



Object detection and segmentation in deep learning

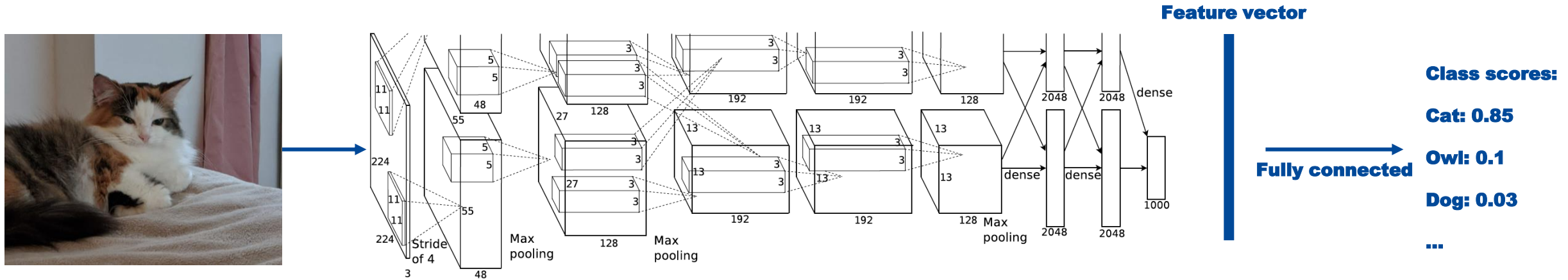
Colin Decourt – ANITI / NXP Semiconductors

08/02/2023

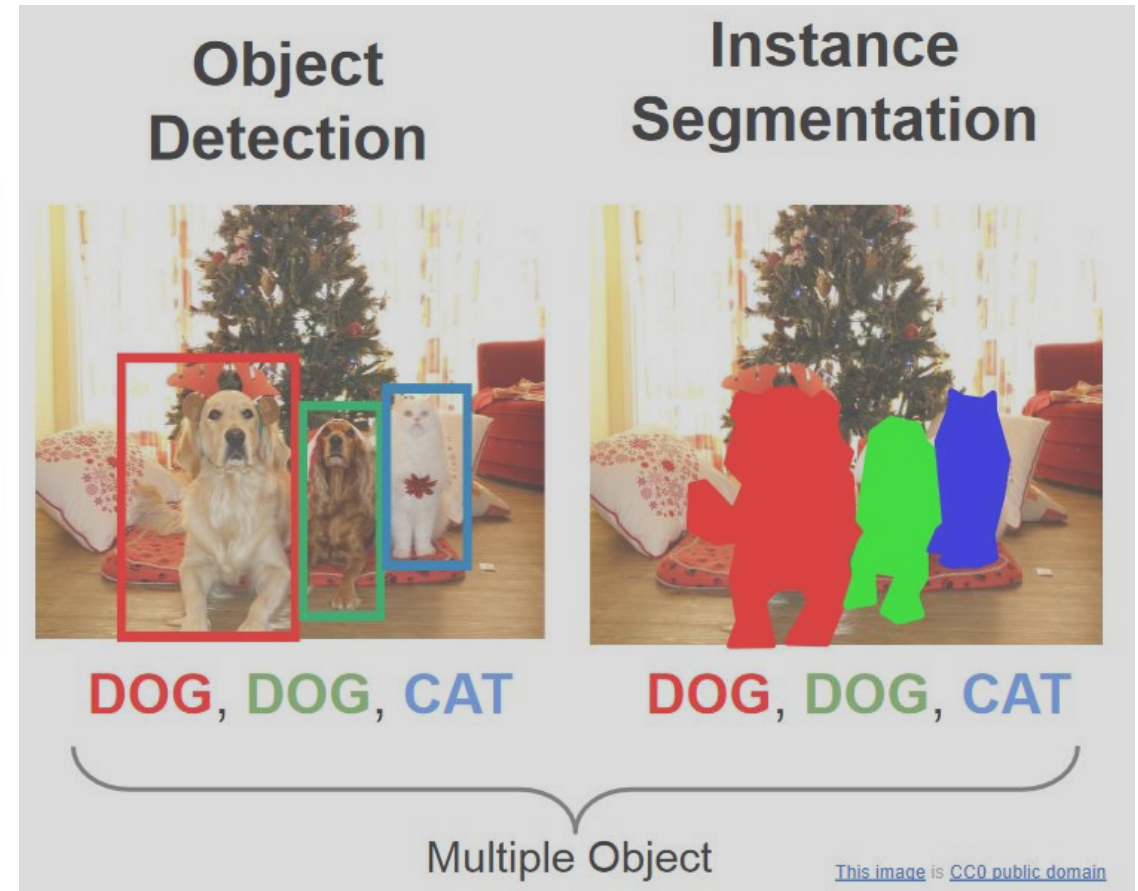
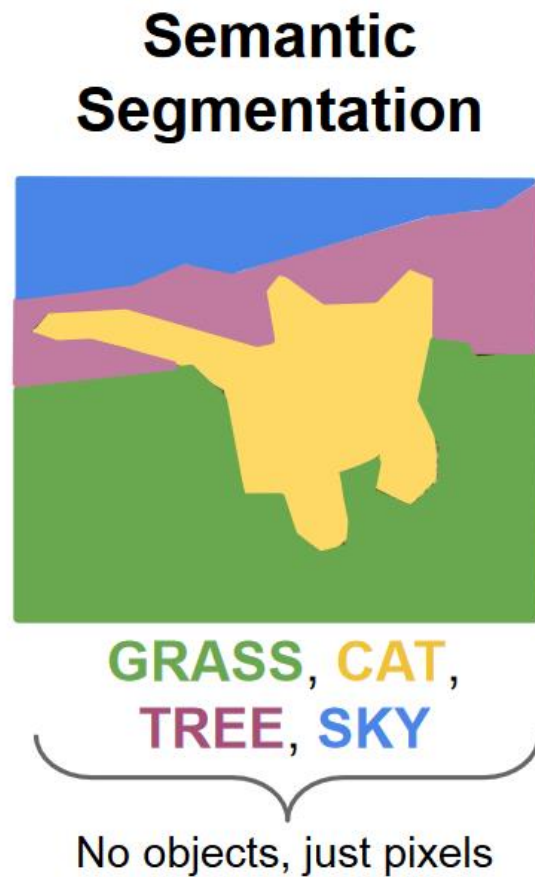
ANITI

Université
Fédérale
Toulouse
Midi-Pyrénées

From classification...

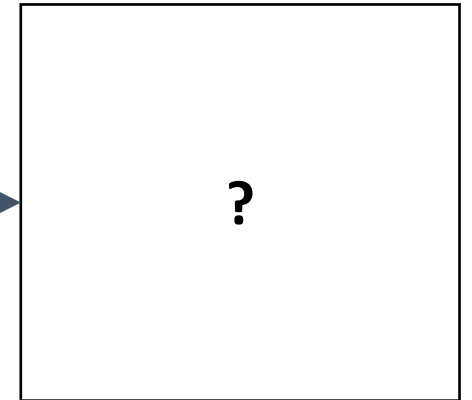
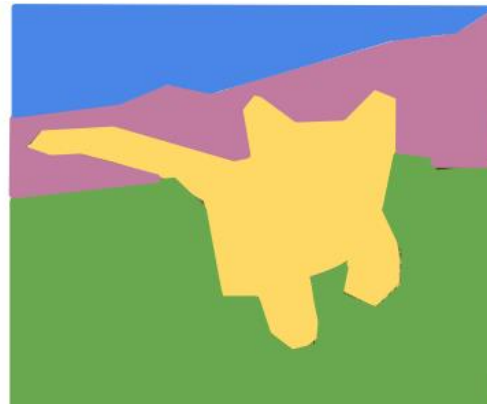


... to object detection and segmentation



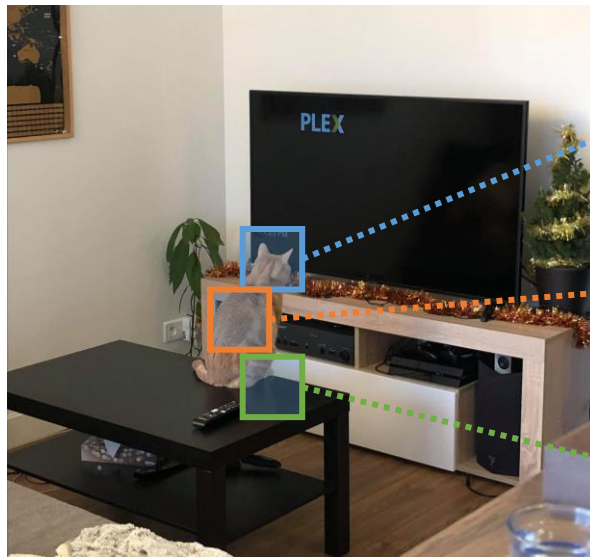
Semantic segmentation: what is it?

- Assign each pixel in the image to a category label (grass, cat, tree, sky...)
- Paired training data: for each training image, each pixel is labelled with a *semantic* category
- At test time, classify each pixel of a new image

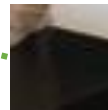


Semantic segmentation idea: sliding window

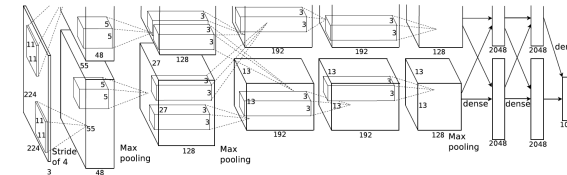
- Classify many crops of the image using CNN (ResNet, VGG...)
- How to classify without context?
- How do we include context?



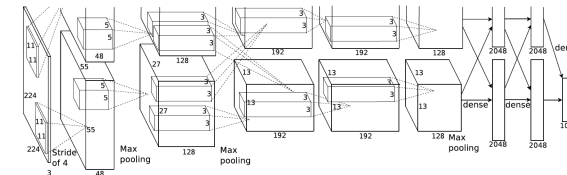
Extract patches



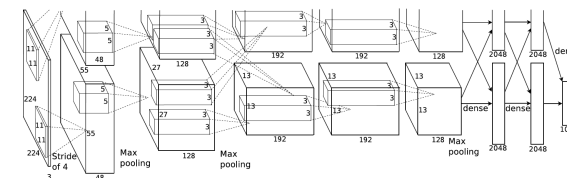
Classify centre pixel with CNN



Cat



Cat

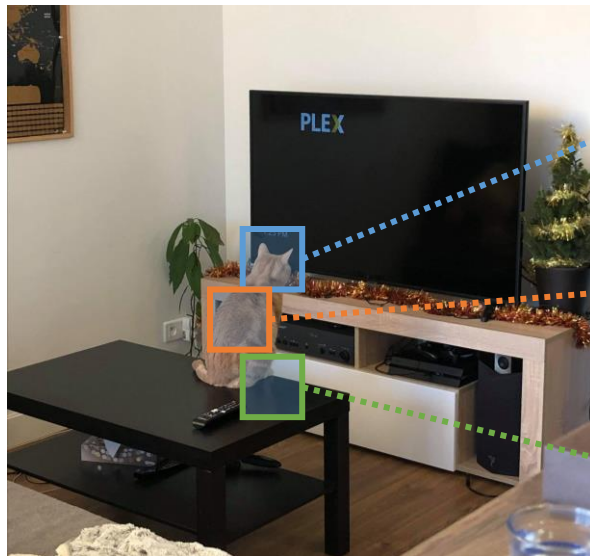


Table

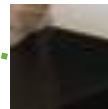
Semantic segmentation idea: sliding window

- Classify many crops of the image using CNN (ResNet, VGG...)
- How to classify without context?
- How do we include context?

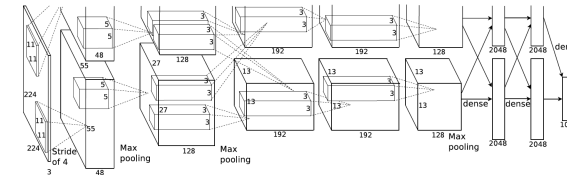
Very inefficient! For a 800x600 images ~ **58M boxes**
Not reusing shared features between overlapping patches



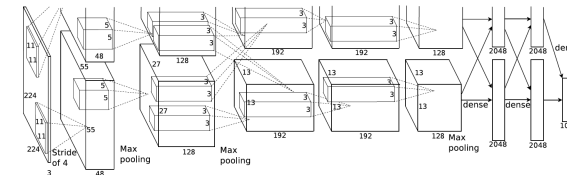
Extract patches



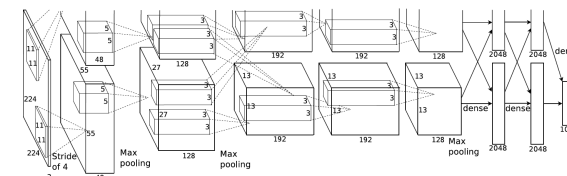
Classify centre pixel with CNN



Cat



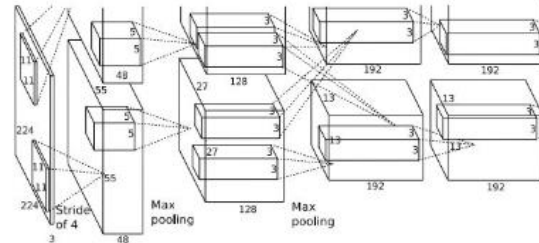
Cat



Table

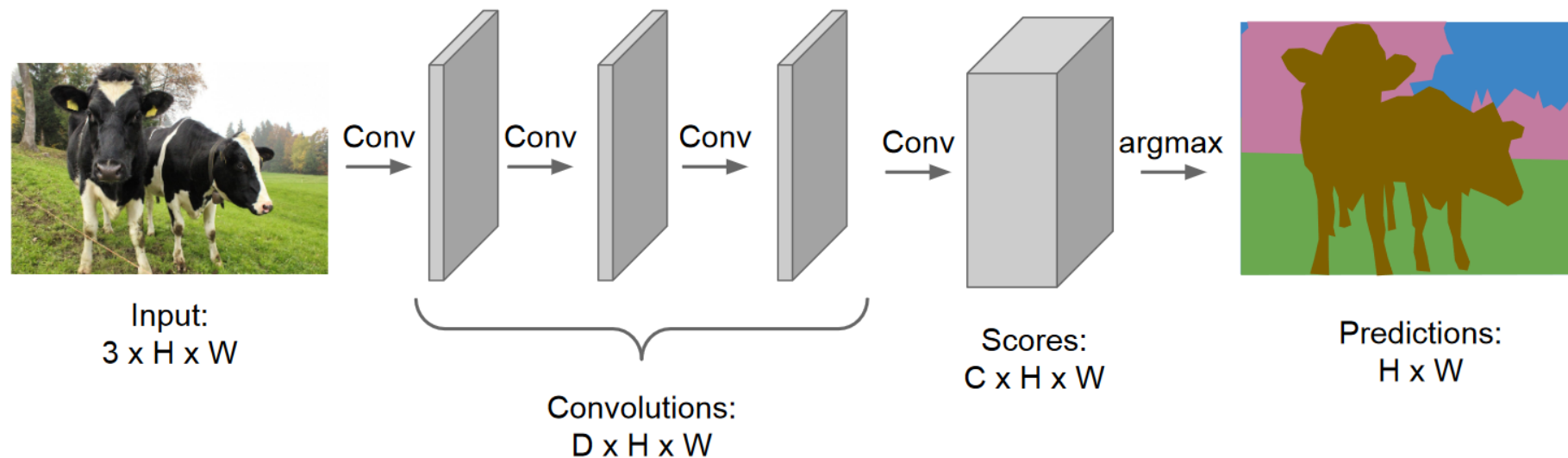
Semantic segmentation idea: convolution

- Idea: encode image features with ConvNets, and perform semantic segmentation on top
- Classification architectures reduce spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as the input size



Semantic segmentation idea: fully convolutional

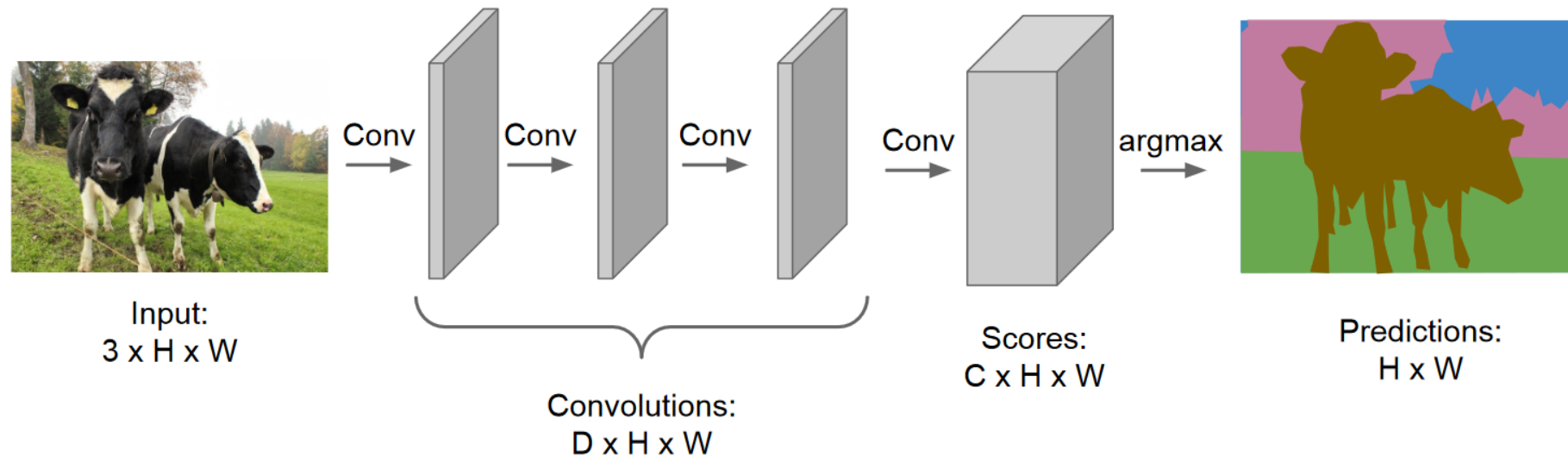
Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at ones



Very expensive as we don't decrease input size!

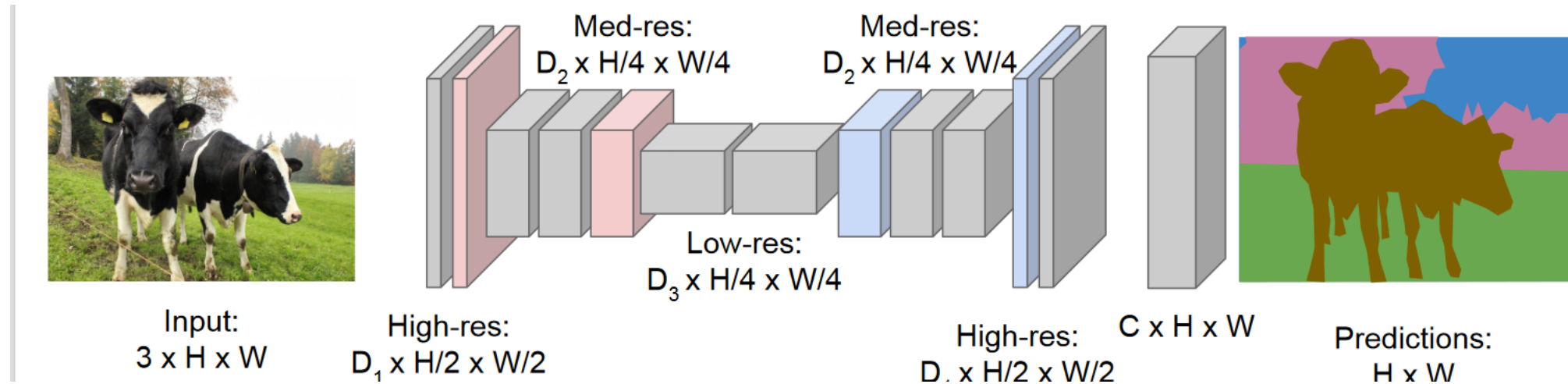
Semantic segmentation idea: fully convolutional

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at ones



Semantic segmentation idea: fully convolutional

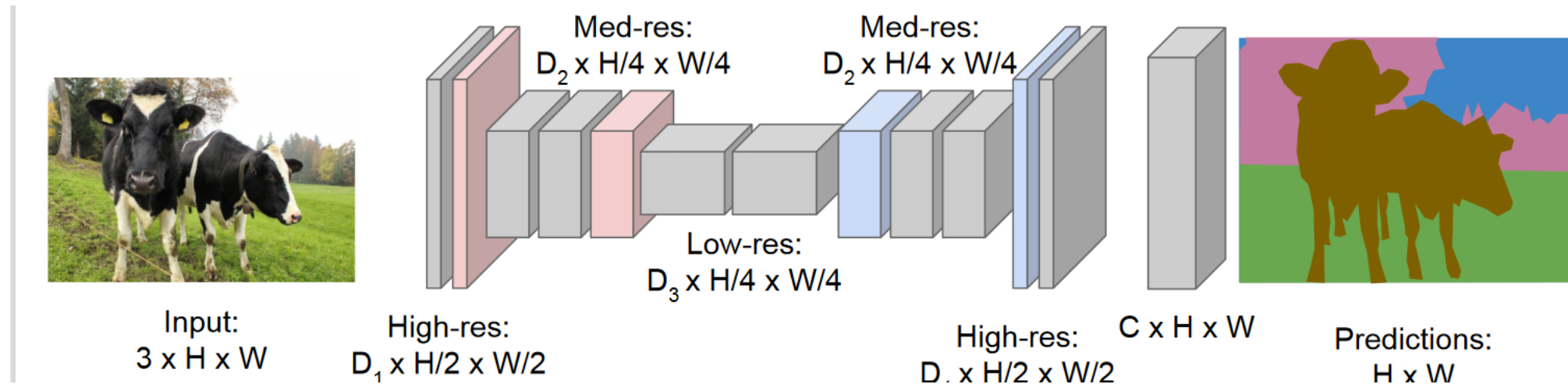
- Design the network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!
- Downsampling can be pooling, strided convolution...
- Upsampling can be nearest neighbour interpolation, unpooling, transposed convolution



Semantic segmentation idea: fully convolutional

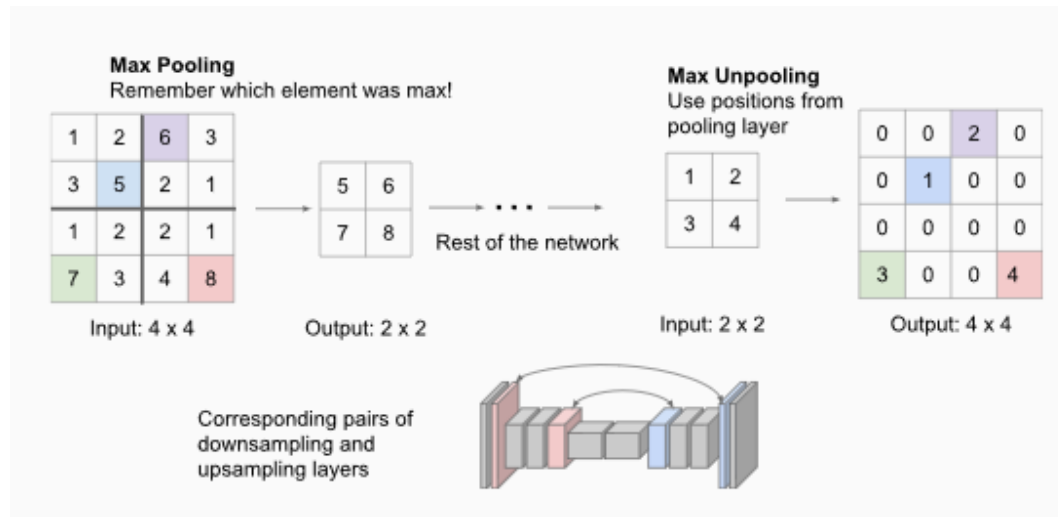
Idea: Learn to project low-resolution features onto pixel level (high-resolution)

- Design the network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!
- Downsampling can be pooling, strided convolution...
- Upsampling can be nearest neighbour interpolation, unpooling, transposed convolution

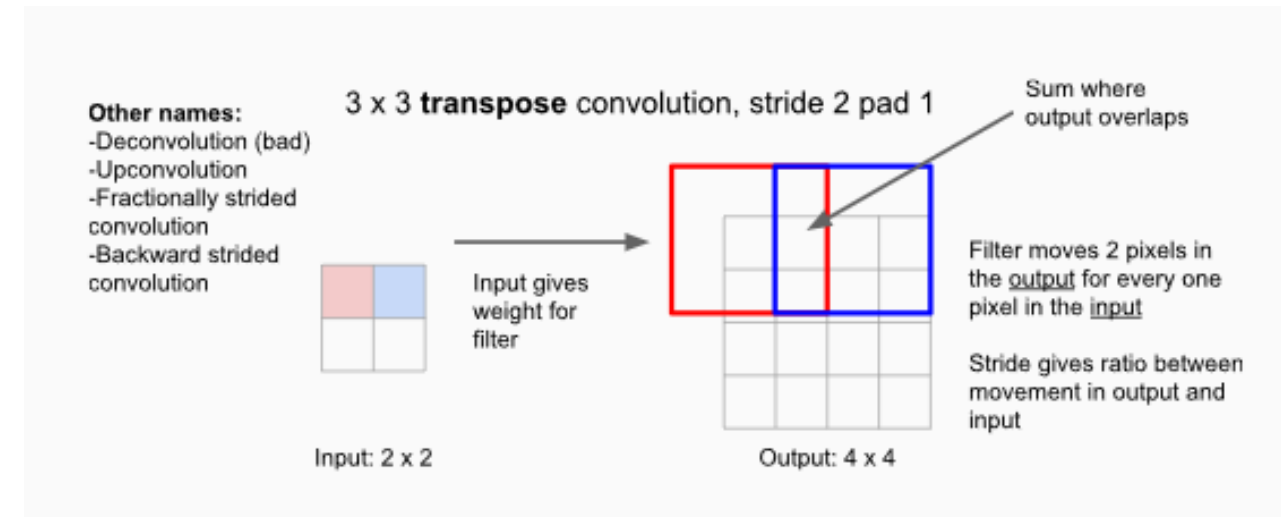


Upsampling low-resolution features

Restore the condensed feature map to the original size of the input image by expanding the feature dimension



1. Unpooling

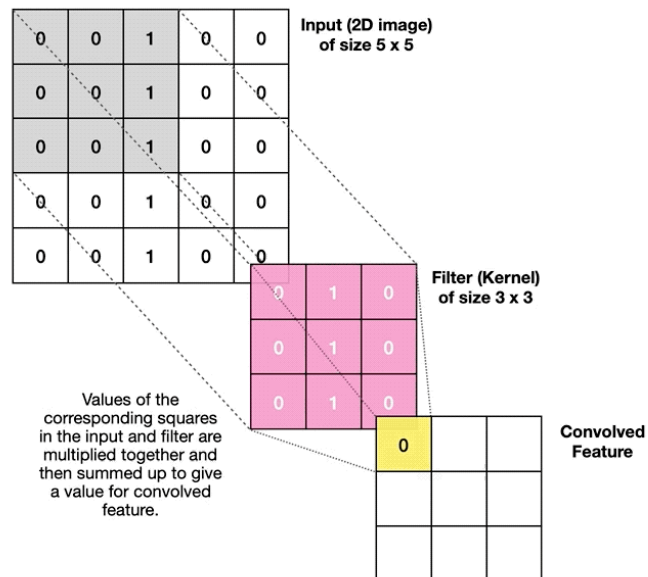


2. Transposed convolution

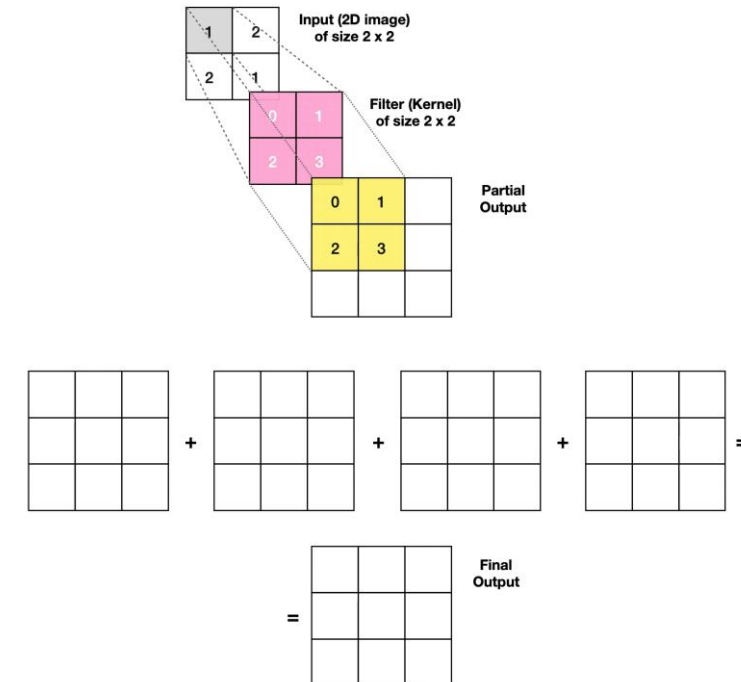
Transposed convolution

Learn to upsample the input feature map to the desired size using **learnable parameters**

Classic convolution

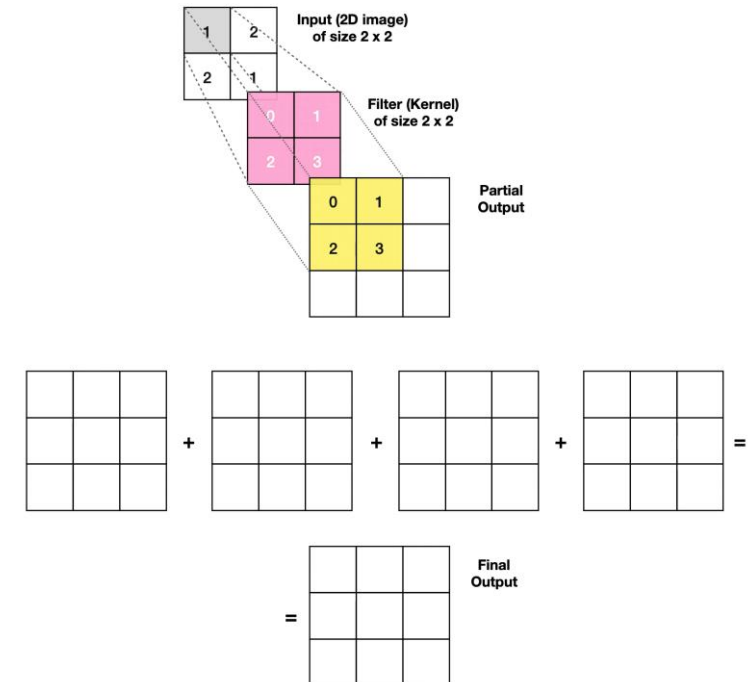
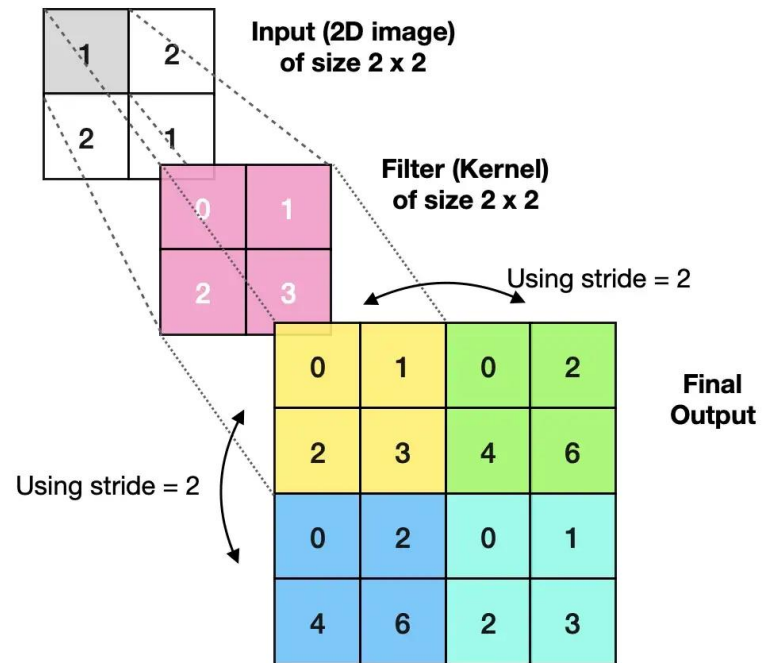


Transposed convolution



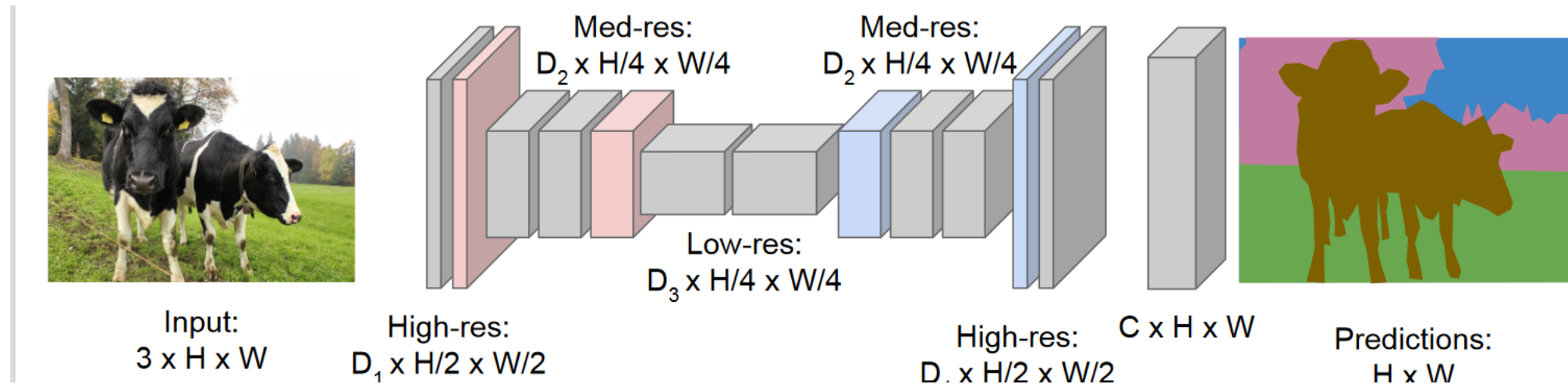
Transposed convolution

Learn to upsample the input feature map to the desired size using **learnable parameters**



Semantic segmentation: fully convolutional

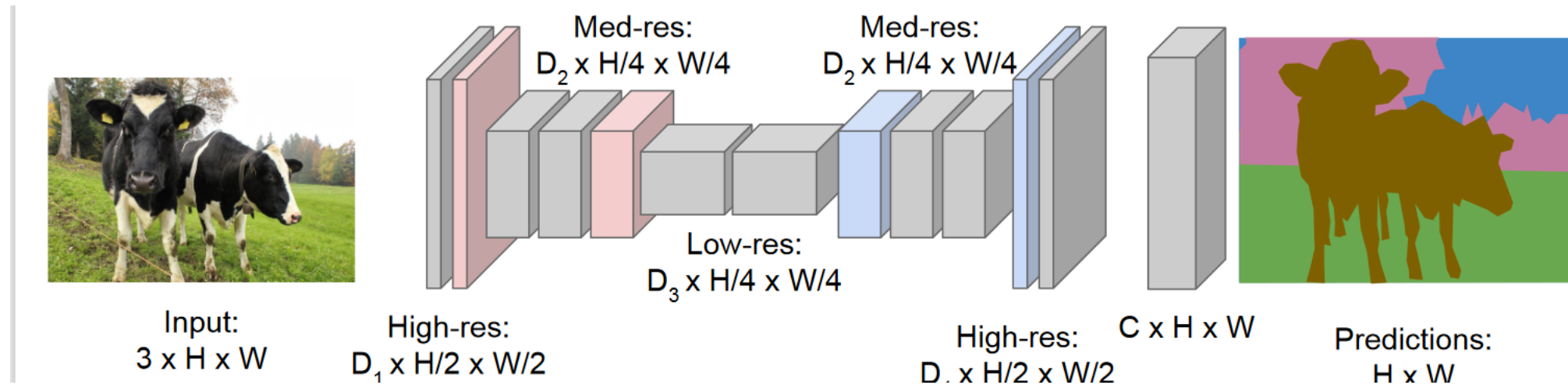
- Design the network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!
- Downsampling can be pooling, strided convolution...
- Upsampling can be nearest neighbour interpolation, unpooling, transposed convolution



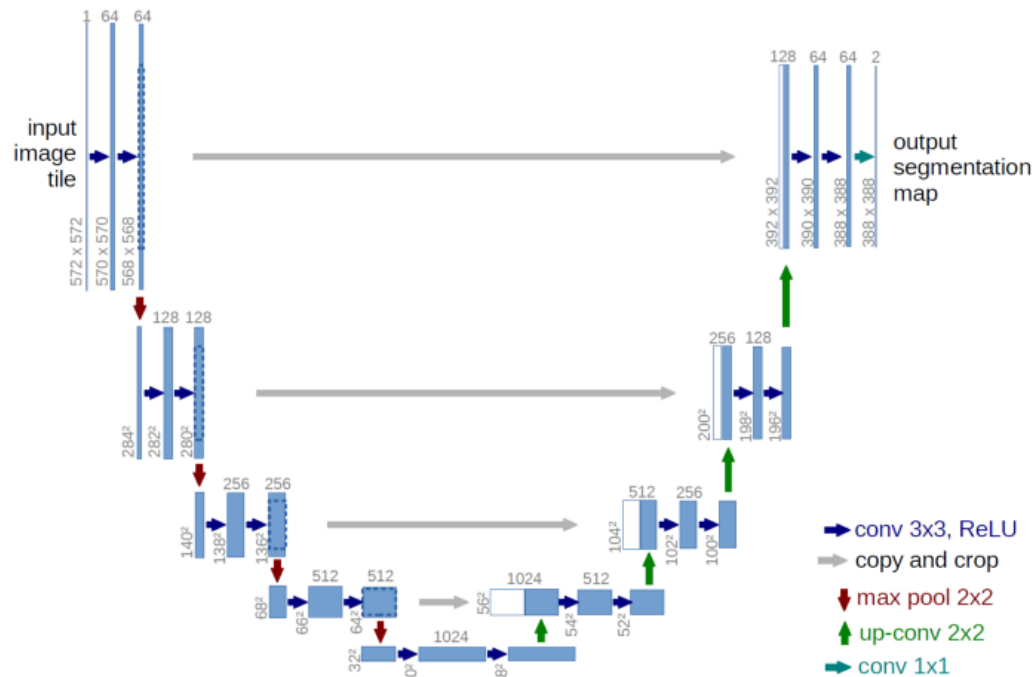
Don't differentiate instances (yet),
only care about pixels

Semantic segmentation: fully convolutional

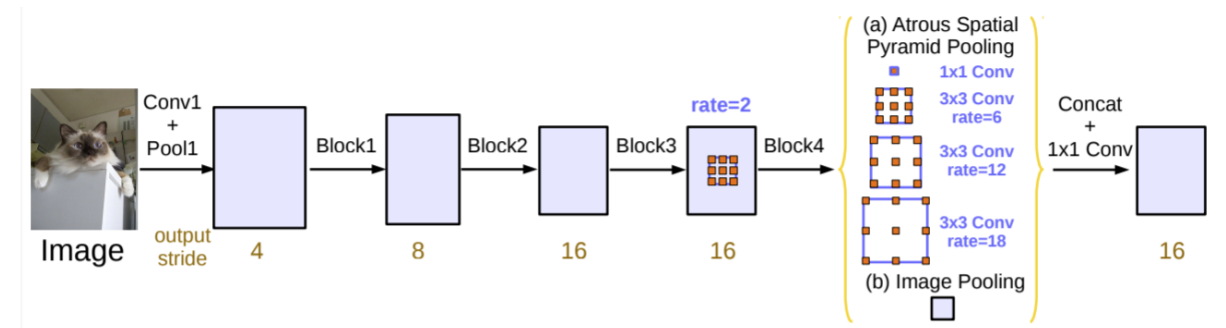
- Design the network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!
- Downsampling can be pooling, strided convolution...
- Upsampling can be nearest neighbour interpolation, unpooling, transposed convolution



Semantic segmentation: some models...

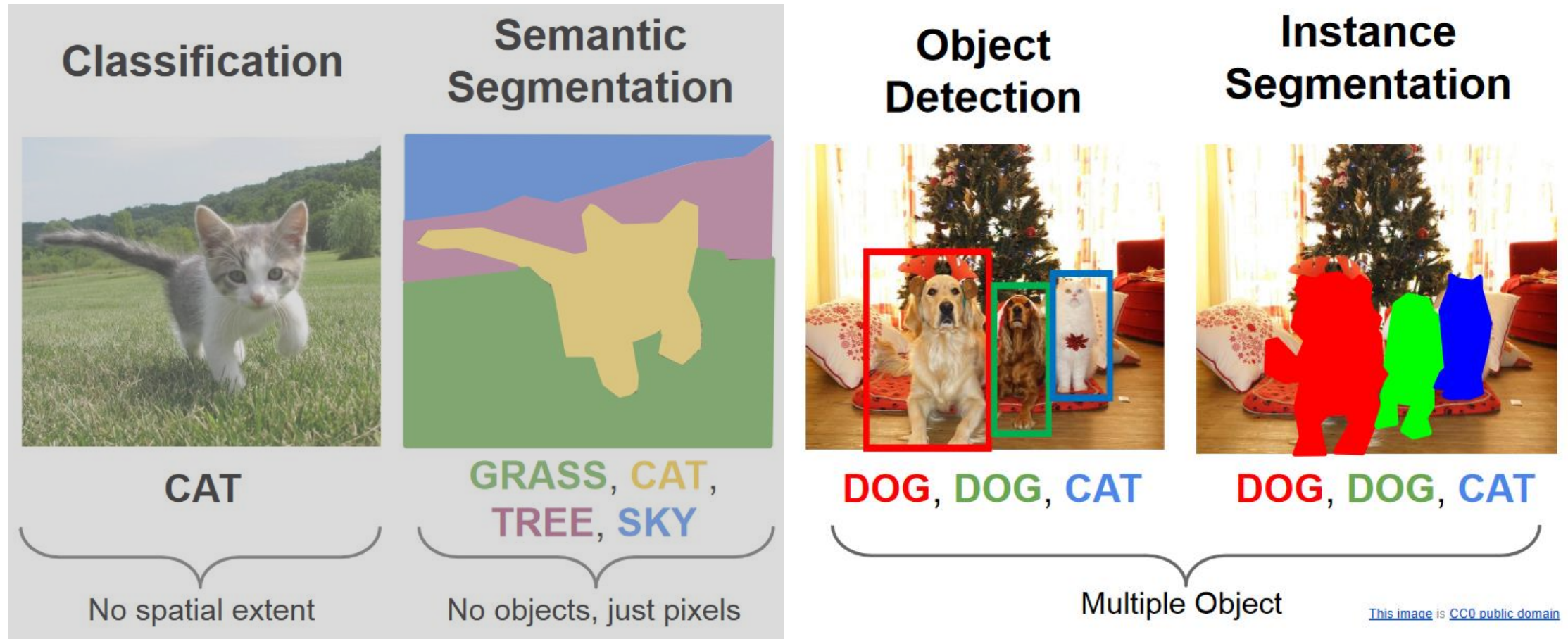


U-Net



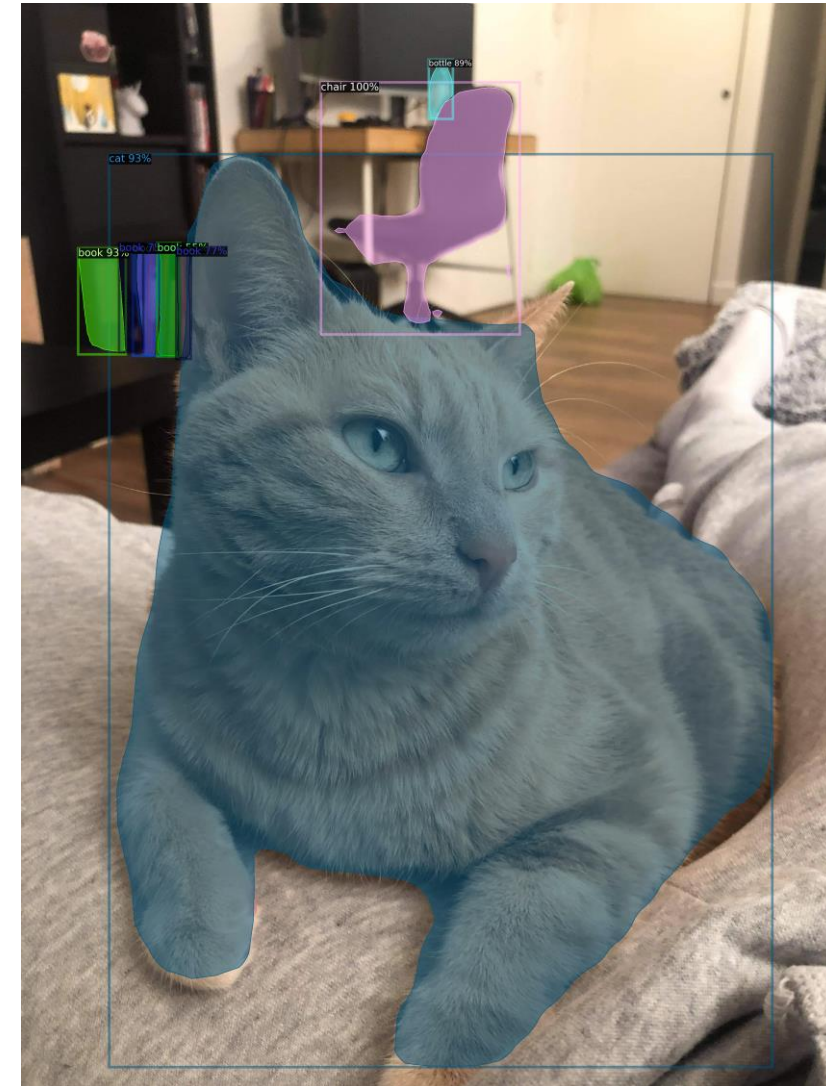
DeepLabv3

From classification to object detection and segmentation



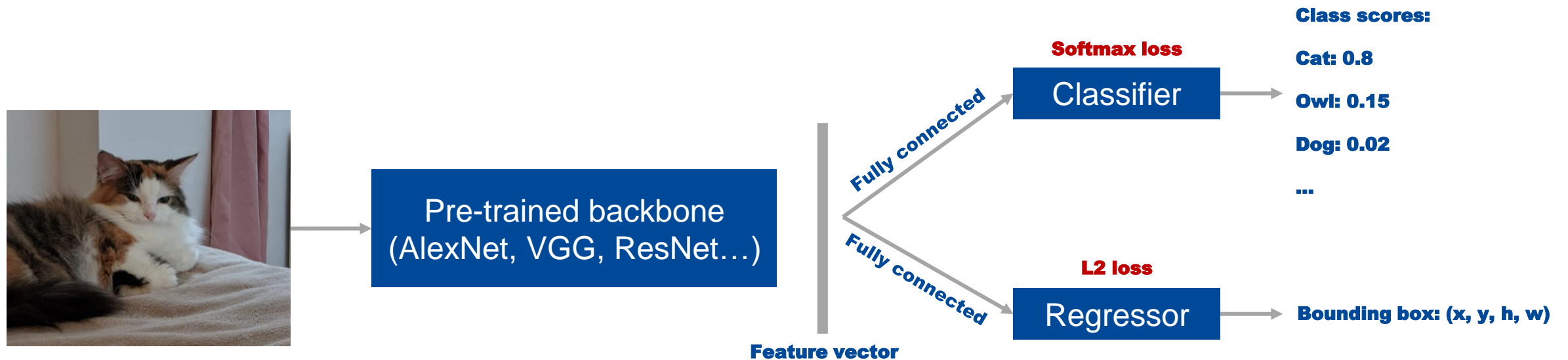
Object detection: what is it?

- The task of assigning a label and a bounding box (or a mask) to all objects in the image
- Semantic segmentation:
 - No objects, just pixels
 - Output a segmentation mask with the same size as the input
- Object detection: for each object output **a label** and a box (or mask)

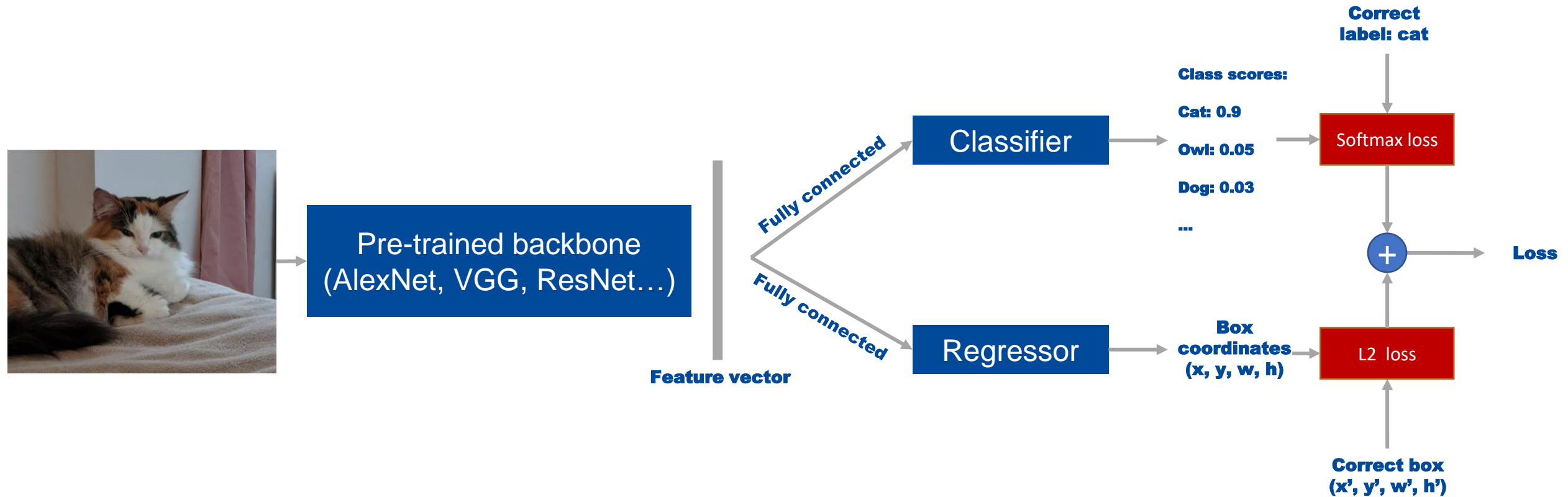


An easy case... Detecting single objects

Object detection is at the same time a **regression** problem and a **classification** task

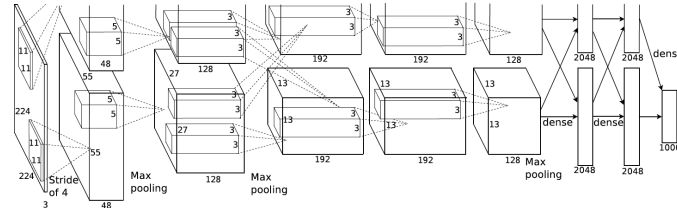
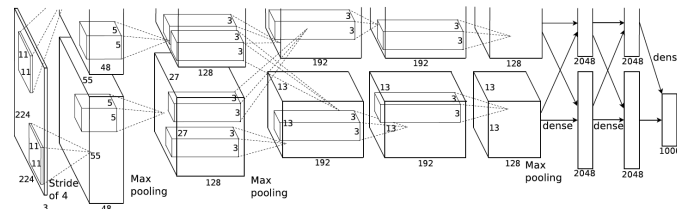
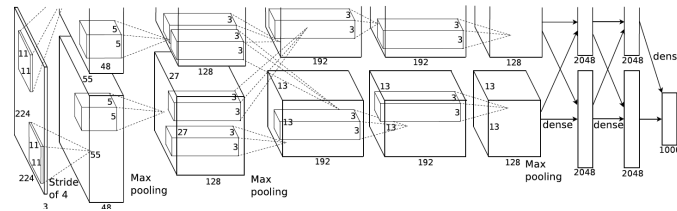


Detecting single objects: multitask loss



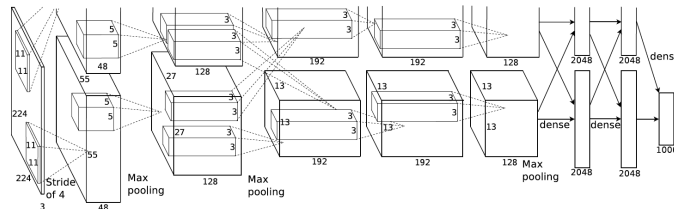
Object detection: challenges

- Multiple outputs: **variable number of objects per image**
- Multiple types of output: category label, bounding boxes (orientation, velocity...)



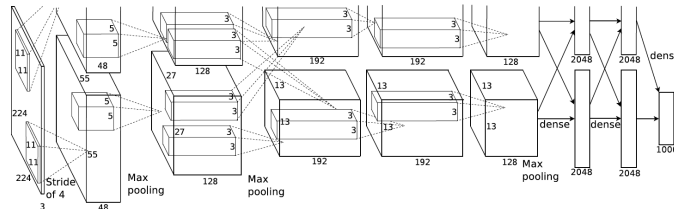
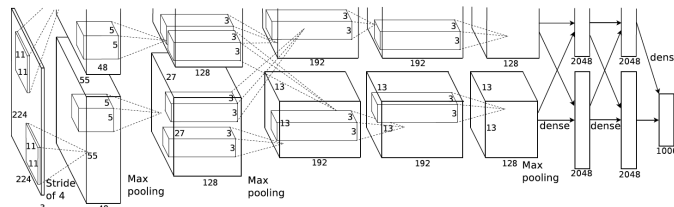
Object detection: challenges

- Multiple outputs: **variable number of objects per image**
- Multiple types of output: category label, bounding boxes (orientation, velocity...)



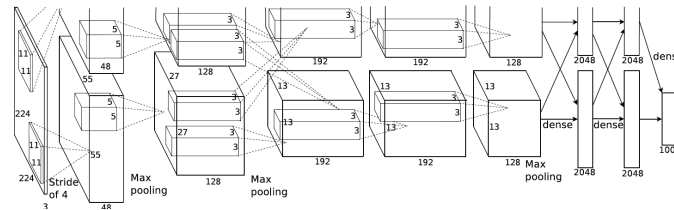
Cat + (x, y, w, h)

4 numbers to predict



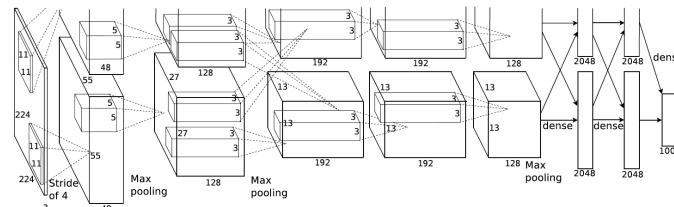
Object detection: challenges

- Multiple outputs: **variable number of objects per image**
- Multiple types of output: category label, bounding boxes (orientation, velocity...)



Cat + (x, y, w, h)

4 numbers to predict

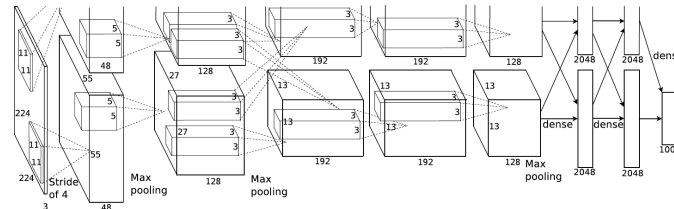


Dog + (x, y, w, h)

Dog + (x, y, w, h)

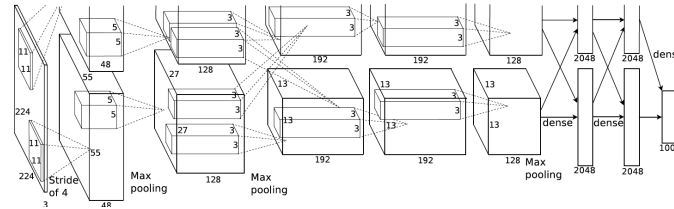
Cat + (x, y, w, h)

12 numbers to predict



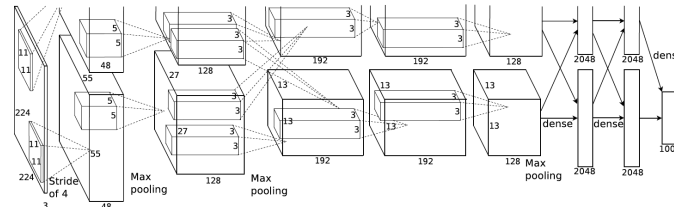
Object detection: challenges

- Multiple outputs: **variable number of objects per image**
- Multiple types of output: category label, bounding boxes (orientation, velocity...)



Cat + (x, y, w, h)

4 numbers to predict

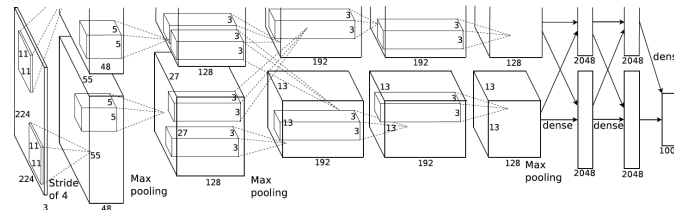


Dog + (x, y, w, h)

Dog + (x, y, w, h)

Cat + (x, y, w, h)

12 numbers to predict



Duck + (x, y, w, h)

Duck + (x, y, w, h)

...

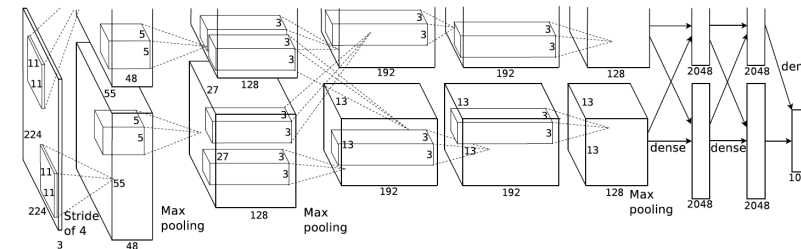
Many numbers to predict!!

Remember... The sliding window idea

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

Very inefficient! For a 800x600 images ~ **58M boxes**
Need to apply CNN to a huge number of locations!

Extract patches



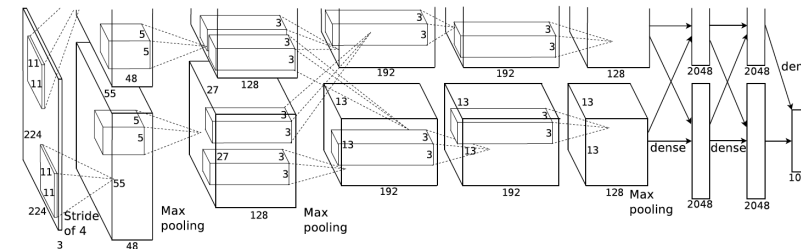
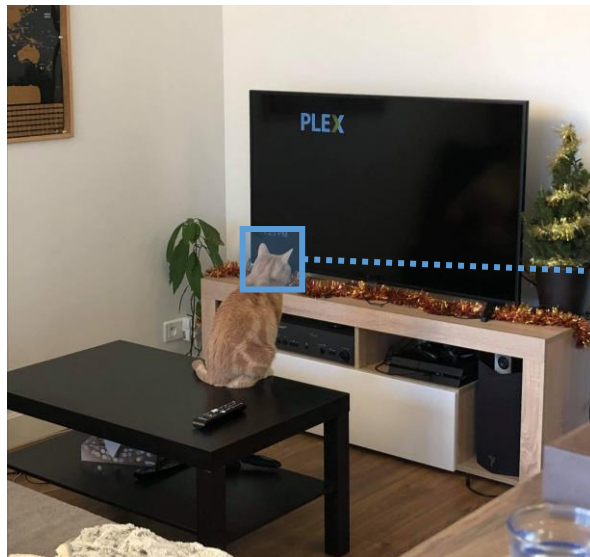
Cat? Yes
Table? No
Background? No

Remember... The sliding window idea

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

Very inefficient! For a 800x600 images ~ **58M boxes**
Need to apply CNN to a huge number of locations!

Extract patches



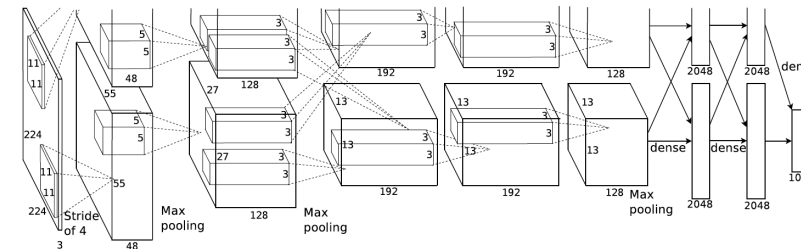
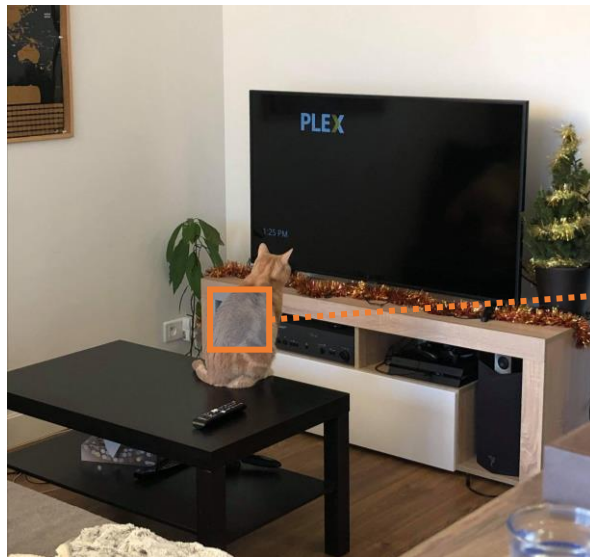
Cat? Yes
Table? No
Background? No

Remember... The sliding window idea

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

Very inefficient! For a 800x600 images ~ **58M boxes**
Need to apply CNN to a huge number of locations!

Extract patches



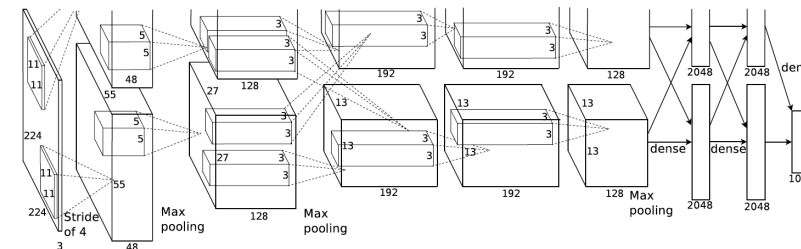
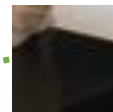
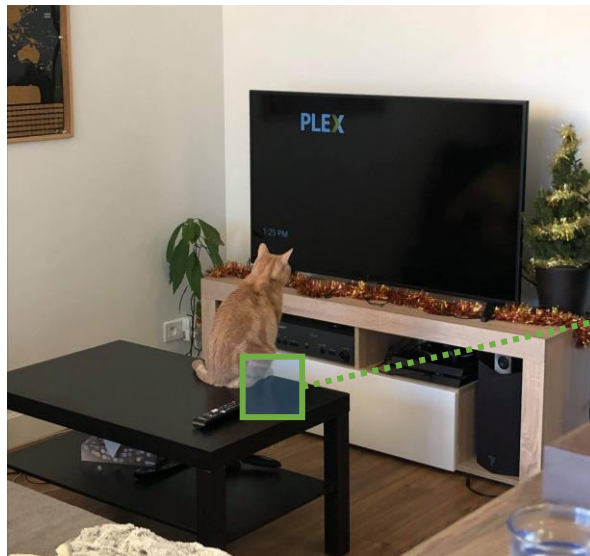
Cat? Yes
Table? No
Background? No

Remember... The sliding window idea

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

Very inefficient! For a 800x600 images ~ **58M boxes**
Need to apply CNN to a huge number of locations!

Extract patches



Cat? No
Table? Yes
Background? No

Object detection: region proposals

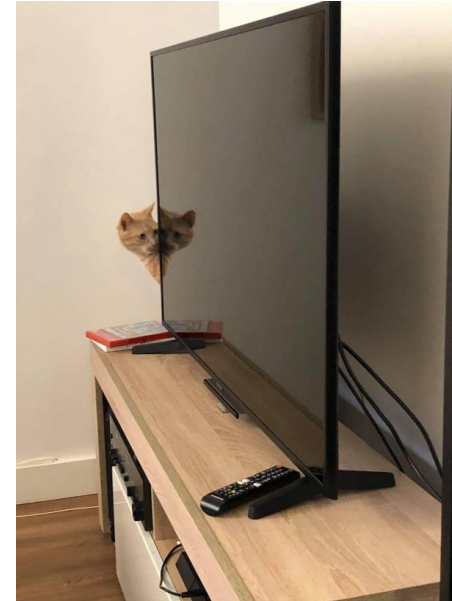
Selective search algorithm :

1. Generate initial sub-segmentation, many candidate regions generation
2. Use greedy algorithm to recursively combine similar region into larger ones
 1. From set of regions, choose two that are most similar.
 2. Combine them into a single, larger region.
 3. Repeat the above steps for multiple iterations.
3. Use the generated regions to produce the final candidate region proposals



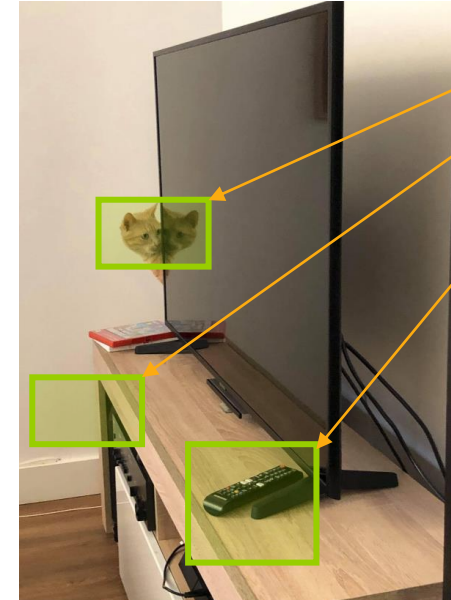
Object detection: R-CNN (Girshick et al., 2013)

- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Object detection: R-CNN (Girshick et al., 2013)

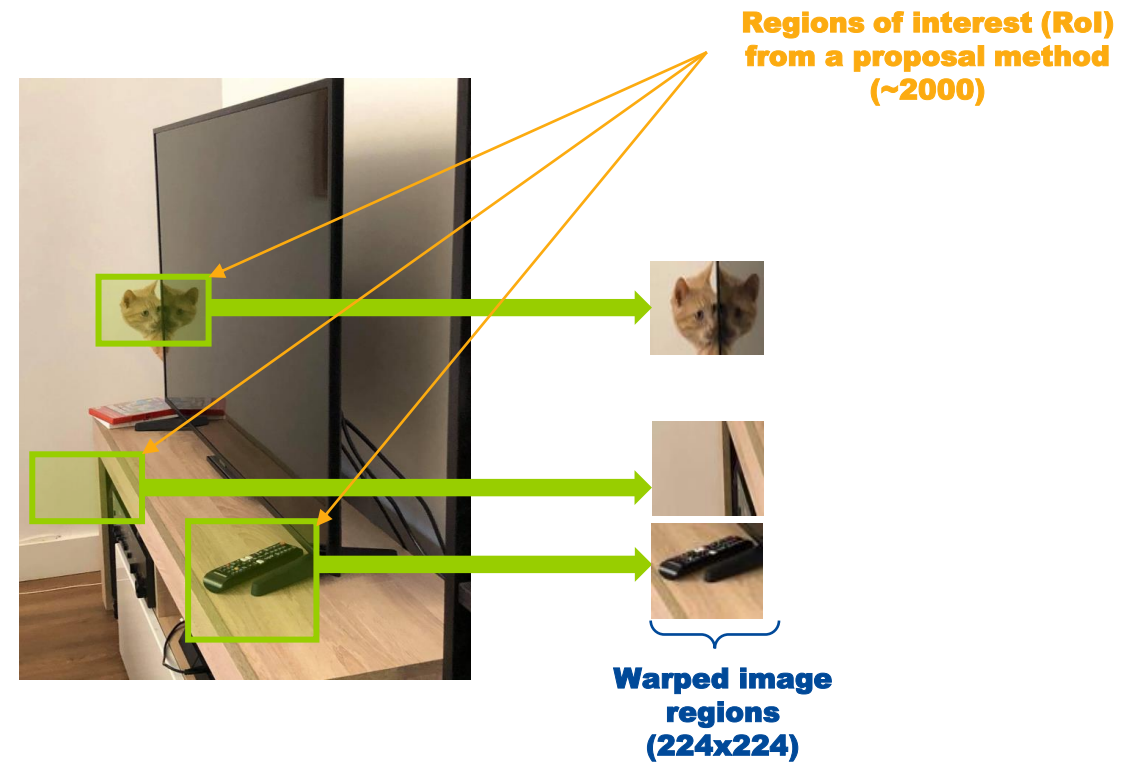
- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Regions of interest (RoI)
from a proposal method
(~2000)

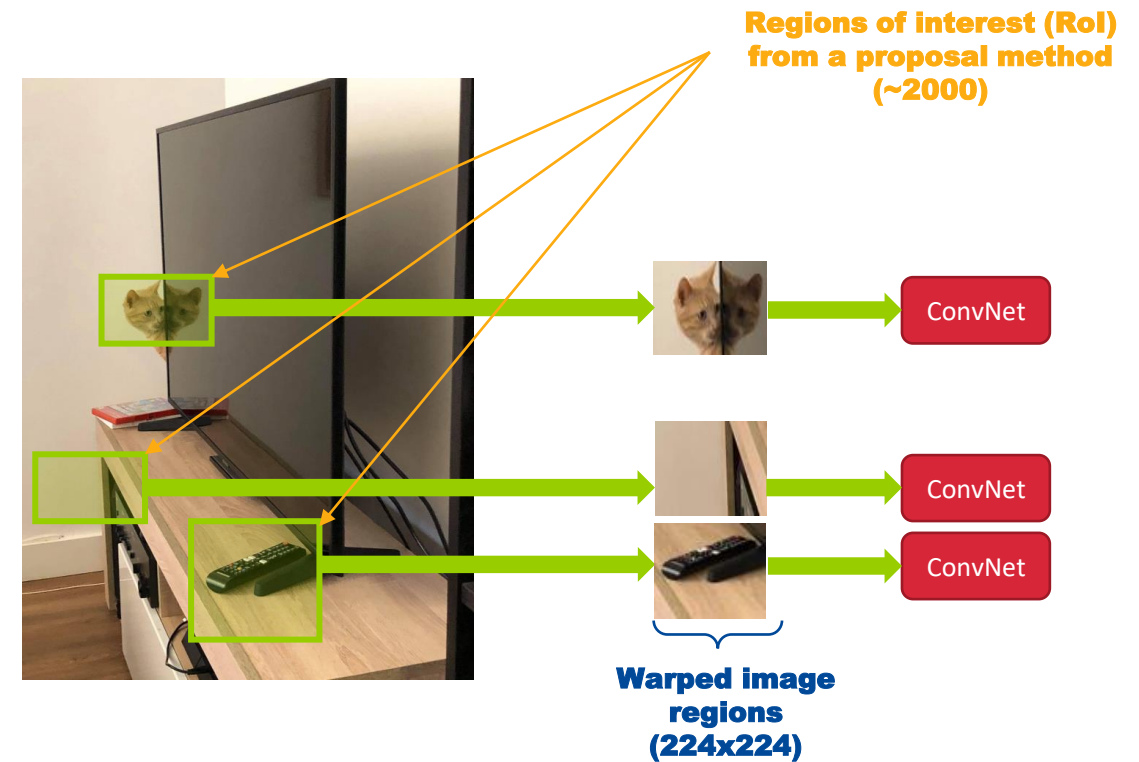
Object detection: R-CNN (Girshick et al., 2013)

- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



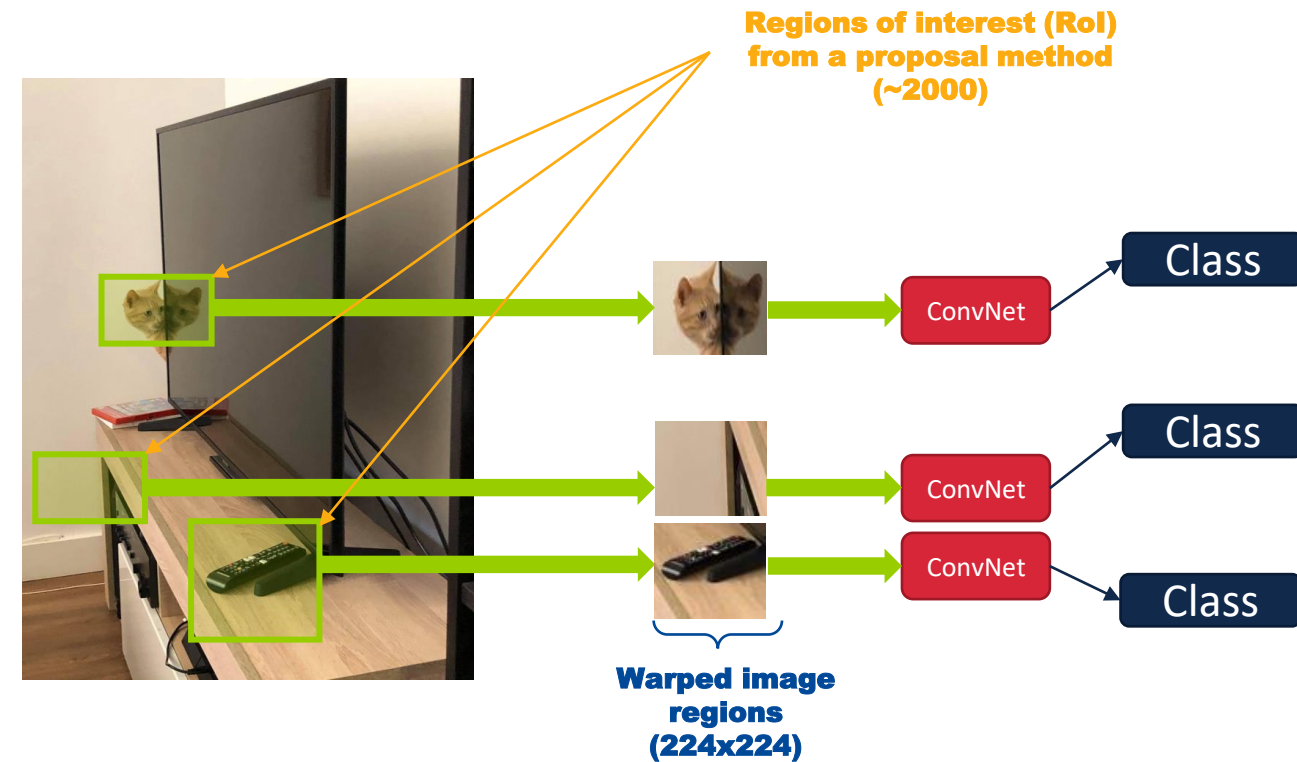
Object detection: R-CNN (Girshick et al., 2013)

- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



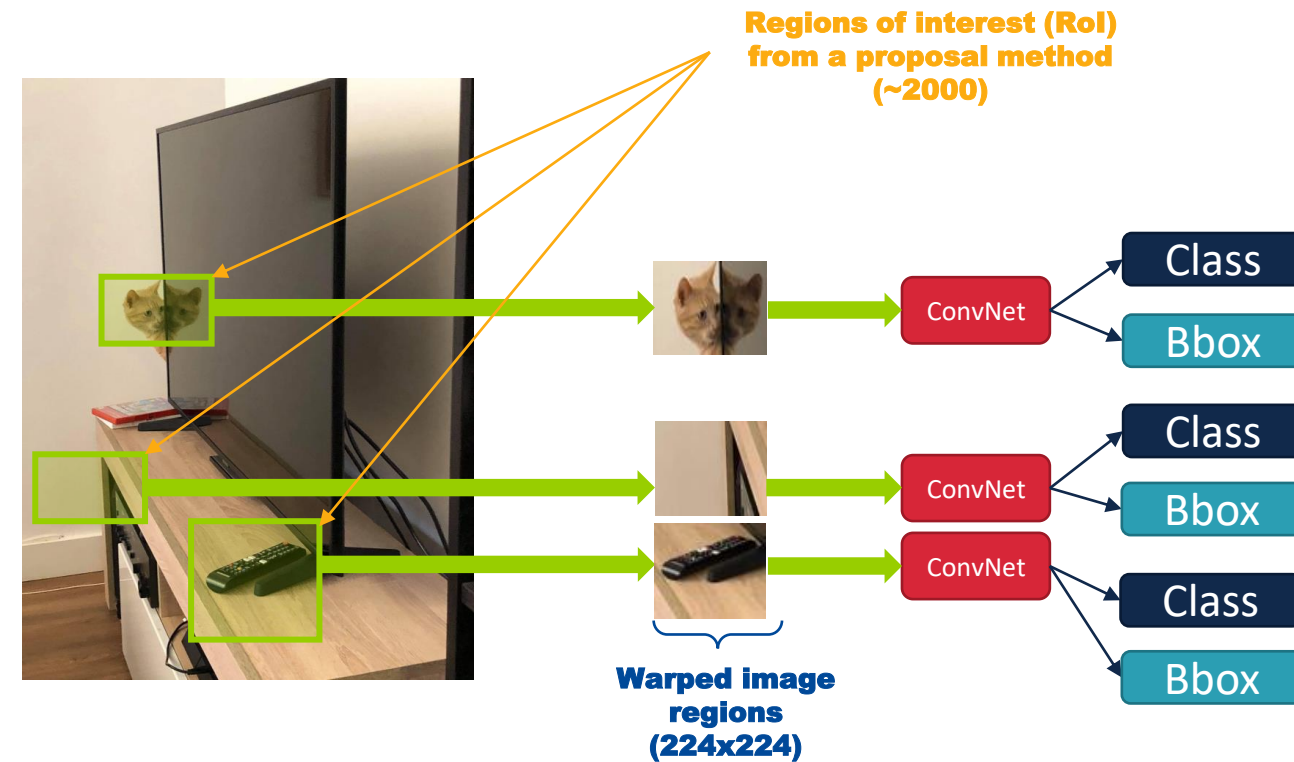
Object detection: R-CNN (Girshick et al., 2013)

- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Object detection: R-CNN (Girshick et al., 2013)

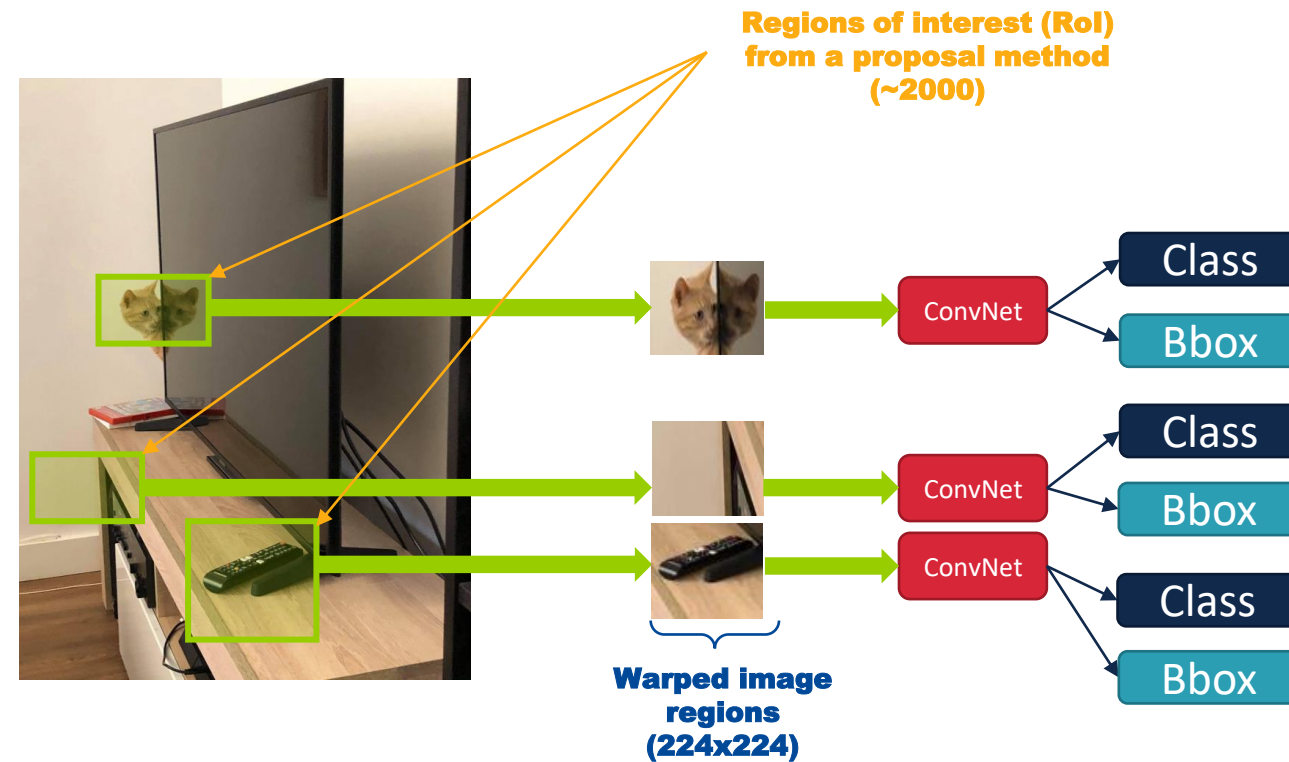
- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Object detection: R-CNN (Girshick et al., 2013)

Problem: very slow! Need 2000 independent forward pass for each image

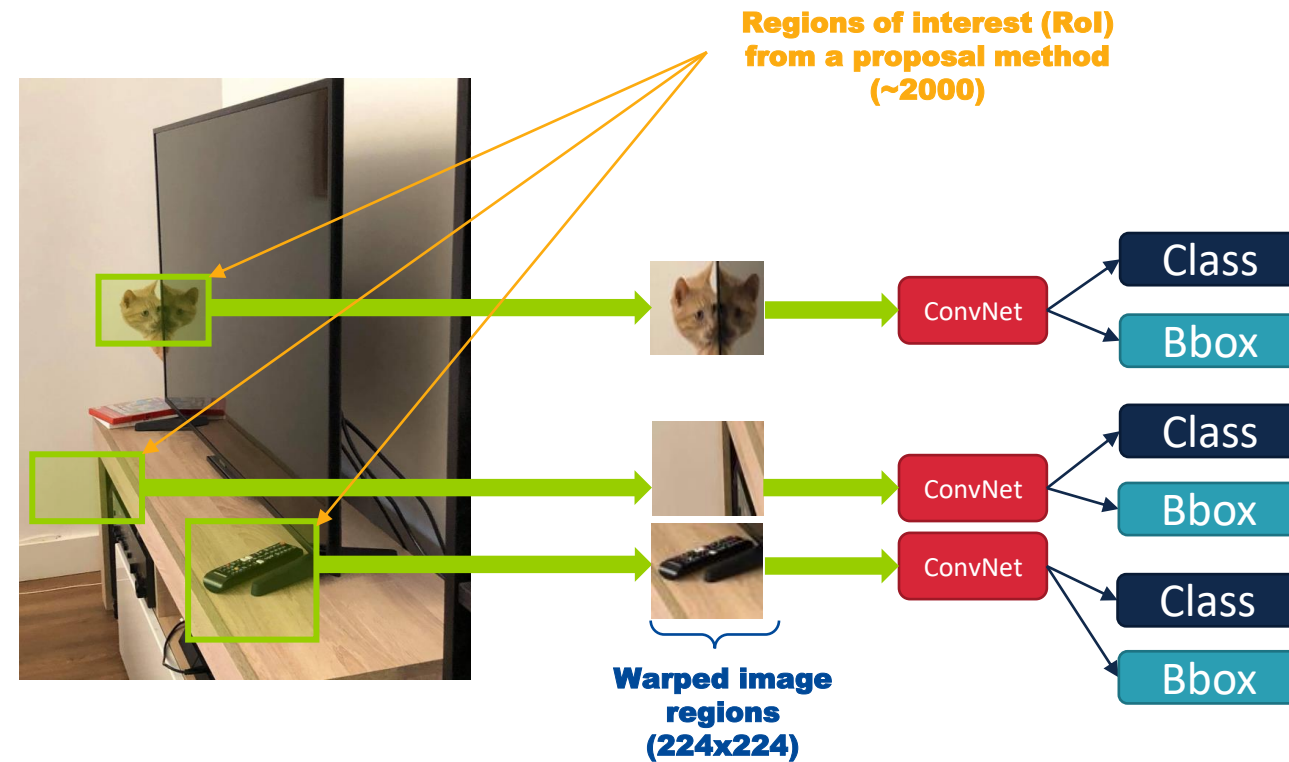
- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Object detection: “*Slow*” R-CNN (Girshick et al., 2013)

Problem: very slow! Need 2000 independent forward pass for each image

- Classify EACH proposed region (SVM)
- Bounding box regression:
 - Predict “transform” to correct the proposed RoI
 - 4 numbers: (t_x, t_y, t_h, t_w)
- Final output:
 - Proposal: (p_x, p_y, p_h, p_w)
 - Transform: (t_x, t_y, t_h, t_w) (a “correction”)
 - Output box: (b_x, b_y, b_h, b_w)
 - $b_x = p_x + p_w t_x$ and $b_y = p_y + p_h t_y$
 - $b_w = p_w e^{t_w}$ and $b_h = p_h e^{t_h}$



Object detection: Fast-RCNN (Girshick et al., 2015)



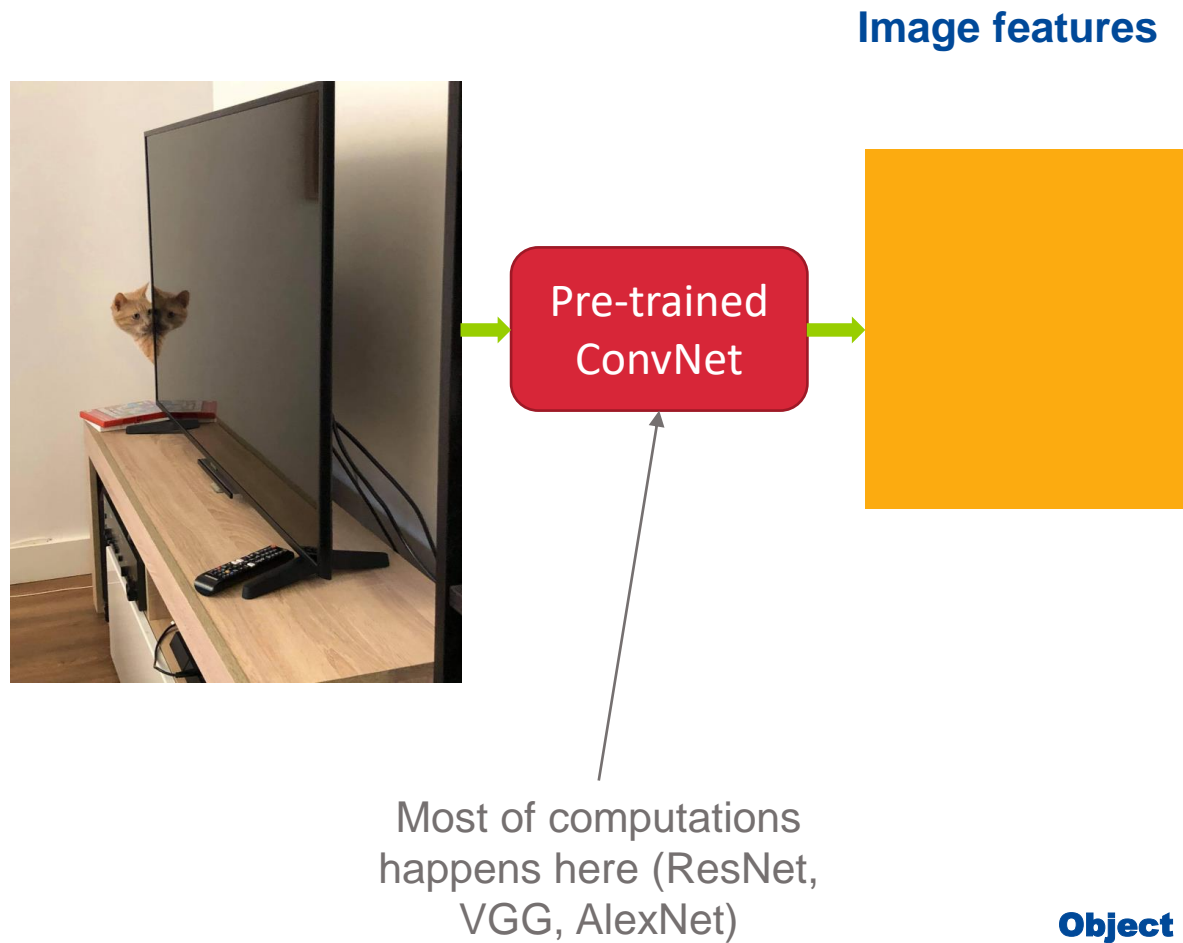
Object detection: Fast-RCNN (Girshick et al., 2015)



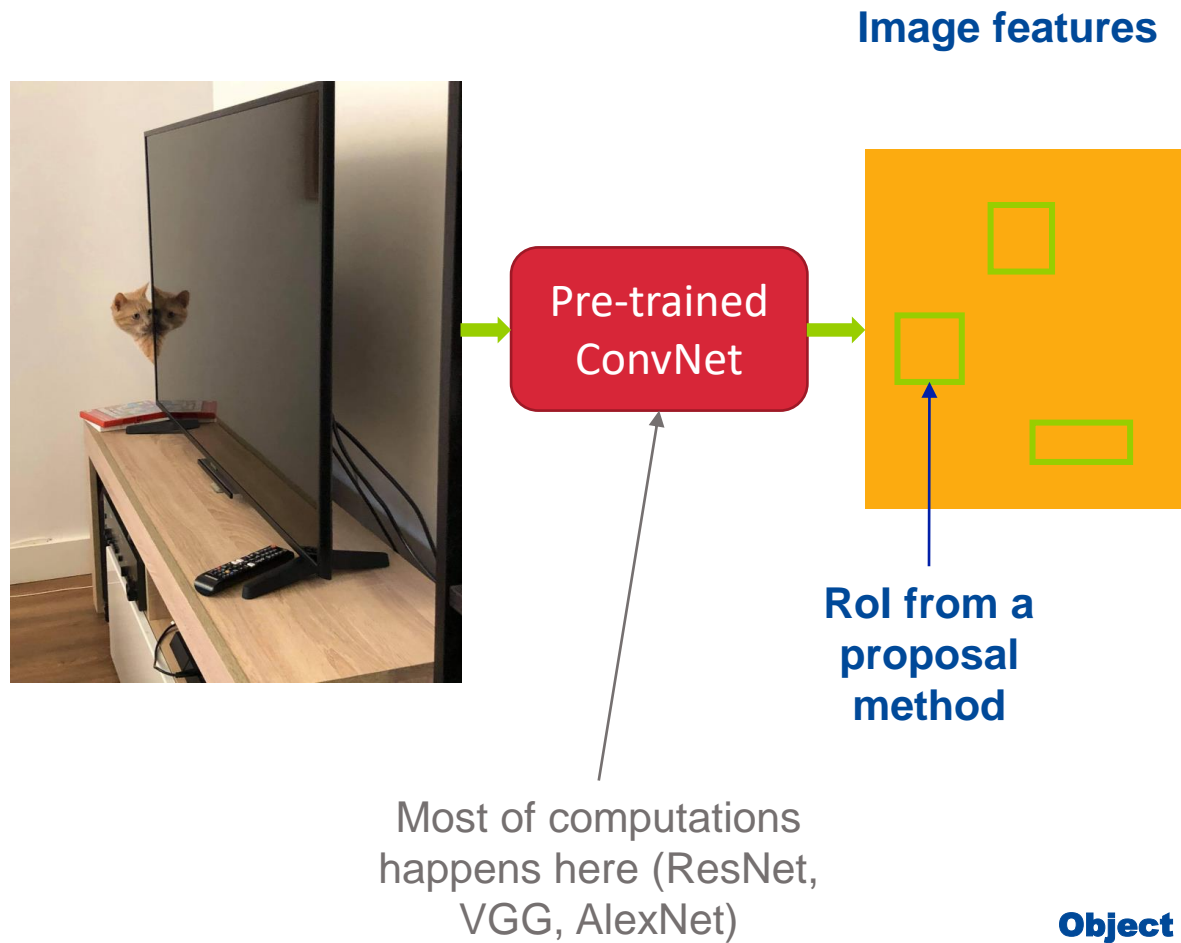
Pre-trained
ConvNet

Most of computations
happens here (ResNet,
VGG, AlexNet)

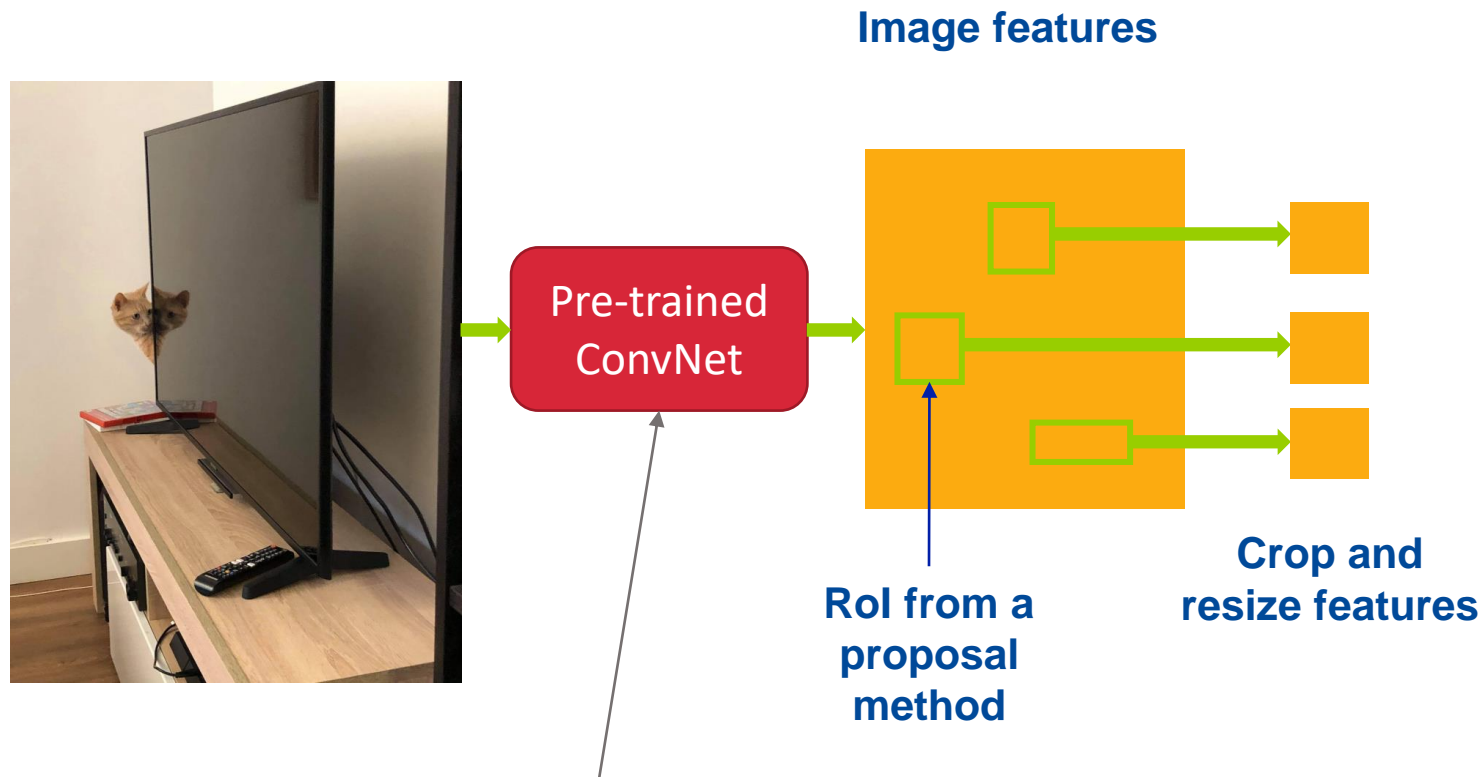
Object detection: Fast-RCNN (Girshick et al., 2015)



Object detection: Fast-RCNN (Girshick et al., 2015)

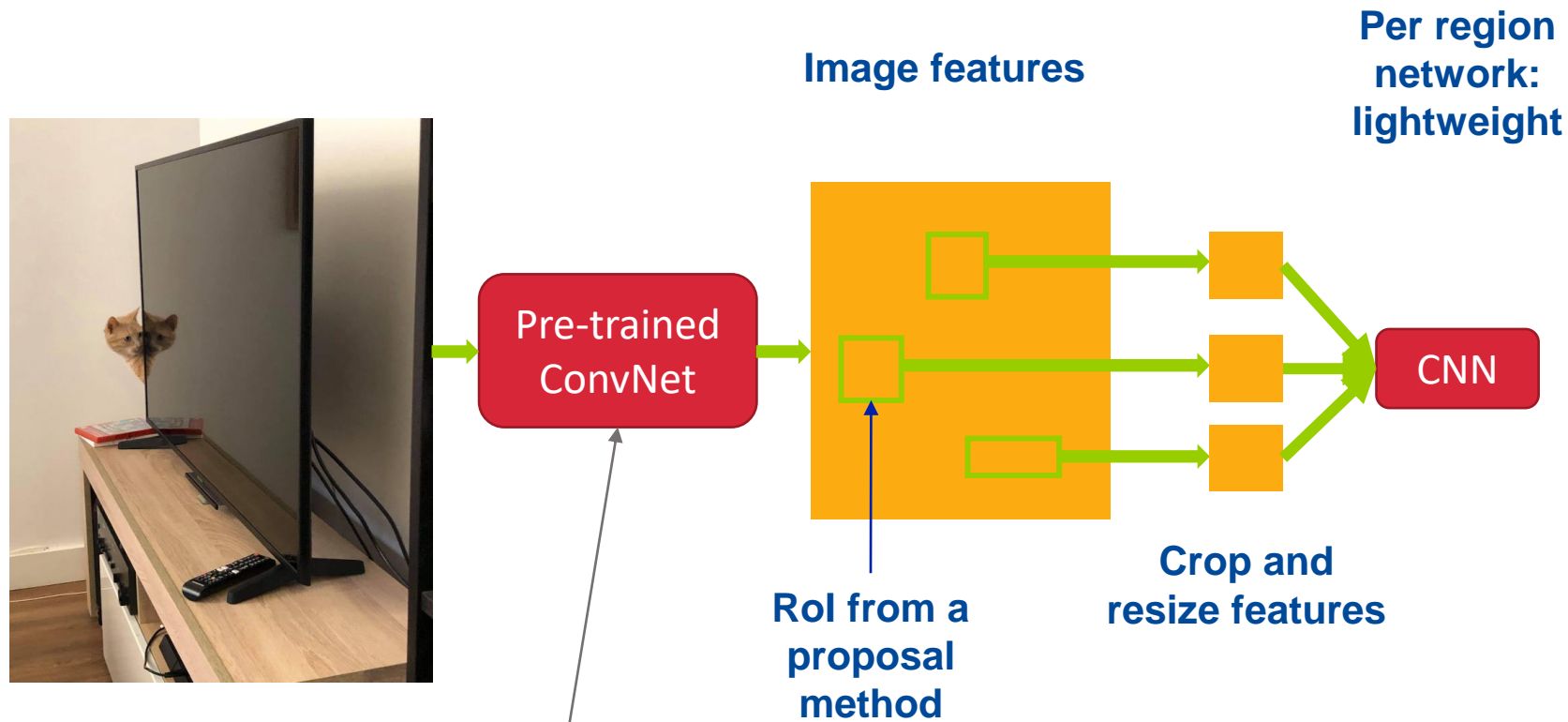


Object detection: Fast-RCNN (Girshick et al., 2015)



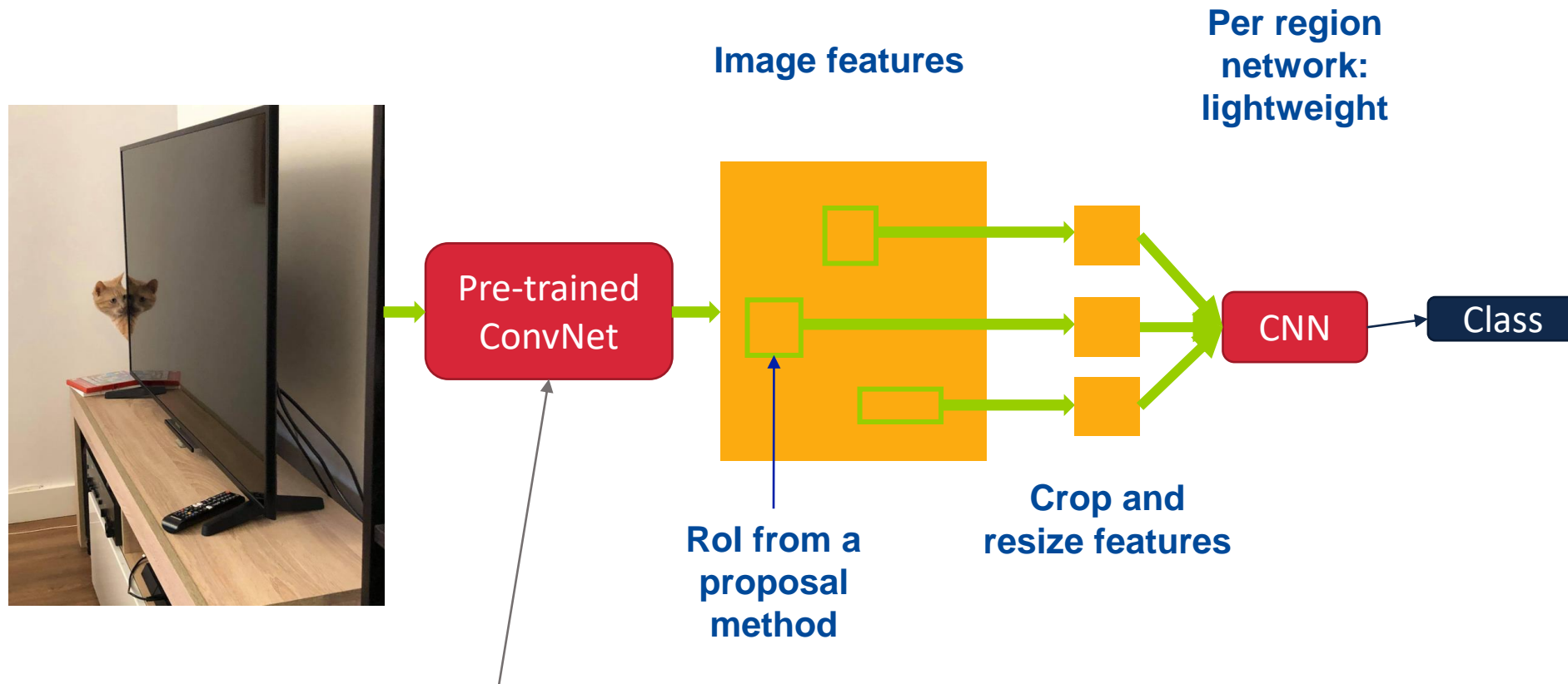
Most of computations happens here (ResNet, VGG, AlexNet)

Object detection: Fast-RCNN (Girshick et al., 2015)



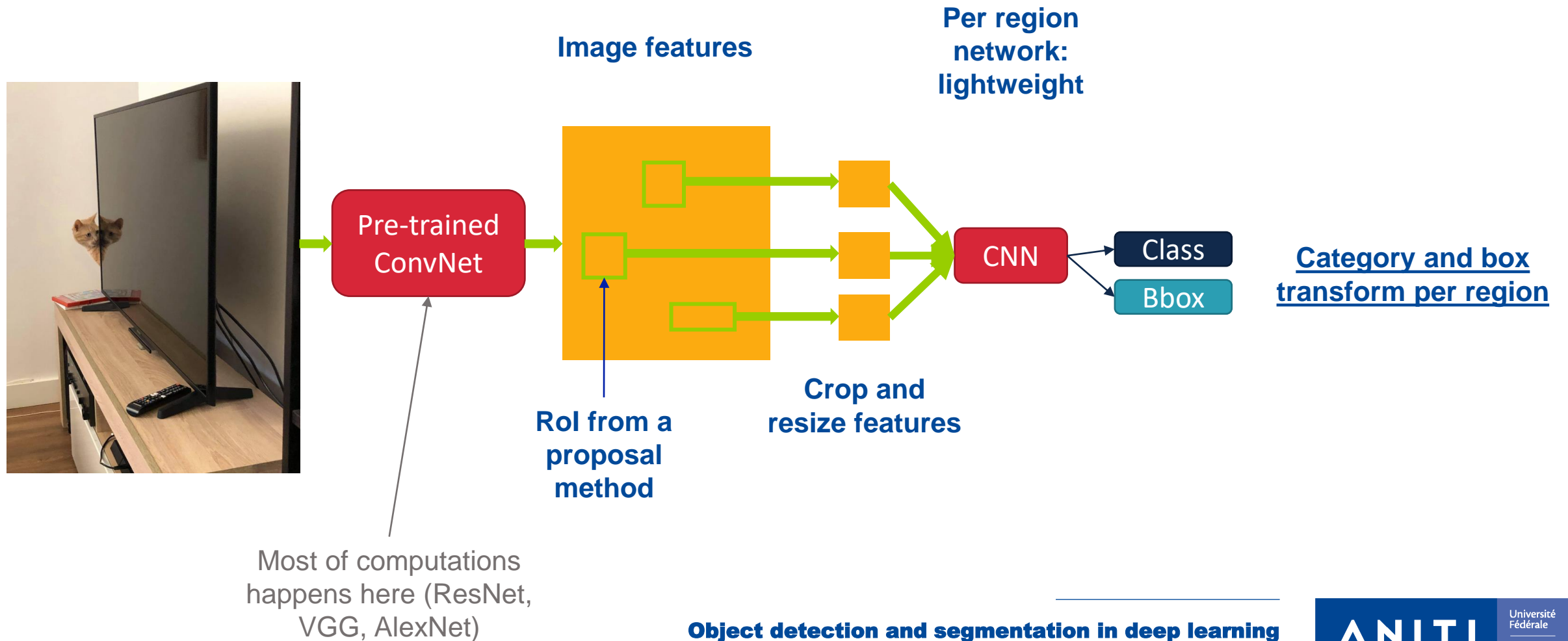
Most of computations happens here (ResNet, VGG, AlexNet)

Object detection: Fast-RCNN (Girshick et al., 2015)

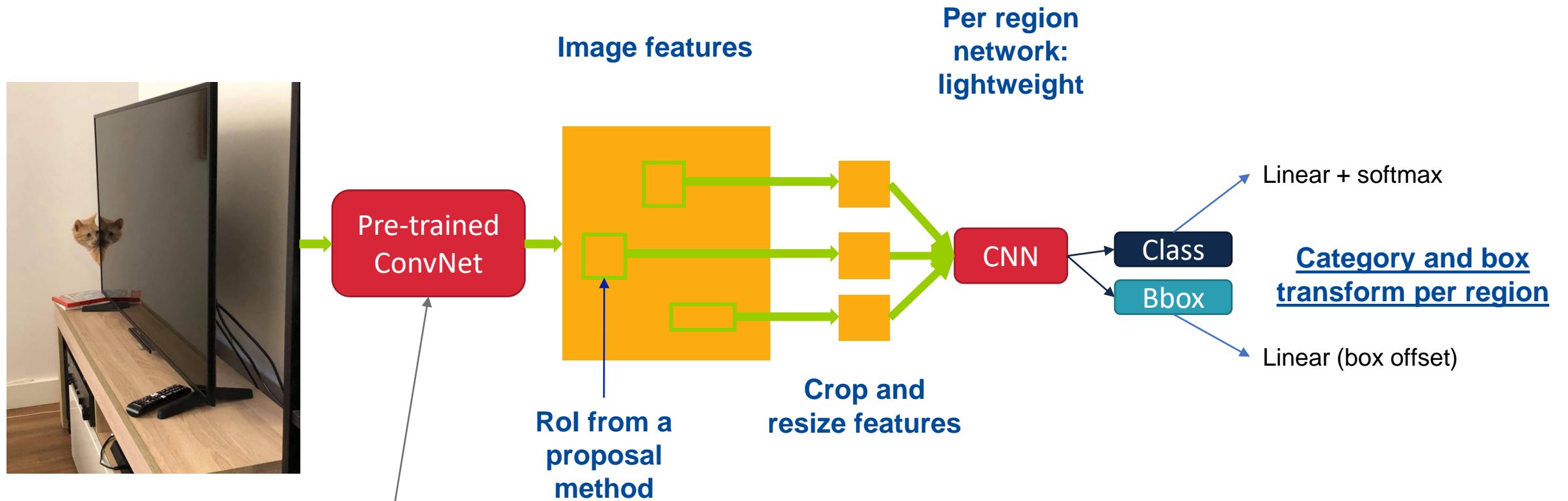


Most of computations happens here (ResNet, VGG, AlexNet)

Object detection: Fast-RCNN (Girshick et al., 2015)



Object detection: Fast-RCNN (Girshick et al., 2015)

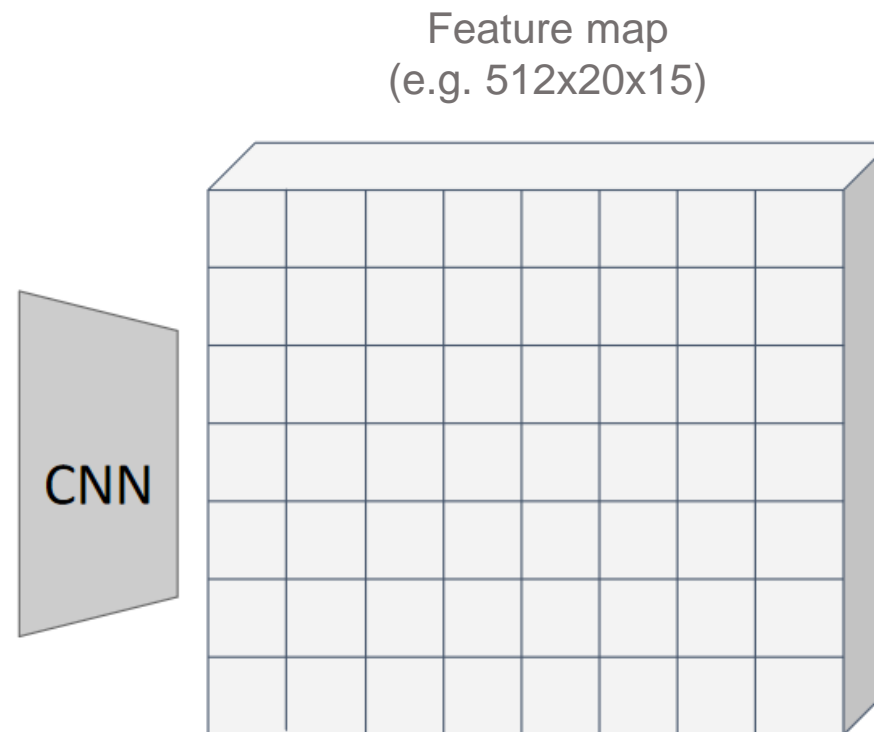


Most of computations happens here (ResNet, VGG, AlexNet)

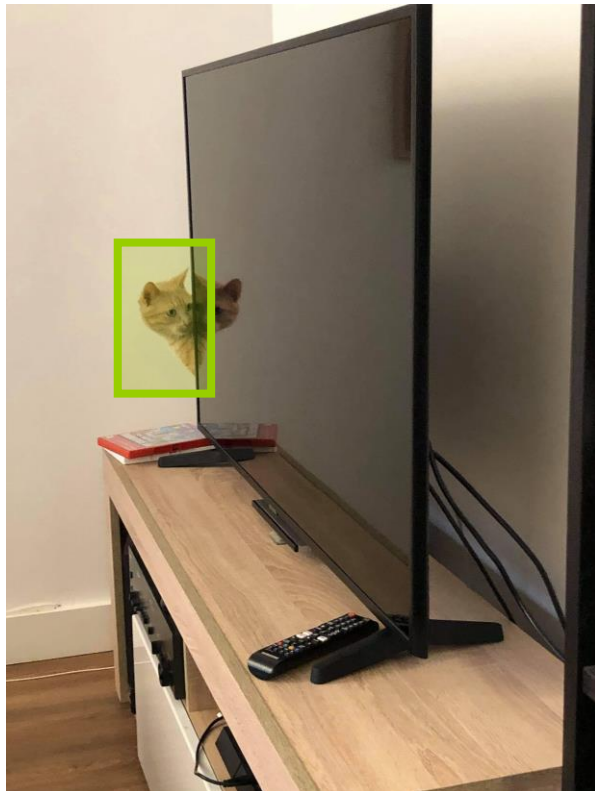
Rol pooling



Input image (e.g.
 $3 \times 640 \times 480$)

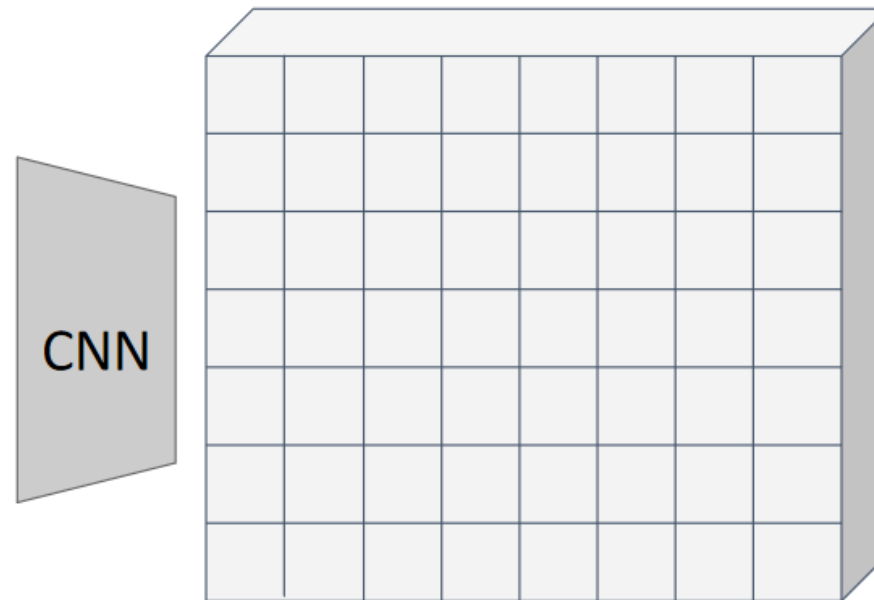


Rol pooling

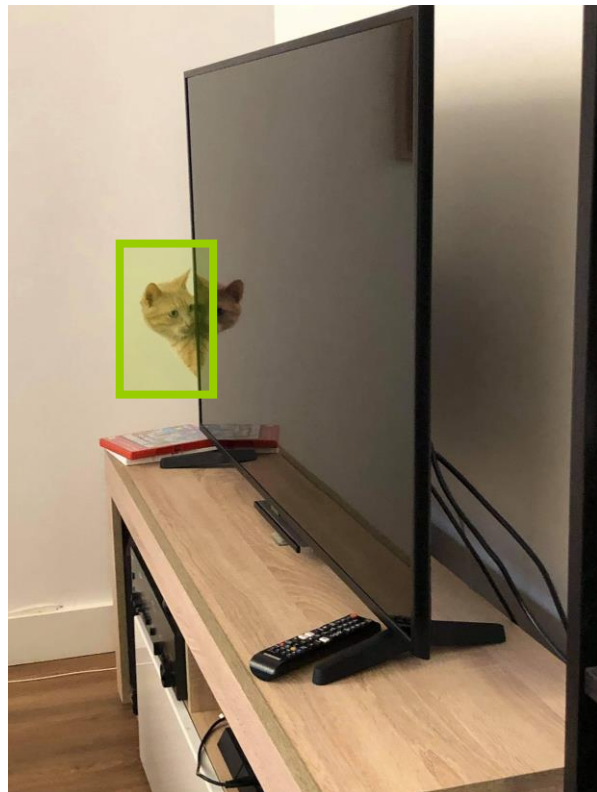


Input image (e.g.
 $3 \times 640 \times 480$)

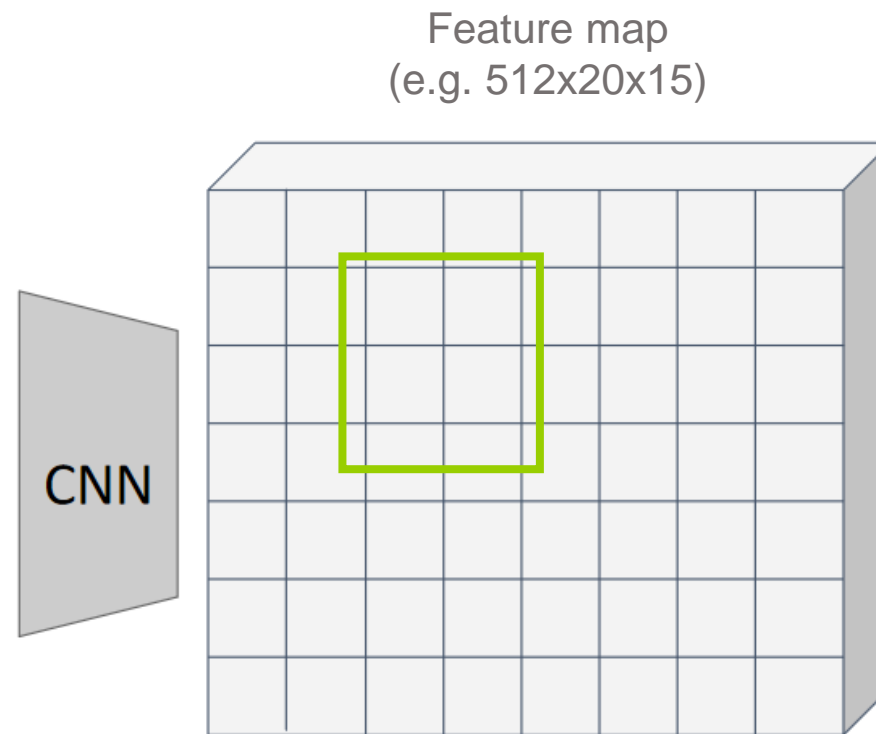
Feature map
(e.g. $512 \times 20 \times 15$)



Rol pooling

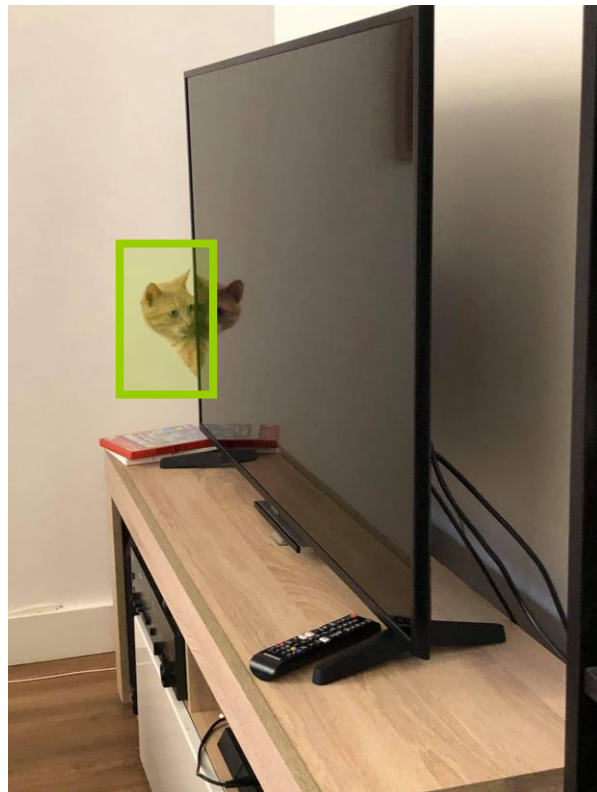


Input image (e.g.
3x640x480)

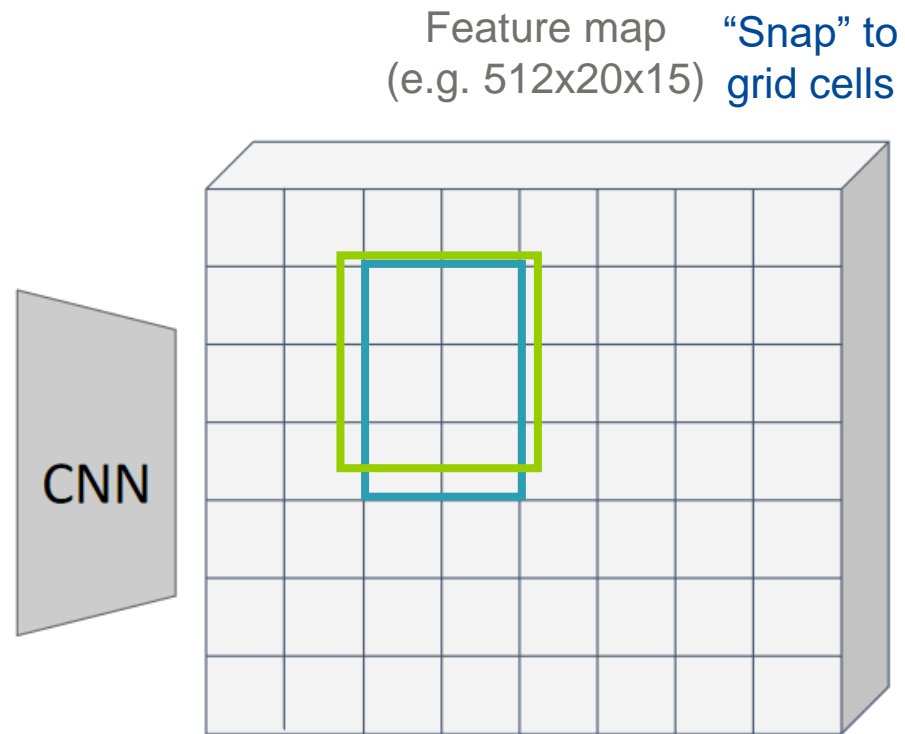


Project proposal onto
features

Rol pooling

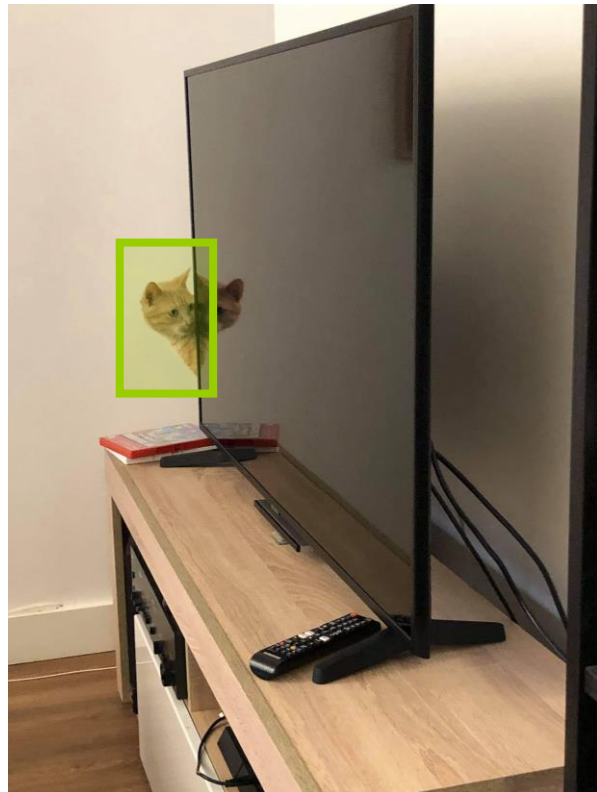


Input image (e.g.
3x640x480)



Project proposal onto
features

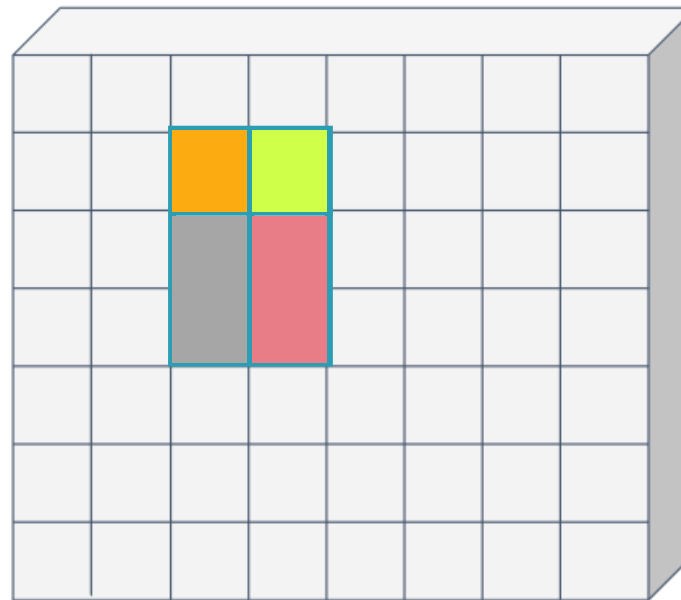
Rol pooling



Input image (e.g.
3x640x480)



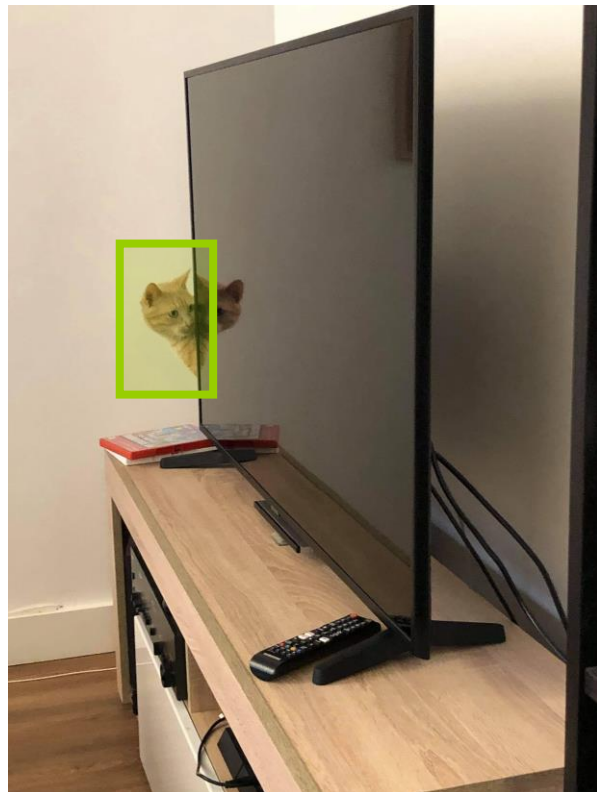
Feature map
(e.g. 512x20x15)



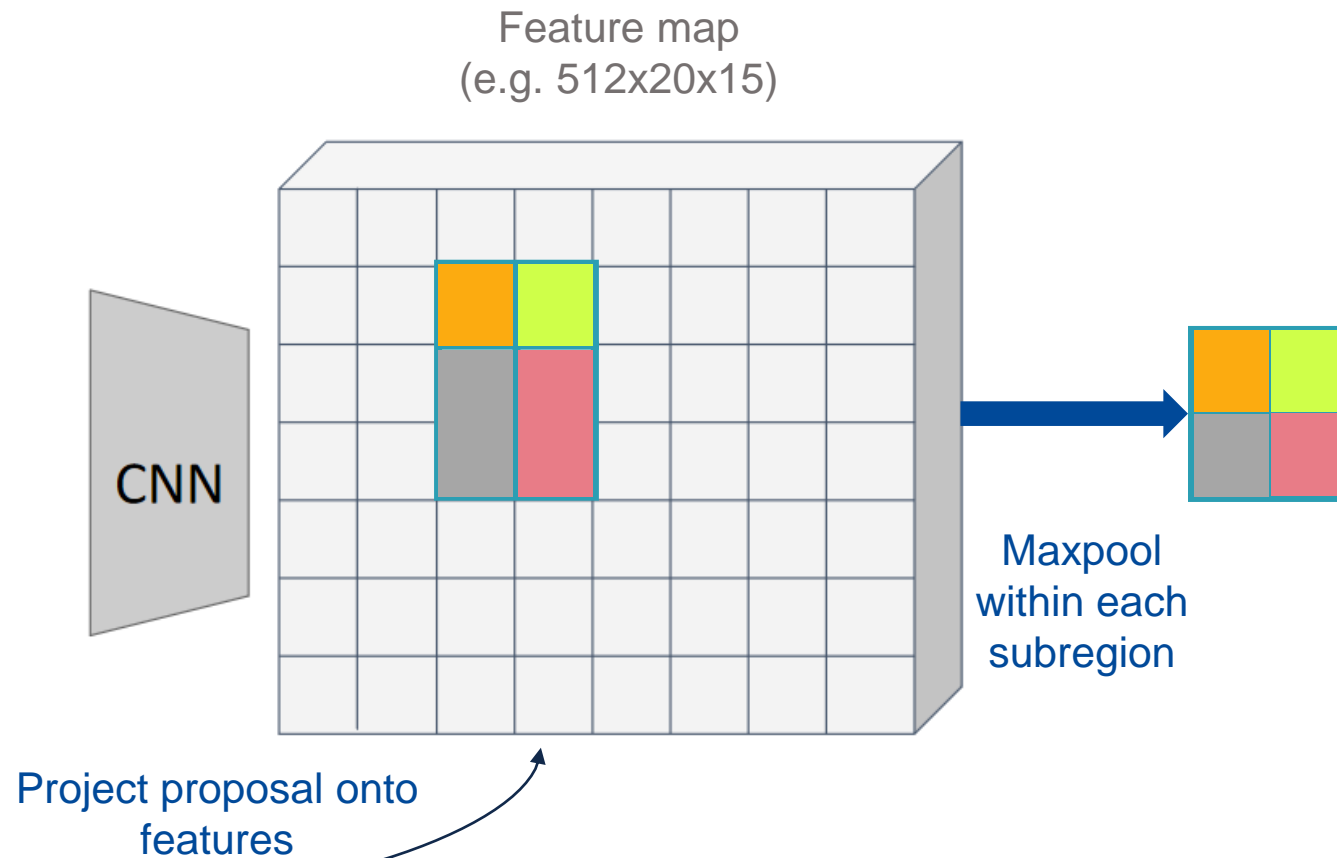
Divide into 2x2 grid

Project proposal onto
features

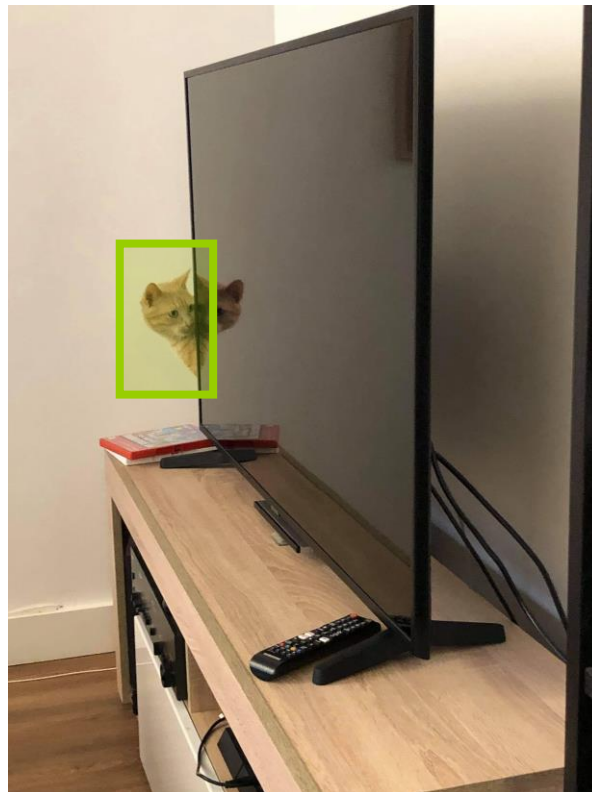
Rol pooling



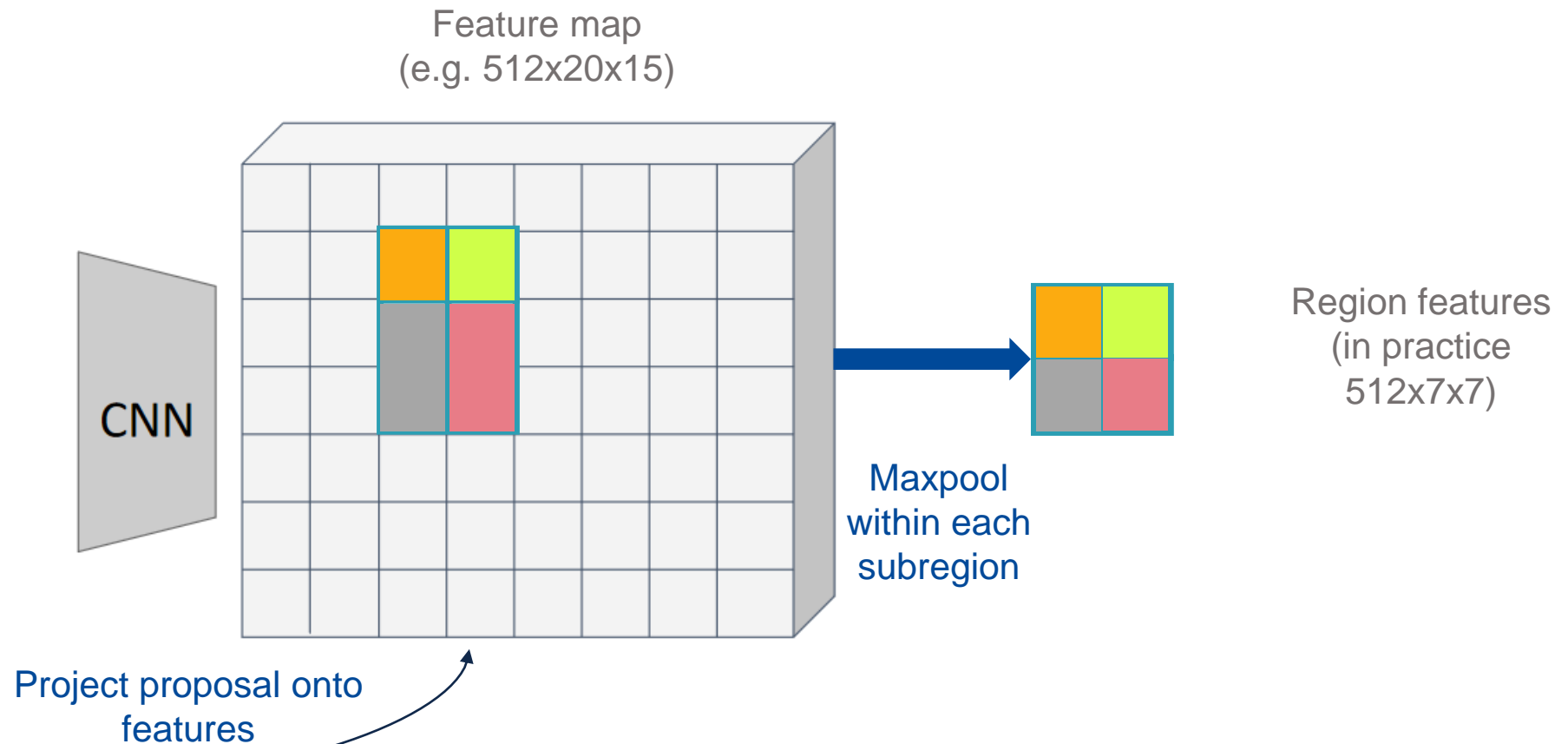
Input image (e.g.
3x640x480)



Rol pooling

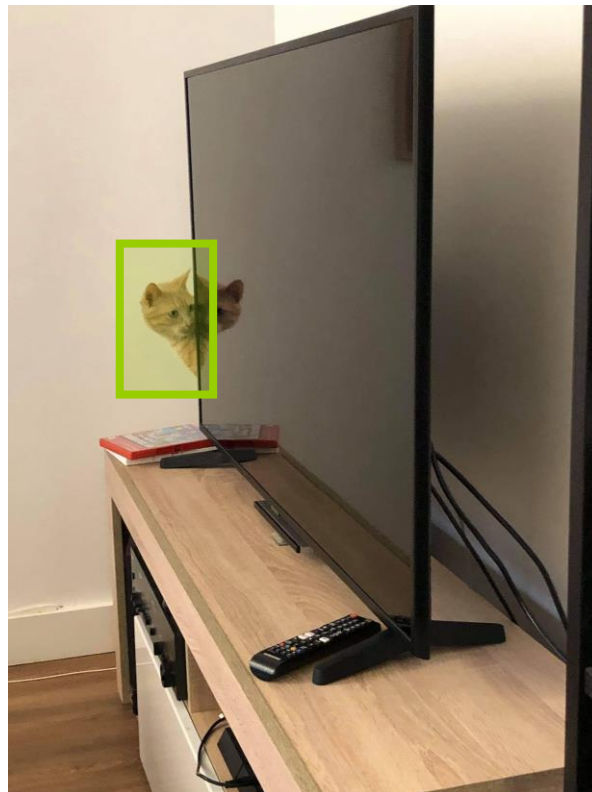


Input image (e.g.
3x640x480)

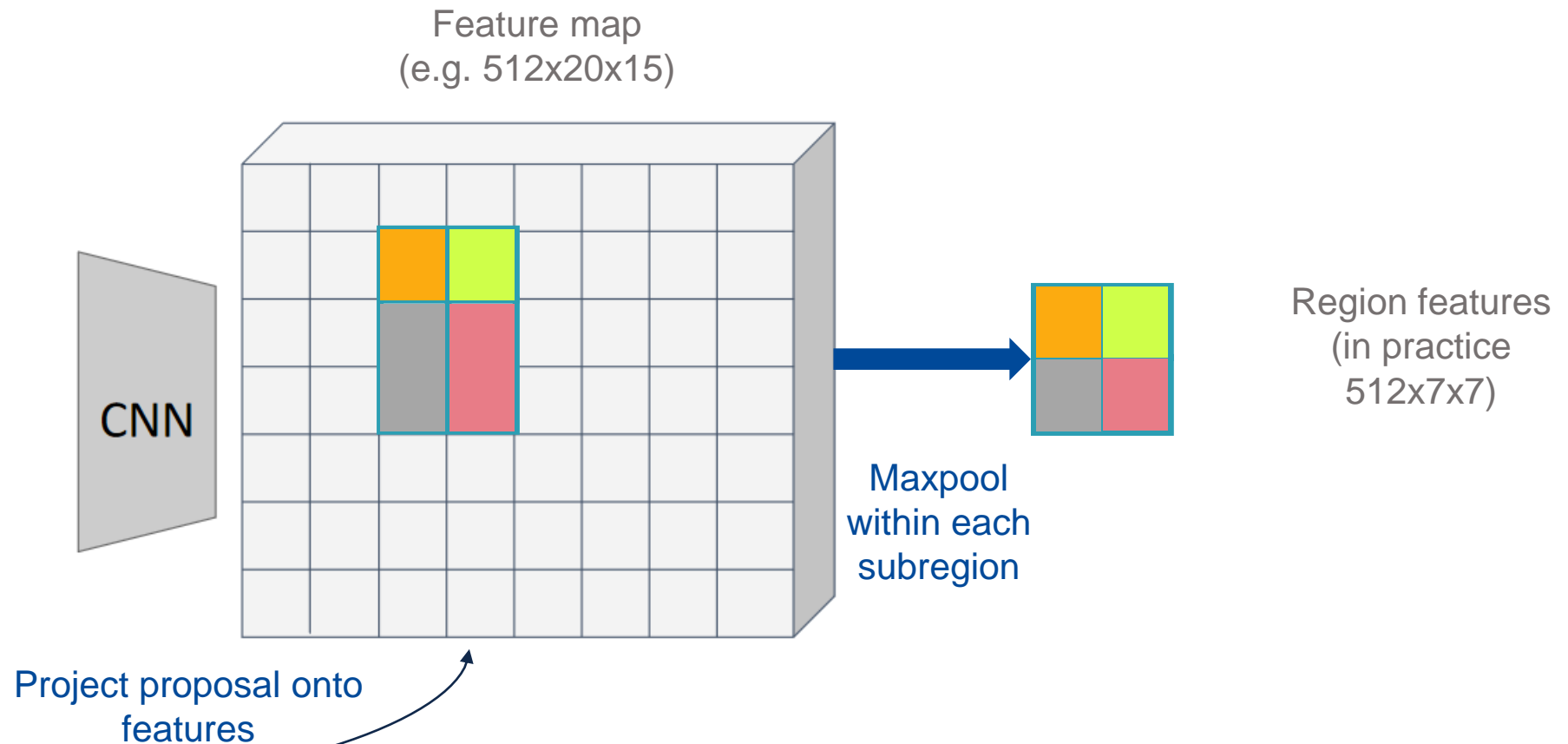


Rol pooling

Problem: region features might be slightly misaligned → Rol align

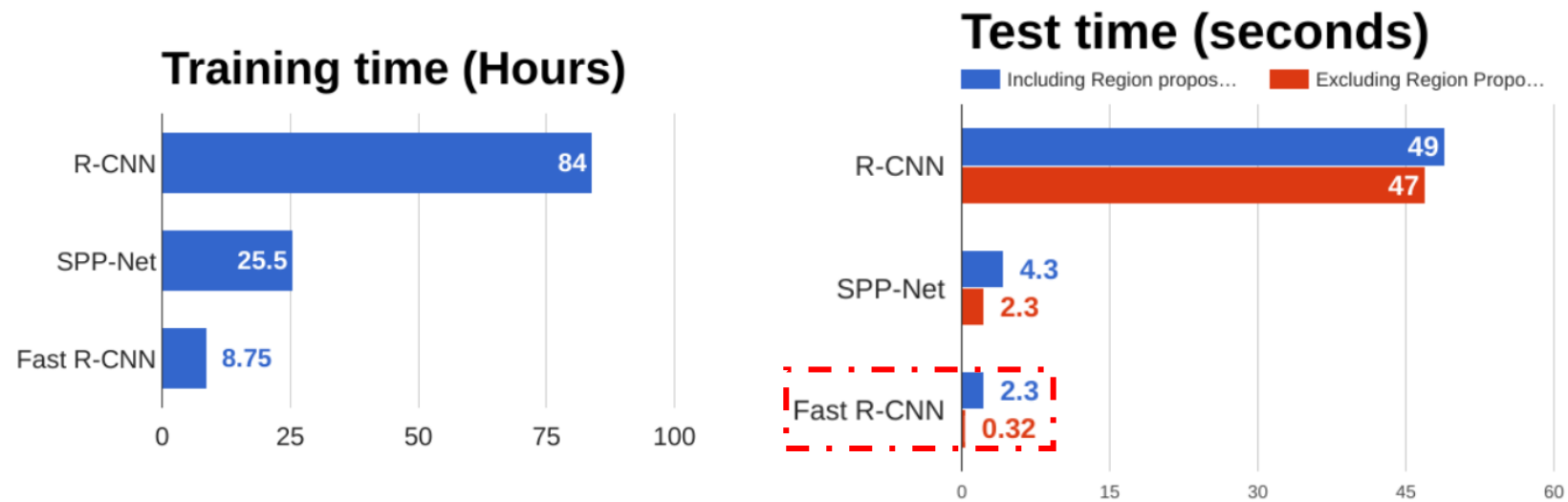


Input image (e.g.
3x640x480)



Object detection: Fast R-CNN vs. “Slow” R-CNN

Runtime dominated by region proposals

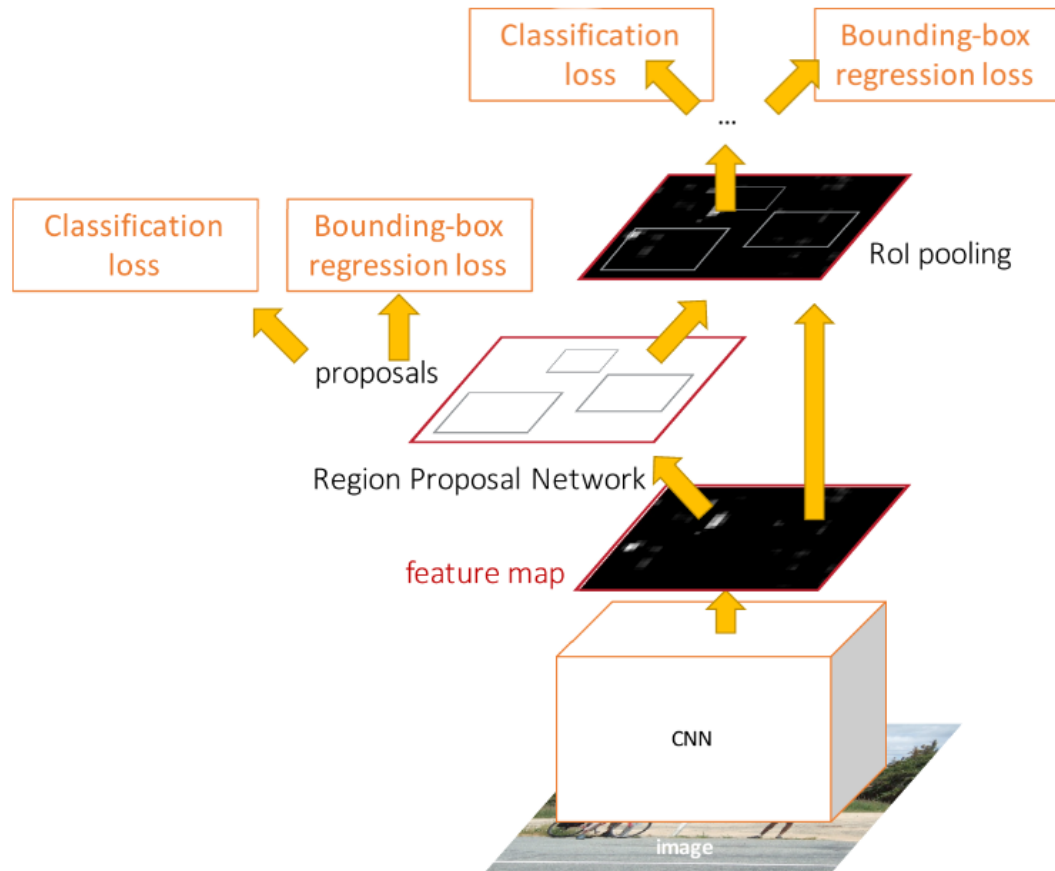


Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

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

Girshick, “Fast R-CNN”, ICCV 2015

Object detection: learn to *propose* regions, Faster R-CNN (Ren et al., 2015)



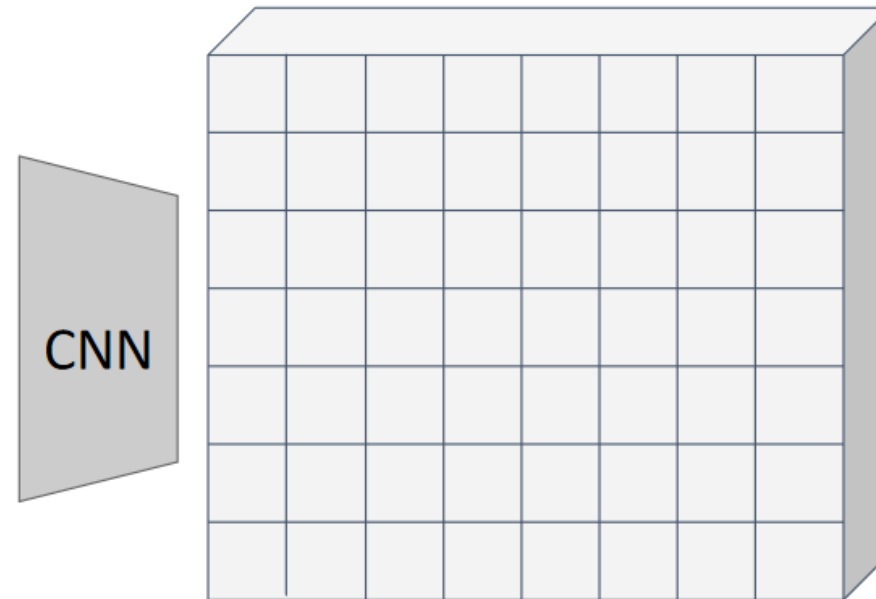
- Idea: insert a **Region Proposal Network (RPN)** to predict proposals from features
- Otherwise same as Fast R-CNN: crop features for each proposal, classify each one

Object detection: Region Proposal Network (RPN)

Imagine an **anchor box** of fixed size at each point in the feature map



Input image (e.g.
3x640x480)

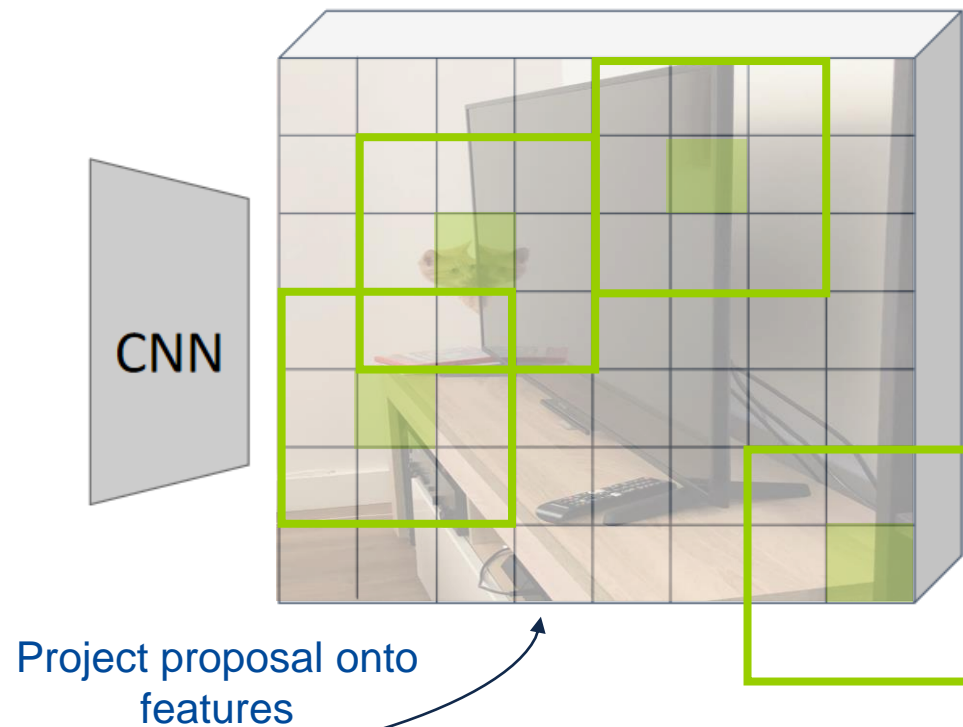


Object detection: Region Proposal Network (RPN)

Imagine an **anchor box** of fixed size at each point in the feature map



Input image (e.g.
3x640x480)

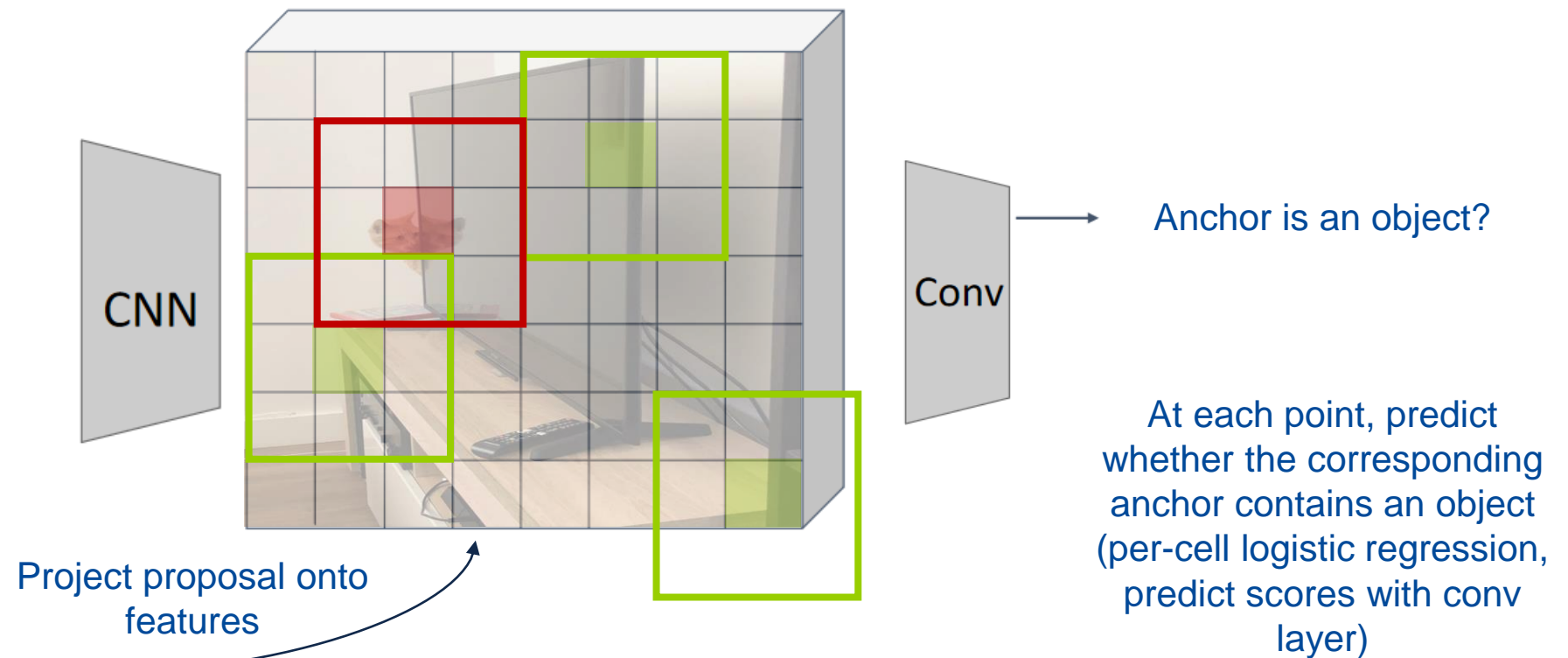


Object detection: Region Proposal Network (RPN)

Imagine an **anchor box** of fixed size at each point in the feature map



Input image (e.g.
3x640x480)

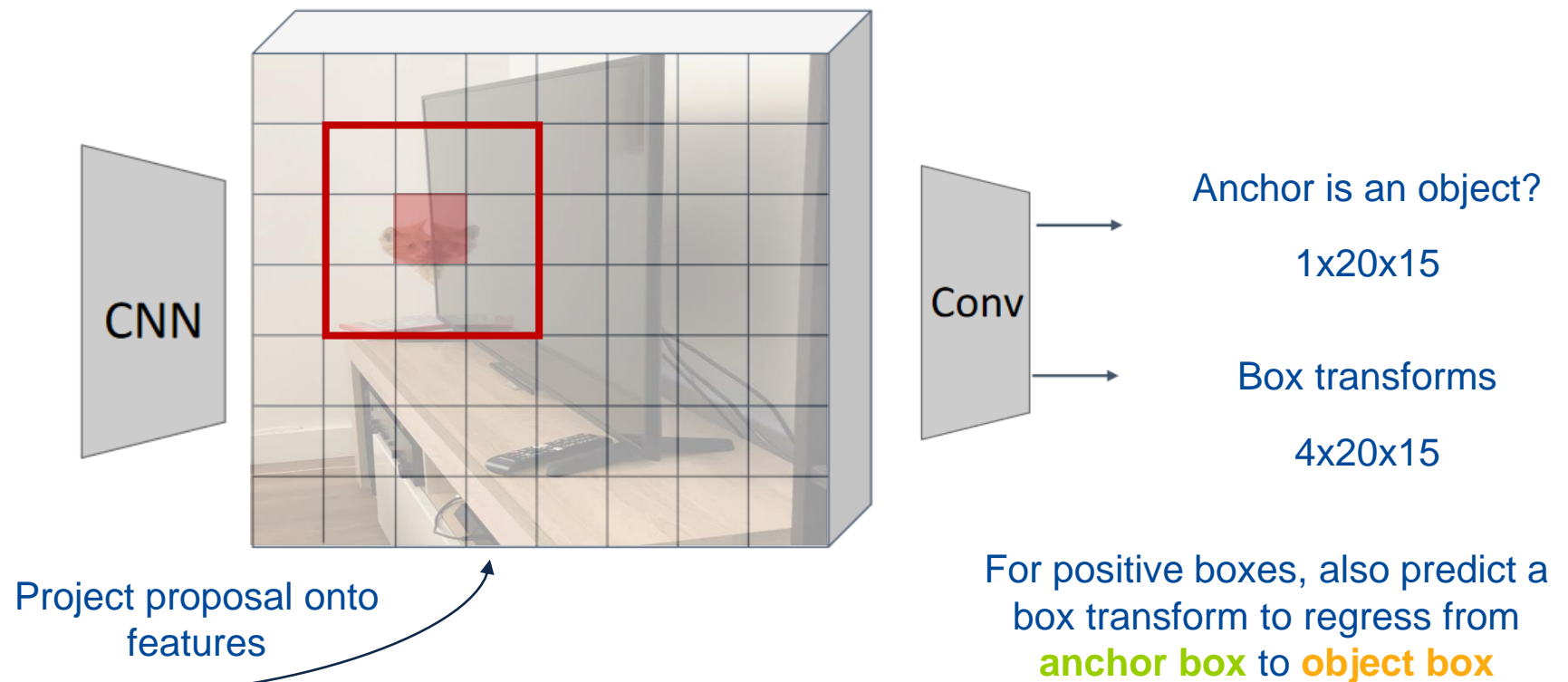


Object detection: Region Proposal Network (RPN)

Imagine an **anchor box** of fixed size at each point in the feature map



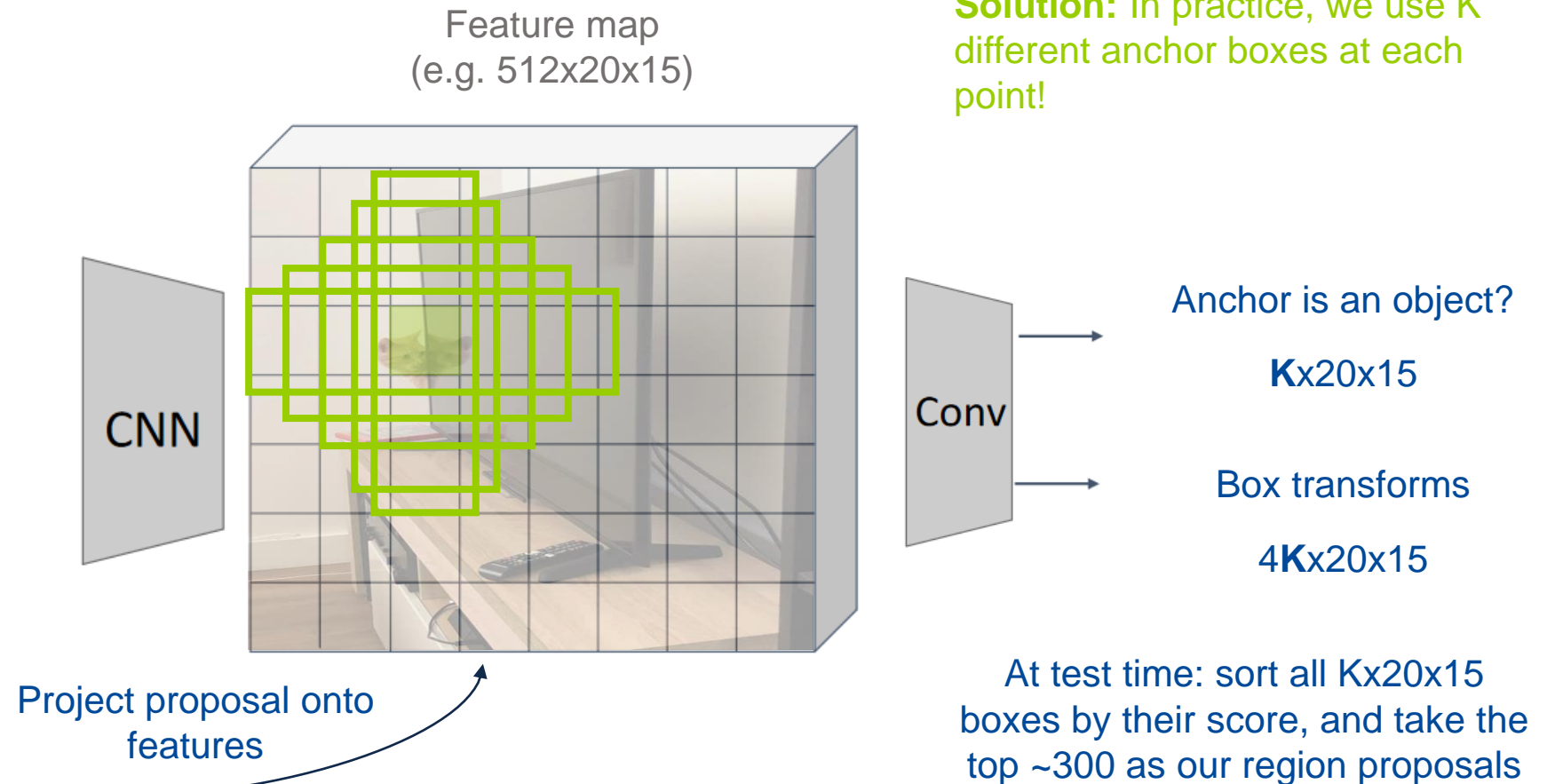
Input image (e.g.
3x640x480)



Object detection: Region Proposal Network (RPN)



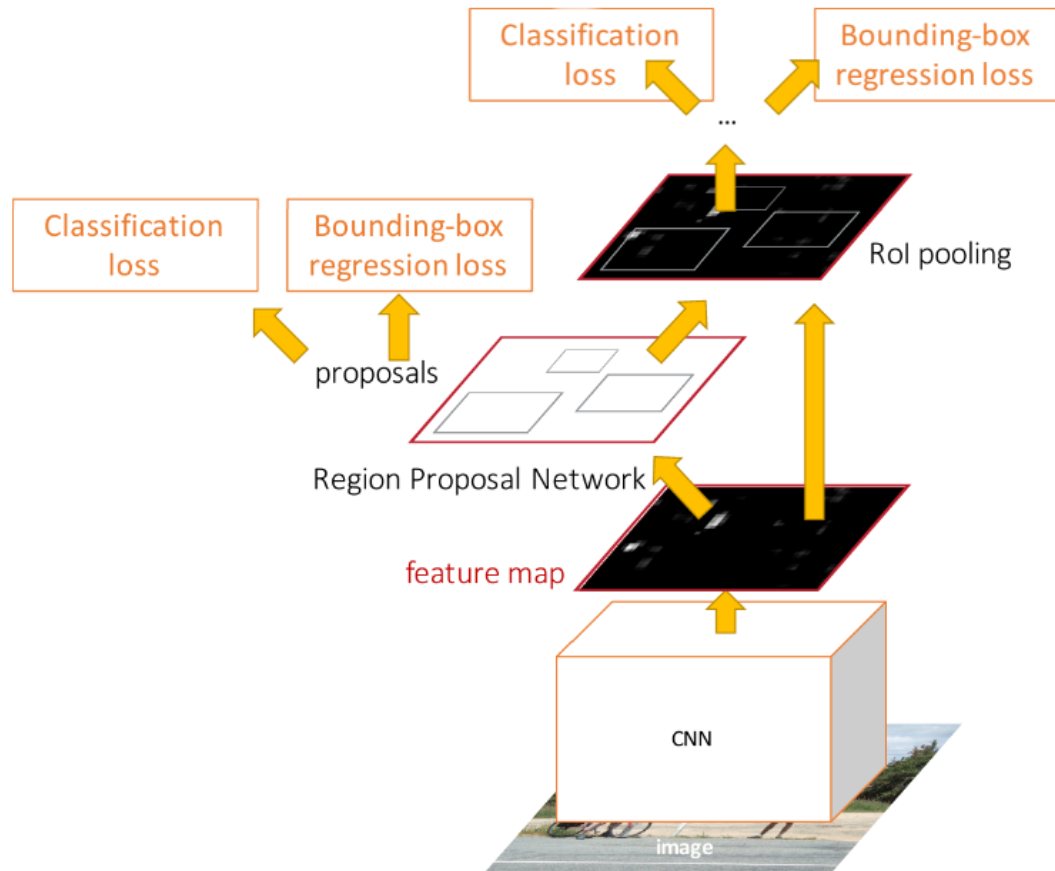
Input image (e.g.
 $3 \times 640 \times 480$)



Problem: Anchor boxes may have the wrong size/shape

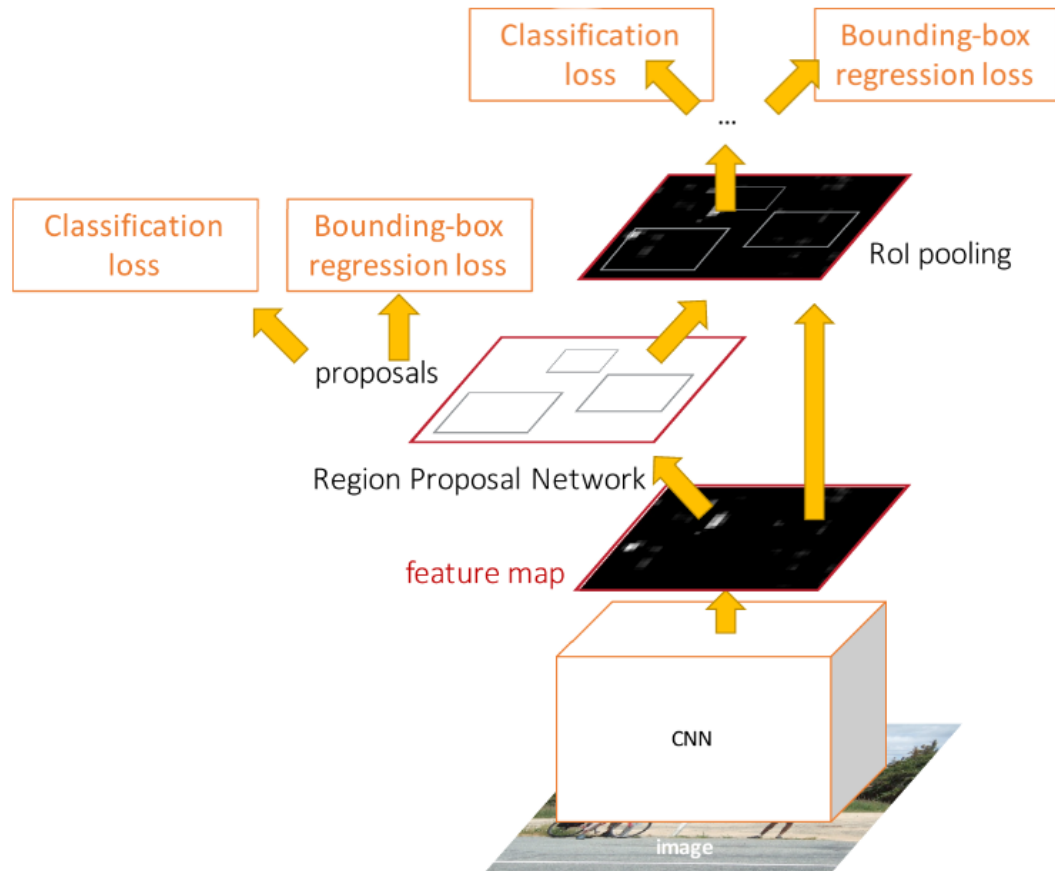
Solution: In practice, we use K different anchor boxes at each point!

Object detection: learn to *propose* regions, Faster R-CNN (Ren et al., 2015)



- Jointly train with 4 losses:
 1. **RPN classification**: anchor box is object / not an object
 2. **RPN regression**: predict transform from anchor box to proposal box
 3. **Object classification**: classify proposals as background / object class
 4. **Object regression**: predict transform from proposal to object box

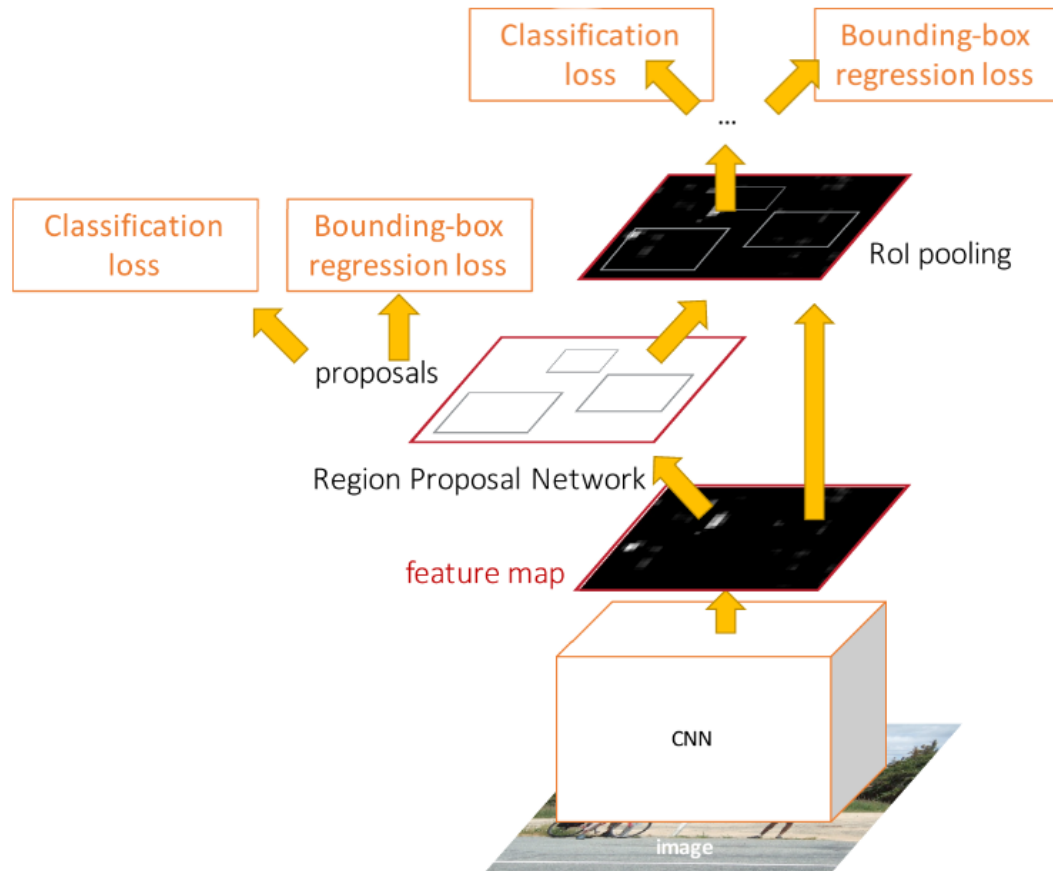
Object detection: learn to *propose* regions, Faster R-CNN (Ren et al., 2015)



- Jointly train with 4 losses:
 - RPN classification**: anchor box is object / not an object
 - RPN regression**: predict transform from anchor box to proposal box
 - Object classification**: classify proposals as background / object class
 - Object regression**: predict transform from proposal to object box

Test time speed: 0.2 seconds vs 2.3 seconds (Fast R-CNN) vs 49 seconds (R-CNN)!

Object detection: learn to *propose* regions, Faster R-CNN (Ren et al., 2015)



- Faster R-CNN is a two-stage object detector:
 1. Extract features using a backbone network and propose regions (mostly background)
 2. For each region: crop features, predict object class and bounding box offset
- Single-stage object detectors: YOLO, SSD, RetinaNet

Object detection: single-stage object detectors (YOLO, SSD, RetinaNet)

- Predict object class and location in ONE single step
- Similar to RPN of Faster R-CNN
- Predict the position of the box AND the class of the object in a box

Object detection: YOLO (Redmond et al., 2015)

Object detection: YOLO (Redmond et al., 2015)

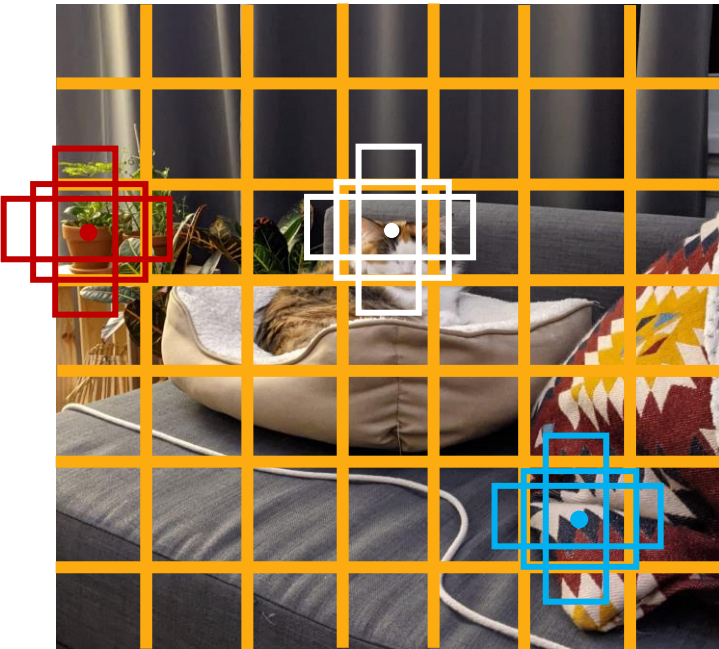
1. Divide the image into cells with an $S \times S$ grid



$S=7$

Object detection: YOLO (Redmond et al., 2015)

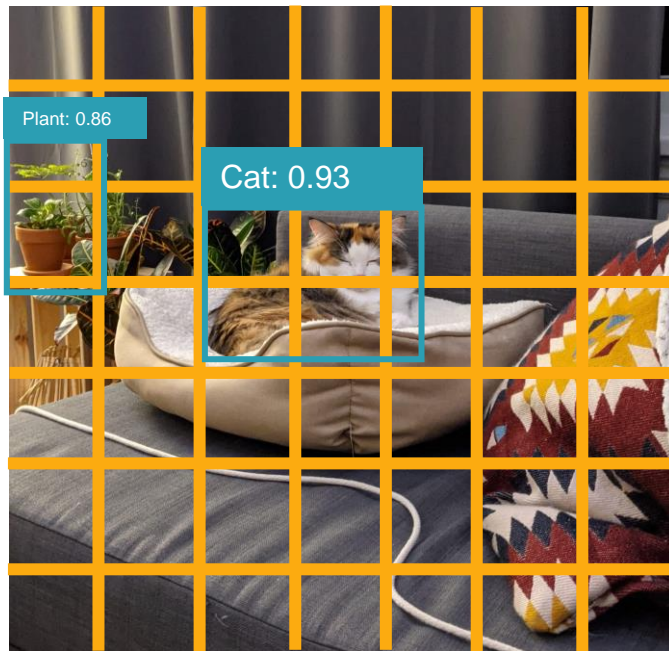
2. Each cell predicts B bounding boxes



Same idea for YOLO v2, v3, v4, ..., v8!

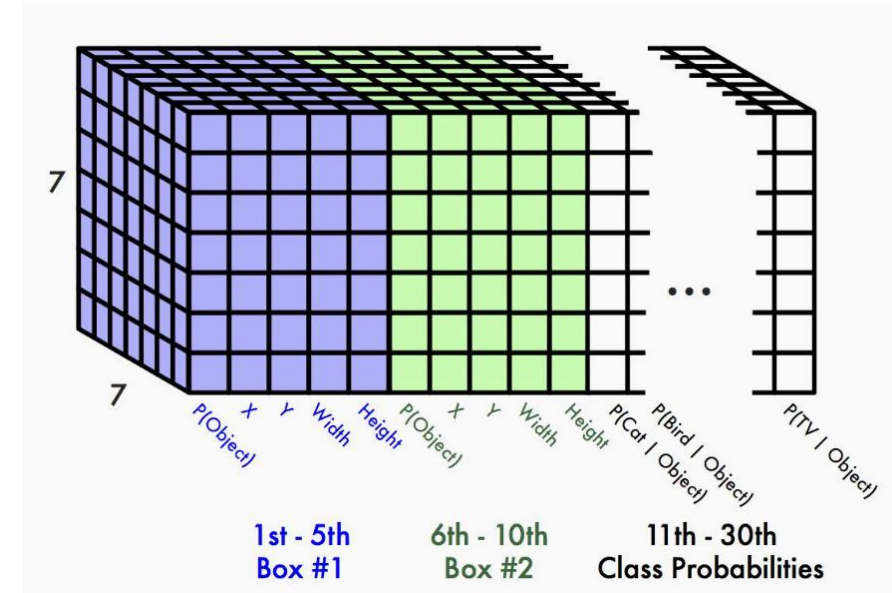
Object detection: YOLO (Redmond et al., 2015)

3. Return bounding boxes above confidence threshold

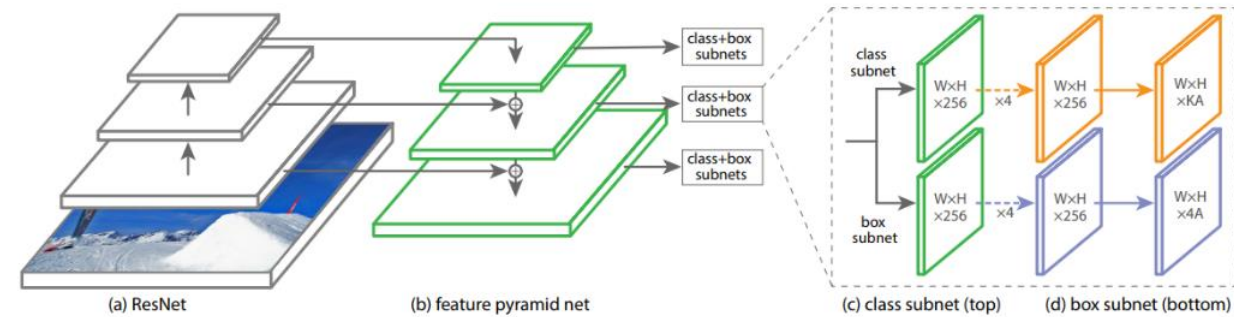
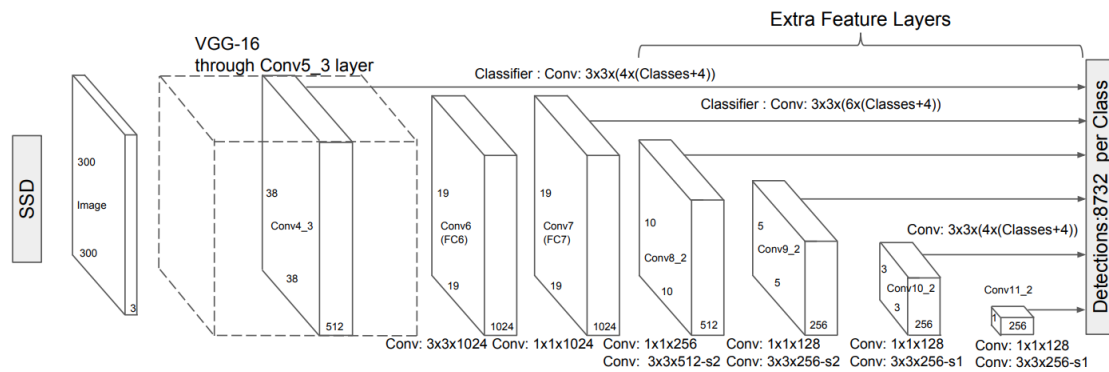


All other bounding boxes have a confidence probability less than the threshold (e.g 0.9) so they are suppressed.

- Each cell predicts:
 - For each anchor box:
 - 4 (box offset) coordinates (dx, dy, dh, dw)
 - 1 confidence value
 - Class probabilities (80 for COCO dataset, 20 for PASCAL-VOC)
- Output: $7 \times 7 \times (5 \cdot B + C)$
- Similar to RPN!



Object detection: single stage object detectors



Object detection (and segmentation): but...

Object detection (and segmentation): but...

Object detection/segmentation is a wide field of research, and many architectures exist.

Object detection (and segmentation): but...

Object detection/segmentation is a wide field of research, and many architectures exist.

- Alternatives to anchor-based methods exist:
 - CenterNet (Duan et al., 2019)
 - FCOS (Tian et al., 2019)
 - R-FCN (Dai et al., 2016)
 - DETR (object detection with transformers) (Carion et al., 2020)

Object detection (and segmentation): but...

Object detection/segmentation is a wide field of research, and many architectures exist.

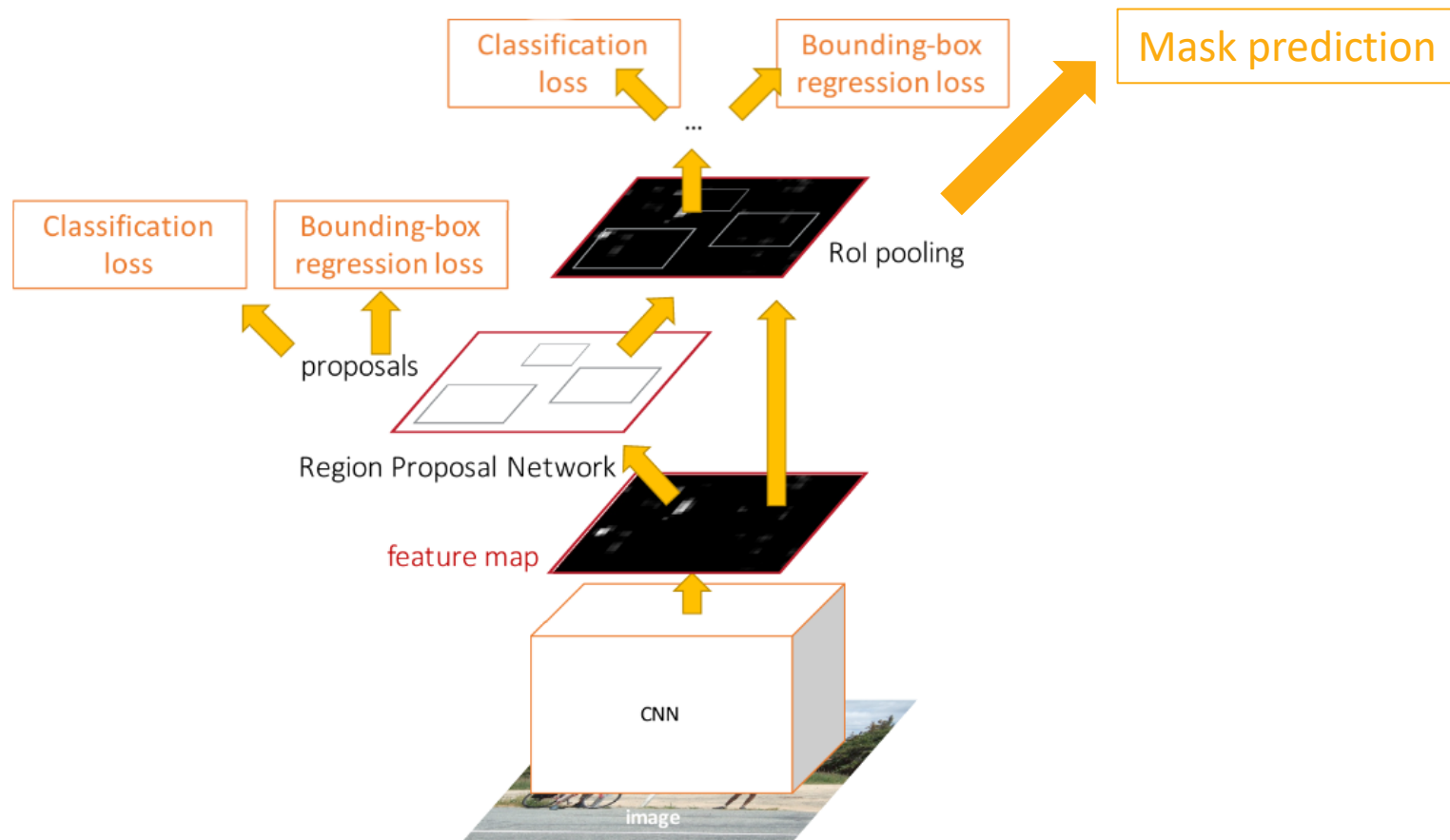
- Alternatives to anchor-based methods exist:
 - CenterNet (Duan et al., 2019)
 - FCOS (Tian et al., 2019)
 - R-FCN (Dai et al., 2016)
 - DETR (object detection with transformers) (Carion et al., 2020)
- Lot of variables:
 - Backbone network: VGG, ResNet, InceptionV2/V3, MobileNet, EfficientNet
 - Architecture style: two-stage, single-stage, hybrid...

Object detection (and segmentation): but...

Object detection/segmentation is a wide field of research, and many architectures exist.

- Alternatives to anchor-based methods exist:
 - CenterNet (Duan et al., 2019)
 - FCOS (Tian et al., 2019)
 - R-FCN (Dai et al., 2016)
 - DETR (object detection with transformers) (Carion et al., 2020)
- Lot of variables:
 - Backbone network: VGG, ResNet, InceptionV2/V3, MobileNet, EfficientNet
 - Architecture style: two-stage, single-stage, hybrid...
- **Takeways:**
 - Two-stage detectors are slower but more accurate
 - Single-stage detectors are faster but not as accurate
 - Bigger / Deeper backbones work better

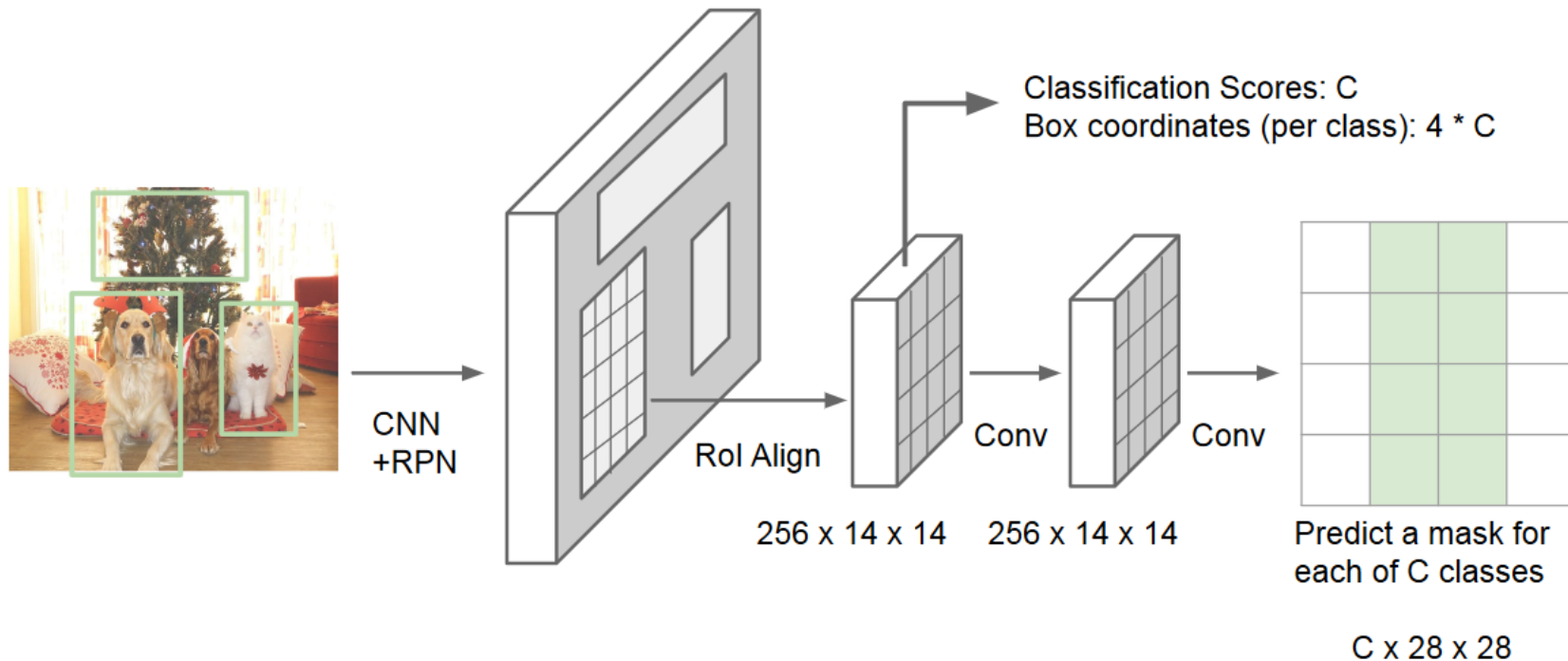
Instance segmentation: Mask R-CNN



Add a small mask network that operates on each RoI and predicts 28x28 binary mask

Instance segmentation: Mask R-CNN

Add a small mask network that operates on each RoI and predicts 28x28 binary mask

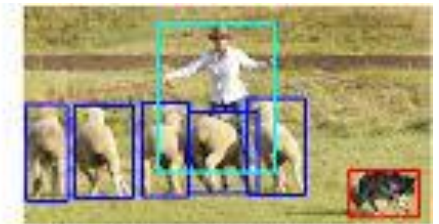


Object detection and segmentation datasets

- Pascal VOC dataset:
 - Detection, classification, segmentation
 - 10000 images with 20 categories
- COCO dataset:
 - Caption generation, object detection, key point detection and object segmentation
 - 120000 images for training / 40000 for validation with 80 categories
- KITTI autonomous driving dataset:
 - Detection, classification, semantic and instance segmentation, tracking...



(a) Image classification



(b) Object localization



(c) Semantic segmentation



(d) This work

Object detection and segmentation: frameworks

Many implementations of the aforementioned model are available on GitHub.

TensorFlow Object Detection API:
https://github.com/tensorflow/models/tree/master/research/object_detection

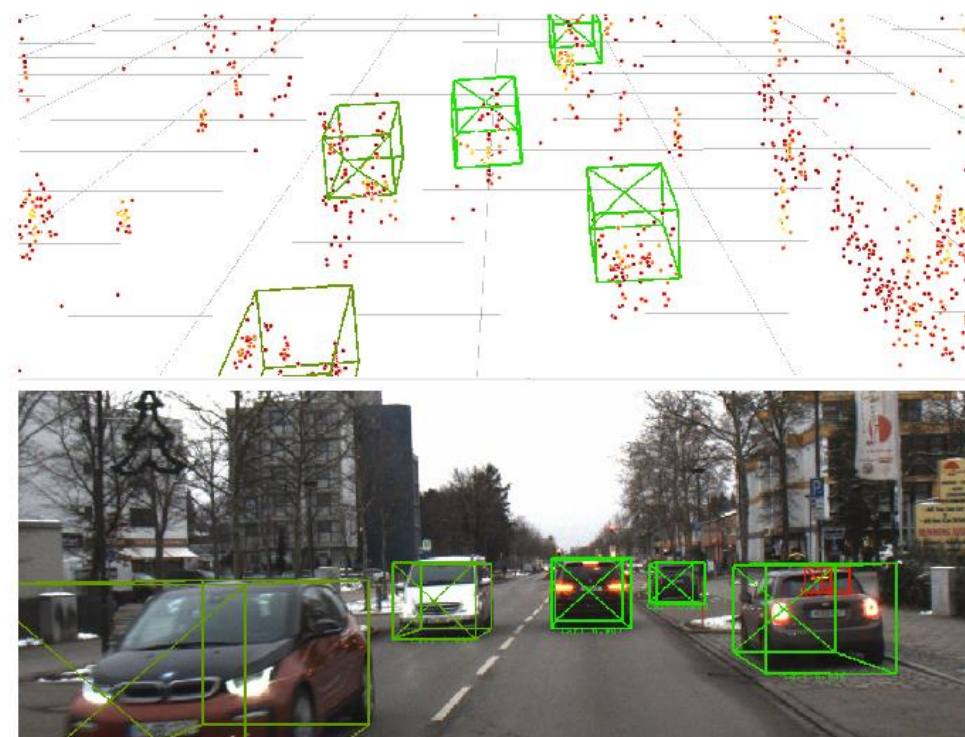
Detectron2 (PyTorch): <https://github.com/facebookresearch/detectron2>

Torchvision (PyTorch): <https://pytorch.org/vision/stable/index.html>

Use pre-trained model to finetune on your own dataset!

Other object detection/segmentation tasks

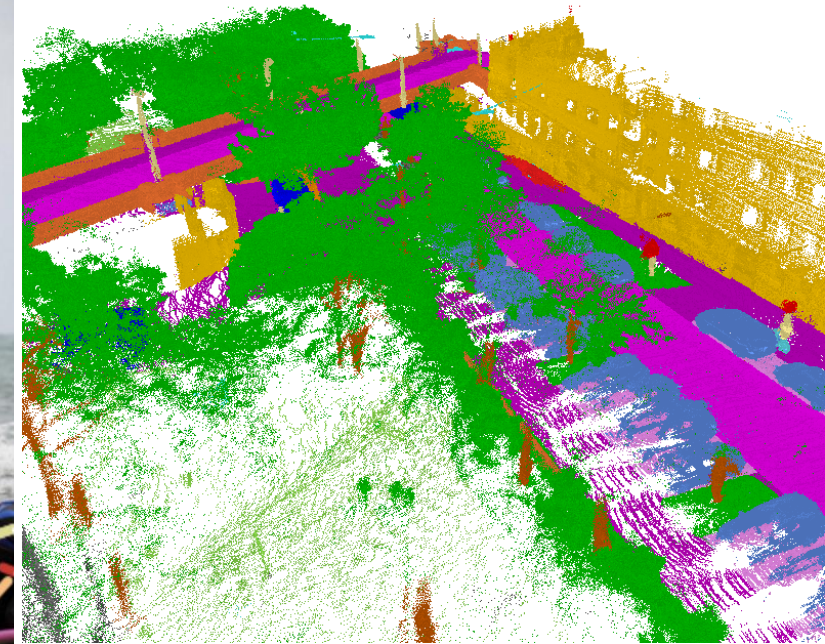
3D object detection



Key point object detection



3D semantic segmentation





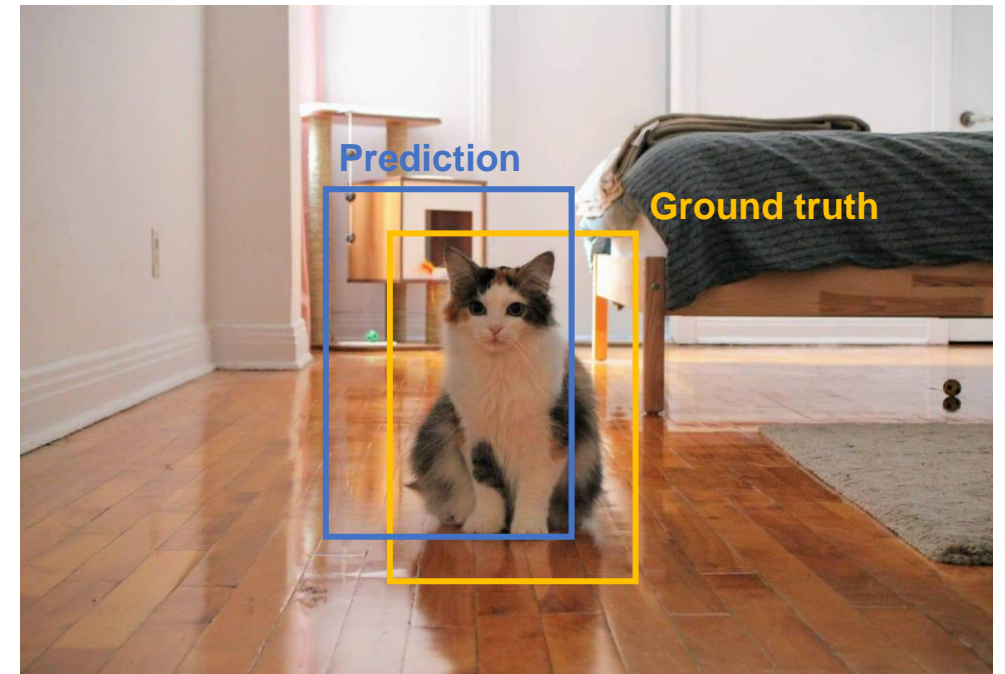
Appendix

ANITI

Université
Fédérale
Toulouse
Midi-Pyrénées

Object detection (or segmentation): comparing boxes (or masks)

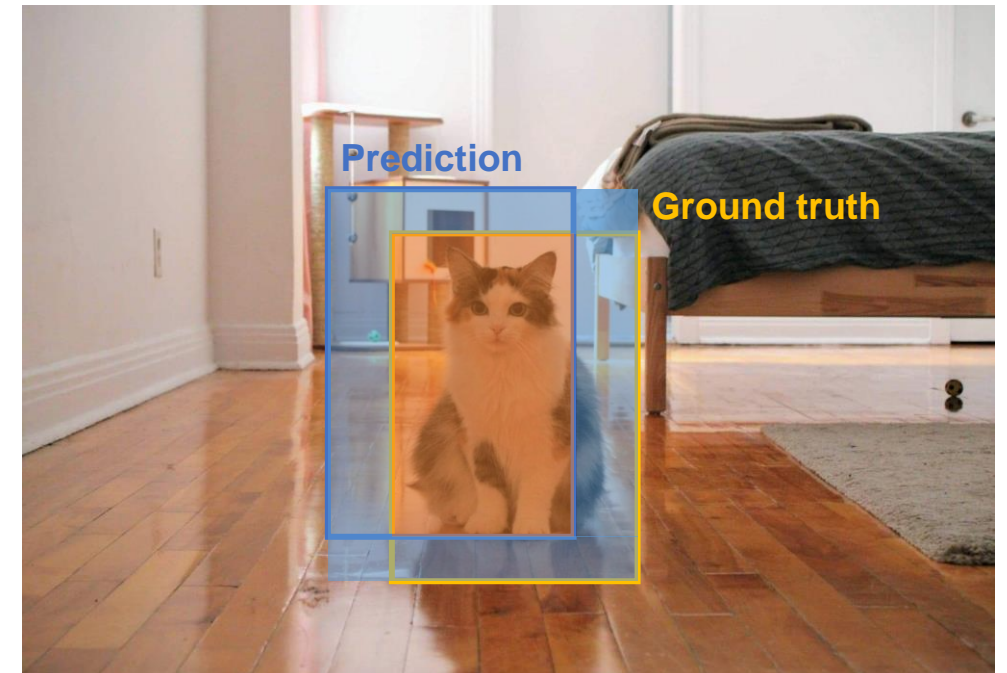
- How can we compare the prediction and the bounding boxes?



Object detection (or segmentation): comparing boxes (or masks)

- How can we compare the prediction and the bounding boxes?
- Use the **Intersection over Union (IoU)**:

$$IoU = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$



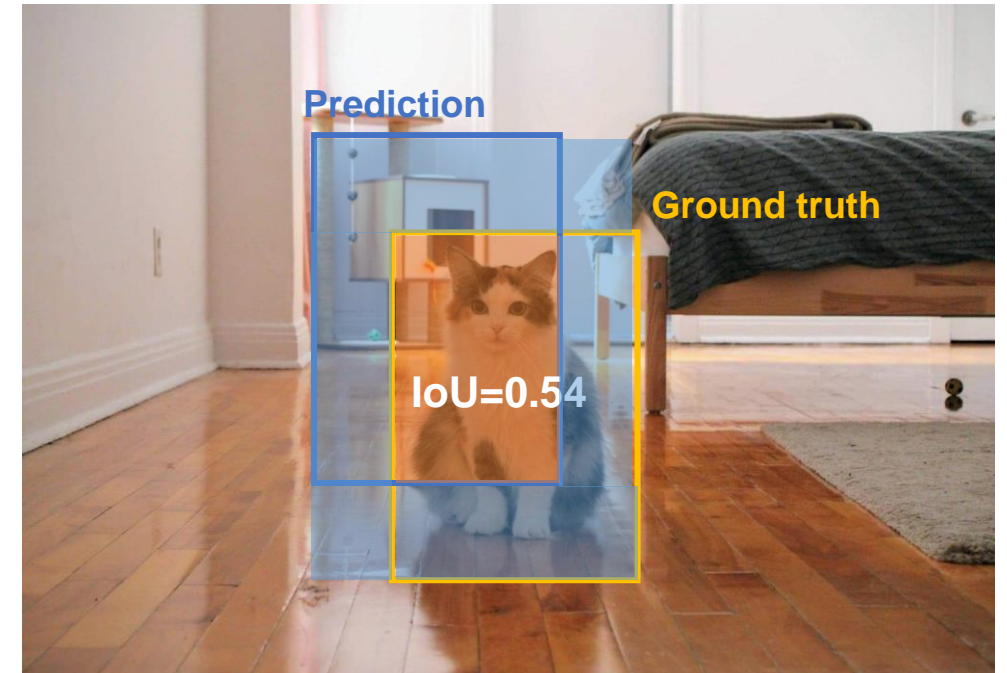
Object detection (or segmentation): comparing boxes (or masks)

- How can we compare the prediction and the bounding boxes?

- Use the **Intersection over Union (IoU)**:

$$IoU = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

- IoU > 0.5 is “decent”



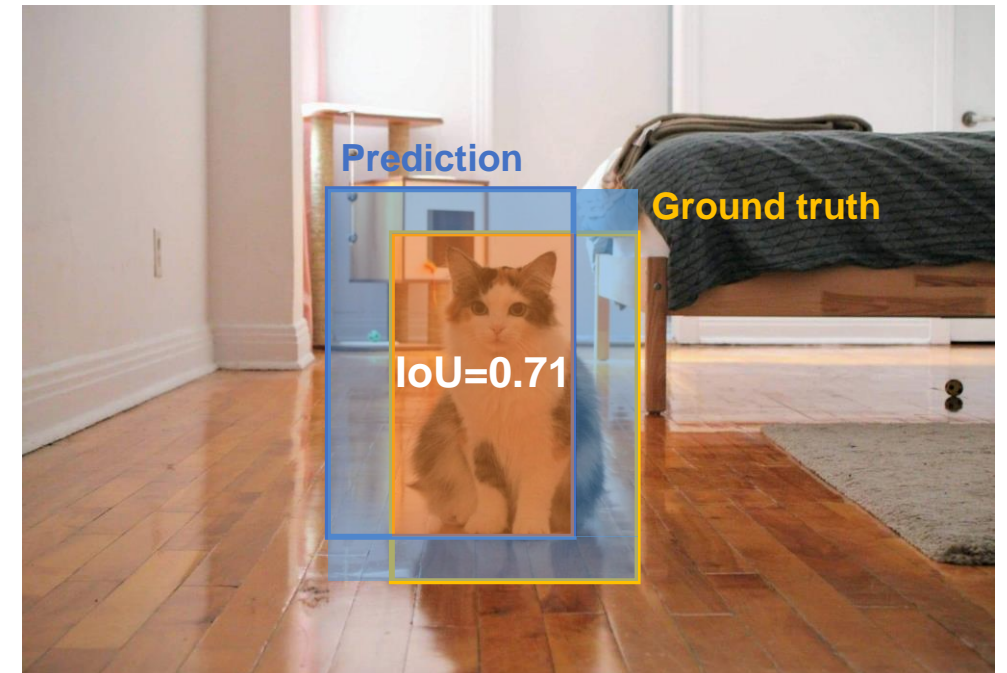
Object detection (or segmentation): comparing boxes (or masks)

- How can we compare the prediction and the bounding boxes?

- Use the **Intersection over Union (IoU)**:

$$IoU = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

- IoU > 0.5 is “decent”
- IoU > 0.7 is “pretty good”

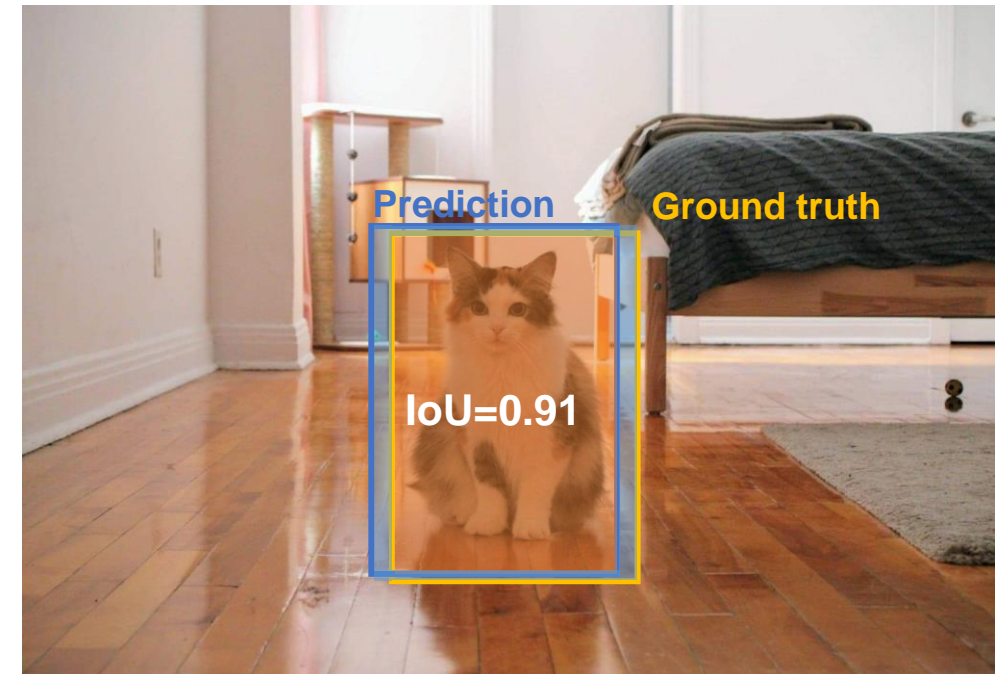


Object detection (or segmentation): comparing boxes (or masks)

- How can we compare the prediction and the bounding boxes?
- Use the **Intersection over Union (IoU)**:

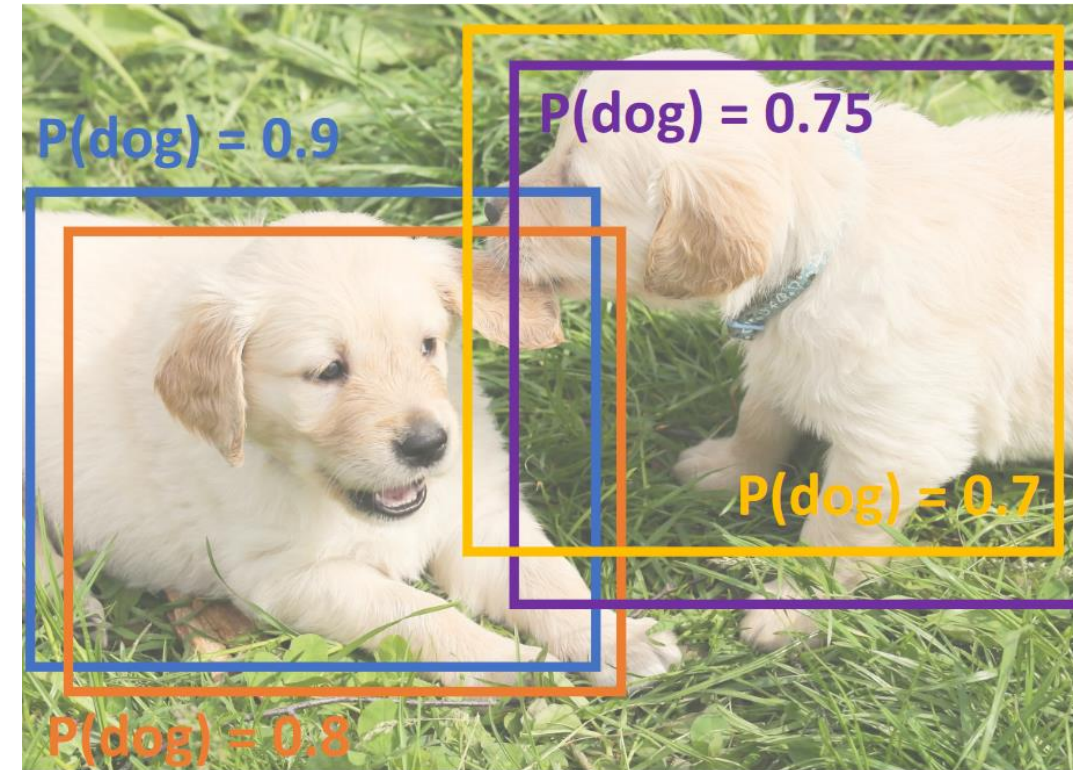
$$IoU = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

- IoU > 0.5 is “decent”
- IoU > 0.7 is “pretty good”
- IoU > 0.9 is “almost perfect”
- For segmentation masks this is done pixel-wise



Object detection: deal with overlapping boxes

- **Problem:** object detectors often output many overlapping detection (due to multiple anchors per pixel)
- **Solution:** post-process raw detections using **Non-Max Suppression (NMS)**
- Algorithm:
 1. Select highest-scoring box
 2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
 3. If any boxes remain, GOTO 1



Puppy image is CC0 Public Domain

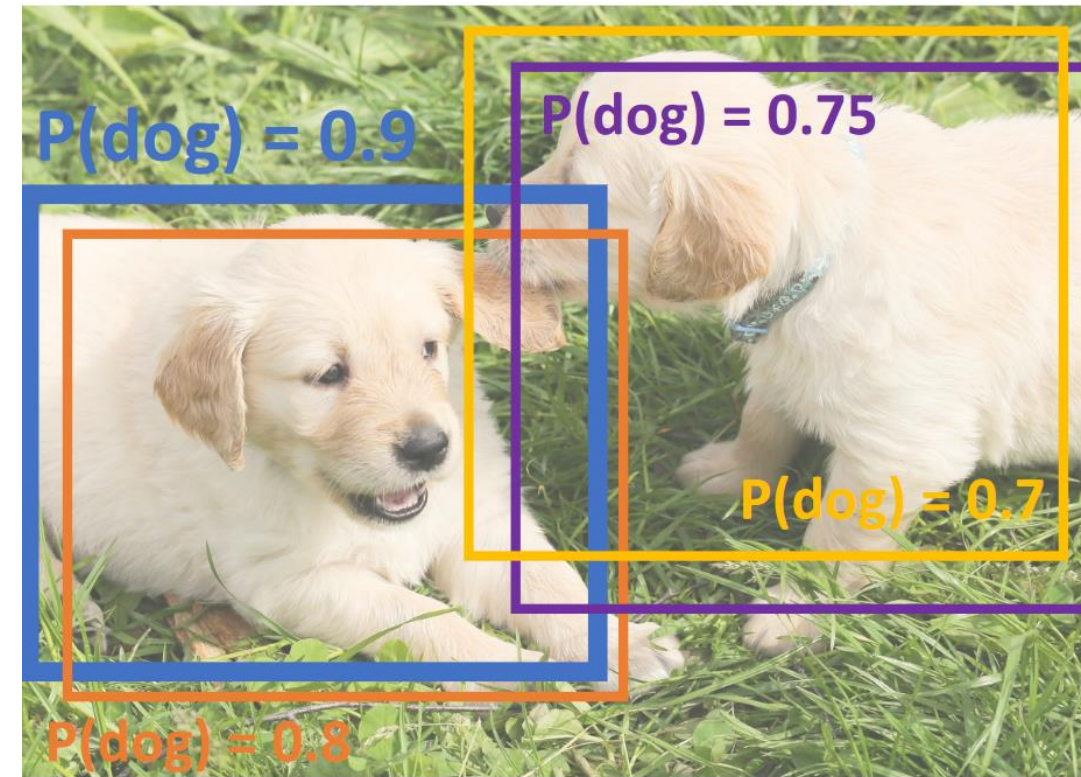
Object detection: deal with overlapping boxes

- **Problem:** object detectors often output many overlapping detection (due to multiple anchors per pixel)
- **Solution:** post-process raw detections using **Non-Max Suppression (NMS)**
- Algorithm:
 1. Select highest-scoring box
 2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
 3. If any boxes remain, GOTO 1

$$\text{IoU}(\text{blue box}, \text{orange box}) = 0.78$$

$$\text{IoU}(\text{blue box}, \text{purple box}) = 0.05$$

$$\text{IoU}(\text{blue box}, \text{yellow box}) = 0.07$$

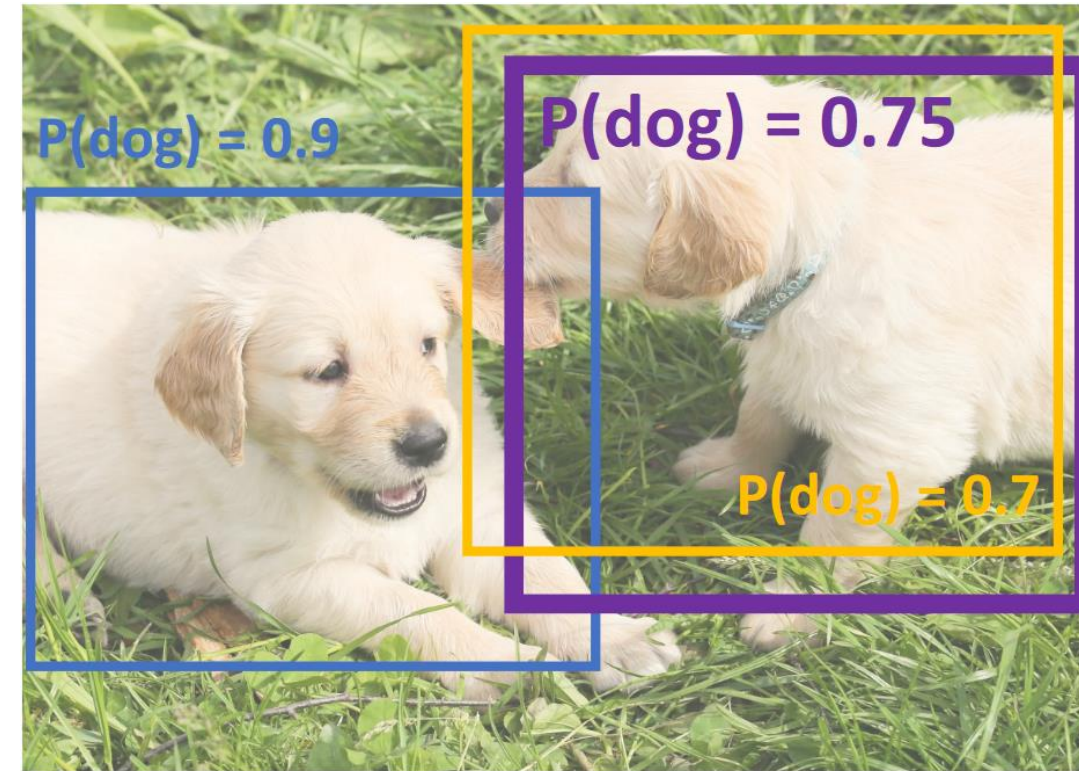


Puppy image is CC0 Public Domain

Object detection: deal with overlapping boxes

- **Problem:** object detectors often output many overlapping detection (due to multiple anchors per pixel)
- **Solution:** post-process raw detections using **Non-Max Suppression (NMS)**
- Algorithm:
 1. Select highest-scoring box
 2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
 3. If any boxes remain, GOTO 1

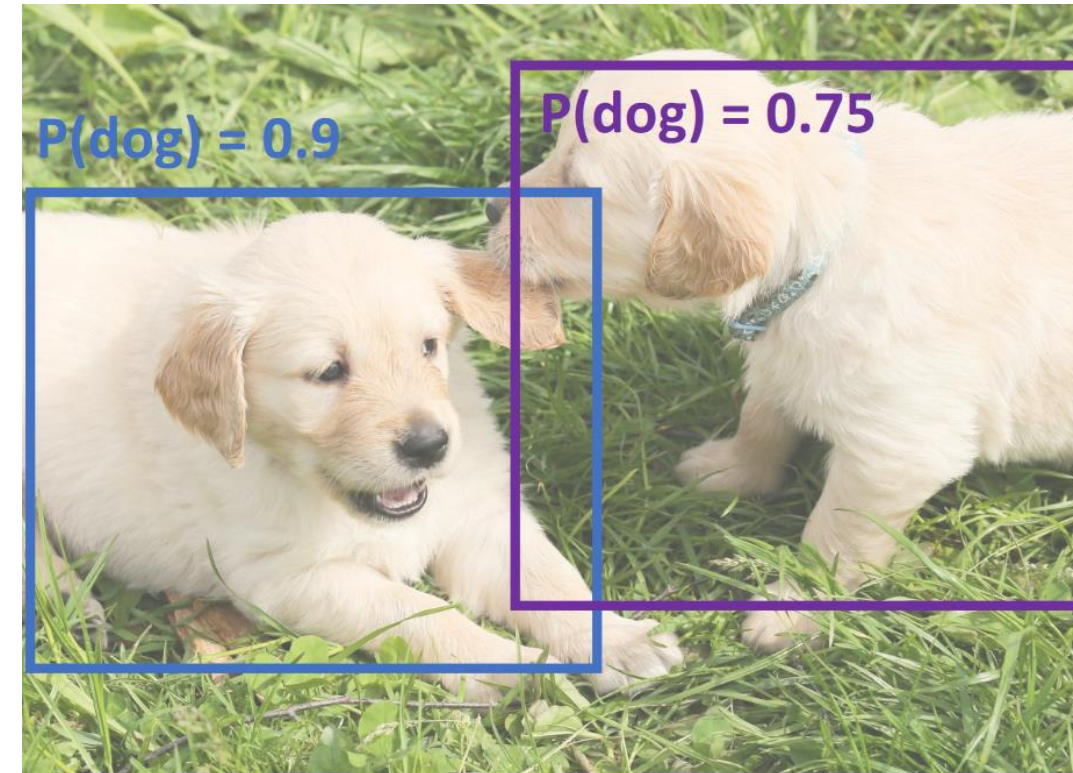
$$\text{IoU}(\blacksquare, \blacksquare) = 0.74$$



Puppy image is CC0 Public Domain

Object detection: deal with overlapping boxes

- **Problem:** object detectors often output many overlapping detection (due to multiple anchors per pixel)
- **Solution:** post-process raw detections using **Non-Max Suppression (NMS)**
- Algorithm:
 1. Select highest-scoring box
 2. Eliminate lower-scoring boxes with IoU > threshold (e.g. 0.7)
 3. If any boxes remain, GOTO 1
- **Problem:** NMS may eliminate “good” boxes when objects are highly overlapping...

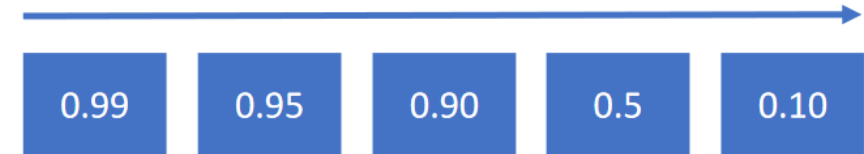


Puppy image is CC0 Public Domain

Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)

All dog detections sorted by score

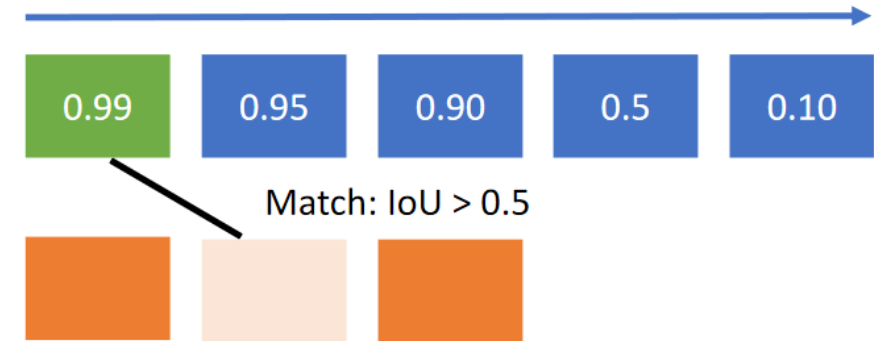


All ground-truth dog boxes

Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative

All dog detections sorted by score

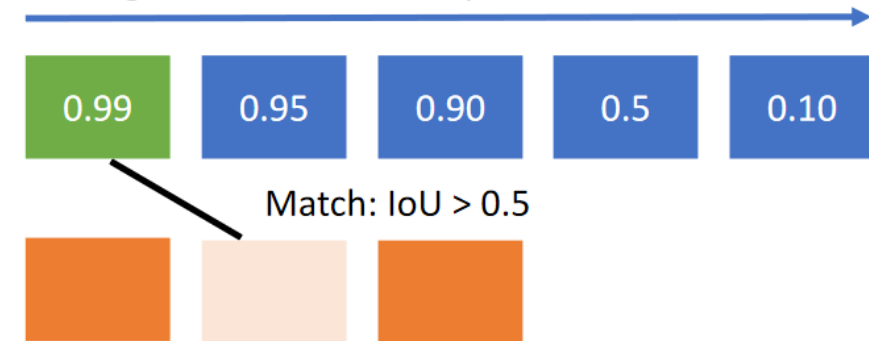


All ground-truth dog boxes

Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve

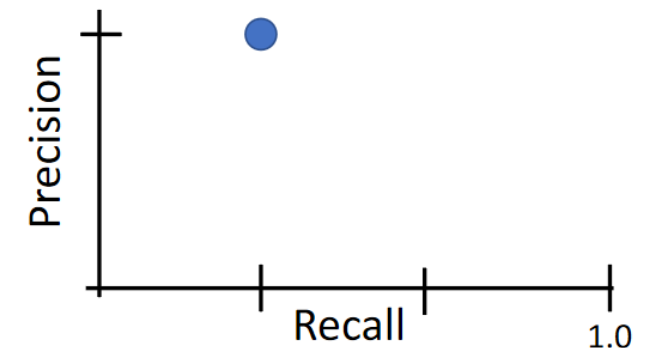
All dog detections sorted by score



All ground-truth dog boxes

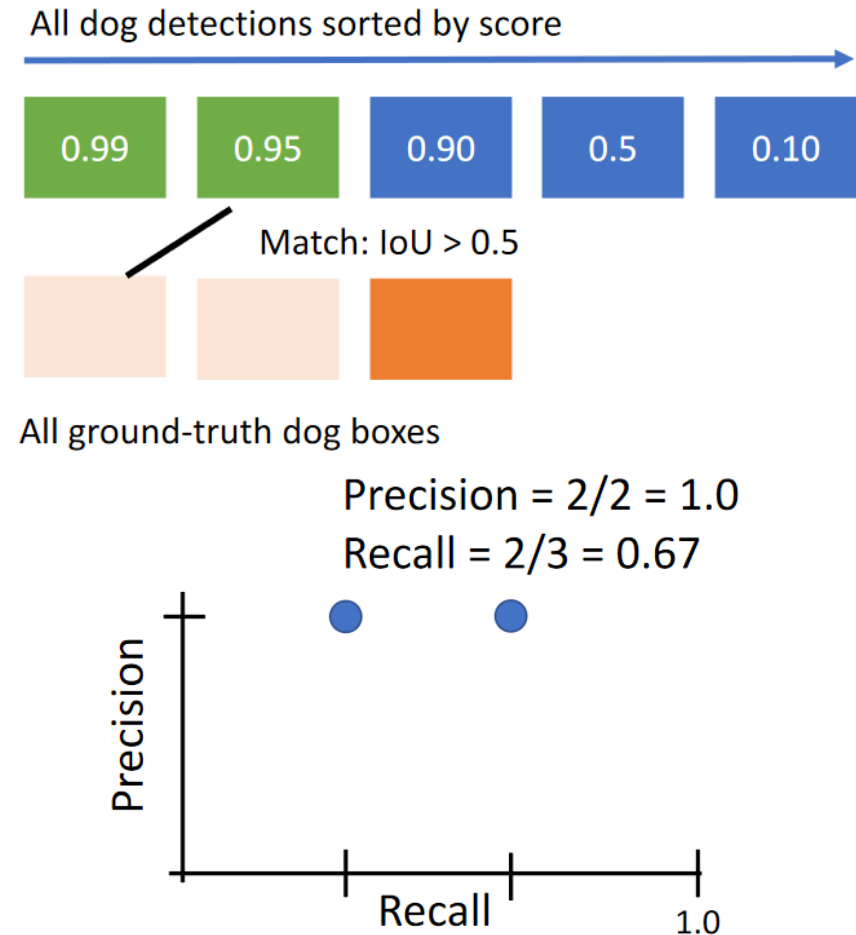
Precision = $1/1 = 1.0$

Recall = $1/3 = 0.33$



Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve



Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve

All dog detections sorted by score



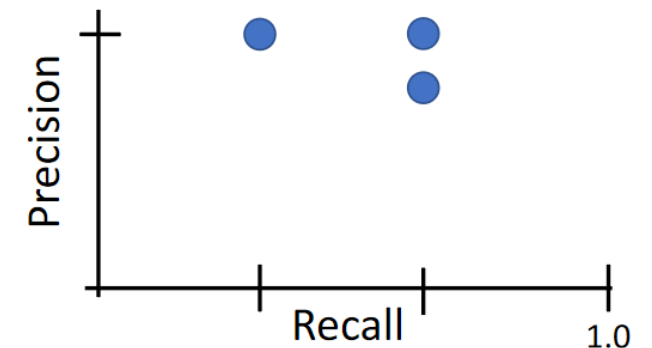
No match > 0.5 IoU with GT



All ground-truth dog boxes

Precision = $2/3 = 0.67$

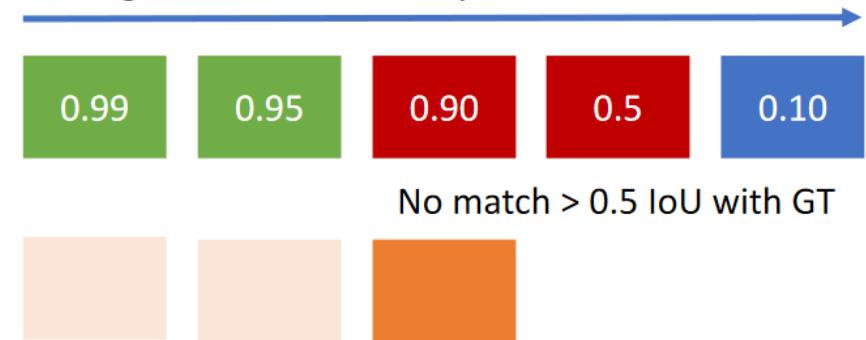
Recall = $2/3 = 0.67$



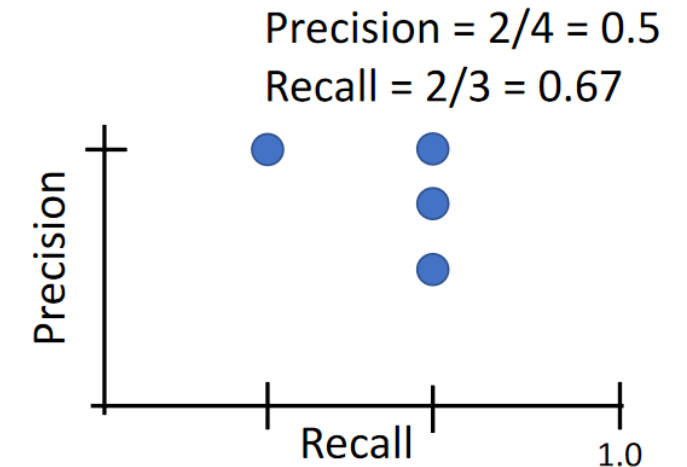
Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve

All dog detections sorted by score



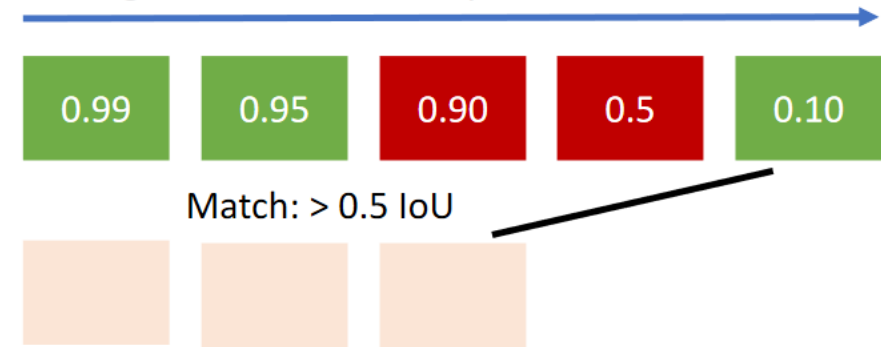
All ground-truth dog boxes



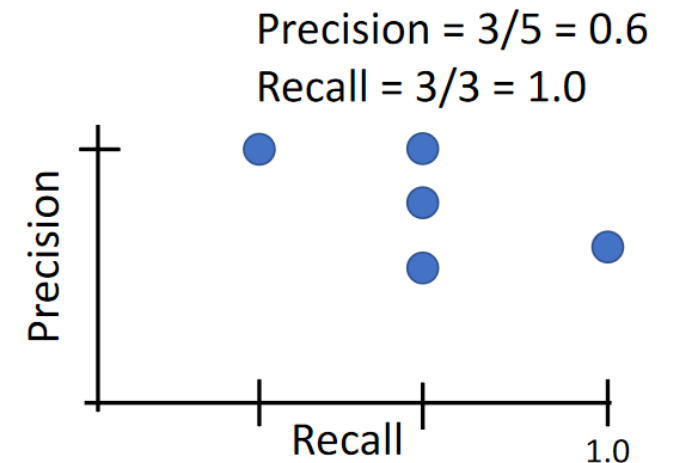
Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve

All dog detections sorted by score



All ground-truth dog boxes



Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve
 2. Average Precision (AP) = Area under PR curve
3. Mean Average Precision (mAP) = average of AP for each category

CarAP = 0.65

Cat AP = 0.80

Dog AP = 0.86

mAP@0.5 = 0.77

Evaluating object detectors – Mean Average Precision (mAP)

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = Area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with IoU > 0.5, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR Curve
 2. Average Precision (AP) = Area under PR curve
3. Mean Average Precision (mAP) = average of AP for each category
4. For “COCO mAP”: Compute mAP@thresh for each IoU threshold (0.5, 0.55, 0.6, ..., 0.95) and take average

mAP@0.50 = 0.77

mAP@0.55 = 0.71

mAP@0.60 = 0.65

...

mAP@0.95 = 0.2

COCO mAP = 0.4