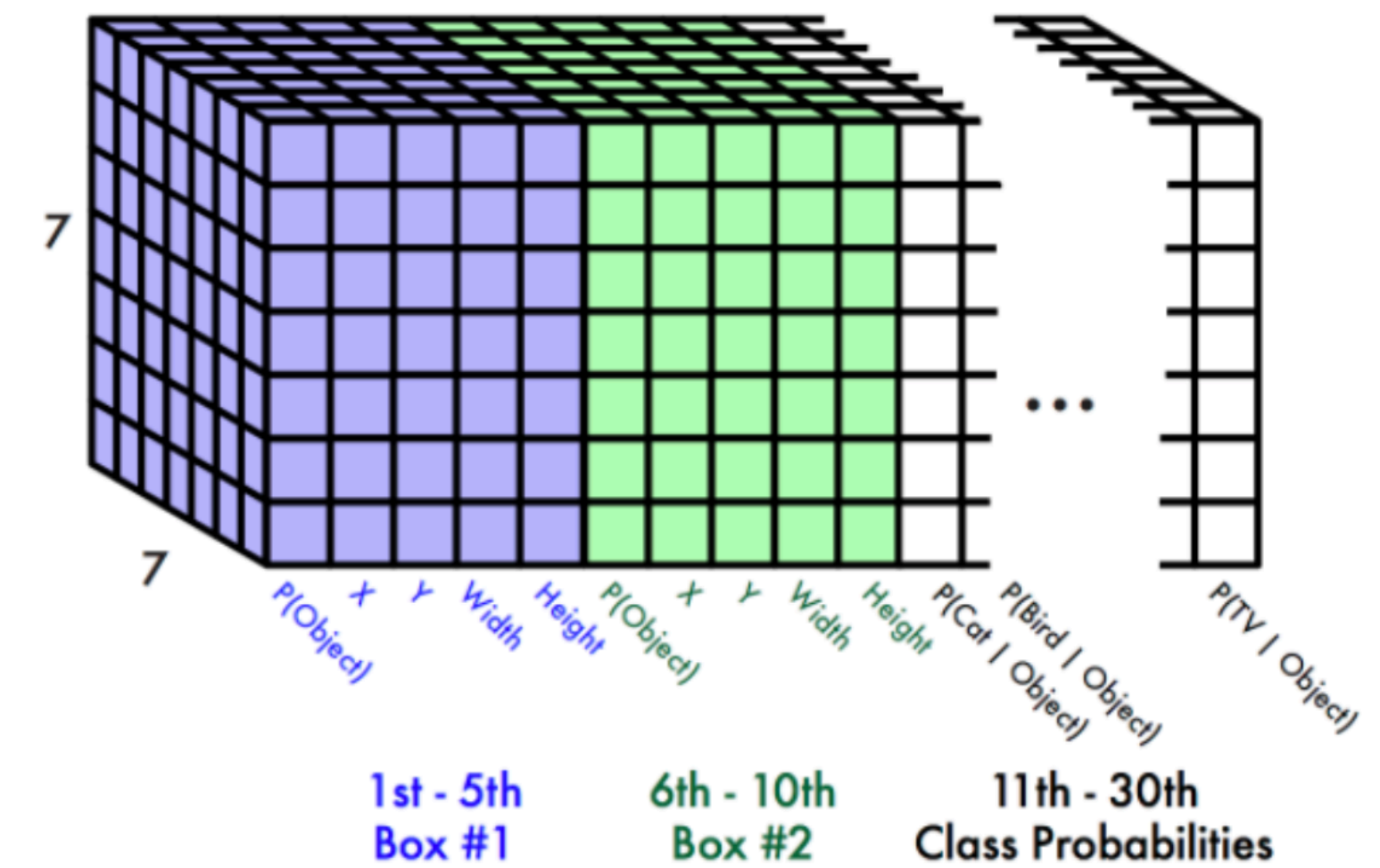


Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$



1st - 5th Box #1  
6th - 10th Box #2  
11th - 30th Class Probabilities





# YOLO v2

<http://pureddie.com/yolo>



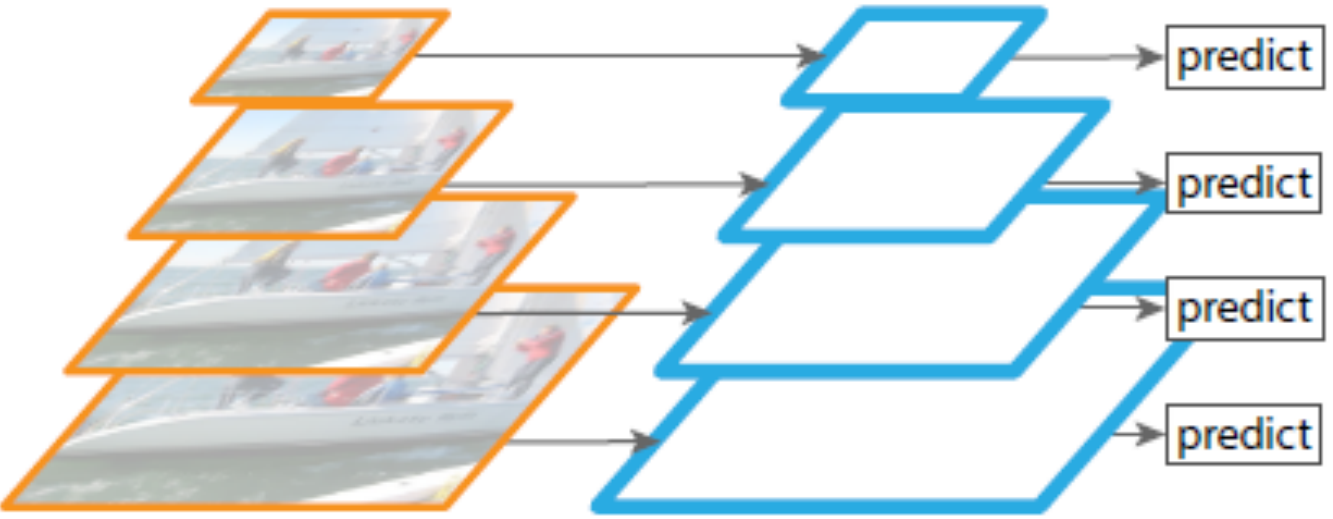


# YOLO v2

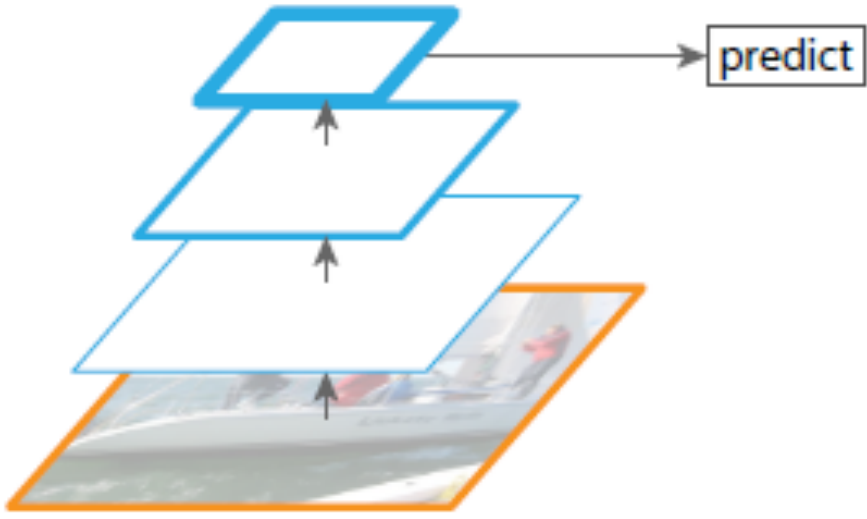
<http://pureddie.com/yolo>



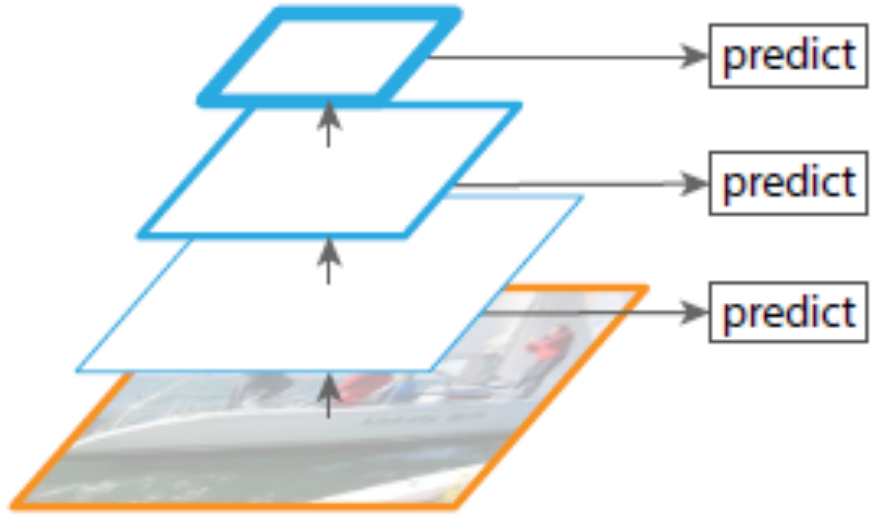
# Feature Pyramid Networks



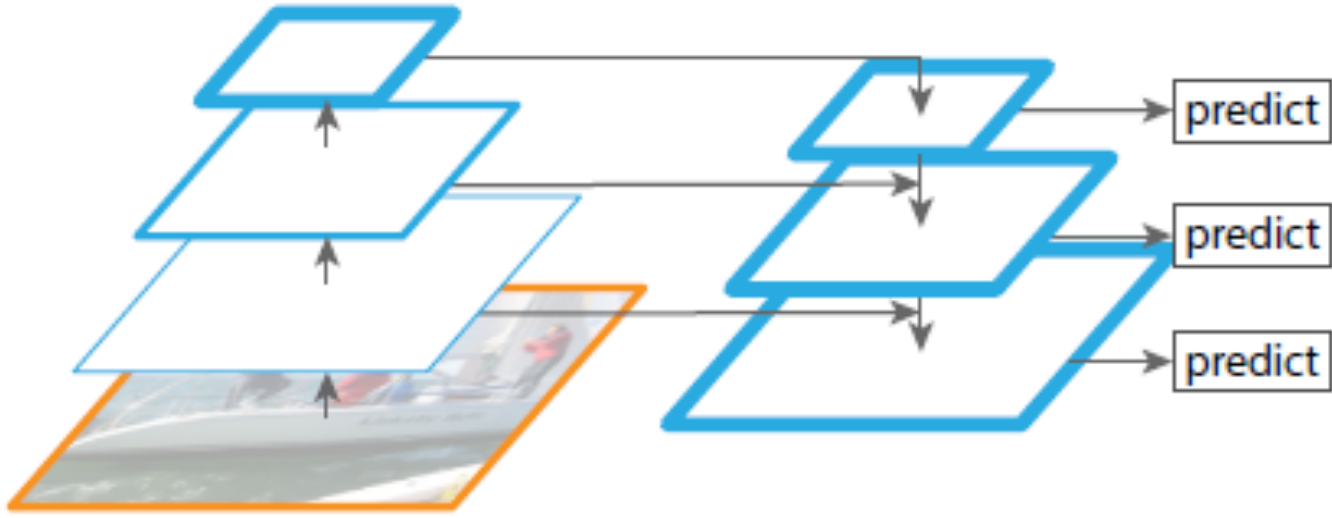
(a) Featurized image pyramid



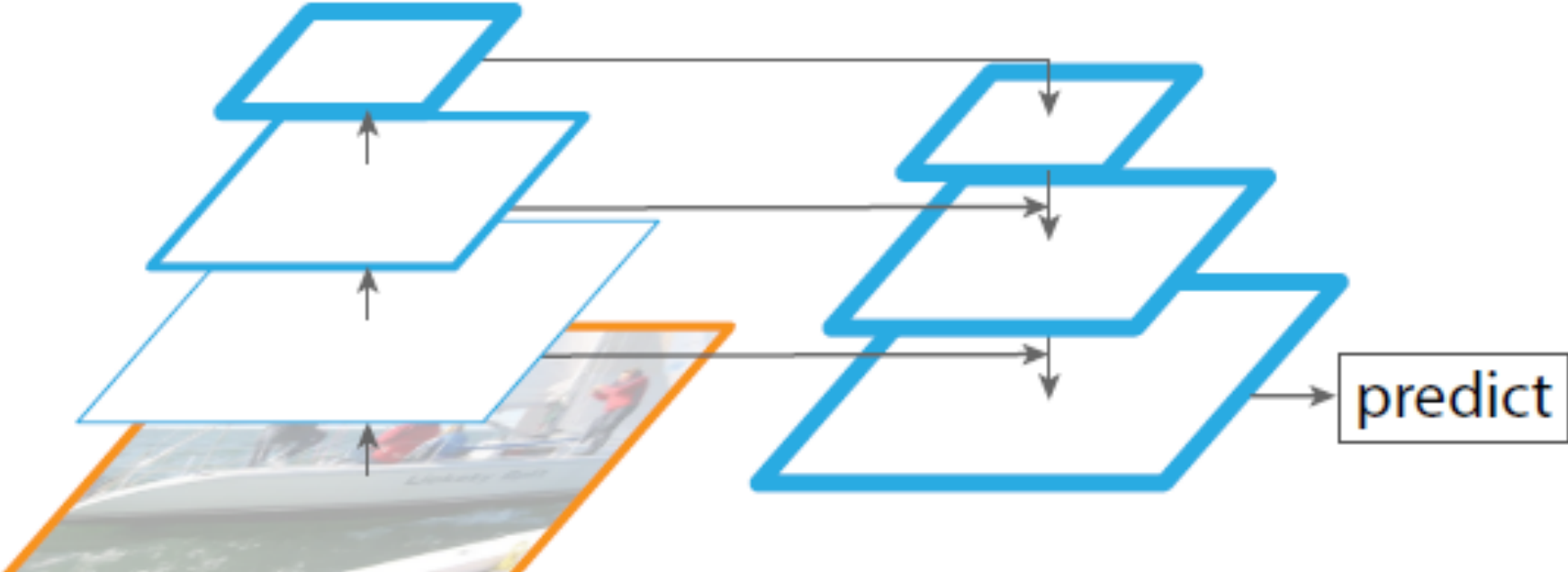
(b) Single feature map



(c) Pyramidal feature hierarchy



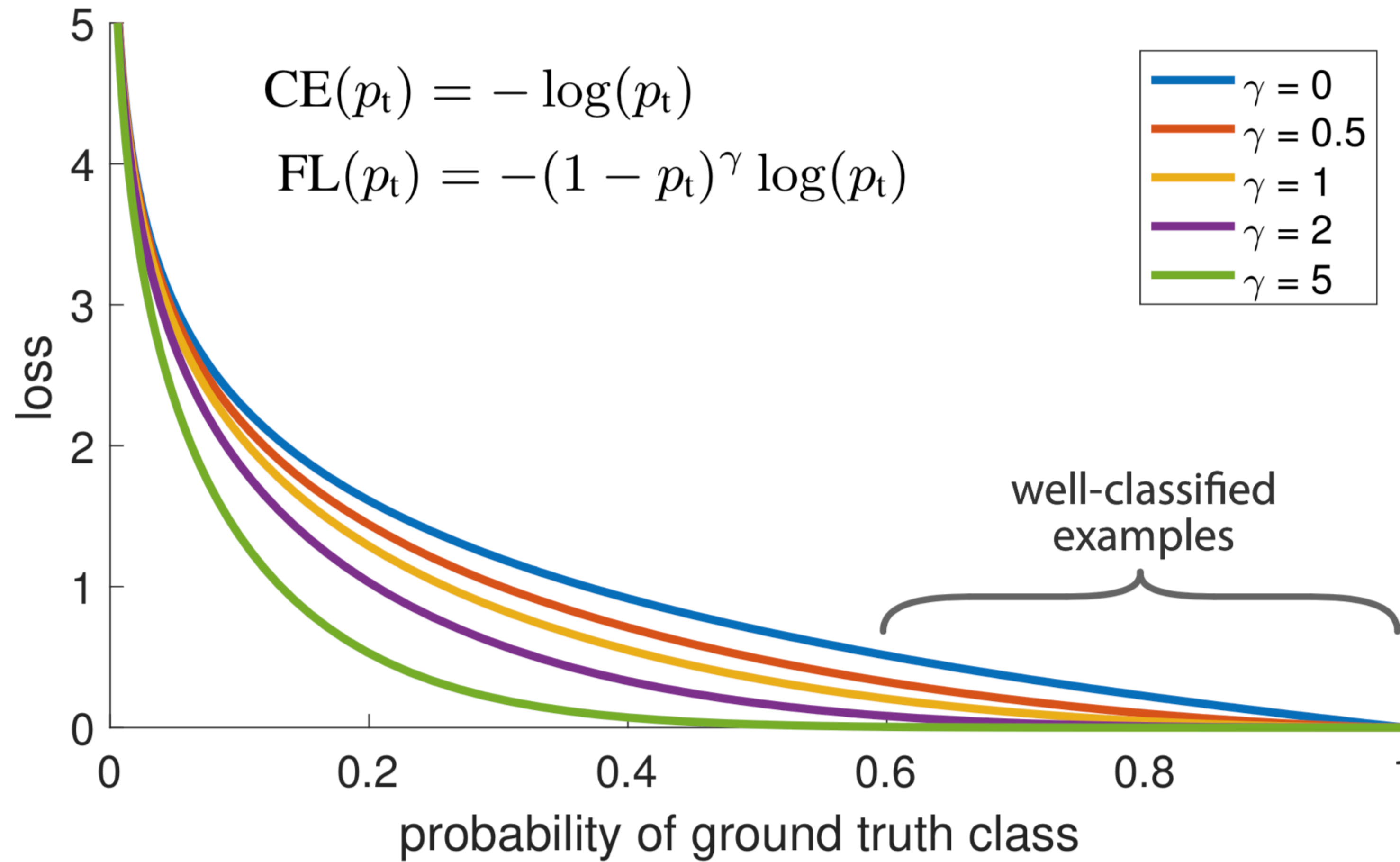
(d) Feature Pyramid Network



(e) Similar Structure with (d)

# Focal Loss

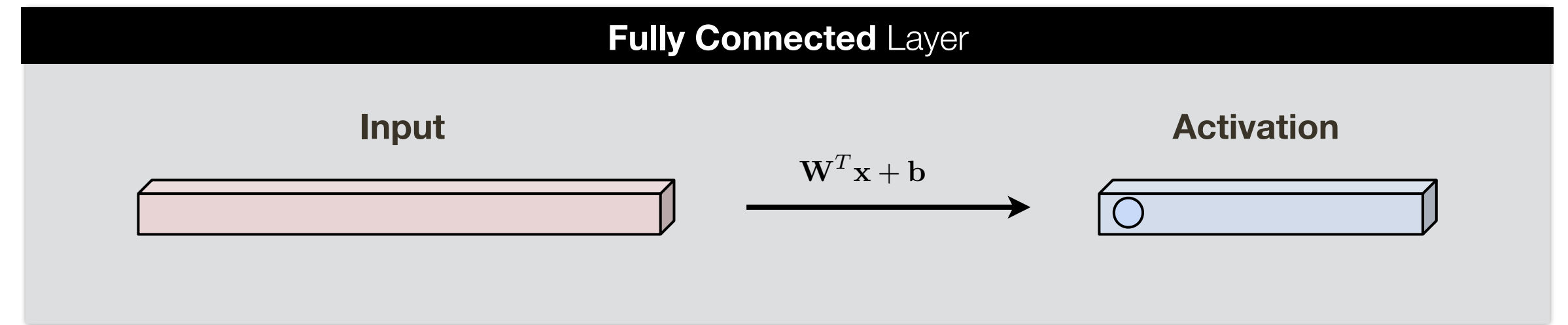
$$p_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{otherwise} \end{cases}$$





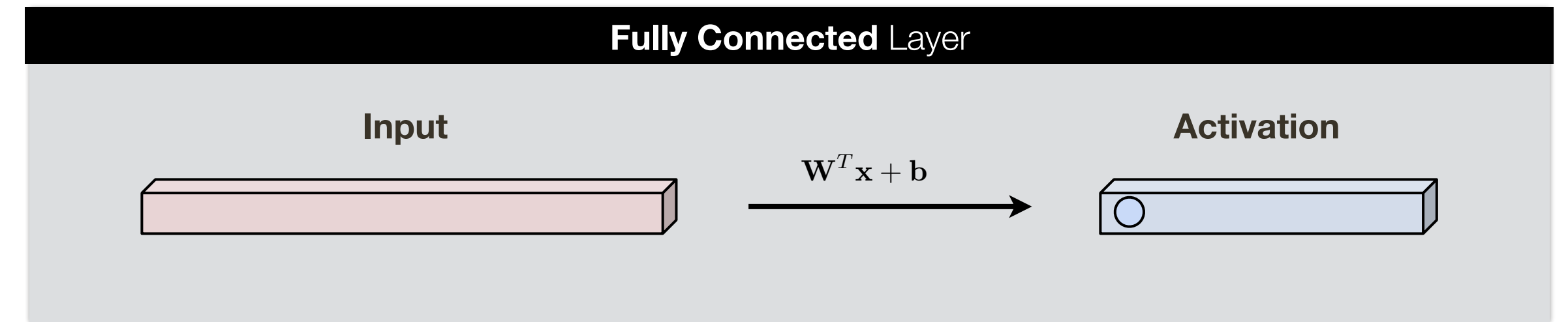
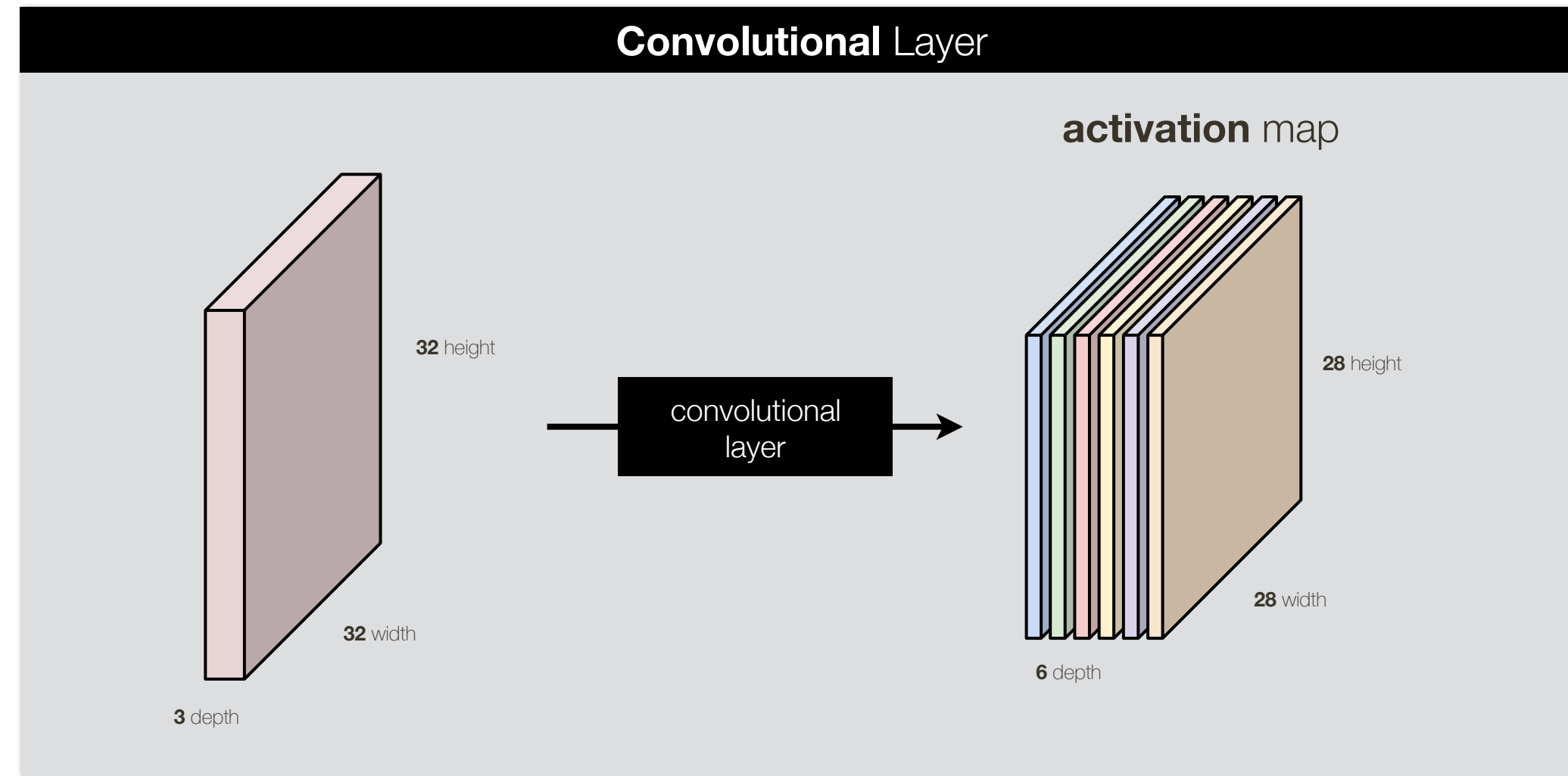
# Review of **CNNs**

# Review of CNNs

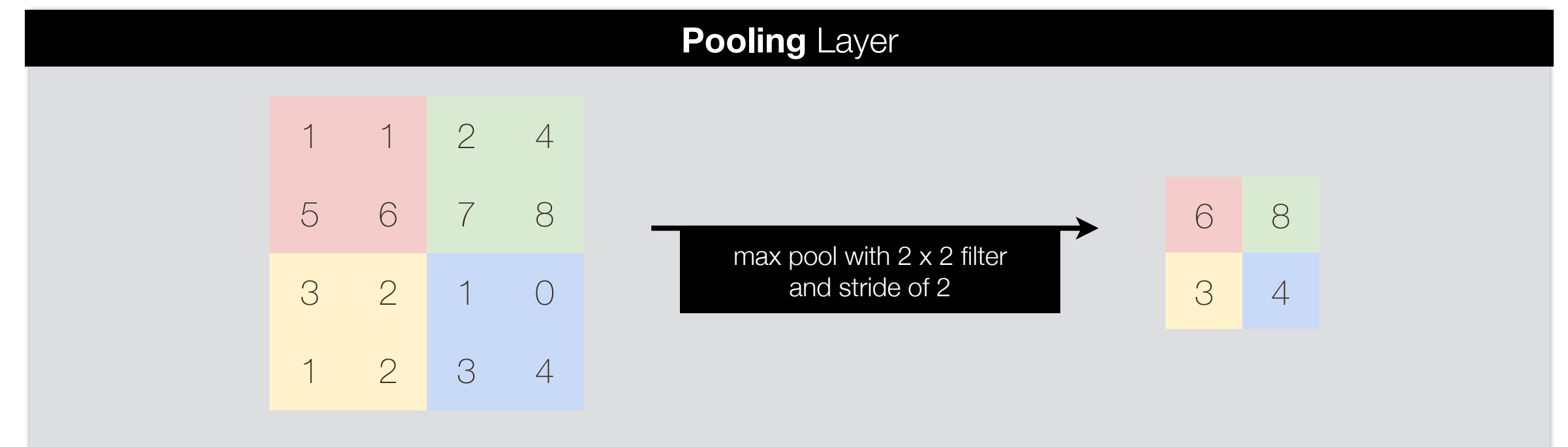
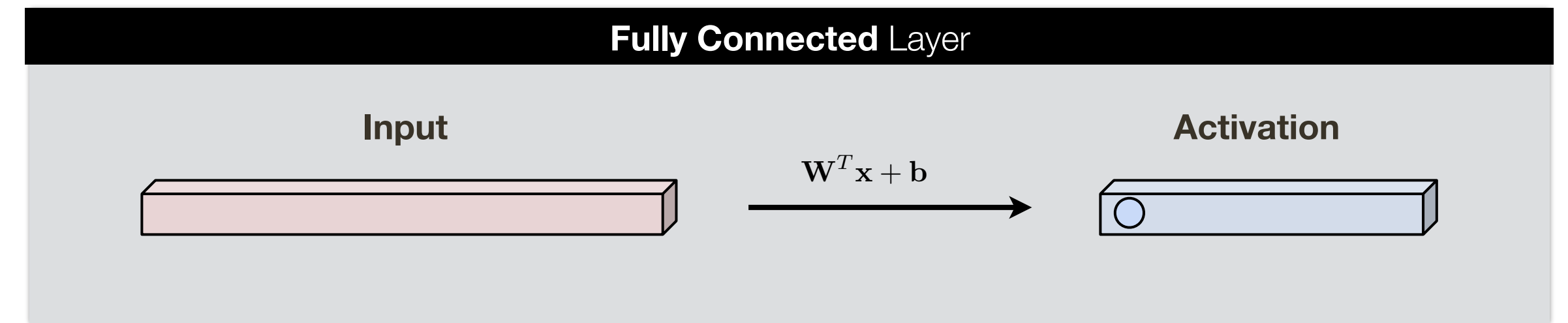
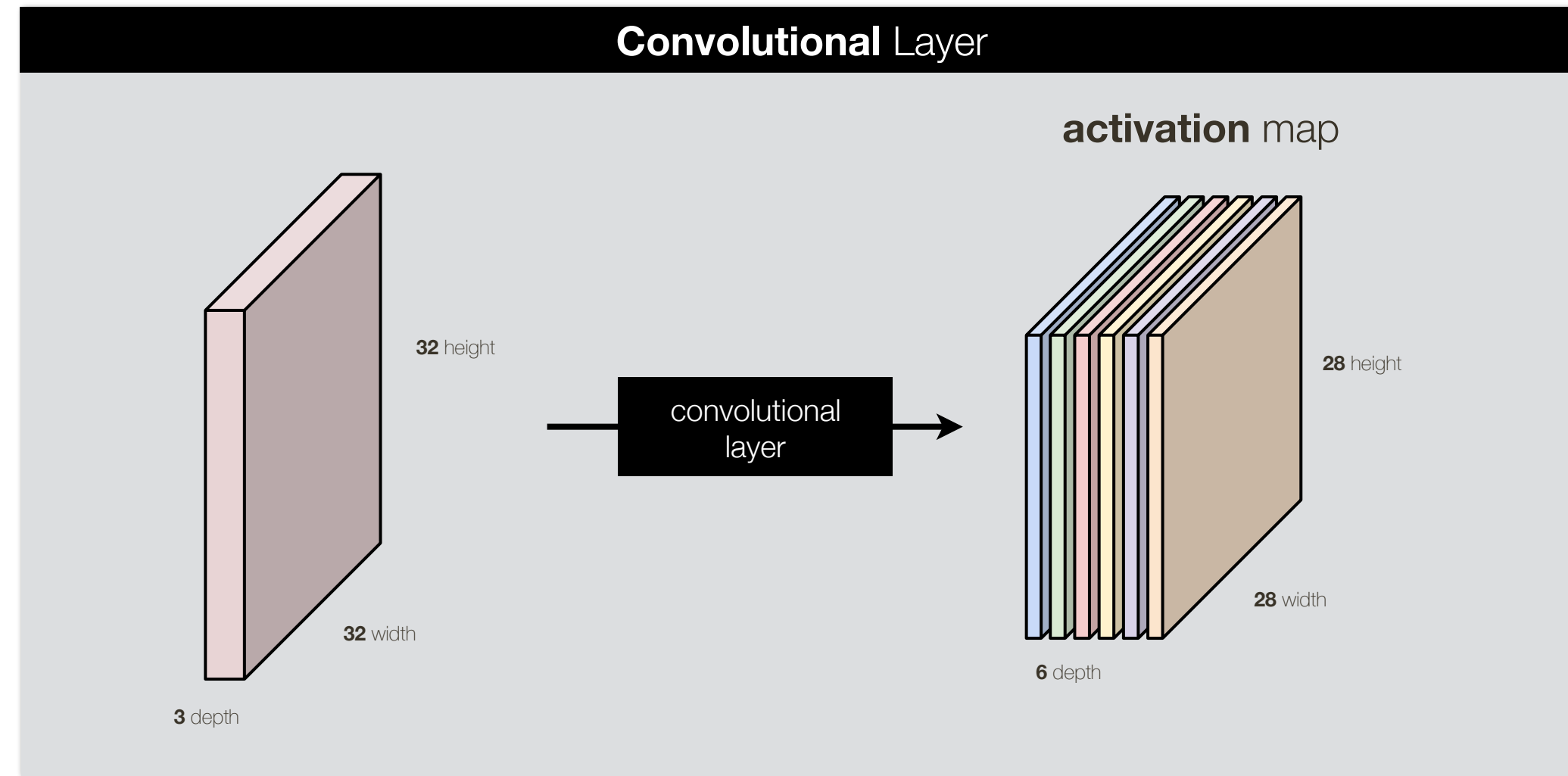




# Review of CNNs

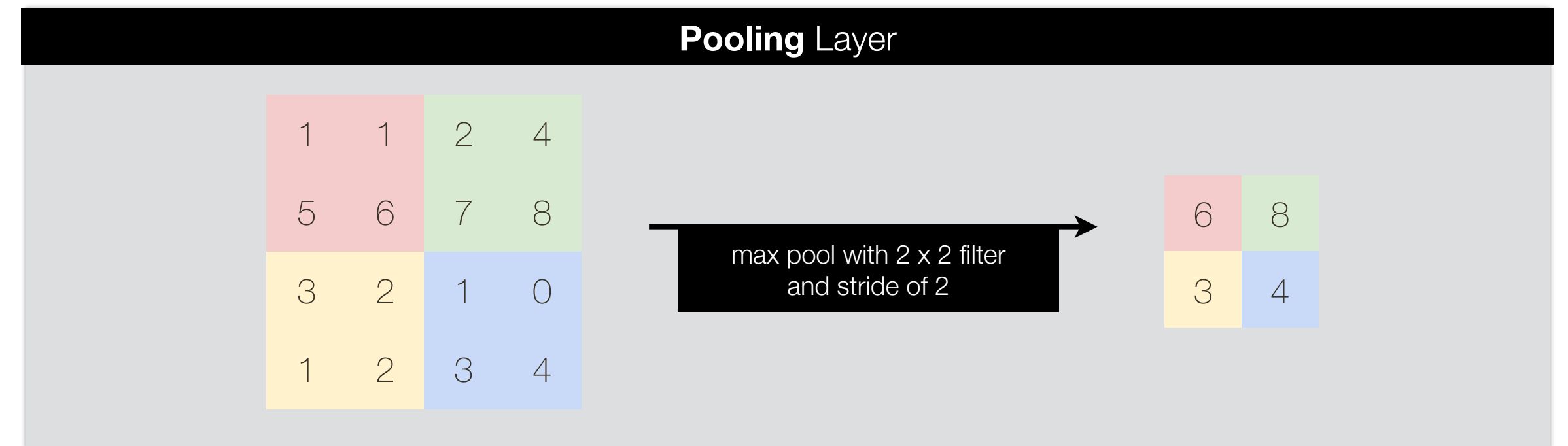
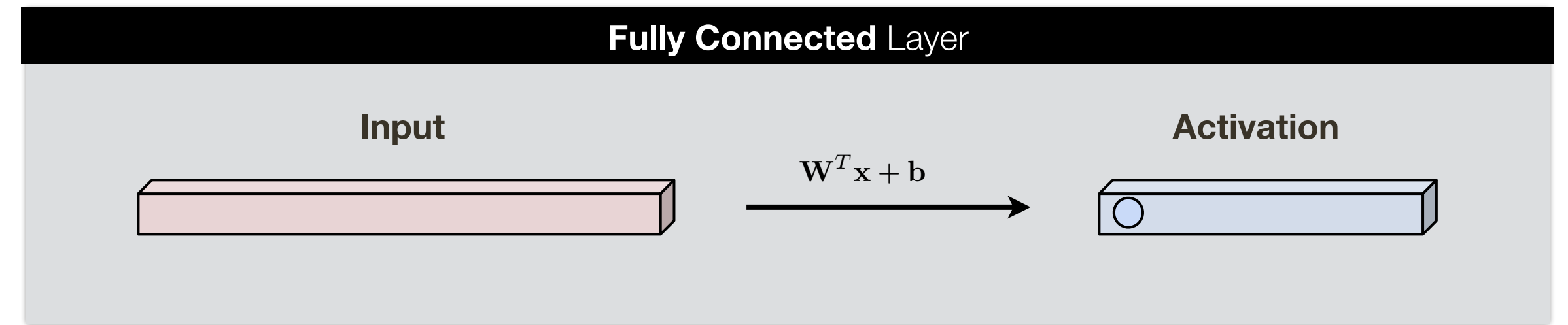
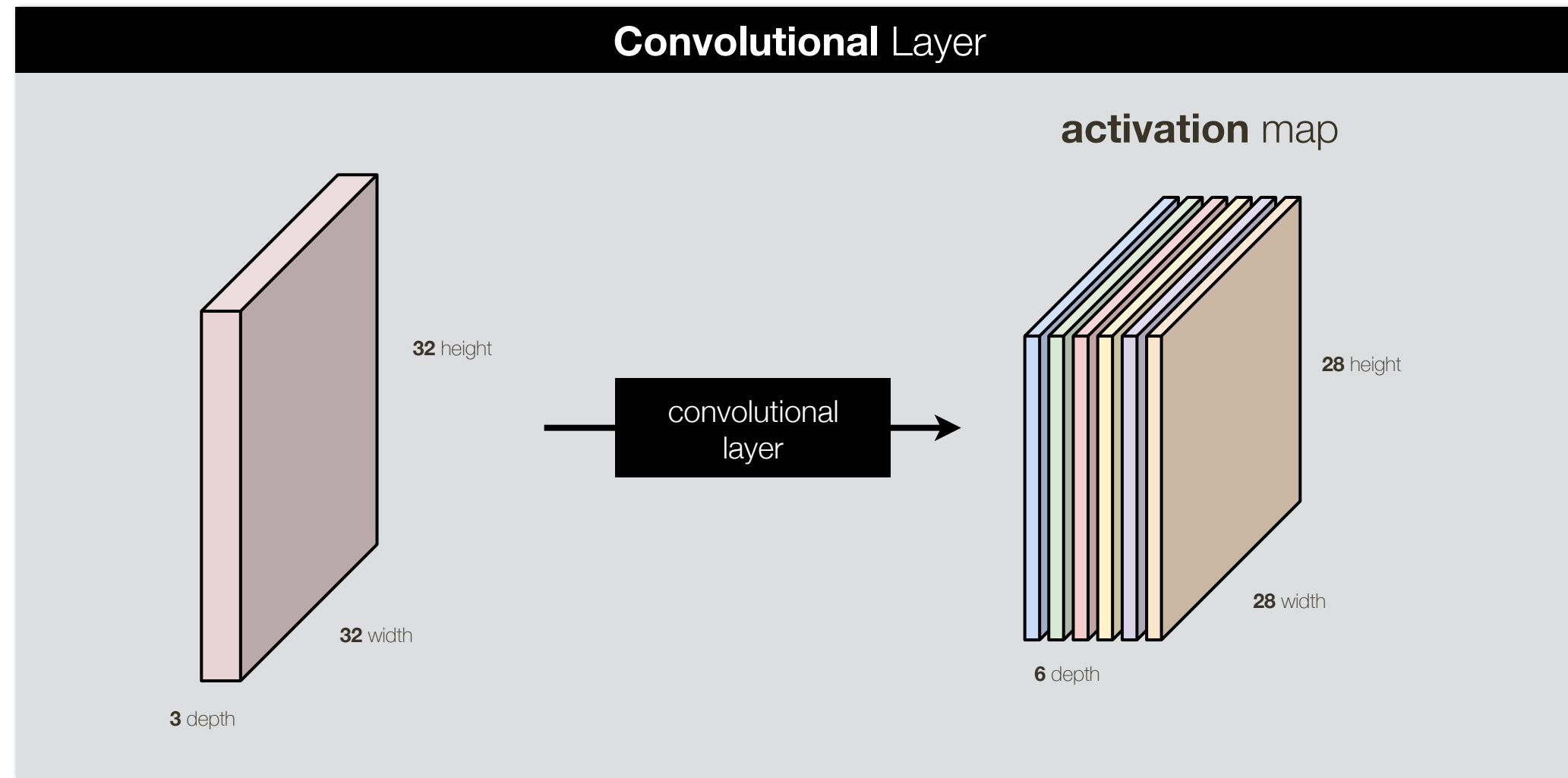


# Review of CNNs





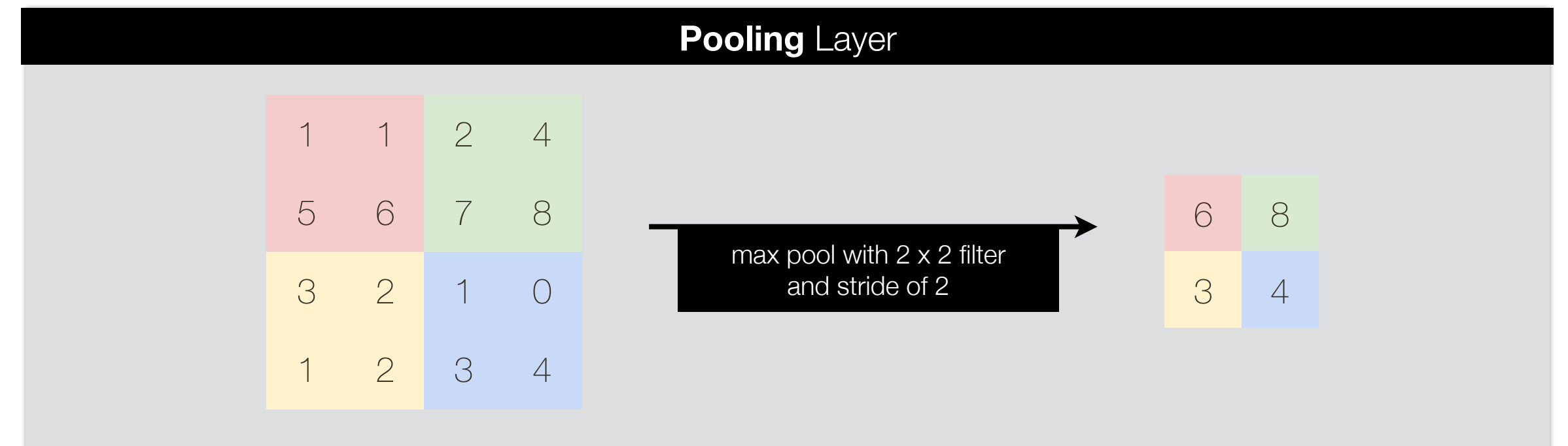
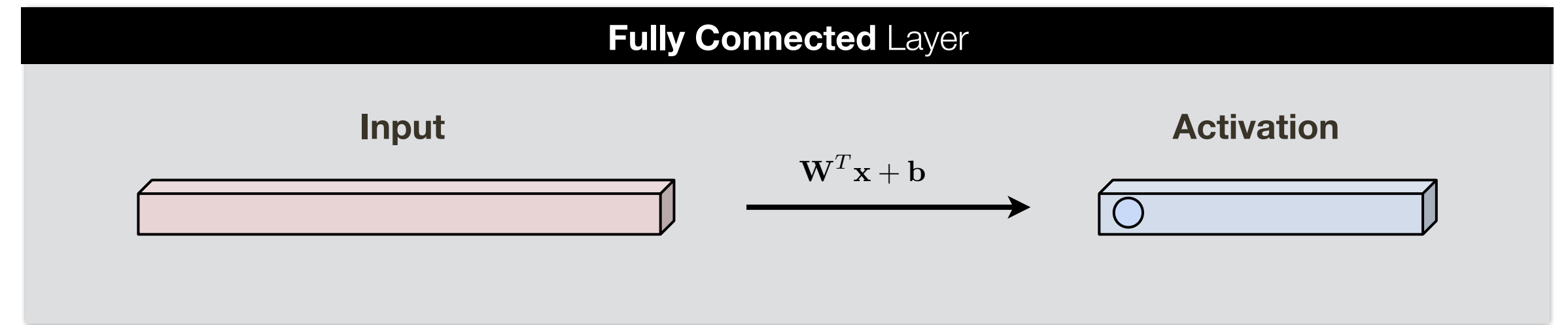
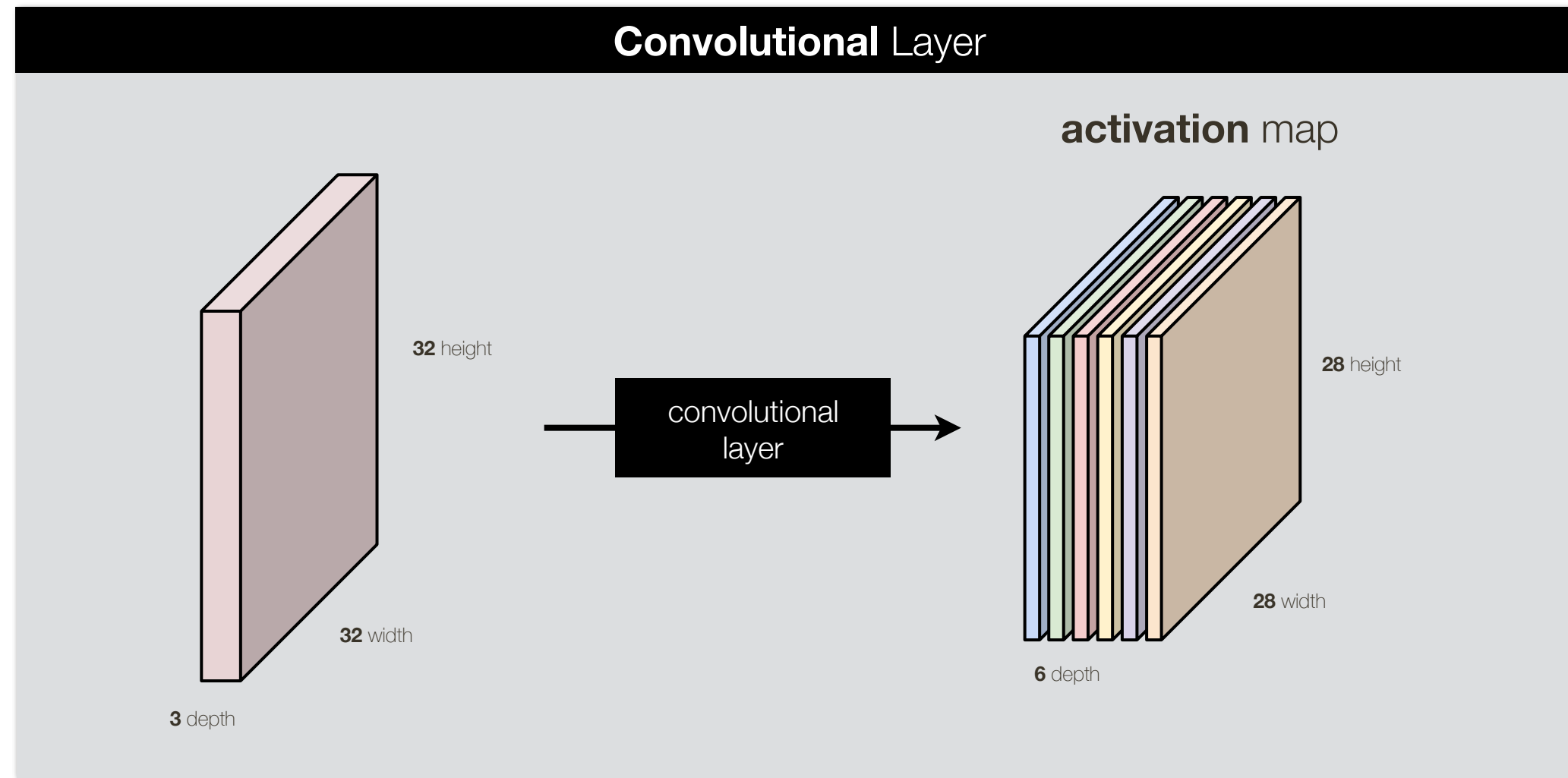
# Review of CNNs



## Effective Techniques for **Training**

- **Regularization:** L1, L2, data augmentation
- **Transfer Learning:** fine-tuning networks

# Review of CNNs



## Effective Techniques for **Training**

- **Regularization:** L1, L2, data augmentation
- **Transfer Learning:** fine-tuning networks

## Vision **Applications** of CNNs

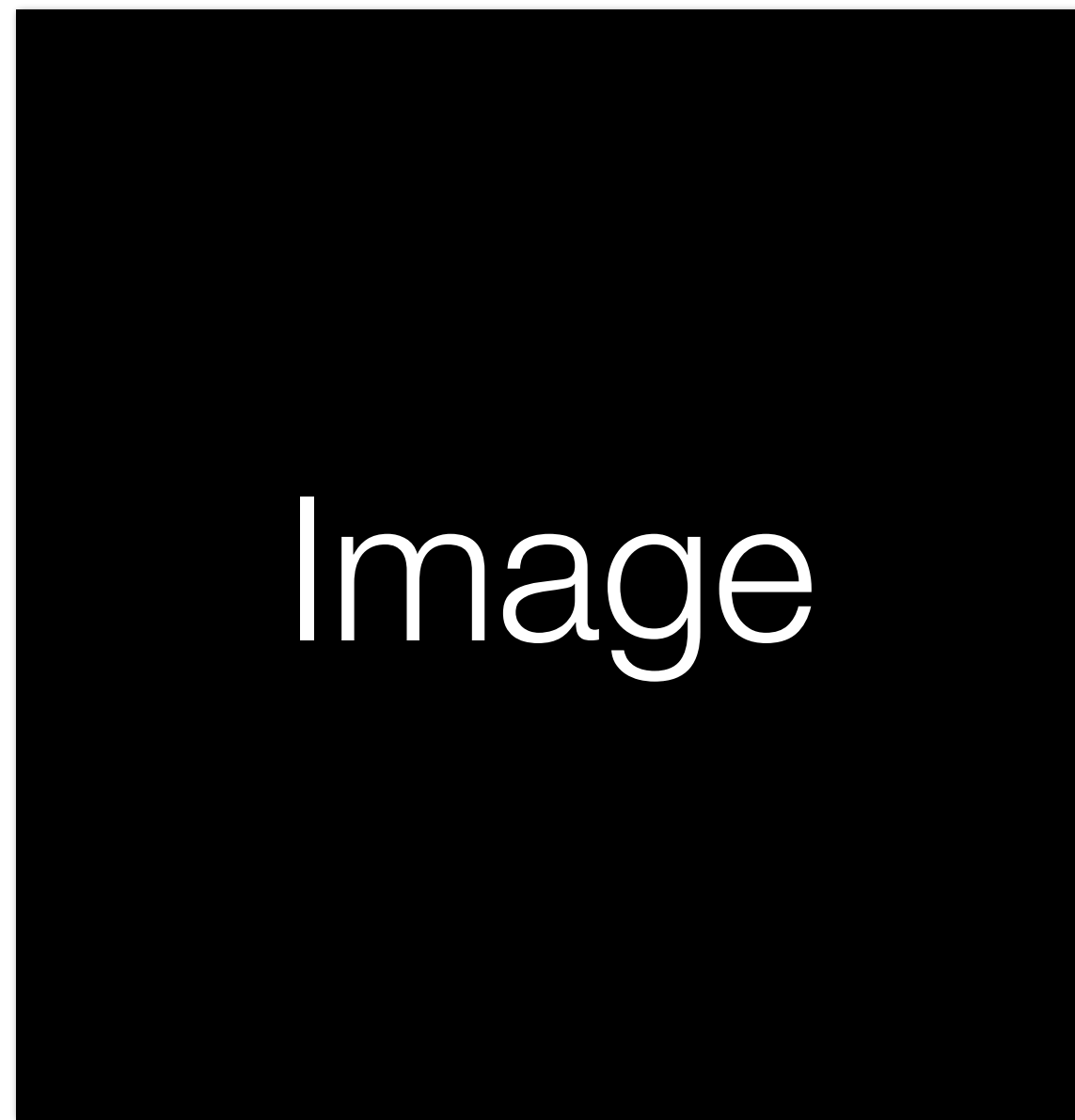
- **Classification:** AlexNet, VGG, GoogleLeNet, ResNet
- **Segmentation:** Fully convolutional CNNs
- **Detection:** R-CNN, Fast R-CNN, Faster R-CNN, YOLO

|                     | Categorization                                 | Detection  | Segmentation    | Instance Segmentation                  |
|---------------------|--|--|-----------------|--|
|                     |  |  |                 |  |
| <b>Multi-class:</b> | Horse<br>Church<br>Toothbrush<br><b>Person</b> | Horse (x, y, w, h)<br>Horse (x, y, w, h)<br>Person (x, y, w, h)<br>Person (x, y, w, h) | Horse<br>Person | Horse1<br>Horse2<br>Person1<br>Person2 |
| <b>Multi-label:</b> | Horse<br>Church<br>Toothbrush<br><b>Person</b> |  |                 |  |

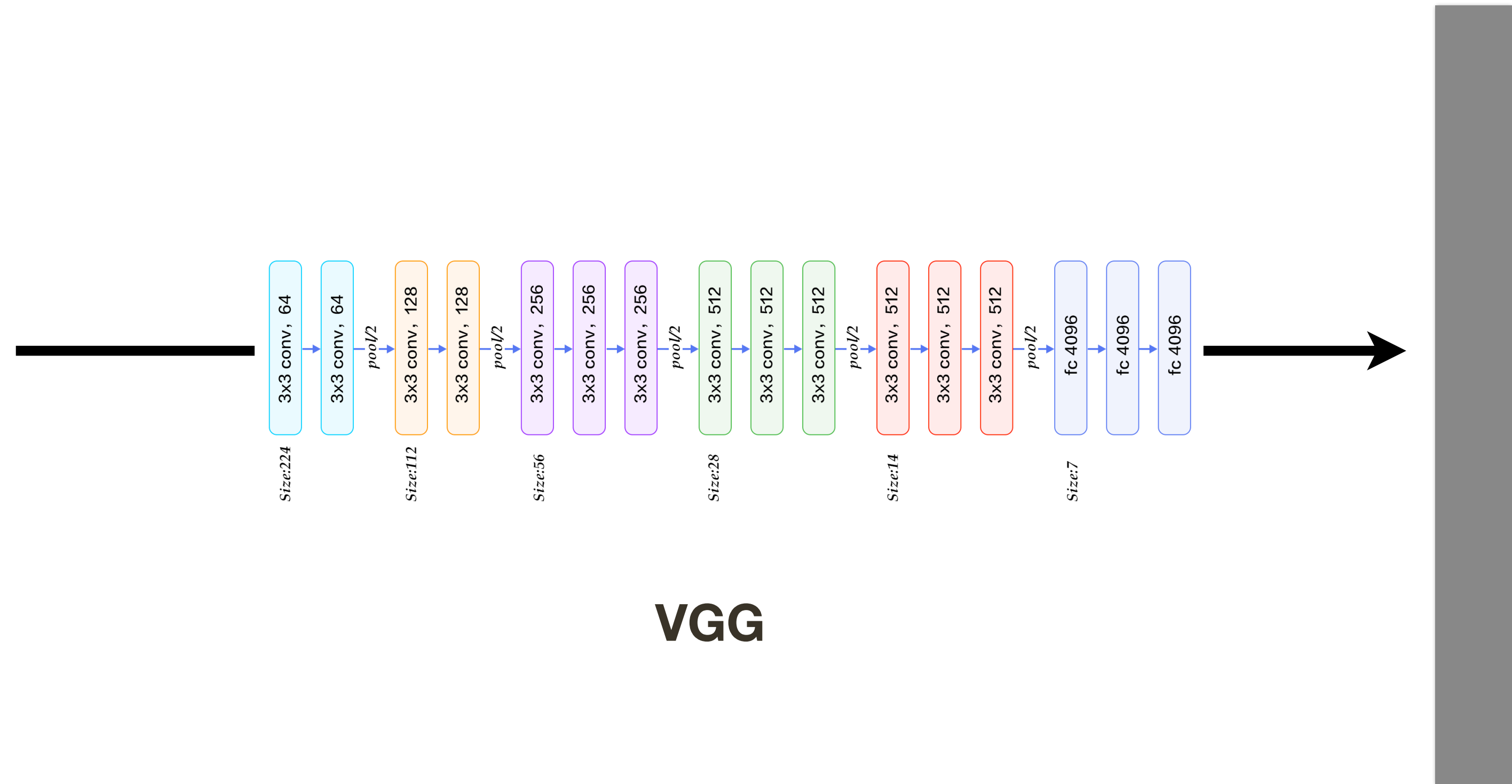
IMAGENET COCO Common Objects in Context



# Any CNN Could be Fully Convolutional

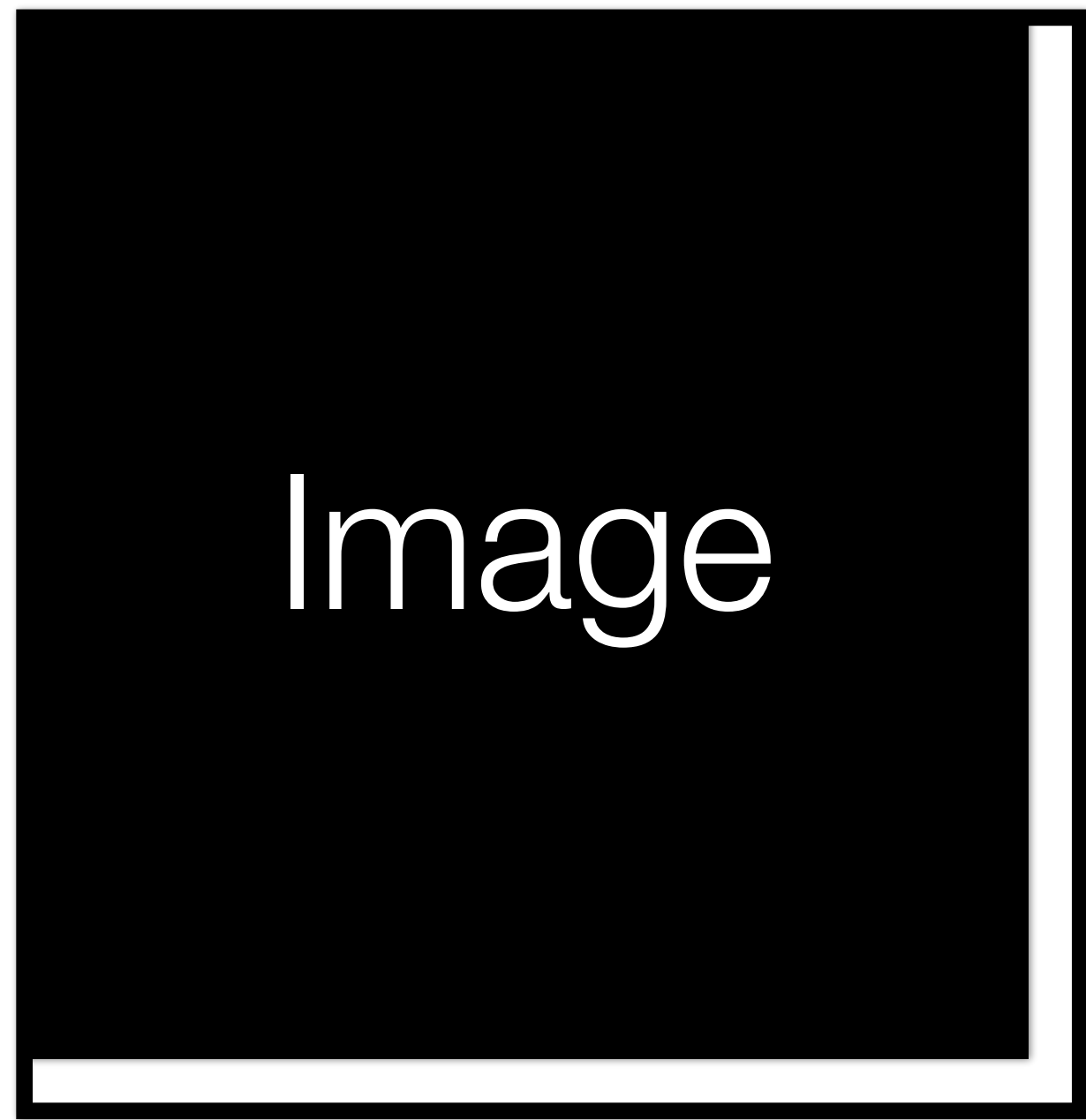


224 x 224

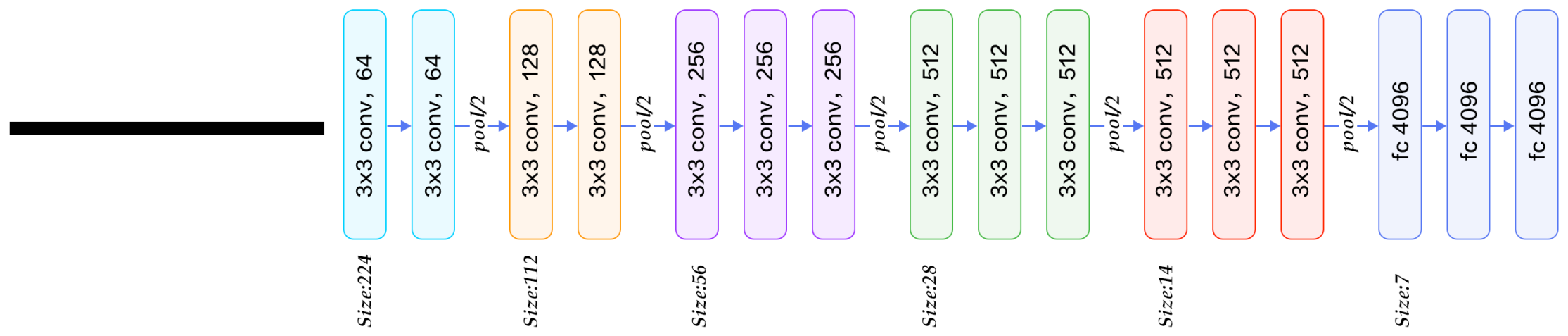


VGG

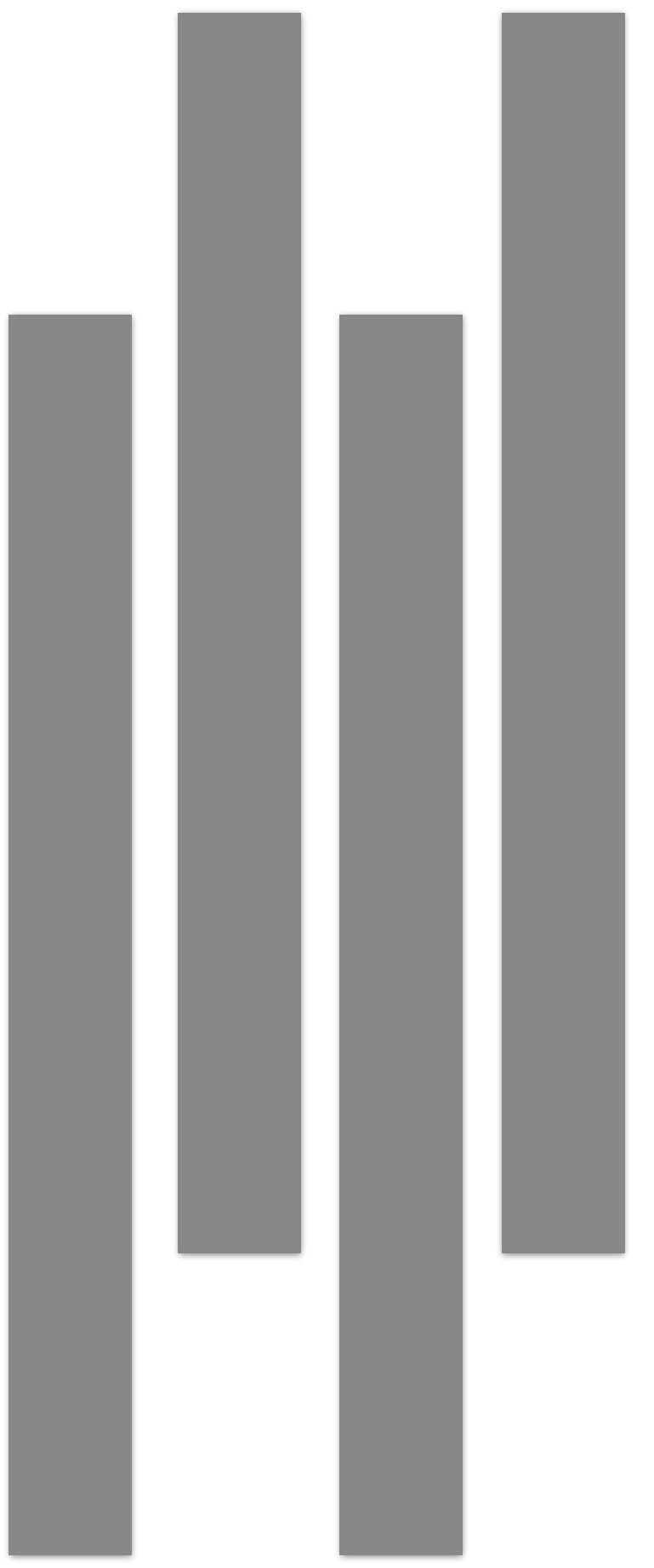
# Any CNN Could be Fully Convolutional



225 x 225



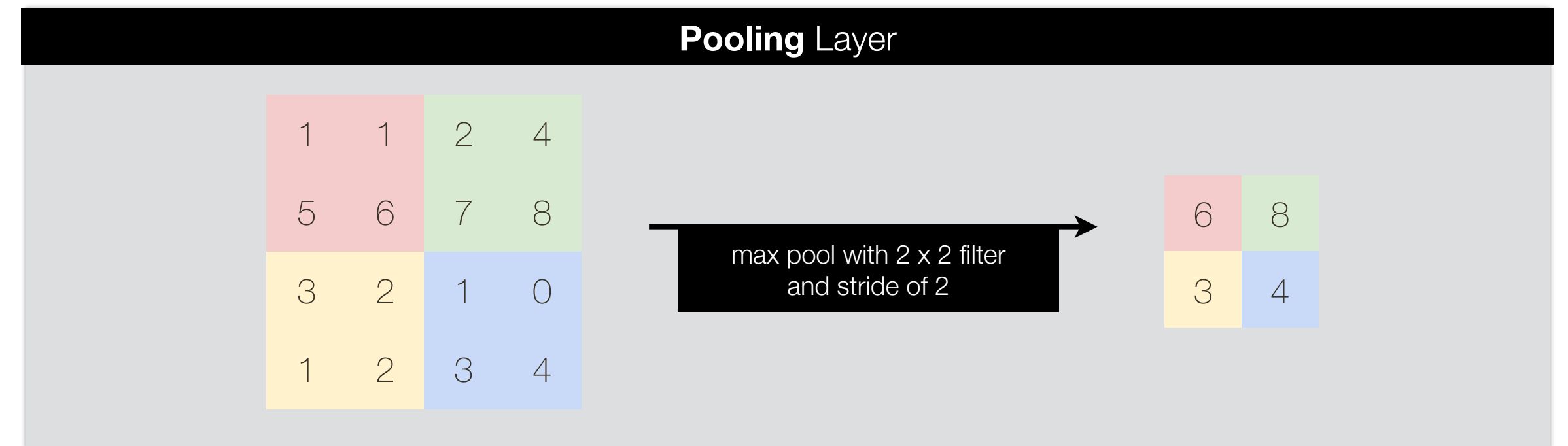
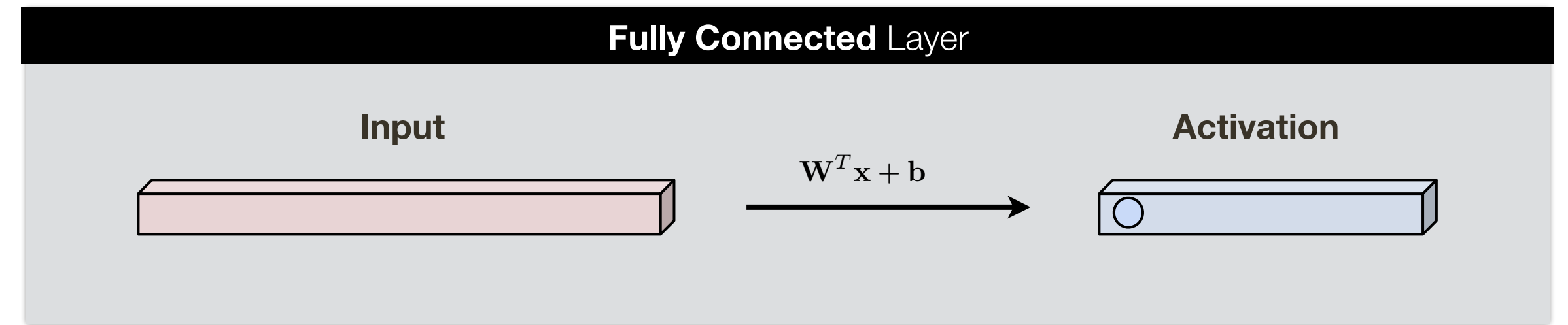
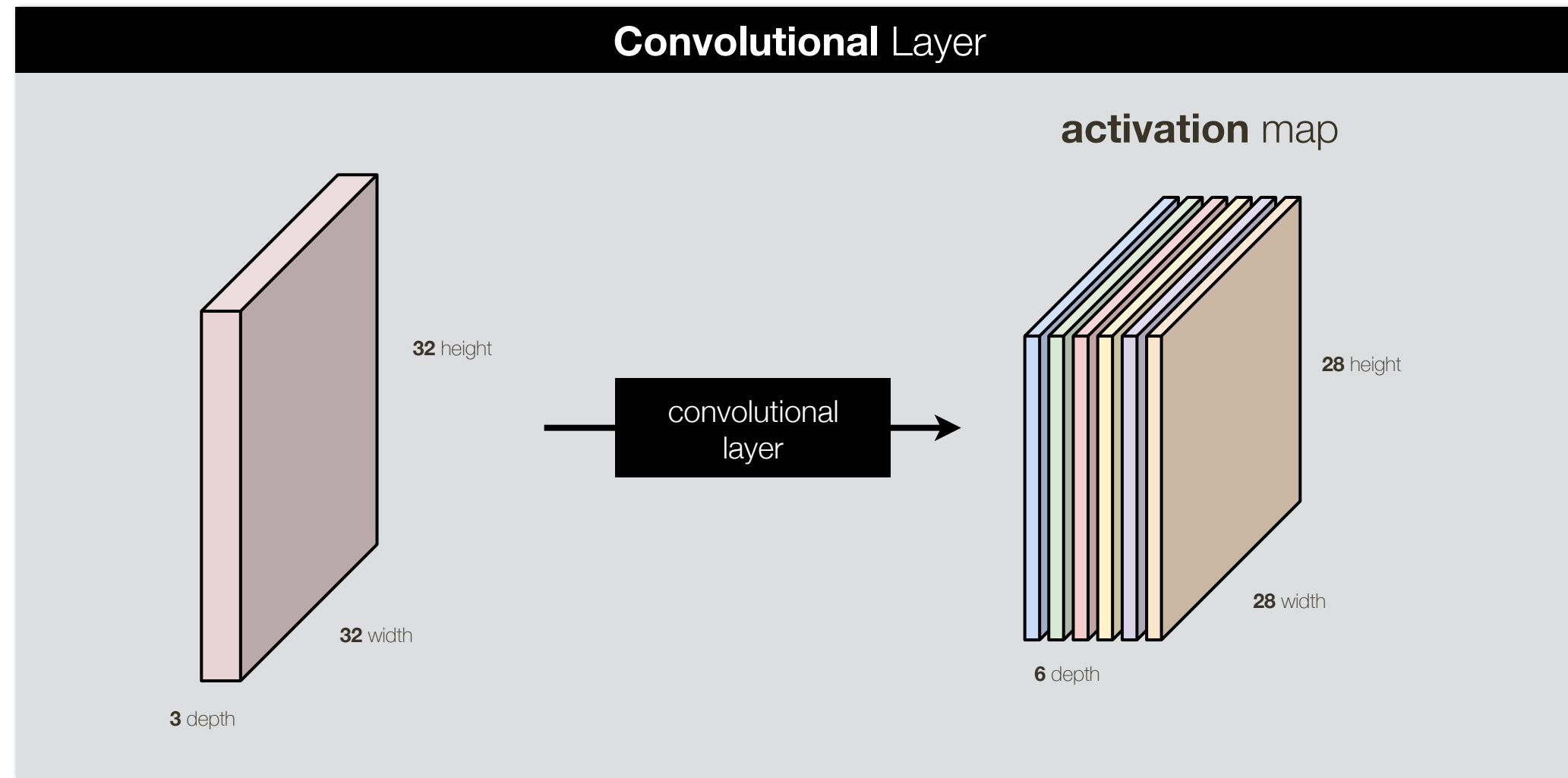
VGG



2 x 2 x 1000



# Review of CNNs



## Effective Techniques for **Training**

- **Regularization:** L1, L2, data augmentation
- **Transfer Learning:** fine-tuning networks

## Vision **Applications** of CNNs

- **Classification:** AlexNet, VGG, GoogleLeNet, ResNet
- **Segmentation:** Fully convolutional CNNs
- **Detection:** R-CNN, Fast R-CNN, Faster R-CNN, YOLO

|                     | Categorization                                 | Detection  | Segmentation    | Instance Segmentation                  |
|---------------------|--|--|-----------------|--|
|                     |  |  |                 |  |
| <b>Multi-class:</b> | Horse<br>Church<br>Toothbrush<br><b>Person</b> | Horse (x, y, w, h)<br>Horse (x, y, w, h)<br>Person (x, y, w, h)<br>Person (x, y, w, h) | Horse<br>Person | Horse1<br>Horse2<br>Person1<br>Person2 |
| <b>Multi-label:</b> | Horse<br>Church<br>Toothbrush<br><b>Person</b> |  |                 |  |

IMAGENET COCO Common Objects in Context