

ROB 498/599: Deep Learning for Robot Perception (DeepRob)

Lecture 12: Object Detection - part 2

02/18/2026



<https://deeprob.org/w26/>

Today

- Feedback and Recap (5min)
- Object Detection (part 2)
 - Fast R-CNN (25min)
 - Receptive Fields
 - ROI Align
 - ROI Pool
 - Faster R-CNN (25min)
 - Region Proposal Network
 - Mask R-CNN (6min)
 - YOLO (6min)
 - Mesh R-CNN (6min)
- Summary and Takeaways (5min)

Recap: Object Detection

- P3 released, Due March 8, 2026
Start NOW!!!
- ~~Also, Object Detection Quiz due Feb 22~~

Classification



"Chocolate Pretzels"

No spatial extent

Semantic Segmentation



Chocolate Pretzels,
Shelf

No objects, just pixels

Object Detection



Flipz, Hershey's, Keese's

Multiple objects

Instance Segmentation

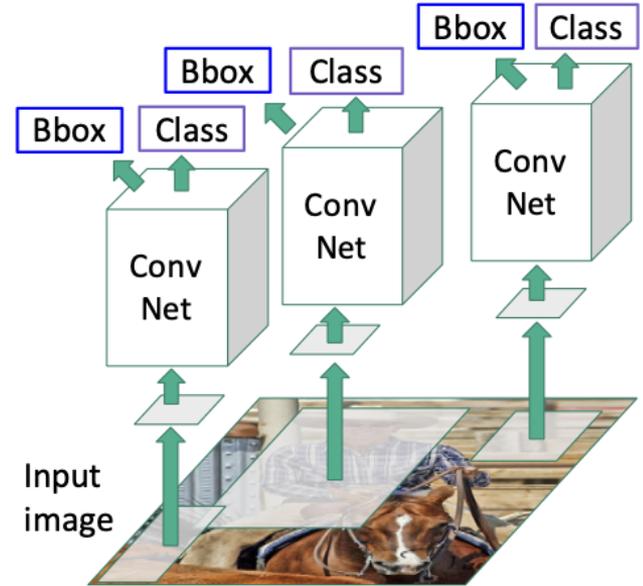


Recap: Slow vs. Fast R-CNN



Input image

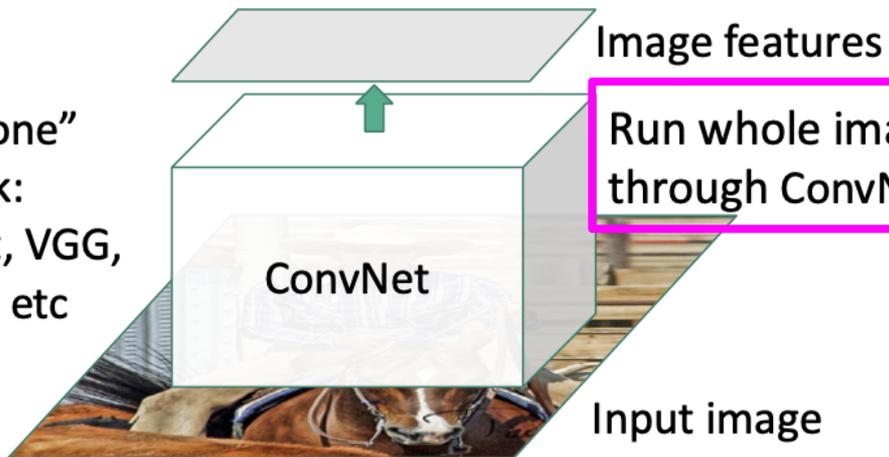
“Slow” R-CNN
Process each region
independently



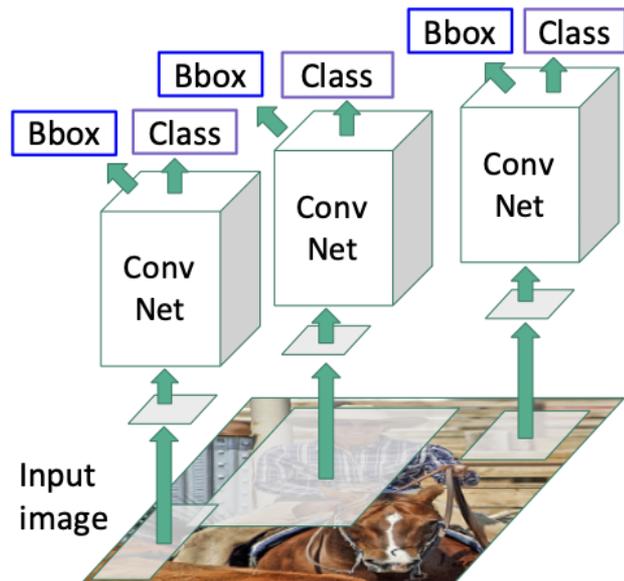
Fast R-CNN

“Heavy duty backbone”

“Backbone”
network:
AlexNet, VGG,
ResNet, etc



“Slow” R-CNN
Process each region
independently

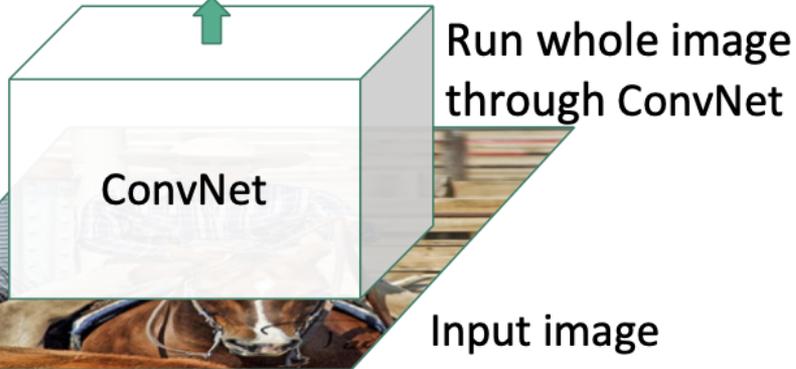


Fast R-CNN

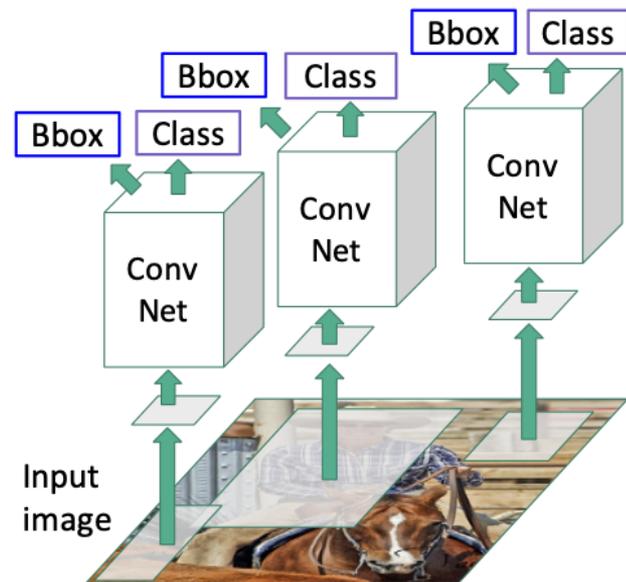
Regions of Interest (RoIs) from a proposal method



“Backbone” network:
AlexNet, VGG,
ResNet, etc



“Slow” R-CNN
Process each region independently

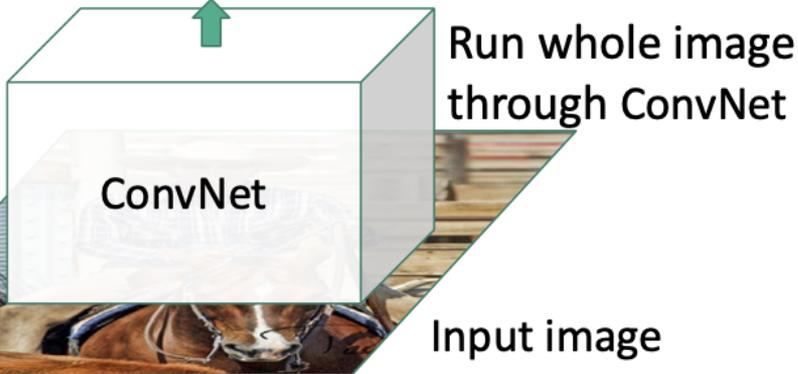


Fast R-CNN

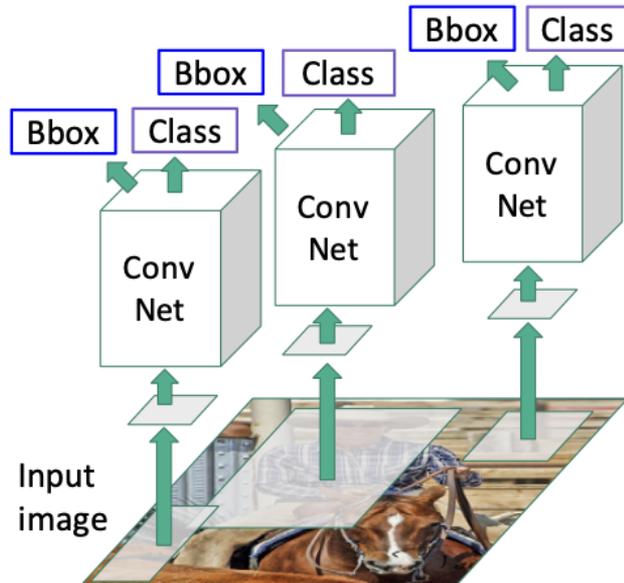
Regions of Interest (Rois) from a proposal method



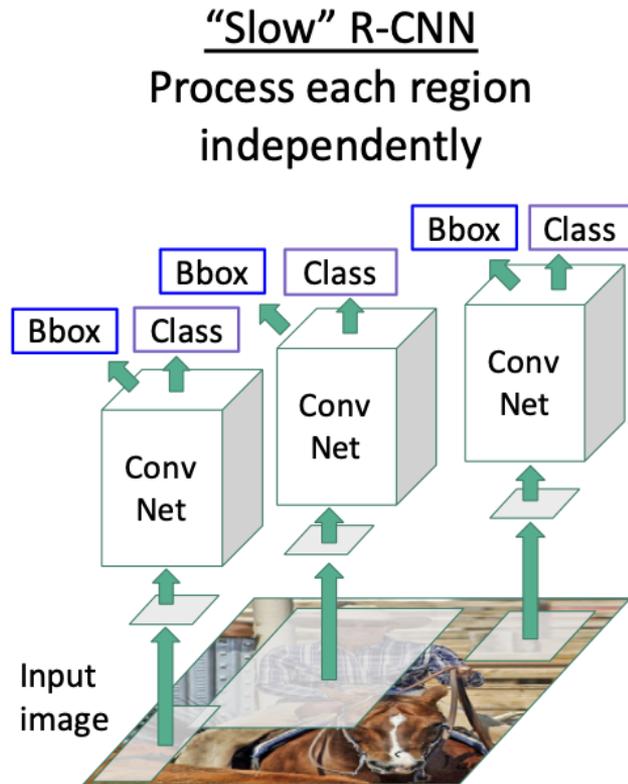
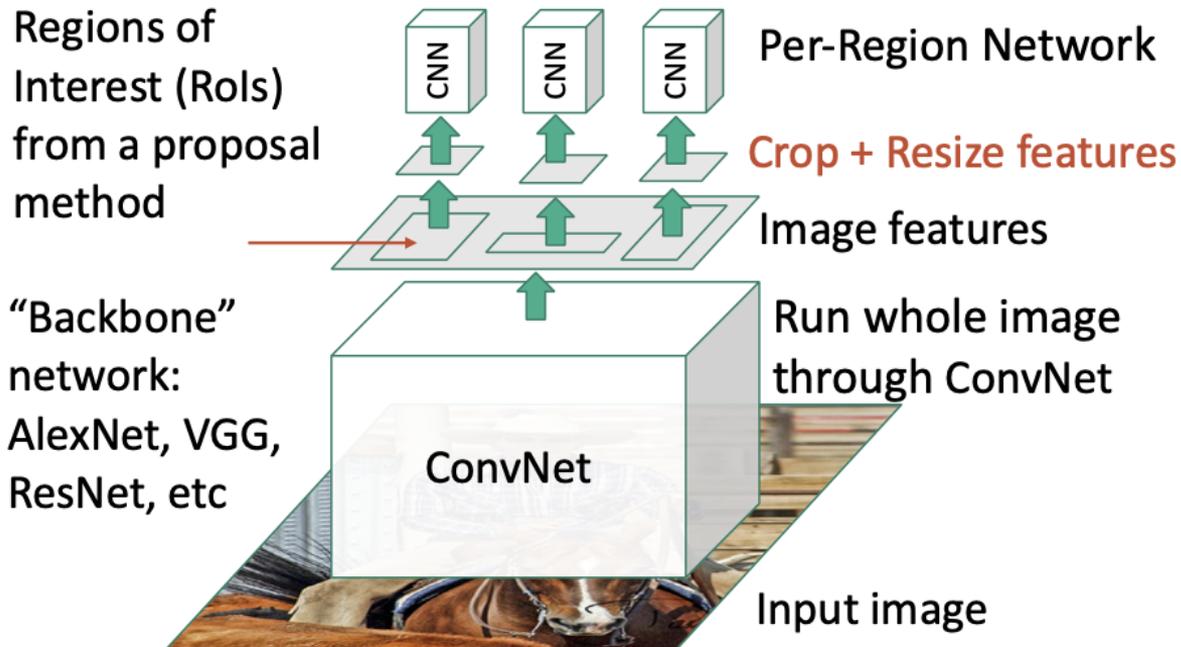
“Backbone” network:
AlexNet, VGG,
ResNet, etc



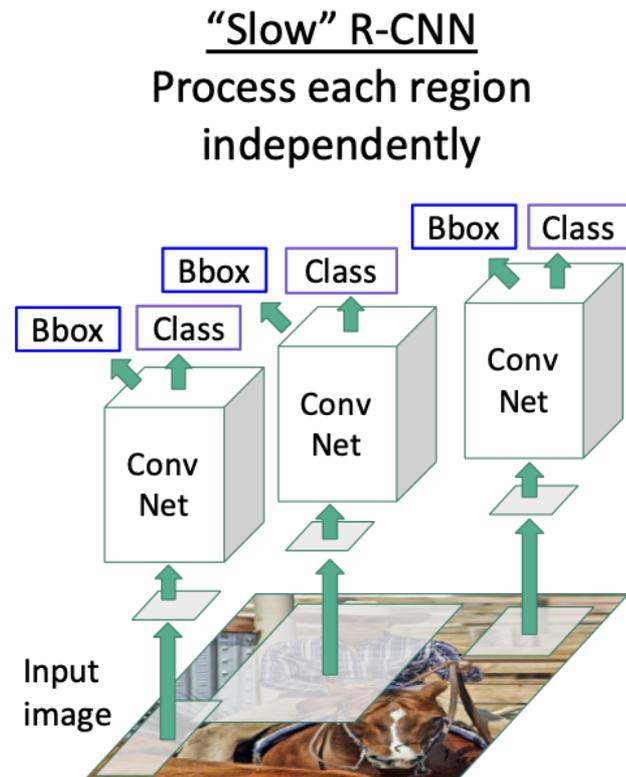
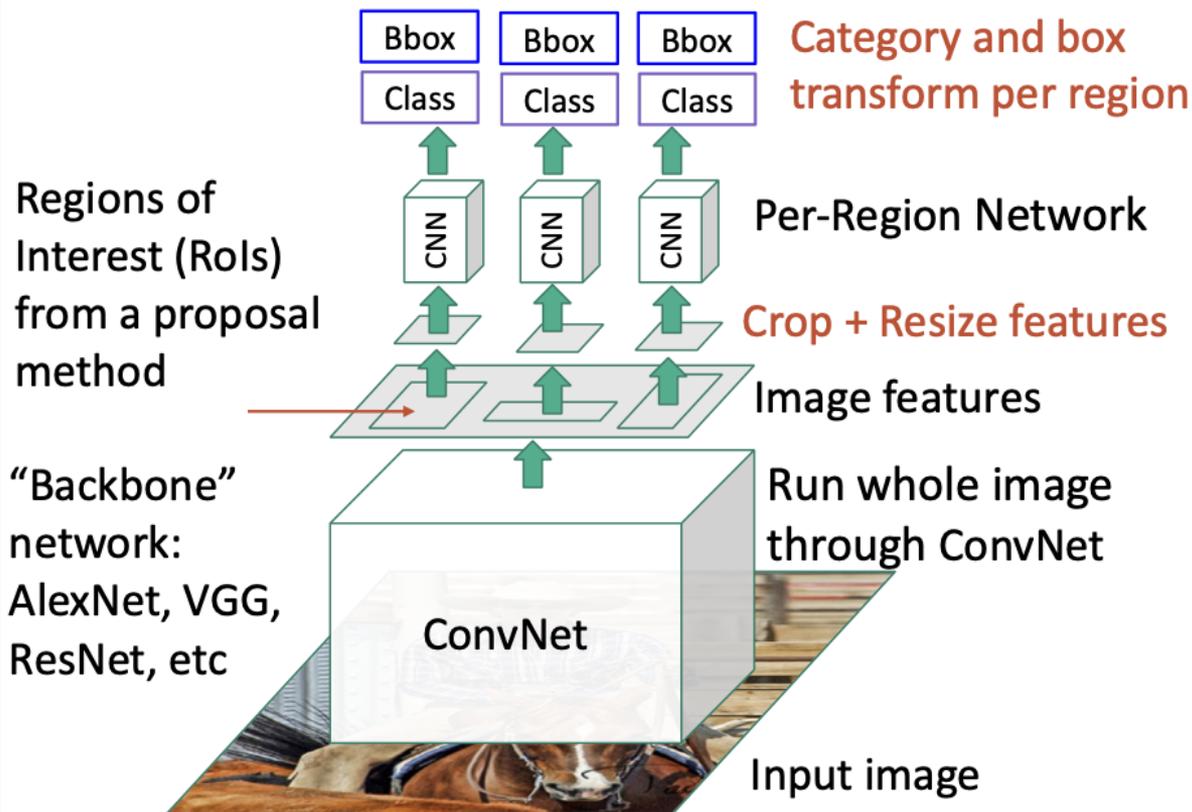
“Slow” R-CNN
Process each region independently



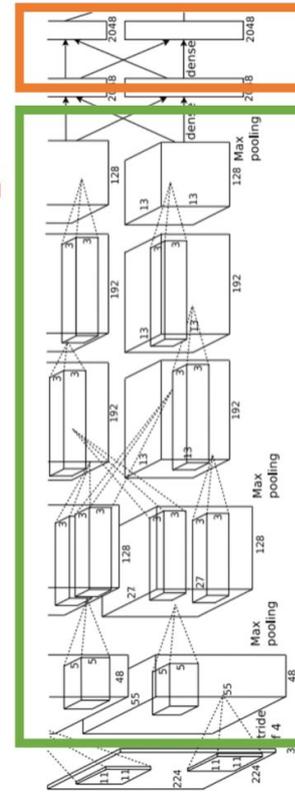
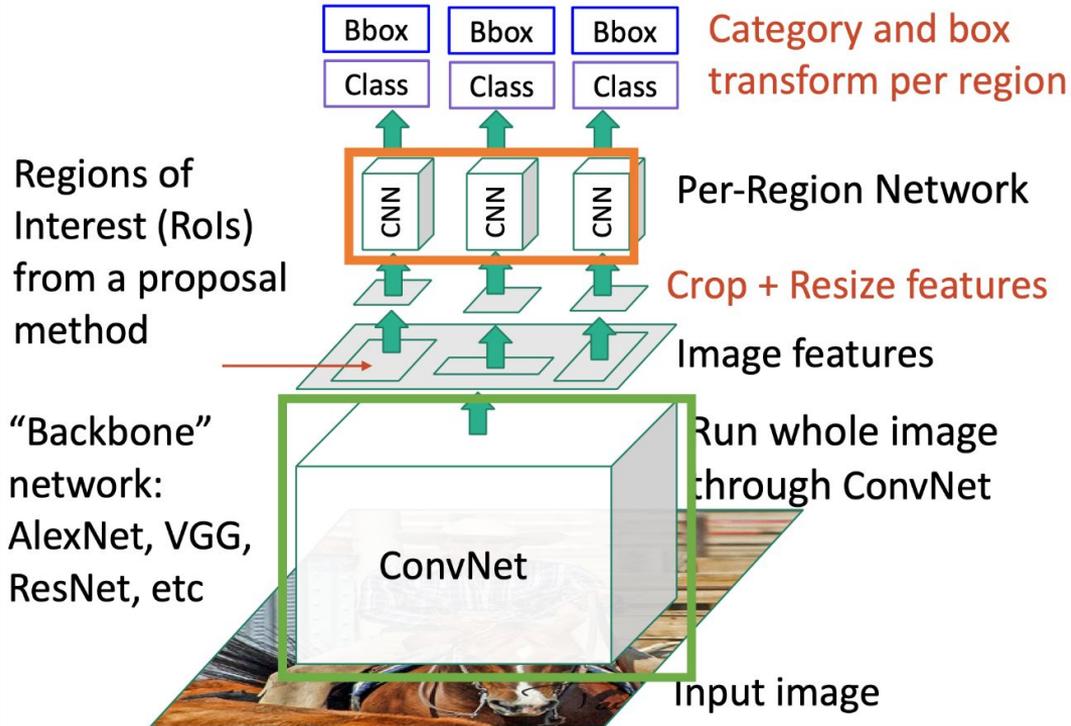
Fast R-CNN



Fast R-CNN

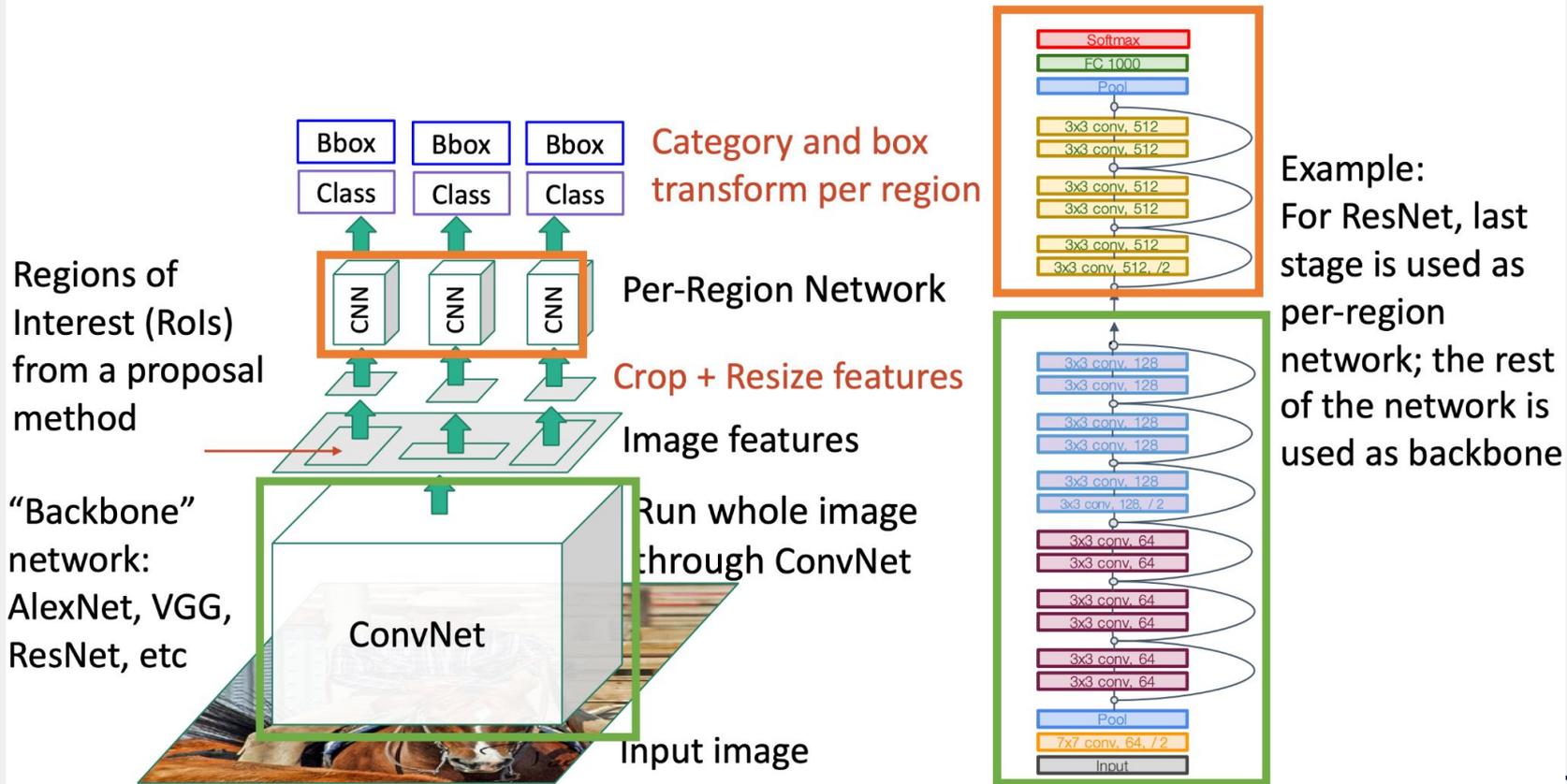


Fast R-CNN

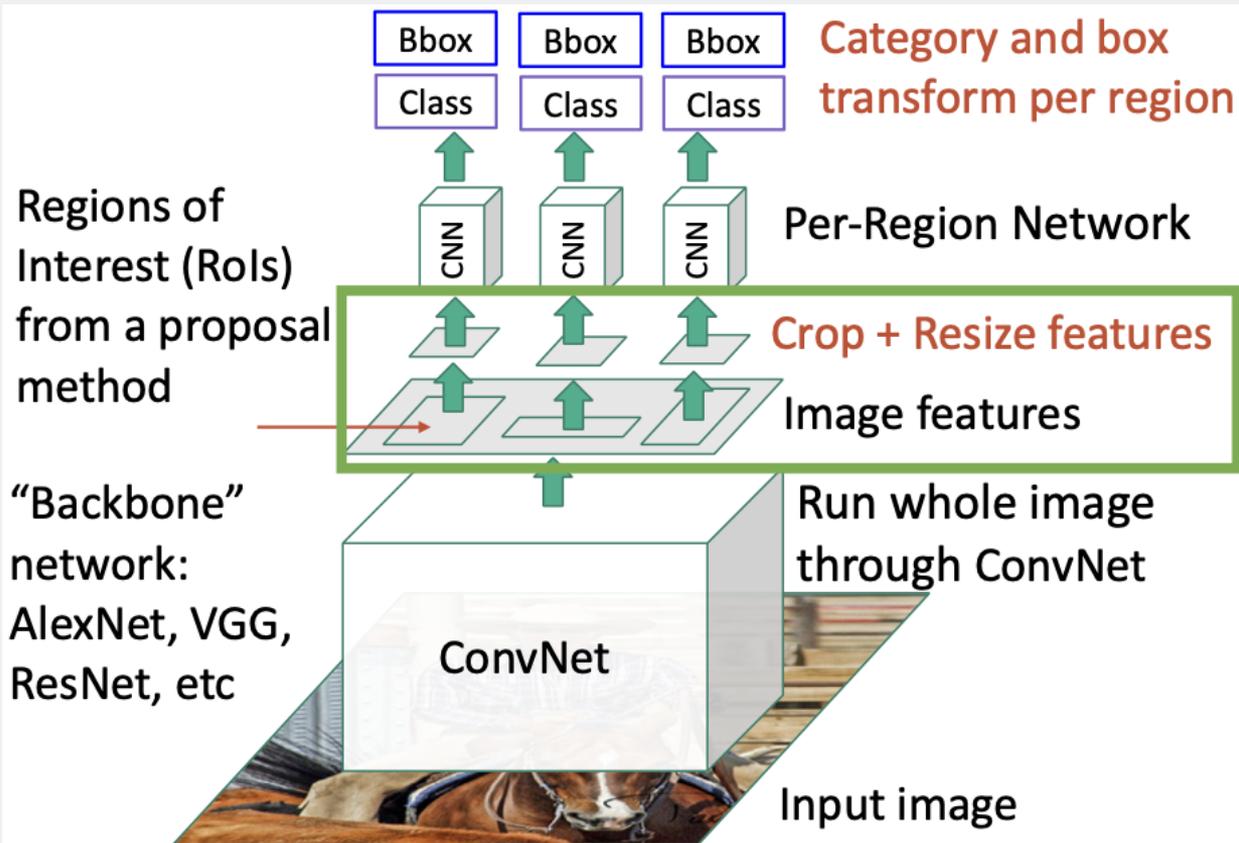


Example:
When using AlexNet for detection, five conv layers are used for backbone and two FC layers are used for per-region network

Fast R-CNN

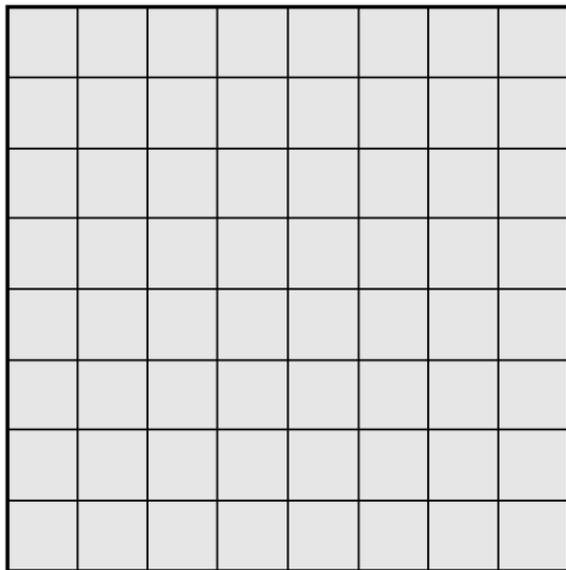


Fast R-CNN



How to crop features?

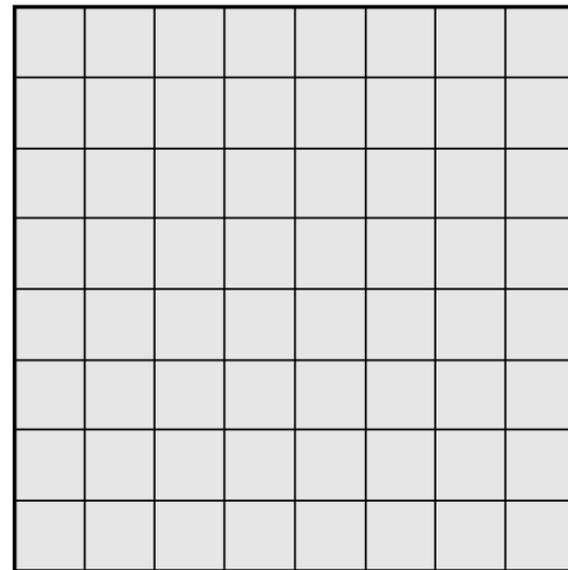
Recall: Receptive Fields



Input Image: 8 x 8

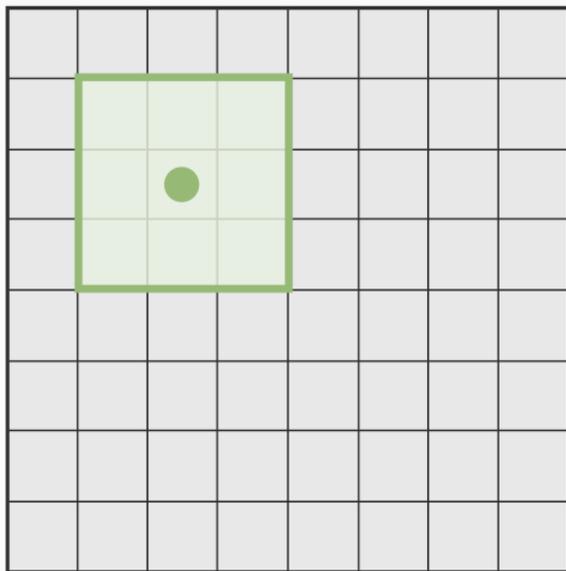
Every position in the
output feature map
depends on a ???
receptive field in the input

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

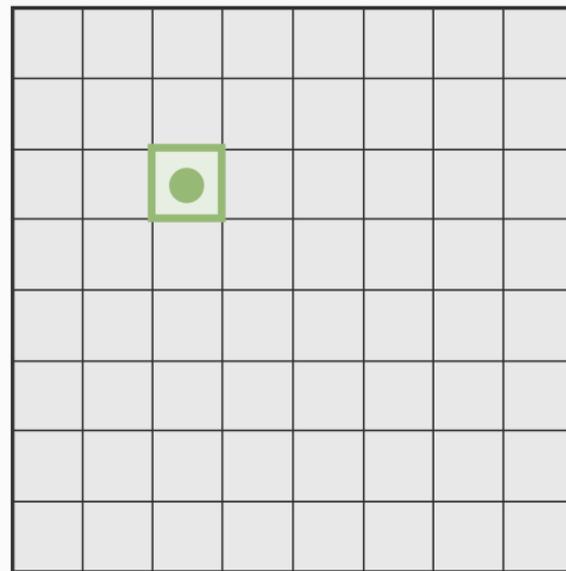
Recall: Receptive Fields



Input Image: 8 x 8

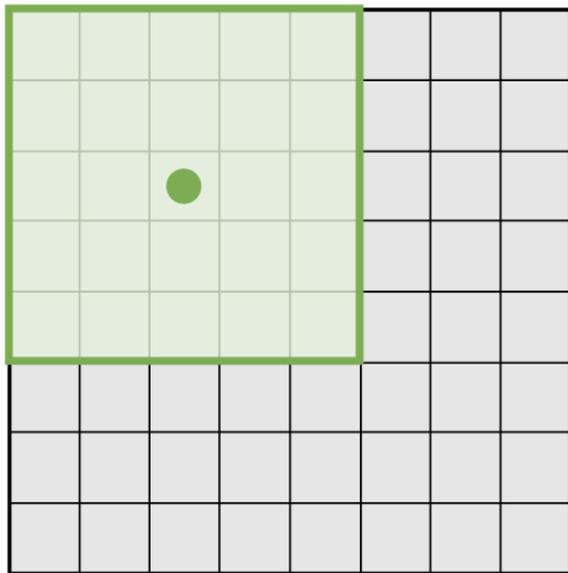
Every position in the output feature map depends on a 3x3 receptive field in the input

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

Recall: Receptive Fields

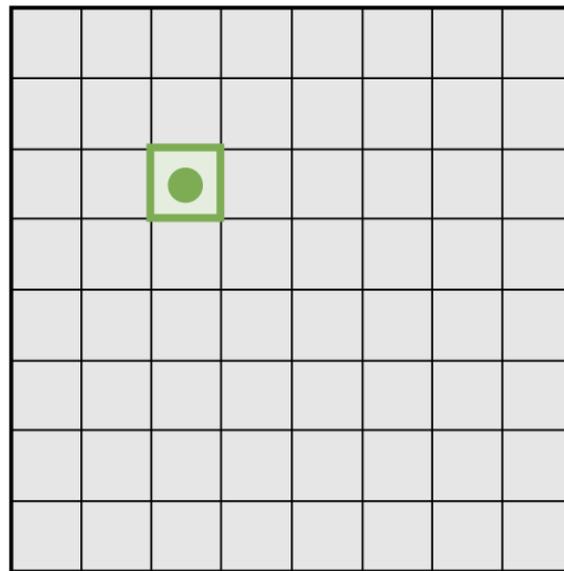


Input Image: 8 x 8

Every position in the output feature map depends on a 5x5 receptive field in the input

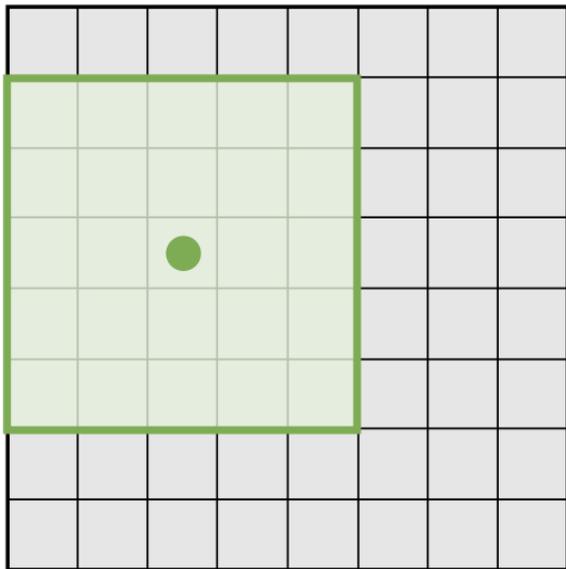
3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8

Recall: Receptive Fields

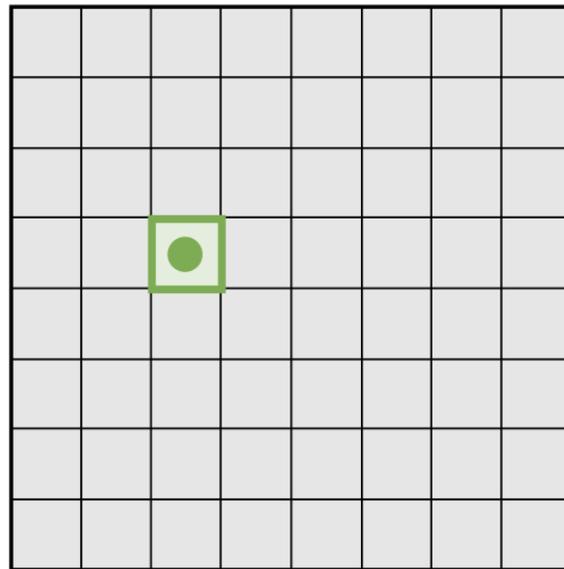


Input Image: 8 x 8

Moving one unit in the output space also moves the receptive field by one

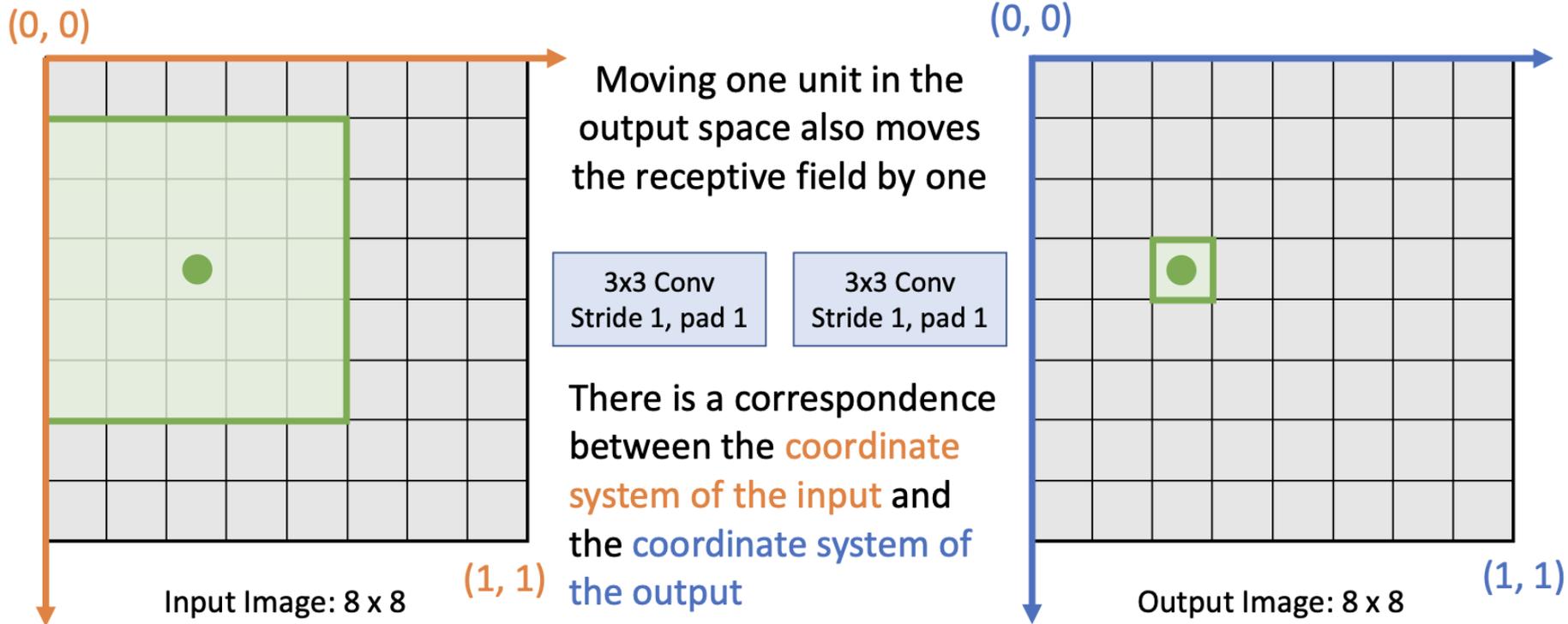
3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1



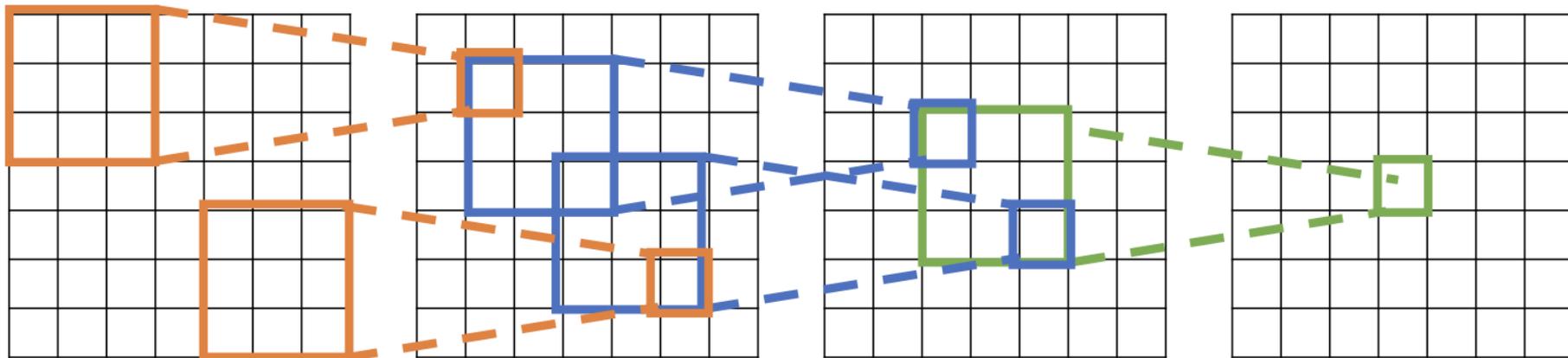
Output Image: 8 x 8

Recall: Receptive Fields



Recall: Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



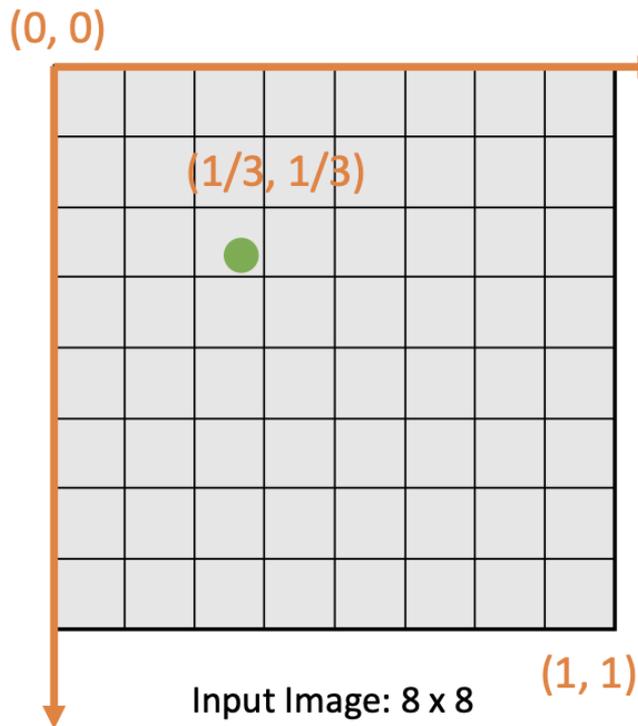
Input

Problem: For large images we need many layers for each output to “see” the whole image image

Solution: Downsample inside the network

Output

Projecting Points

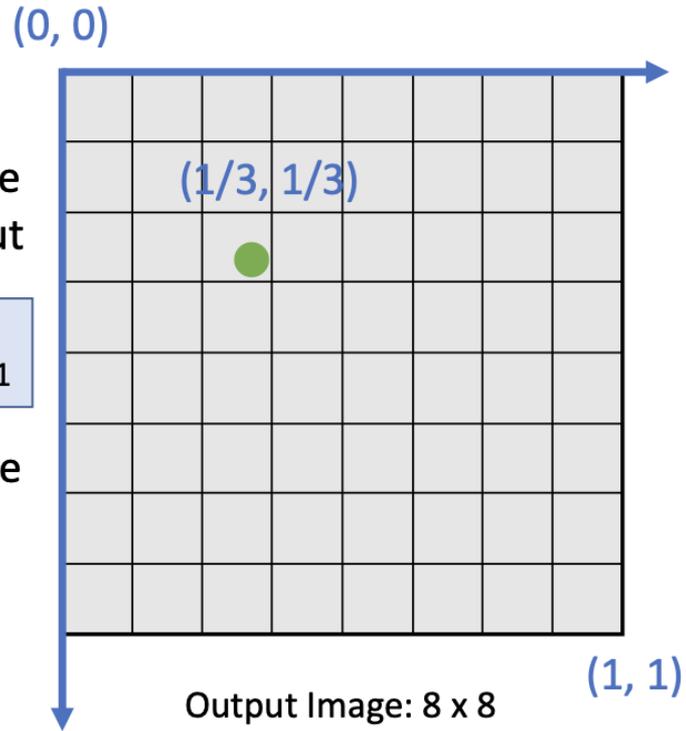


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

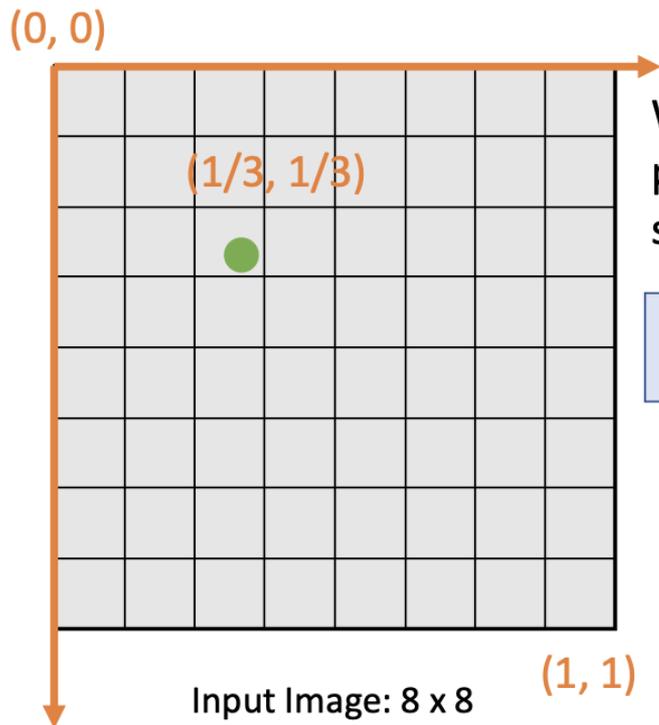
3x3 Conv
Stride 1, pad 1

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Projecting Points

Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different

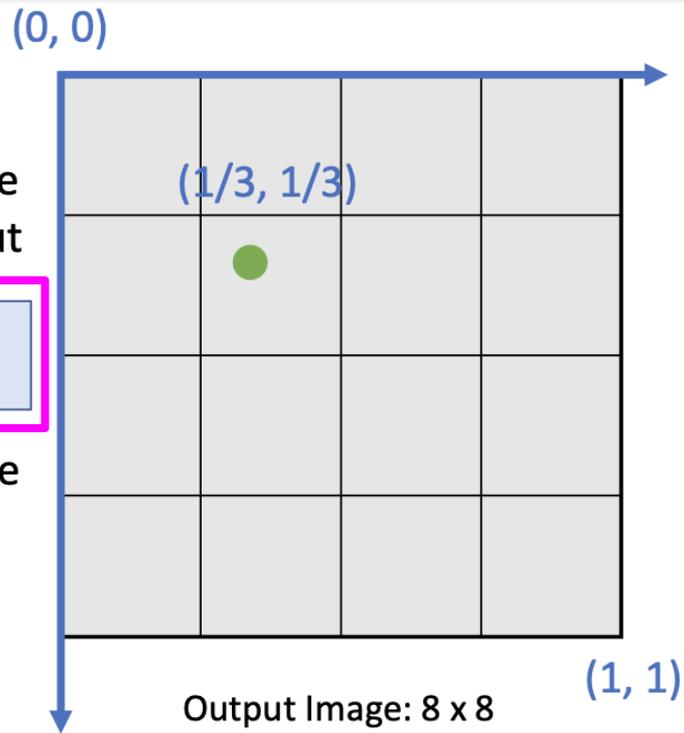


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

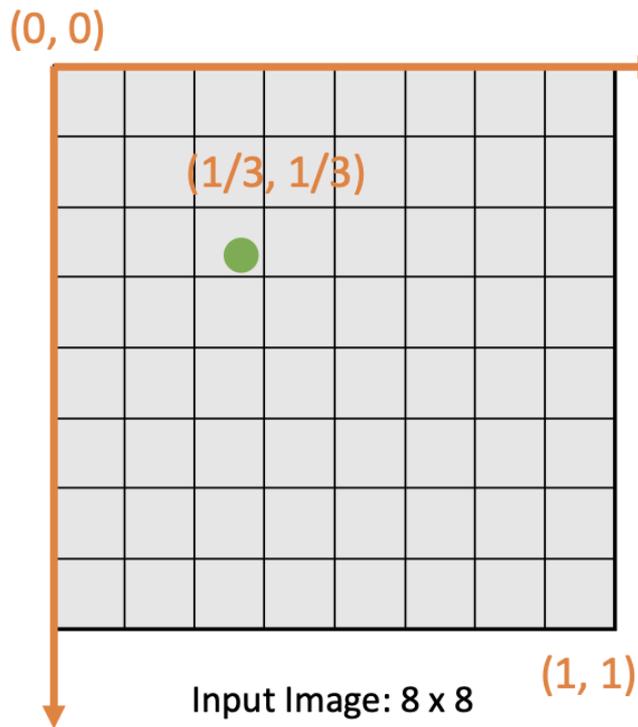
2x2 MaxPool
Stride 2

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Projecting Points

Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different



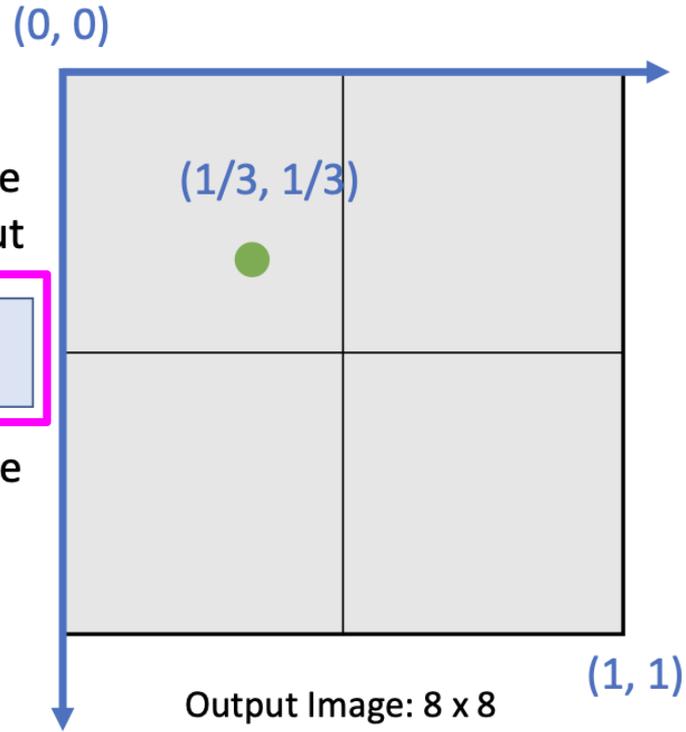
Input Image: 8 x 8

We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

4x4 MaxPool
Stride 4

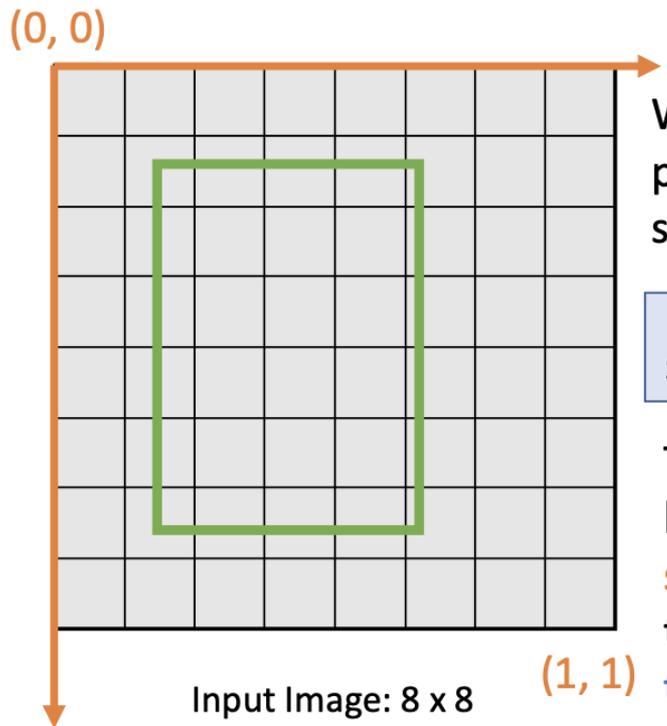
There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Output Image: 8 x 8

Projecting Points

We can use this idea to project **bounding boxes** between an input image and a feature map

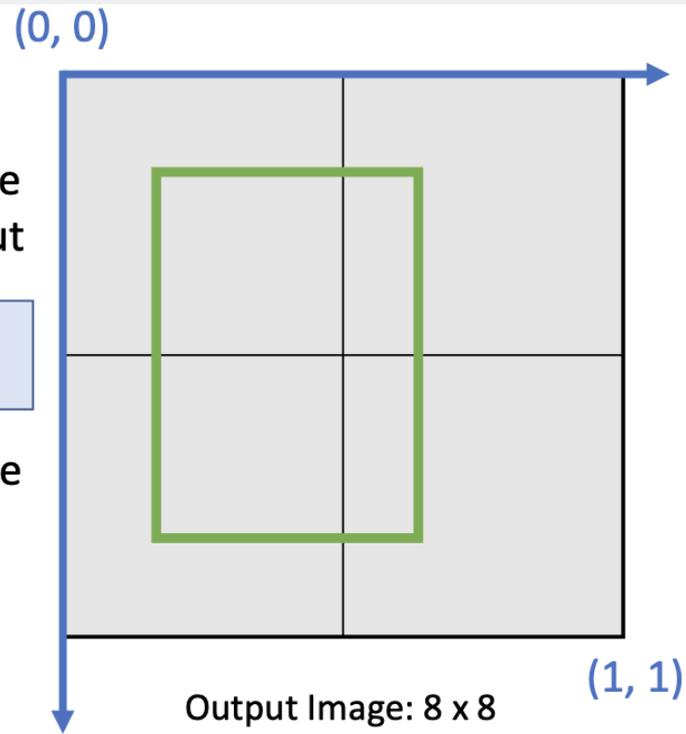


We can align arbitrary points between coordinate system of input and output

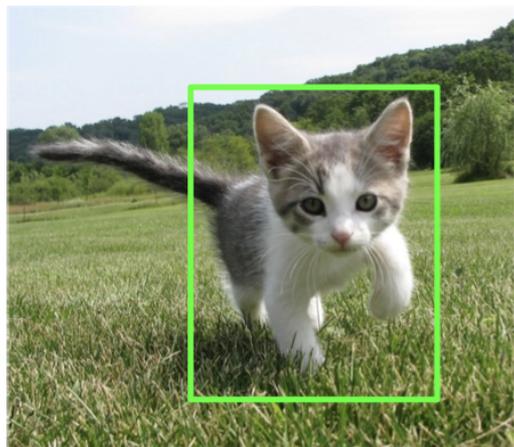
3x3 Conv
Stride 1, pad 1

4x4 MaxPool
Stride 4

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**



Cropping Features: RoI Pool



Input Image
(e.g. 3 x 640 x 480)

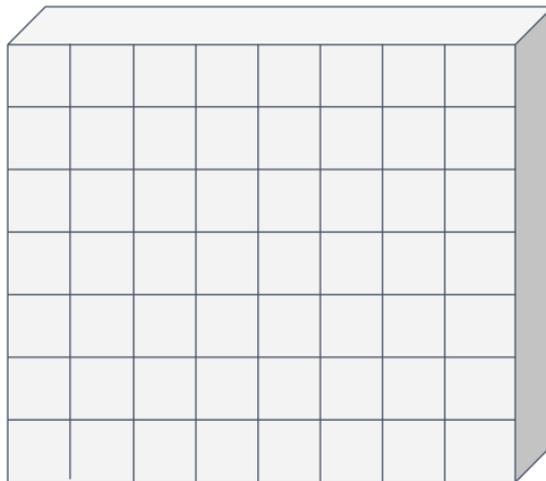
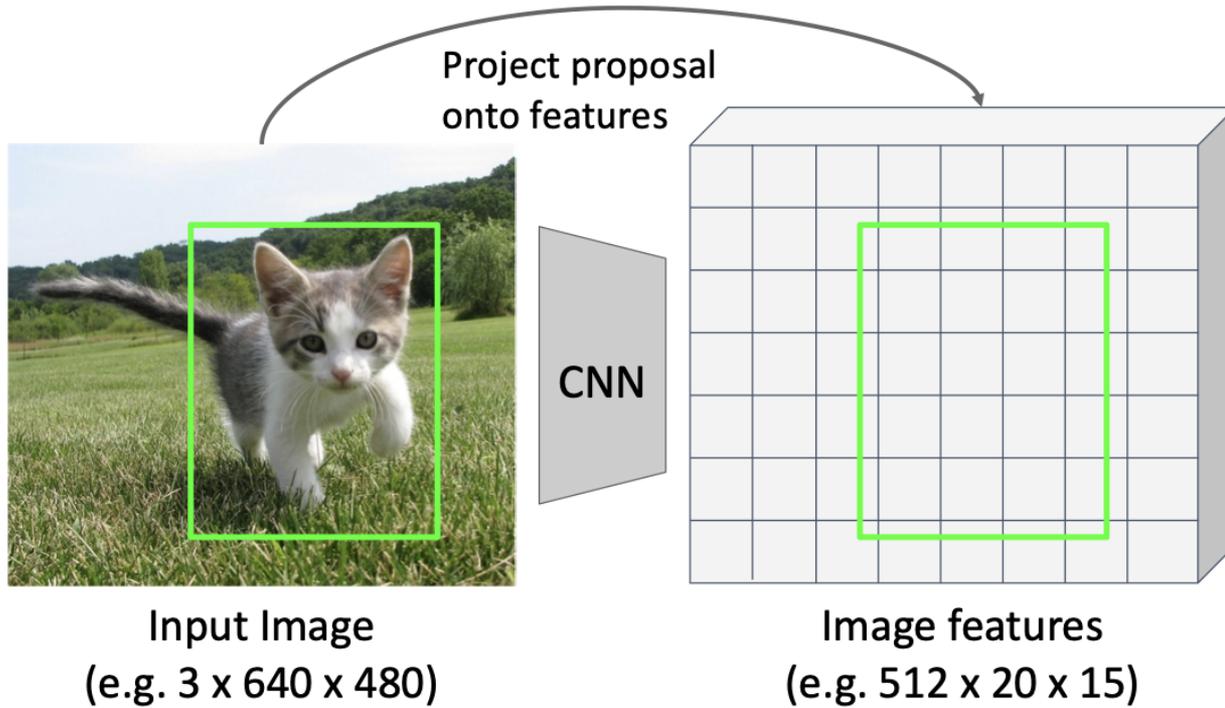


Image features
(e.g. 512 x 20 x 15)

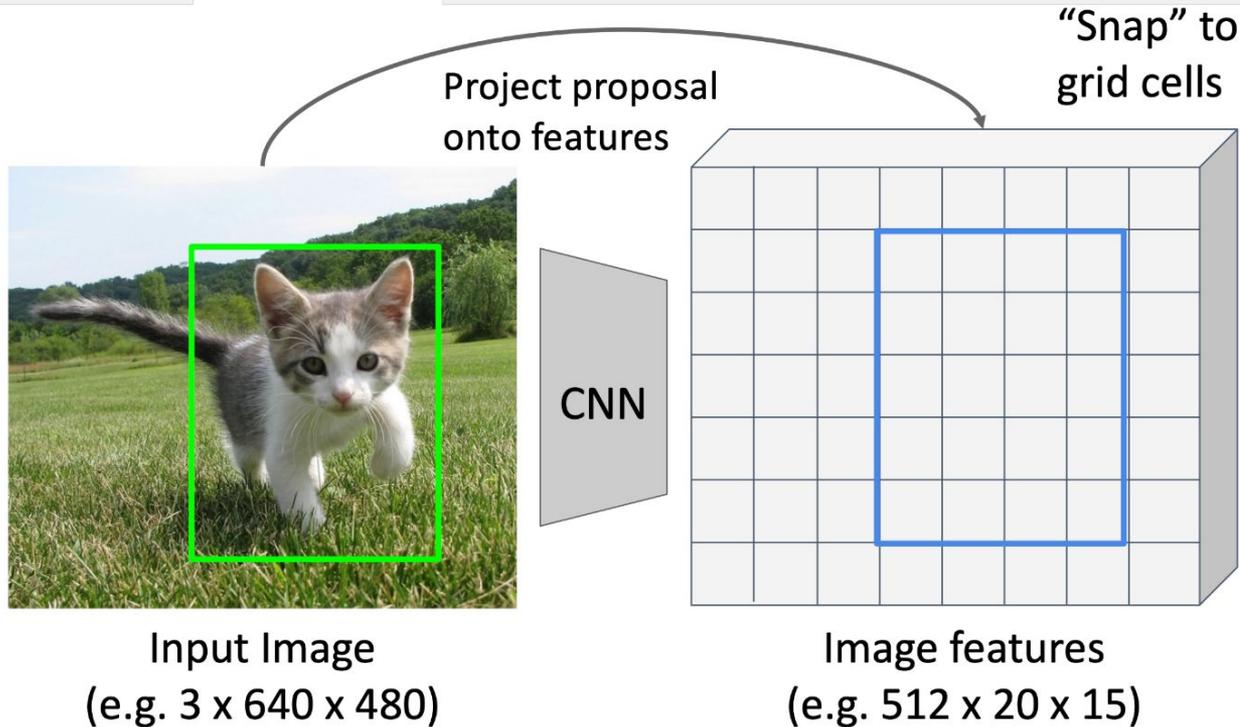
Want features for the
box of a fixed size
(2x2 in this example,
7x7 or 14x14 in practice)

Cropping Features: RoI Pool



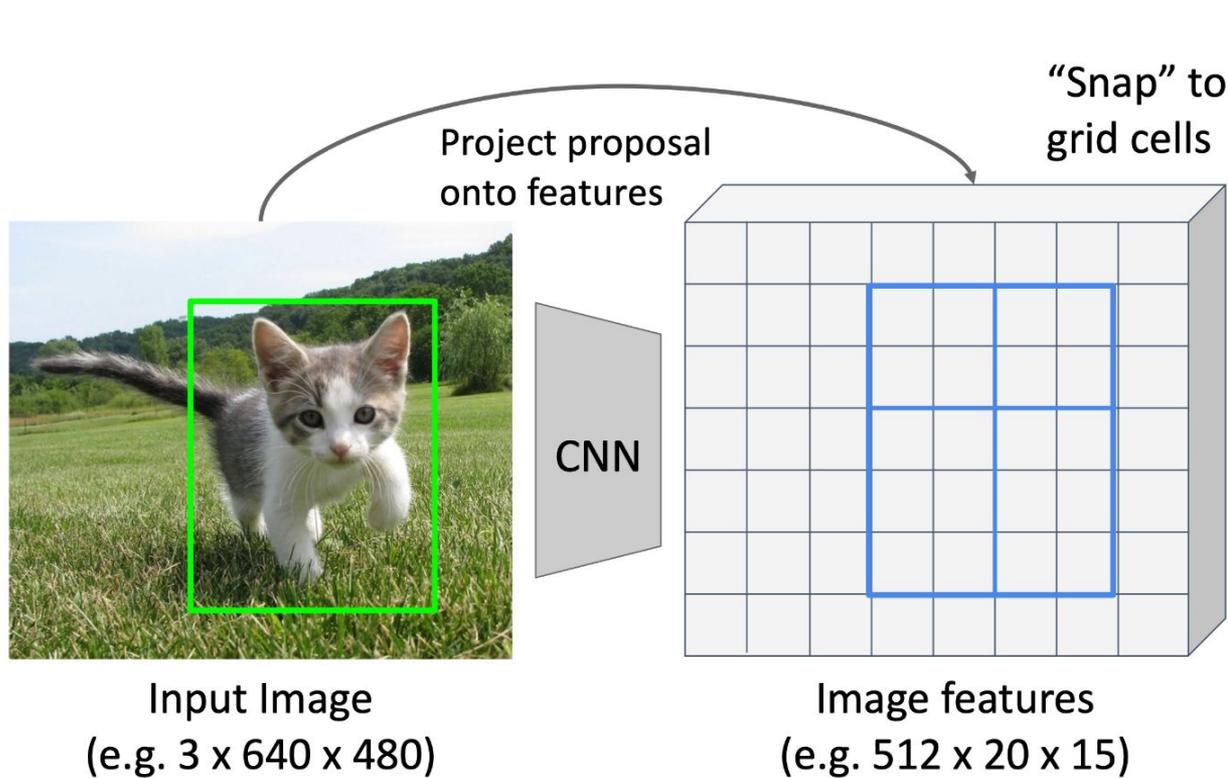
Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

Cropping Features: RoI Pool



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

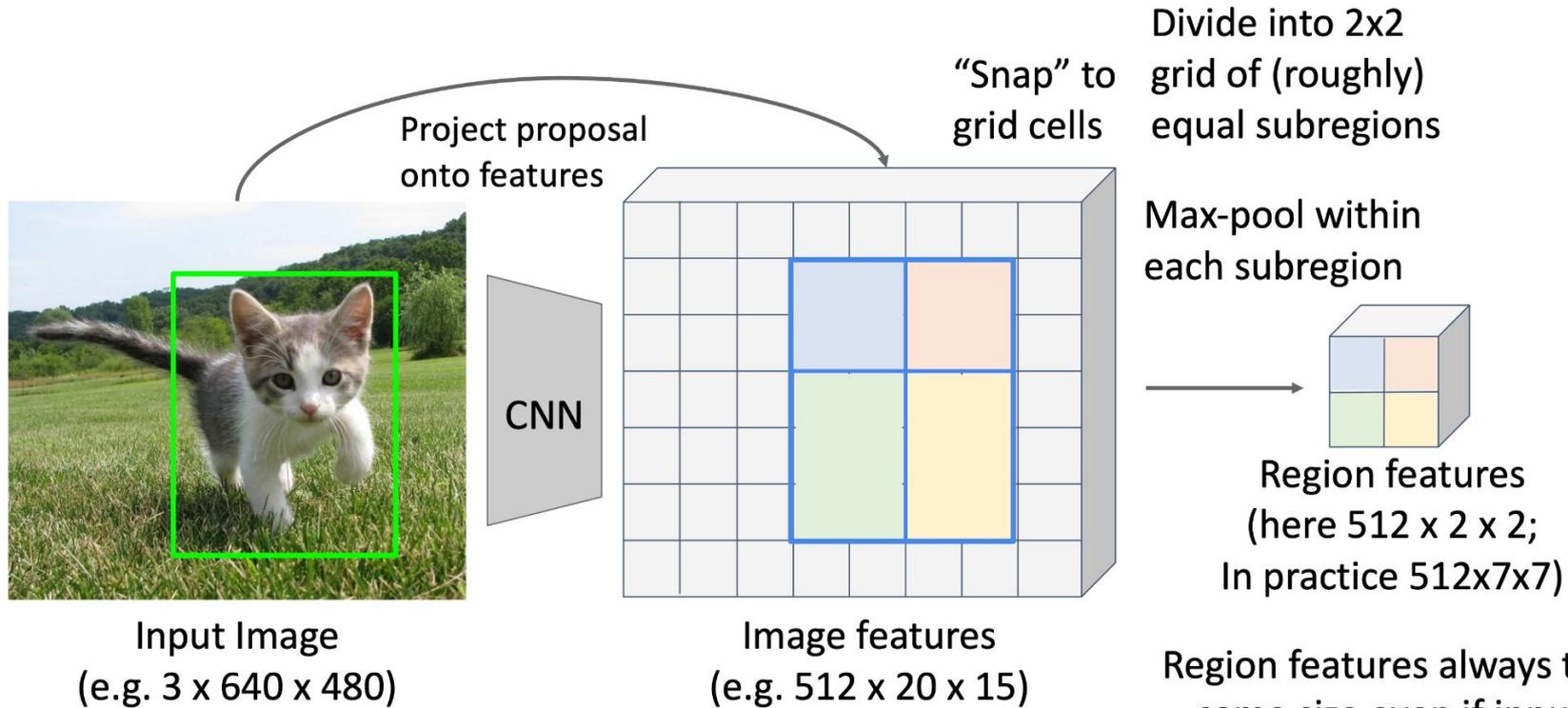
Cropping Features: RoI Pool



Divide into 2x2
grid of (roughly)
equal subregions

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

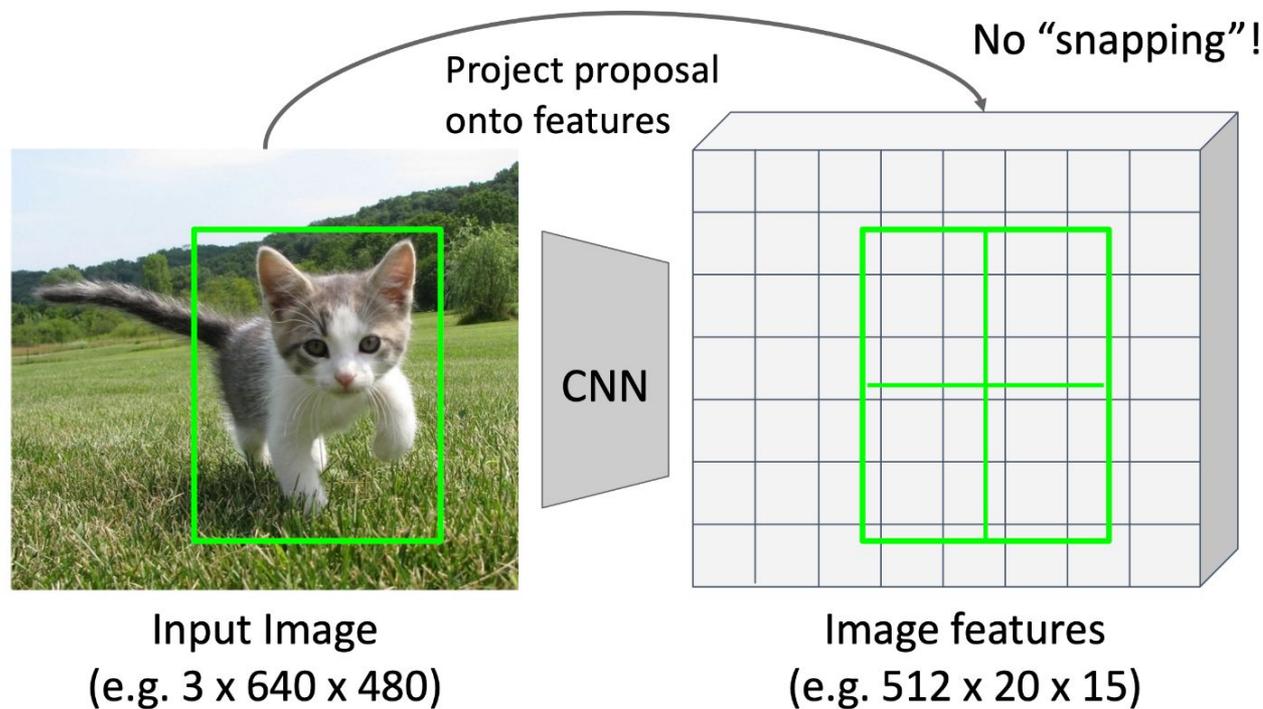
Cropping Features: RoI Pool



Problem: Slight misalignment due to snapping; different-sized subregions is weird

Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



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

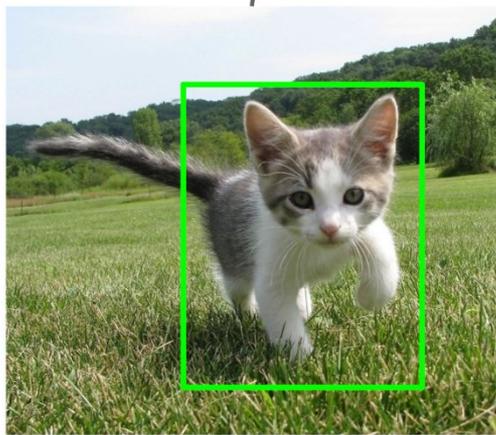
Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)

Project proposal
onto features

No “snapping”!

Sample features at
regularly-spaced points
in each subregion using
bilinear interpolation



Input Image
(e.g. 3 x 640 x 480)

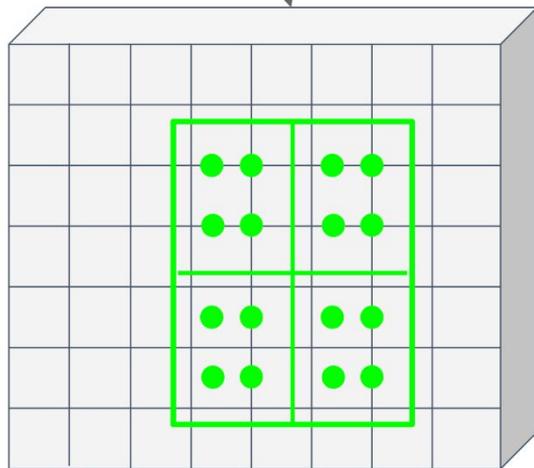
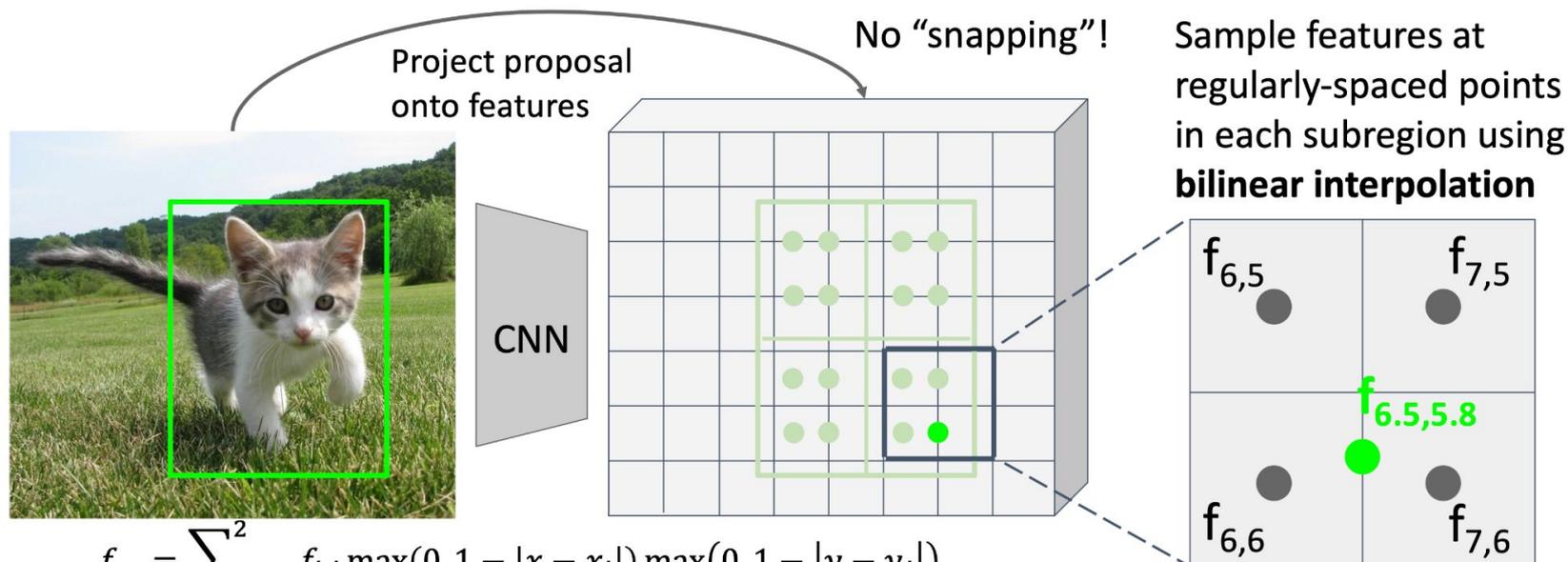


Image features
(e.g. 512 x 20 x 15)

Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

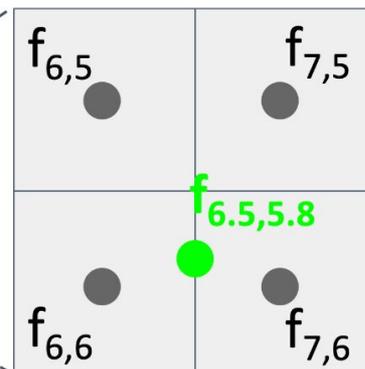
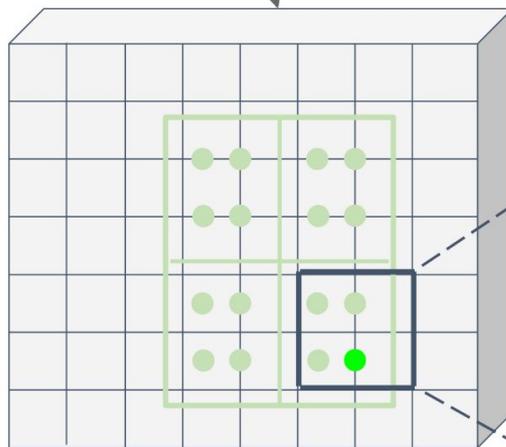
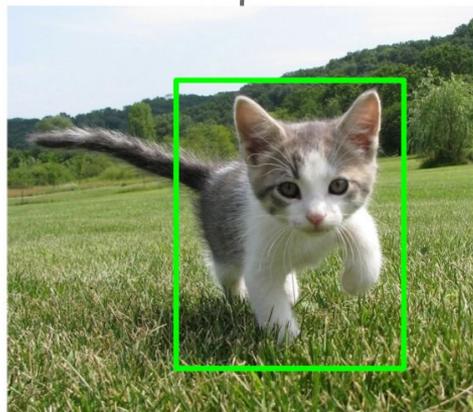
Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)

Project proposal
onto features

No “snapping”!

Sample features at
regularly-spaced points
in each subregion using
bilinear interpolation

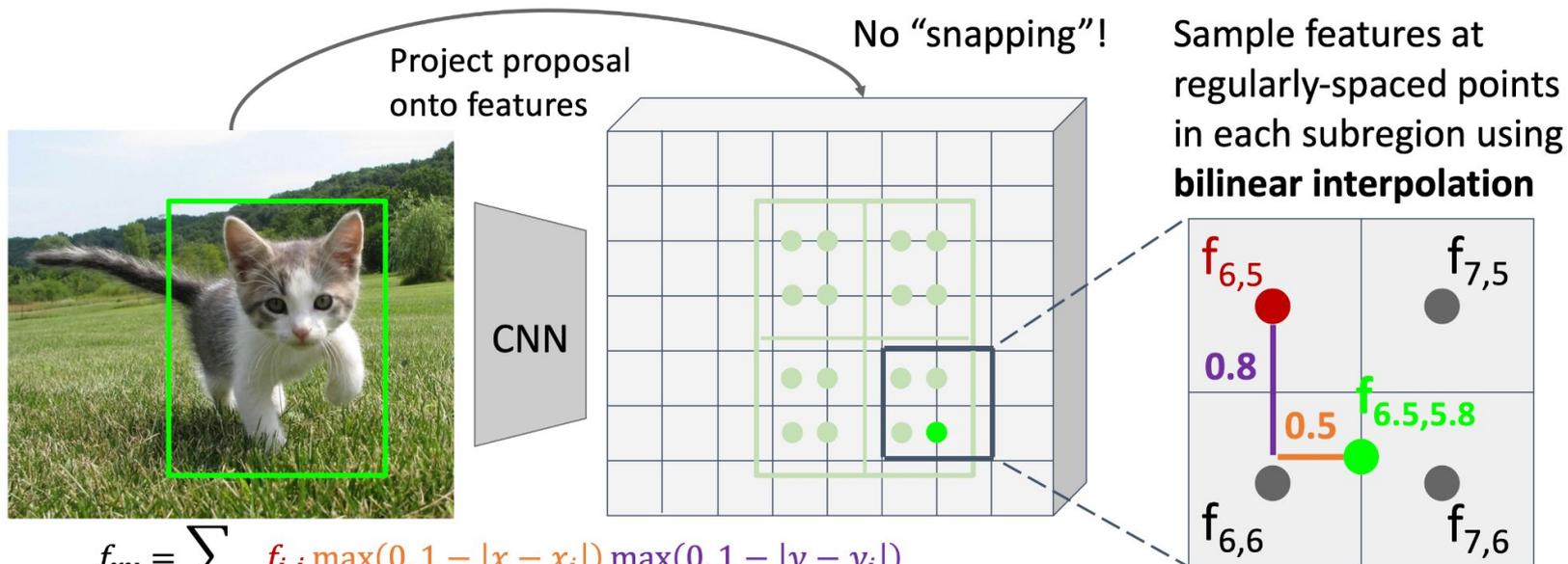


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &\quad + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

Feature f_{xy} for point (x, y) is a
linear combination of features
at its four neighboring grid cells:

Cropping Features: RoI Align

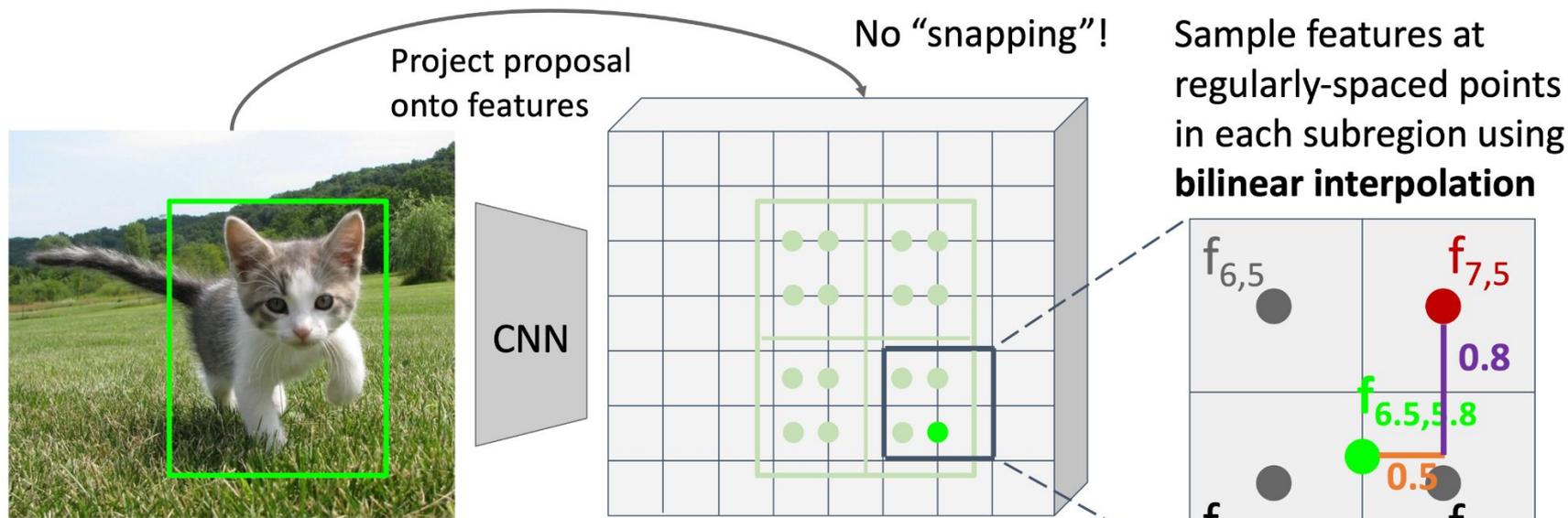


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$\begin{aligned} f_{6.5,5.8} &= (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) \\ &+ (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8) \end{aligned}$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align

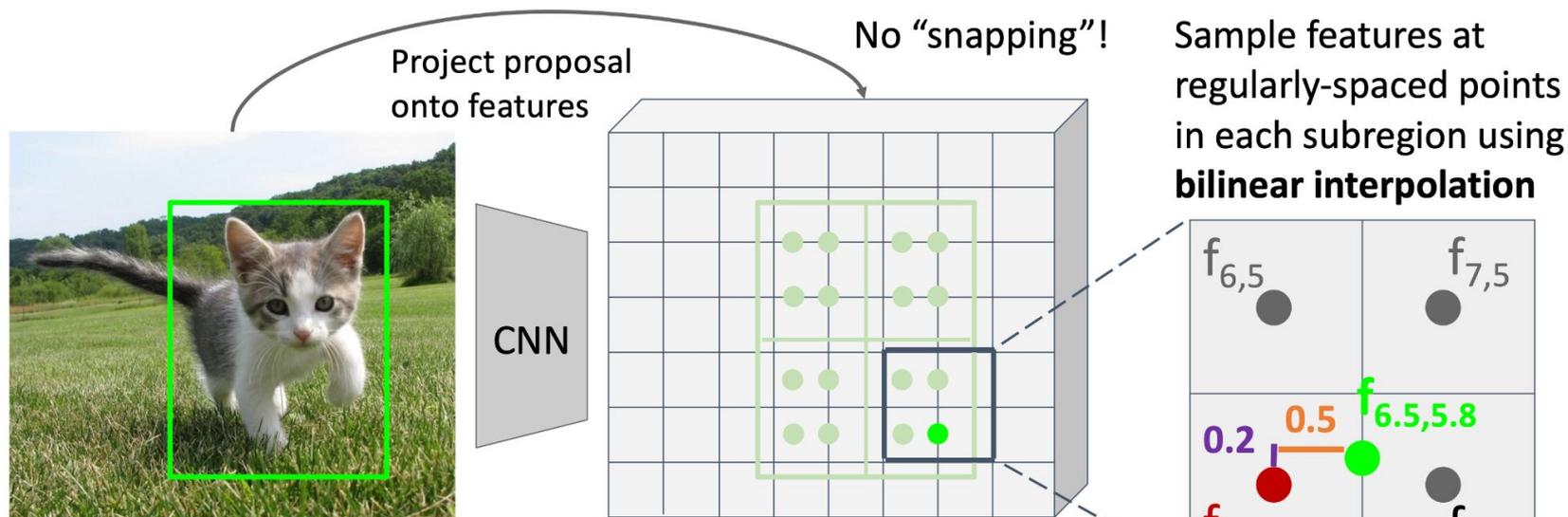


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align

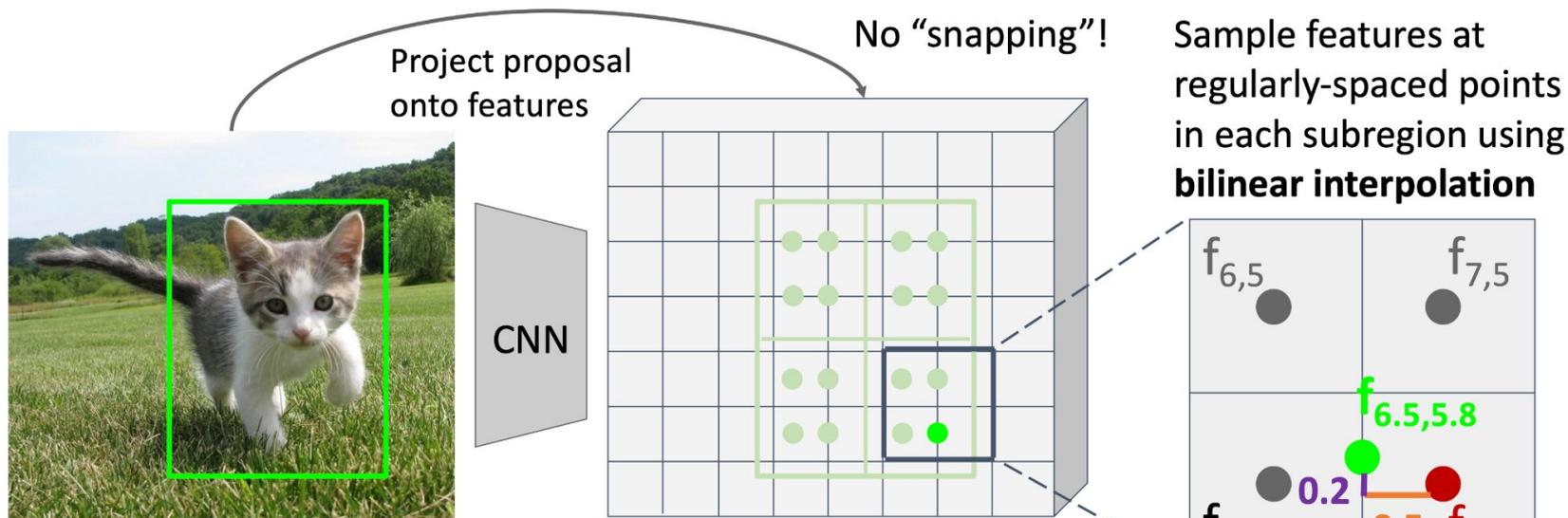


$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

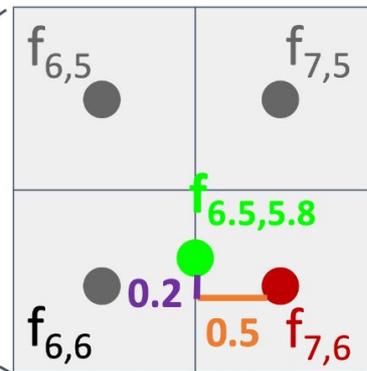
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



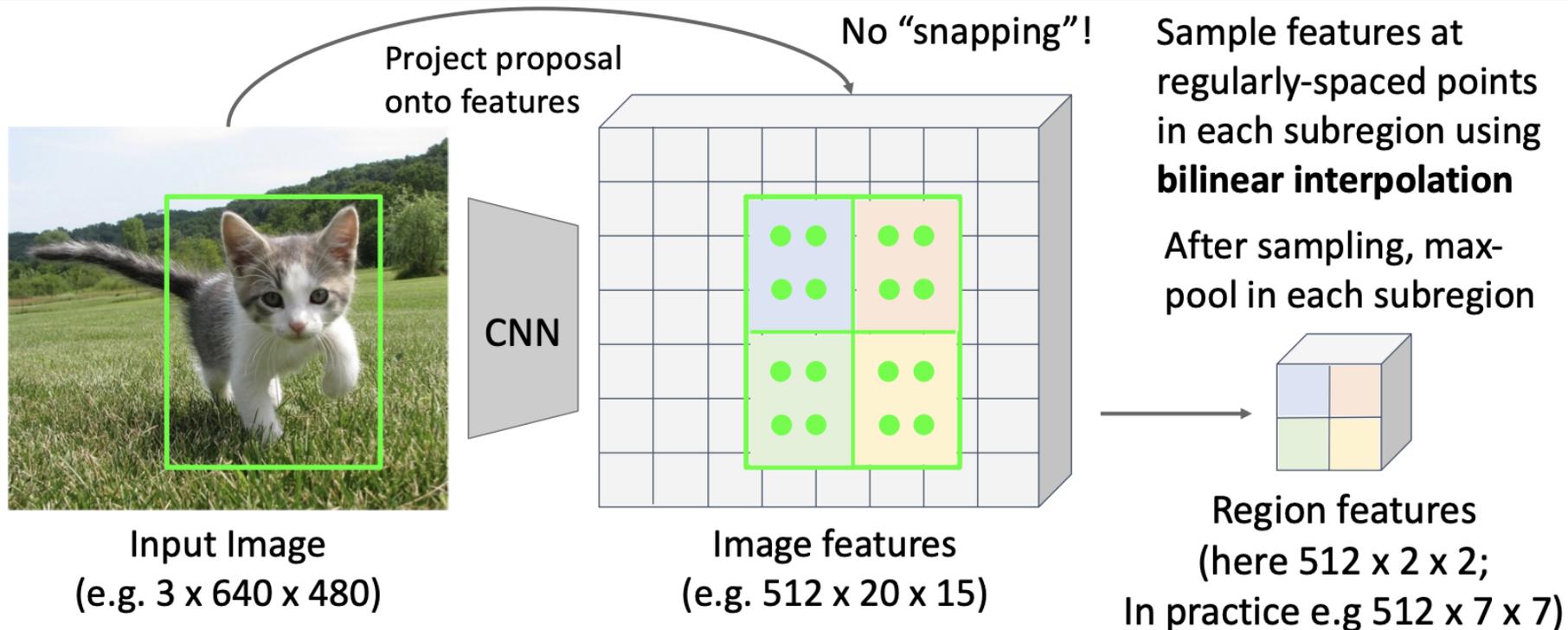
$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$



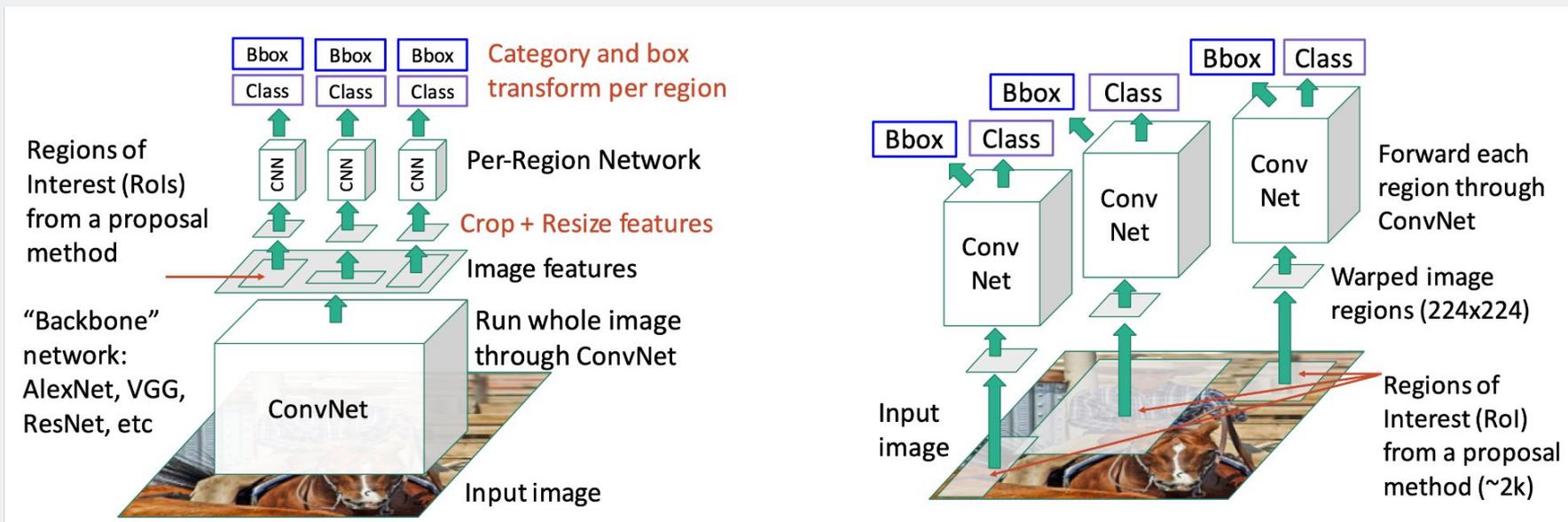
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



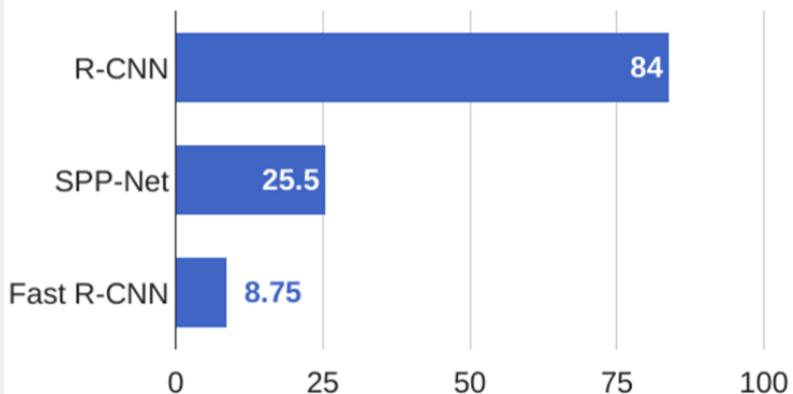
Fast R-CNN vs “Slow” R-CNN

Fast R-CNN: Apply differentiable cropping to shared image features

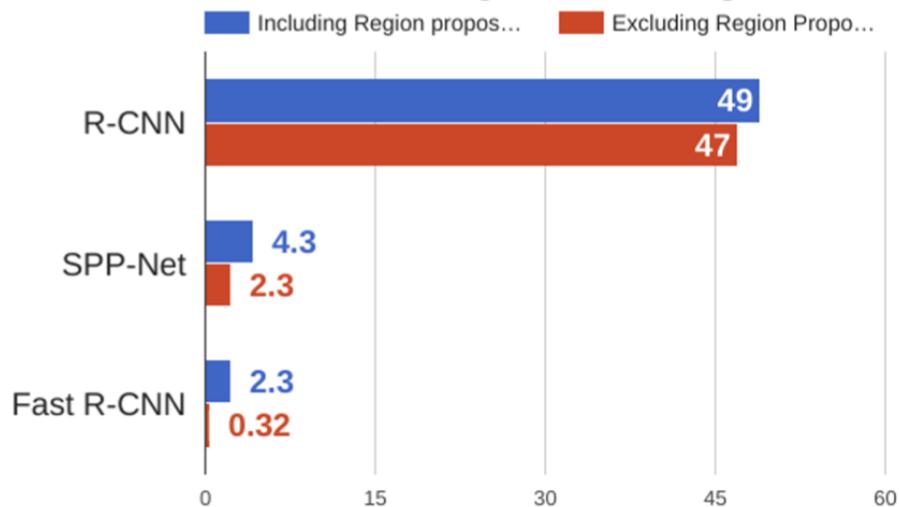


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)



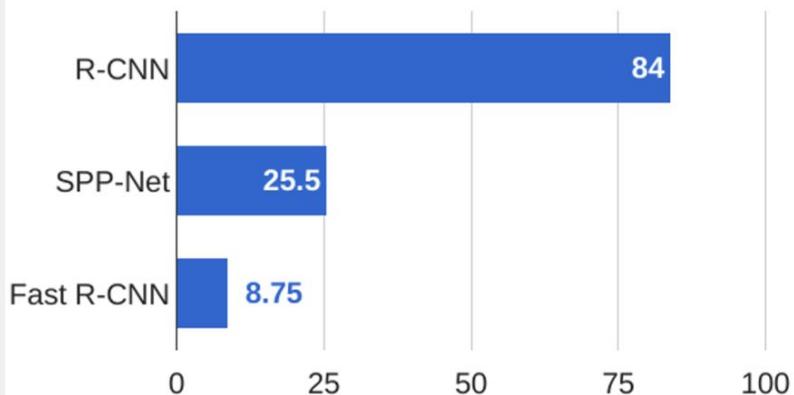
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

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

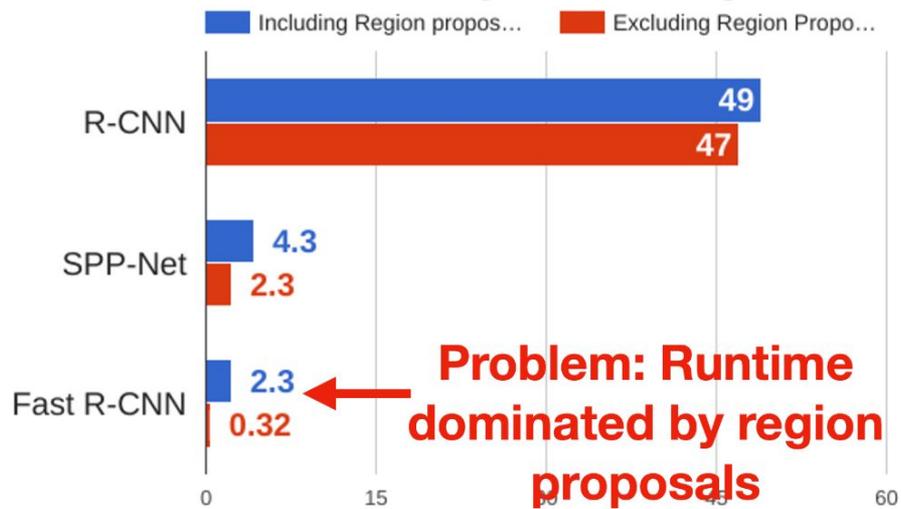
Girshick, “Fast R-CNN”, ICCV 2015

Fast R-CNN vs “Slow” R-CNN

Training time (Hours)

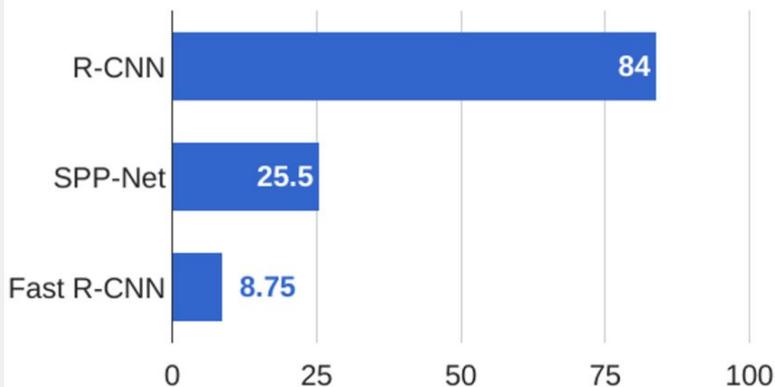


Test time (seconds)

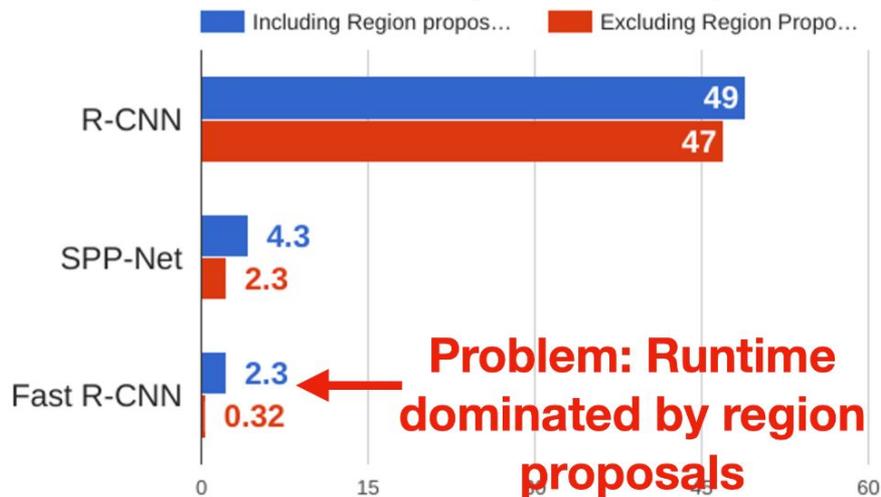


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)



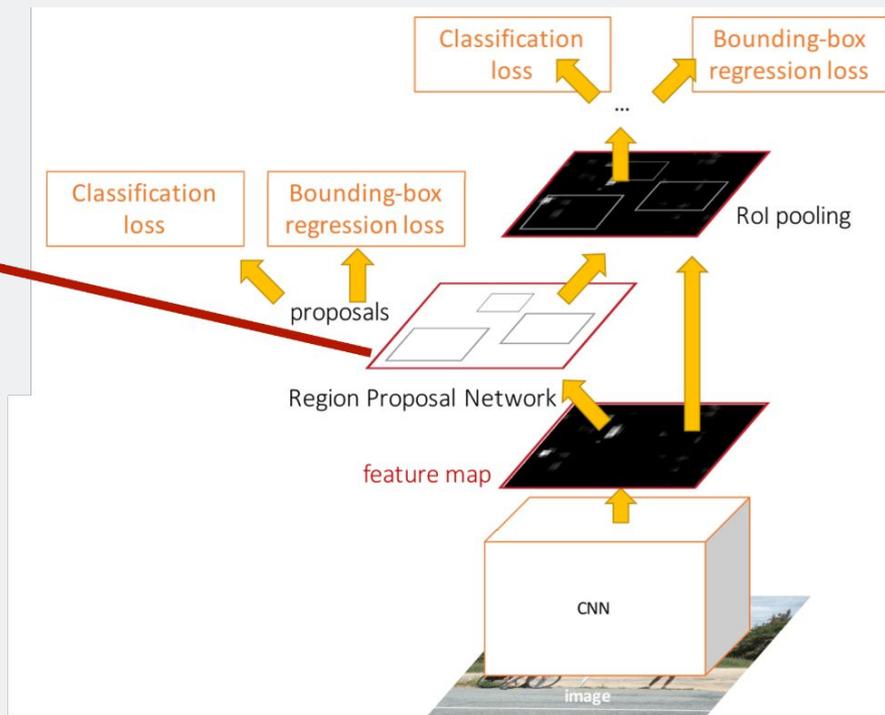
Problem: Runtime dominated by region proposals

Recall: Region proposals computed by heuristic “Selective search” algorithm on CPU — let’s learn them with a CNN

Fast_{er} R-CNN: Learnable Region Proposals

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

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

Figure copyright 2015, Ross Girshick; reproduced with permission

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

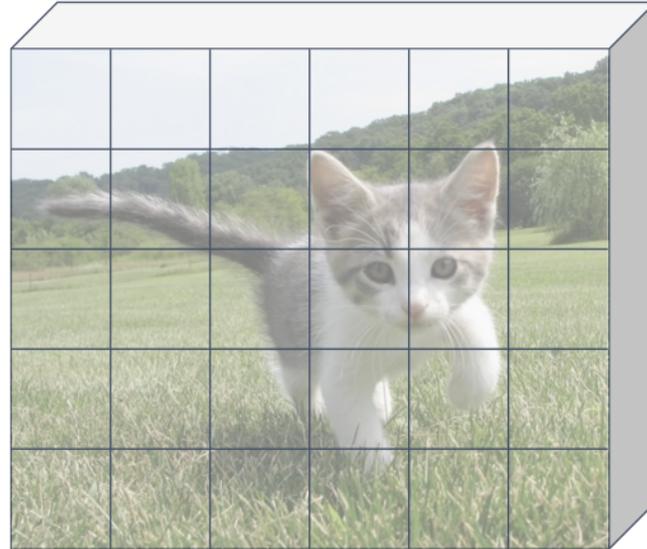


Image features
(e.g. 512 x 5 x 6)

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

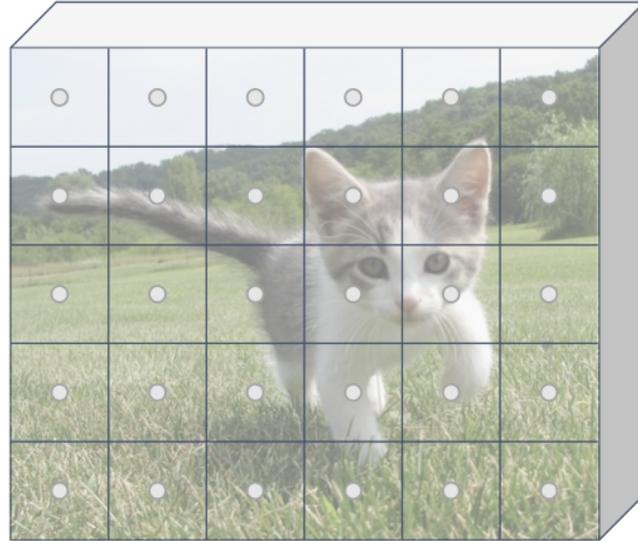
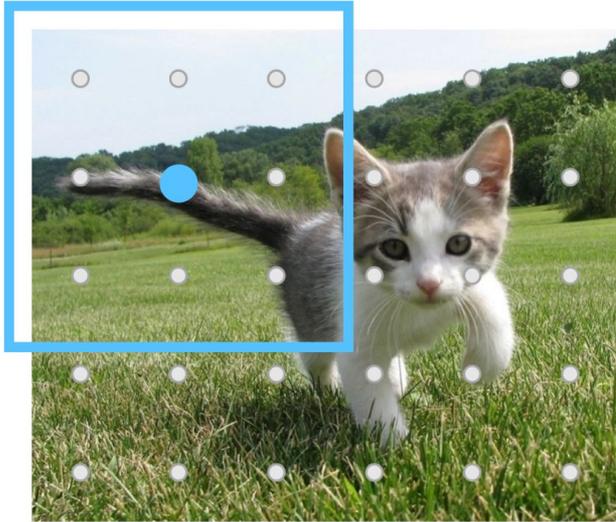


Image features
(e.g. 512 x 5 x 6)

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

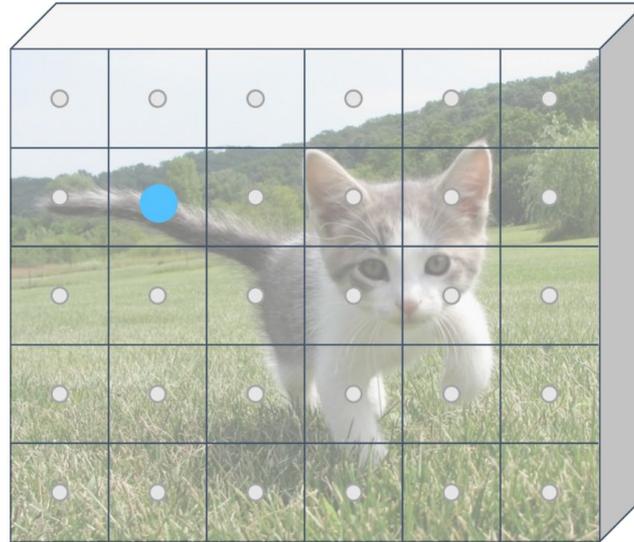
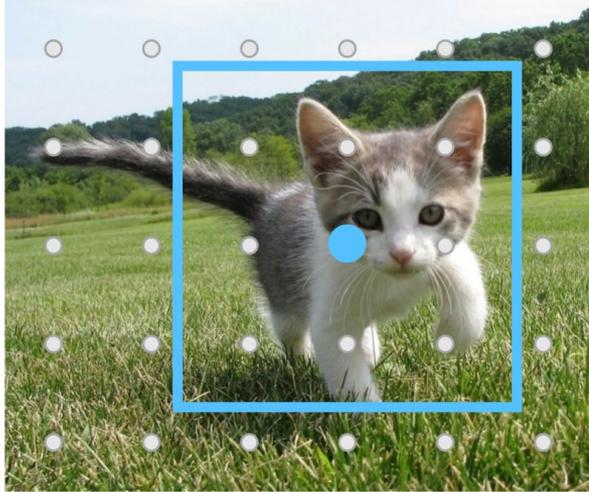


Image features
(e.g. 512 x 5 x 6)

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

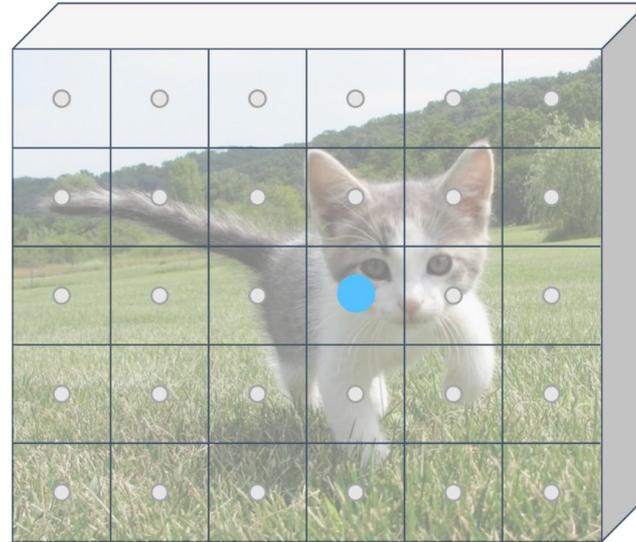
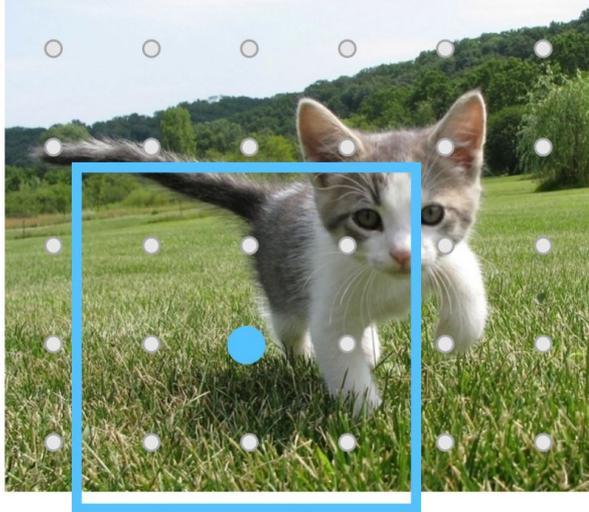


Image features
(e.g. 512 x 5 x 6)

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

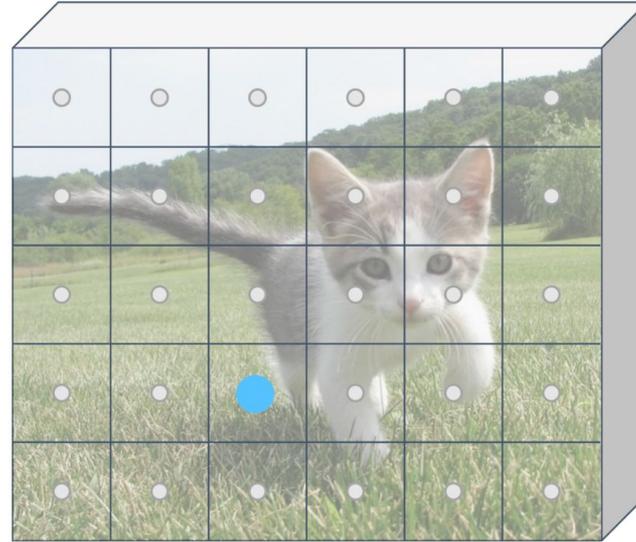


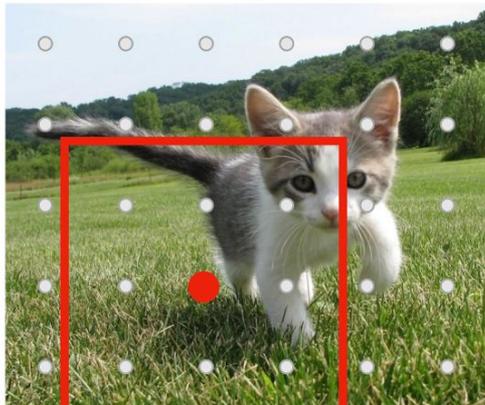
Image features
(e.g. 512 x 5 x 6)

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image

Each feature corresponds to a point in the input



Input Image
(e.g. 3 x 640 x 480)

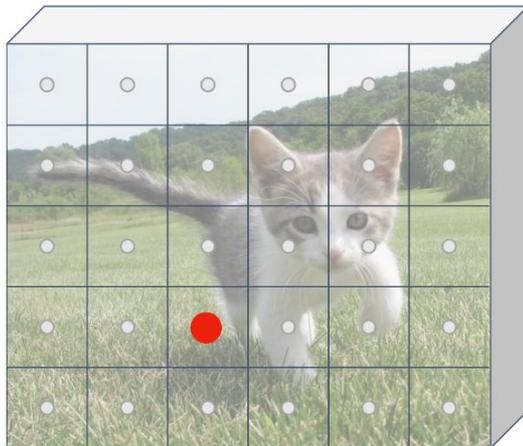


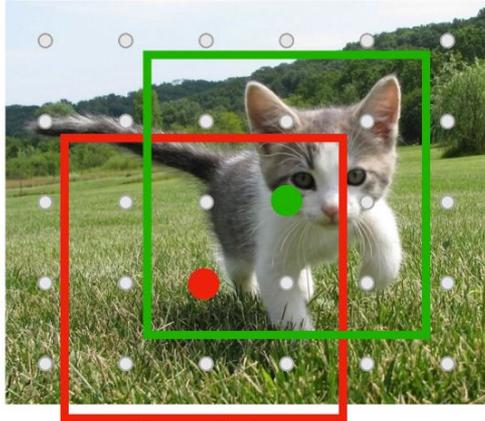
Image features
(e.g. 512 x 5 x 6)

Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

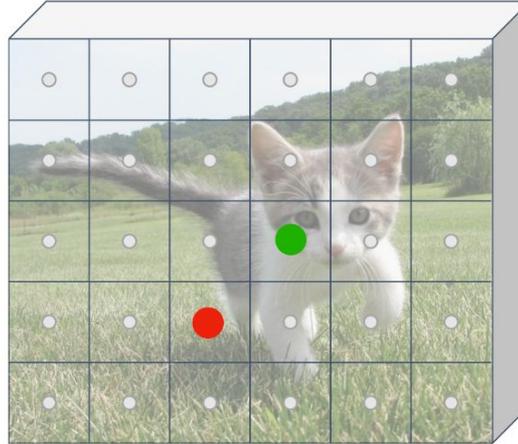


Image features
(e.g. 512 x 5 x 6)

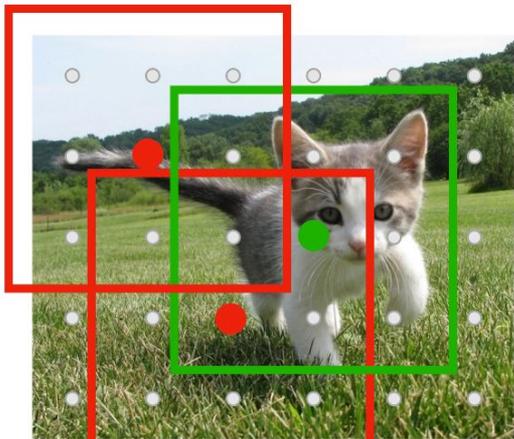


Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

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

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

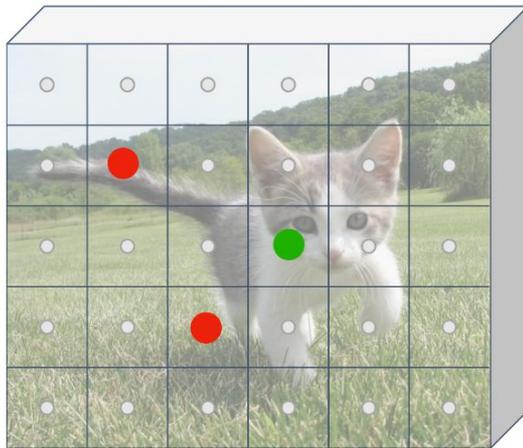
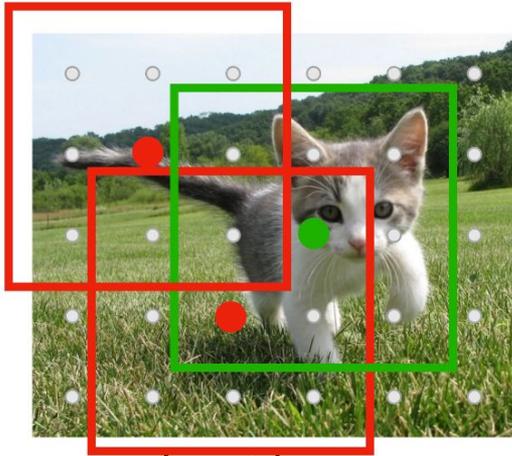


Image features
(e.g. 512 x 5 x 6)

Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

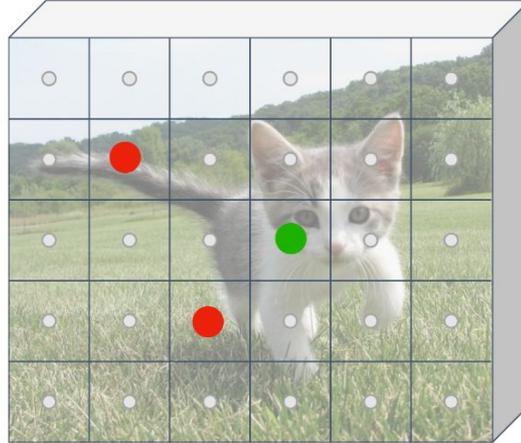
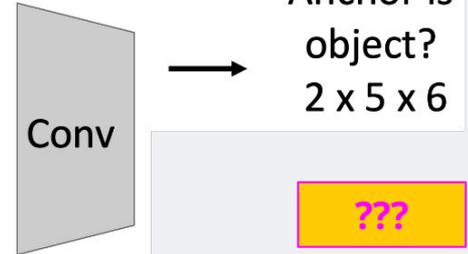


Image features
(e.g. 512 x 5 x 6)

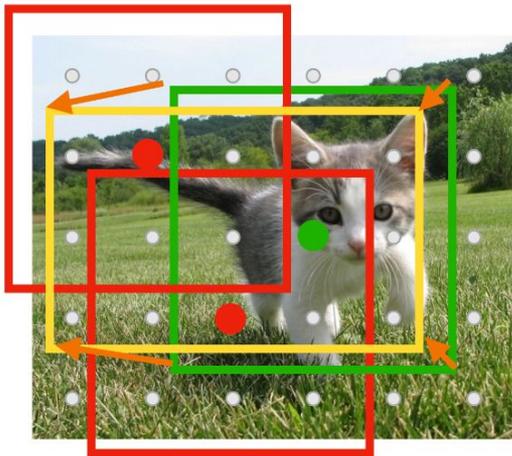
Predict object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)



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 x 640 x 480)



Each feature corresponds to a point in the input

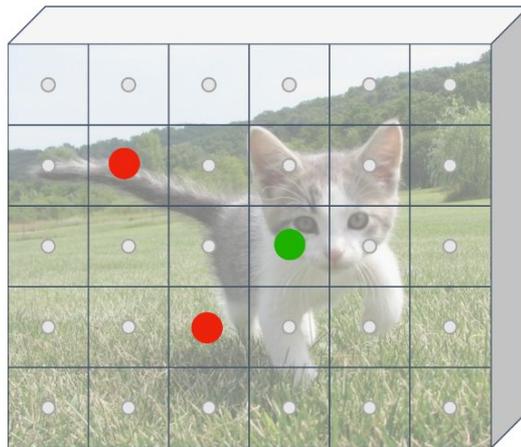
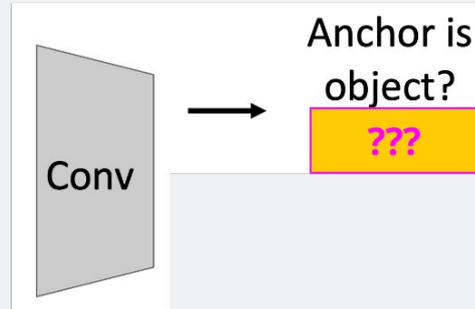


Image features
(e.g. 512 x 5 x 6)

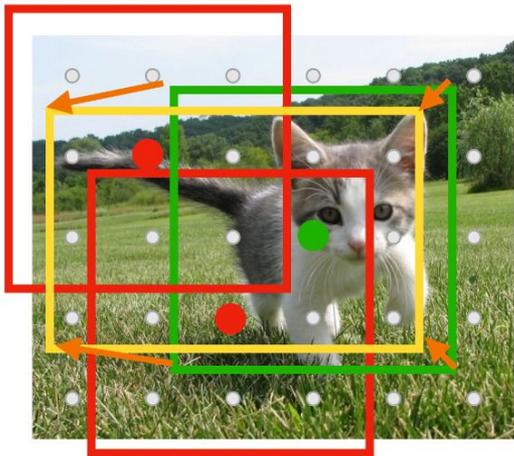
For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)



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 x 640 x 480)



Each feature corresponds to a point in the input

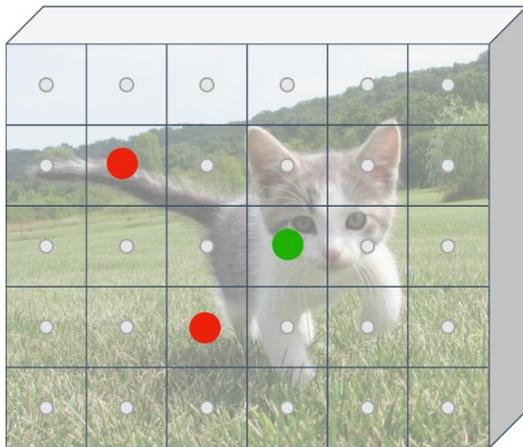


Image features
(e.g. 512 x 5 x 6)

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

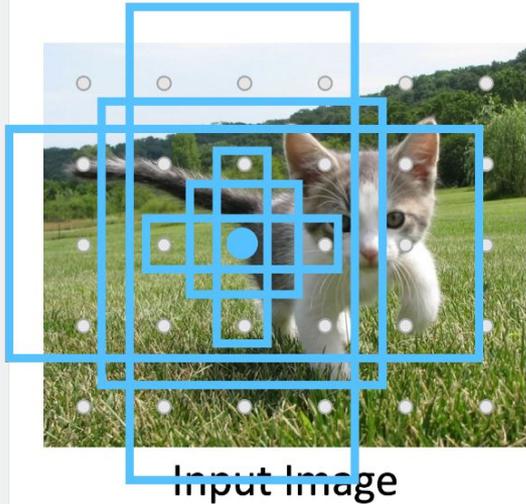


Anchor is object?
???
Anchor transforms
4 x 5 x 6

Classify each anchor as **positive (object)** or **negative (no object)**

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

CNN

Each feature corresponds to a point in the input

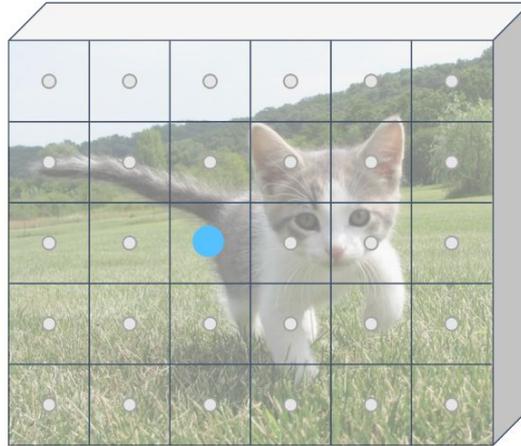


Image features
(e.g. 512 x 5 x 6)

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



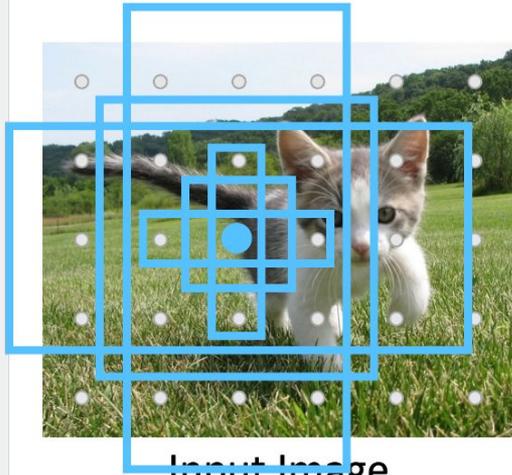
Anchor is object?
 $2K \times 5 \times 6$



Anchor transforms
 $4K \times 5 \times 6$

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

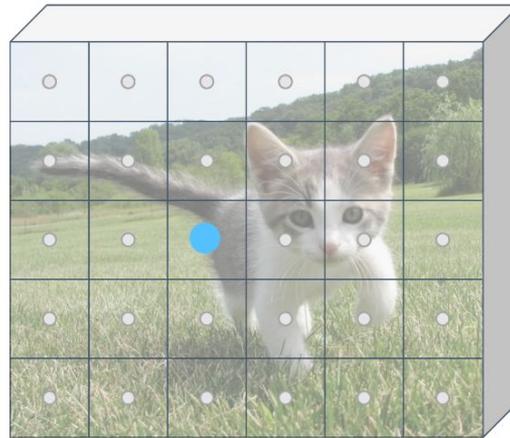


Image features
(e.g. 512 x 5 x 6)

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



Anchor is
object?
2K x 5 x 6

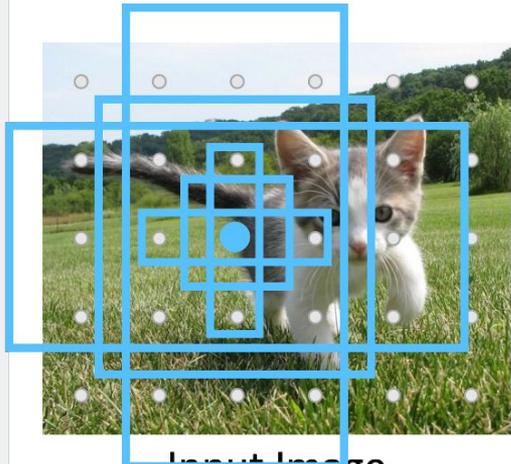


Anchor
transforms
4K x 5 x 6

During training, supervised positive / negative anchors and box transforms like R-CNN

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

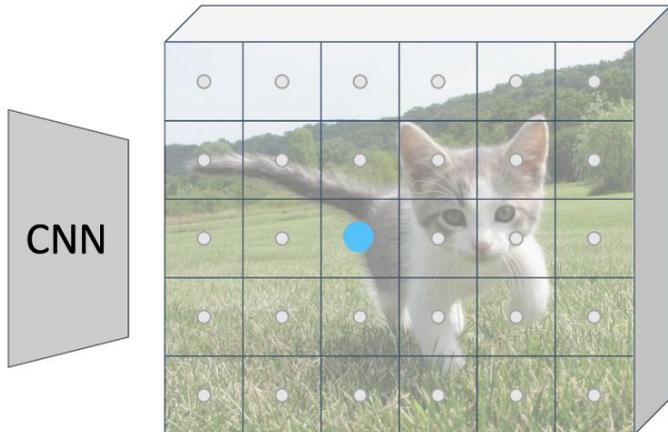


Image features
(e.g. 512 x 5 x 6)

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



Anchor is object?
2K x 5 x 6

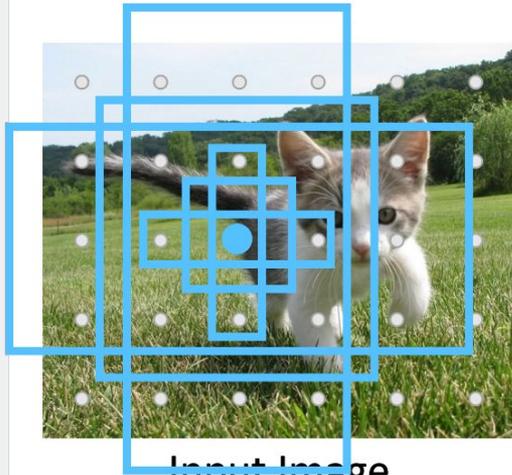


Anchor transforms
4K x 5 x 6

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

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

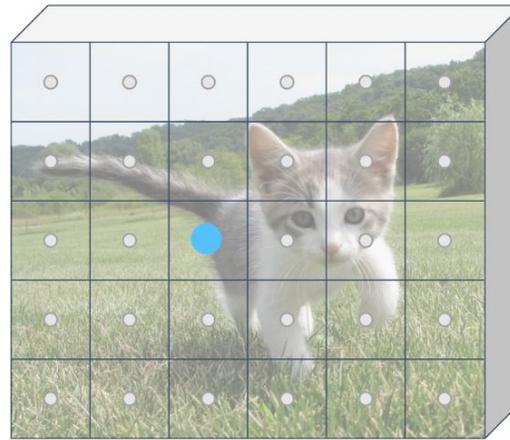


Image features
(e.g. 512 x 5 x 6)



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



Anchor is object?
2K x 5 x 6

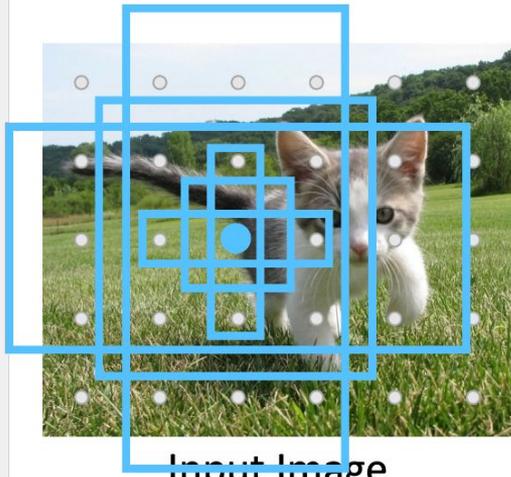


Anchor transforms
4K x 5 x 6

Negative anchors: < 0.3 IoU with all GT boxes. Don't supervised transforms for negative boxes.

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

Each feature corresponds to a point in the input

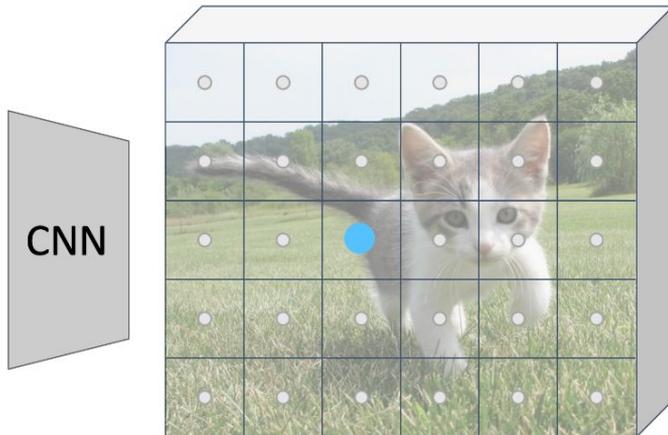
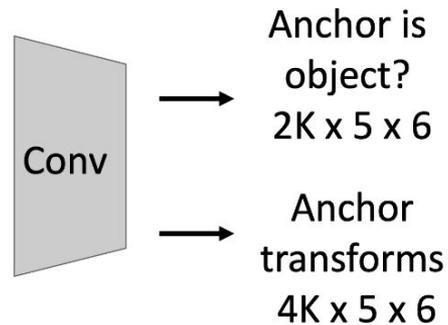


Image features
(e.g. 512 x 5 x 6)

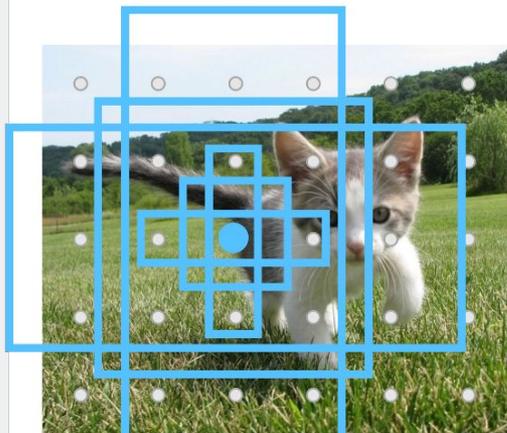
In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



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

Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

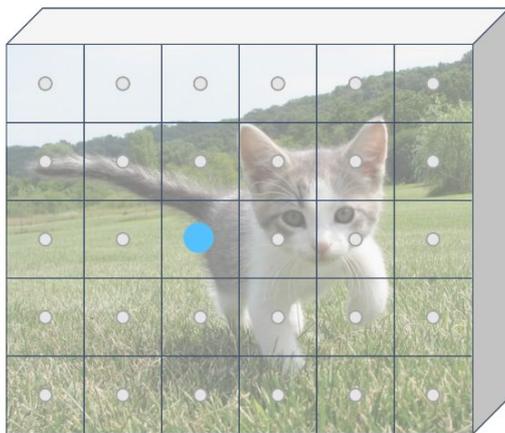


Image features
(e.g. 512 x 5 x 6)

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



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



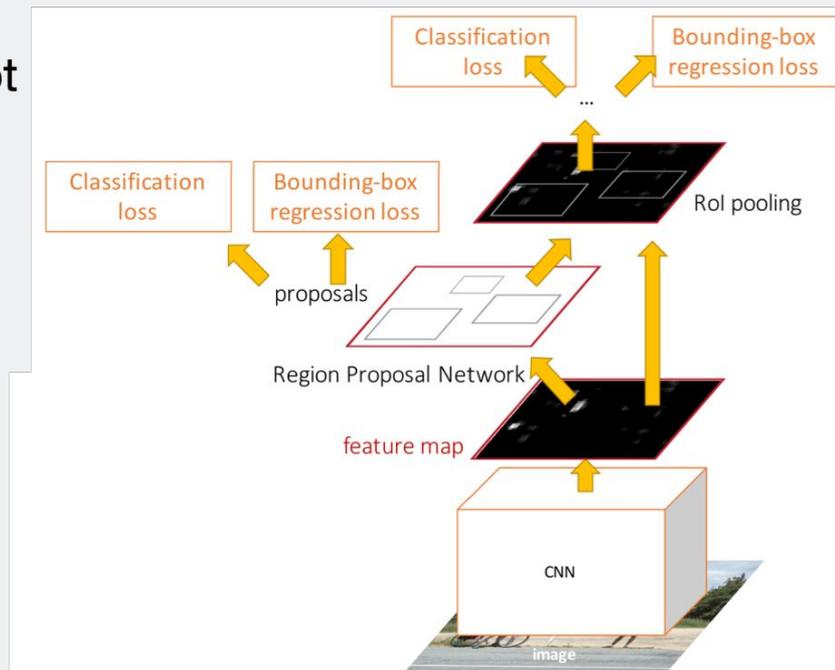
Anchor transforms
 $4K \times 5 \times 6$

At test-time, sort all $K \cdot 5 \cdot 6$ boxes by their positive score, take top 300 as our region proposals

Faster R-CNN: Learnable Region Proposals

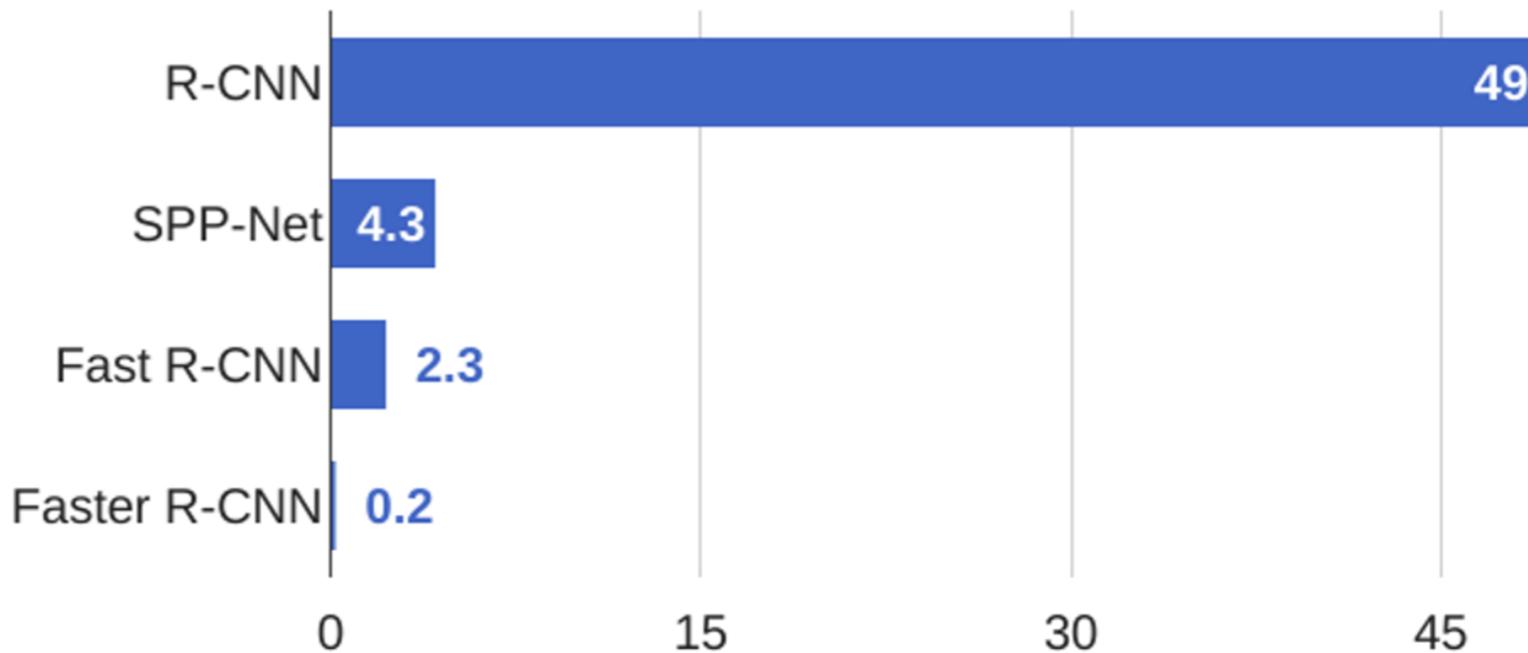
Jointly train four losses:

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



Faster R-CNN: Learnable Region Proposals

R-CNN Test-Time Speed



Extend Faster R-CNN to Image Segmentation: Mask R-CNN

Classification



"Chocolate Pretzels"

No spatial extent

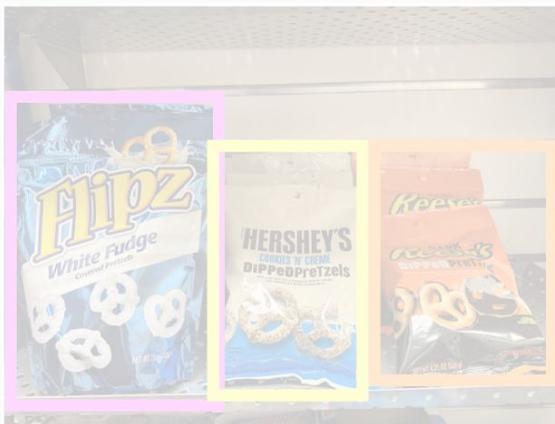
Semantic Segmentation



Chocolate Pretzels,
Shelf

No objects, just pixels

Object Detection



Flipz, Hershey's, Keese's

Multiple objects

Instance Segmentation



Extend Faster R-CNN to Image Segmentation: Mask R-CNN

Instance Segmentation

Detect all objects in the image and identify the pixels that belong to each object (Only things!)

Approach

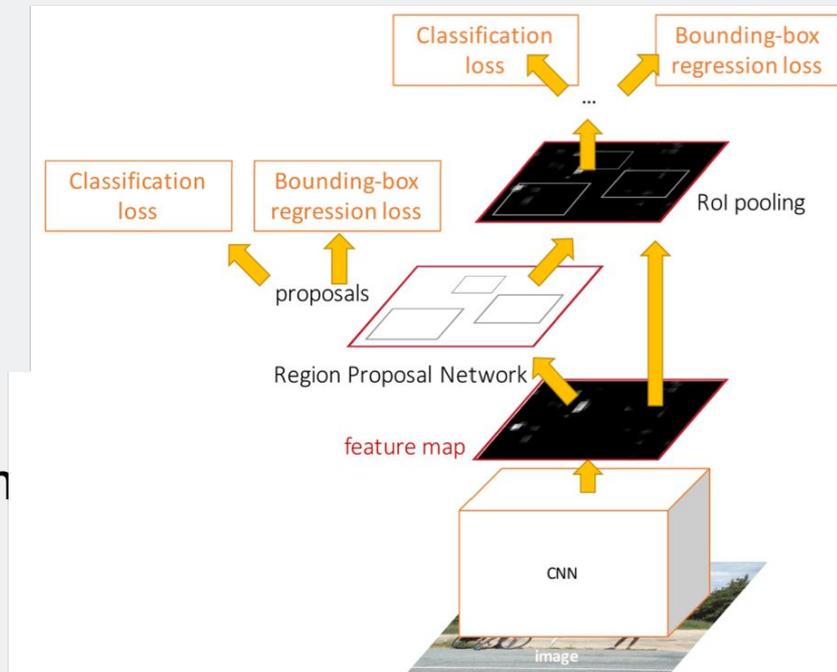
Perform object detection then predict a segmentation mask for each object detected!



Extend Faster R-CNN to Mask R-CNN

Faster R-CNN

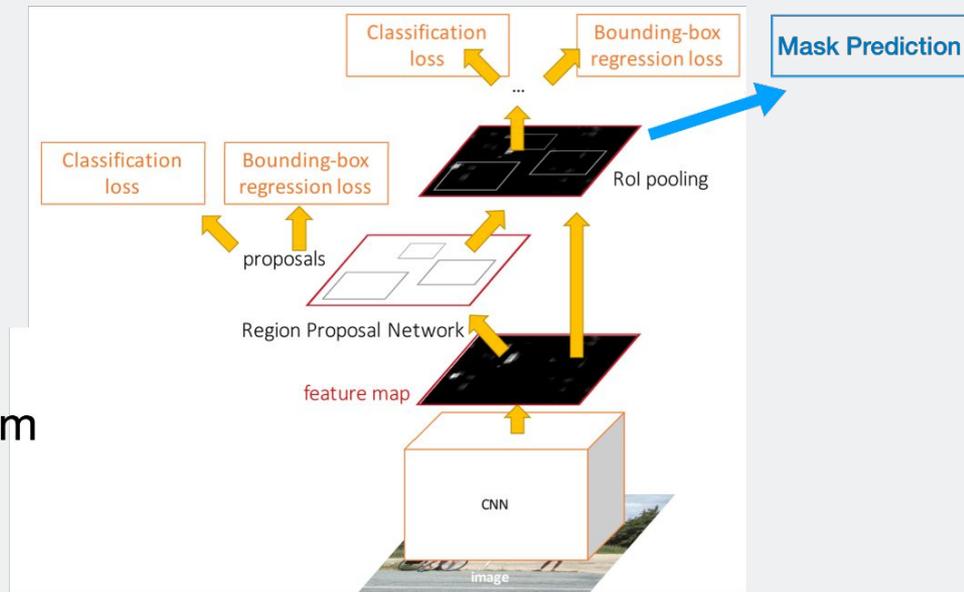
1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 1. **Object classification:** classify proposals
 2. **Object regression:** predict transform from proposal box to object box



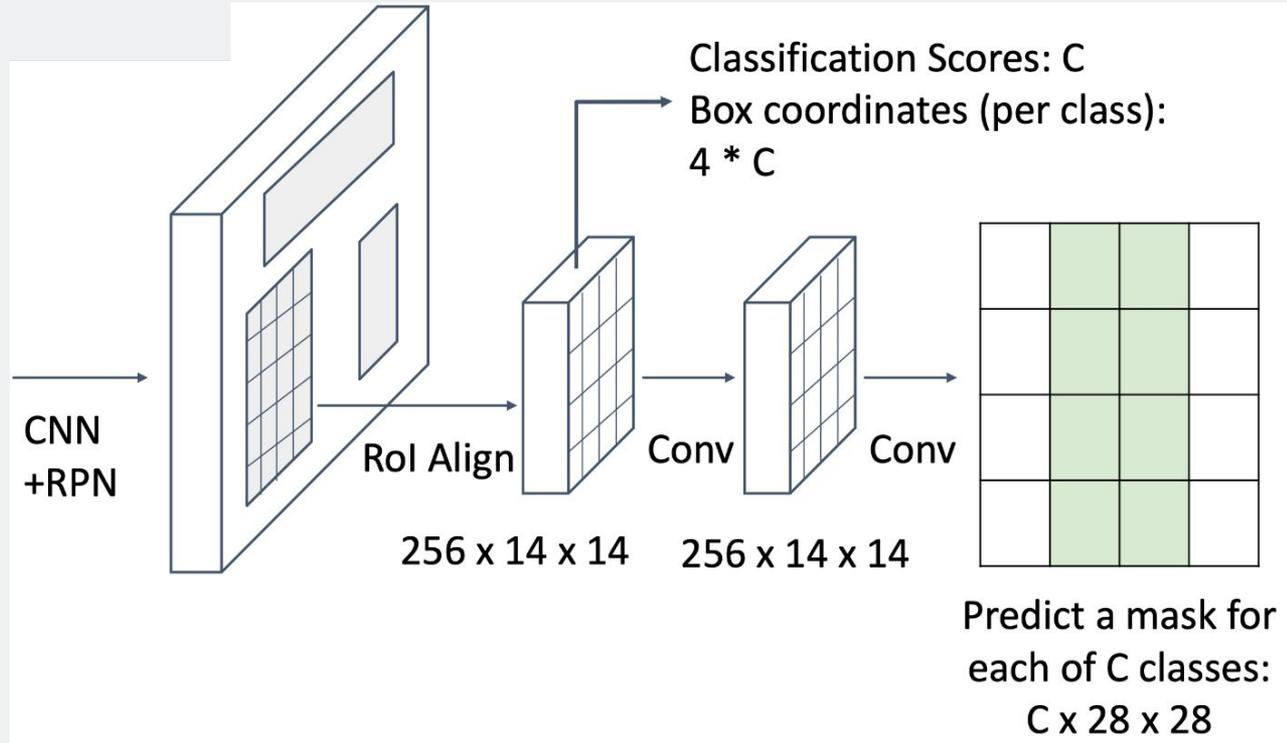
Extend Faster R-CNN to Mask R-CNN

Mask R-CNN

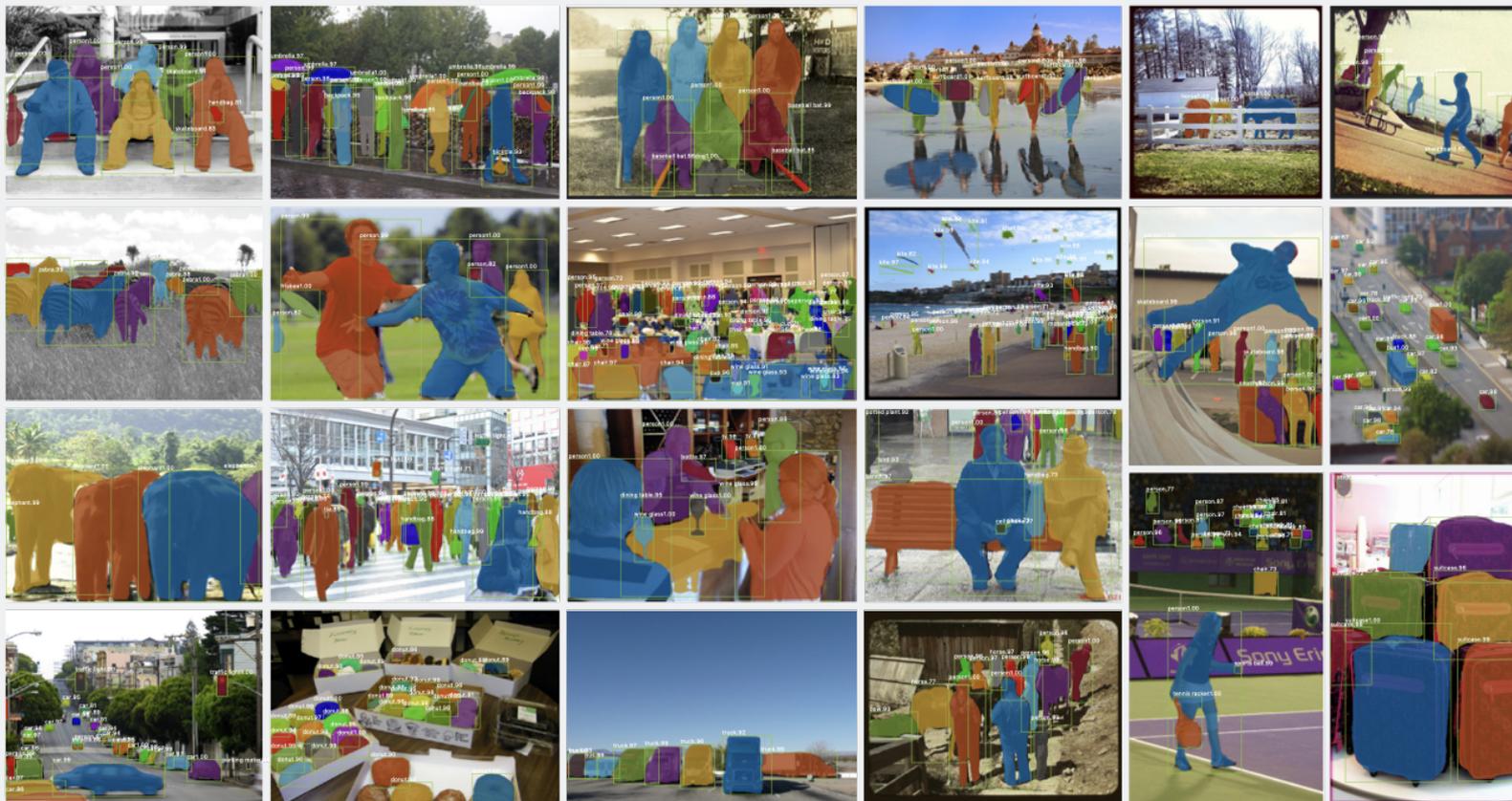
1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region



Mask R-CNN



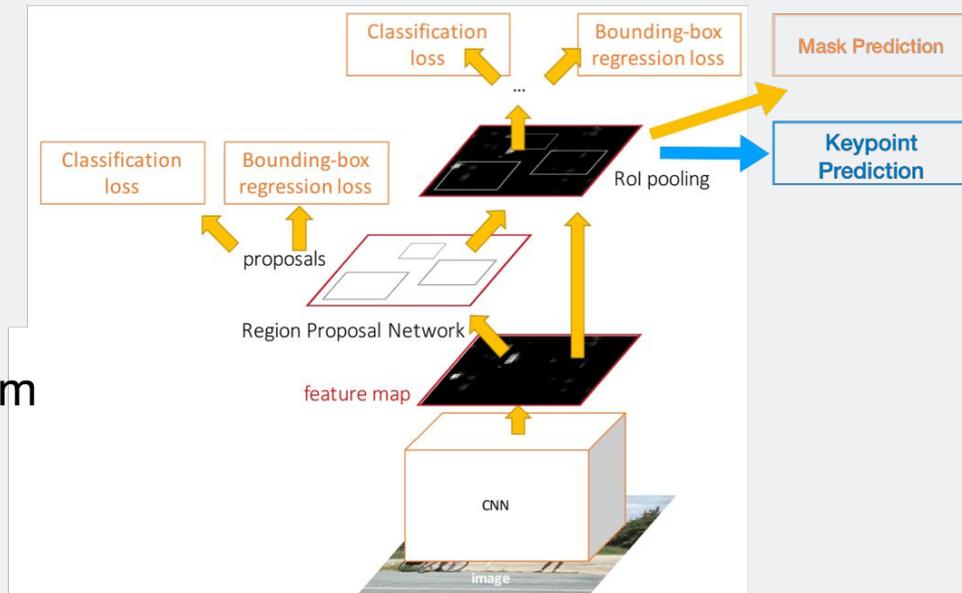
Mask R-CNN



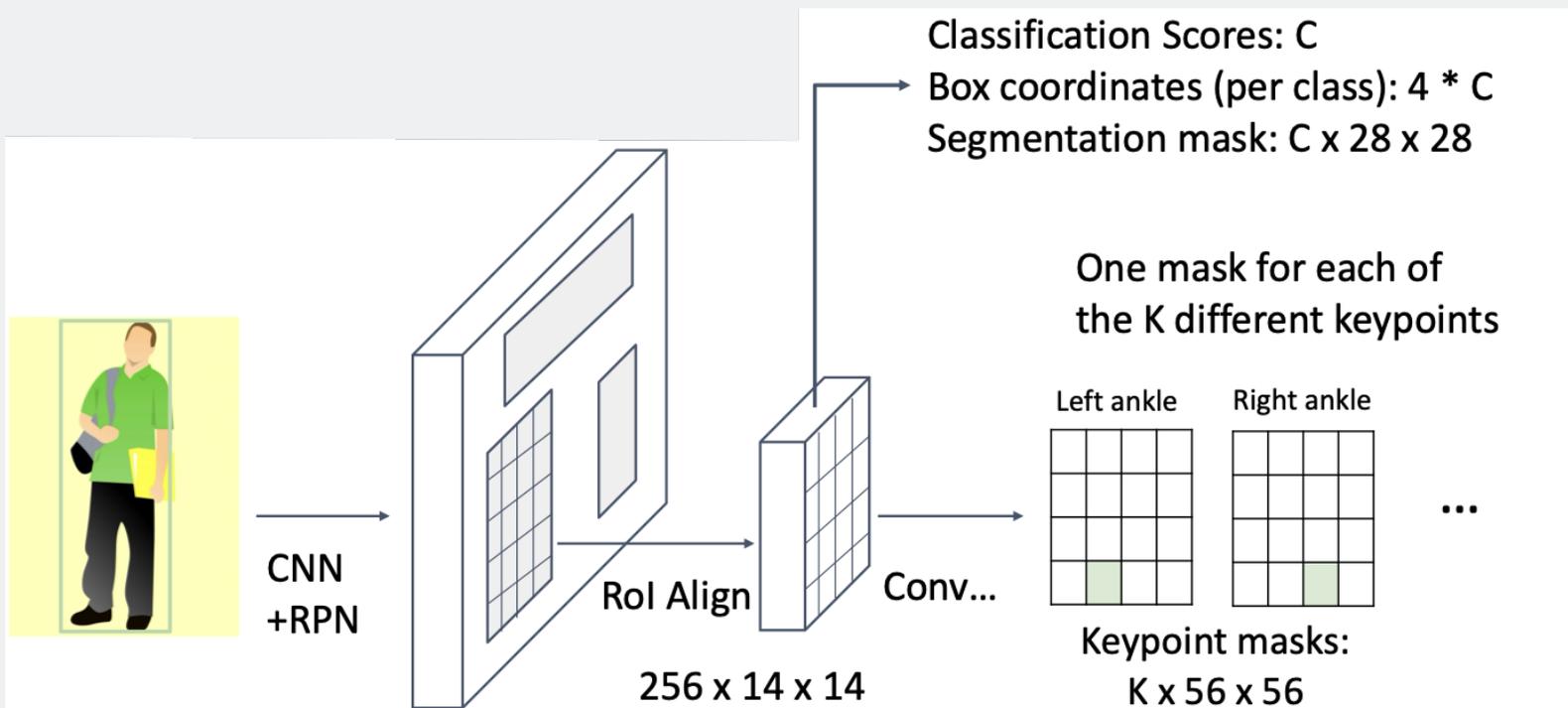
Mask R-CNN for Human Pose Estimation

Mask R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region
 - d. **Keypoint Prediction:** predict binary mask for human key points

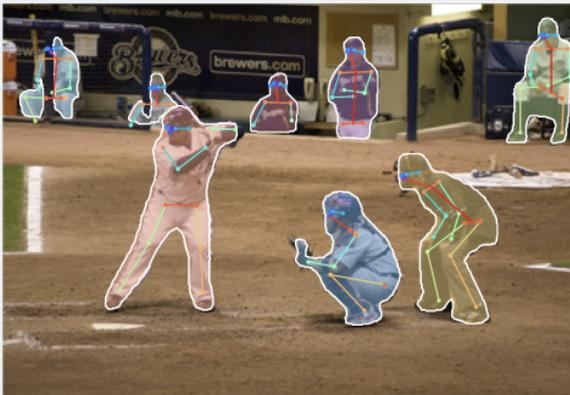
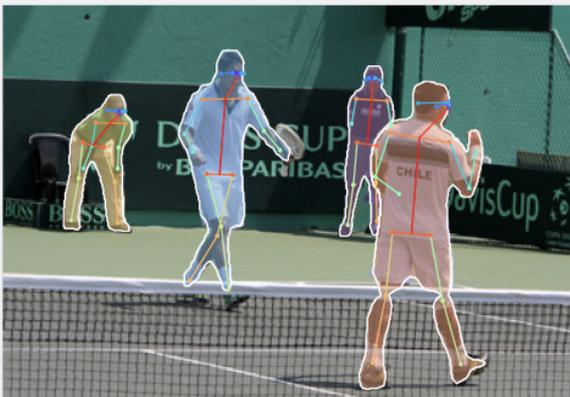


Mask R-CNN for Human Pose Estimation



Ground-truth has one “pixel” turned on per keypoint. Train with softmax loss

Mask R-CNN for Human Pose Estimation



Detection without Anchors - Transformers

(more details later)

- + No RPN
- + No Anchors
- + No NMS

- Slow to train
- Worse than Faster R-CNN for smaller objects

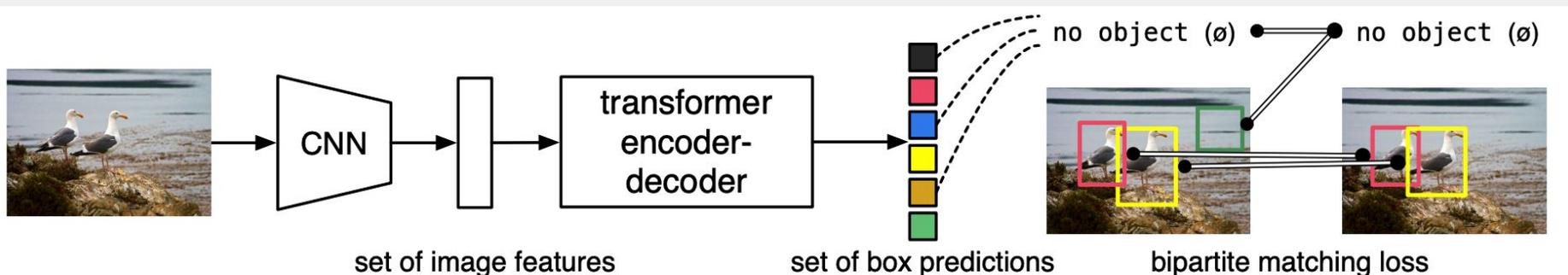


Fig. 1: DETR directly predicts (in parallel) the final set of detections by combining a common CNN with a transformer architecture. During training, bipartite matching uniquely assigns predictions with ground truth boxes. Prediction with no match should yield a “no object” (\emptyset) class prediction.

Two Stage vs One Stage Detectors

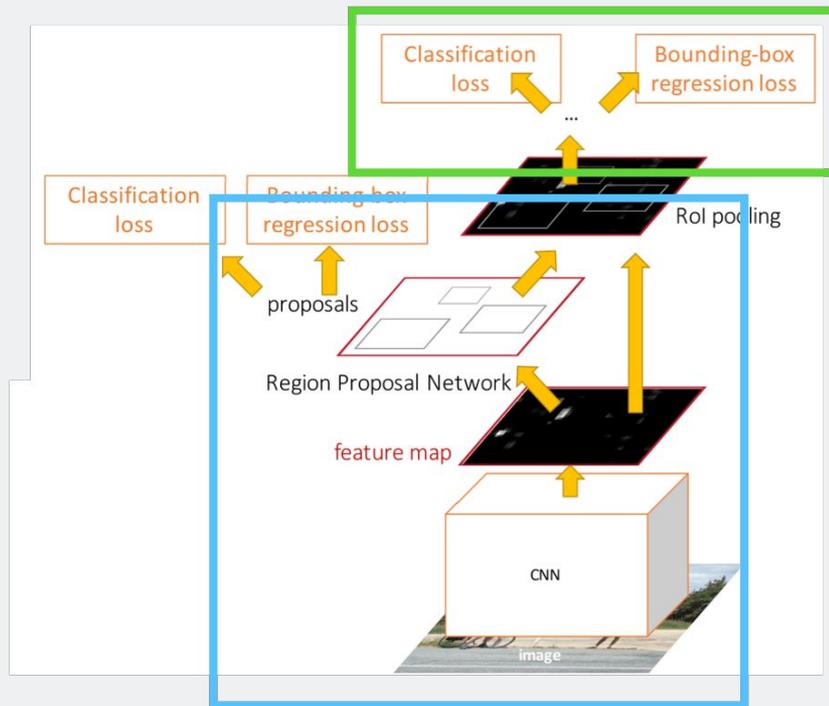
Faster R-CNN is a two-stage object detector

First stage: Run once per image

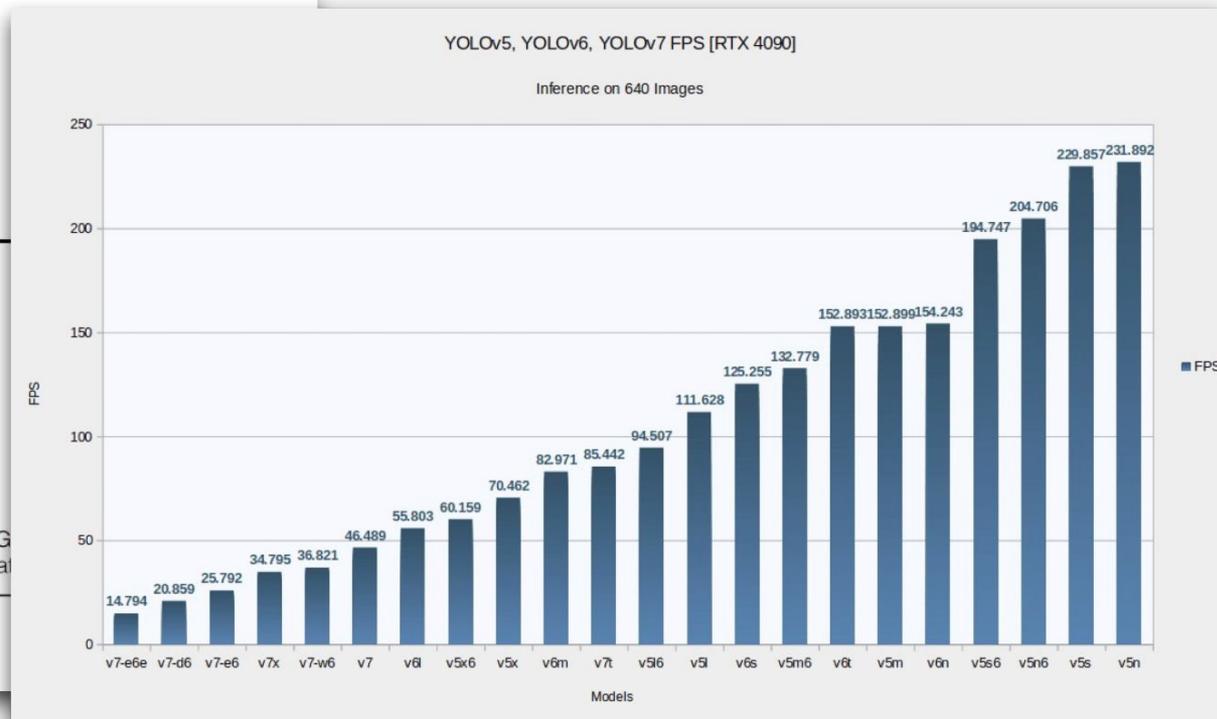
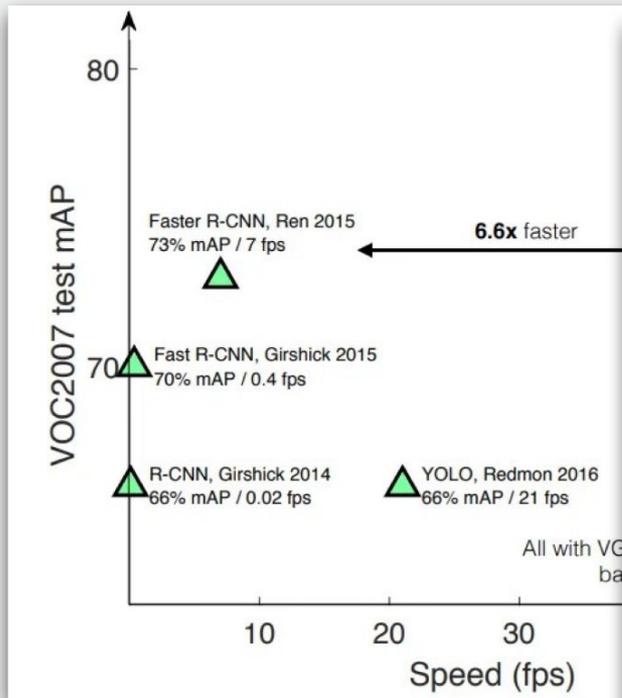
- Backbone Network
- Region Proposal Network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict Object Class
- Prediction bbox offset

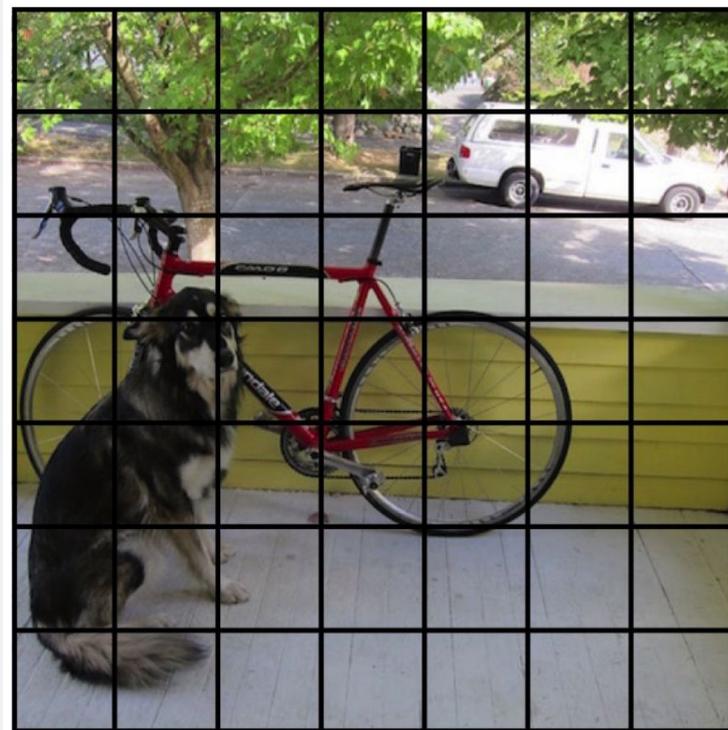


Two Stage vs One Stage Detectors



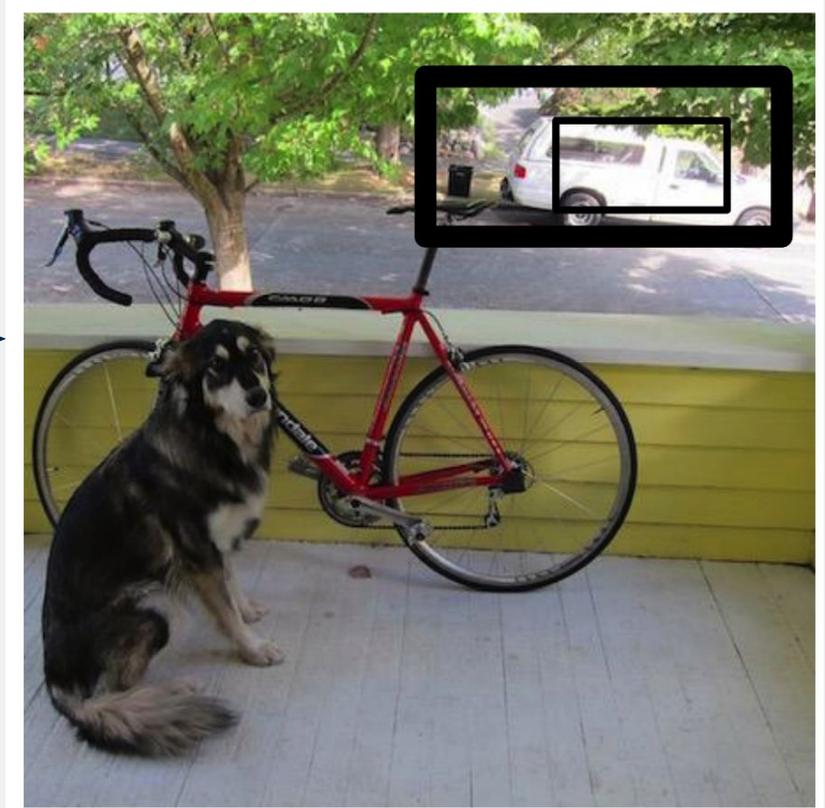
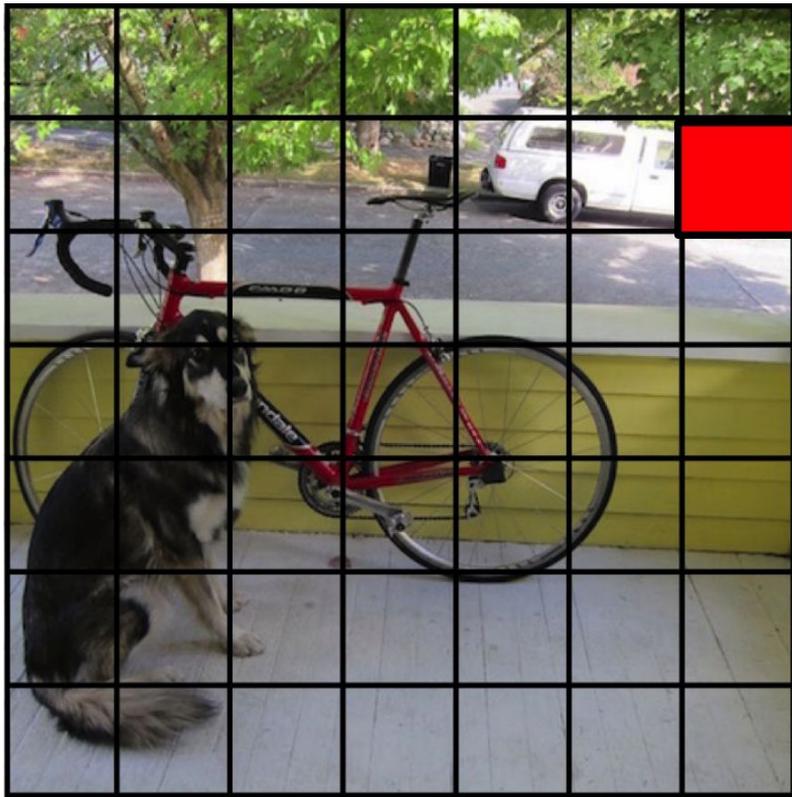
YOLO: You Only Look Once

Split into $S \times S$ grid



YOLO

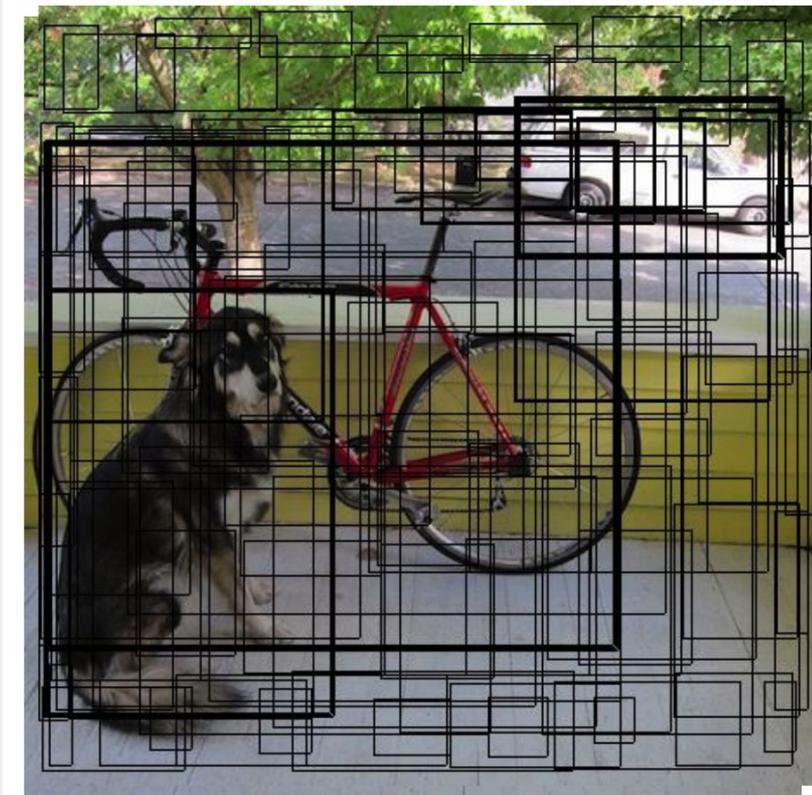
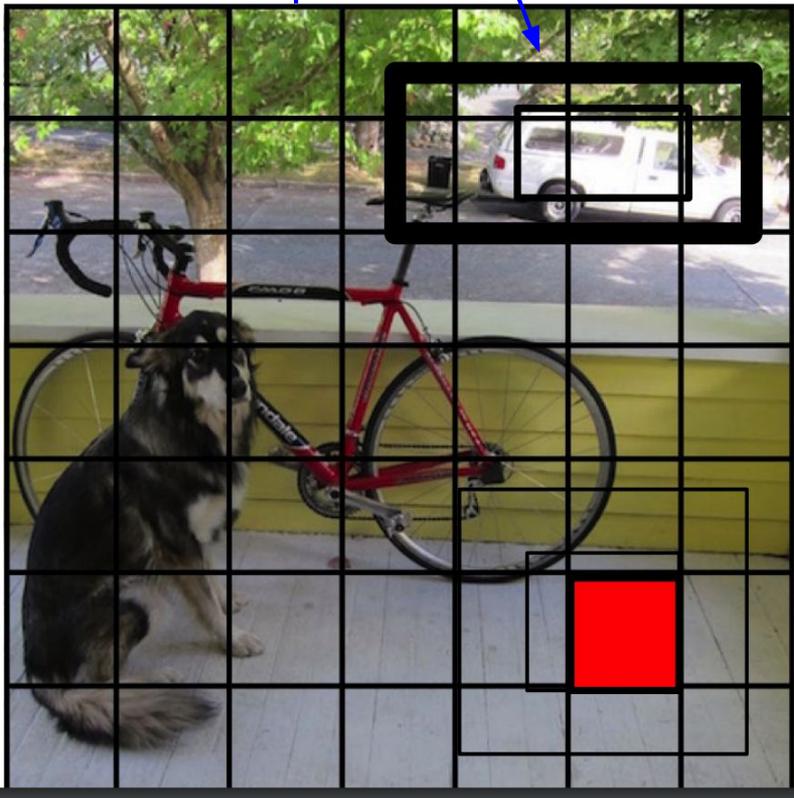
Each cell predicts **B boxes**(x,y,w,h) and **confidences of each box**: $P(\text{Object})$



YOLO

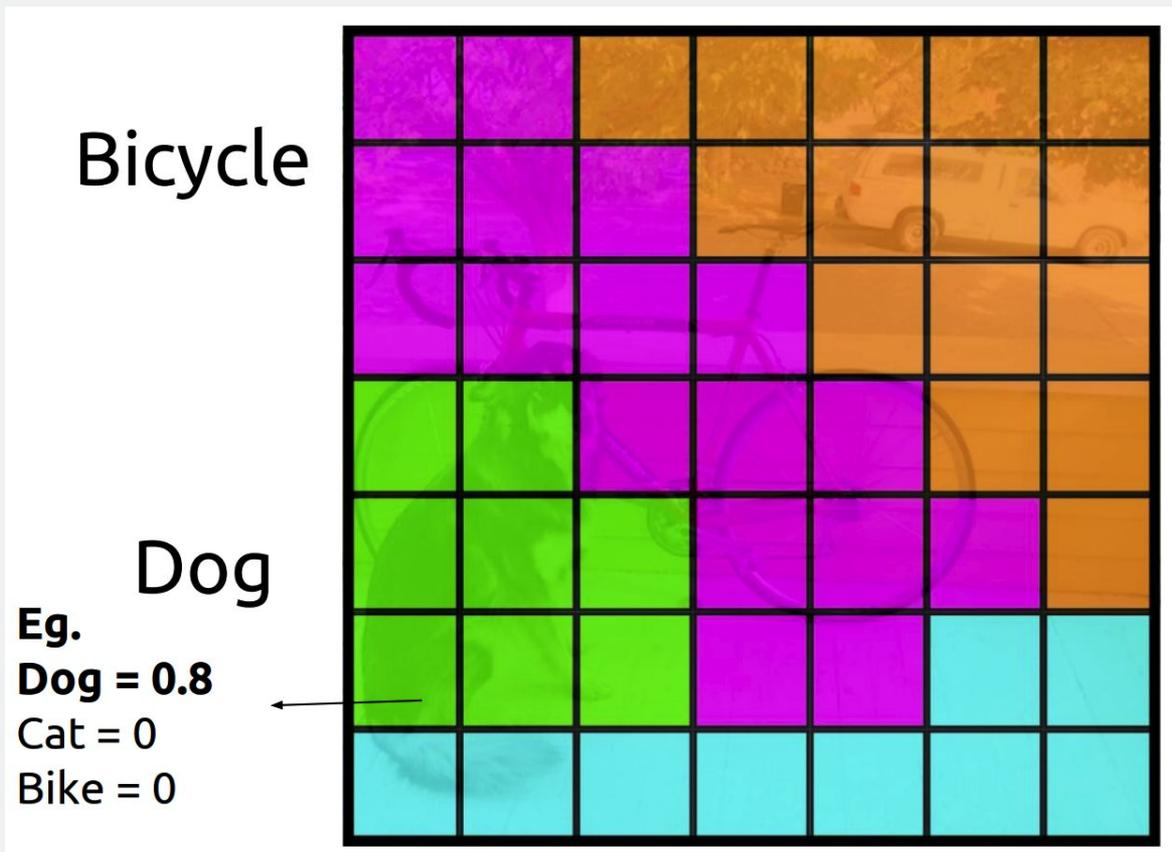
(example: $B=2$)

Each cell predicts boxes and confidences



YOLO

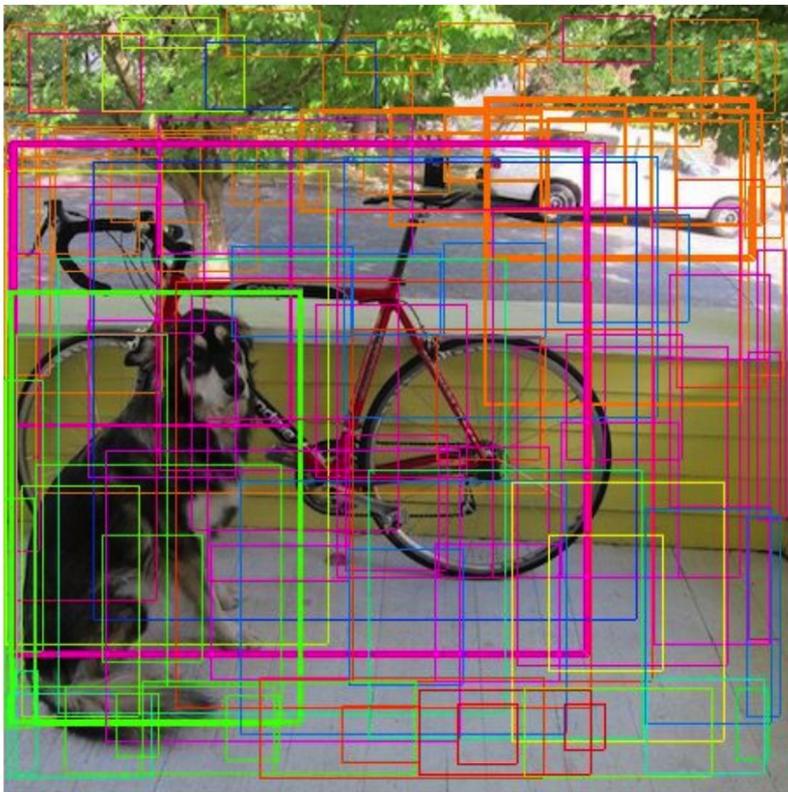
Each cell also predicts a class probability
Conditioned on object: $P(\text{Class} \mid \text{Object})$



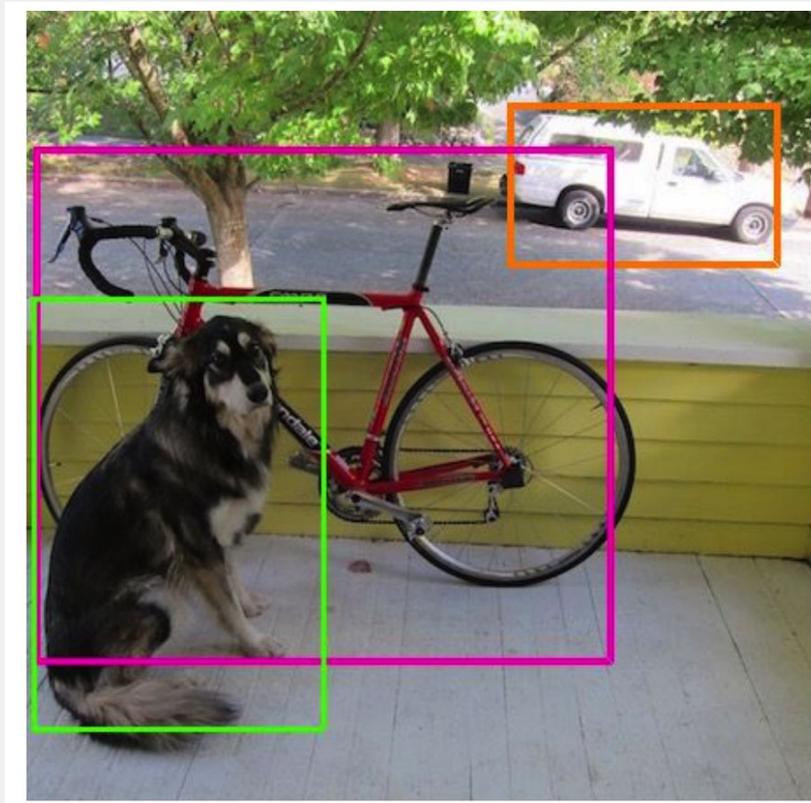
YOLO

Combine boxes and class probabilities
Apply NMS

$$P(\text{class}|\text{Object}) * P(\text{Object}) \\ = P(\text{class})$$



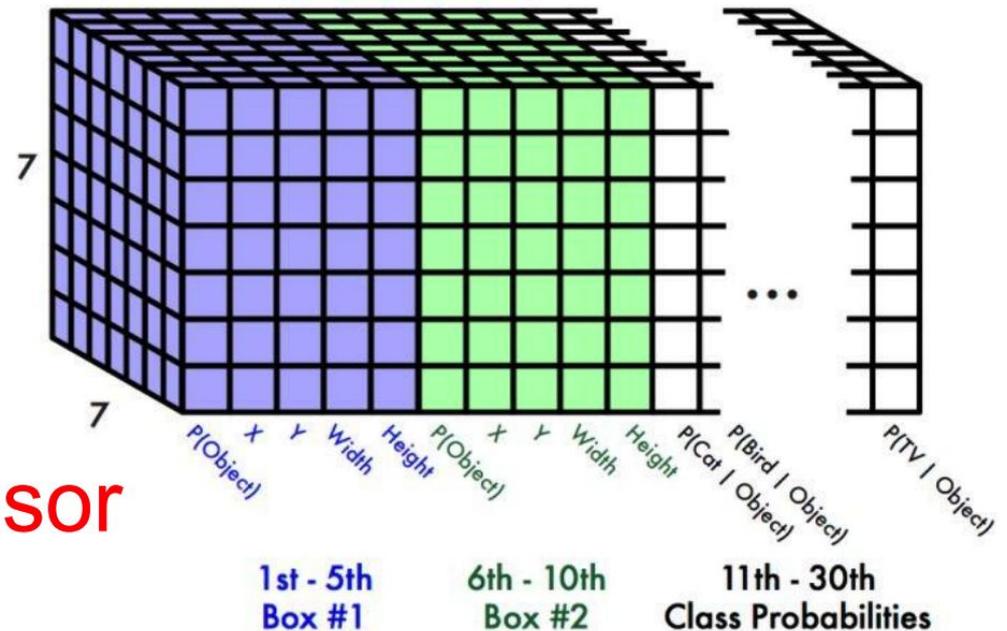
NMS
→



YOLO

Each cell predicts:

- For each bounding box:
 - 4 coordinates (x, y, w, h)
 - 1 confidence value
- Some number of class probabilities



$S * S * (B * 5 + C)$ tensor

What's the size of the YOLO output

YOLO Variants

<https://arxiv.org/abs/1506.02640> (original YOLO paper, 2015)

<https://arxiv.org/abs/2304.00501> (from YOLOv1 to YOLOv8 and YOLO-NAS, 2023)

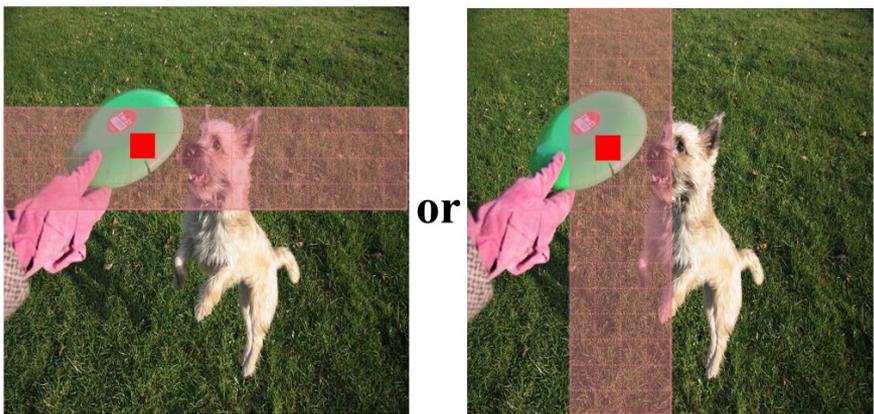
<https://arxiv.org/abs/2402.13616> (YOLOv9, Feb. 2024)

<https://arxiv.org/abs/2405.14458> (YOLOv10, May 2024)

<https://arxiv.org/abs/2410.17725> (YOLOv11, Oct. 2024)

<https://www.arxiv.org/abs/2502.12524> (YOLOv12, Feb.18, 2025)

<https://github.com/sunsmarterjje/yolov12> (YOLOv12 pytorch, Attention-centric YOLO)



Area attention (Ours)



input

YOLOv10

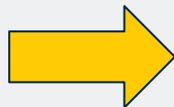
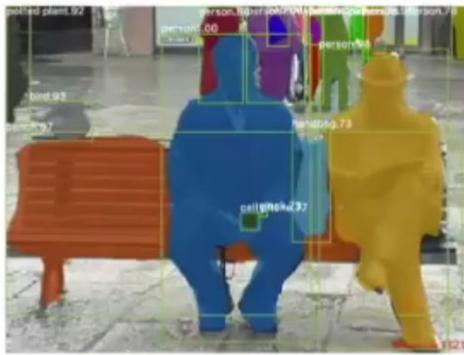
YOLOv11

YOLOv12

heatmap
comparison

3D Shape Prediction: Mesh R-CNN

Mask R-CNN:
2D Image -> 2D shapes



Mesh R-CNN:
2D image -> 3D triangle meshes

Input Image



2D Recognition



3D Meshes



3D Voxels

Figure 1. Mesh R-CNN takes an input image, predicts object instances in that image and infers their 3D shape. To capture diversity in geometries and topologies, it first predicts coarse voxels which are refined for accurate mesh predictions.

3D Shape Prediction: Mesh R-CNN

Input: Single RGB image

Output:

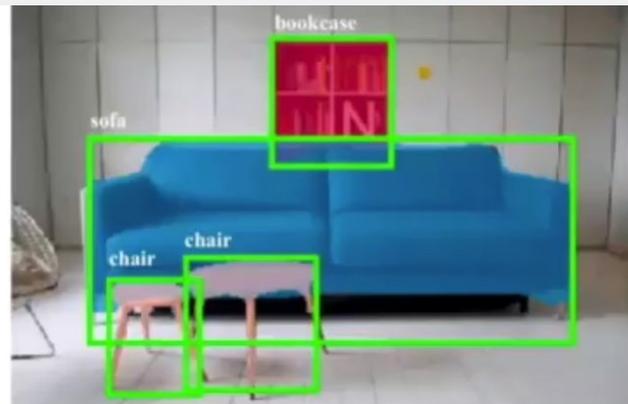
A set of detected objects

For each object:

- Bounding box
- Category label
- Instance segmentation
- 3D triangle mesh

Mask R-CNN

Mesh head



Attach a customized head that operates on each RoI coming out of Mask R-CNN

3D Shape Prediction: Mesh R-CNN

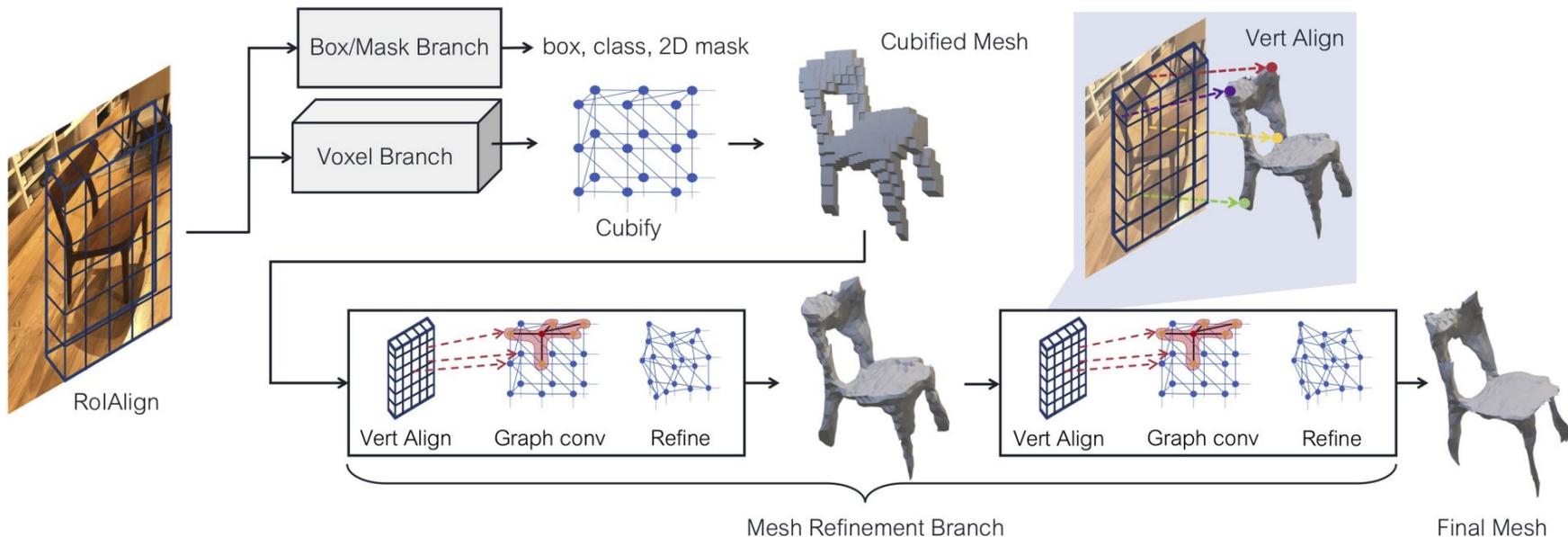
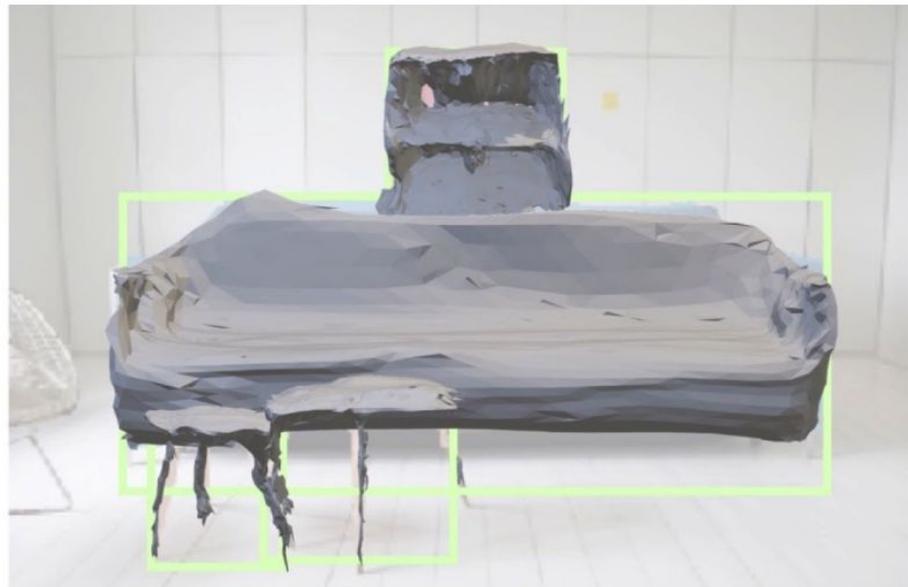
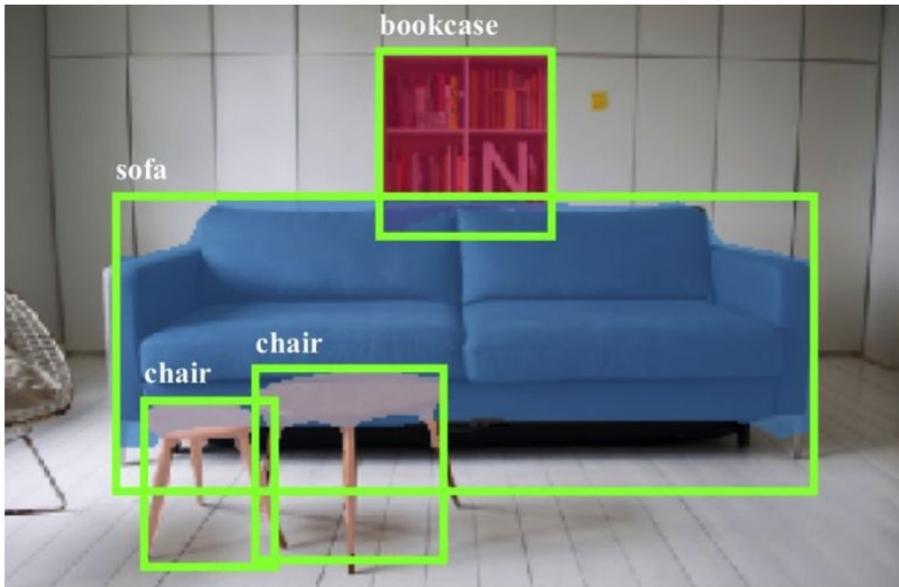


Figure 3. System overview of Mesh R-CNN. We augment Mask R-CNN with 3D shape inference. The *voxel branch* predicts a coarse shape for each detected object which is further deformed with a sequence of refinement stages in the *mesh refinement branch*.

3D Shape Prediction: 2024



https://proceedings.neurips.cc/paper_files/paper/2024/file/29c8c615b3187ee995029284702d3f43-Paper-Conference.pdf (2024 NIPS, 3D Scene Diffusion)



3D Shape Prediction: 2025

<https://github.com/facebookresearch/sam-3d-objects> (SAM 3D: 3Dfy Anything in Images, Nov 2025)

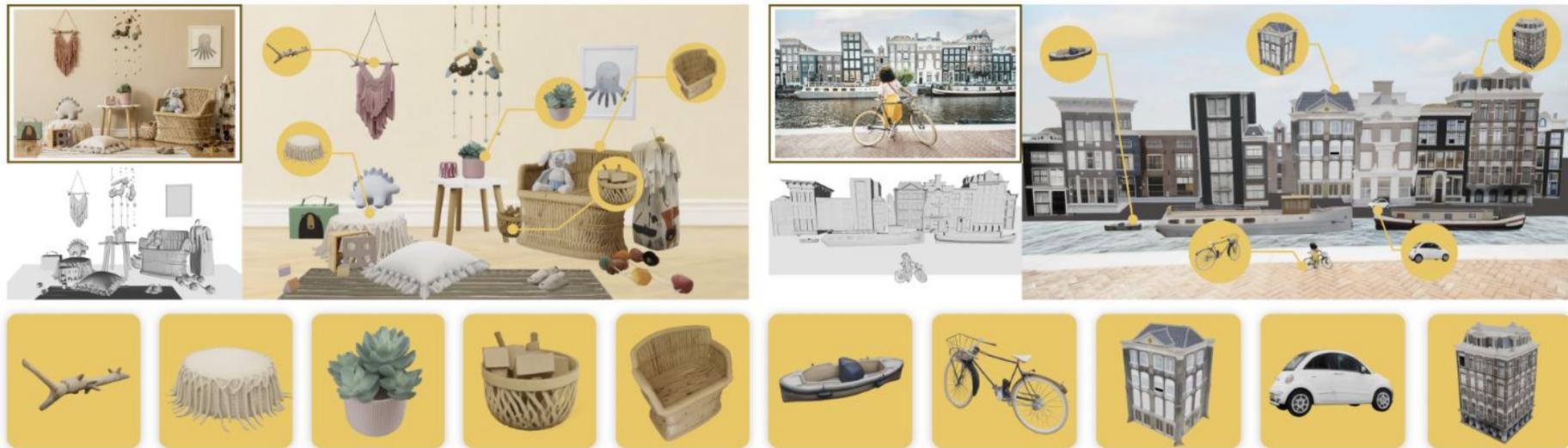


Figure 1 SAM 3D converts a single image into a composable 3D scene made of individual objects: Our method predicts per-object geometry, texture, and layout, enabling full scene reconstruction. Bottom: high-quality 3D assets recovered for each object.

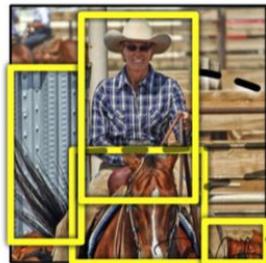
Summary

- P3 released, Due March 8, 2026
Start NOW!!!
- Also, Canvas Quizzes

“R-CNN Family”

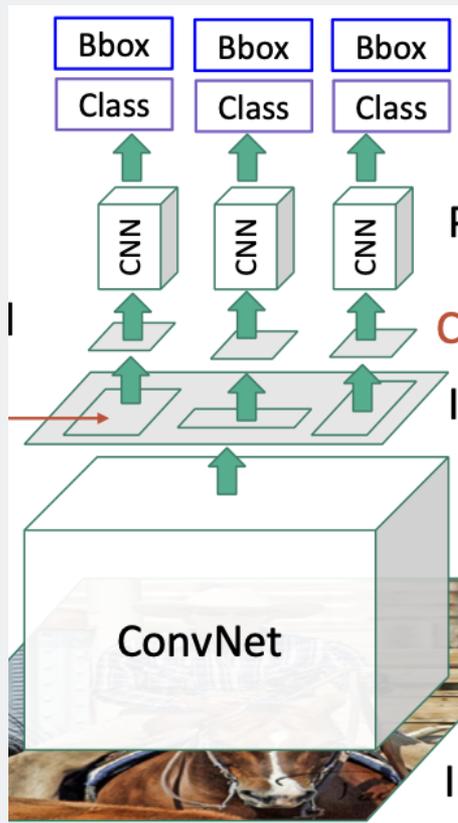


1. Input images



2. Extract region proposals (~2k)

Region proposals



ROI Pool
ROI Align

RPNs

Instance Segmentation

