

Introduction to Computer Vision: Object detection (part 2)

Natalia Neverova

Last lecture's recap

- What is the object detection task?
- How is this task formalized?
- How is this task evaluated?
- Why is the task difficult?
- What are important datasets to know about?



Today's focus

- Principles of multi-object detection
- Region proposals
- Bounding box regression
- R-CNN training and inference
- Non-Max suppression
- Modern object detectors

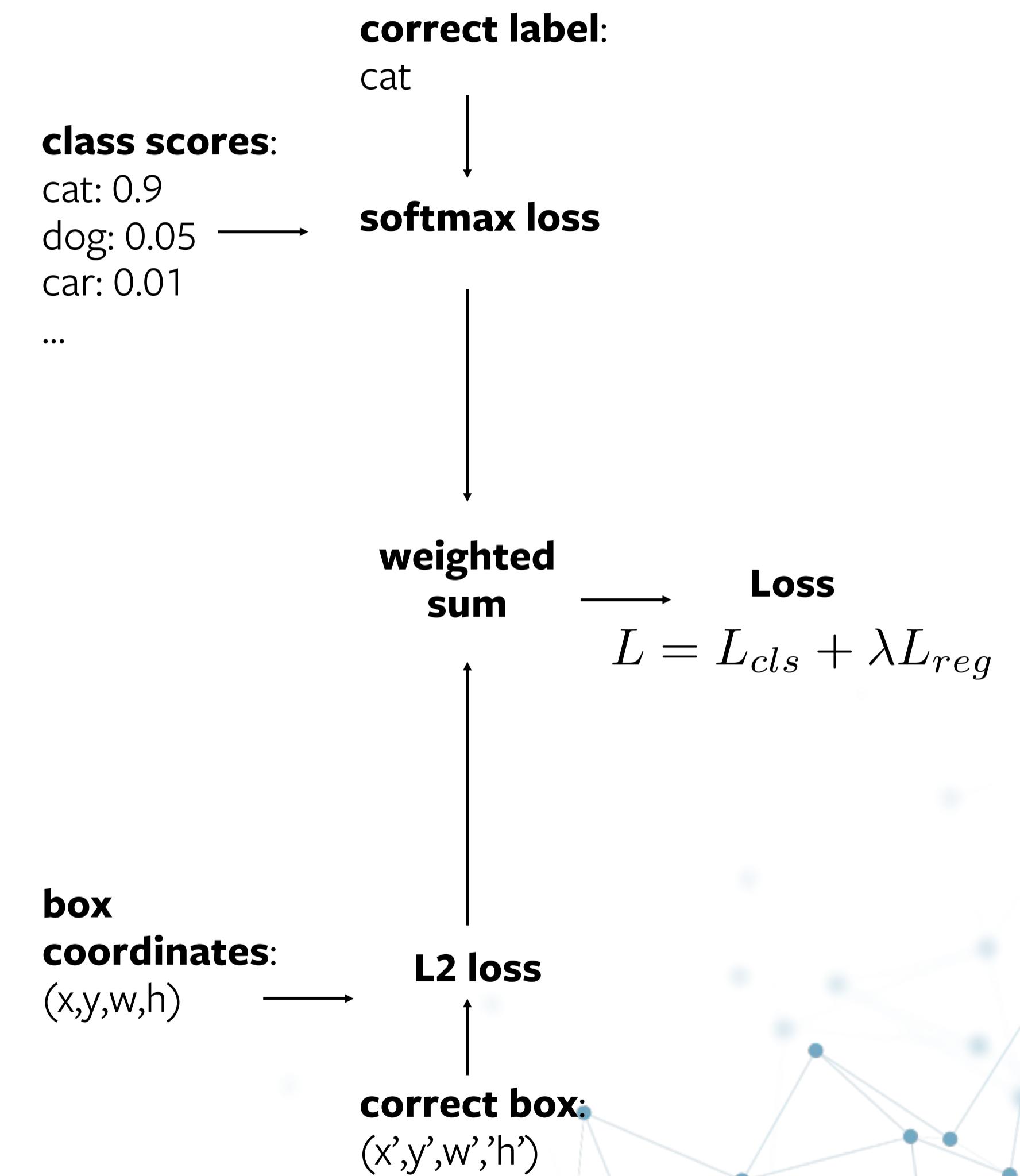
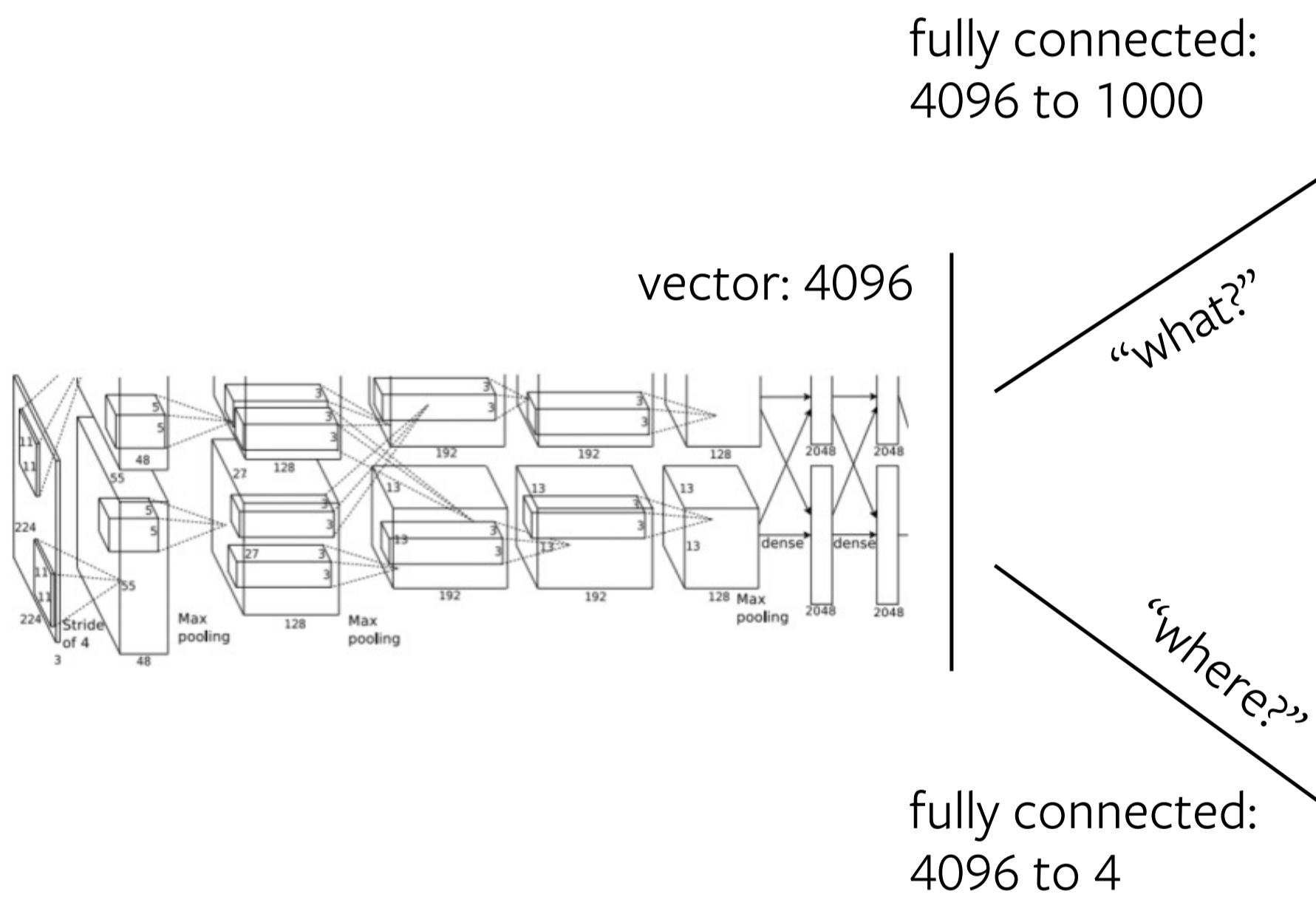
Recap: object detection task

What objects are in the image and **where** are they?

Detecting a single object



treat localization as a regression problem!

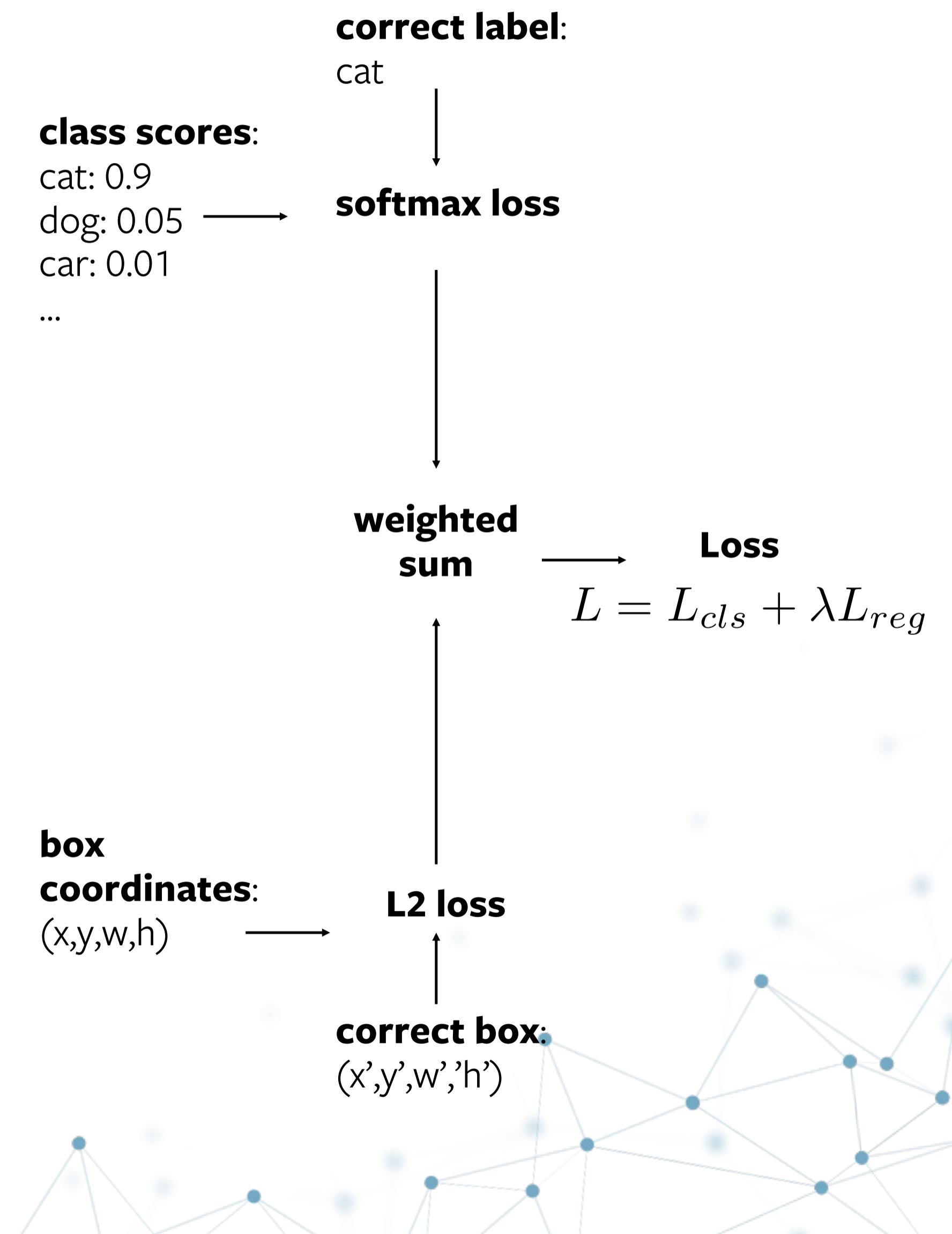
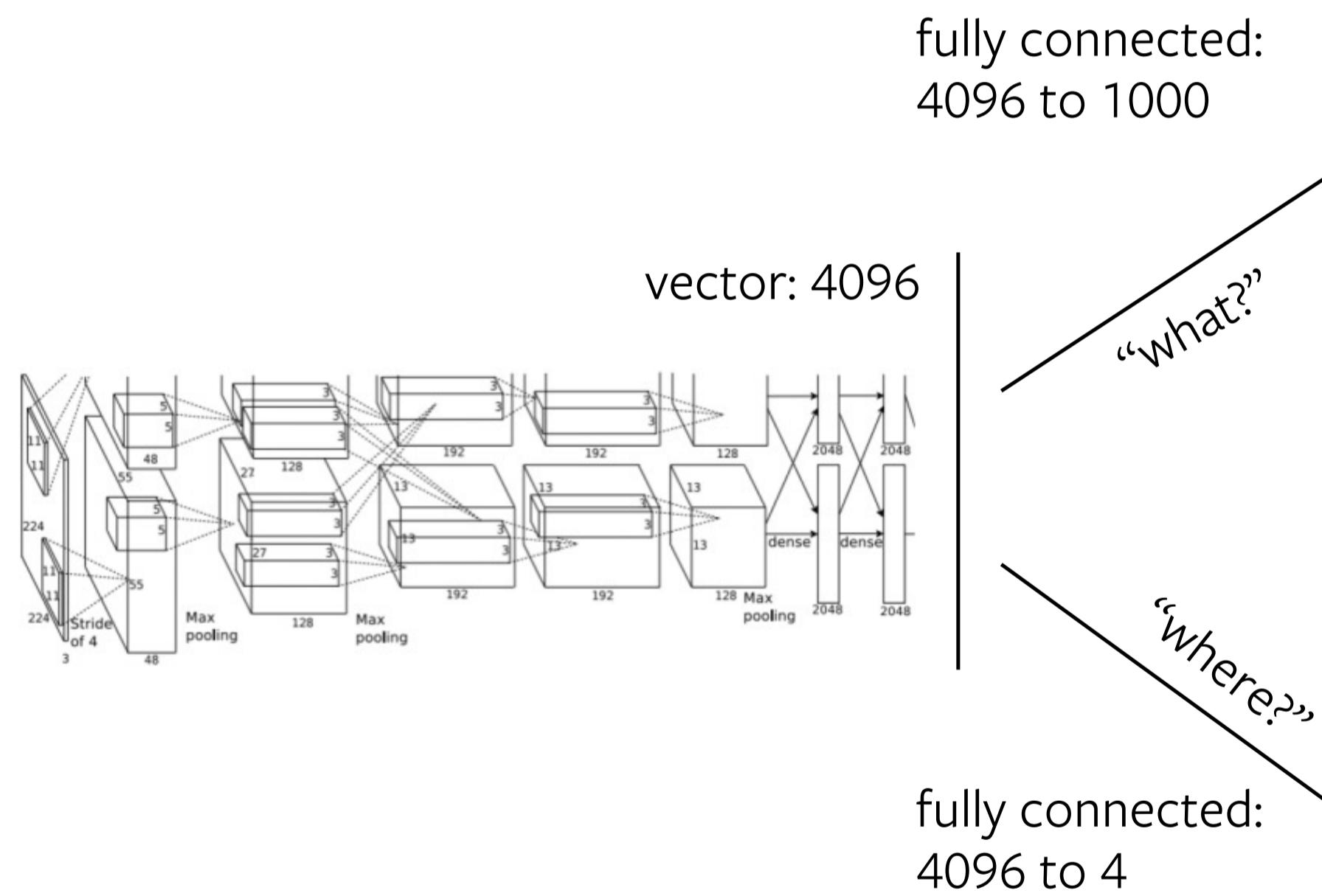


Detecting a single object

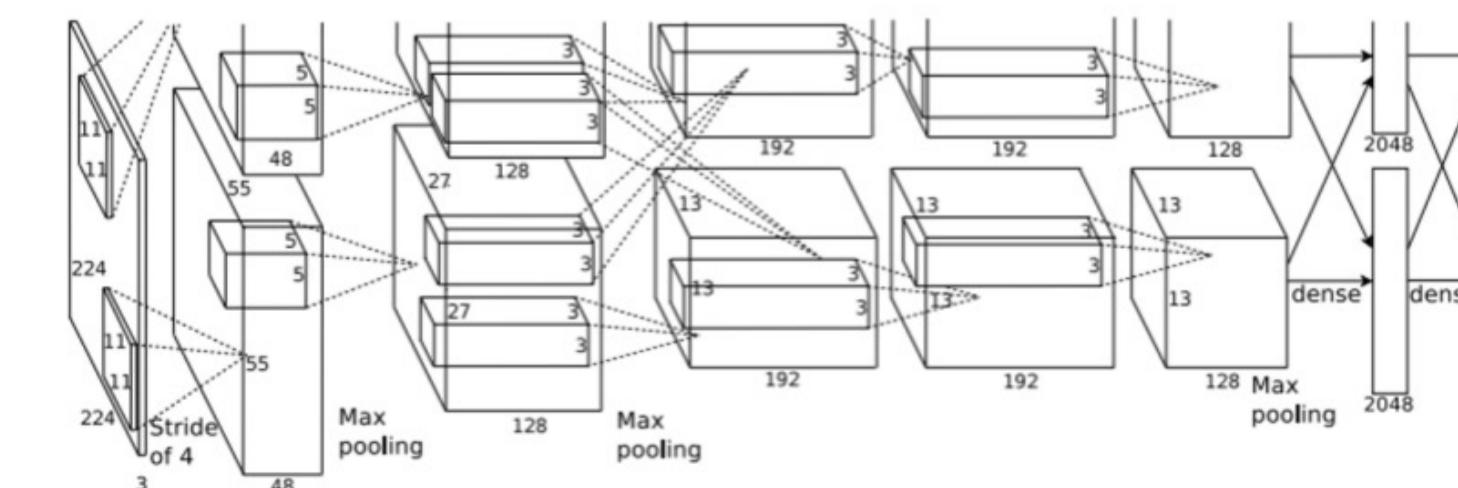
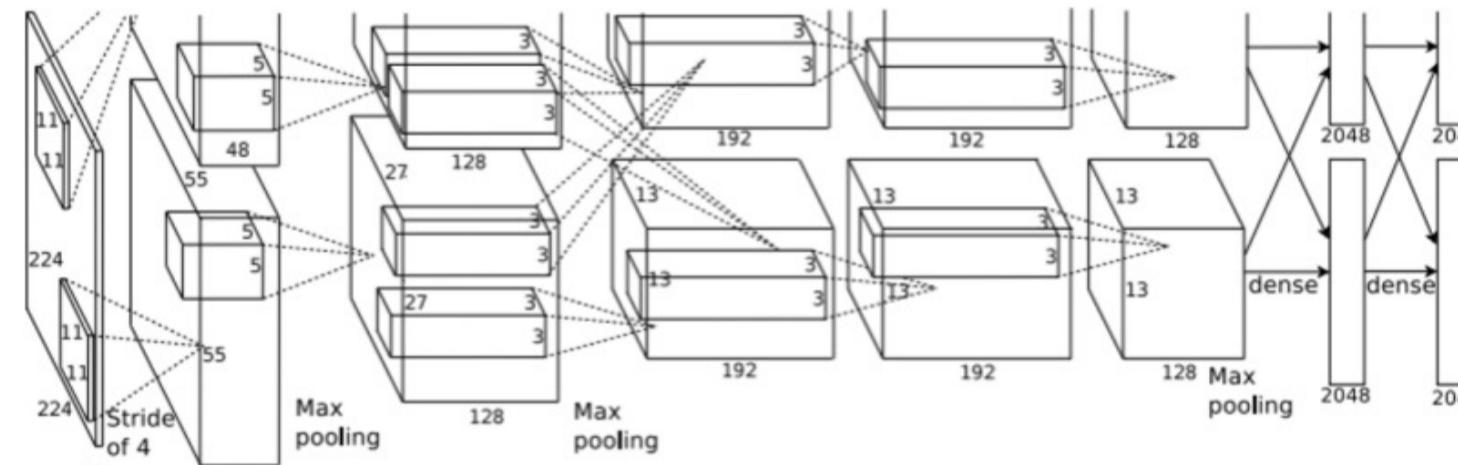
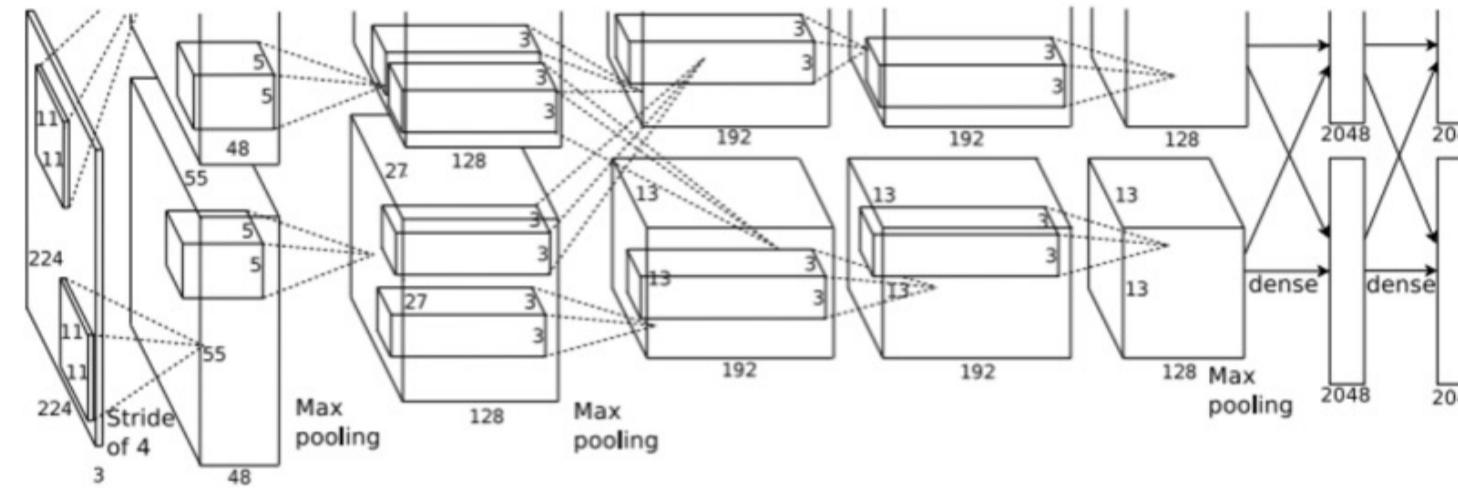


treat localization as a regression problem!

**things get much more complicated with multiple objects!
(today!)**



Detecting multiple objects



CAT: (x, y, w, h)

DOG: (x, y, w, h)

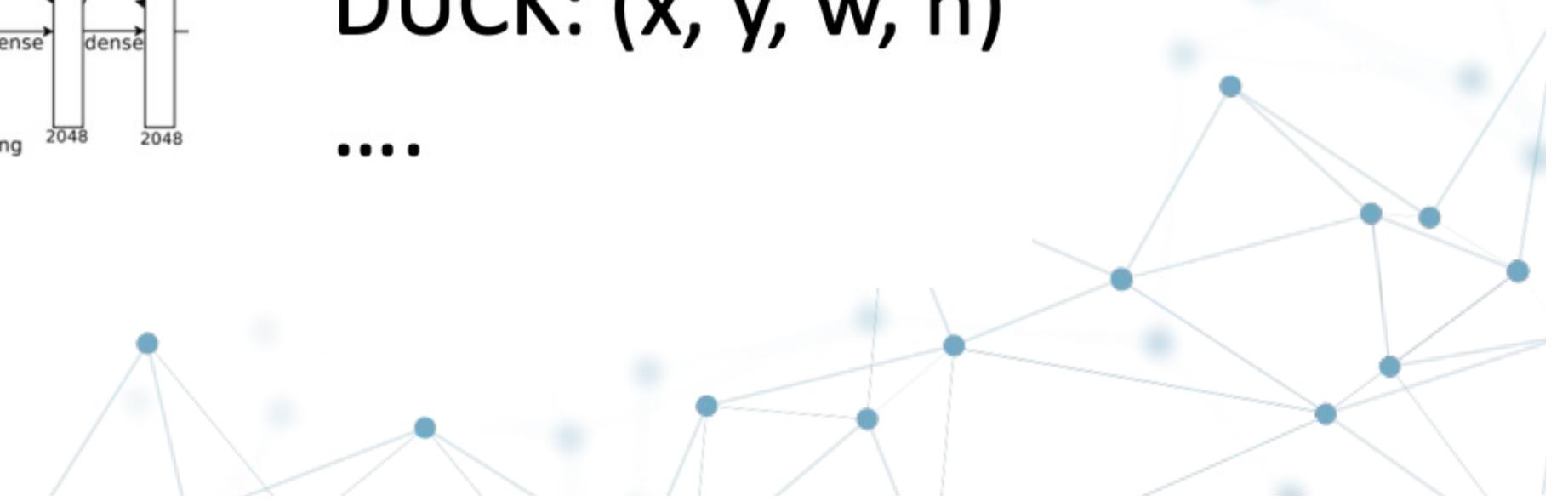
DOG: (x, y, w, h)

CAT: (x, y, w, h)

DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

...



Multiple object detection as ML problem

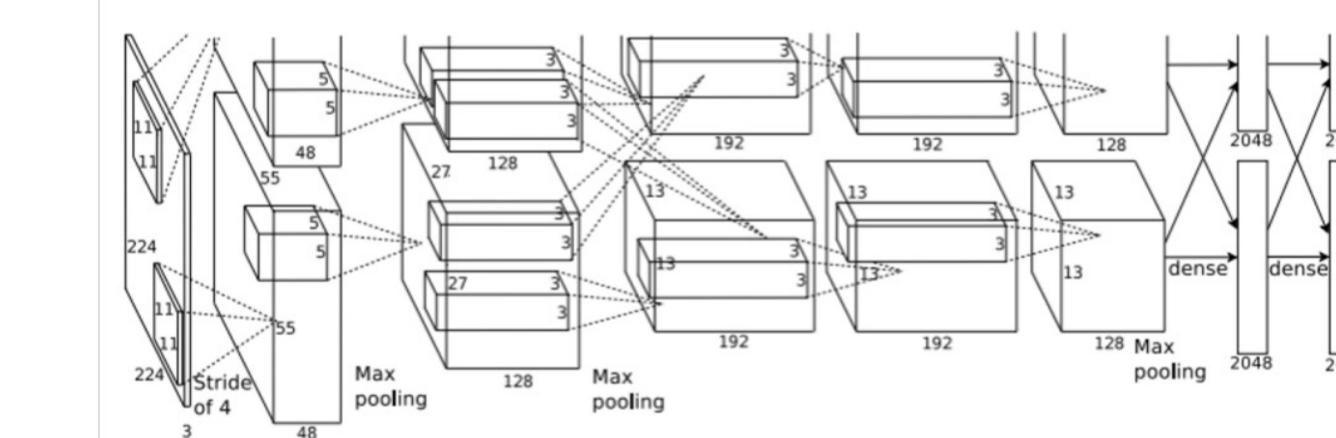
Detection := the classification of boxes

Intuition: train a classifier to learn to classify every possible box
(box's class is 0 if it's background)

Foundational idea: sliding window



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

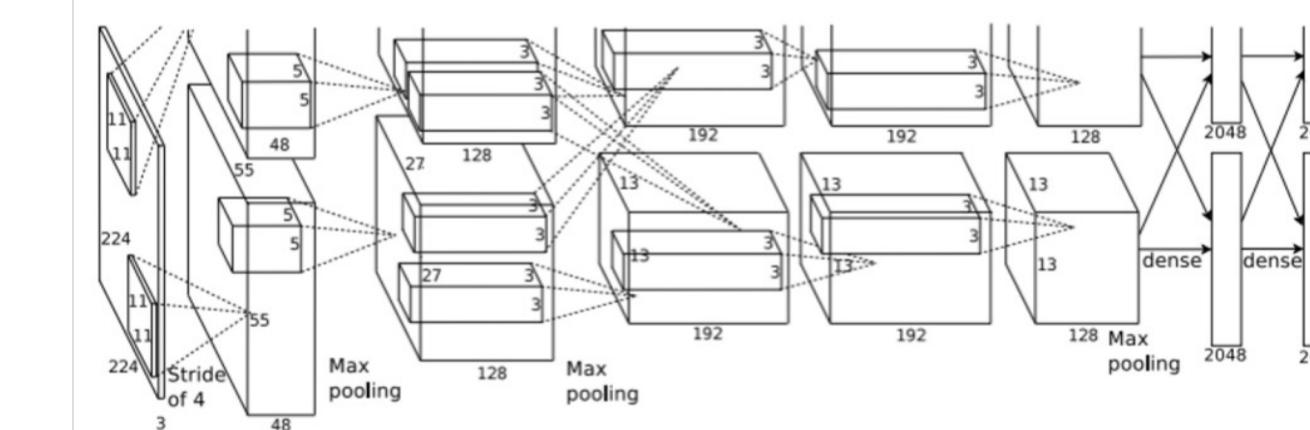


Dog? No
Cat? No
Background? Yes

Foundational idea: sliding window



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

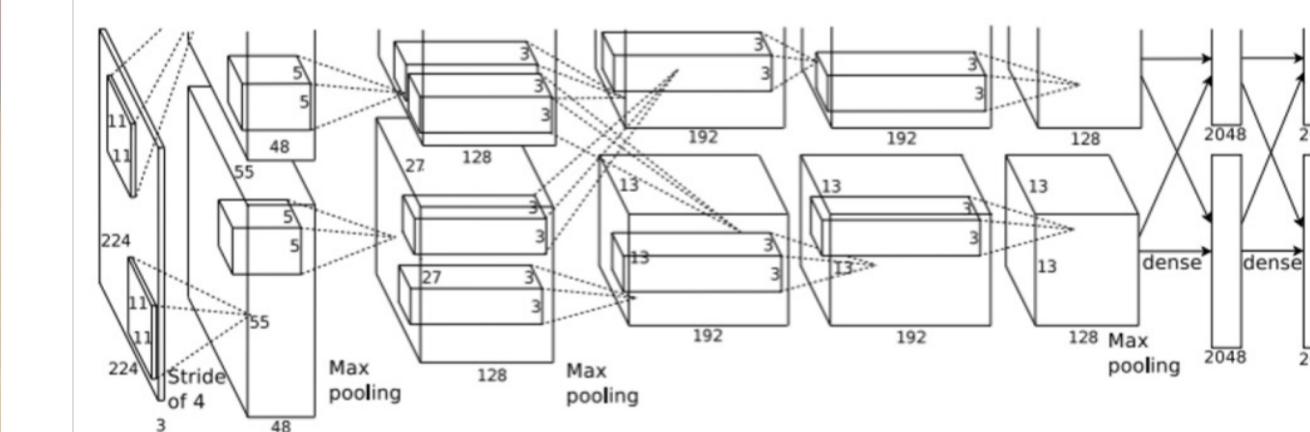


Dog? Yes
Cat? No
Background? No

Foundational idea: sliding window



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

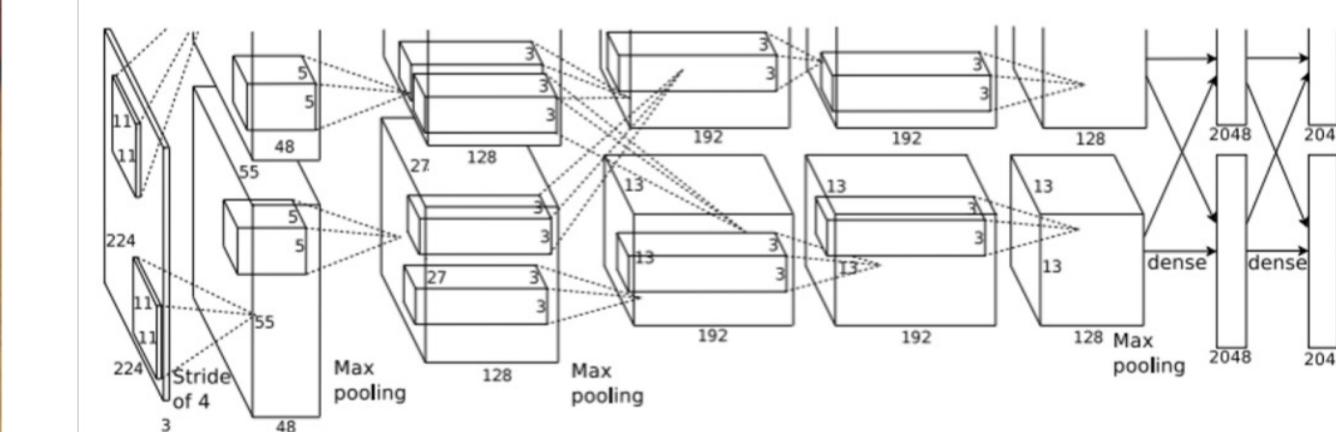


Dog? Yes
Cat? No
Background? No

Foundational idea: sliding window



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



Dog? No
Cat? Yes
Background? No

Foundational idea: sliding window



But: How many possible boxes are there in an image size $H \times W$?

Consider a box of size $h \times w$:
possible x positions: $W - w + 1$
possible y positions: $H - h + 1$
possible positions:
 $(W - w + 1) * (H - h + 1)$

Foundational idea: sliding window

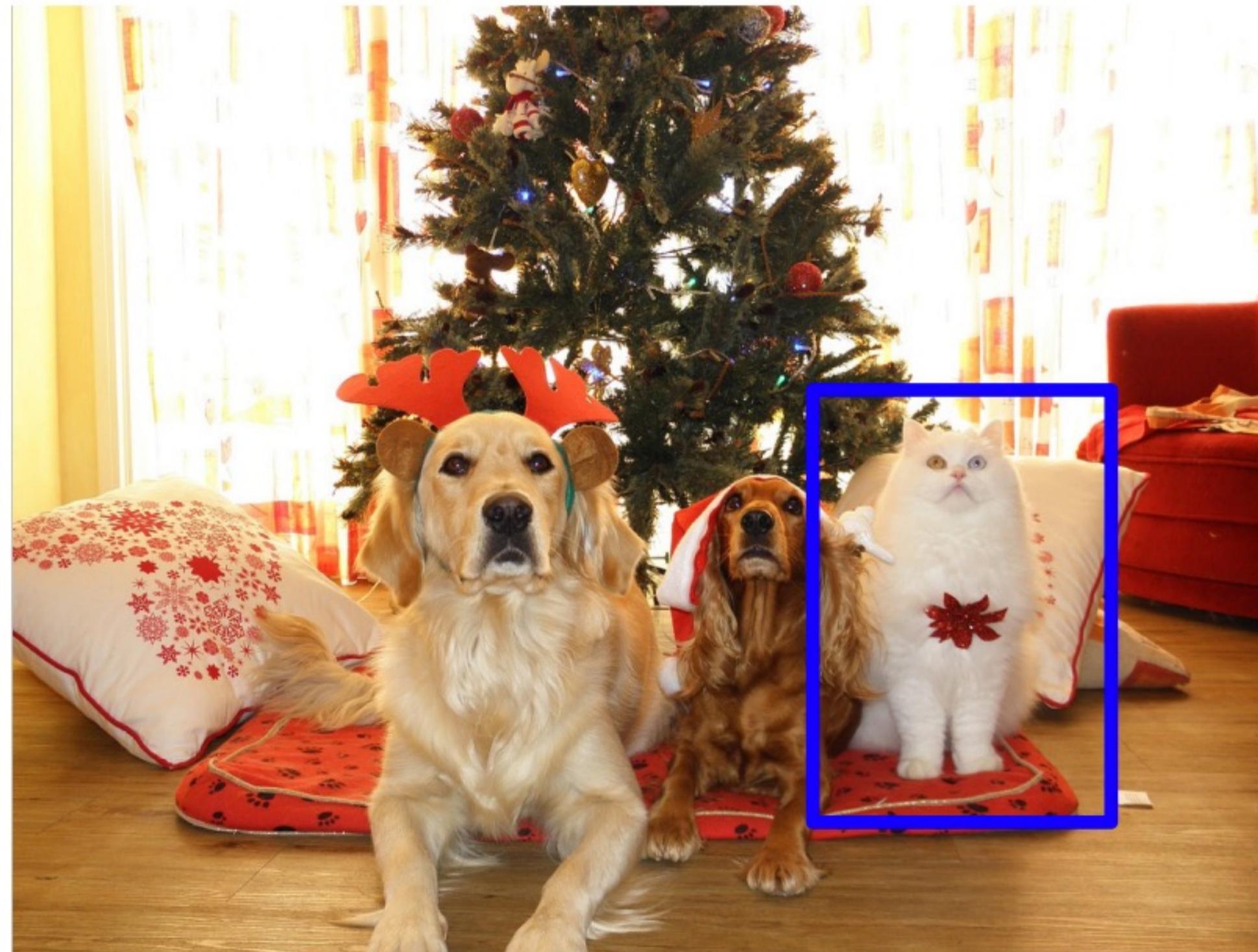


But: How many possible boxes are there in an image size $H \times W$?

Total possible boxes:

$$\sum_{h=1}^H \sum_{w=1}^W (W - w + 1)(H - h + 1) = \frac{H(H + 1)}{2} \frac{W(W + 1)}{2}$$

Foundational idea: sliding window



But: How many possible boxes are there in an image size $H \times W$?

Total possible boxes:

$$\sum_{h=1}^H \sum_{w=1}^W (W - w + 1)(H - h + 1) = \frac{H(H + 1)}{2} \frac{W(W + 1)}{2}$$

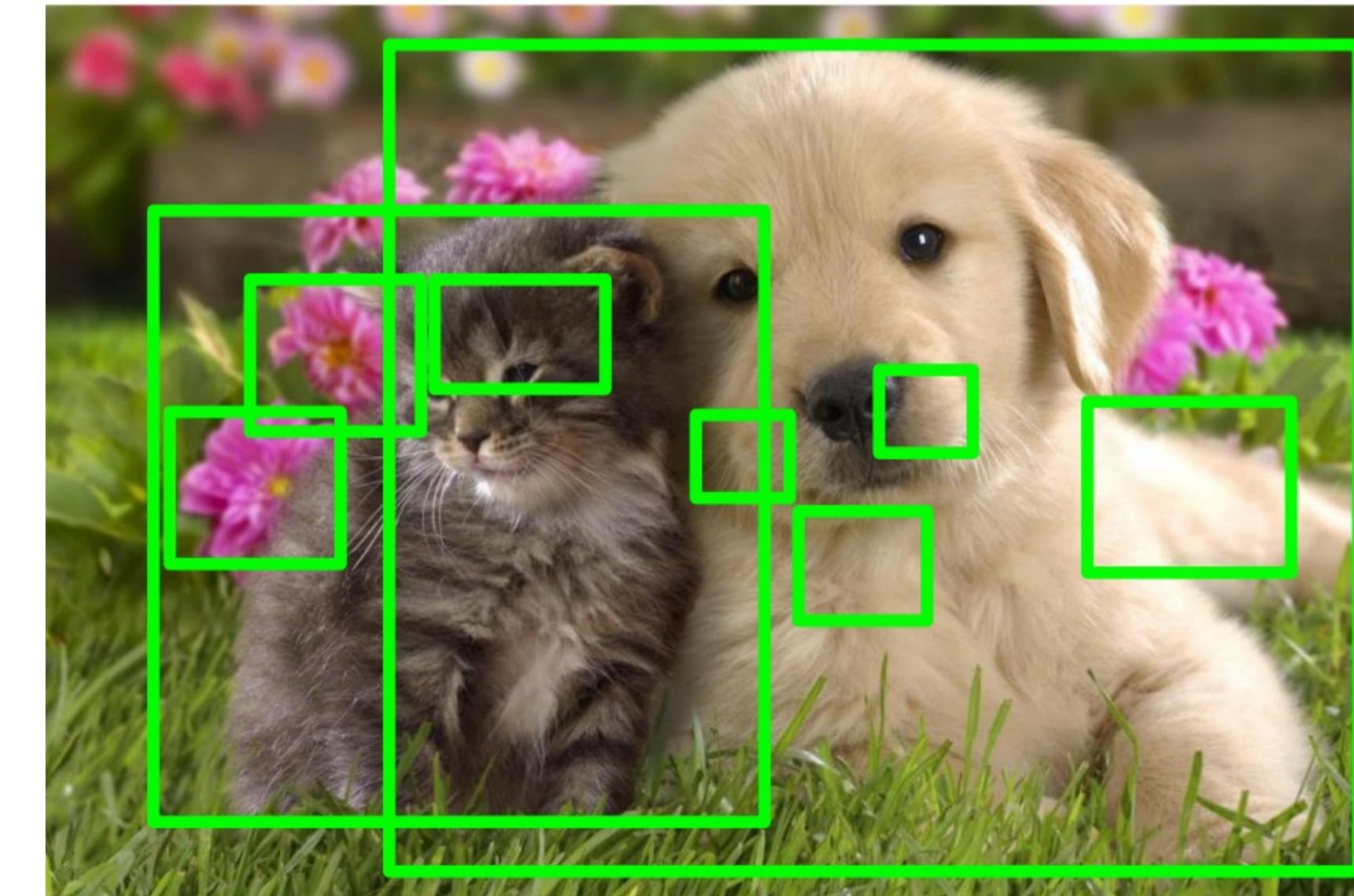
800x600 image has ~58M boxes!
No way we can evaluate them all

Region proposals

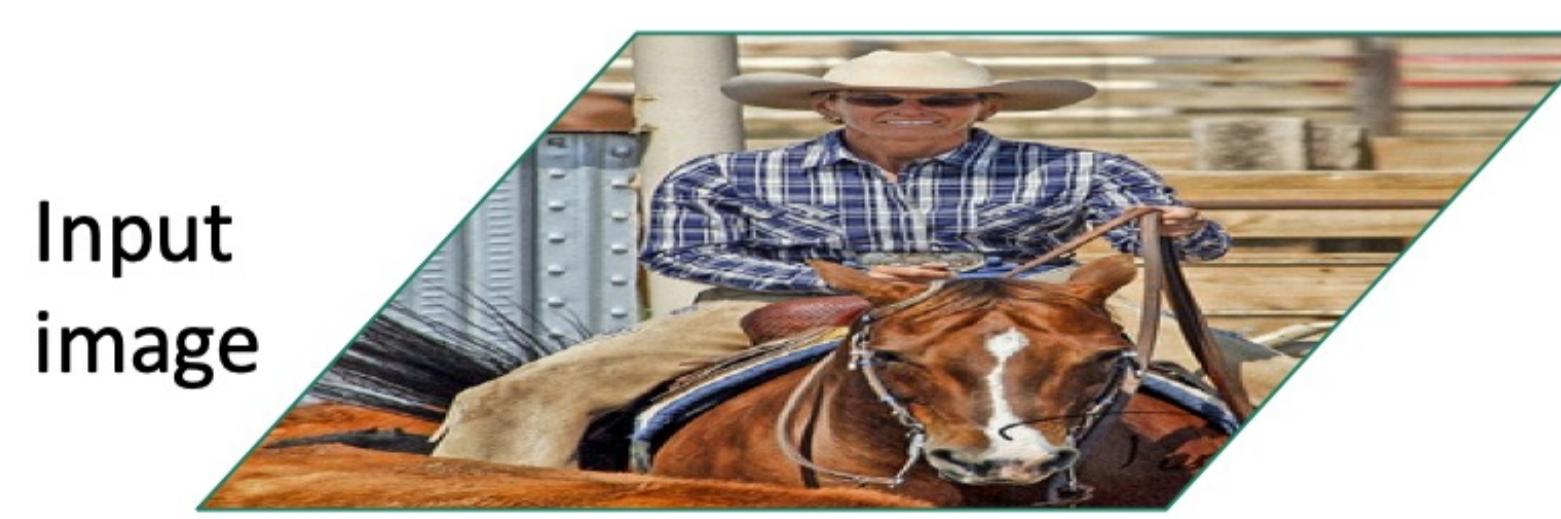
Idea: Find a small set of boxes that are likely to cover all objects.

Often based on heuristics: e.g. look for “blob-like” image regions

Relatively fast to run; e.g. Selective Search [Uijings et al., 2013] gives 2000 regions proposals in a few seconds on CPU



R-CNN: Region-based CNN



Input
image

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



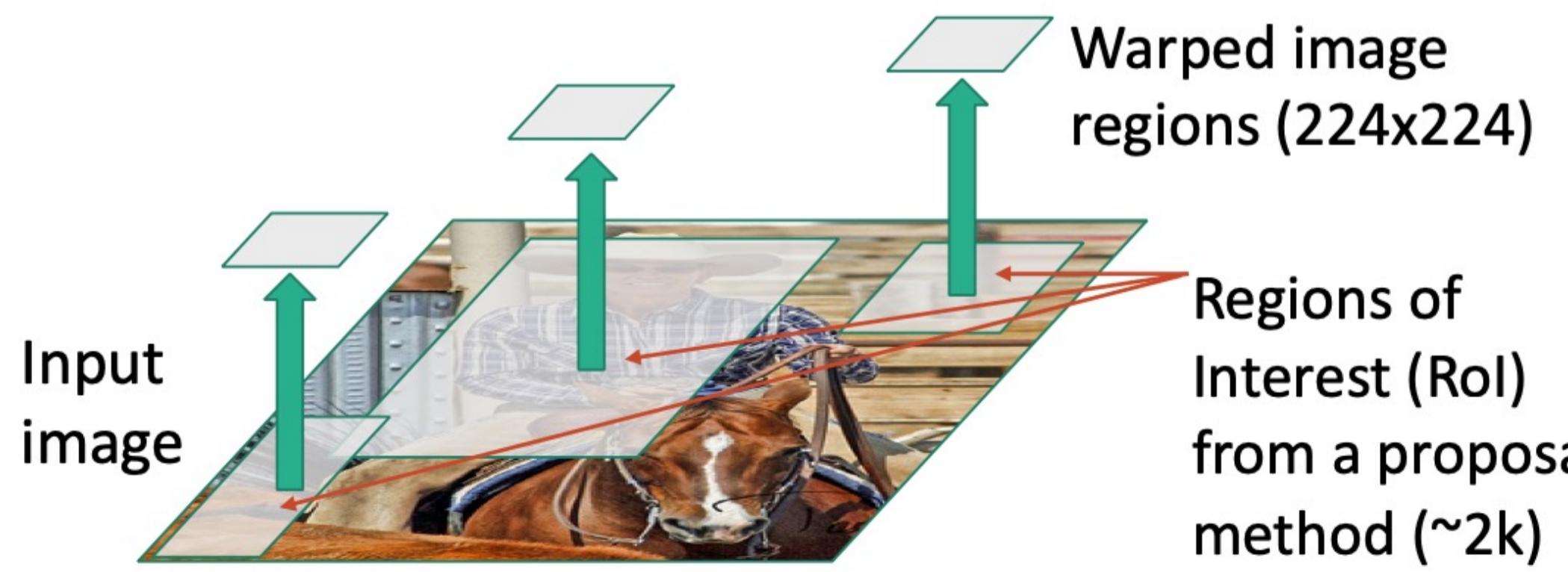
R-CNN: Region-based CNN

(1) select region proposals



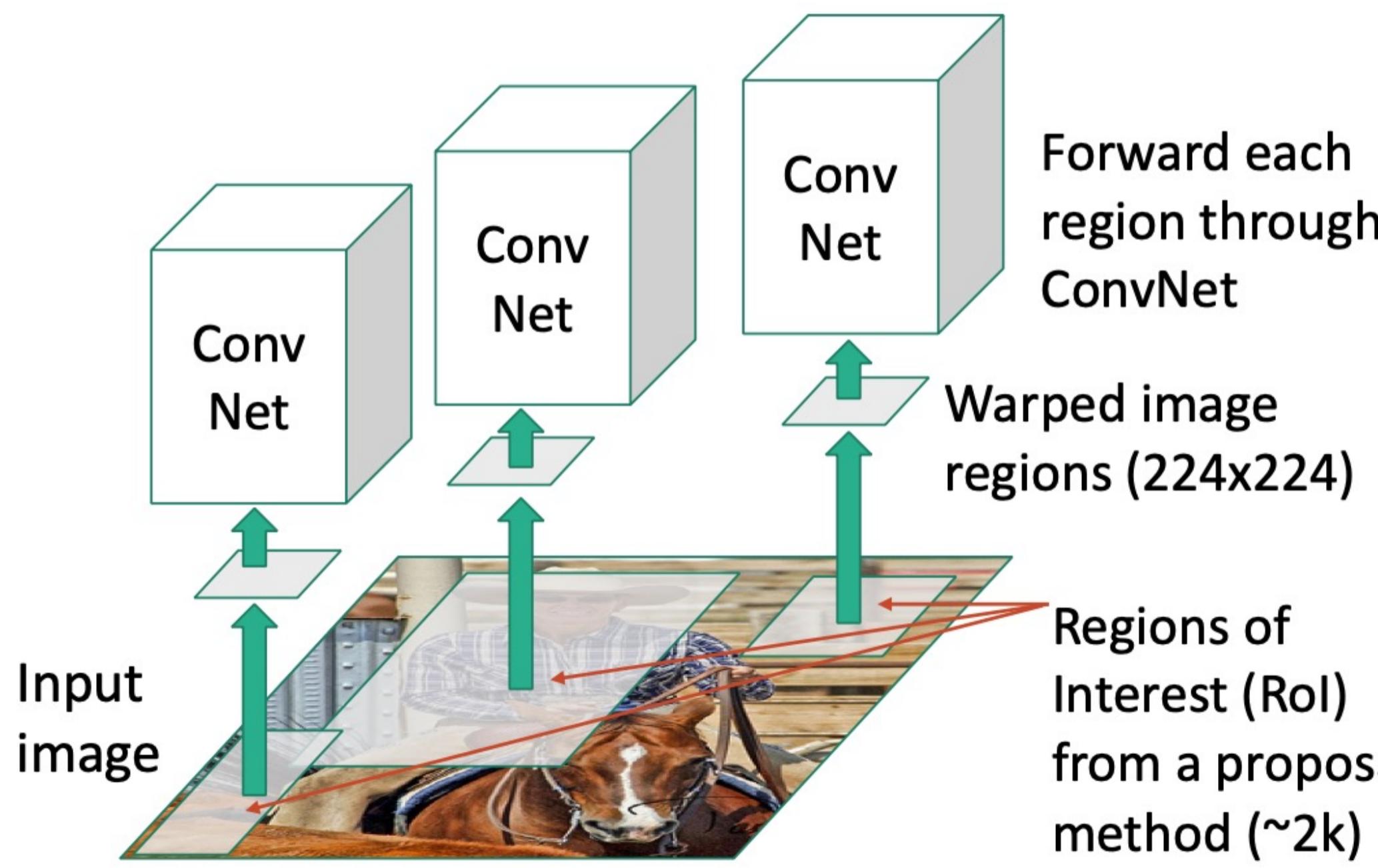
R-CNN: Region-based CNN

- (1) select region proposals
- (2) normalize in size

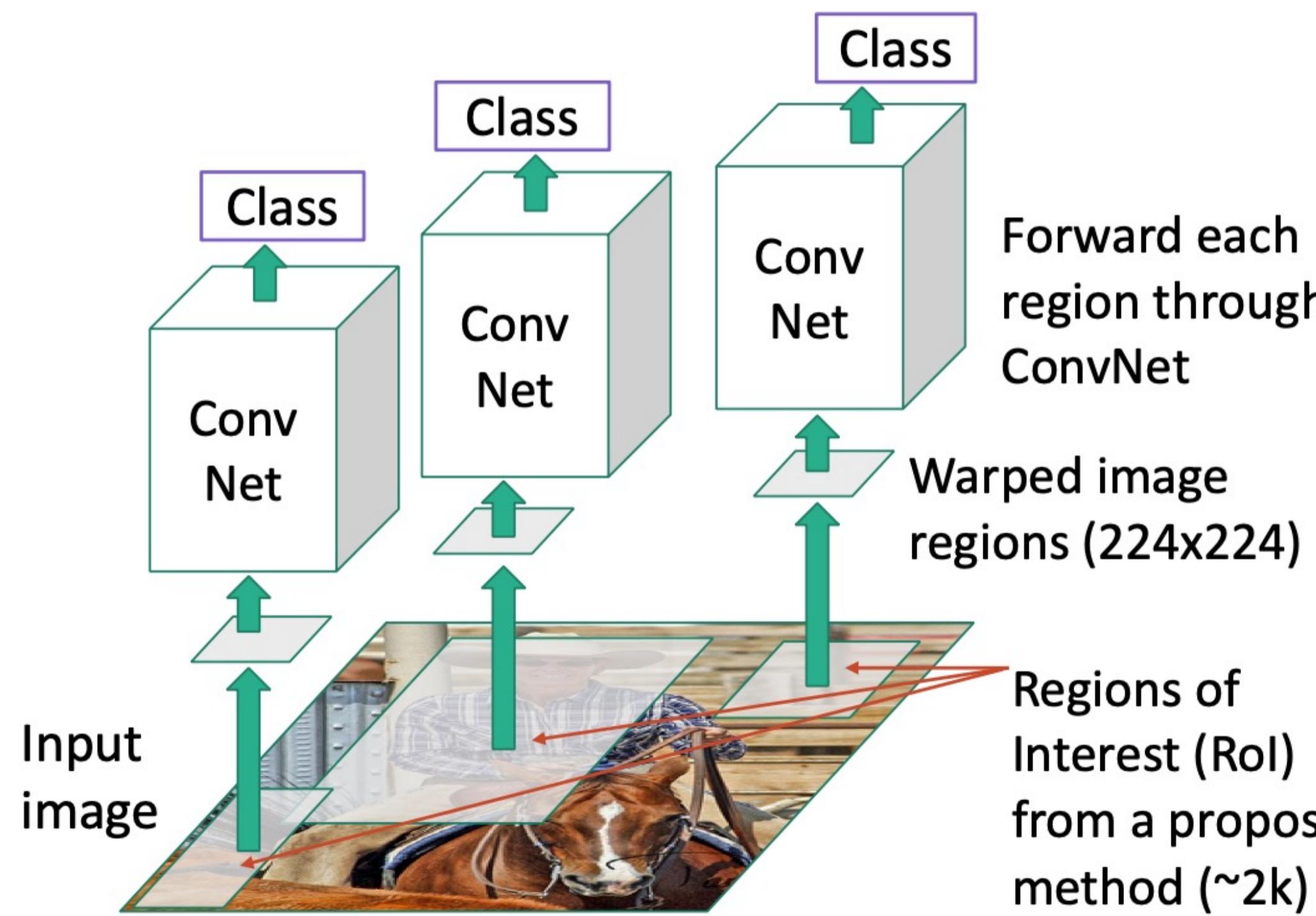


R-CNN: Region-based CNN

- (1) select region proposals
- (2) normalize in size
- (3) process each region

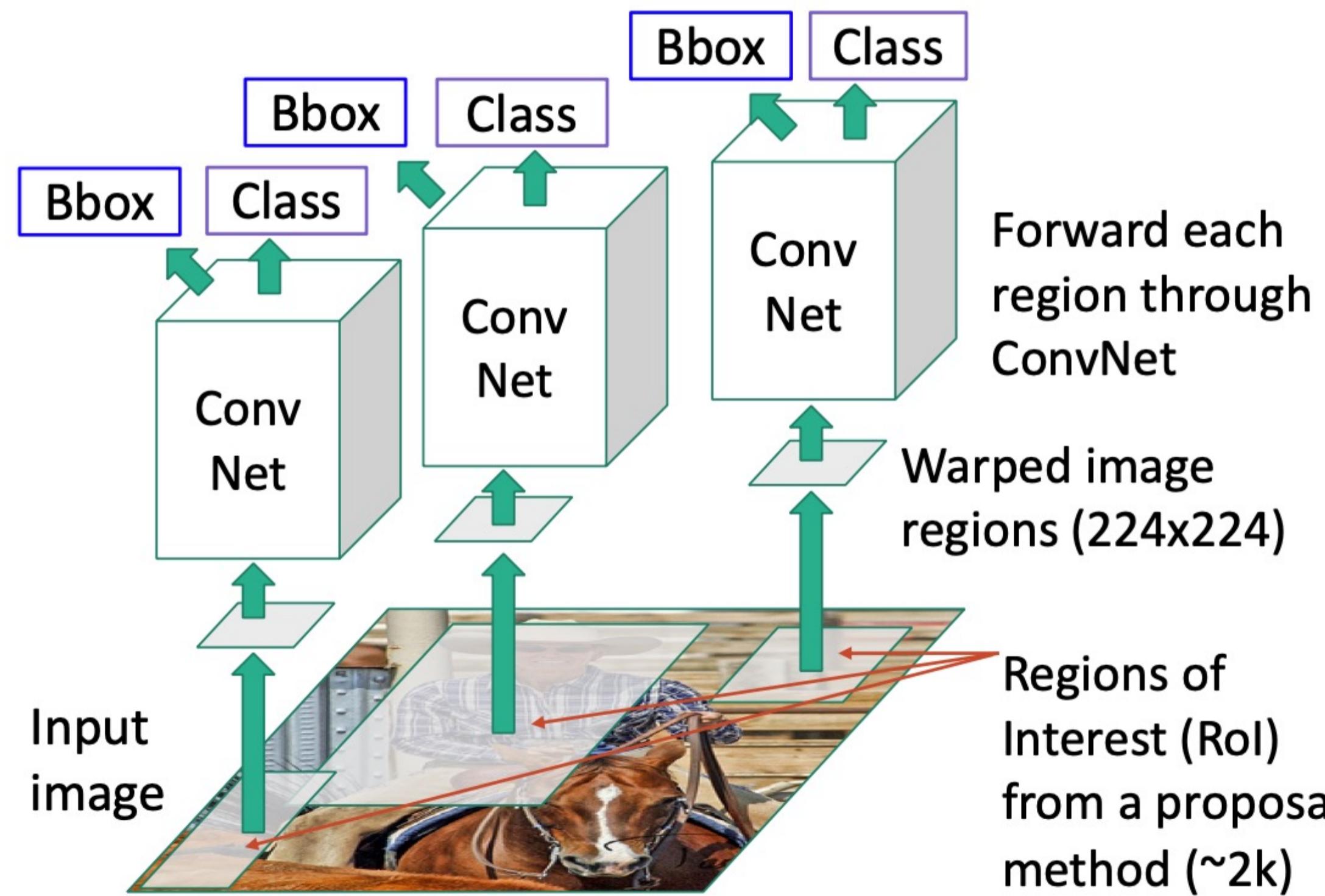


R-CNN: Region-based CNN



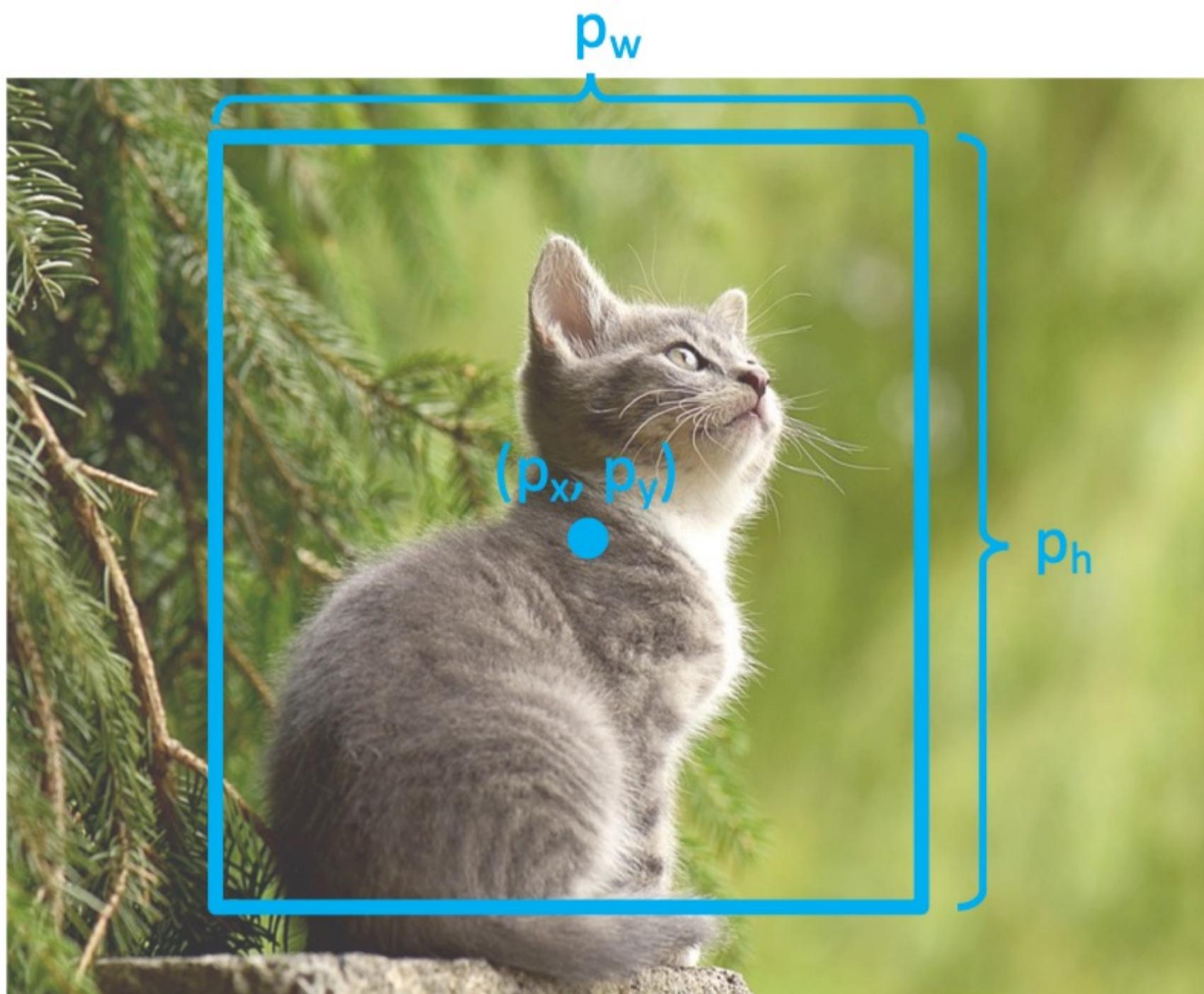
- (1) select region proposals
- (2) normalize in size
- (3) process each region
- (4) classify each region

R-CNN: Region-based CNN



- (1) select region proposals
- (2) normalize in size
- (3) process each region
- (4) classify each region
- (5) bounding box regression:
predict “transform” to
correct the ROI:
 (t_x, t_y, t_h, t_w)

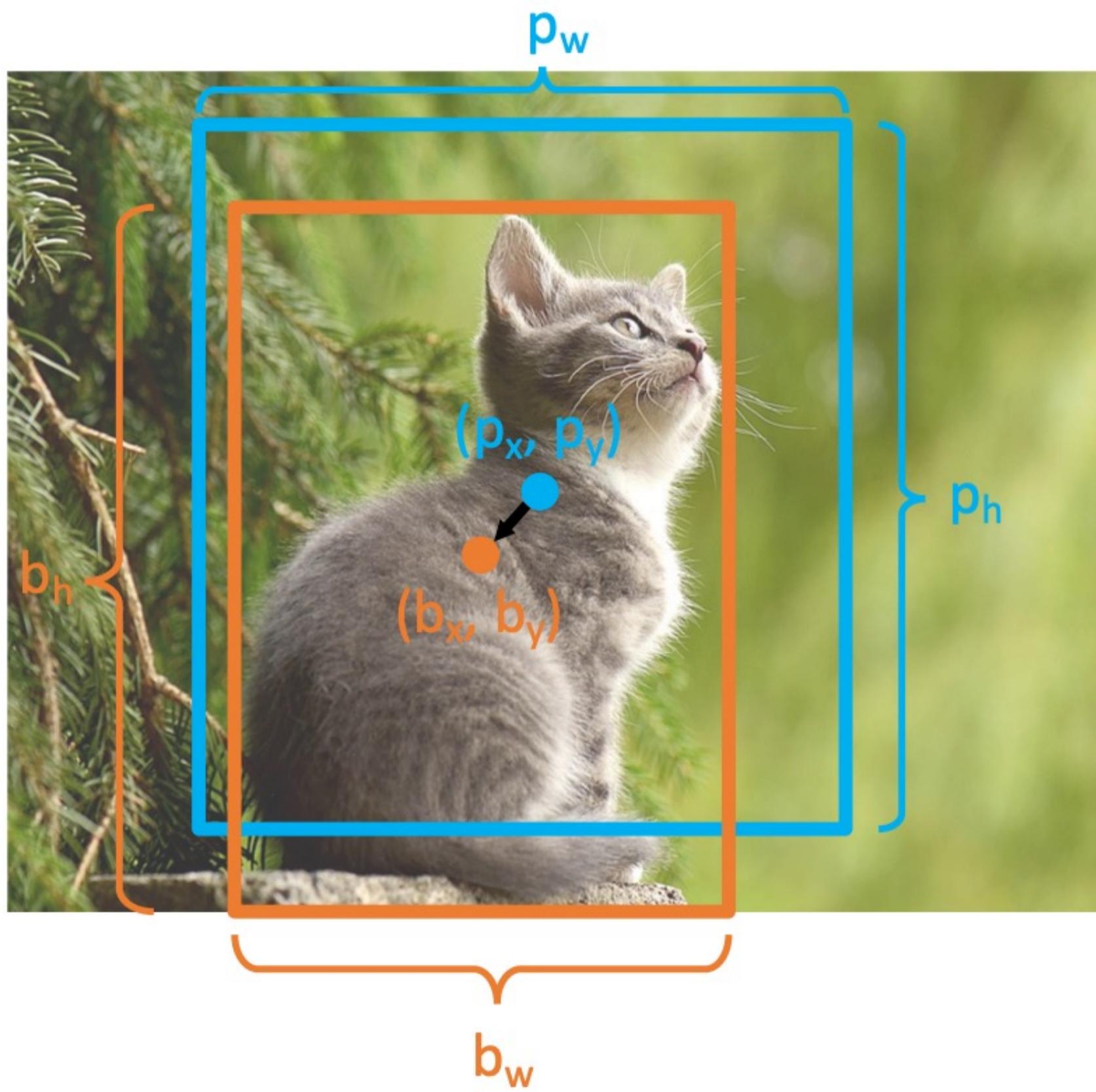
R-CNN: Bounding box regression



Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts **a transform** (t_x, t_y, t_w, t_h) to correct the region proposal

R-CNN: Bounding box regression



Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts **a transform** (t_x, t_y, t_w, t_h) to correct the region proposal

The **output box** is defined by:

$$b_x = p_x + p_w t_x$$

$$b_y = p_y + p_h t_y$$

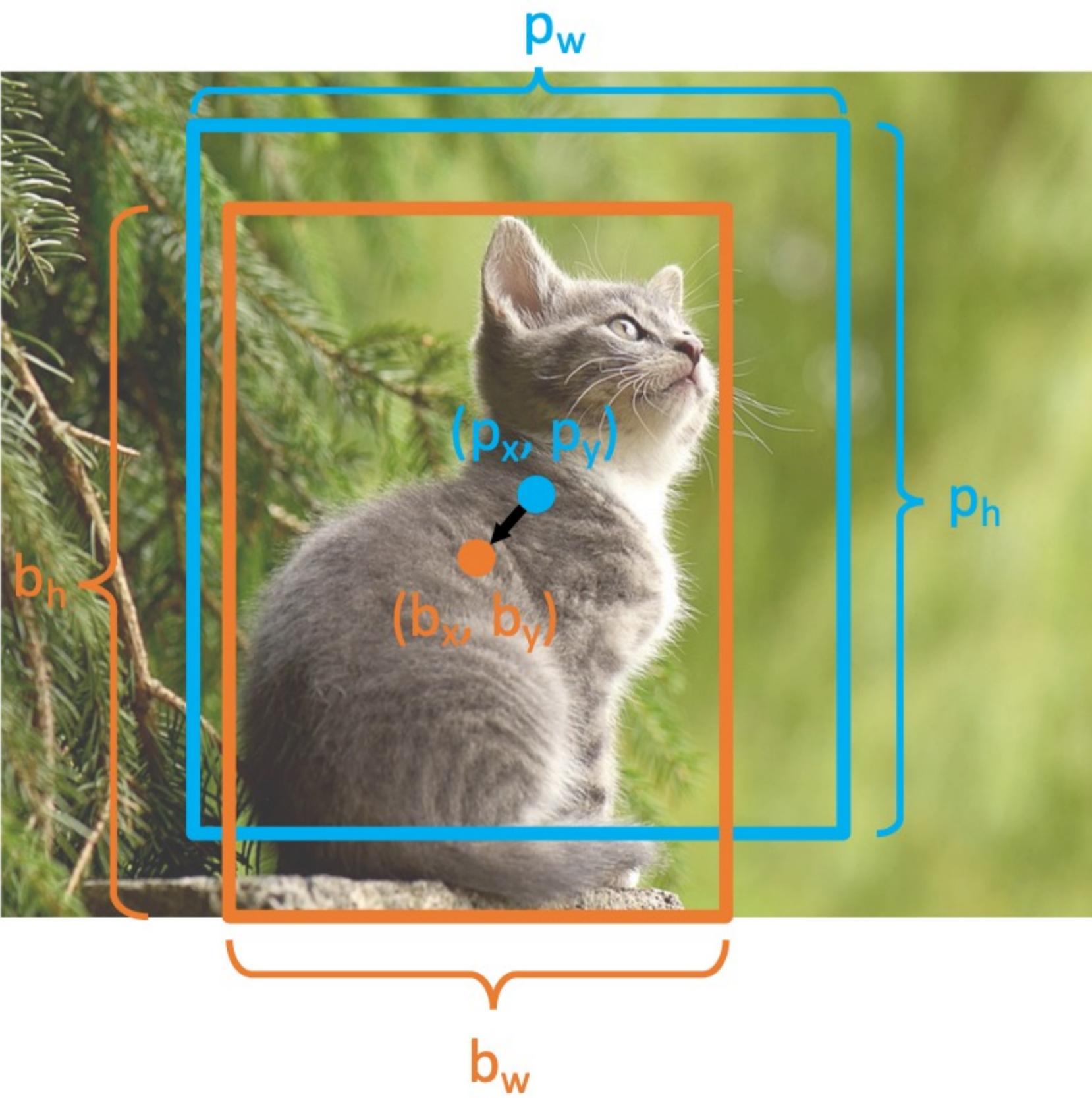
$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

shift center by amount relative to proposal size

scale proposal; exp ensures that scaling factor is >0

R-CNN: Bounding box regression



Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts **a transform** (t_x, t_y, t_w, t_h) to correct the region proposal

The **output box** is defined by:

$$b_x = p_x + p_w t_x$$

$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

Scale / translation invariance:

Transform encodes relative difference between proposal and output; important since CNN doesn't see absolute size or position after cropping

R-CNN: Bounding box regression

The **output box** is defined by:

$$b_x = p_x + p_w t_x$$

$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

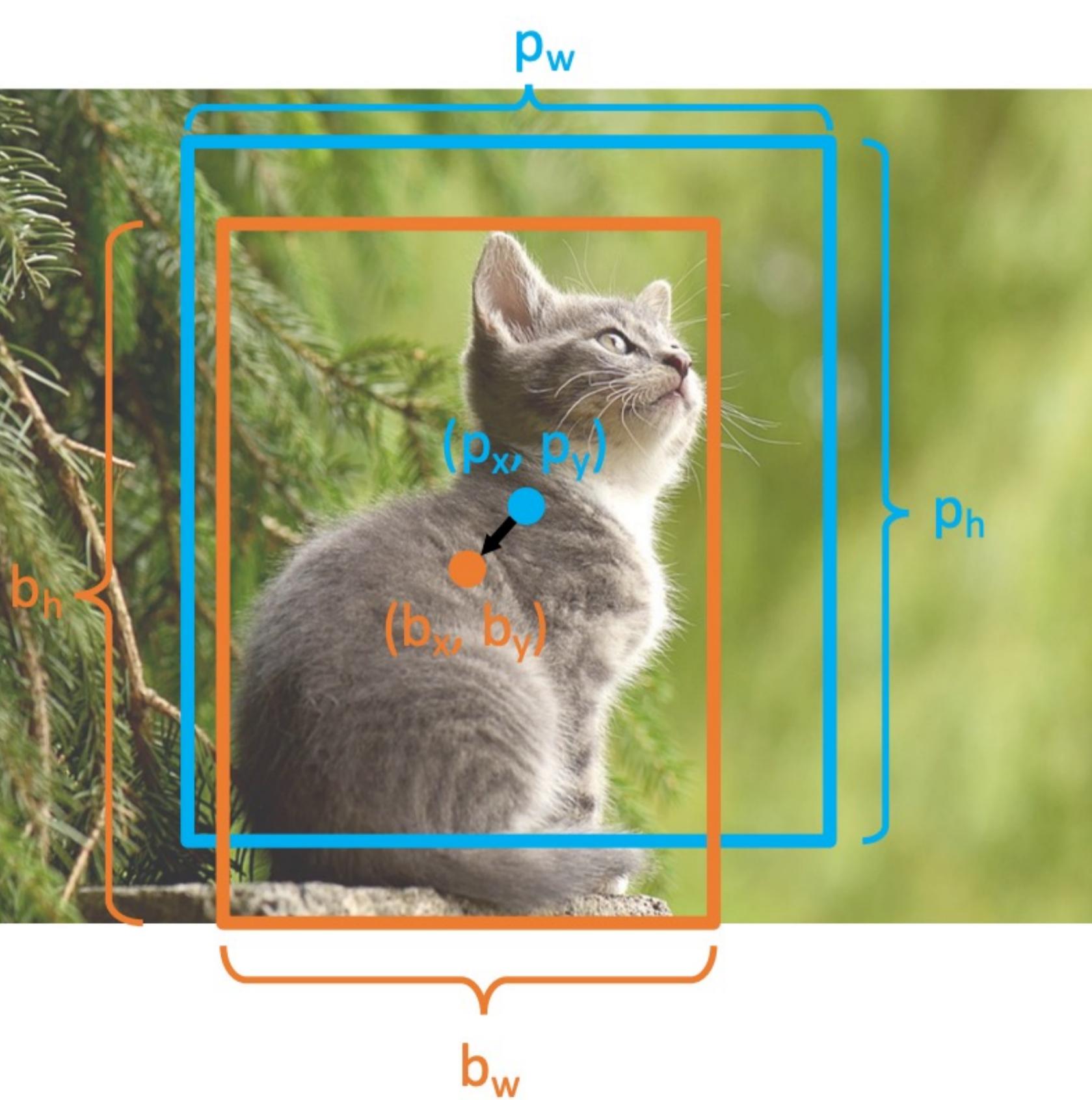
Given **proposal** and **target output**, we can solve for the **transform** the network should output.

$$t_x = (b_x - p_x)/p_w$$

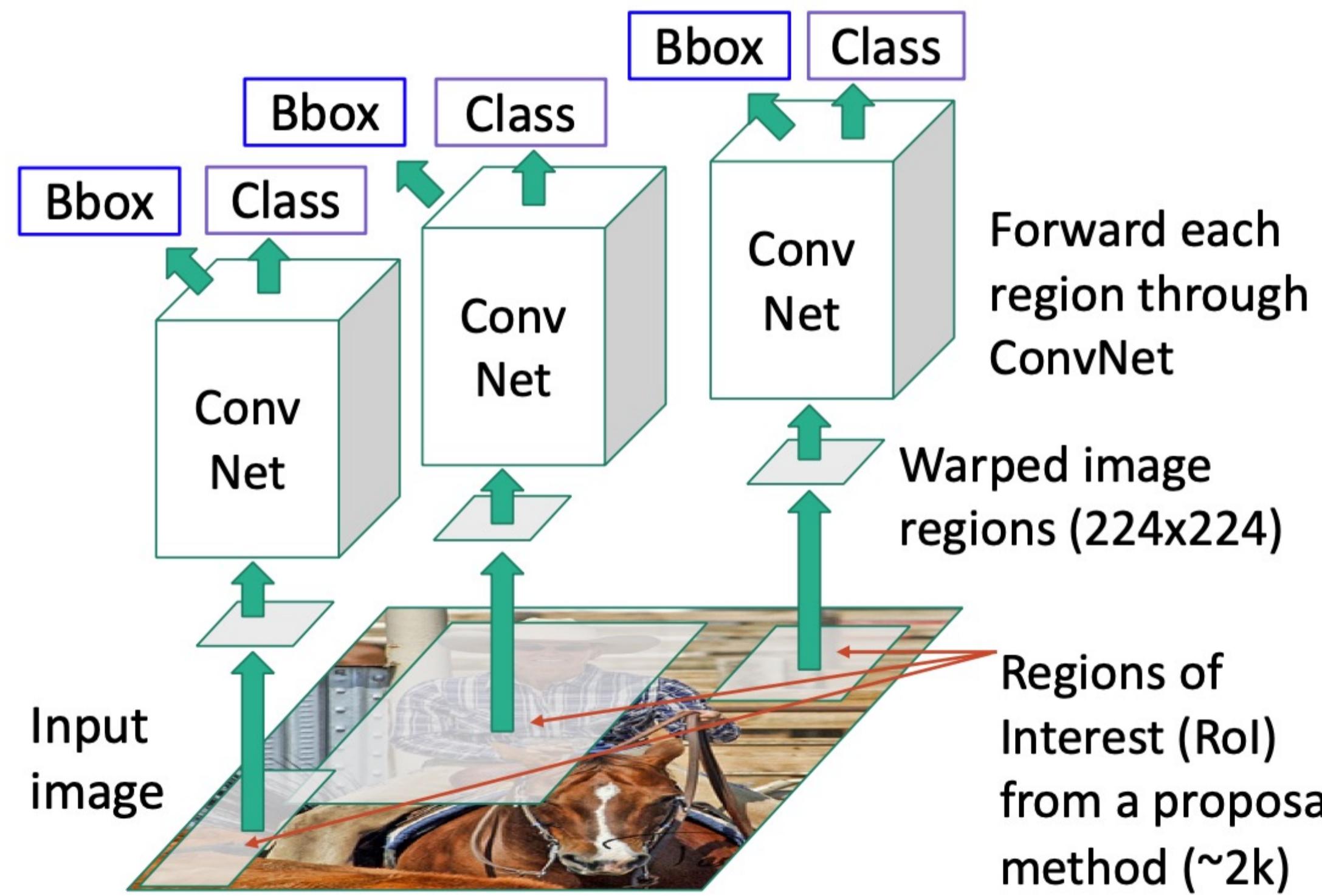
$$t_y = (b_y - p_y)/p_h$$

$$t_w = \log(b_w/p_w)$$

$$t_h = \log(b_h/p_h)$$

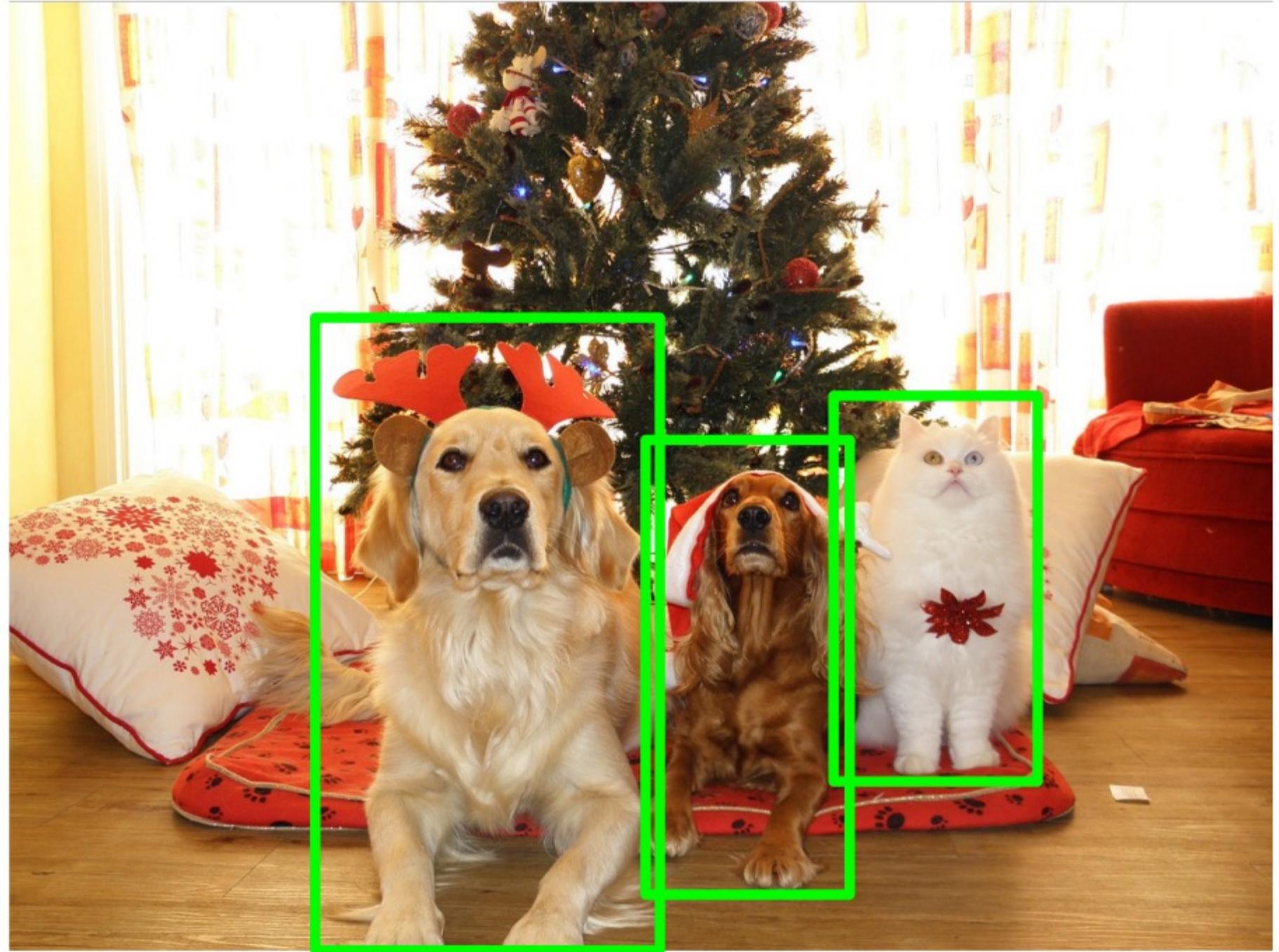


R-CNN: Region-based CNN



- (1) select region proposals
- (2) normalize in size
- (3) process each region
- (4) classify each region
- (5) bounding box regression:
predict “transform” to
correct the ROI:
 (t_x, t_y, t_h, t_w)

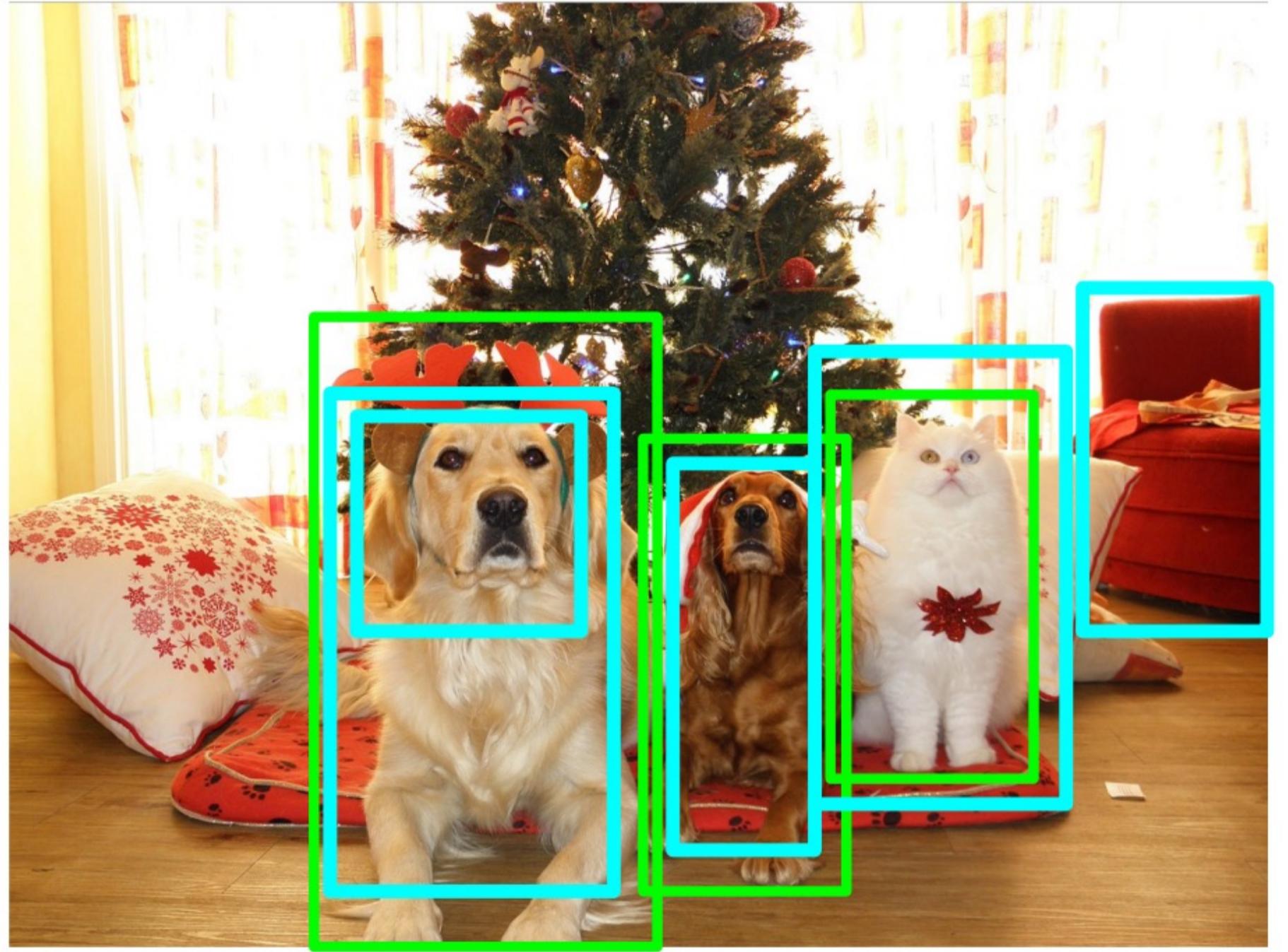
R-CNN training



Ground-Truth boxes



R-CNN training

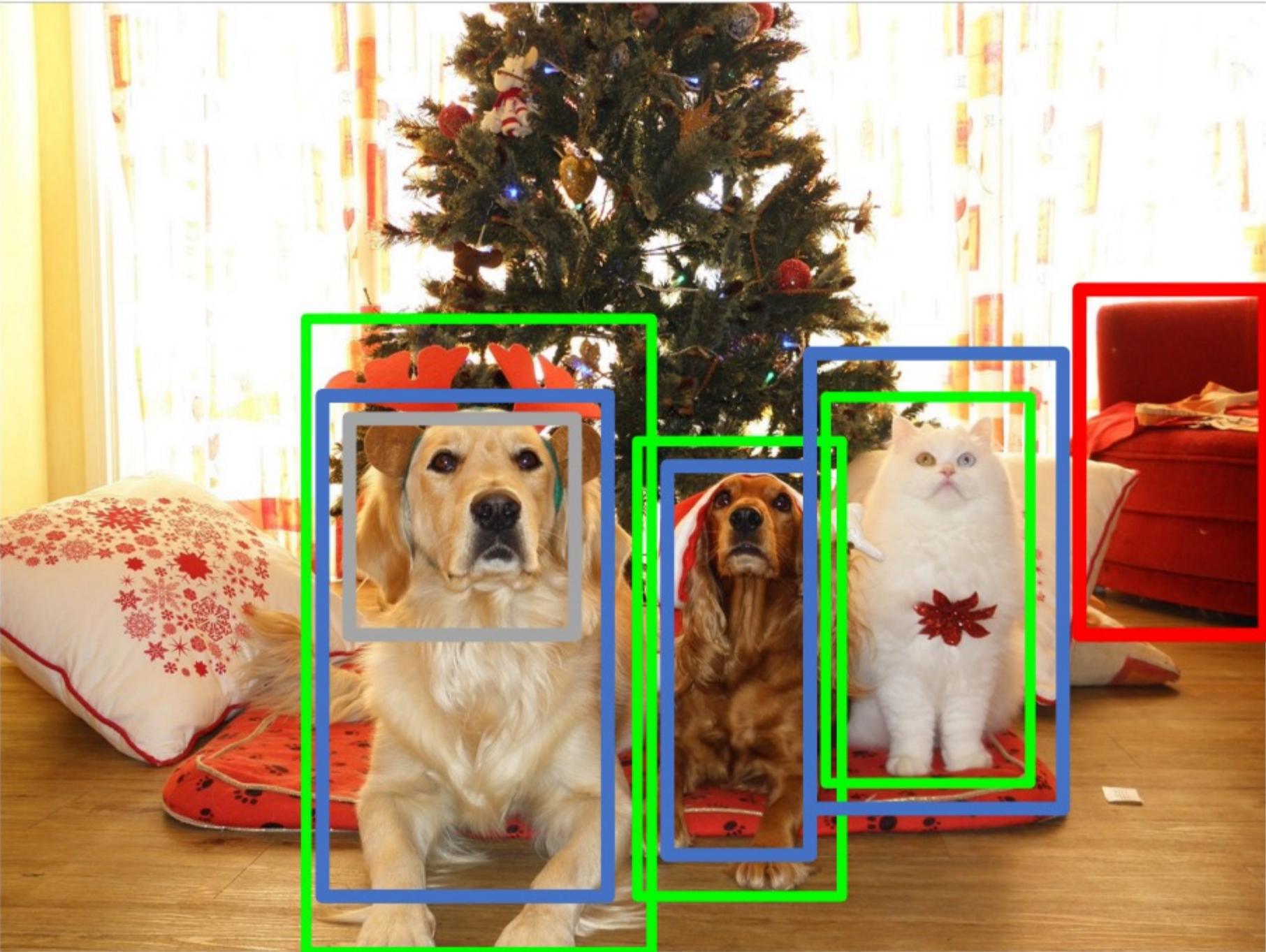


Ground-Truth boxes

Region Proposals



R-CNN training



Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes:

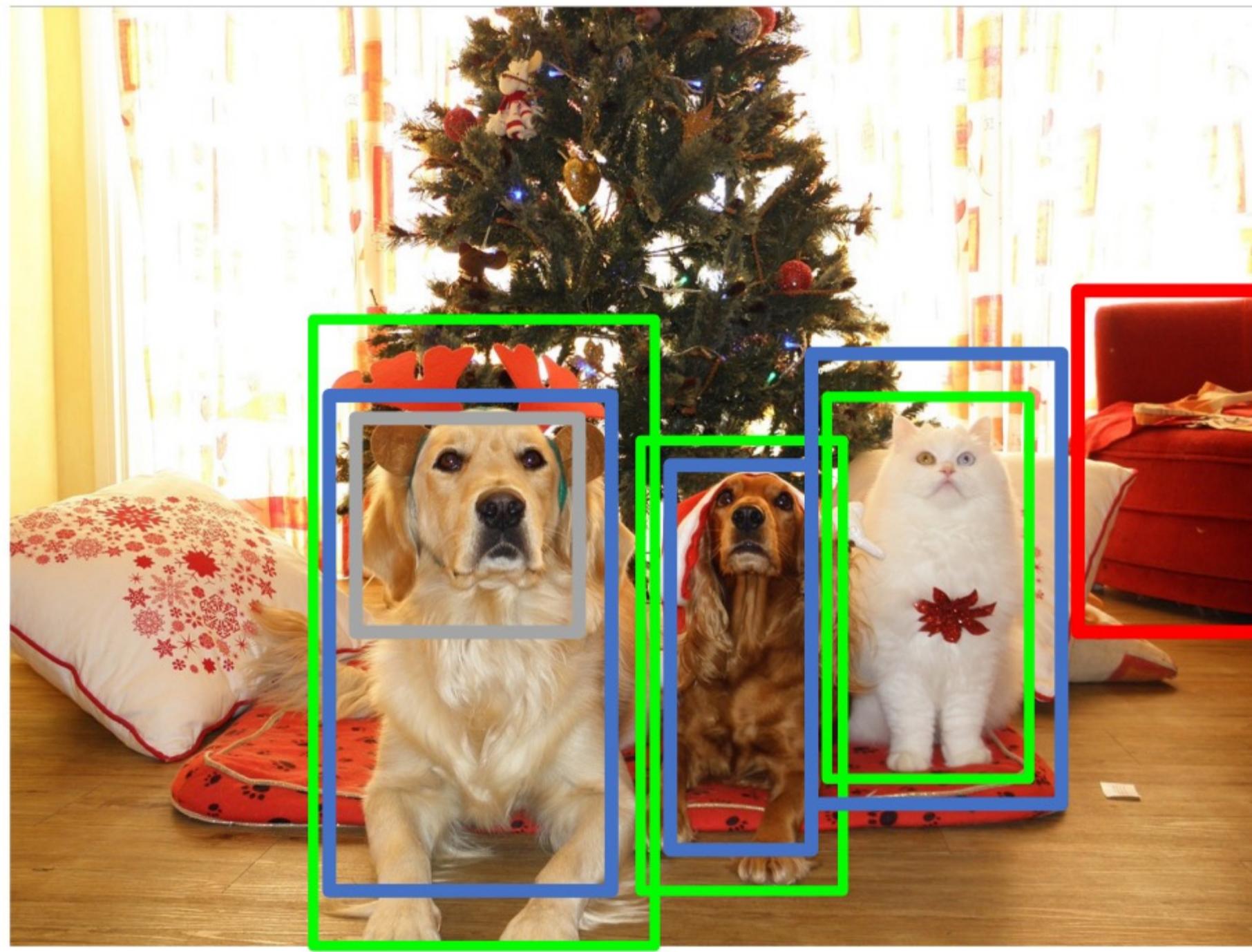
positive: > 0.5 IoU with GT box

negative: < 0.3 IoU with all GT boxes

neutral: between 0.3 and 0.5 IoU with GT boxes



R-CNN training



GT Boxes

Positive

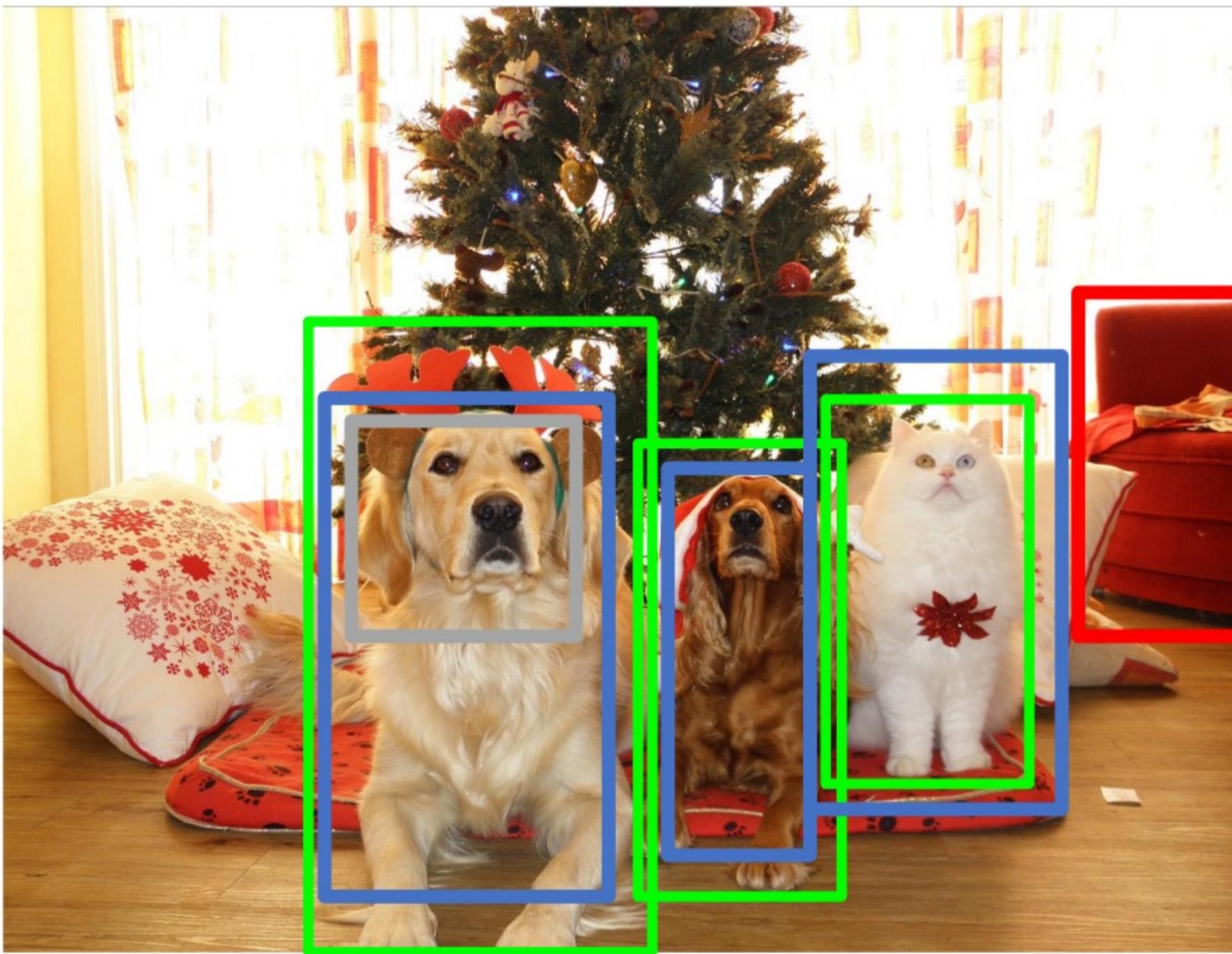
Neutral

Negative



Crop pixels from each positive and negative proposal, resize to 224x224

R-CNN training



GT Boxes

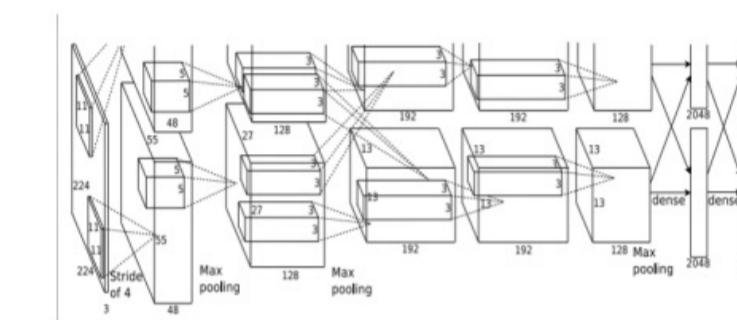
Neutral

Positive

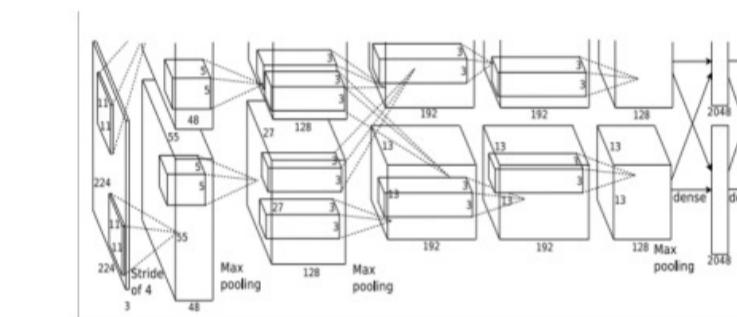
Negative



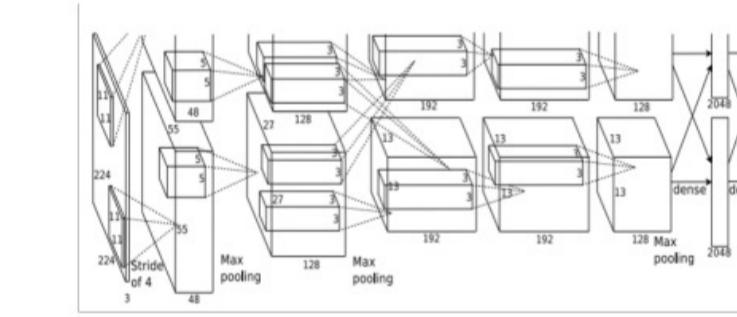
Run each region through CNN
Positive regions: predict class and transform
Negative regions: just predict class



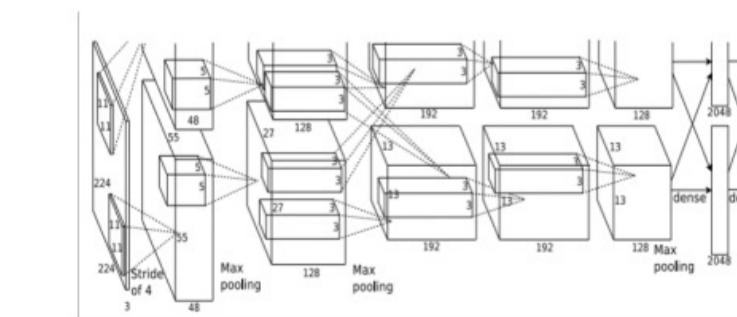
class target: dog
box target: →



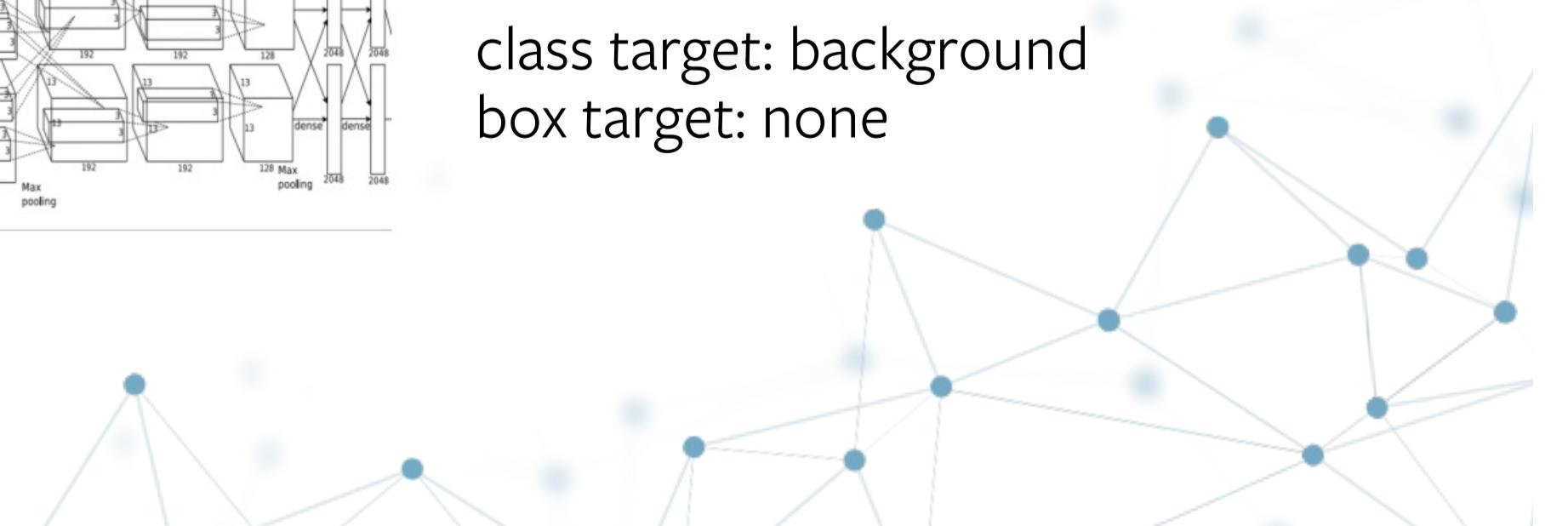
class target: cat
box target: →



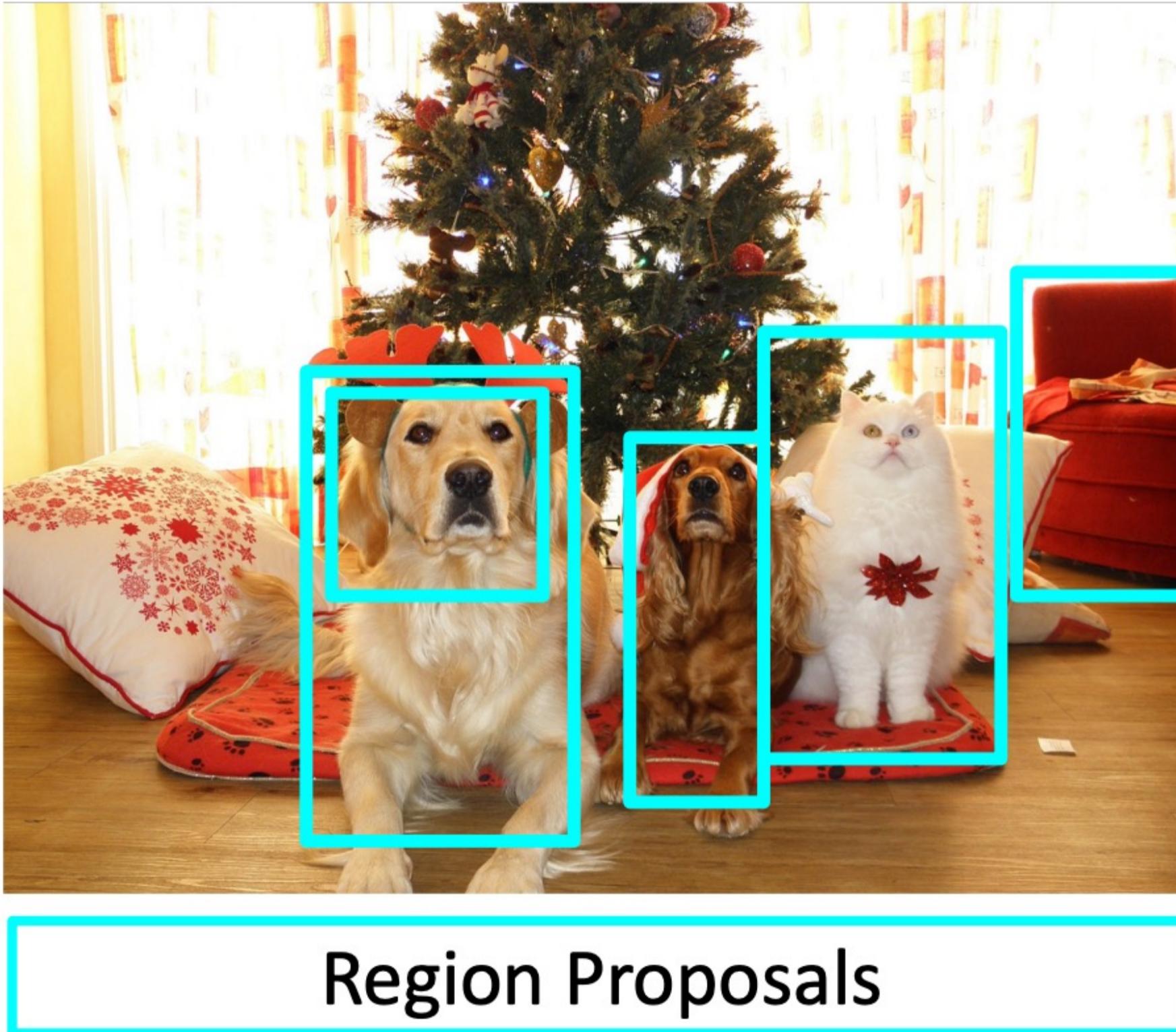
class target: dog
box target: →



class target: background
box target: none



R-CNN testing

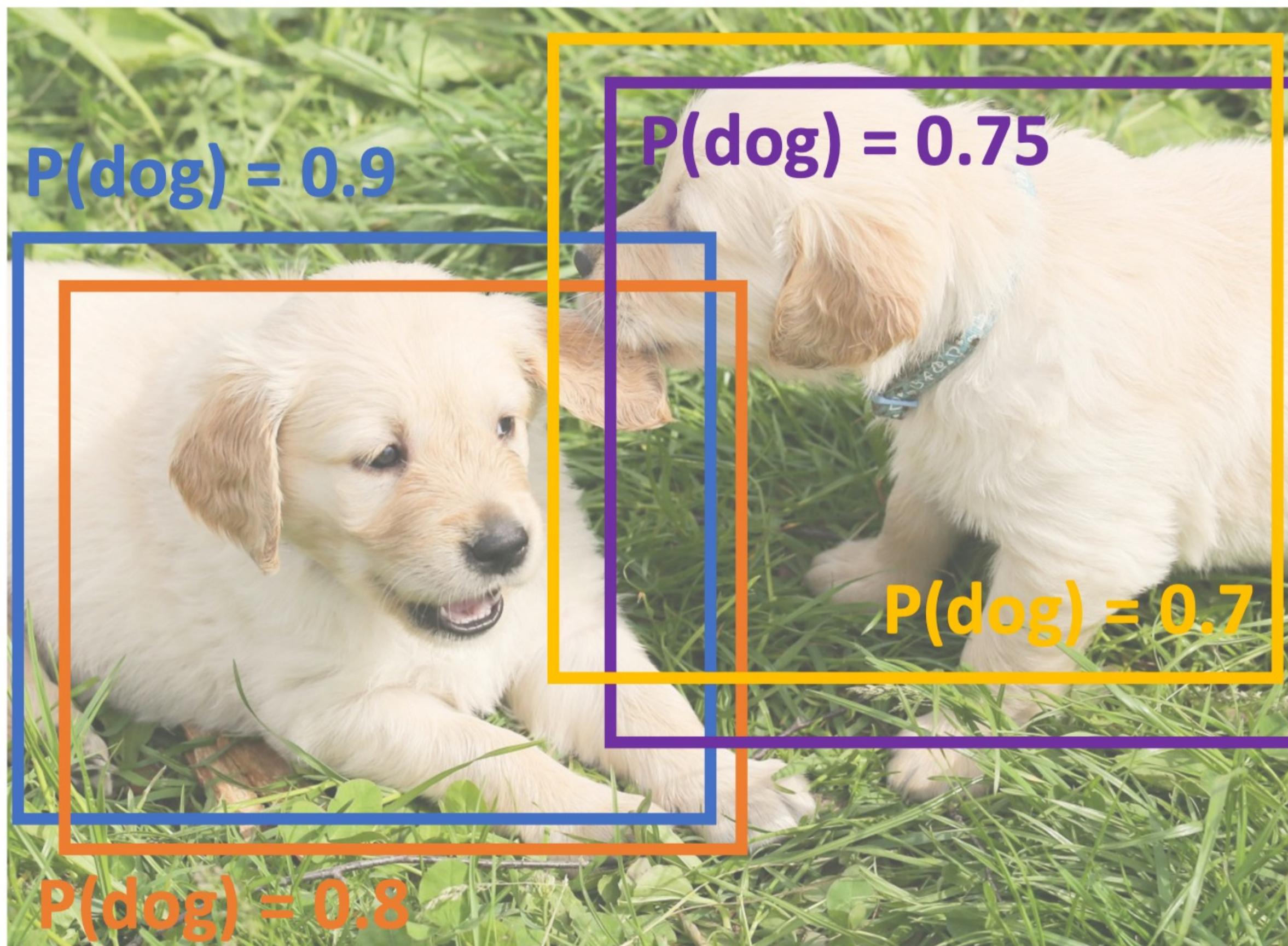


- (1) Run proposal method
- (2) Run CNN for each proposal to get class scores, transforms
- (3) Threshold class scores to get a set of detections

2 problems:

- CNN often outputs overlapping boxes
- How to set thresholds?

Overlapping boxes

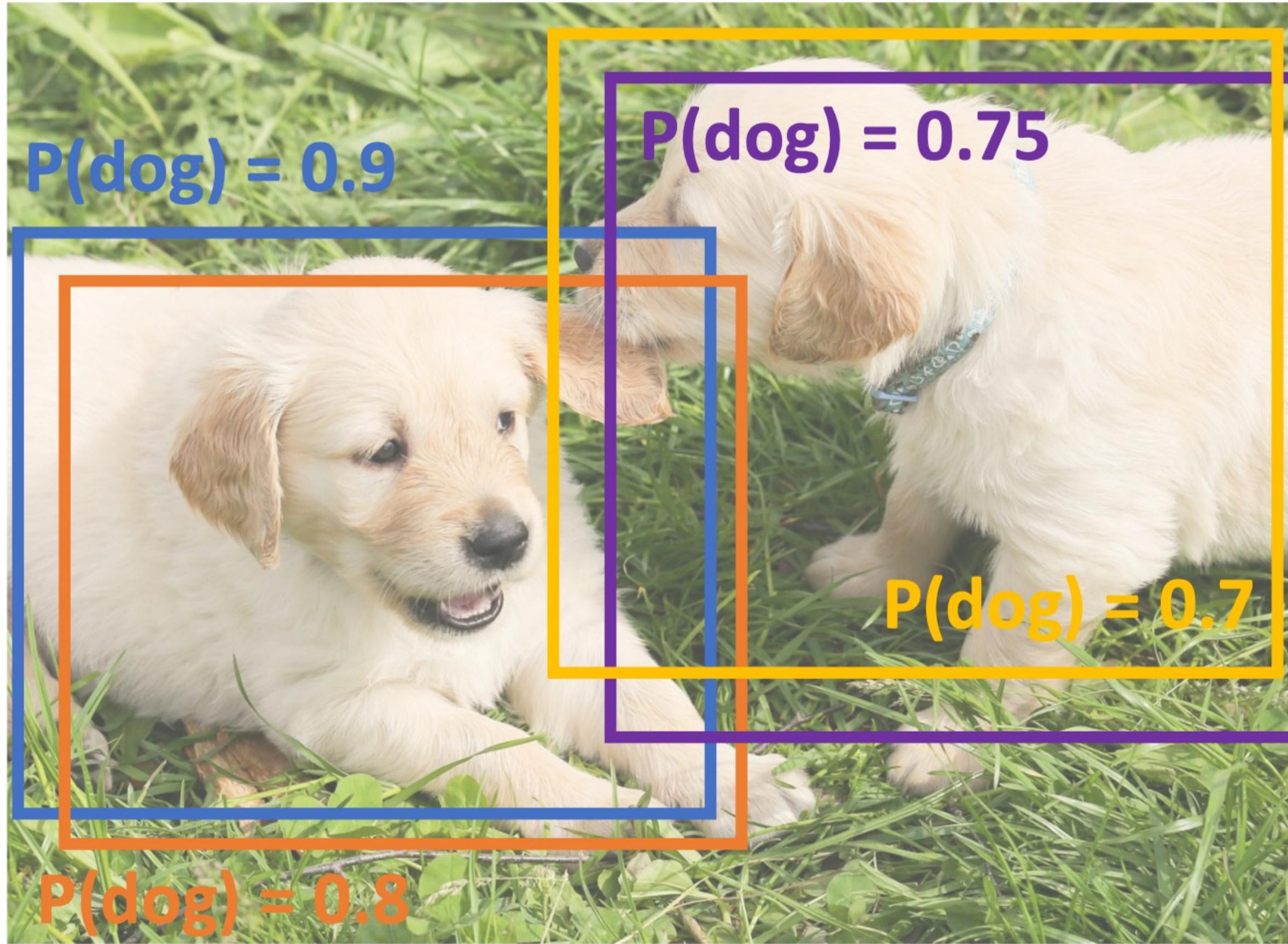


Problem: Object detectors often output many overlapping detections

Solution: post-process raw detections using **Non-Max Suppression (NMS)**

- (1) select highest-scoring box
- (2) eliminate lower-scoring boxes with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
- (3) If any boxes remain, GOTO 1

Overlapping boxes



- (1) select highest-scoring box
- (2) eliminate lower-scoring boxes
with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
- (3) If any boxes remain, GOTO 1

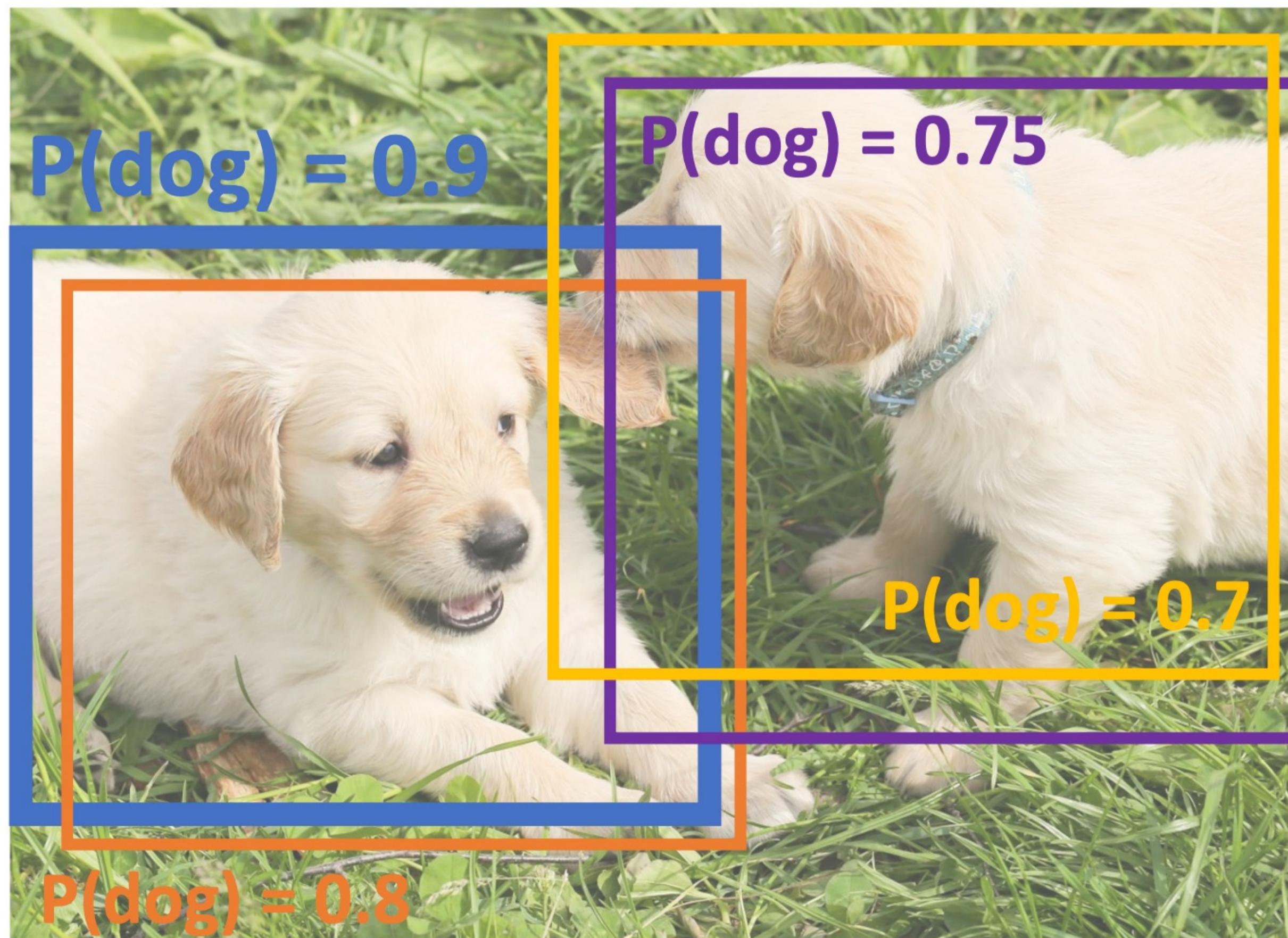
$$\text{IoU}(\square, \blacksquare) = \mathbf{0.78}$$

$$\text{IoU}(\square, \blacksquare) = 0.05$$

$$\text{IoU}(\square, \blacksquare) = 0.07$$



Overlapping boxes



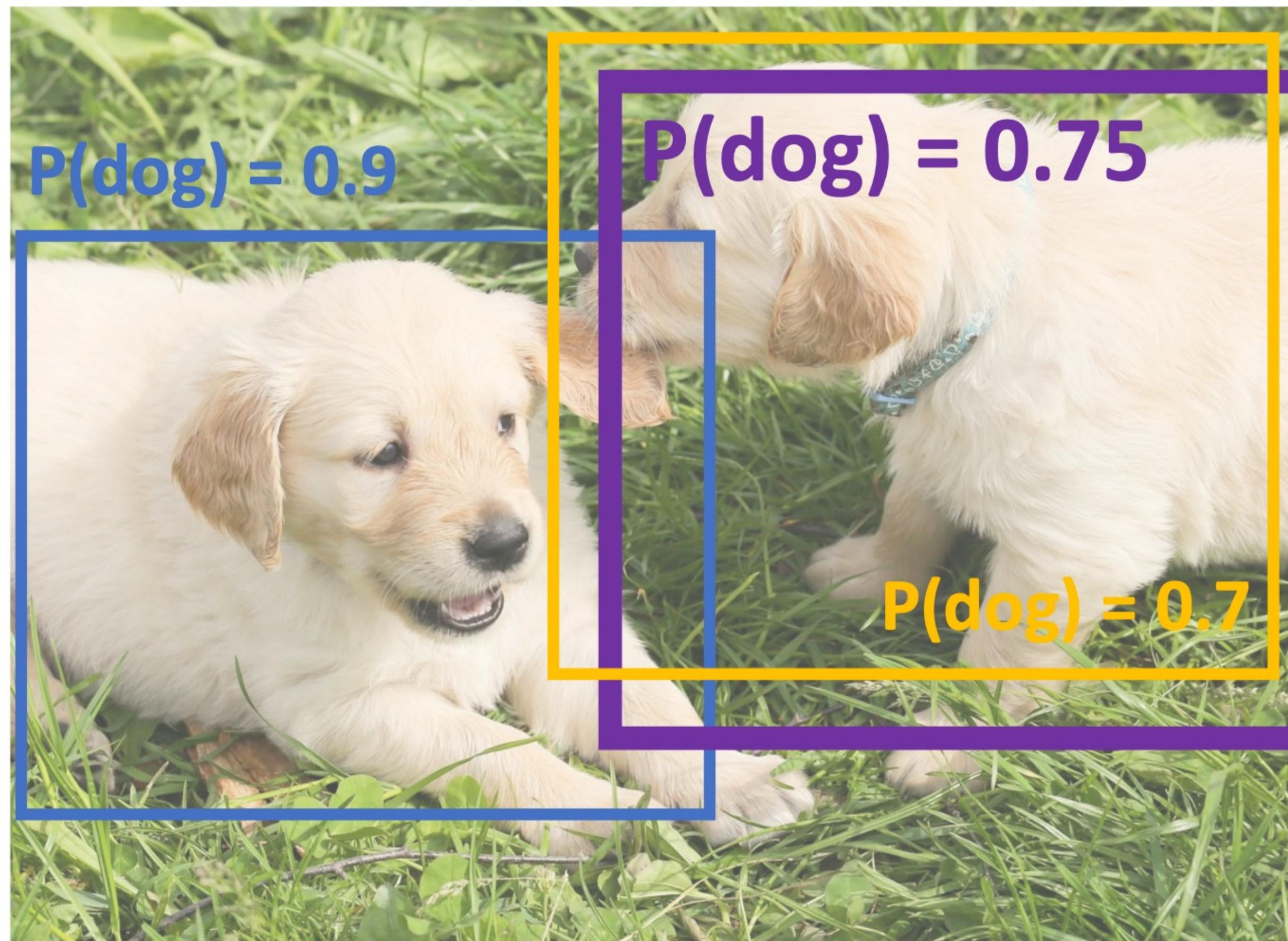
- (1) select highest-scoring box
- (2) eliminate lower-scoring boxes
with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
- (3) If any boxes remain, GOTO 1

$$\text{IoU}(\square, \blacksquare) = \mathbf{0.78}$$

$$\text{IoU}(\square, \blacksquare) = 0.05$$

$$\text{IoU}(\square, \blacksquare) = 0.07$$

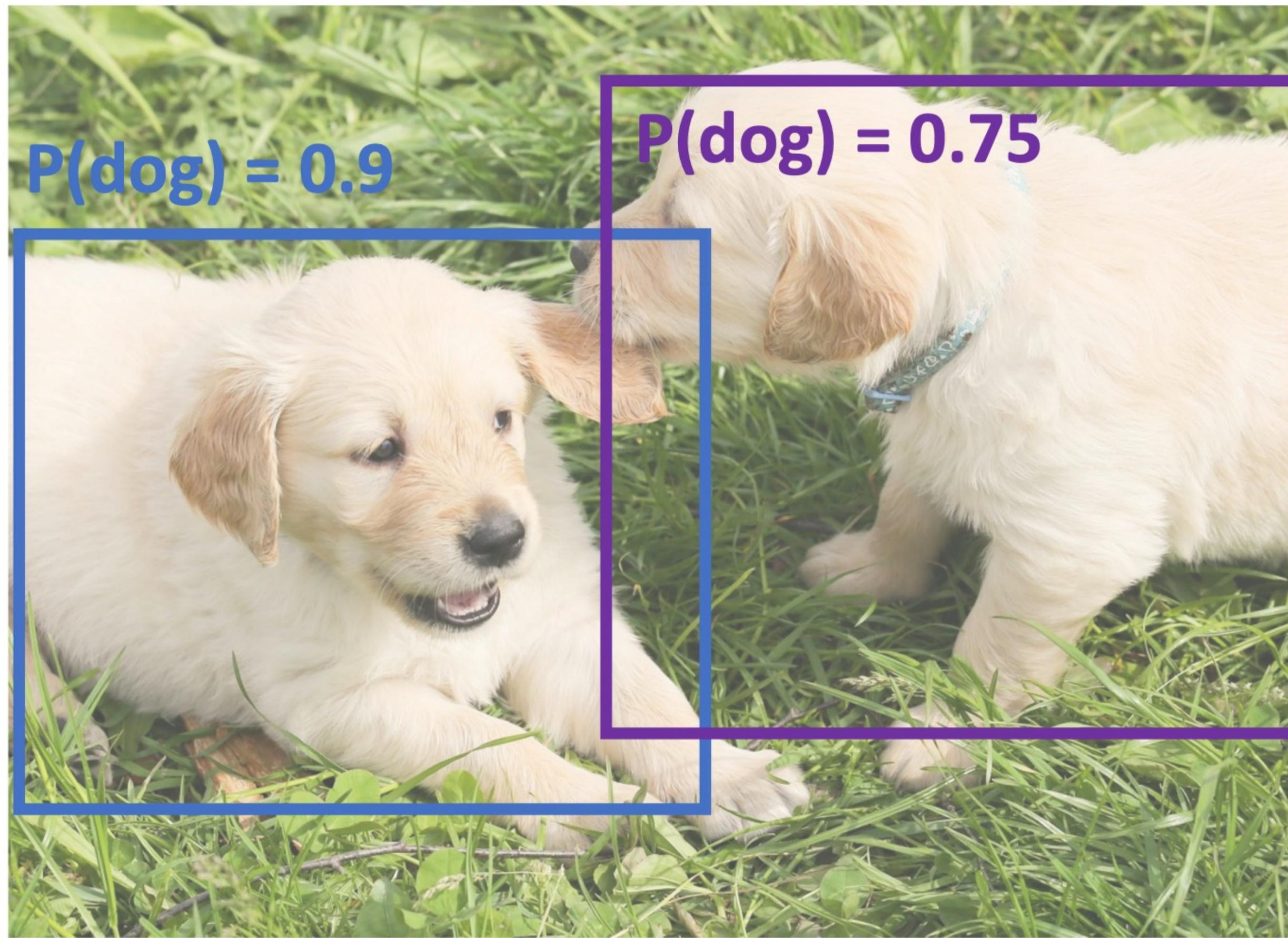
Overlapping boxes



- (1) select highest-scoring box
- (2) eliminate lower-scoring boxes
with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
- (3) If any boxes remain, GOTO 1

$$\text{IoU}(\text{purple box}, \text{yellow box}) = \mathbf{0.74}$$

Overlapping boxes



- (1) select highest-scoring box
- (2) eliminate lower-scoring boxes
with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
- (3) If any boxes remain, GOTO 1



Overlapping boxes

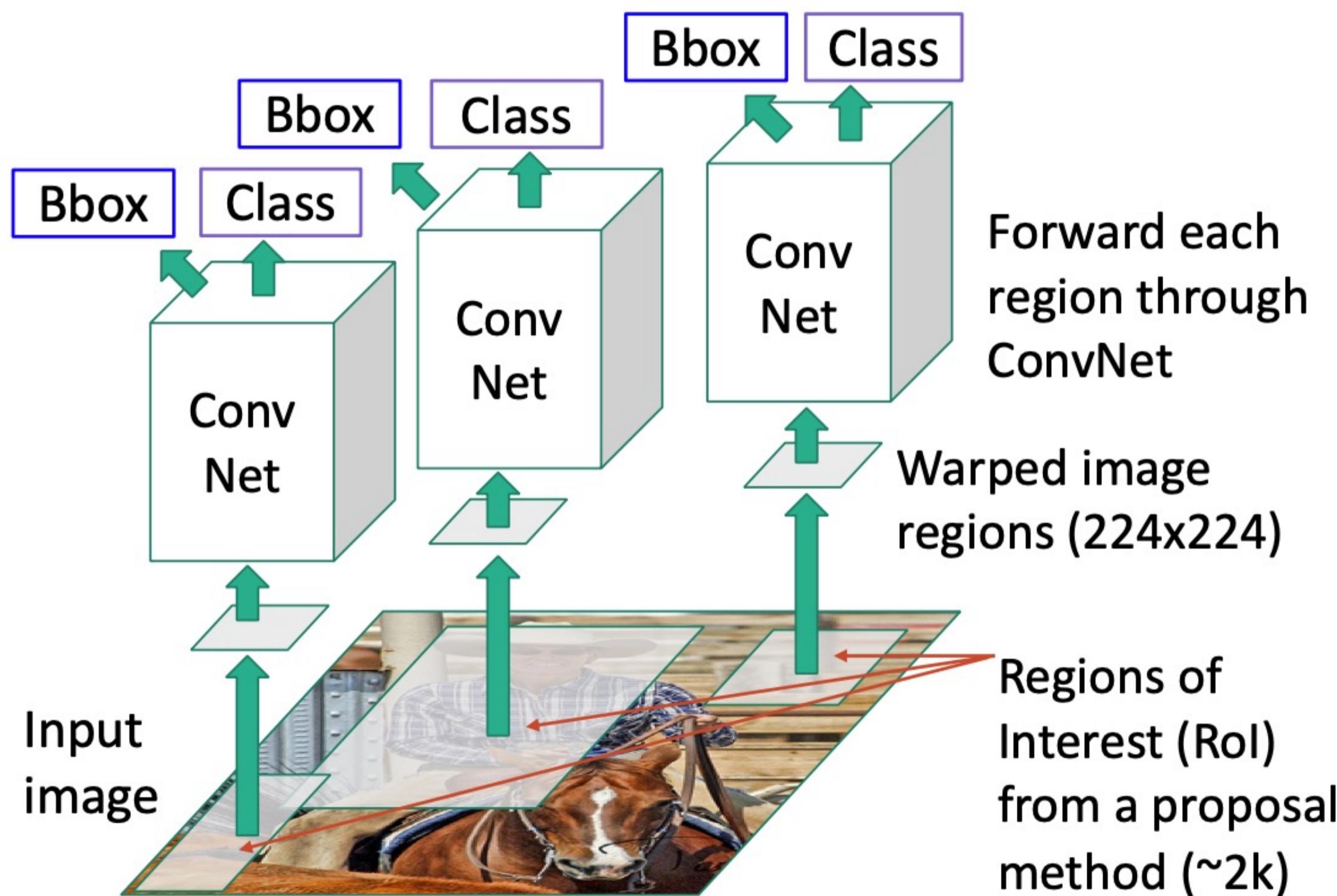


Remaining problem:
NMS may eliminate
“good” boxes when
objects are highly
overlapping – there is still
no good solution!

Today's focus

- Principles of multi-object detection
- Region proposals
- Bounding box regression
- R-CNN training and inference
- Non-Max suppression
- Modern object detectors

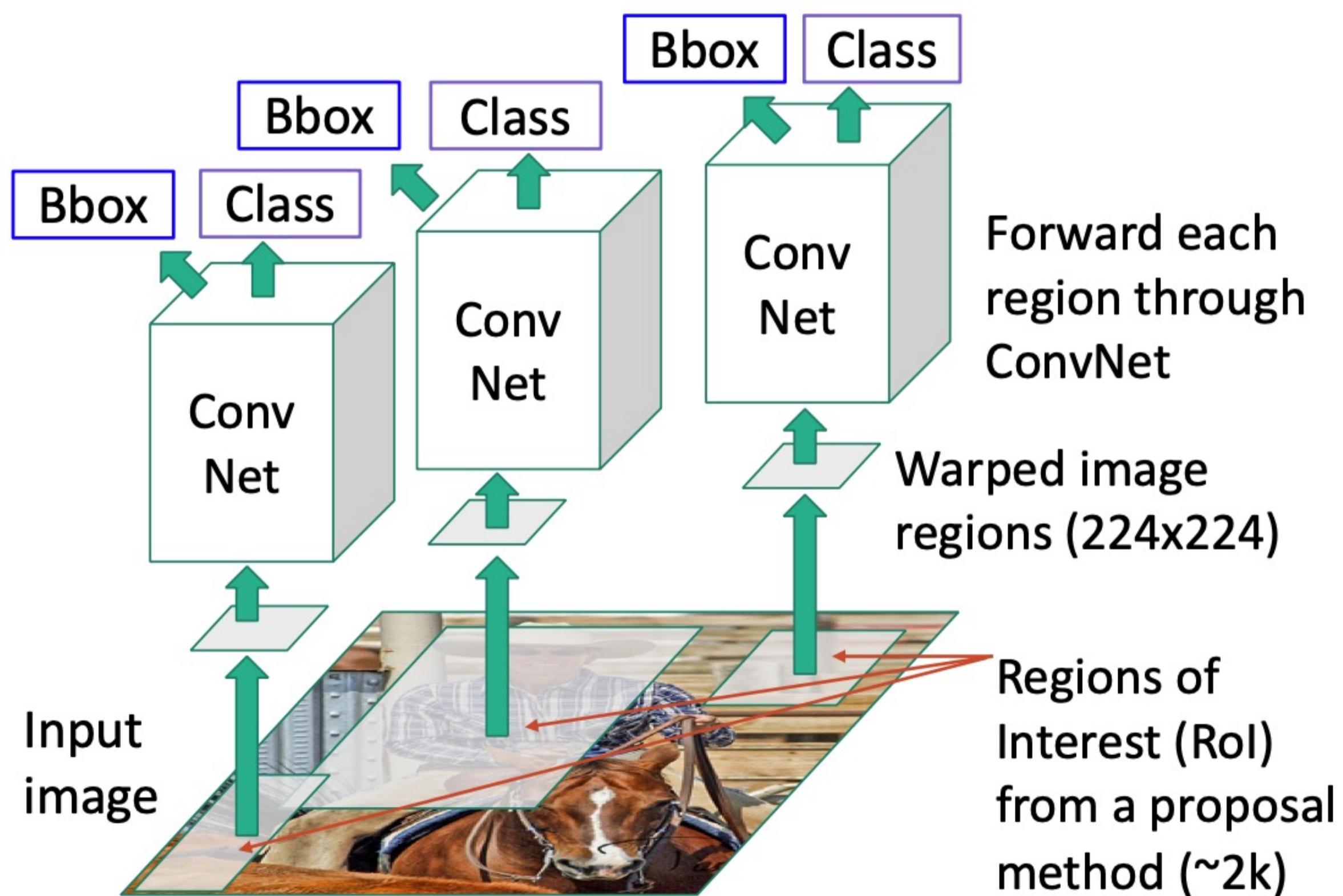
Recap: R-CNN



Problem: Very slow!

Need to do 2000 forward passes through CNN per image

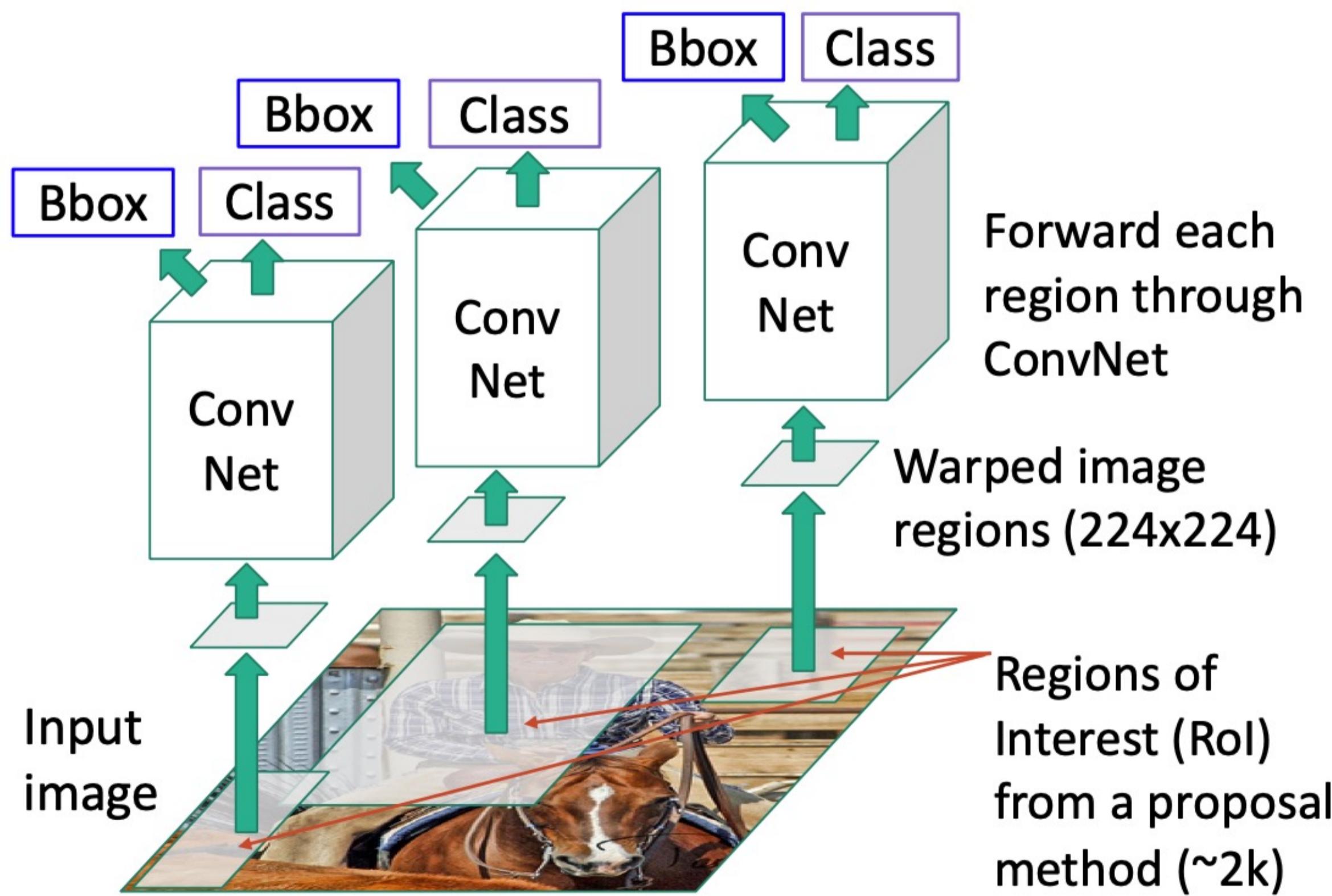
Recap: R-CNN



Problem: Very slow!

Need to do 2000 forward passes through CNN per image

Recap: R-CNN



Problem: Very slow!

Need to do 2000 forward passes through CNN per image

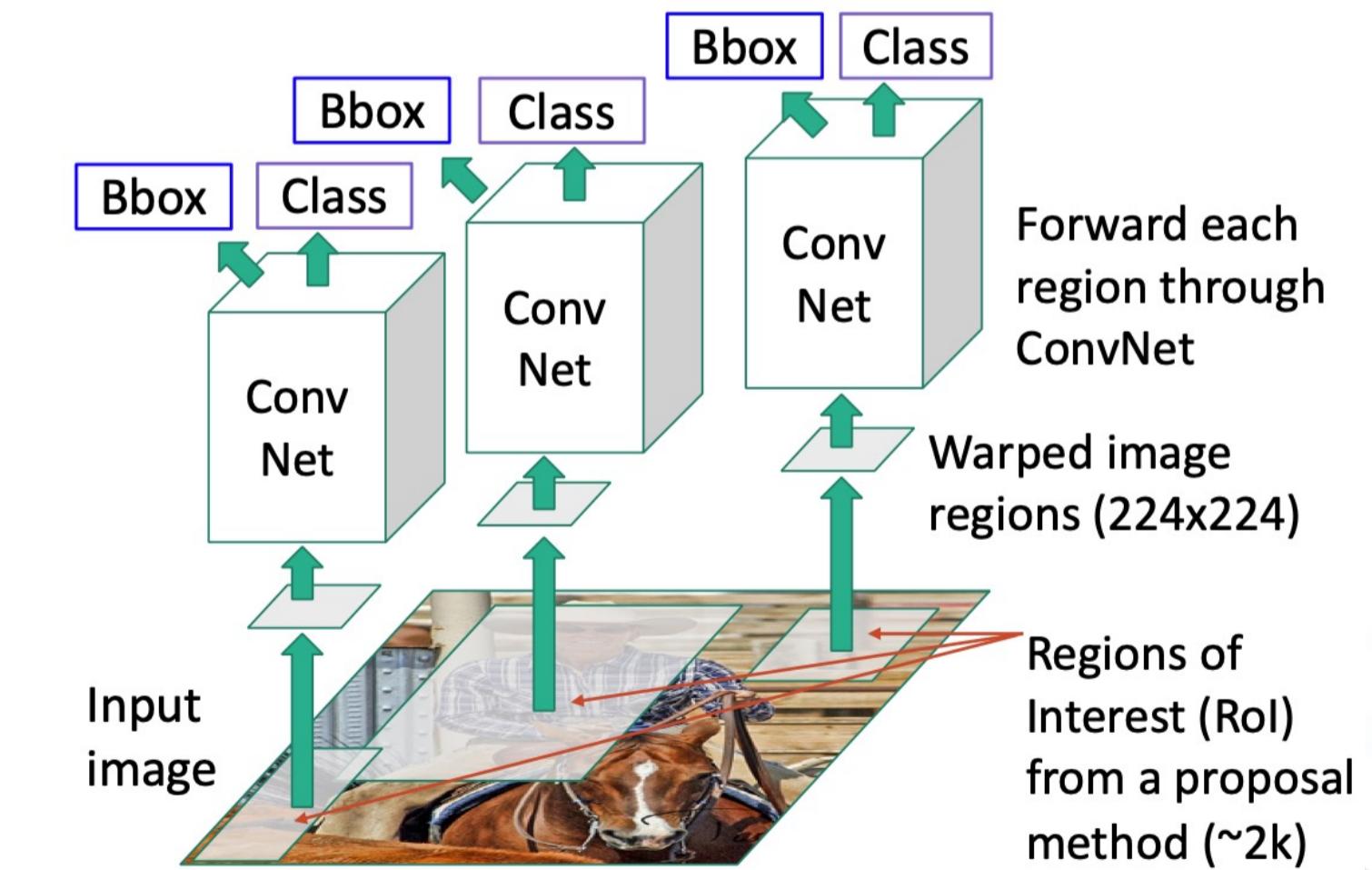
Idea: Overlapping proposals cause a lot of repeated work – same pixels processed many times. Can we avoid this?

Fast R-CNN

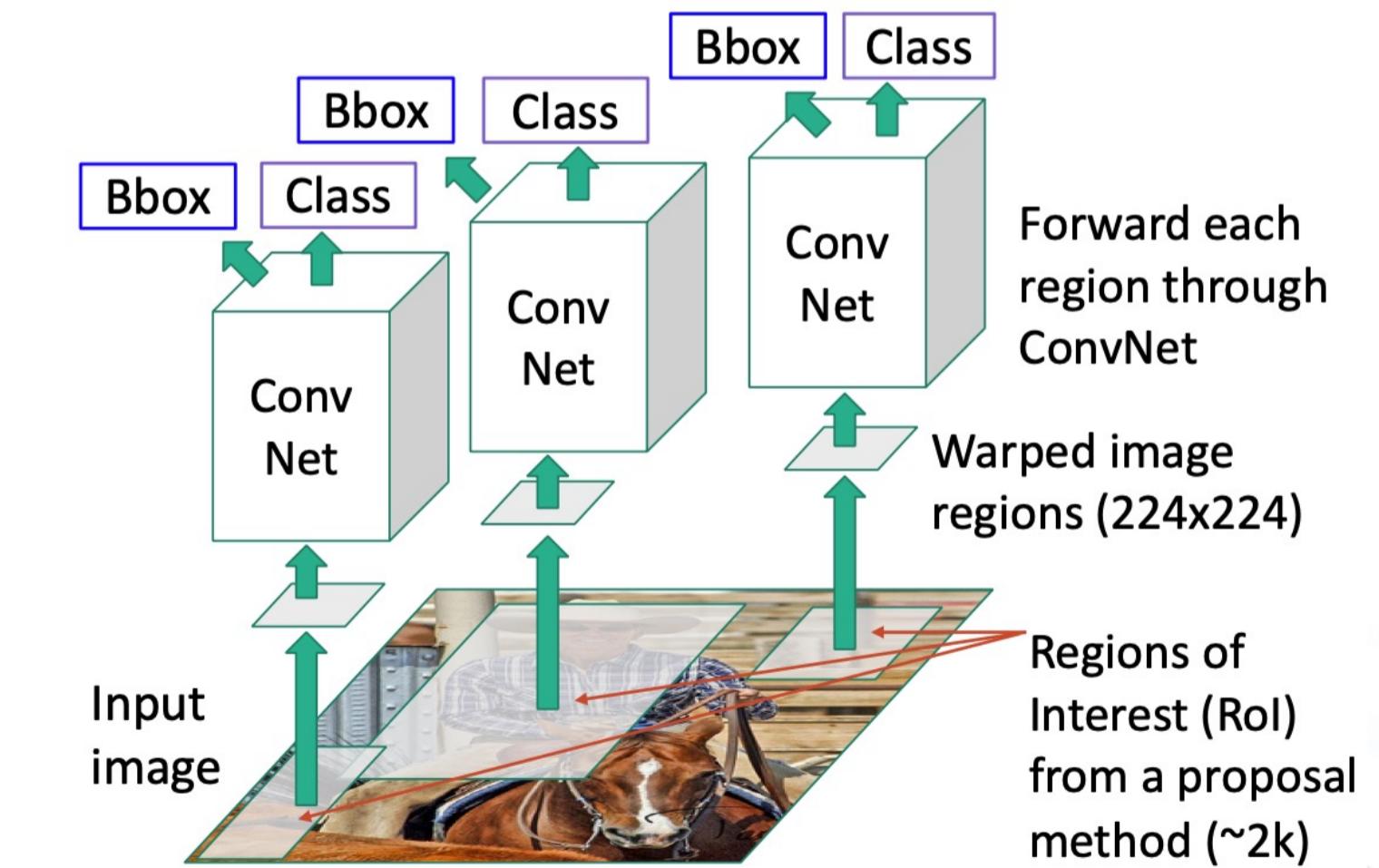
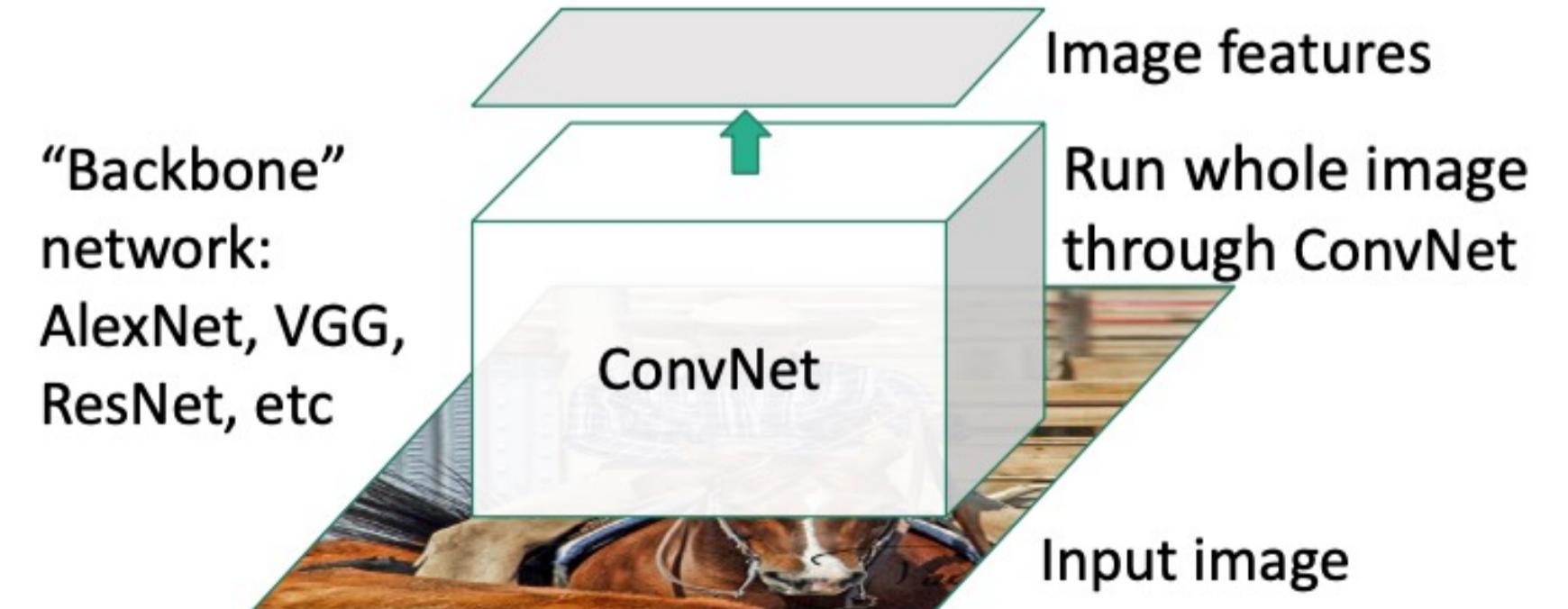


Input image

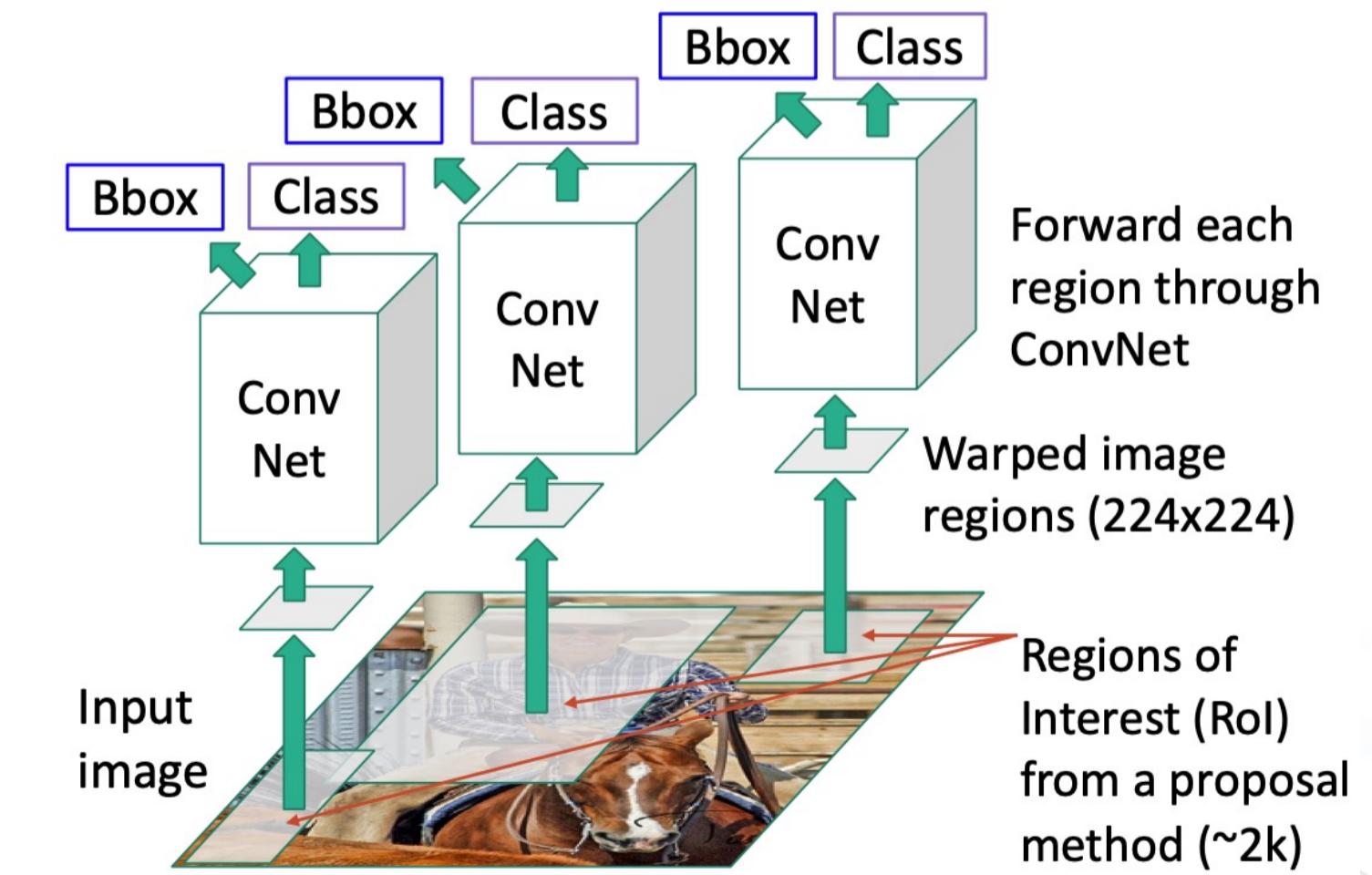
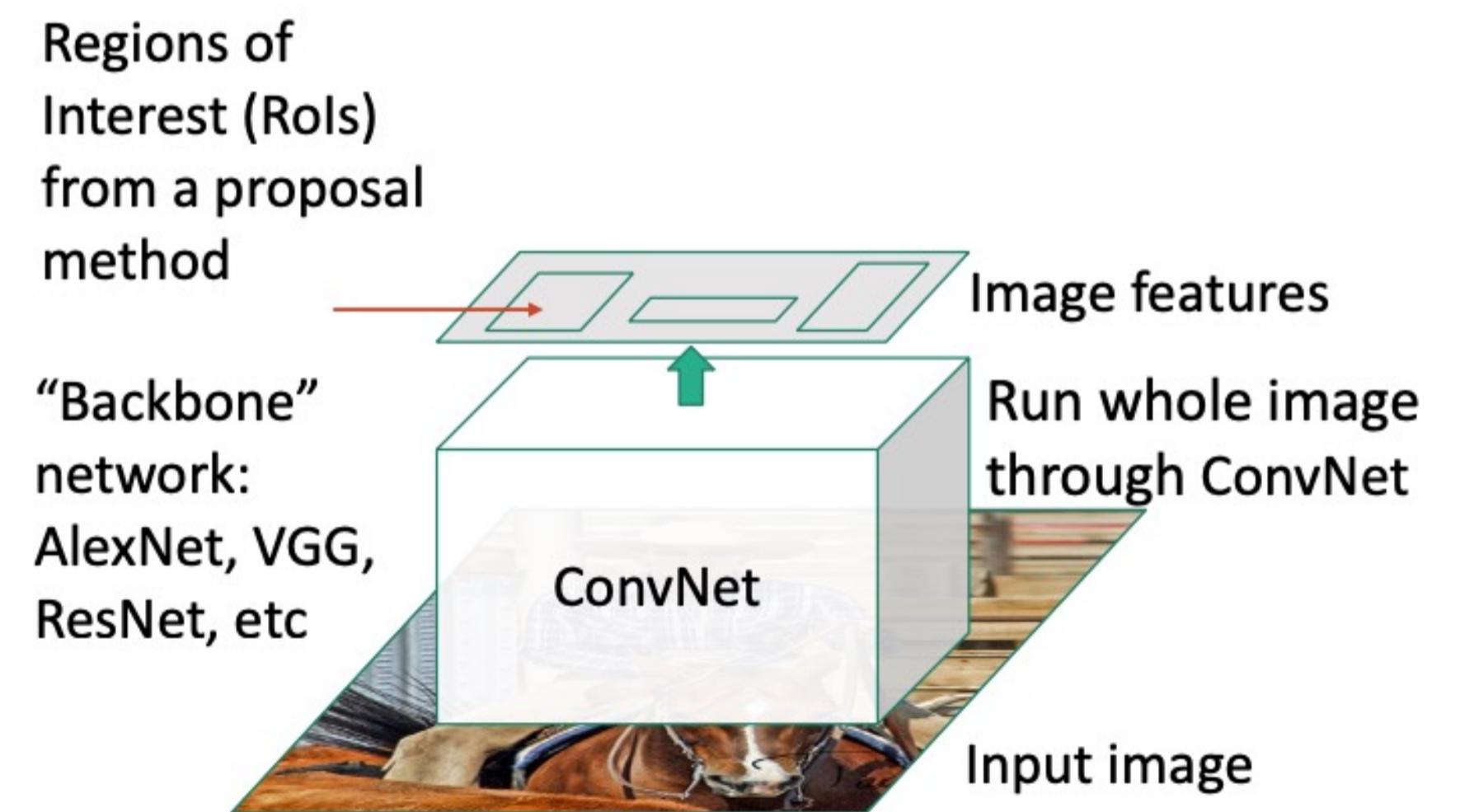
R. Girshick, "Fast R-CNN", ICCV 2015



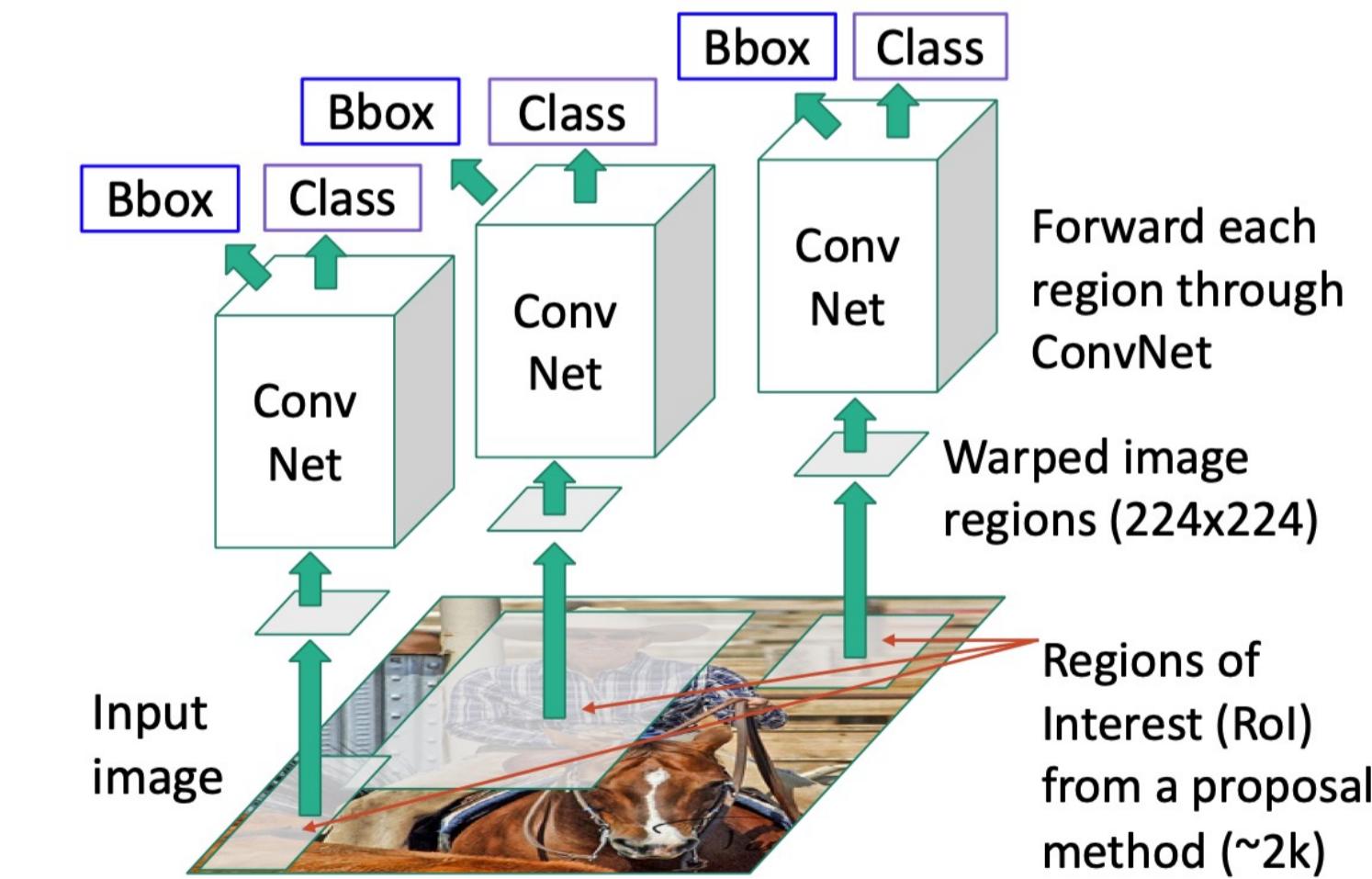
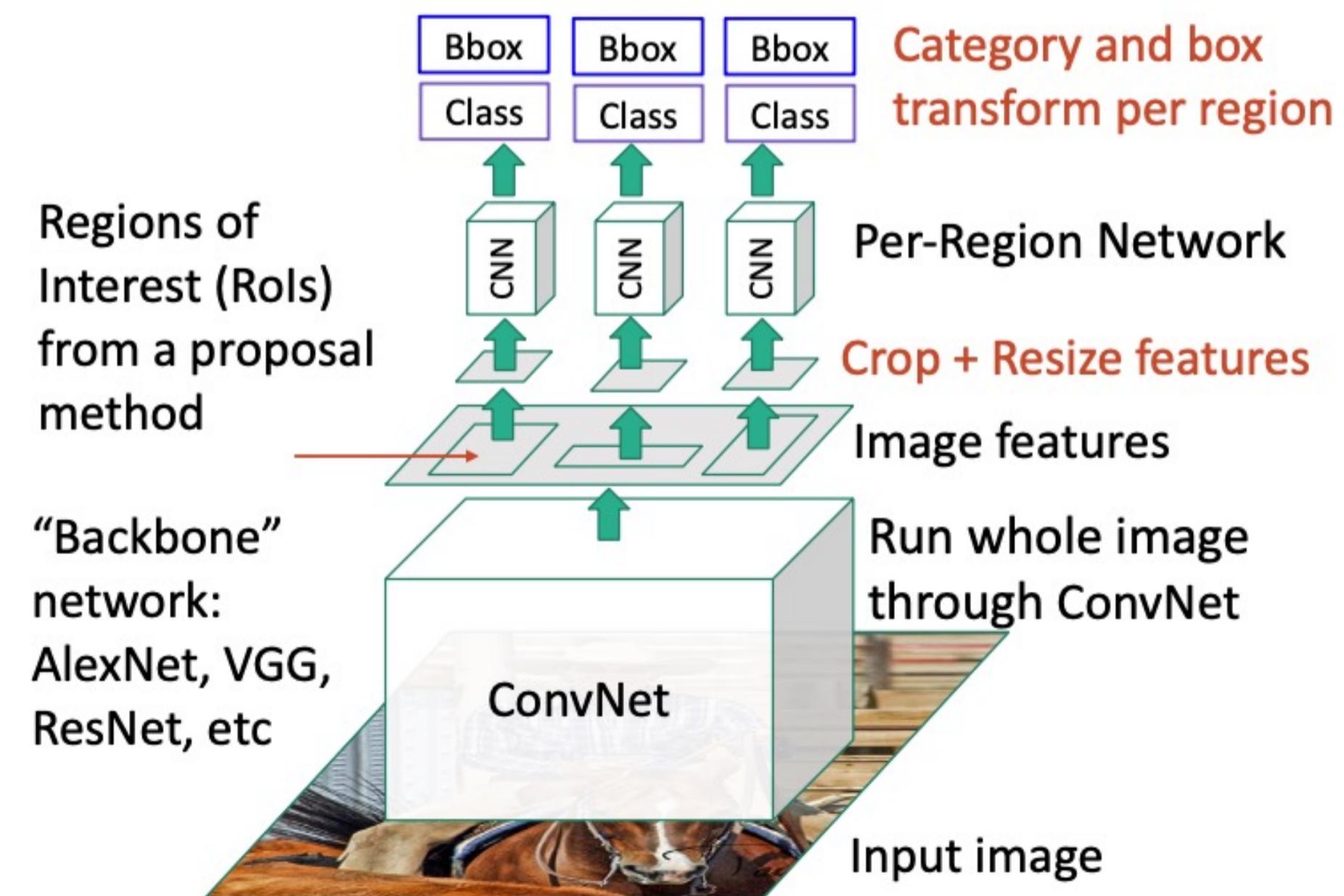
Fast R-CNN



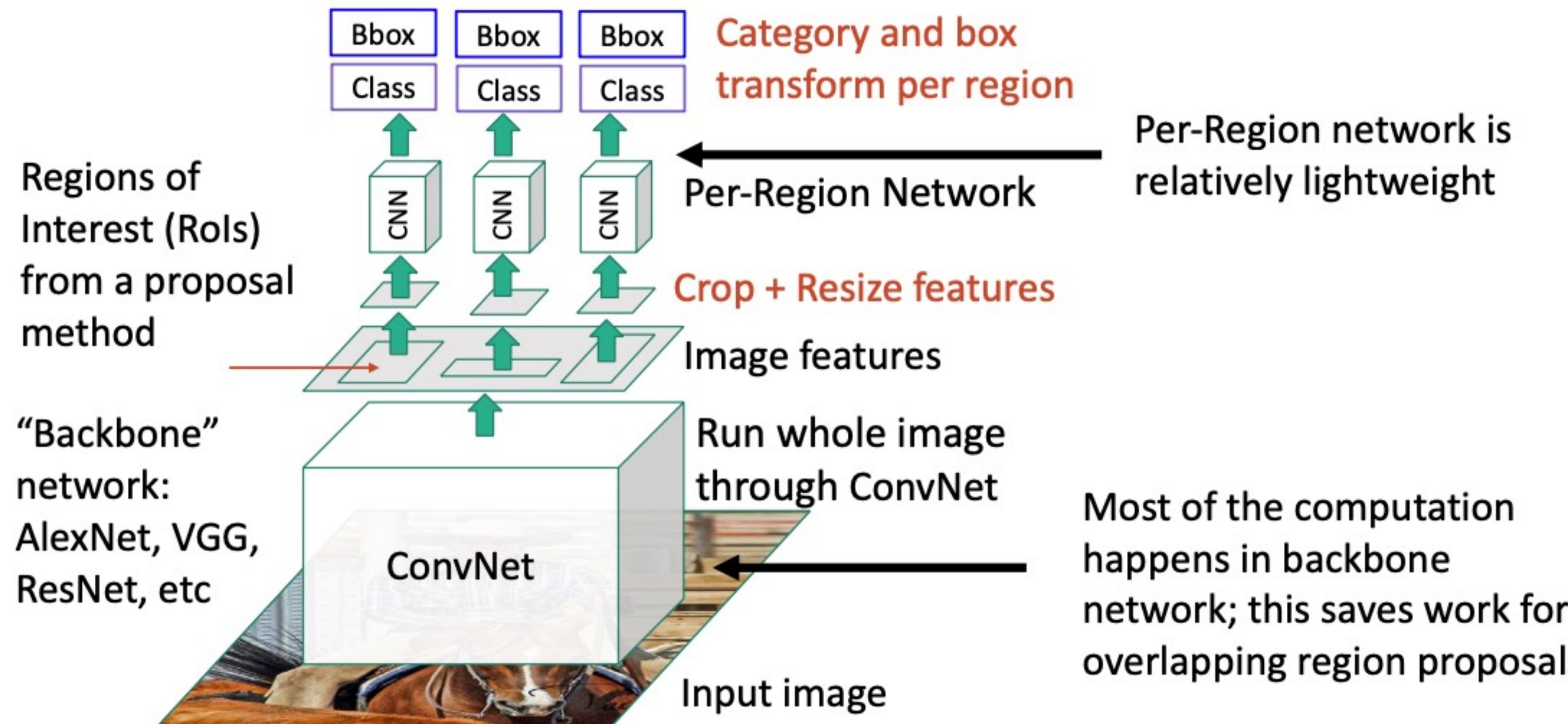
Fast R-CNN



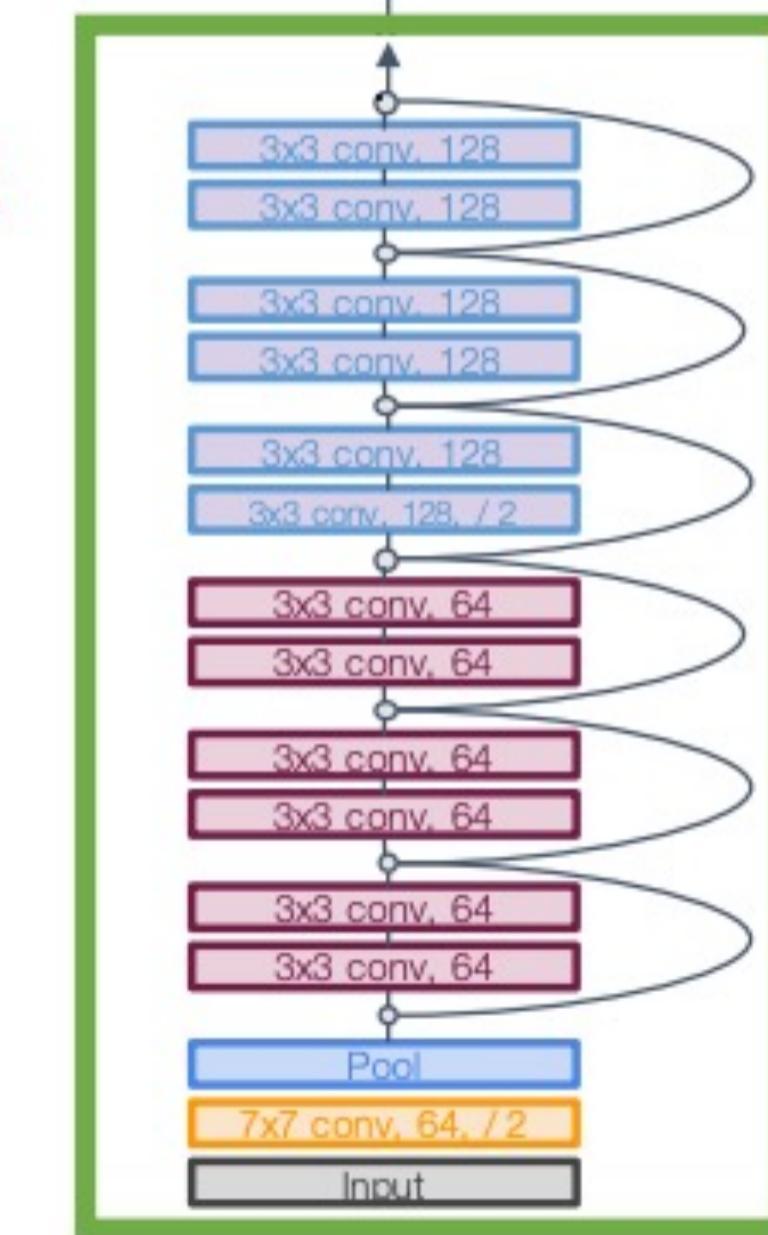
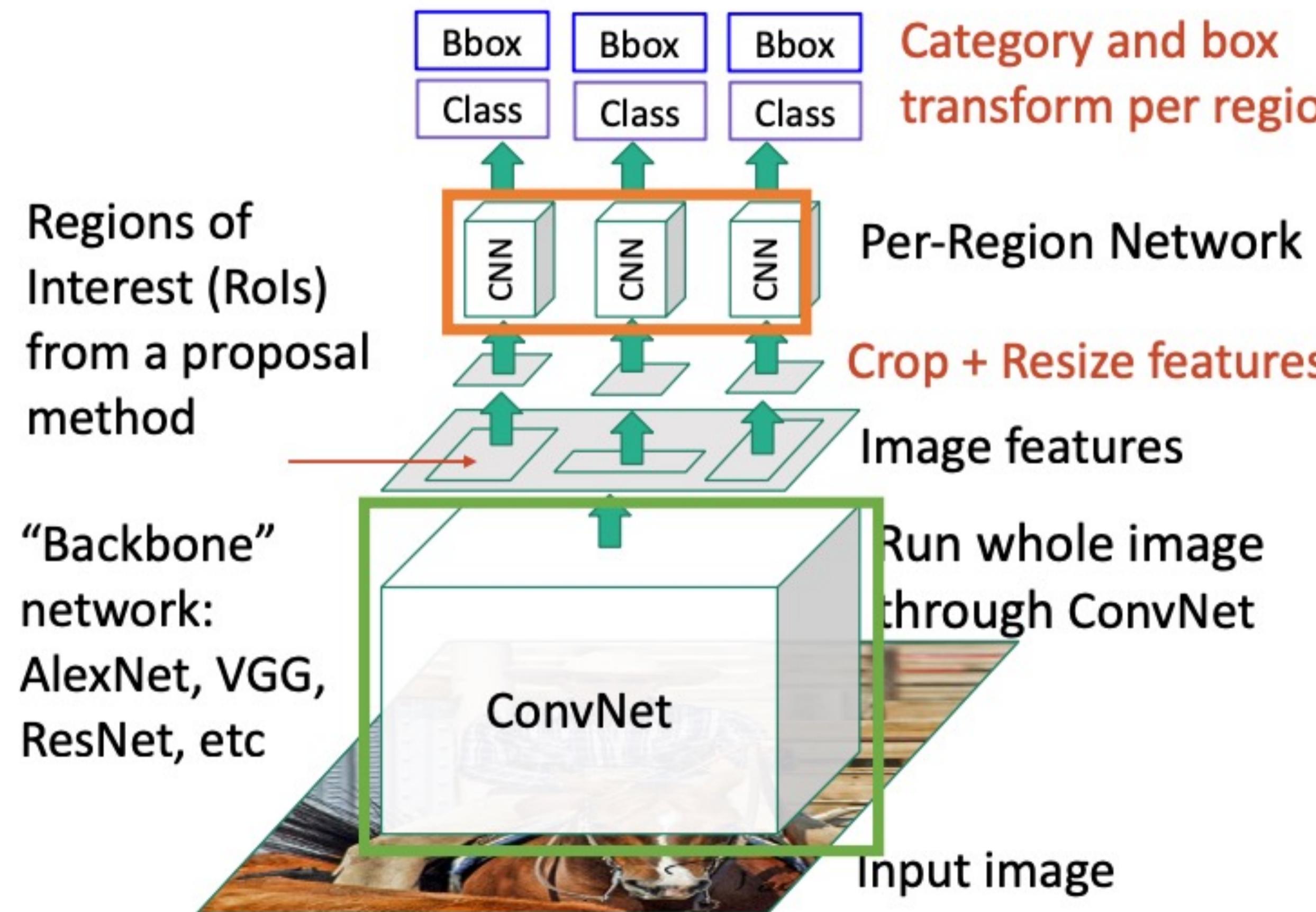
Fast R-CNN



Fast R-CNN

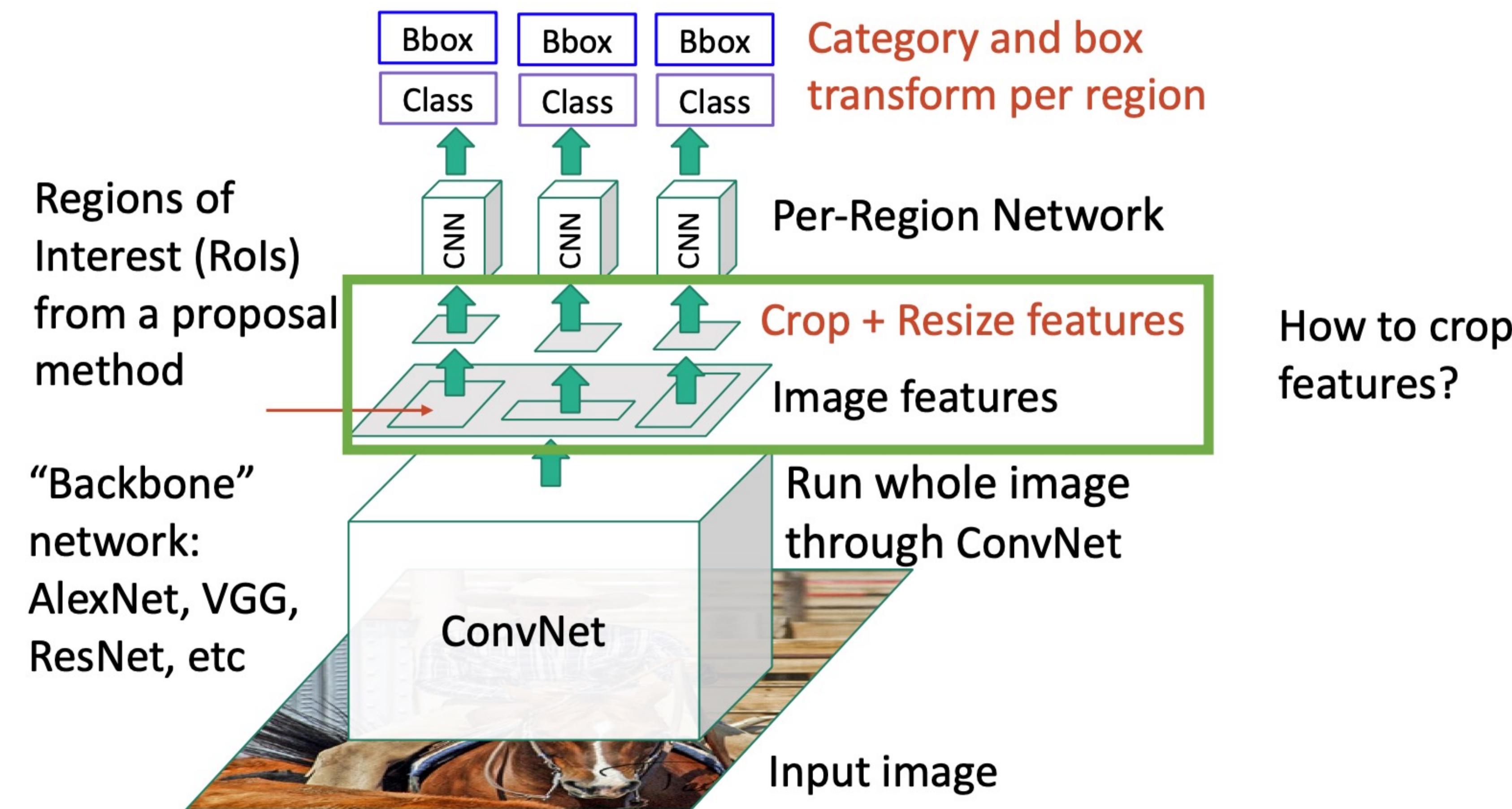


Fast R-CNN

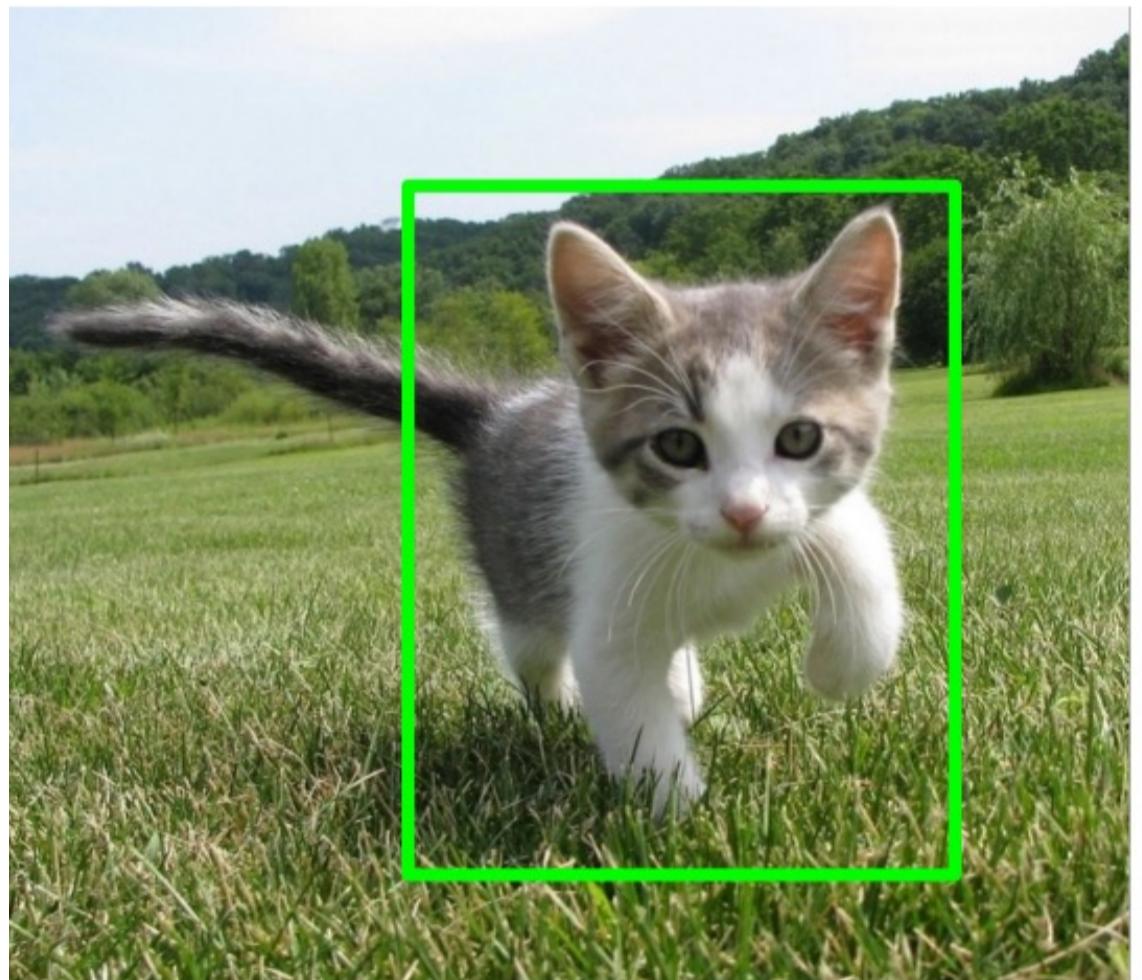


Example:
For ResNet, last stage is used as per-region network; the rest of the network is used as backbone

Fast R-CNN

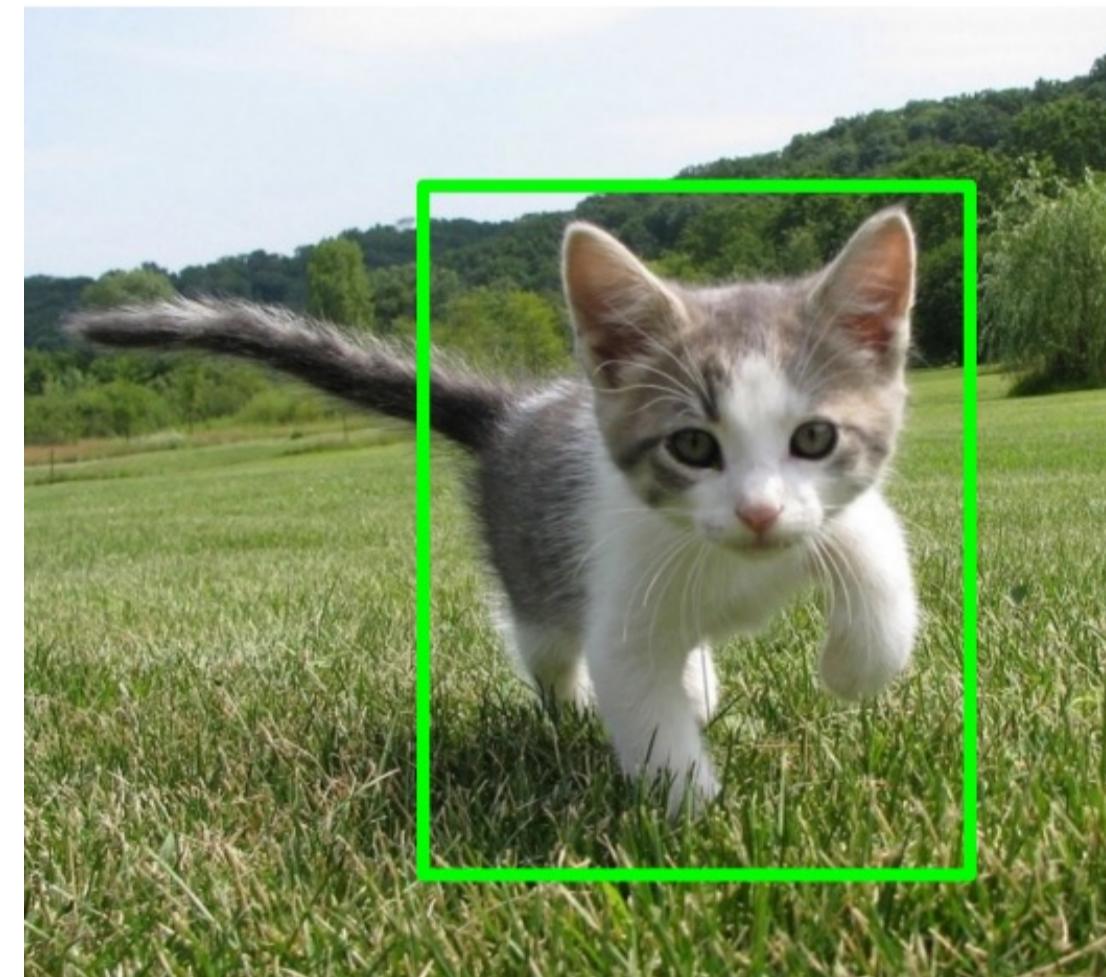


Cropping features: ROI Pool



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

Cropping features: ROI Pool



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

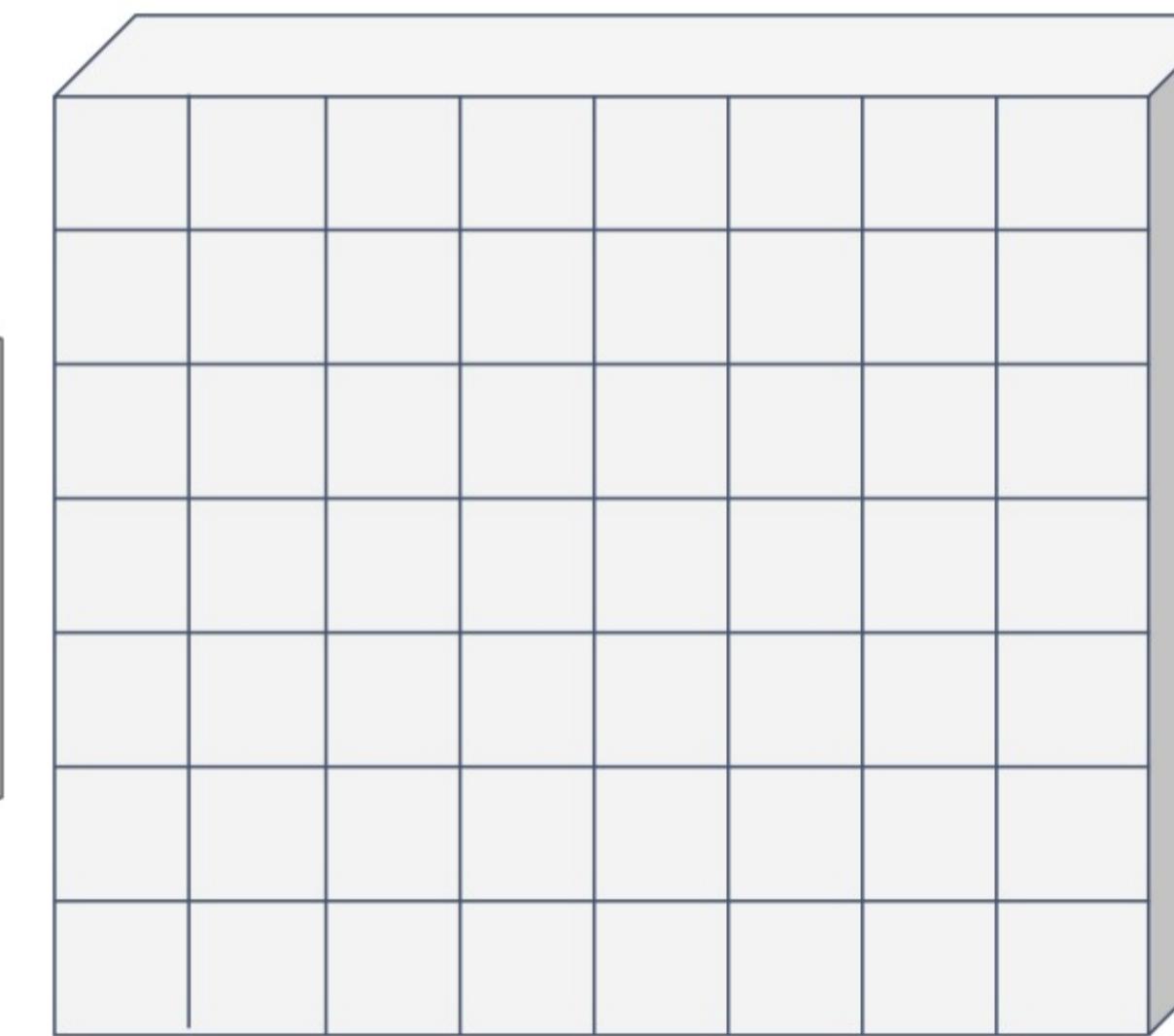
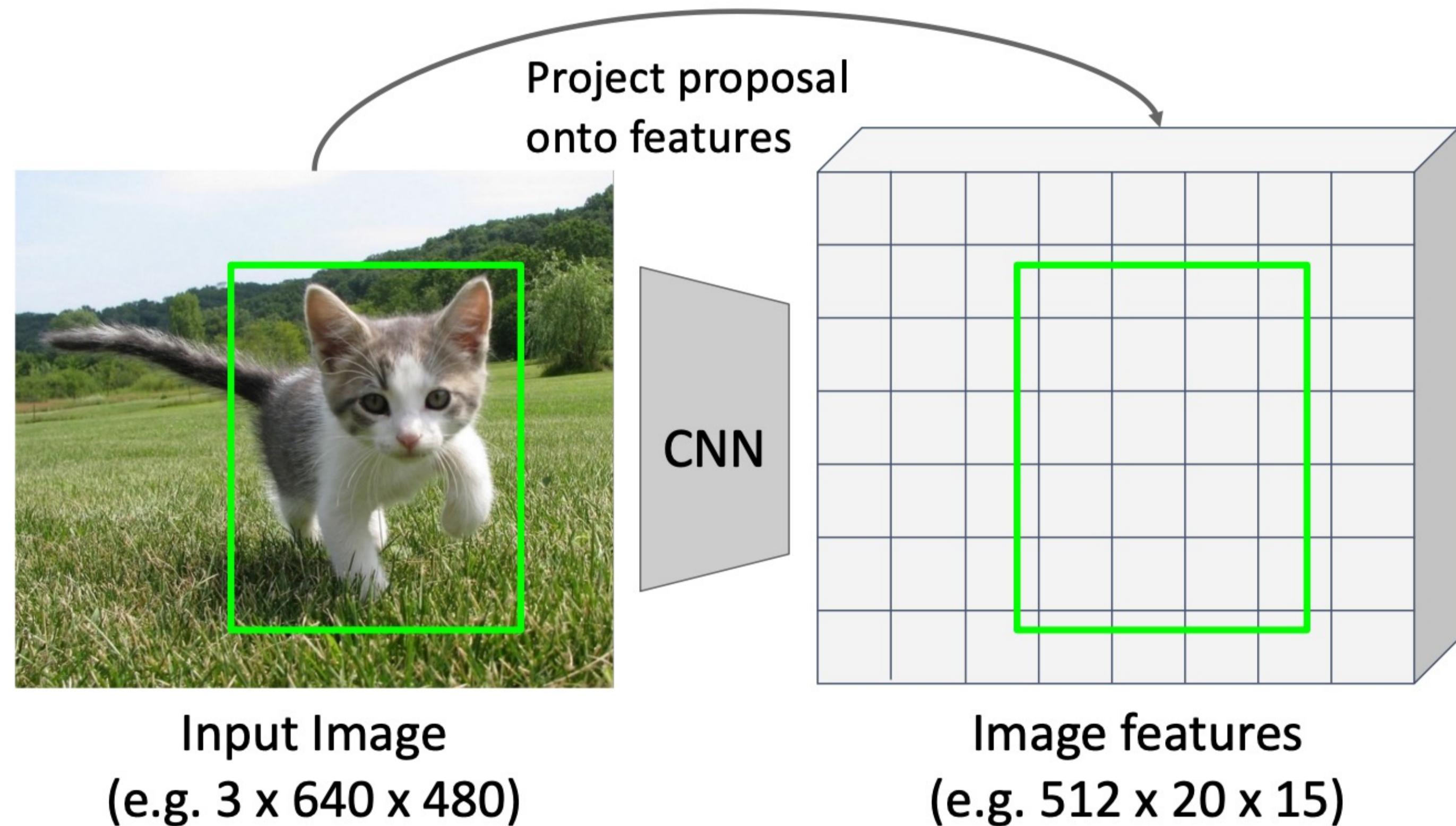


Image features
(e.g. $512 \times 20 \times 15$)

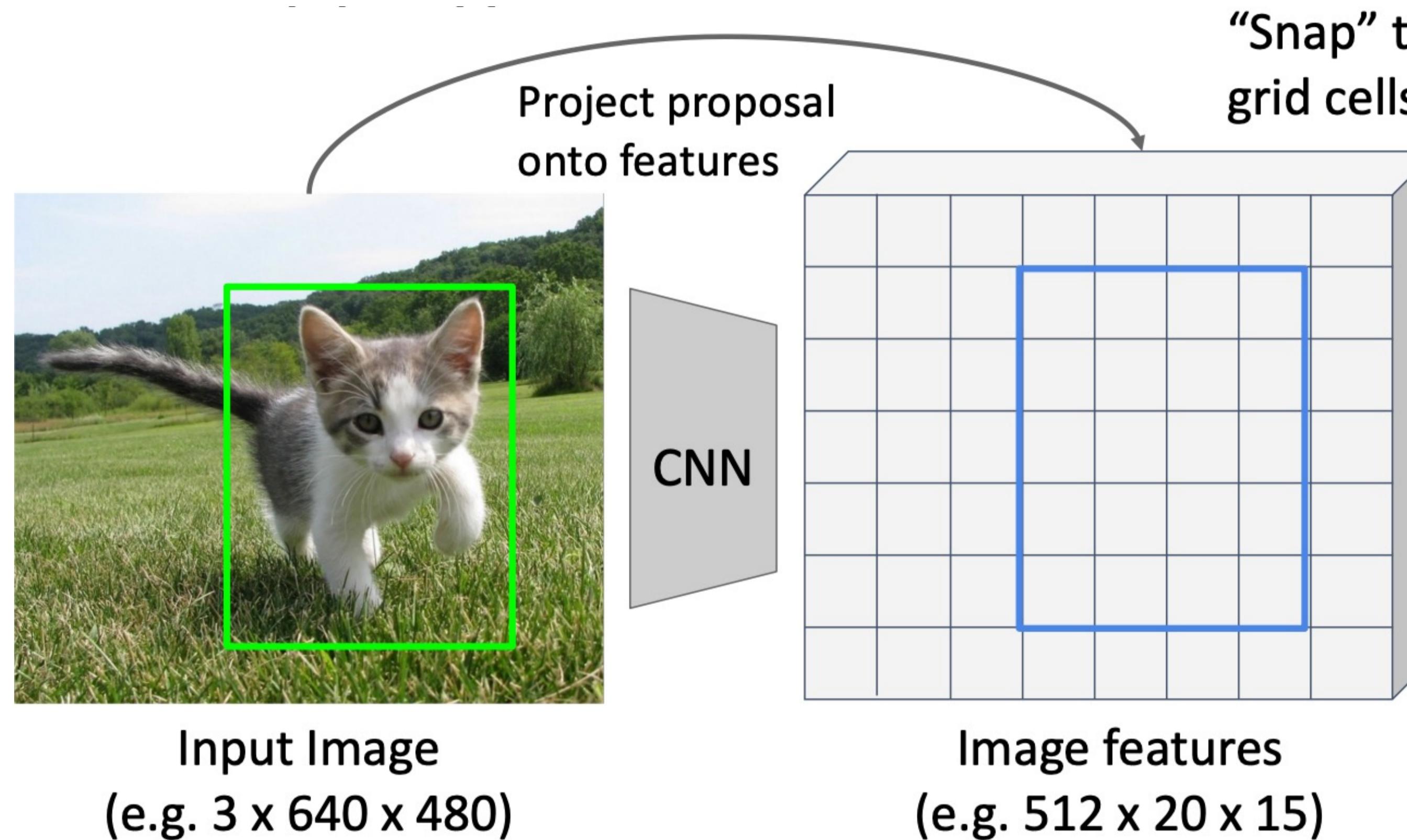
Want features for the
box of a fixed size
(2×2 in this example,
but 7×7 or 14×14 in
practice)

Cropping features: ROI Pool



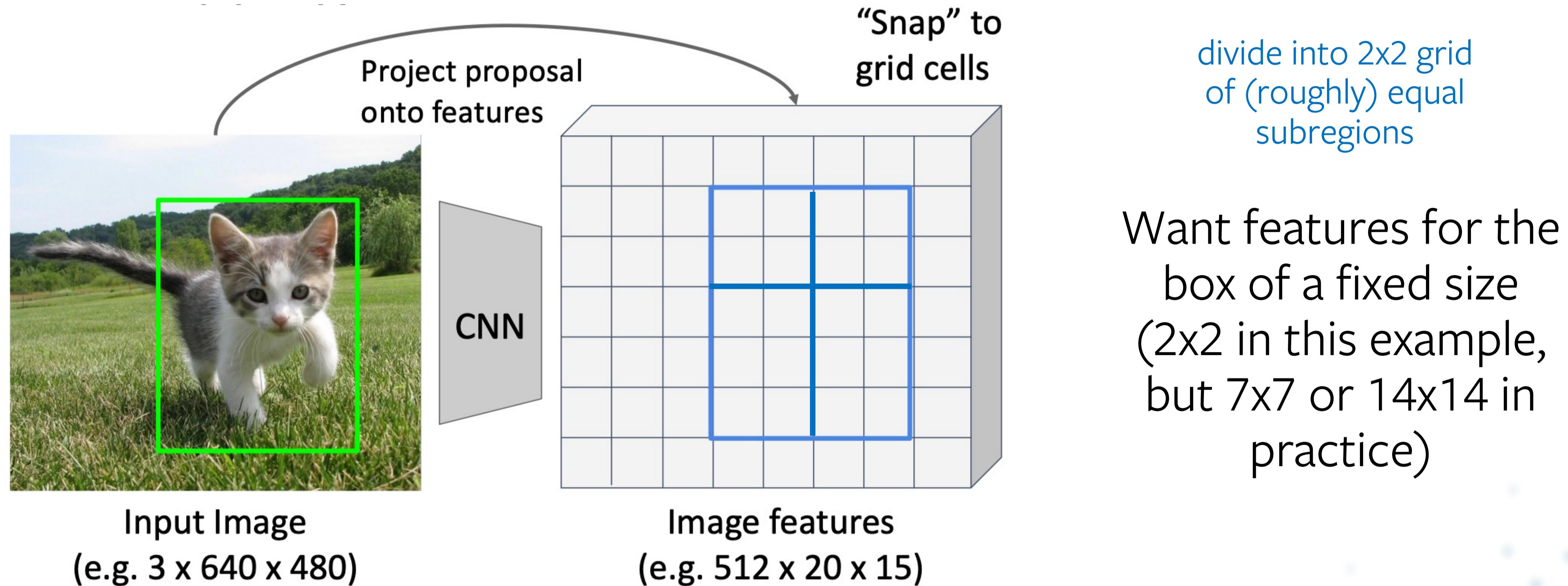
Want features for the
box of a fixed size
(2×2 in this example,
but 7×7 or 14×14 in
practice)

Cropping features: ROI Pool

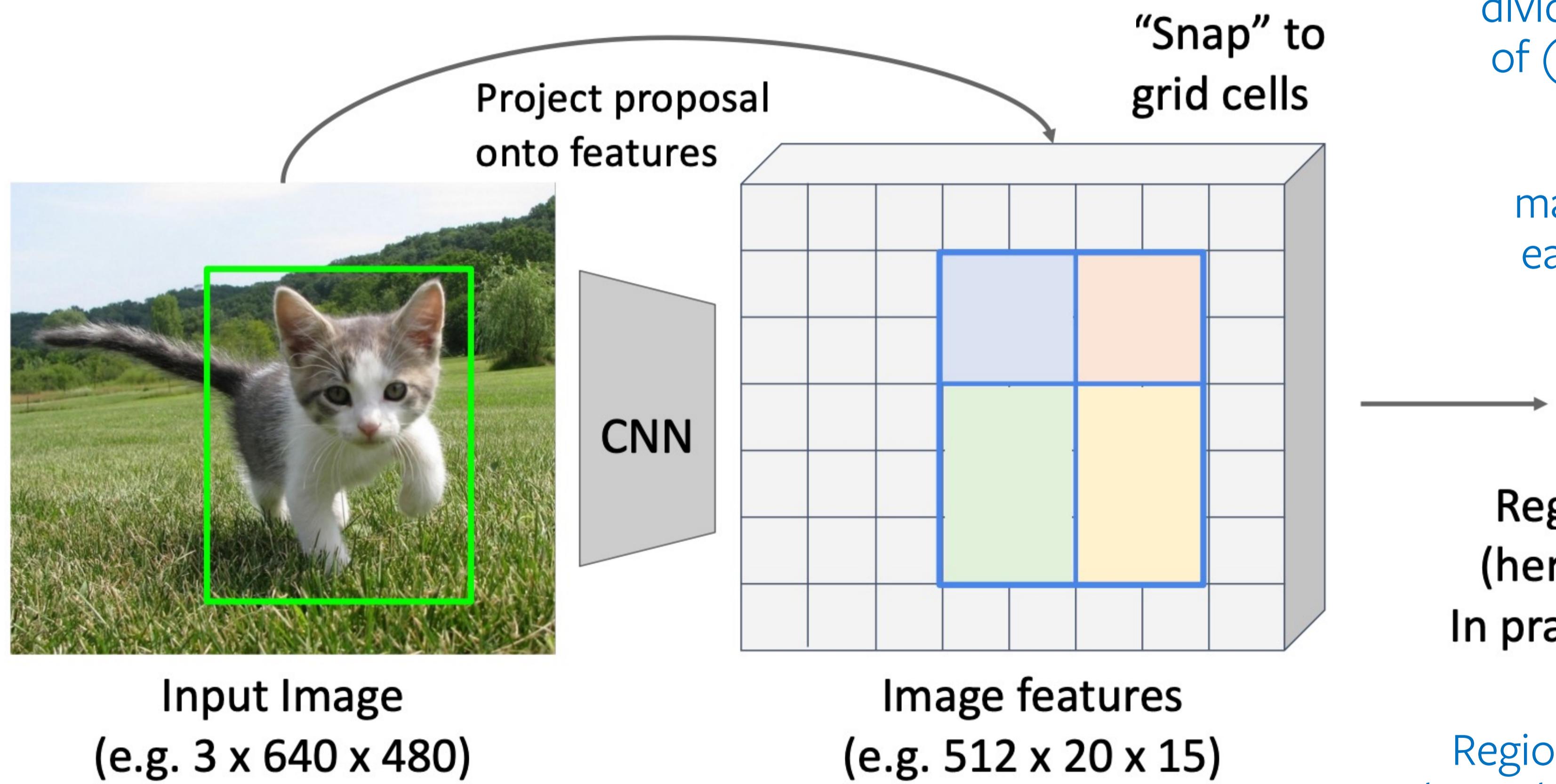


Want features for the box of a fixed size (2x2 in this example, but 7x7 or 14x14 in practice)

Cropping features: ROI Pool

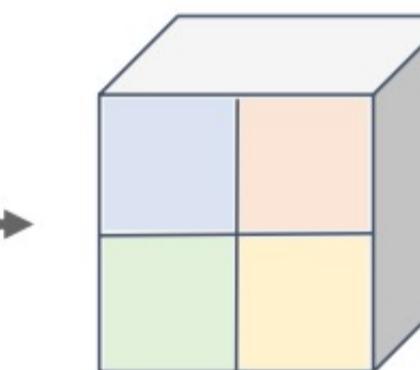


Cropping features: ROI Pool



divide into 2×2 grid of (roughly) equal subregions

max-pool within each subregion

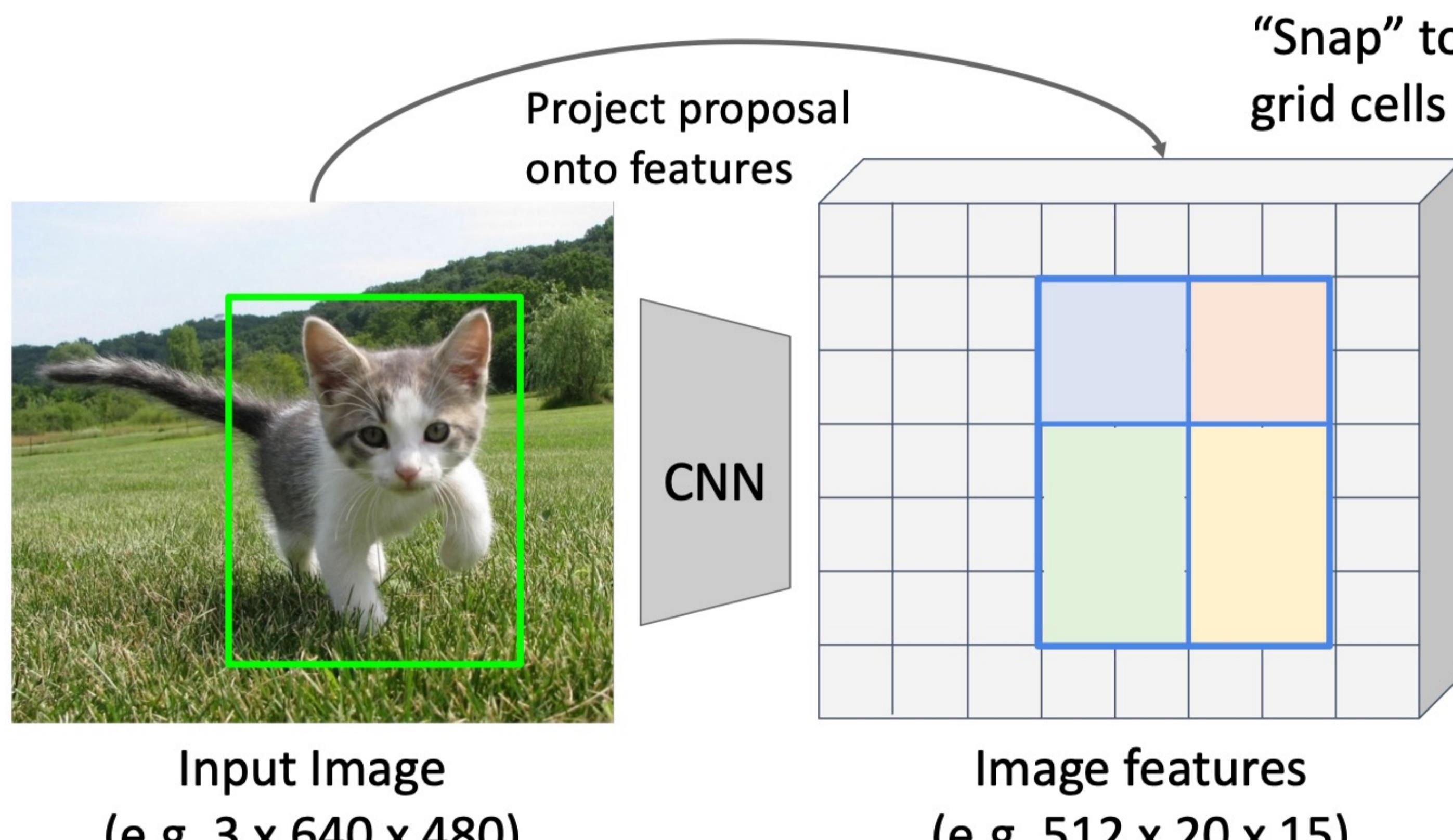


Region features
(here $512 \times 2 \times 2$;
In practice $512 \times 7 \times 7$)

Region features always have the same size even if input regions of different sizes



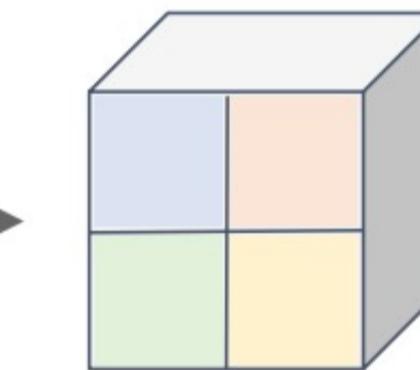
Cropping features: ROI Pool



Problem: lots of approximation and misalignment

divide into 2×2 grid of (roughly) equal subregions

max-pool within each subregion

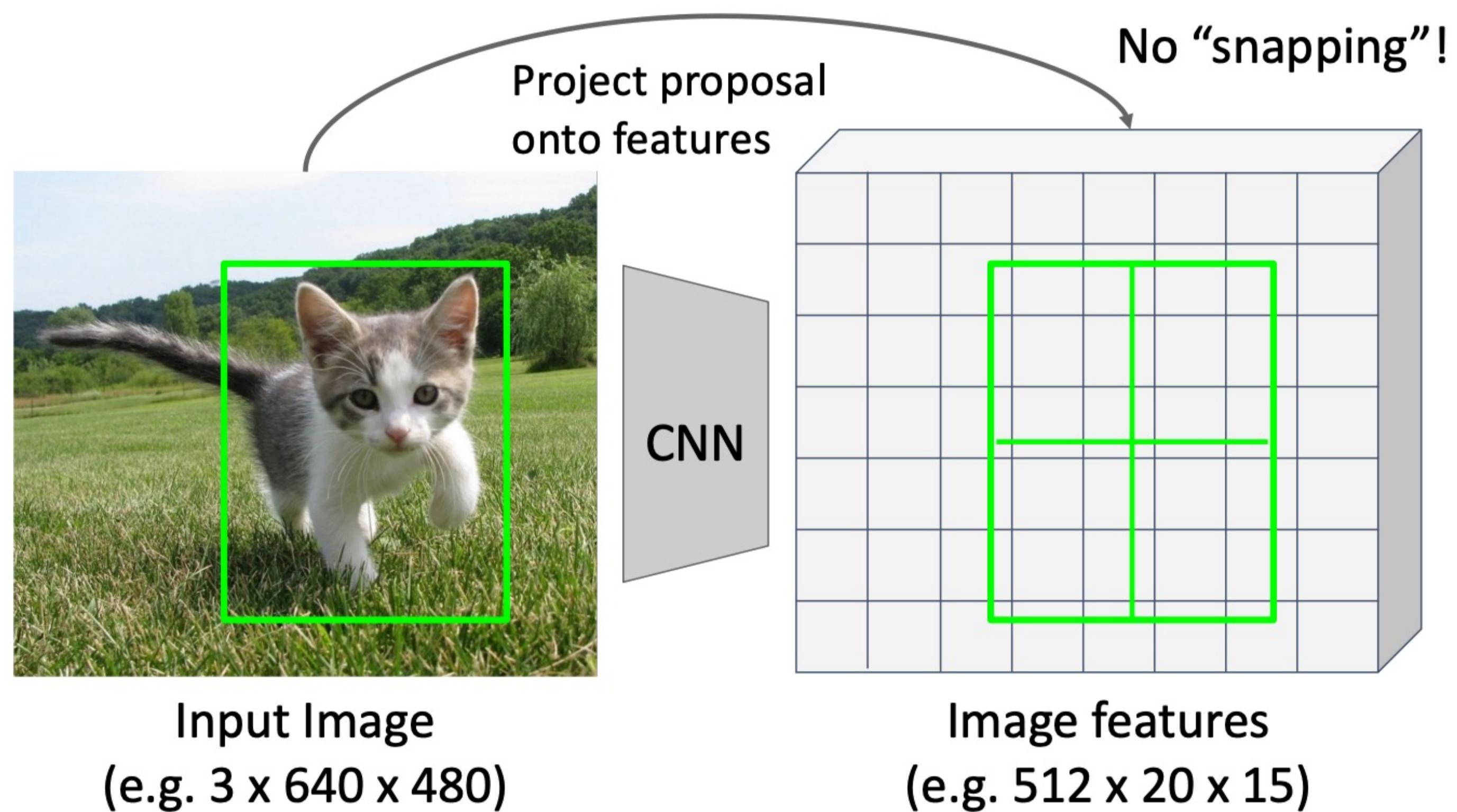


Region features
(here $512 \times 2 \times 2$;
In practice $512 \times 7 \times 7$)

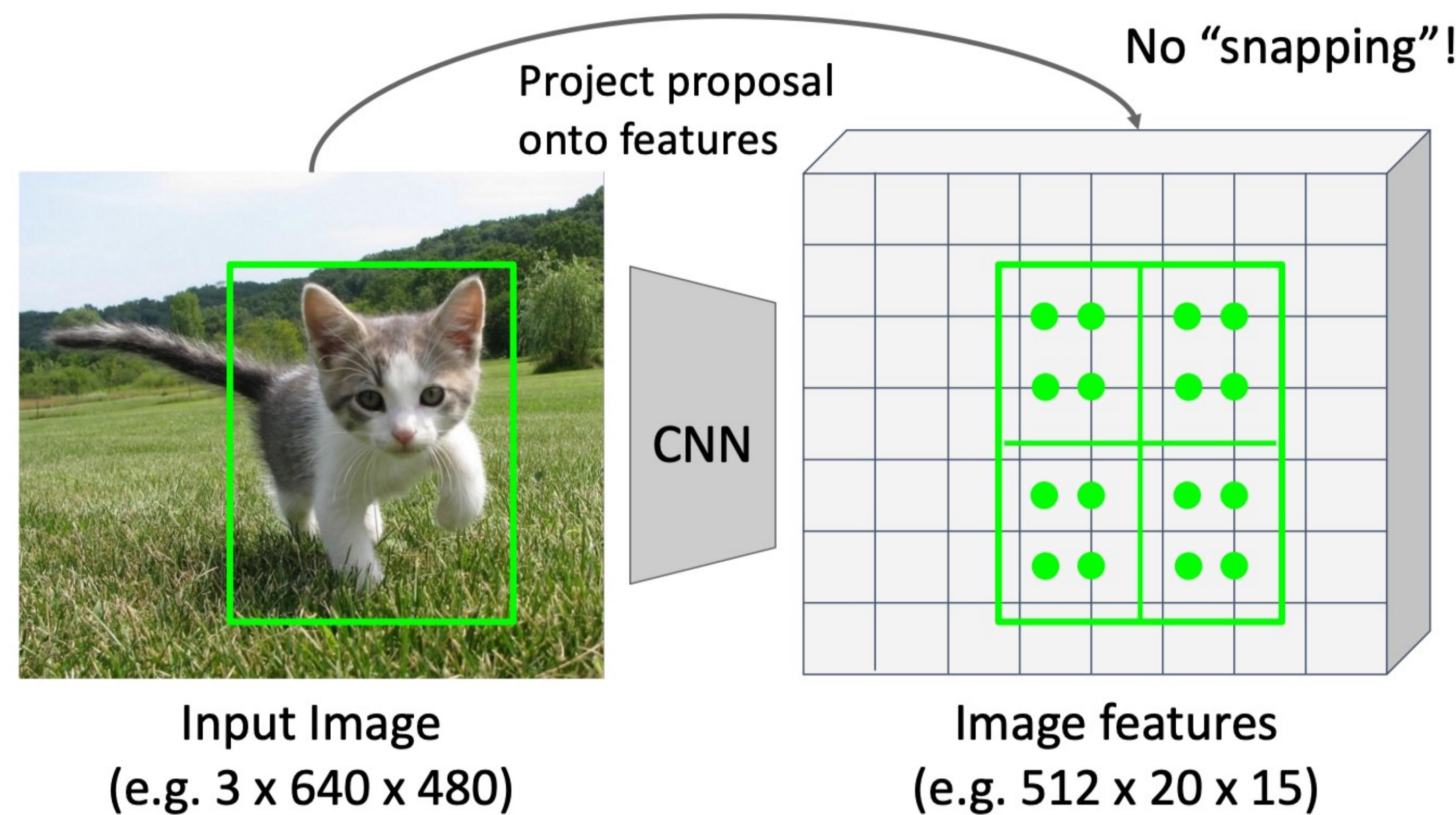
Region features always have the same size even if input regions of different sizes



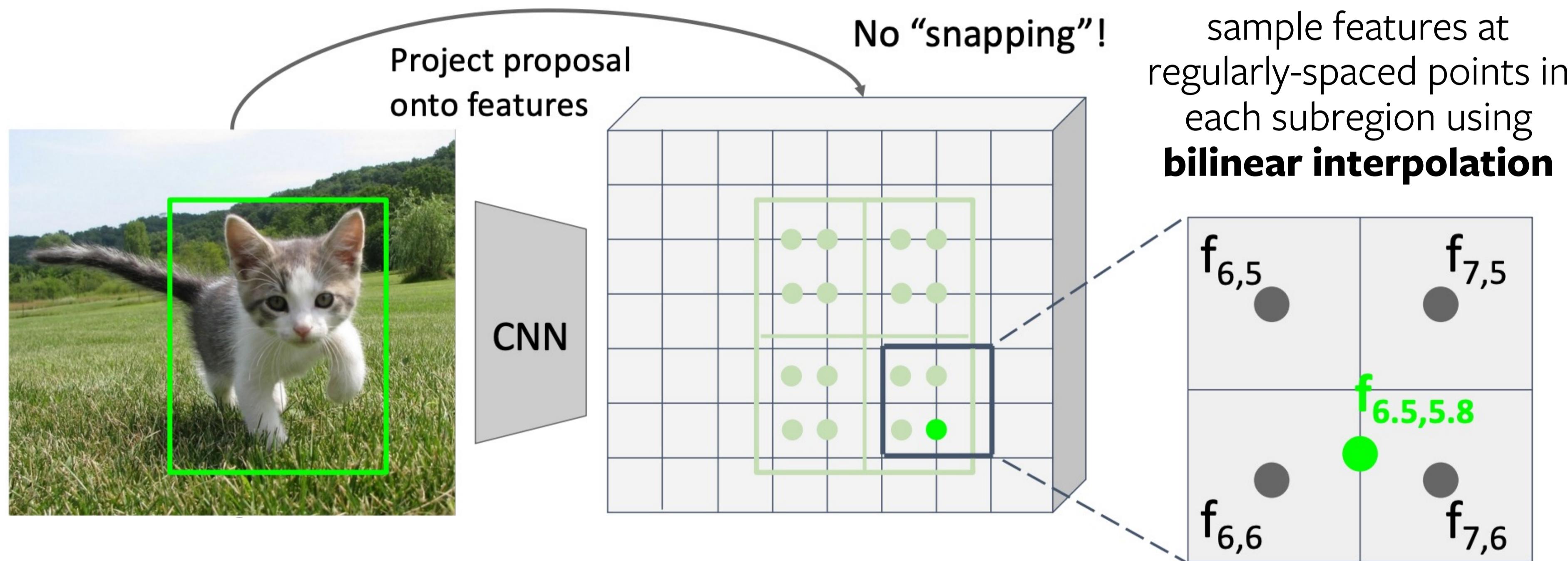
Cropping features: ROI Align



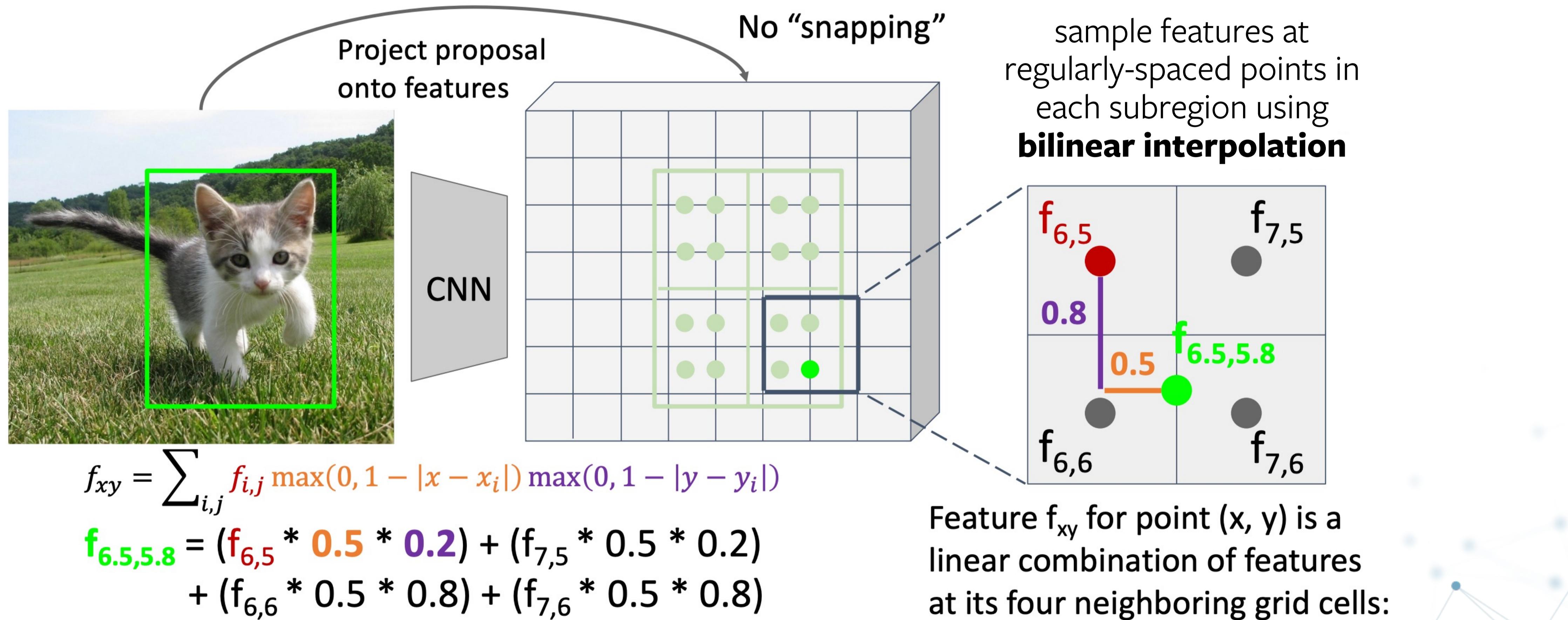
Cropping features: ROI Align



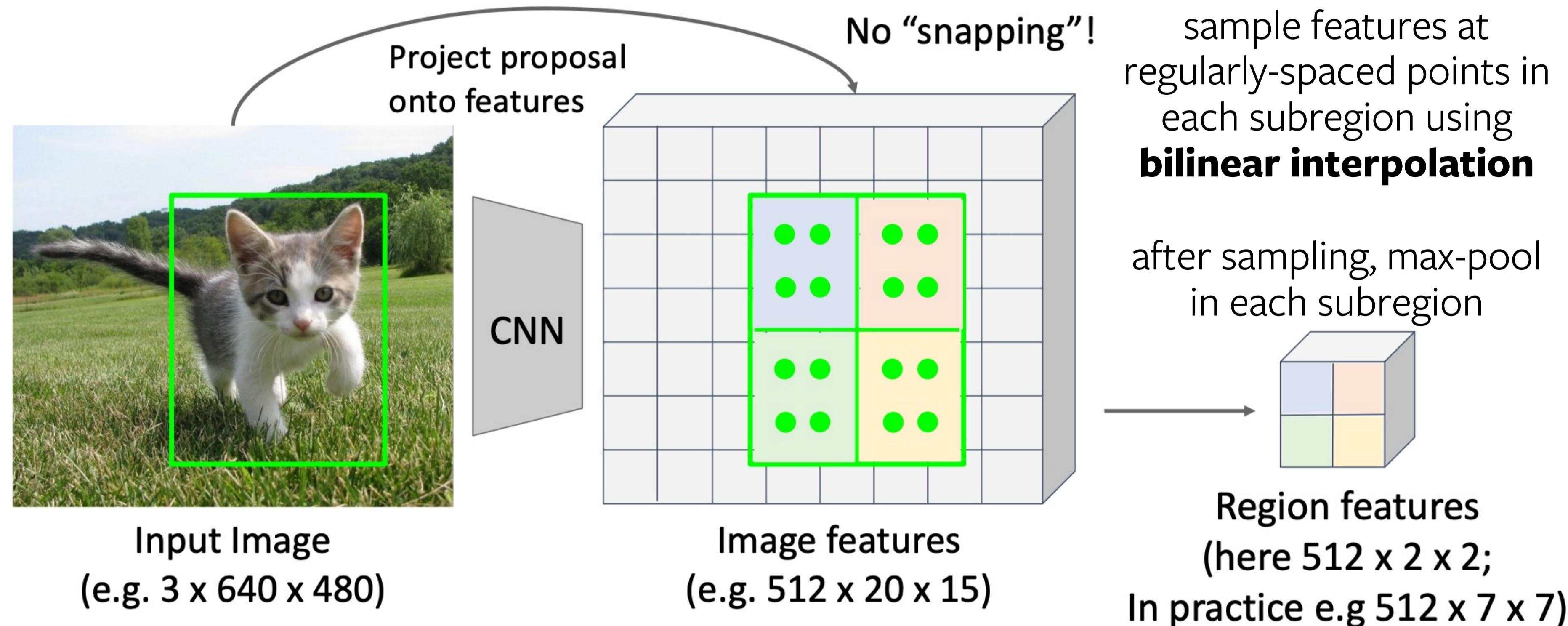
Cropping features: ROI Align



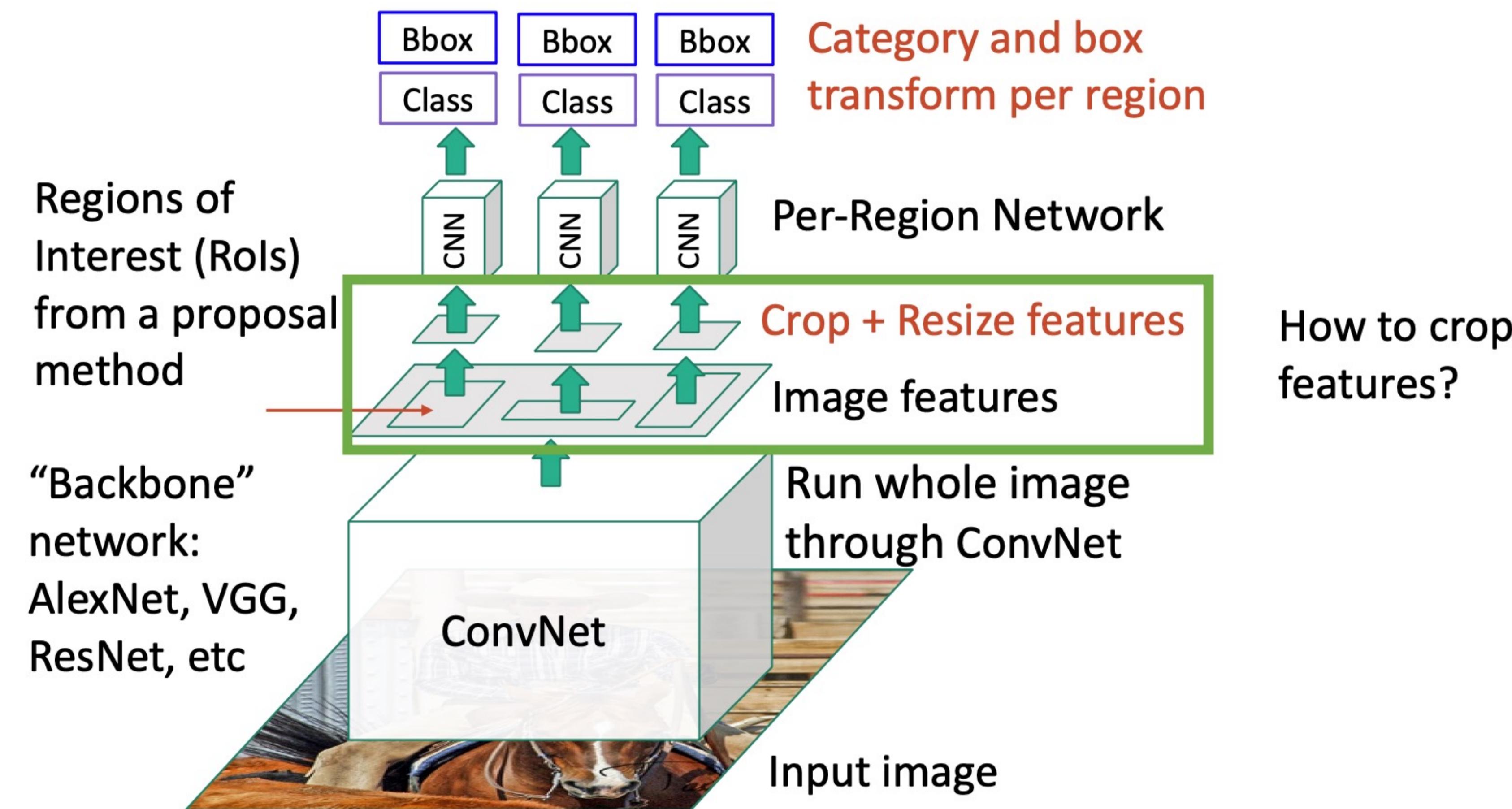
Cropping features: ROI Align



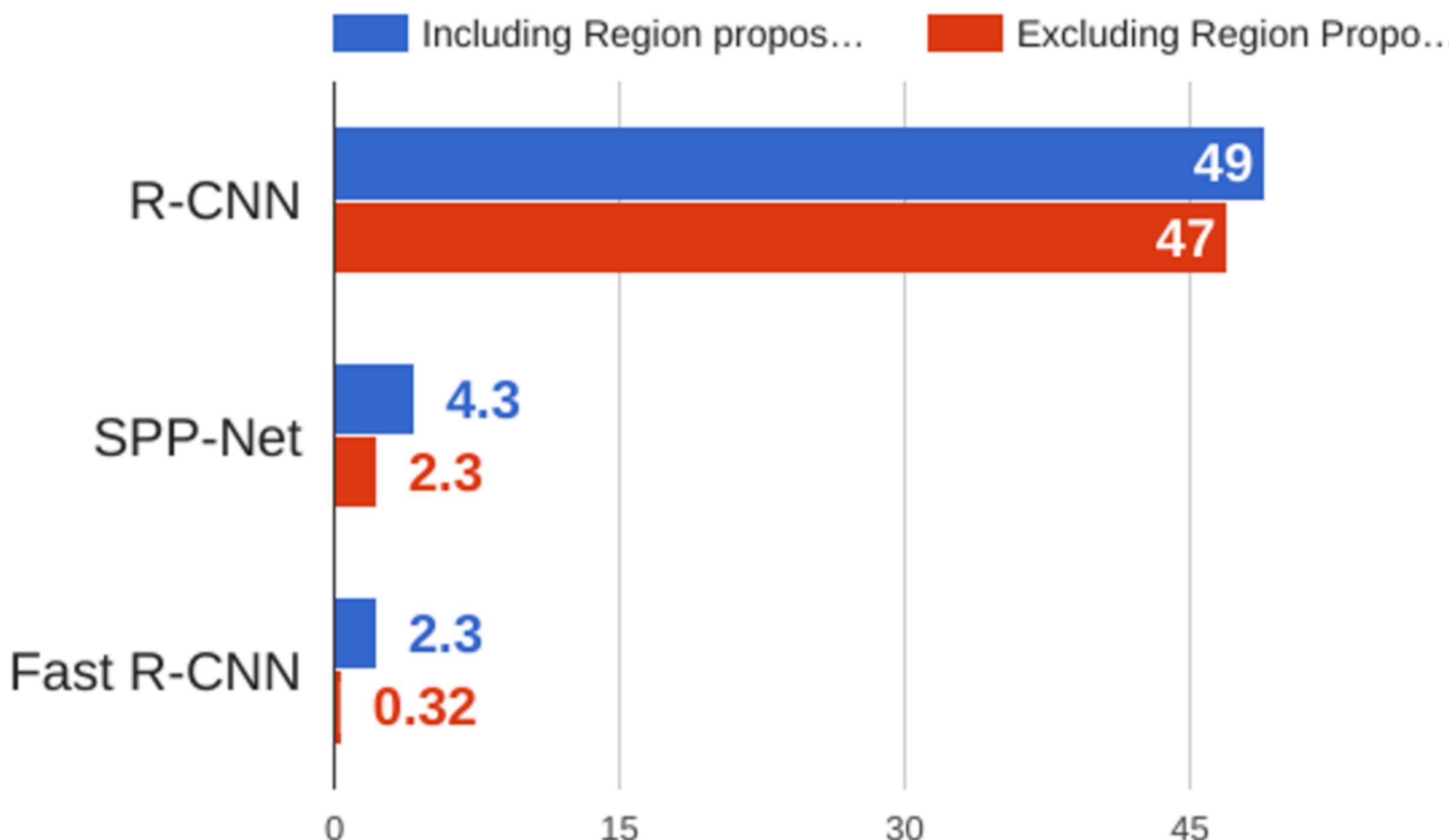
Cropping features: ROI Align



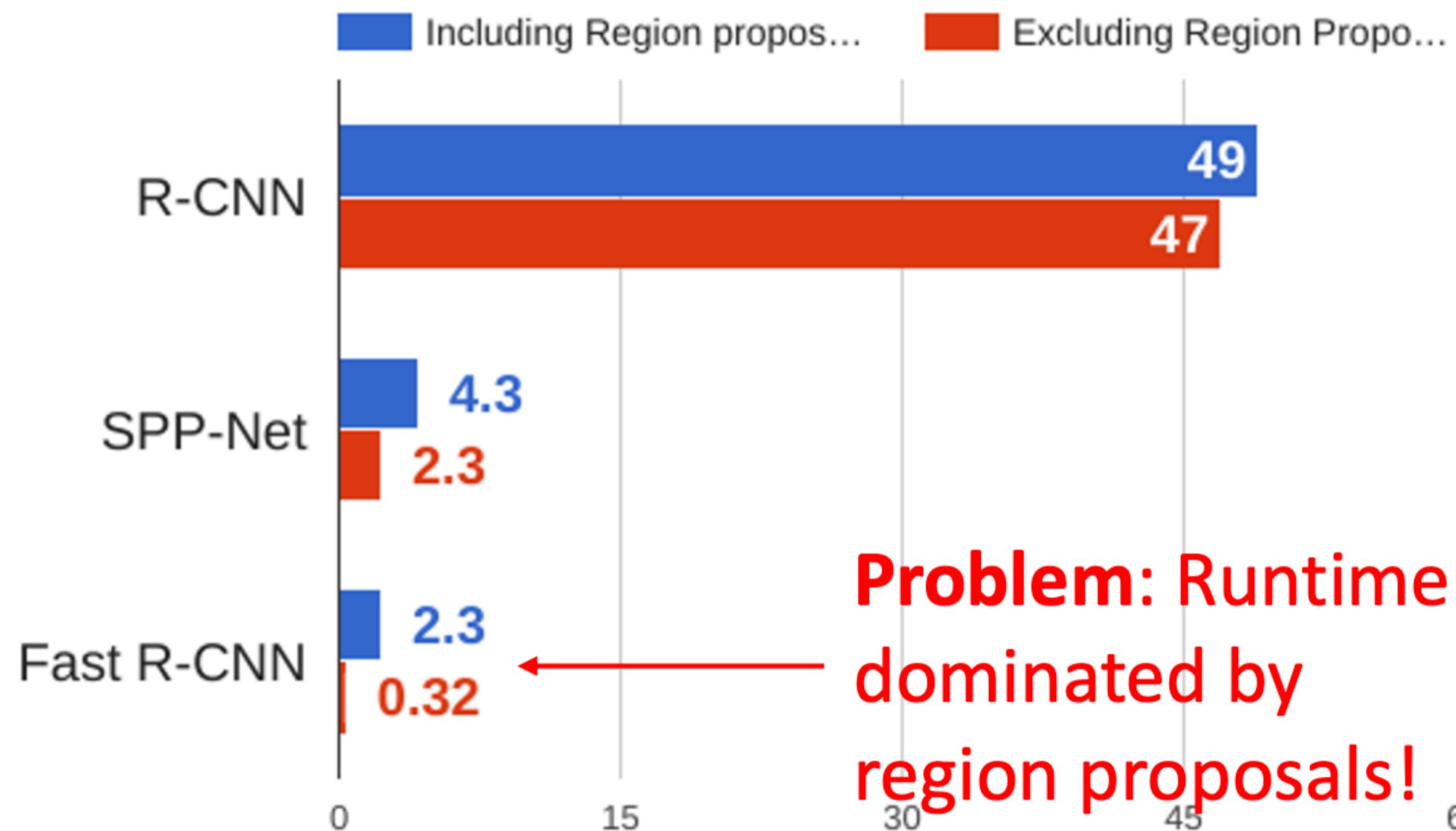
Fast R-CNN



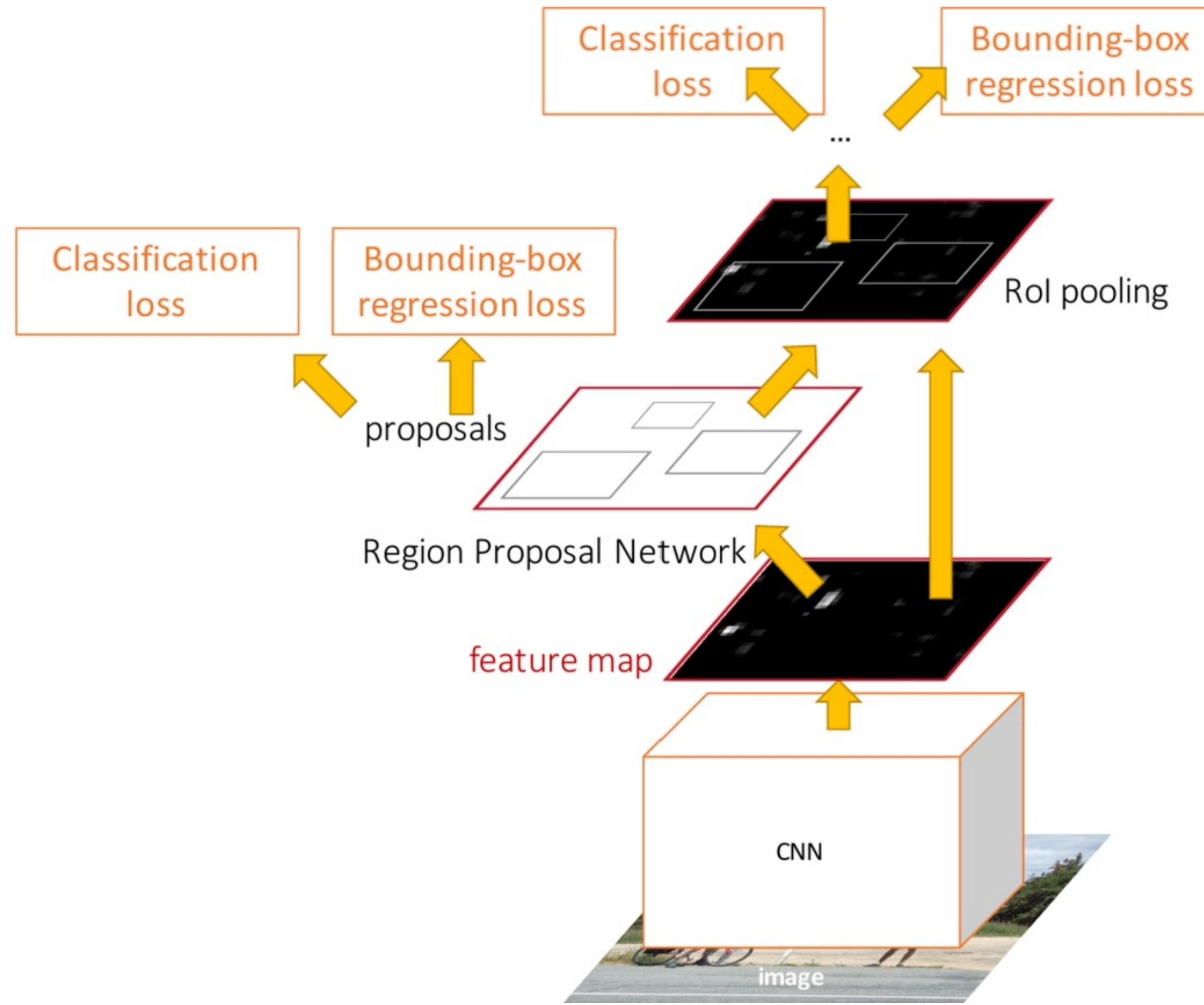
Test time (seconds)



Test time (seconds)



Faster R-CNN: Learnable Region Proposals



Idea: Calculating region proposals separately is slow – let's learn them jointly instead!

Implementation: insert Region Proposal Network (RPN) to predict proposals from features

Region proposal network (RPN)

Run backbone CNN to get
features aligned to input image



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

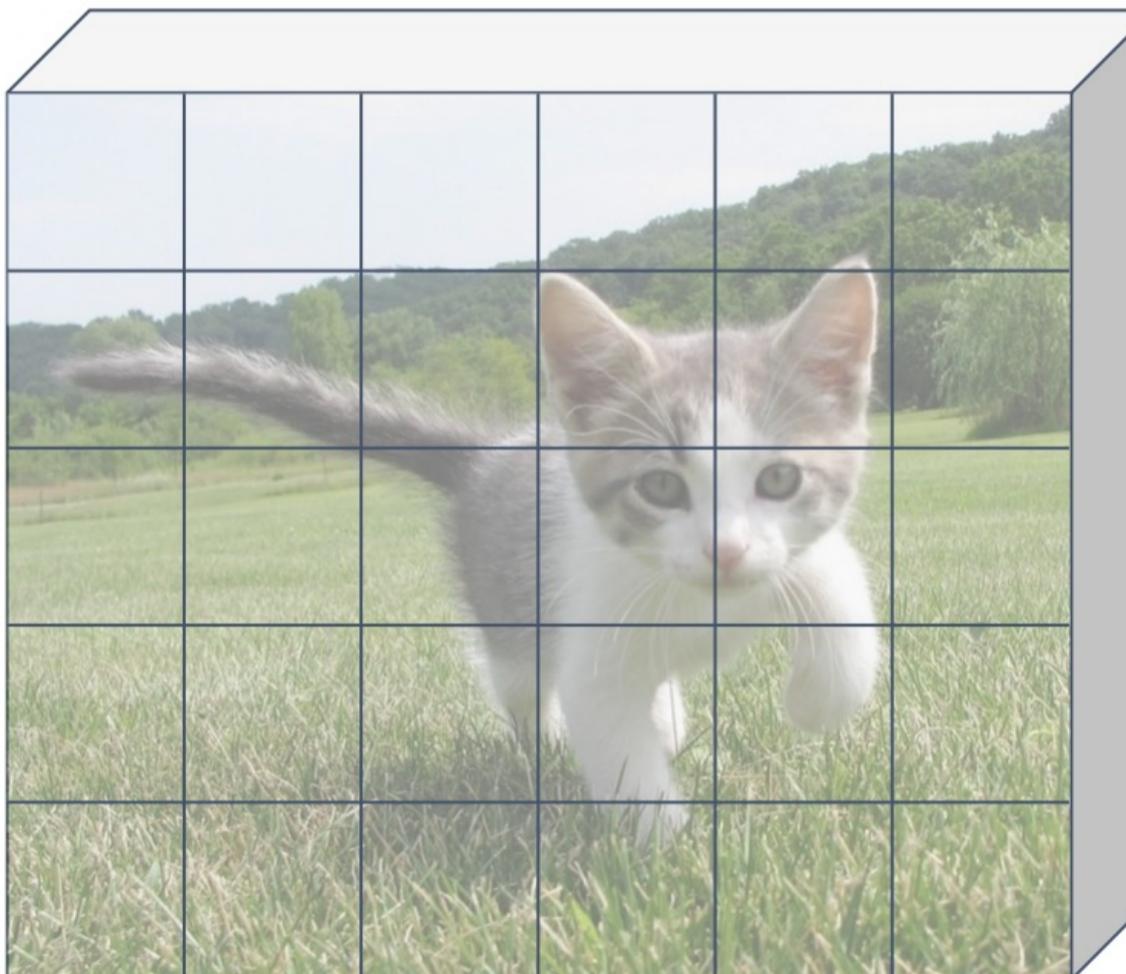
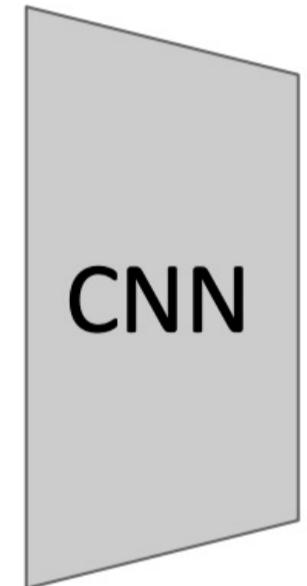


Image features
(e.g. $512 \times 5 \times 6$)



Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

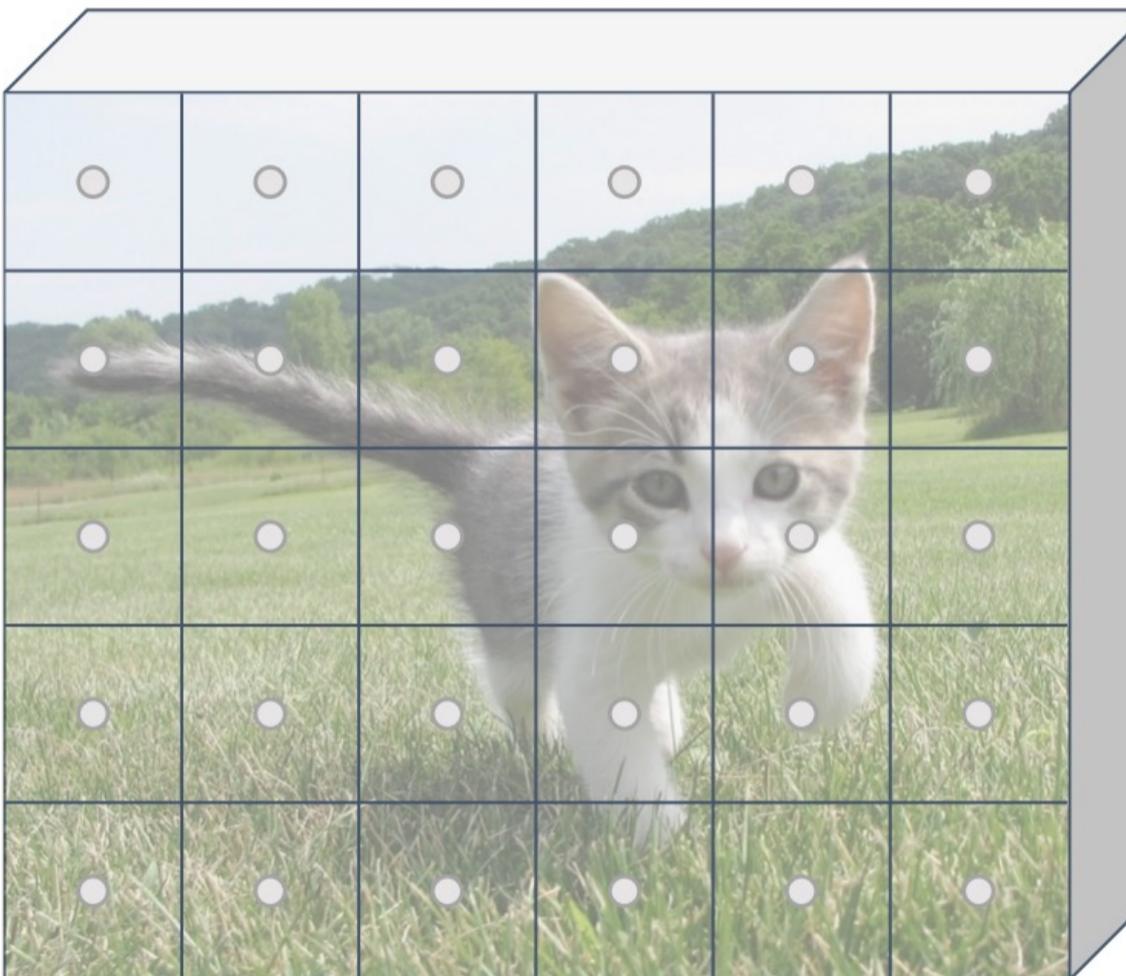
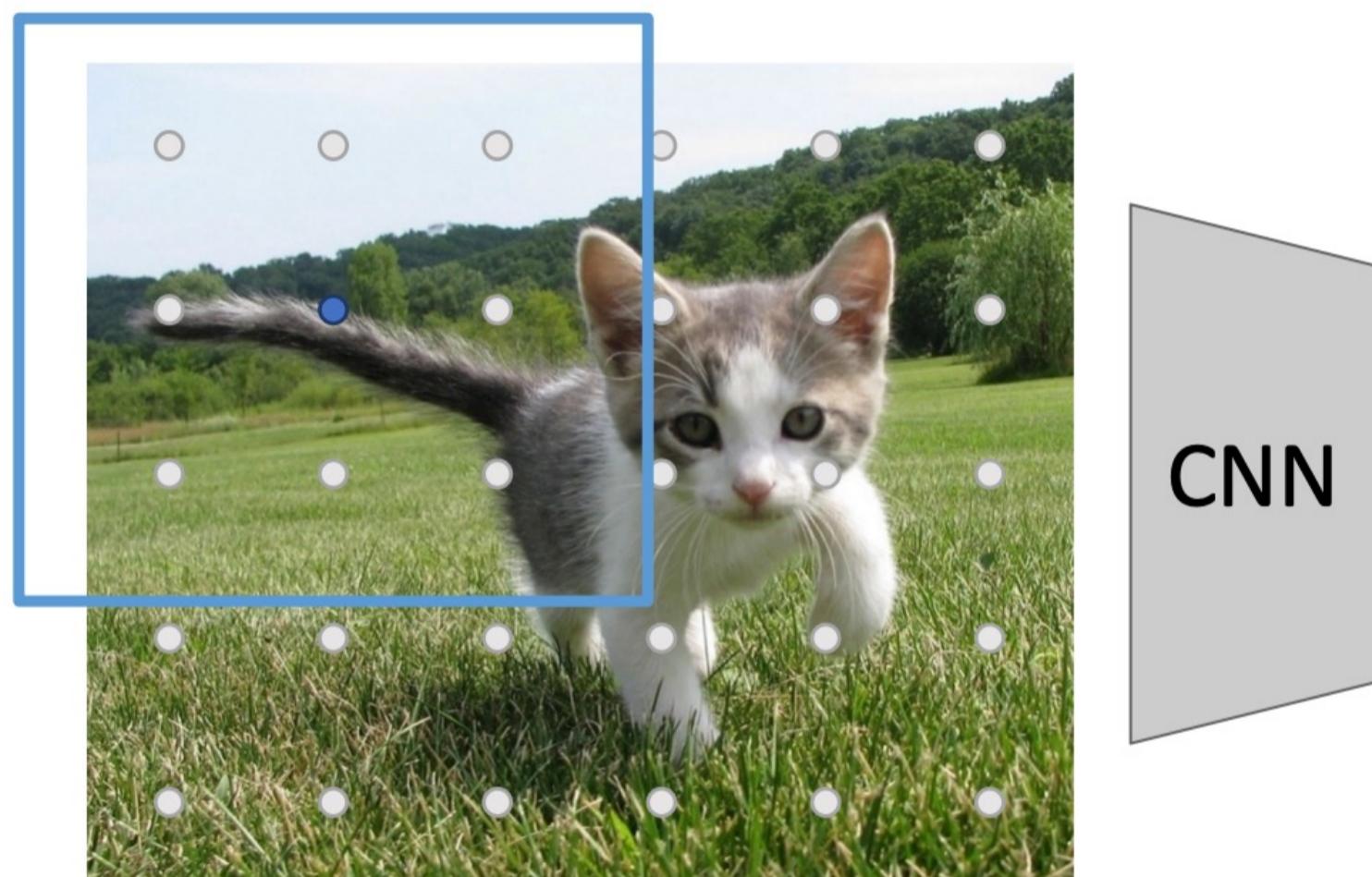


Image features
(e.g. $512 \times 5 \times 6$)



Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

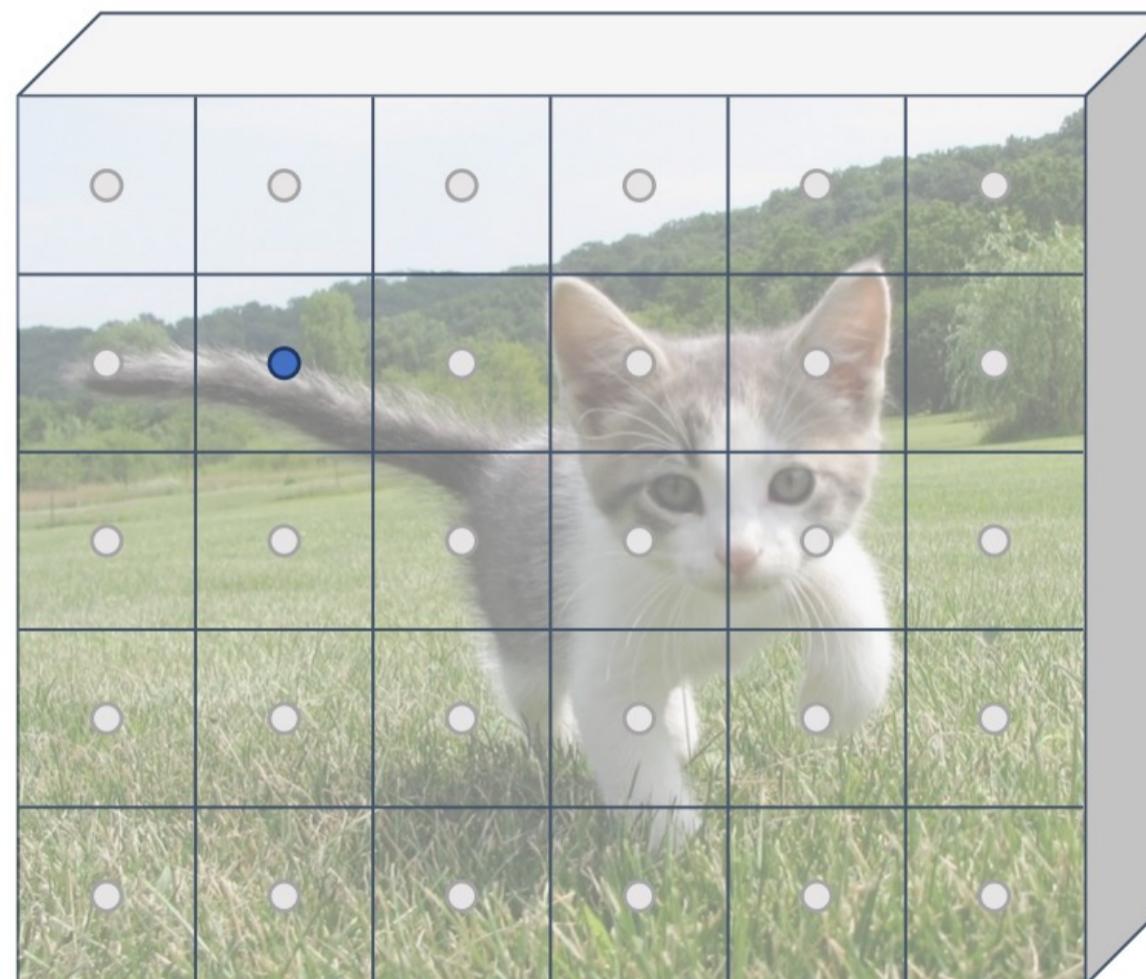
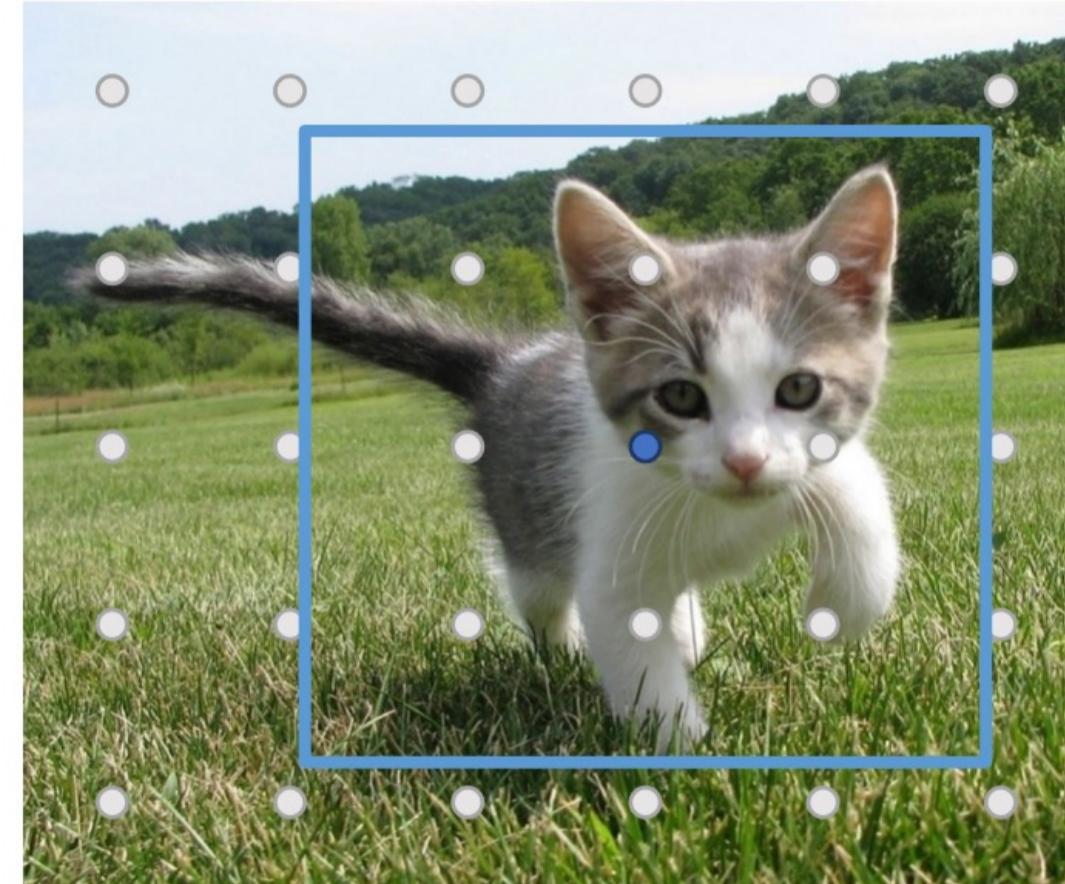


Image features
(e.g. $512 \times 5 \times 6$)

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

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

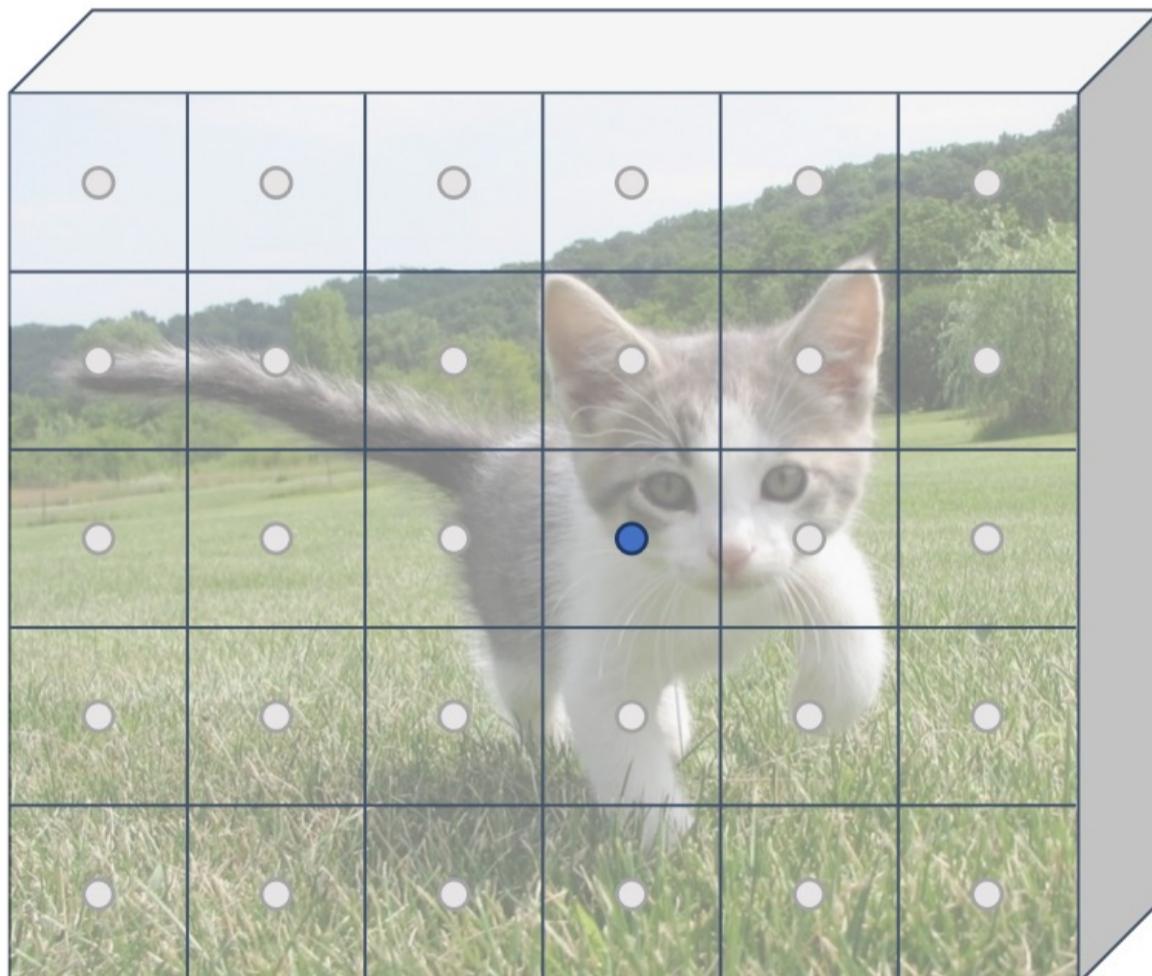
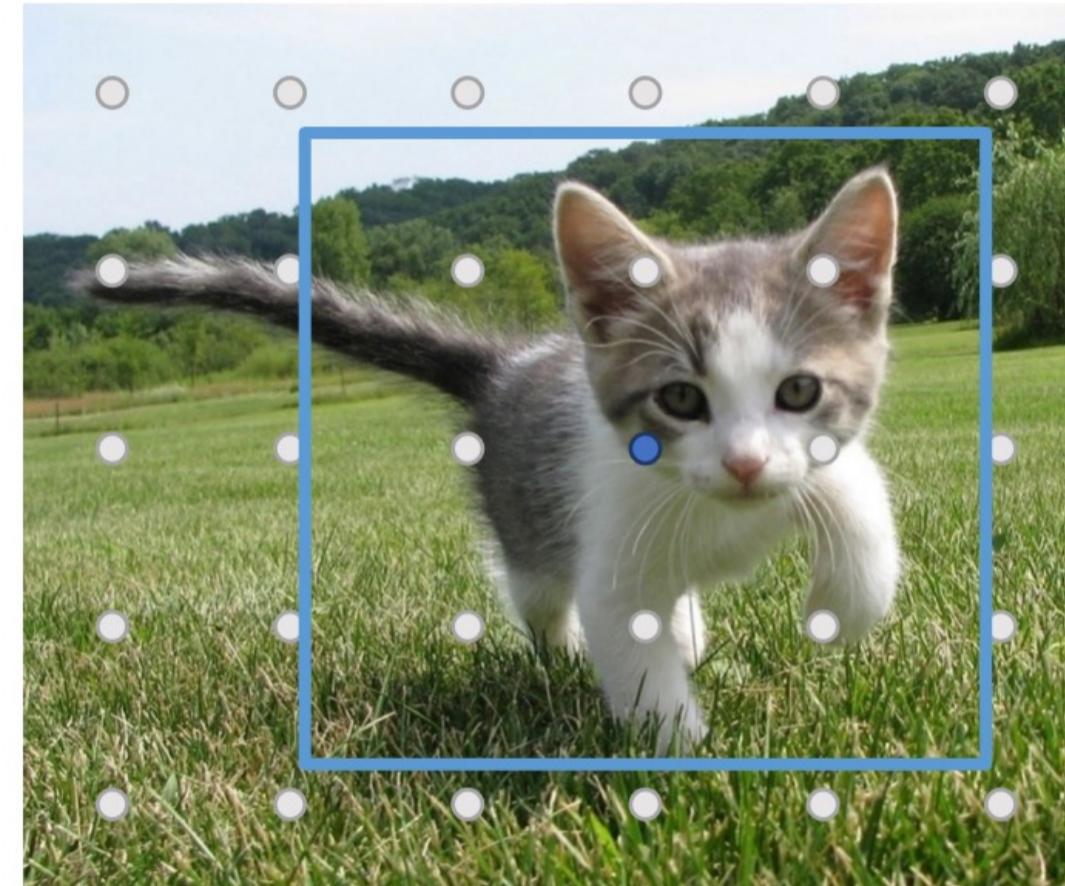


Image features
(e.g. $512 \times 5 \times 6$)

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

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

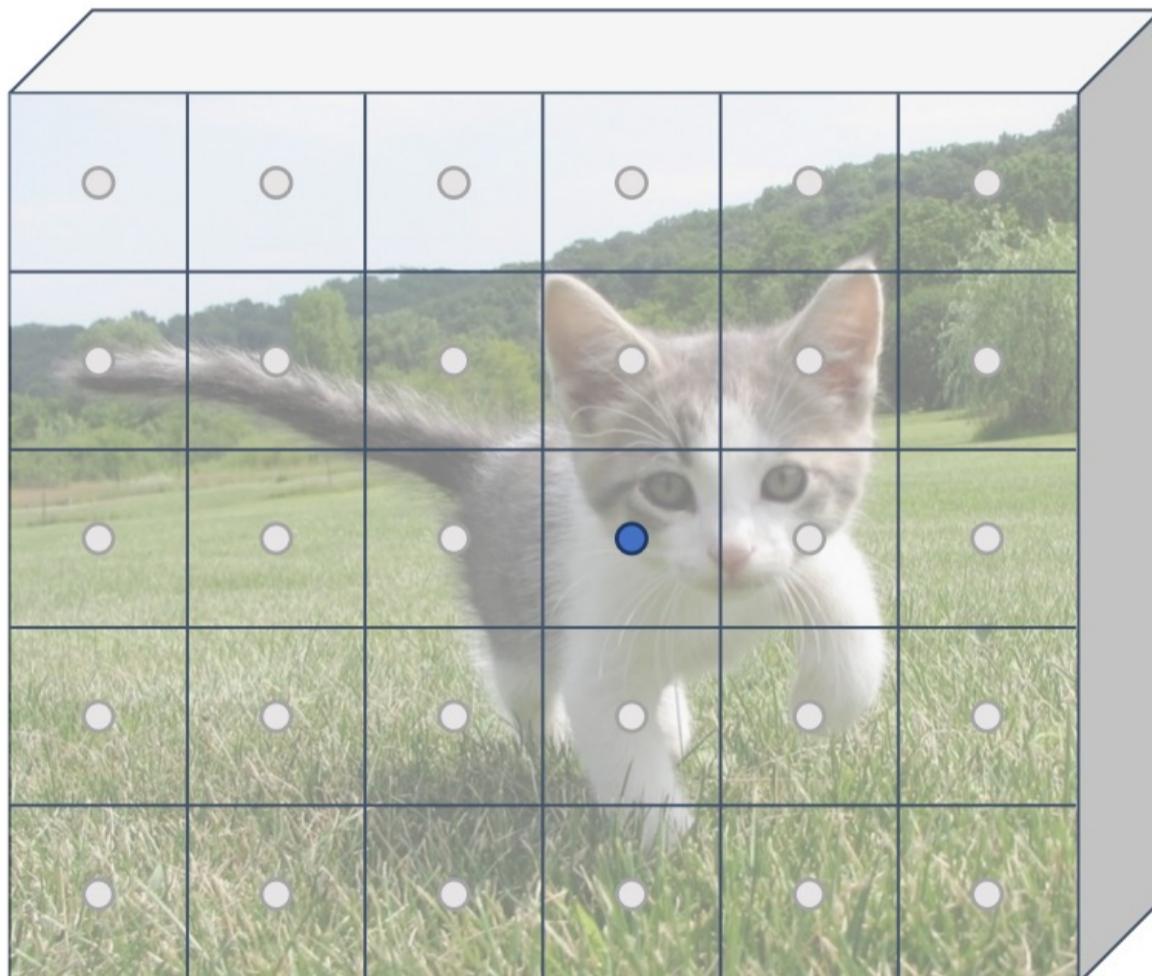
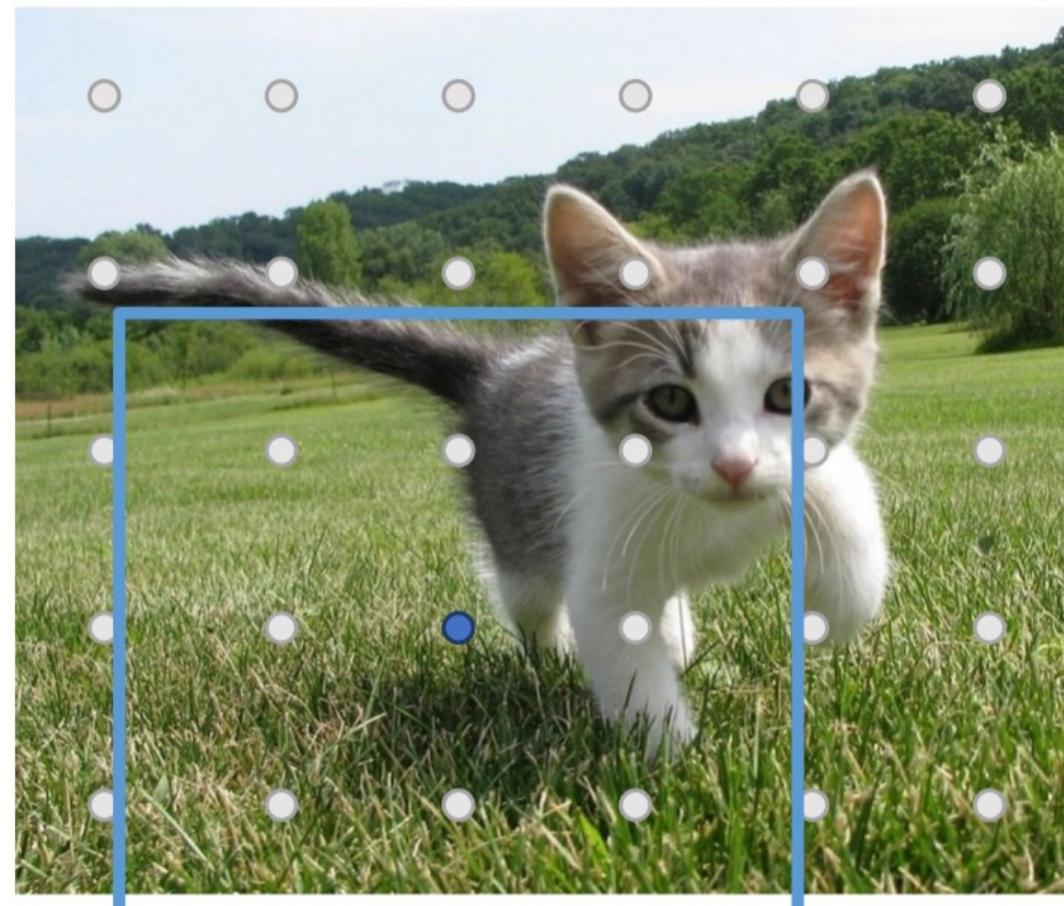


Image features
(e.g. $512 \times 5 \times 6$)

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

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

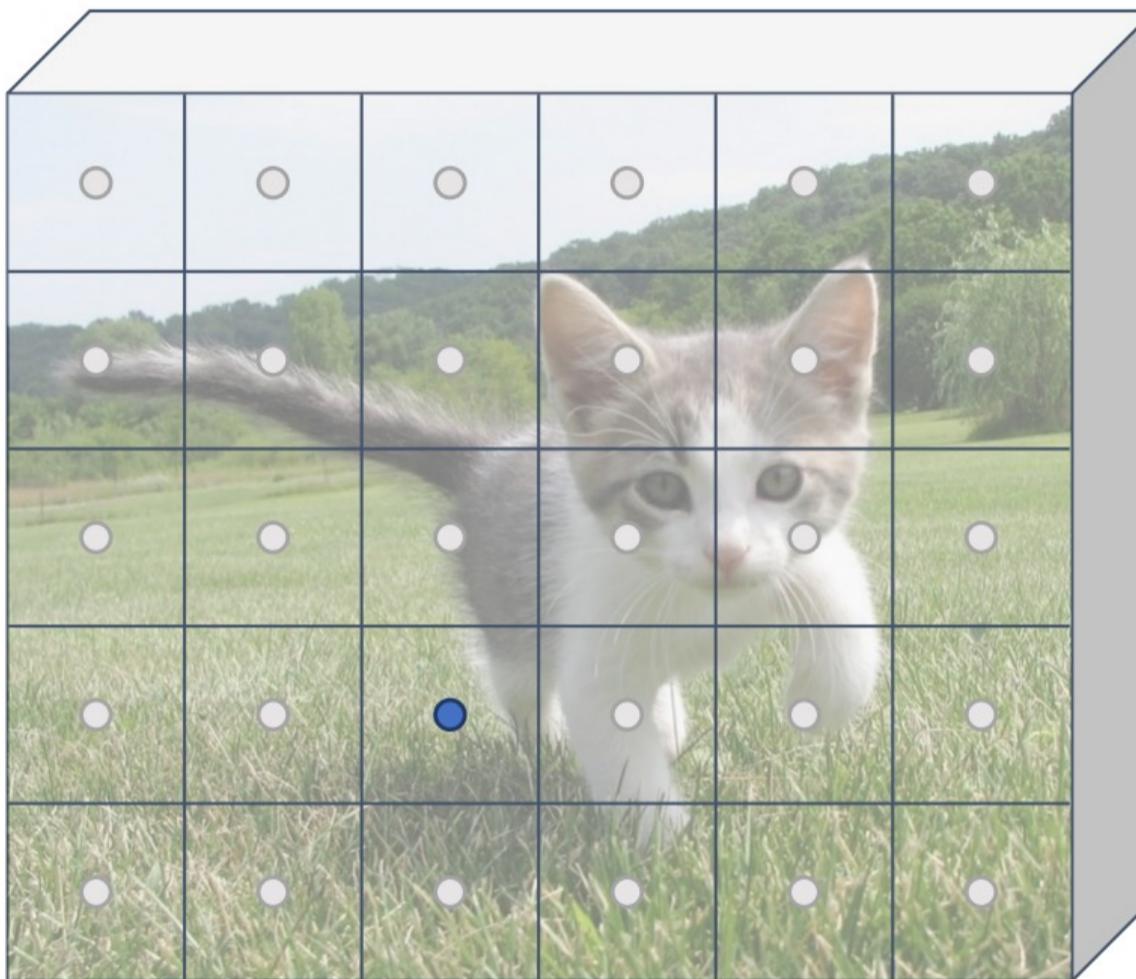
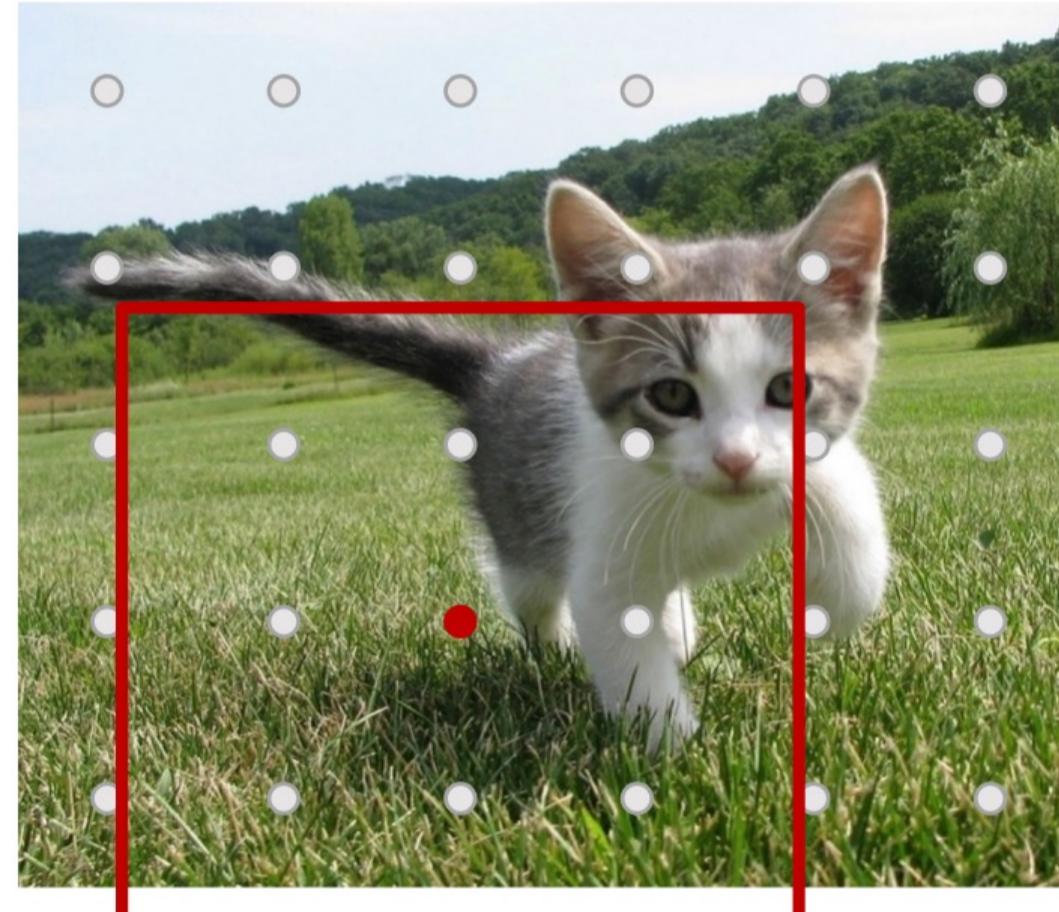


Image features
(e.g. $512 \times 5 \times 6$)

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

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



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

Each feature corresponds to point in the input

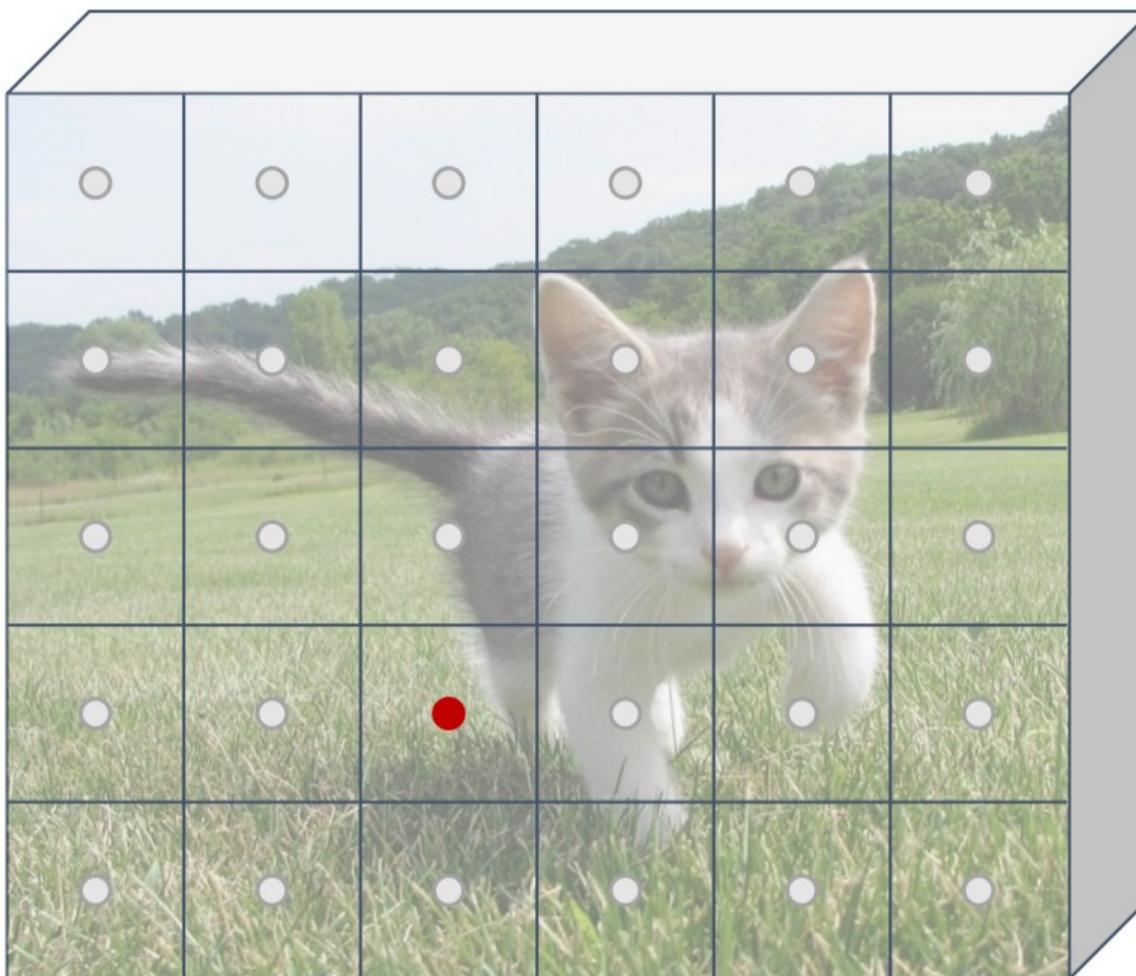
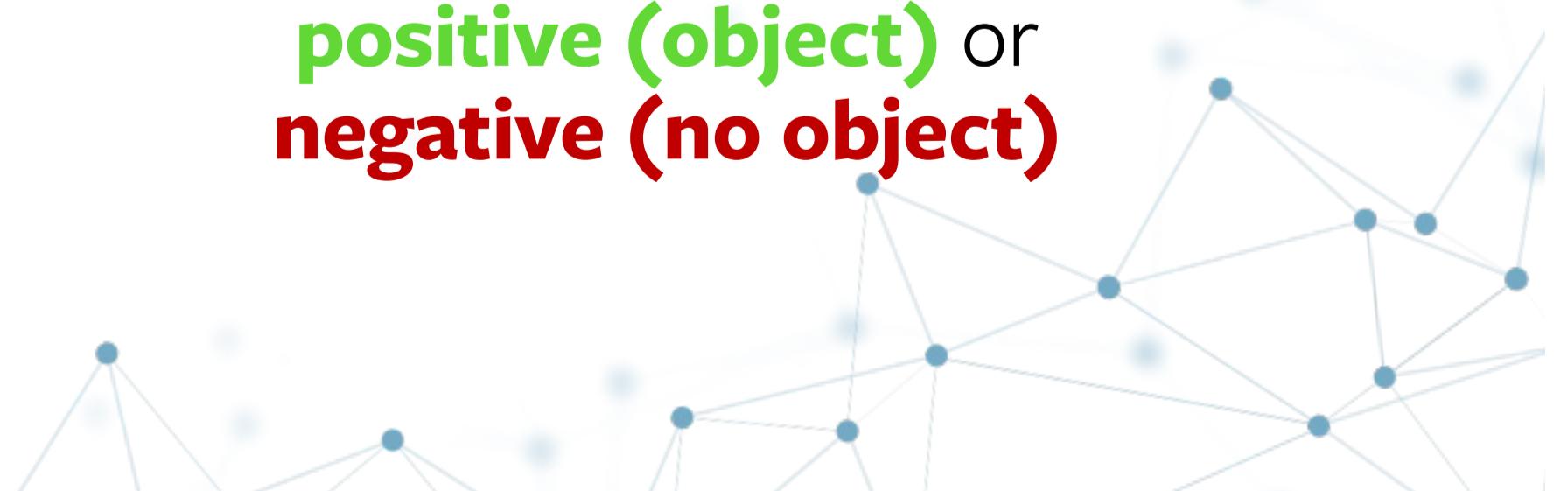


Image features
(e.g. $512 \times 5 \times 6$)

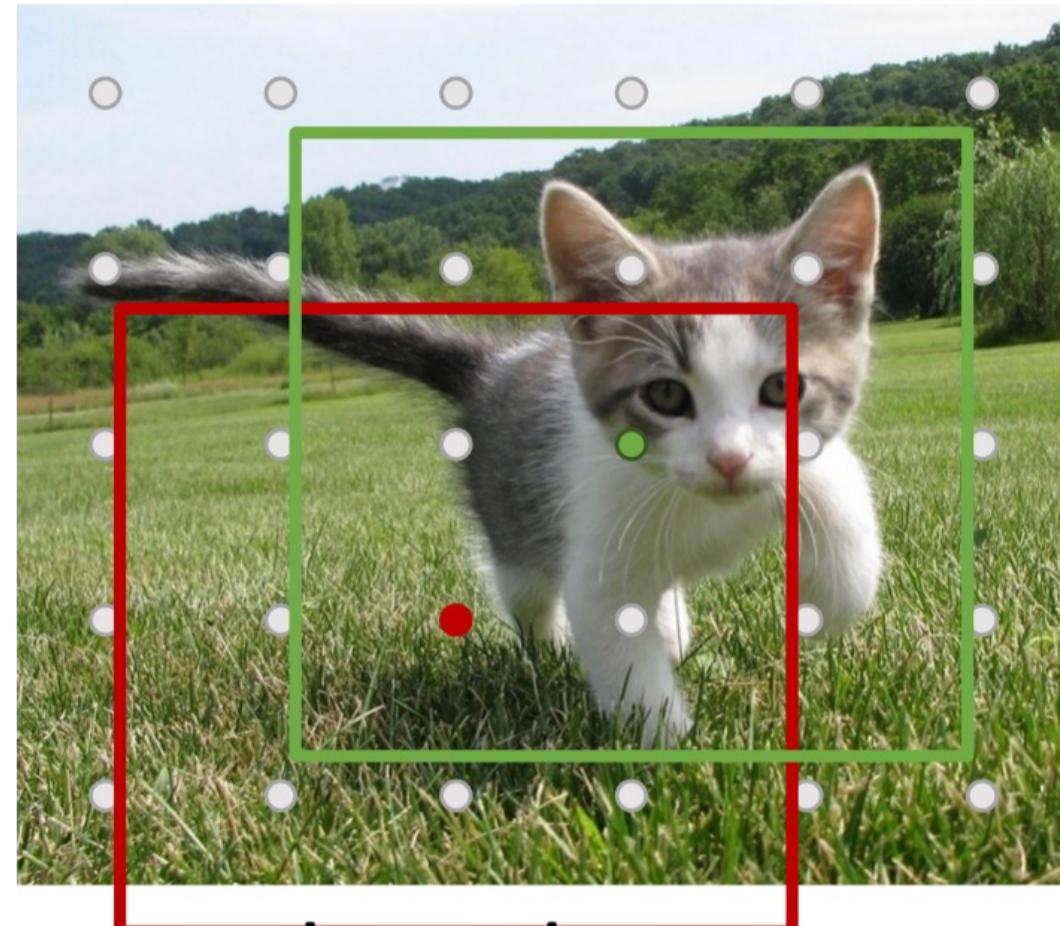
imagine an **anchor box** of fixed size at each point in the feature map

classify each anchor as
positive (object) or
negative (no object)



Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



Input Image

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

Each feature corresponds to point in the input

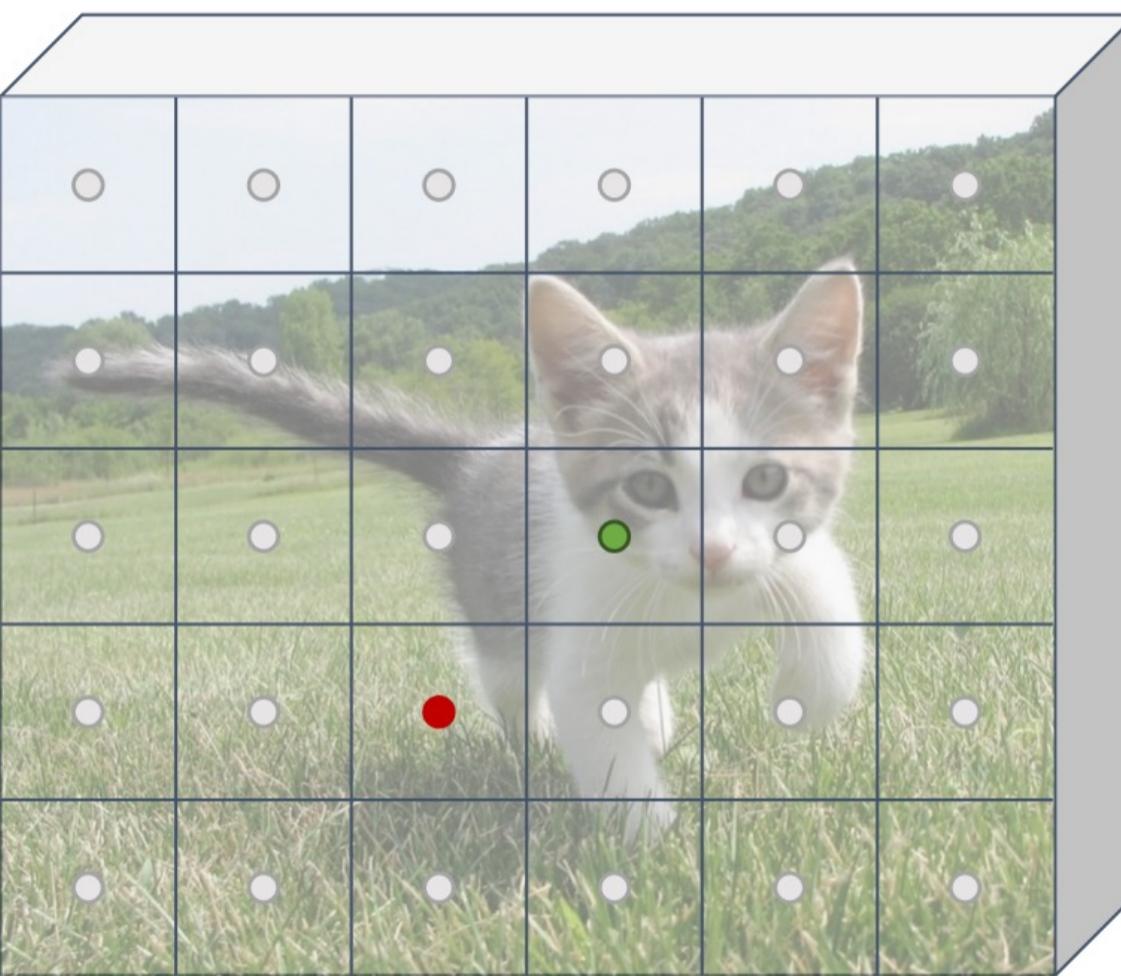


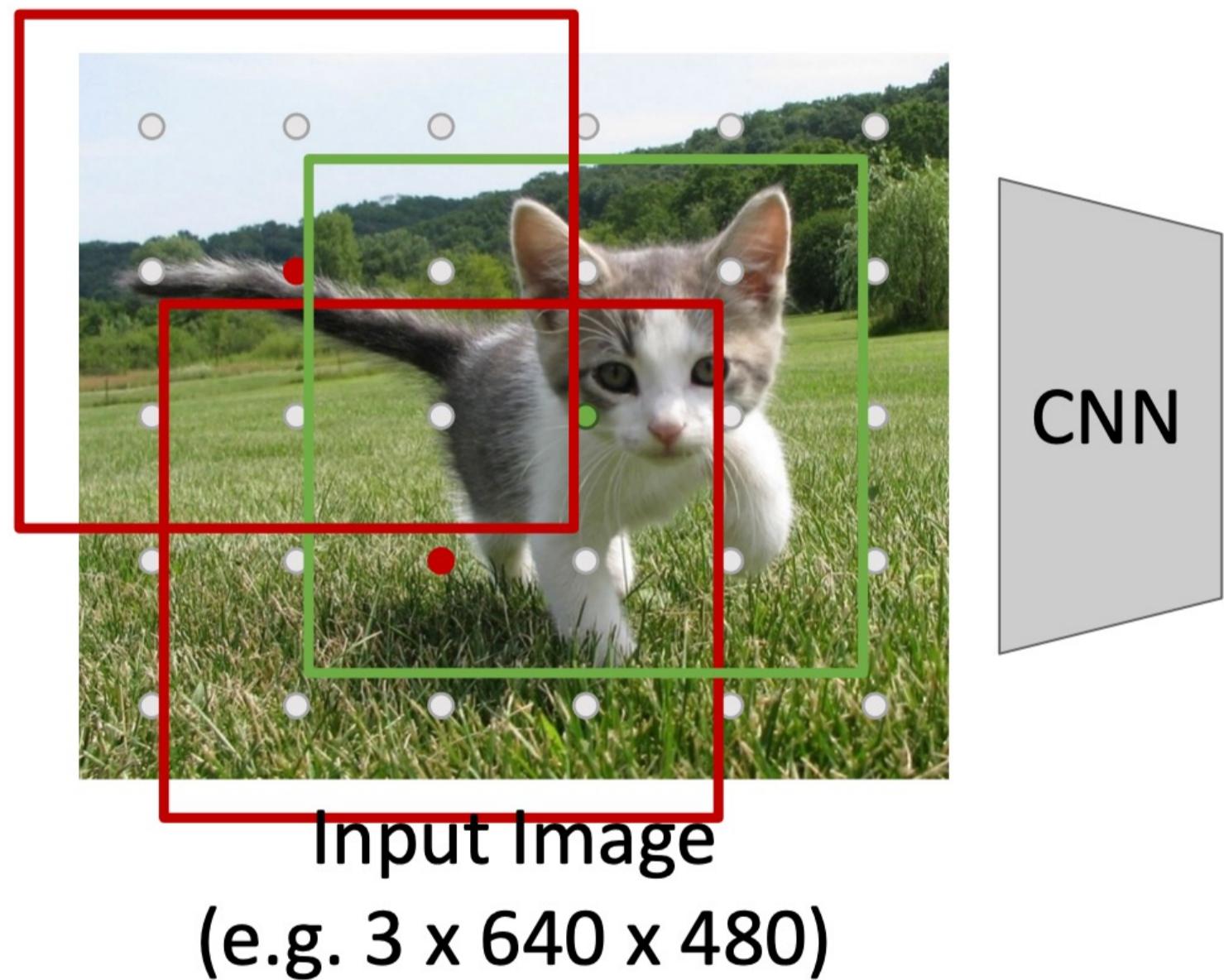
Image features
(e.g. $512 \times 5 \times 6$)

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

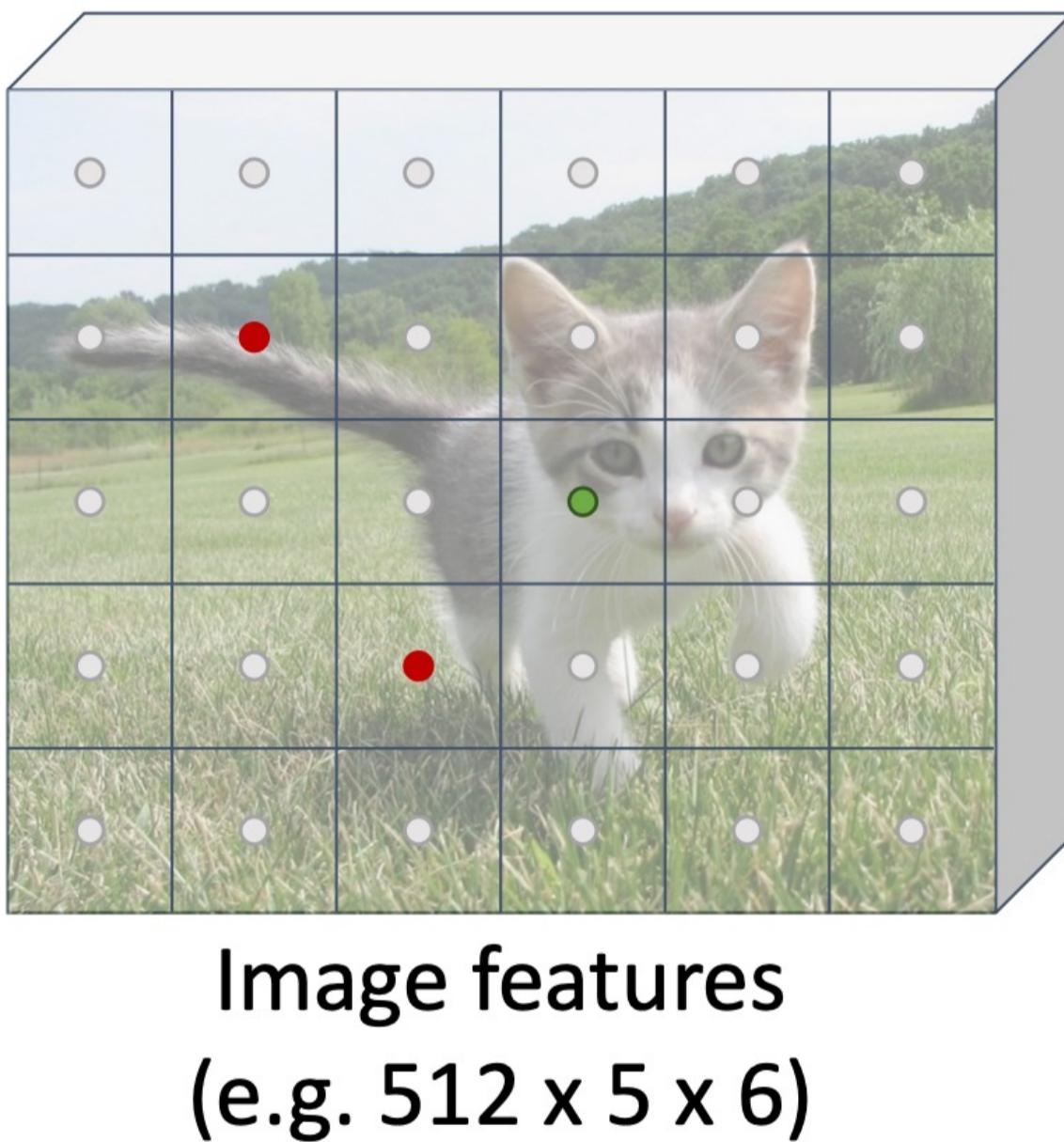
classify each anchor as
positive (object) or
negative (no object)

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to point in the input

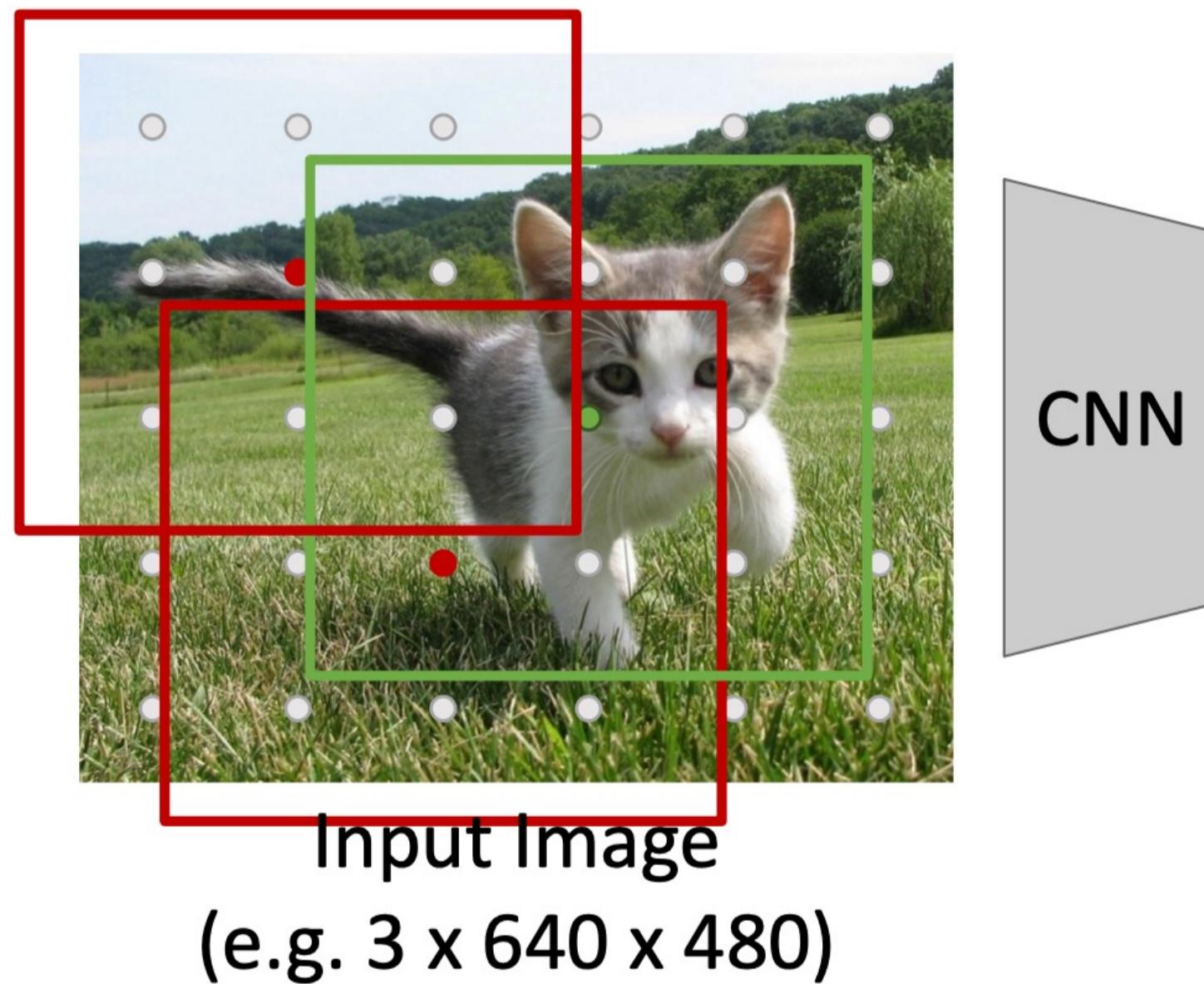


imagine an **anchor box** of fixed size at each point in the feature map

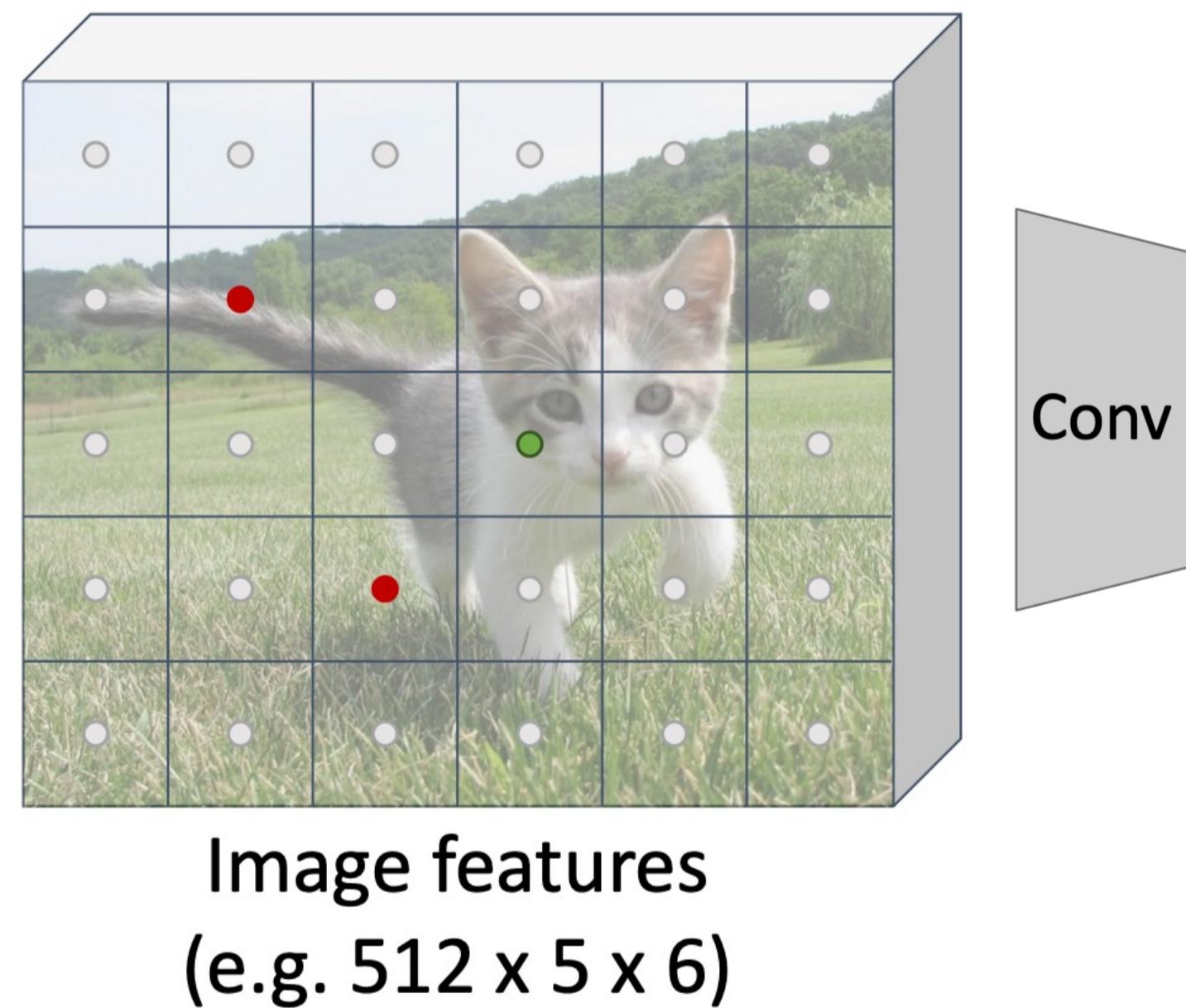
classify each anchor as
positive (object) or
negative (no object)

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to point in the input



project object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)

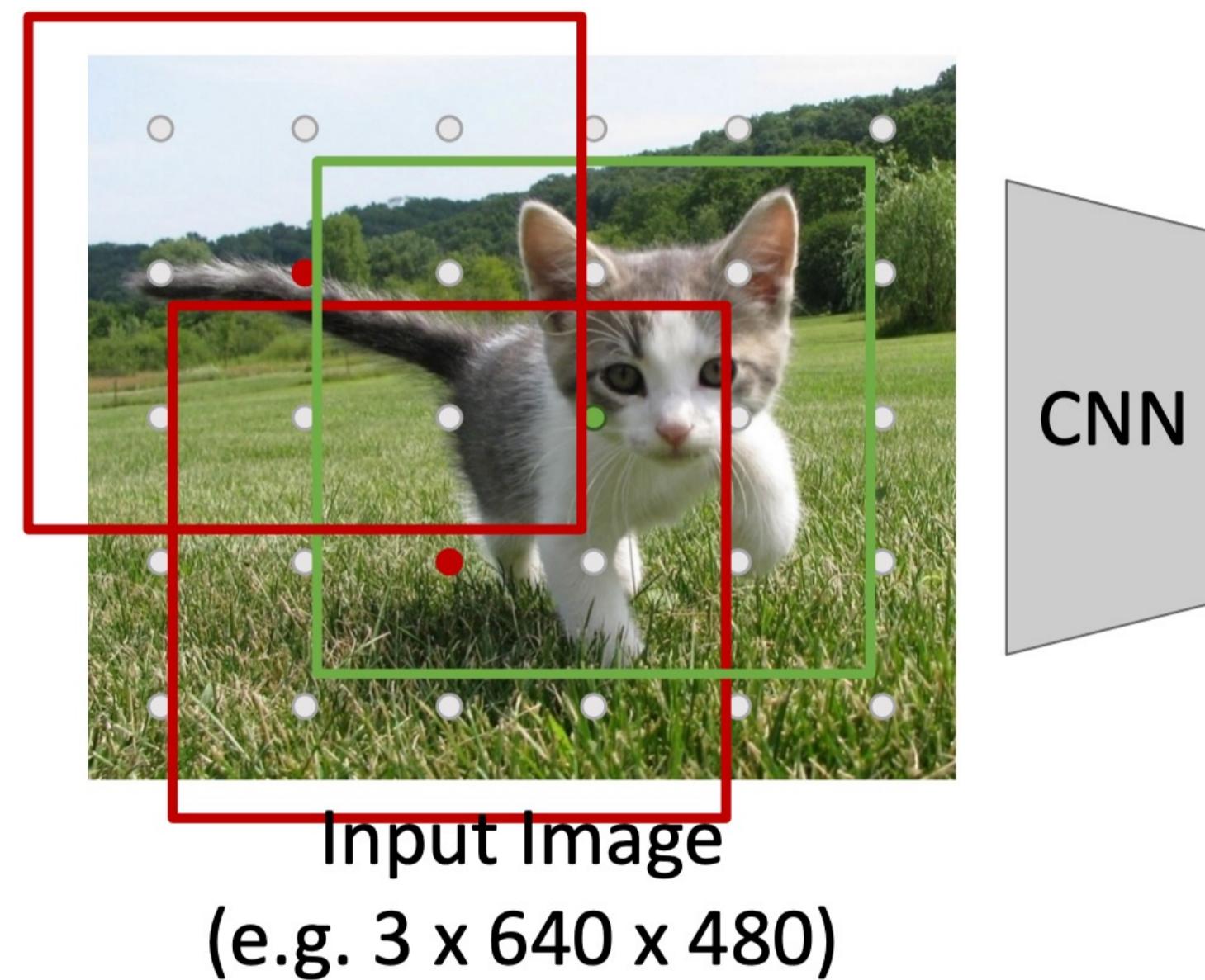
Anchor is object?
 $2 \times 5 \times 6$

classify each anchor as
positive (object) or
negative (no object)

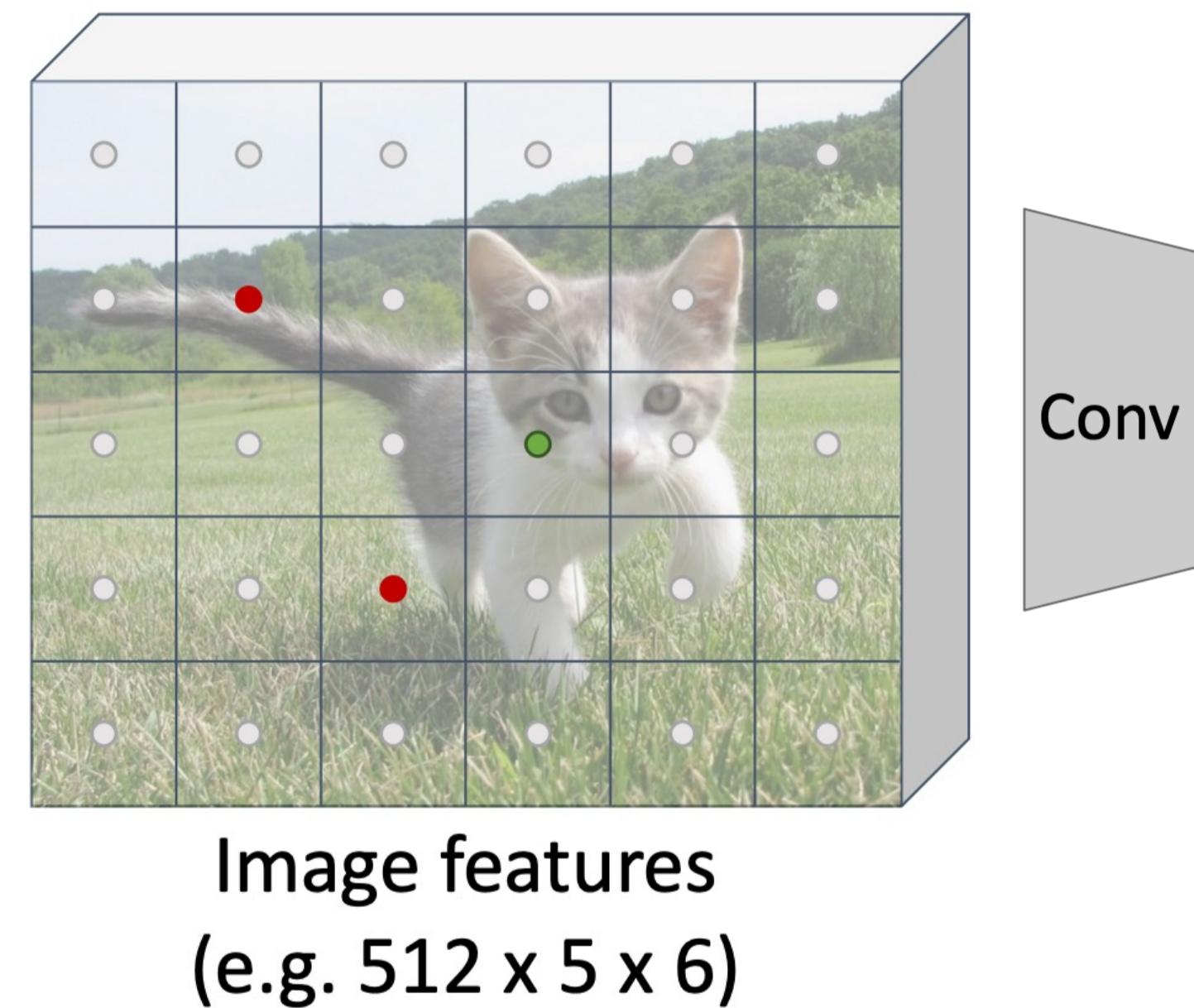


Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



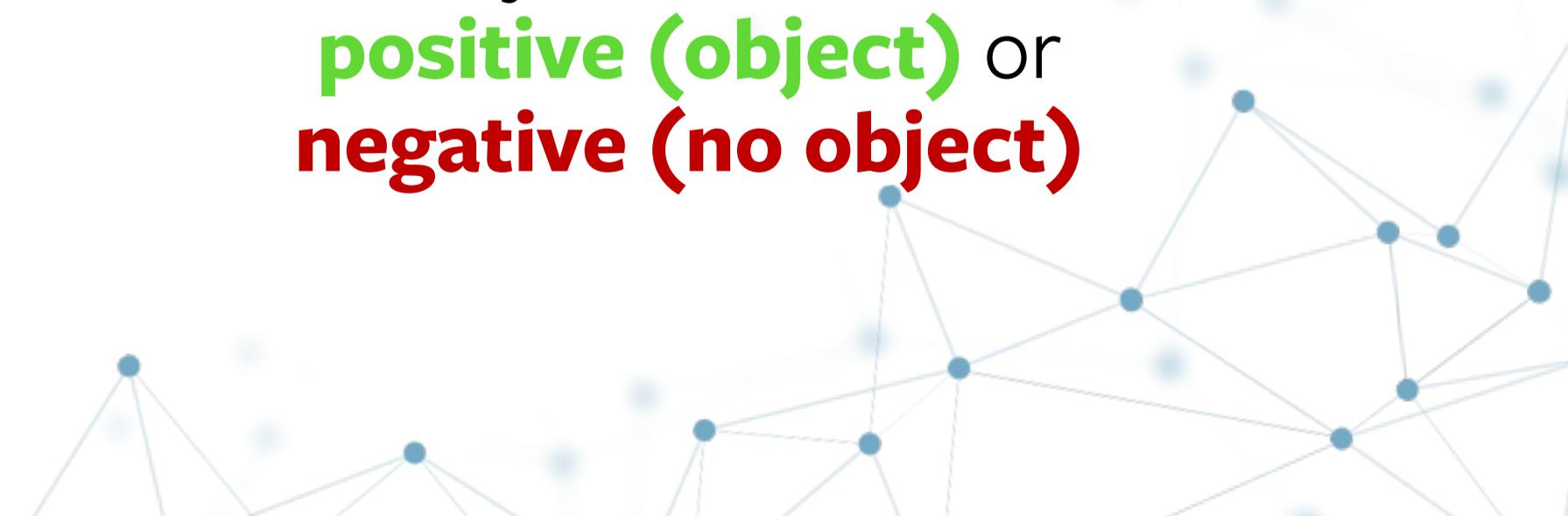
Each feature corresponds to point in the input



project object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)

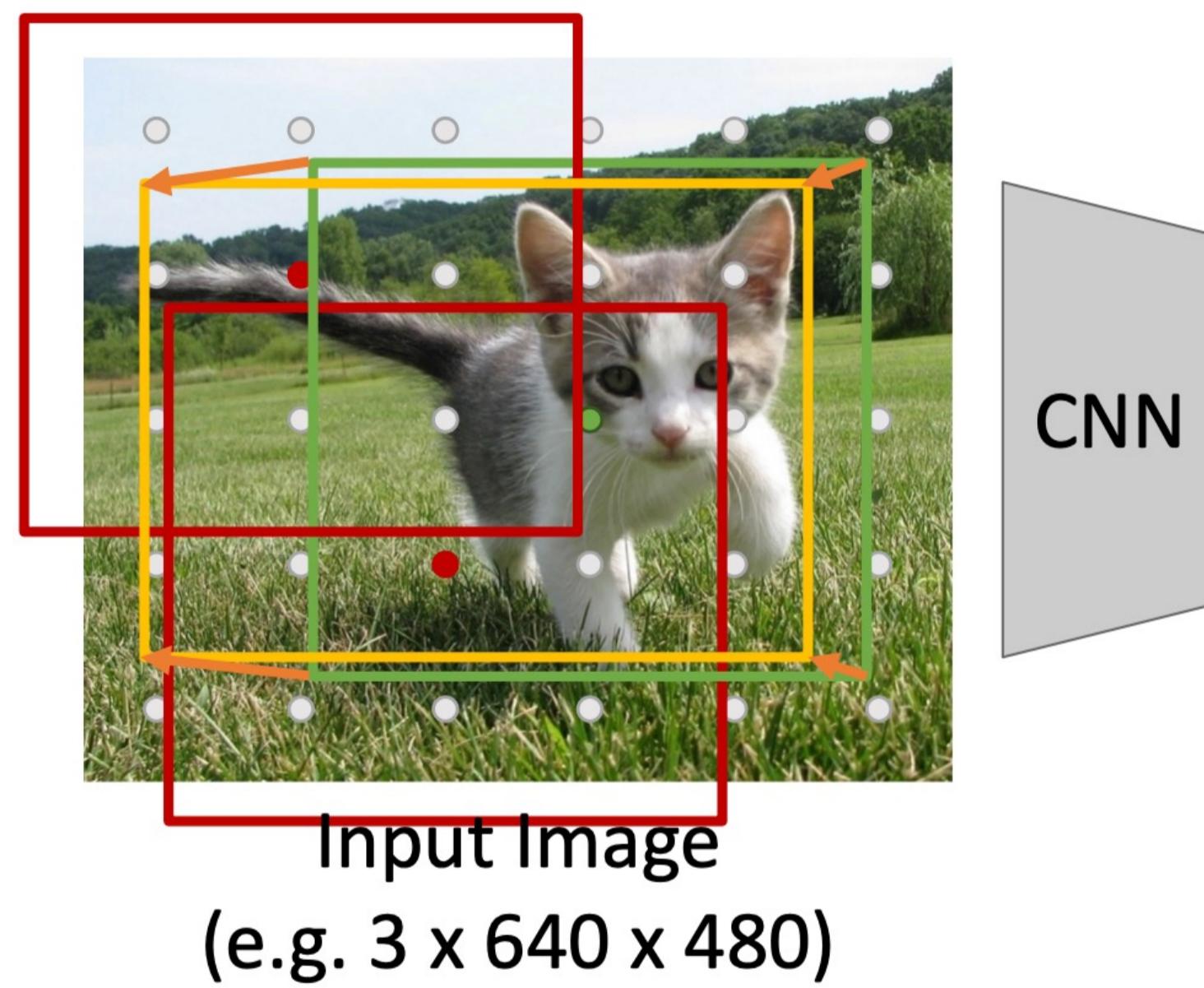
Anchor is object?
 $2 \times 5 \times 6$

classify each anchor as
positive (object) or
negative (no object)

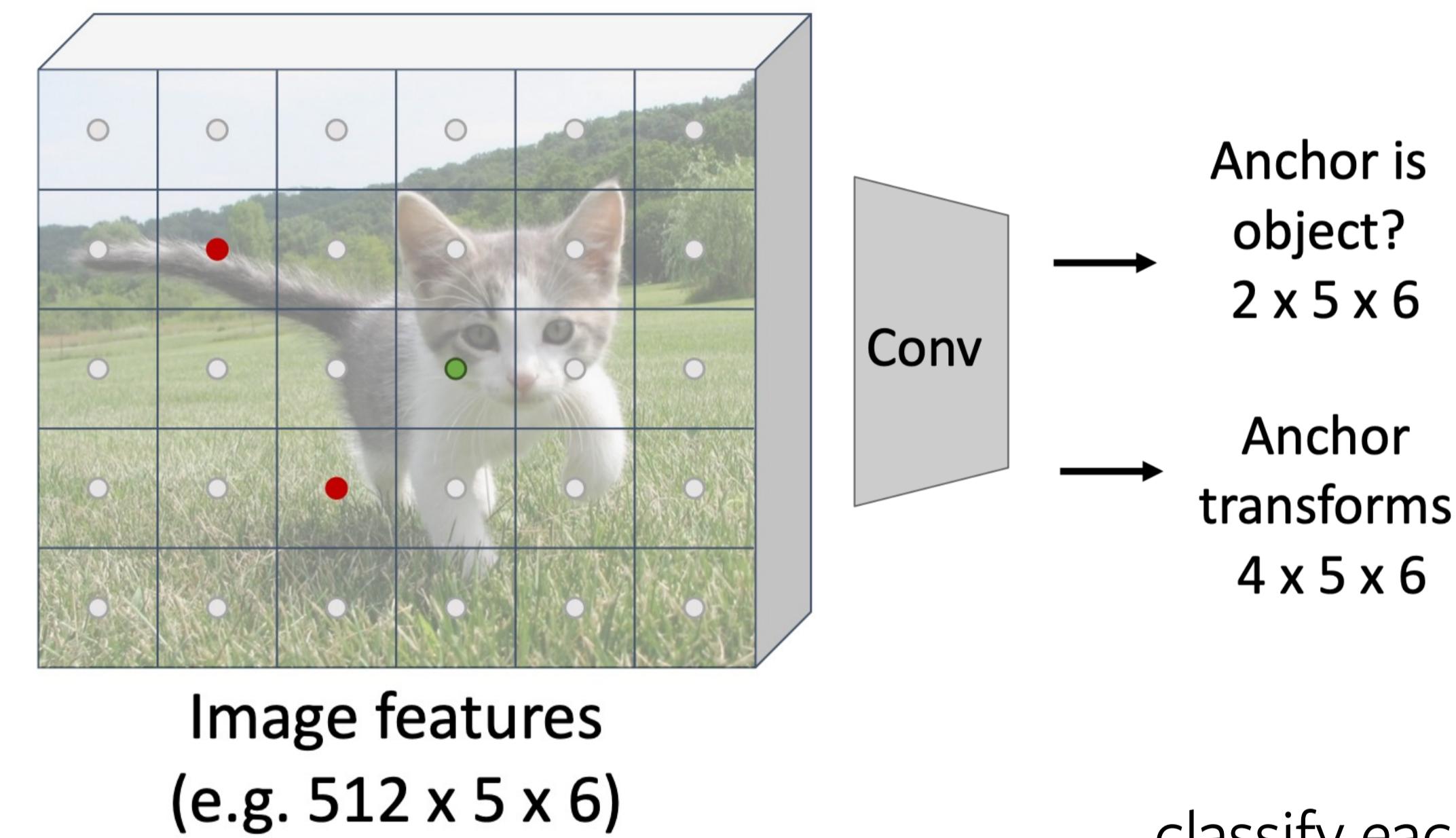


Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to point in the input



for **positive anchors**, also predict a **transform** that is converting the anchor to the **GT box** (like R-CNN)

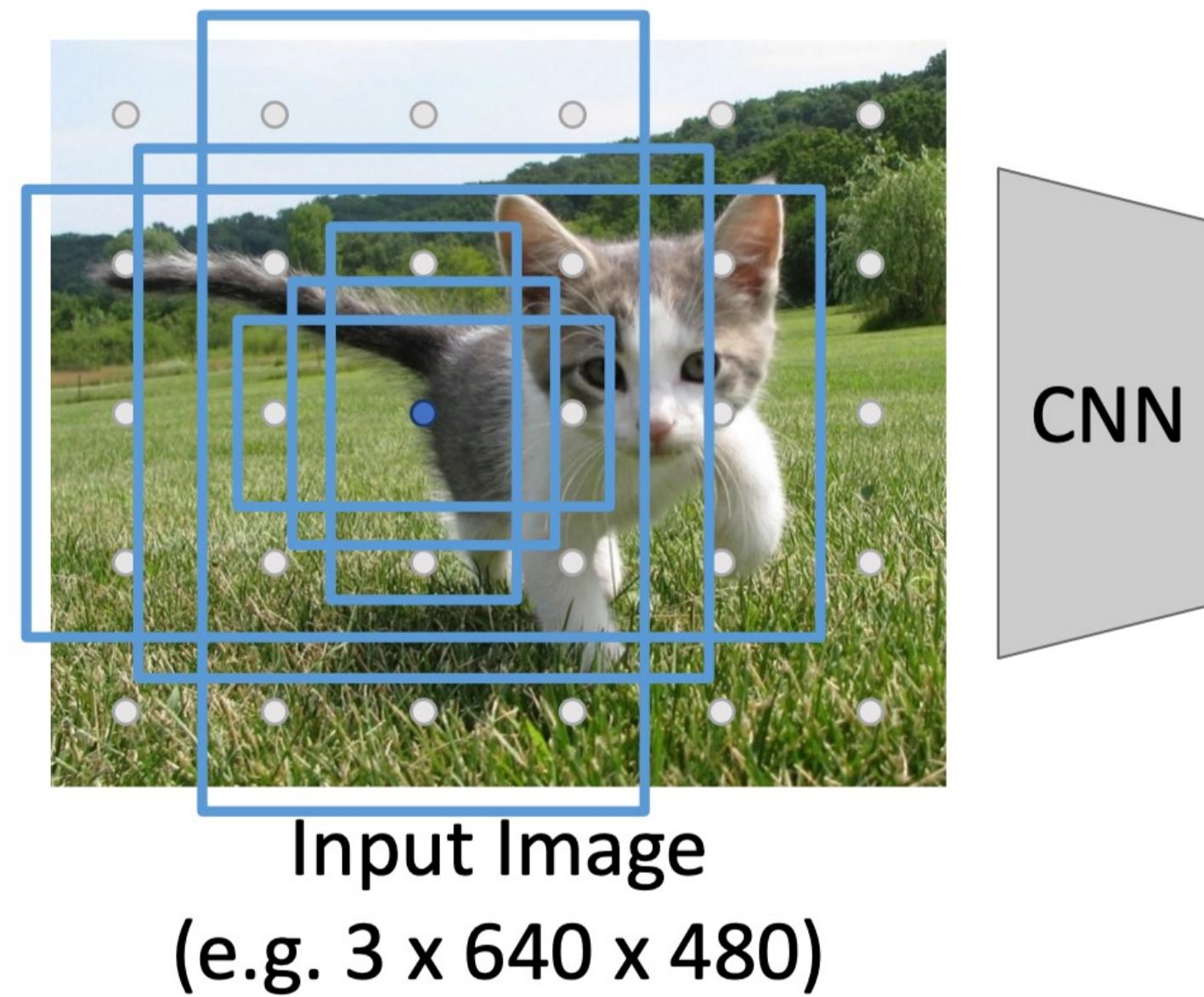
Anchor is object?
 $2 \times 5 \times 6$

Anchor transforms
 $4 \times 5 \times 6$

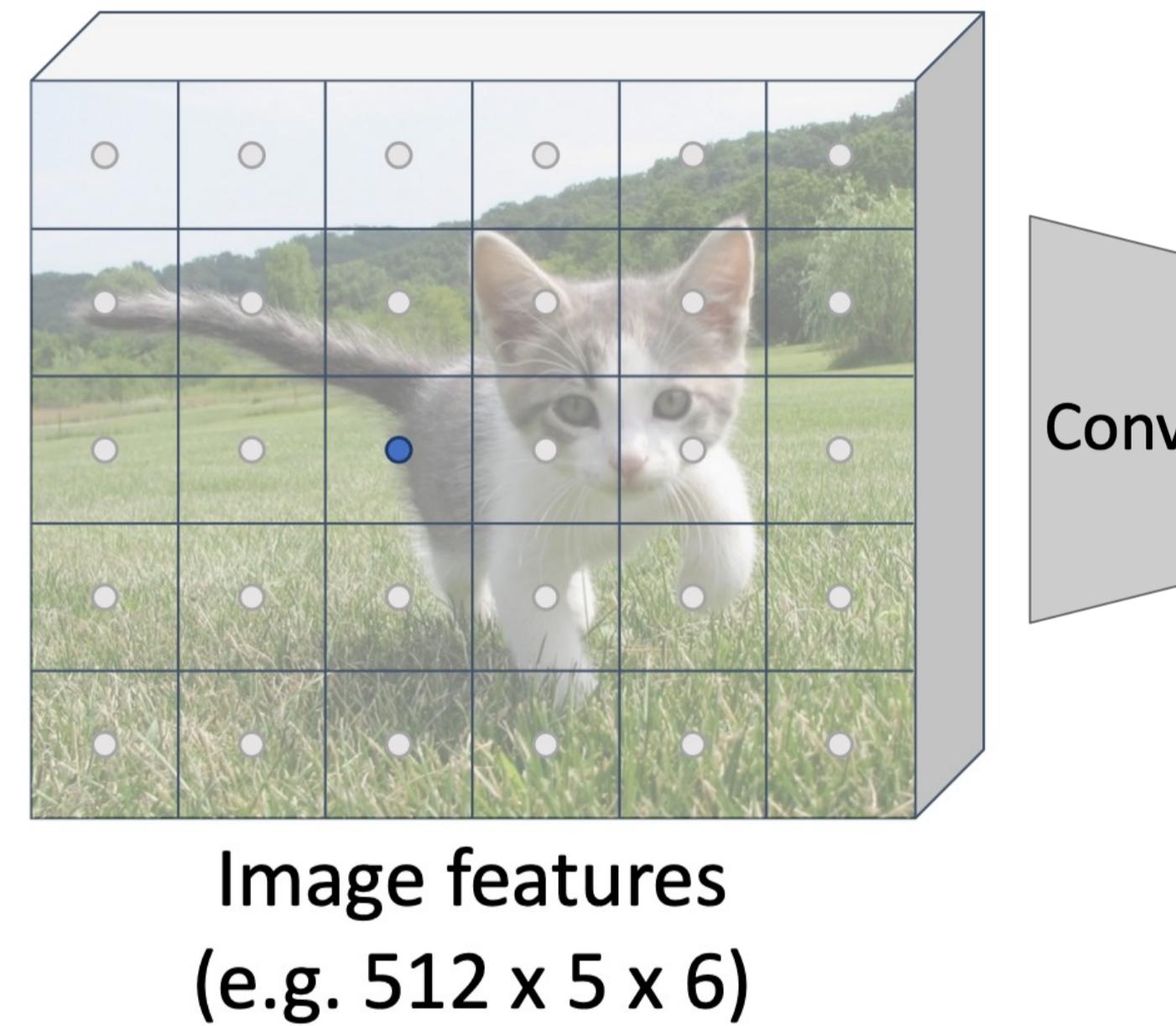
classify each anchor as
positive (object) or
negative (no object)

Region proposal network (RPN)

Run backbone CNN to get features aligned to input image



Each feature corresponds to point in the input



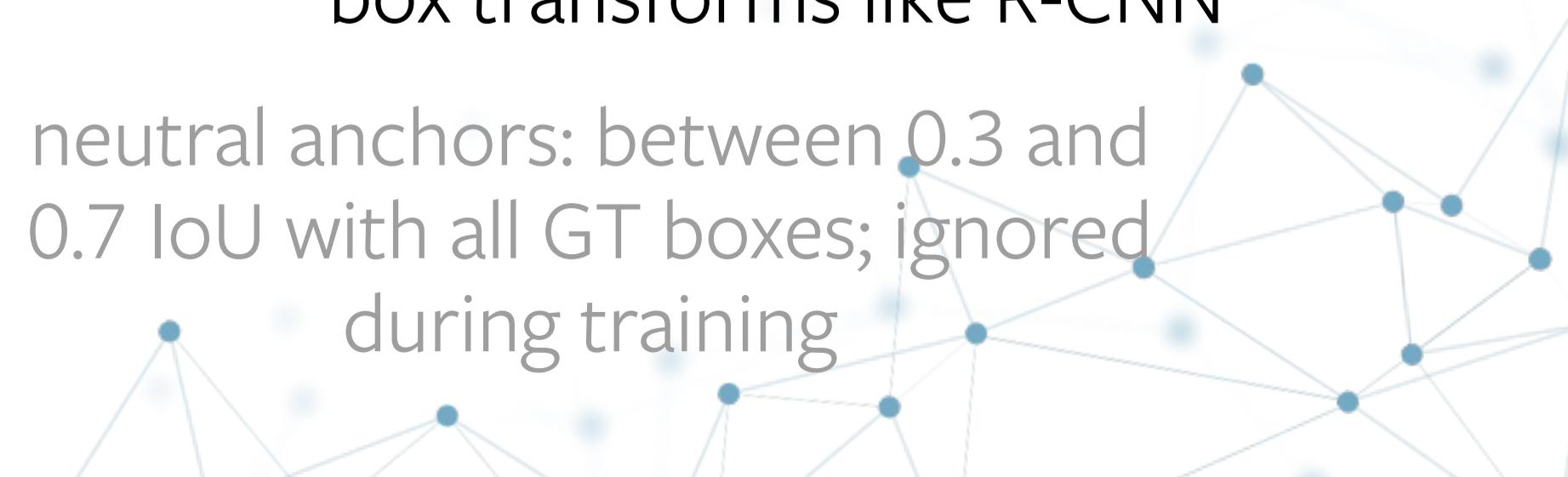
in practice: rather than using one anchor per point, instead consider K different anchors with different size and scale (here K=6)

positive anchors: ≥ 0.7 IoU with some GT box (plus highest IoU to each GT)

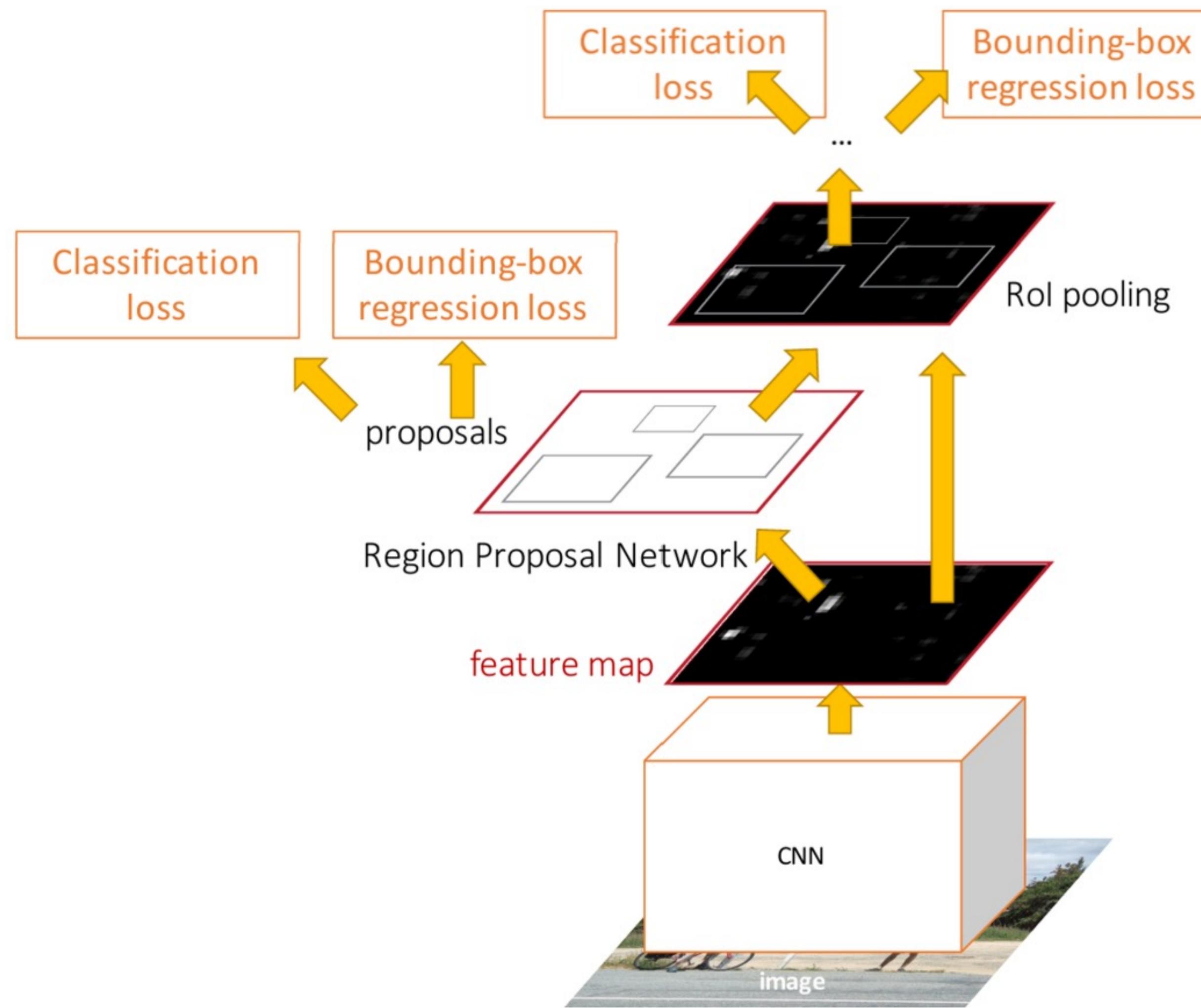
negative anchors: ≤ 0.3 IoU with all GT boxes. Don't supervise transforms for negative boxes

during training supervise positive / negative anchors and box transforms like R-CNN

neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training



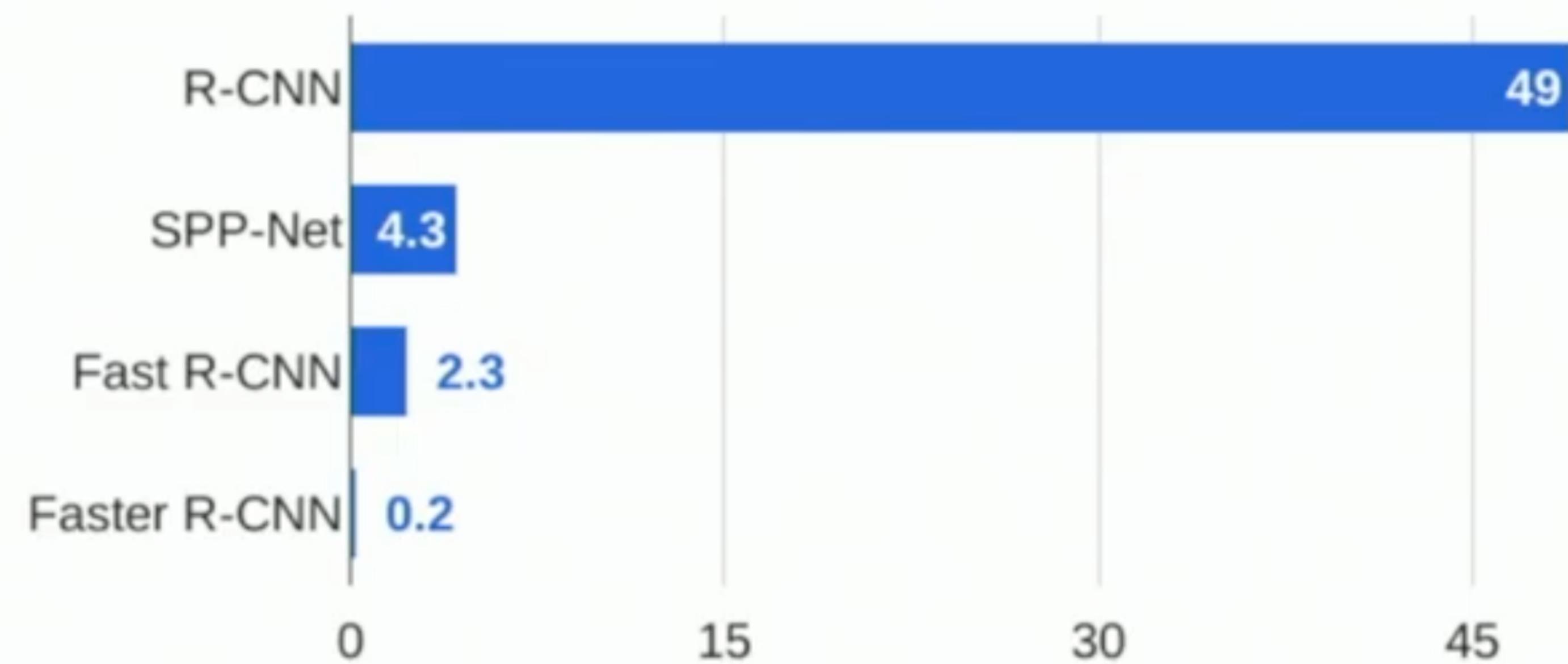
Faster R-CNN: Learnable Region Proposals



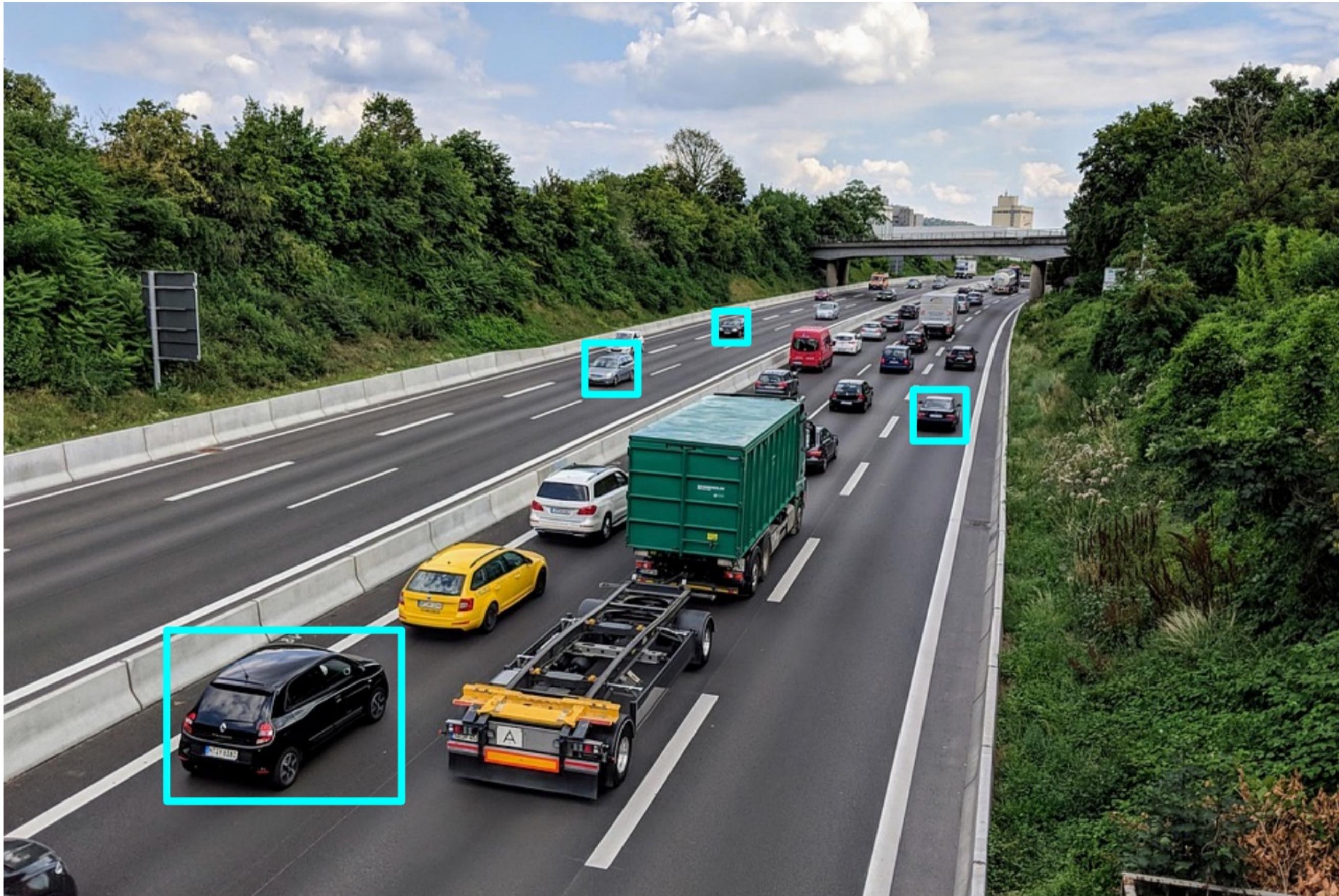
Jointly train with 4 losses:

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 box to object box

R-CNN Test-Time Speed



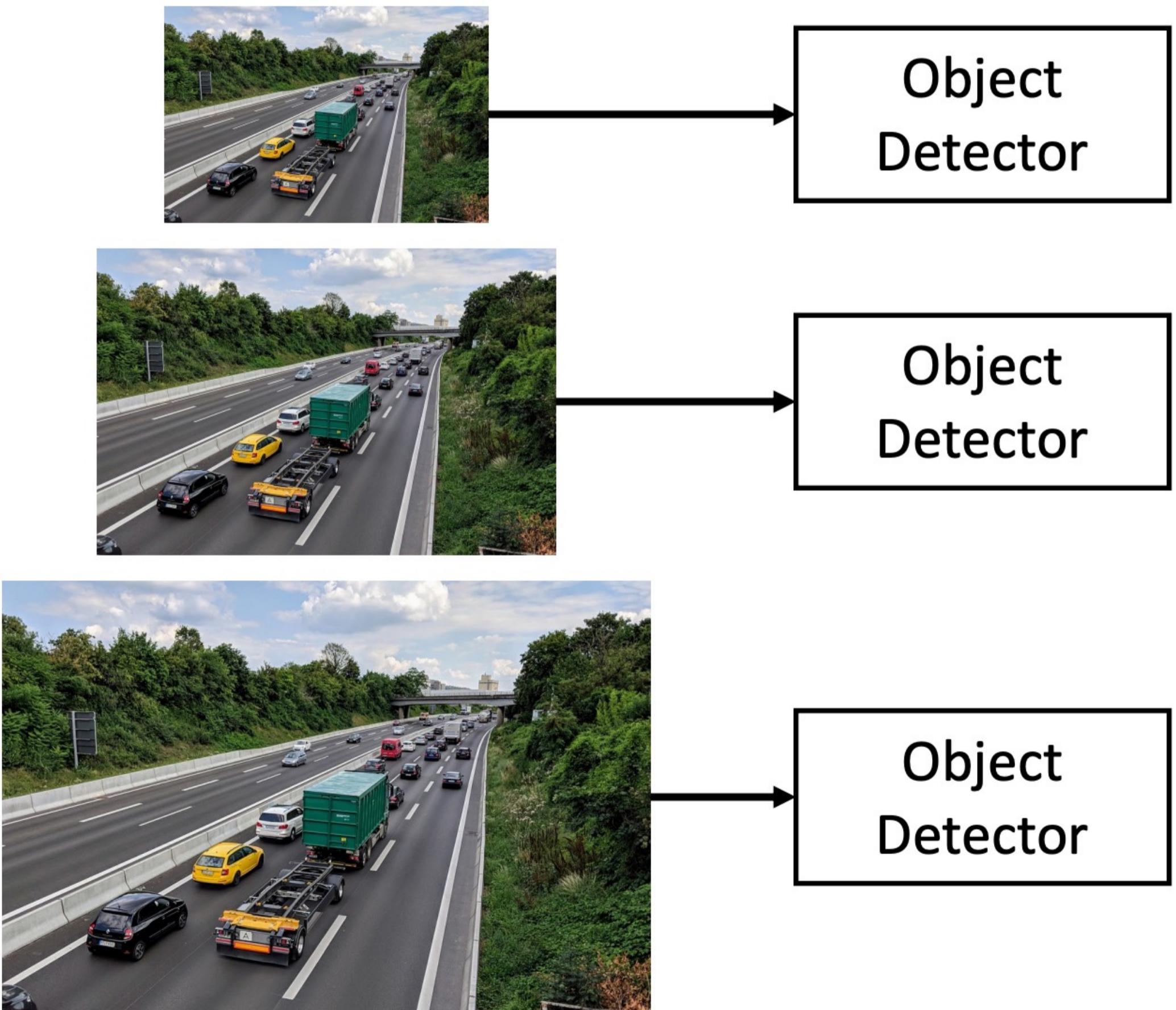
Dealing with scale



We need to detect objects of many different scales. How to improve scale invariance of the detector?



Classical idea: image pyramid

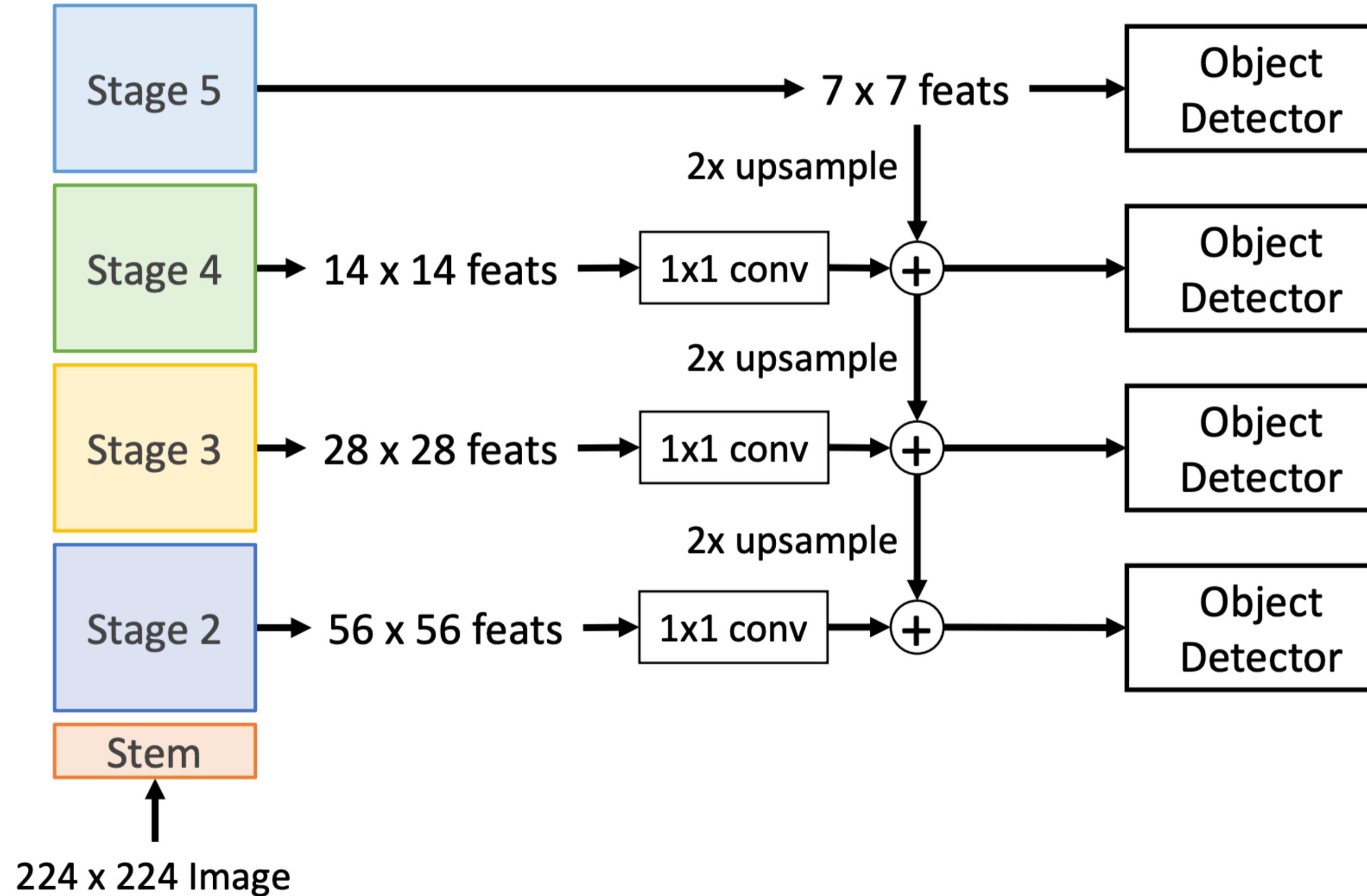


Build an image pyramid by resizing the image to different scales, then process each image scale independently.

Problem: expensive!



Feature Pyramid Network (FPN)



Idea: add top down connections that feed information from high level features back down to lower level features

Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice

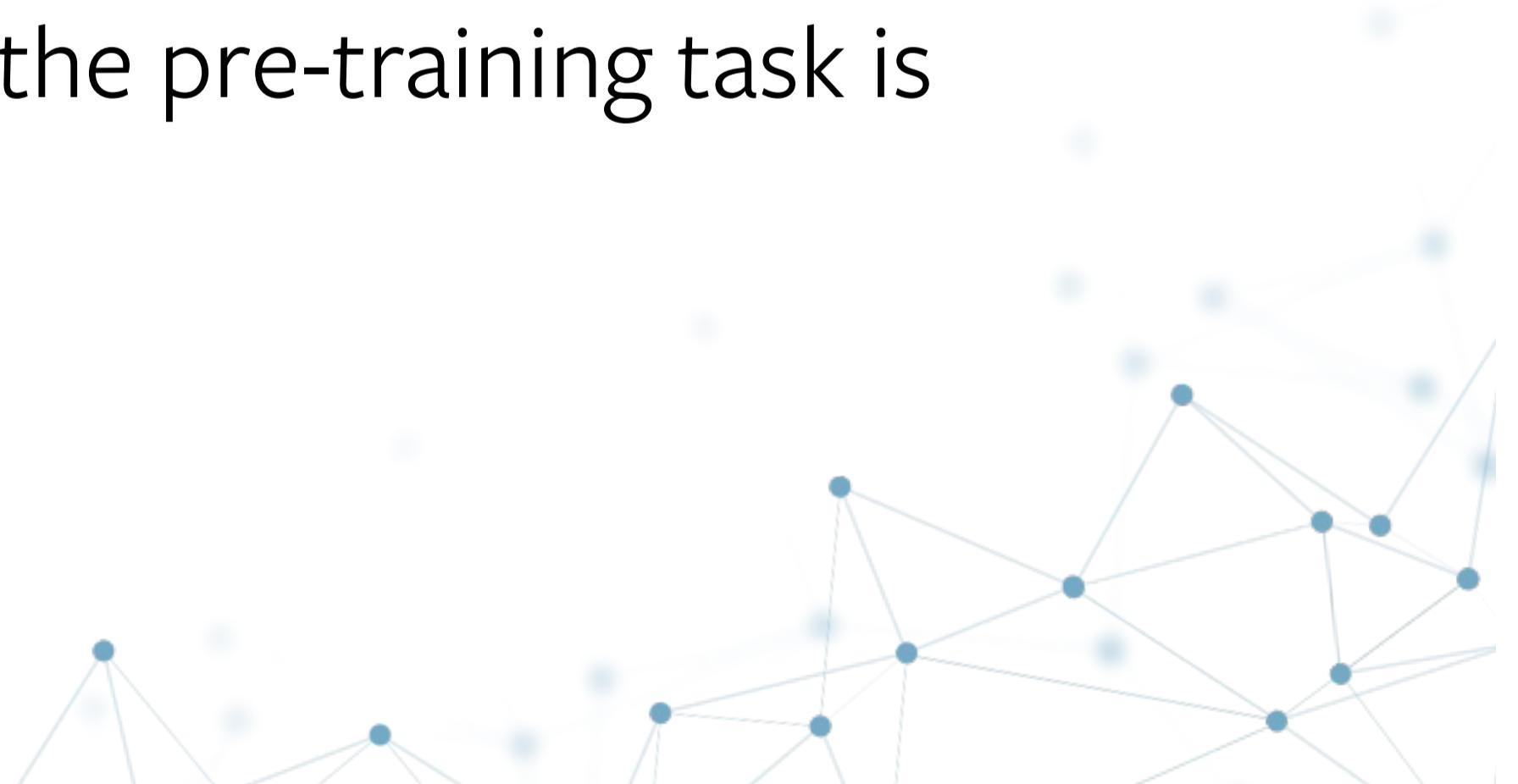
Pre-training (transfer learning)

Neural networks are typically initialized randomly.

Pre-training – an alternative to random initialization:

- (1) train a neural network for some time (e.g., classification on ImageNet, detection on COCO)
- (2) use this pre-trained network as an initialization for your task

This is often called **transfer learning**: information from the pre-training task is transferred to your downstream task.



Backbones

Architecture:

- CNNs such as AlexNet, VGG16, VGG19 etc
- ResNets such as ResNet34, ResNet50 etc
- Transformers such as ViT

Pretraining:

- On ImageNet for image classification
- Self-supervised learning, e.g. MAE

Modifications:

- Remove final global pooling layer and MLP (the part of the network that is designed for classification)
→ output spatial dims are proportional to input spatial dims



What we didn't talk about:

Two-stage vs single-stage detectors

Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

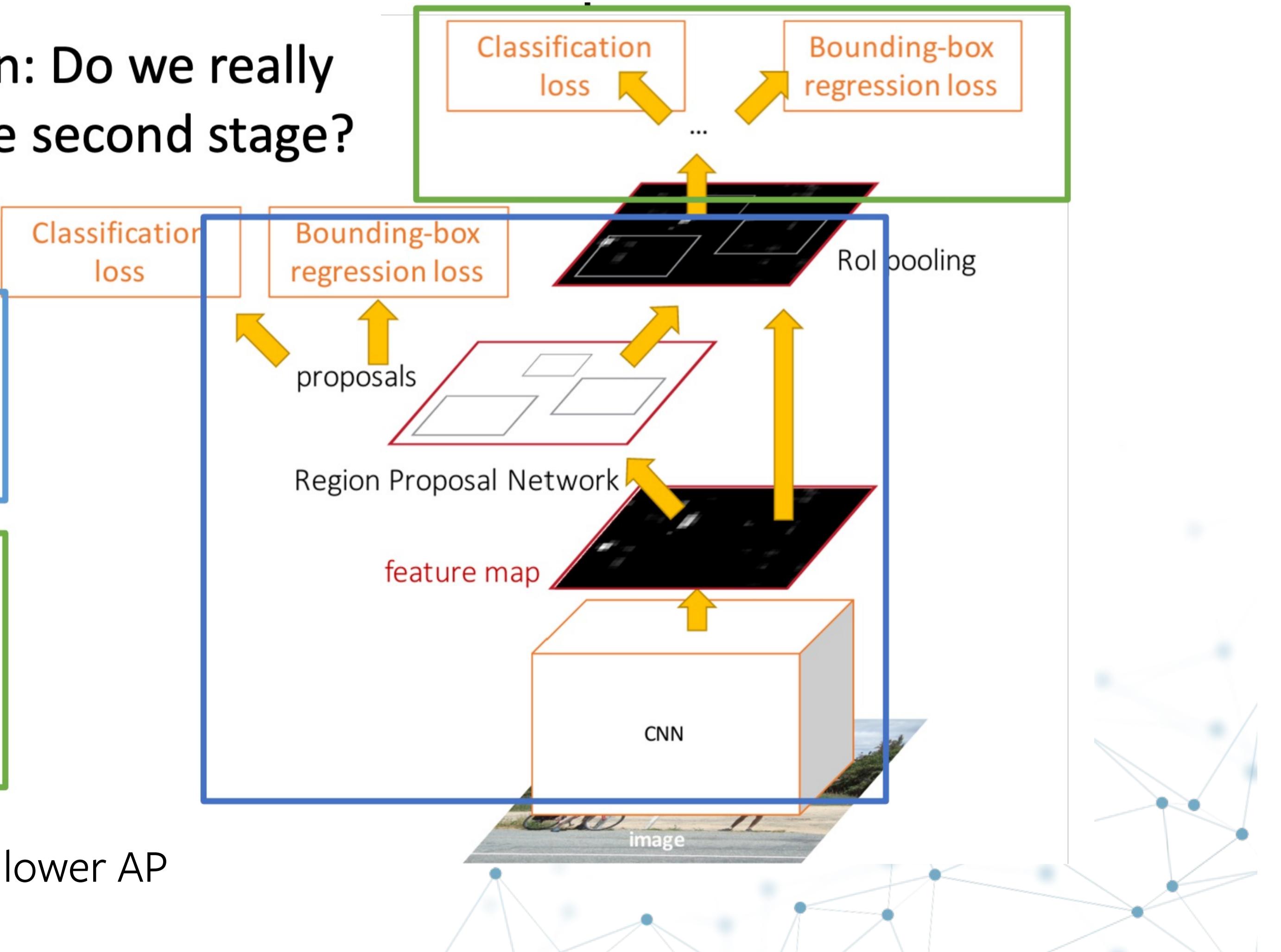
- Backbone network
- Region proposal network

Second stage: Run once per region

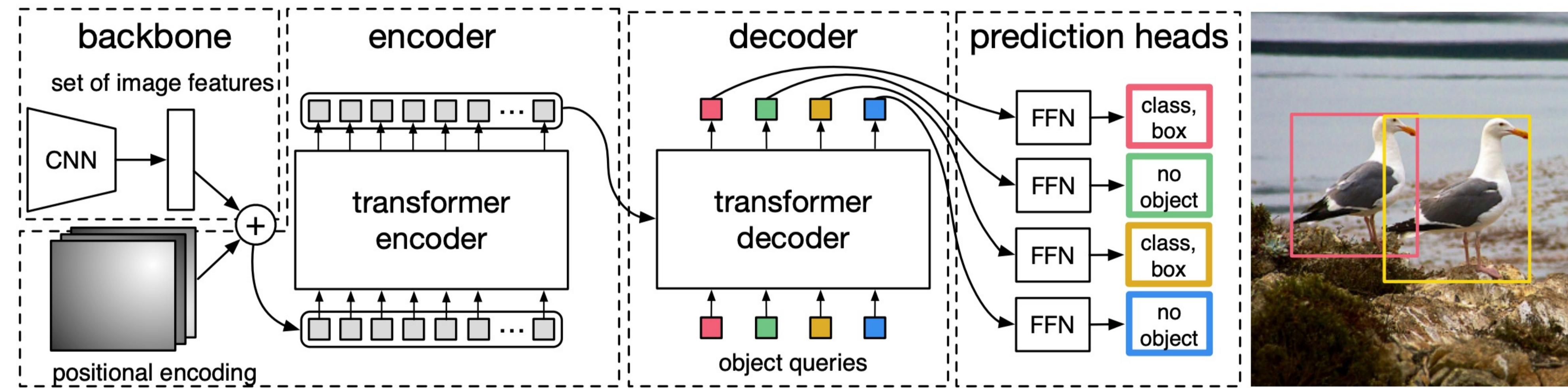
- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset

Single-stage: RetinaNet, SSD, FCOS – faster but lower AP

Question: Do we really
need the second stage?



What we didn't talk about: detection with transformers



Carion, Massa et al. End-to-End Object Detection with Transformers

Today's focus

- Principles of multi-object detection
- Region proposals
- Bounding box regression
- R-CNN training and inference
- Non-Max suppression
- Modern object detectors

Questions?