

Computer Vision

Naeemullah Khan
naeemullah.khan@kaust.edu.sa



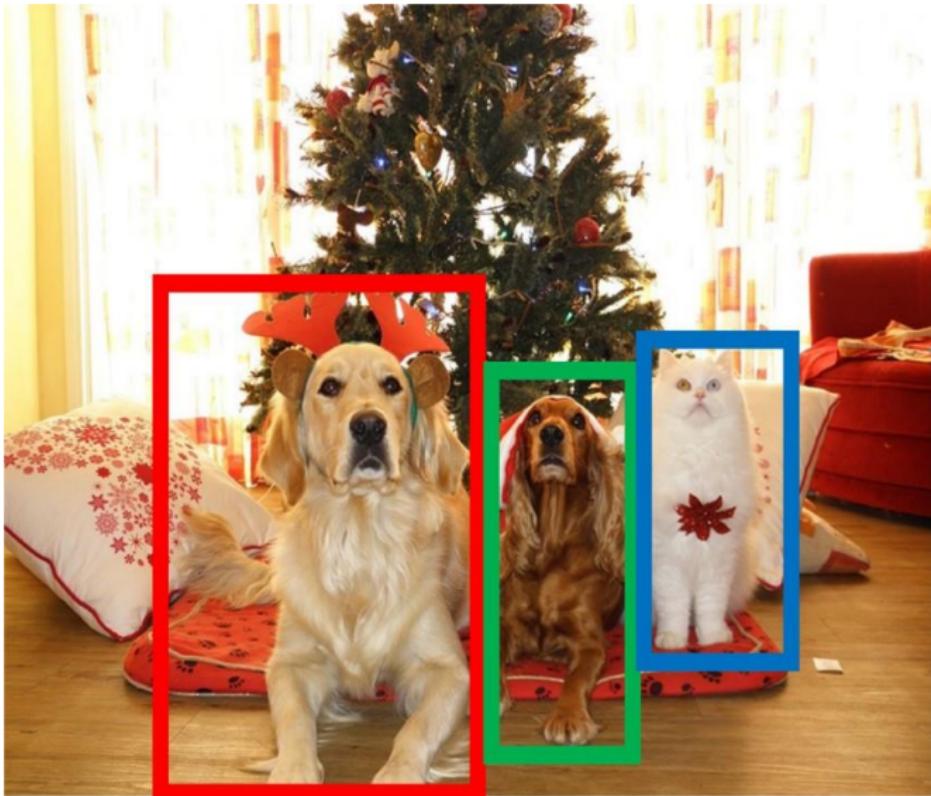
جامعة الملك عبد الله
لعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

January 07, 2024

- ▶ **Input:** Single RGB Image
- ▶ **Output:** A set of detected objects. For each object predict:
 - Category label (from fixed, known set of categories)
 - Bounding box (four numbers: x, y, width, height)

Object Detection (cont.)



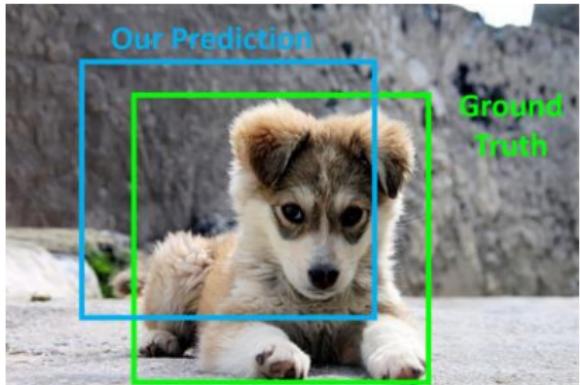
Object Detection: Challenges



- ▶ **Multiple outputs:** Need to output variable numbers of objects per image
- ▶ **Multiple types of output:** Need to predict "what" (category label) as well as "where" (bounding box)
- ▶ **Large images:** Classification works at 224x224; need higher resolution for detection, often $\sim 800 \times 600$

Comparing Boxes: Intersection over Union (IoU)

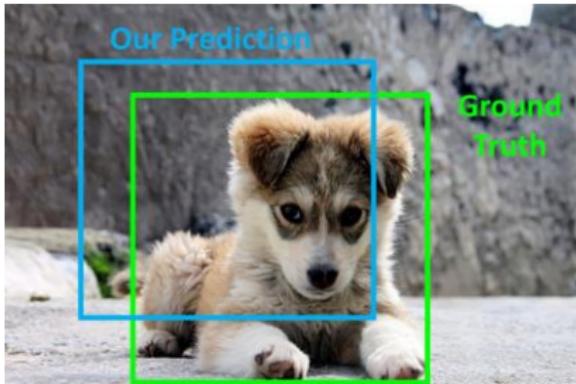
- ▶ How can we compare our prediction to the ground-truth box?



Comparing Boxes: Intersection over Union (IoU)

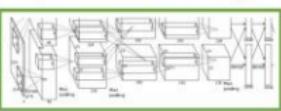
- ▶ How can we compare our prediction to the ground-truth box?
- ▶ **Intersection over Union (IoU)**
(Also called "Jaccard similarity" or "Jaccard index"):

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



Detecting a Single Object

Often pretrained
on ImageNet
(Transfer learning)

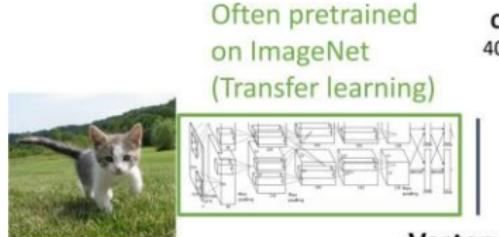


[This image](#) is CC0 public domain

Vector:
4096

Treat localization as a
regression problem!

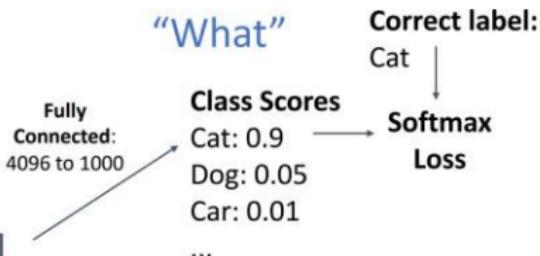
Detecting a Single Object (cont.)



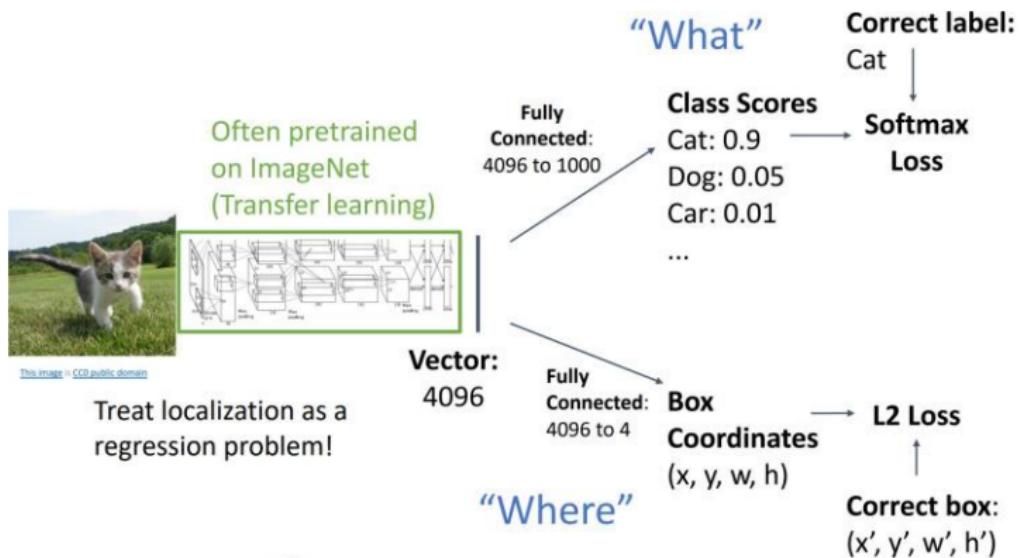
[This image is CC0 public domain](#)

Treat localization as a regression problem!

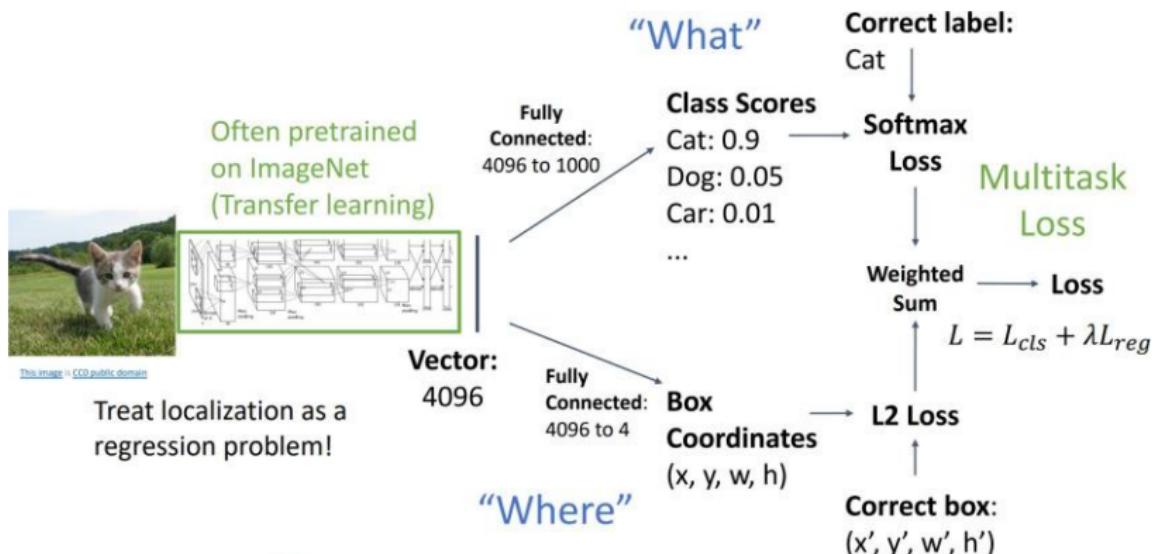
Vector:
4096



Detecting a Single Object (cont.)



Detecting a Single Object (cont.)



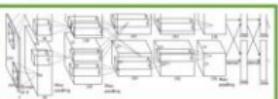
This image is CC0 public domain

Treat localization as a regression problem!

Detecting a Single Object (cont.)



Often pretrained
on ImageNet
(Transfer learning)



This image is CC0 public domain

Treat localization as a
regression problem!

Problem: Images can have
more than one object!

Vector:
4096

Fully
Connected:
4096 to 1000

“What”

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

Fully
Connected:
4096 to 4

“Where”

Box
Coordinates
(x, y, w, h)

Correct label:
Cat

Softmax
Loss

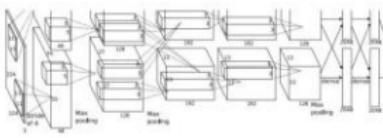
Multitask
Loss

Weighted
Sum

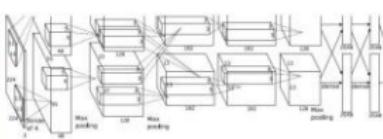
$$L = L_{cls} + \lambda L_{reg}$$

Correct box:
(x', y', w', h')

Multiple Objects



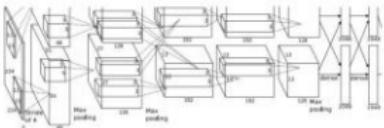
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



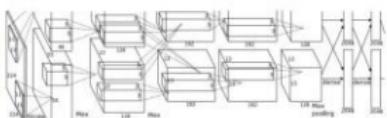
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

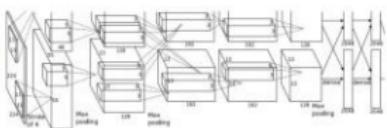
Multiple Objects (cont.)



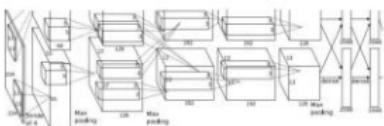
Each image needs a different number of outputs!



CAT: (x, y, w, h) **4 numbers**



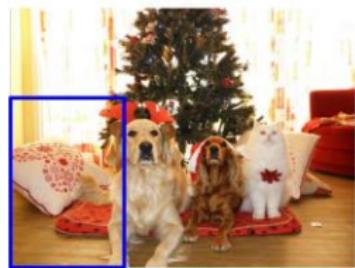
DOG: (x, y, w, h)
DOG: (x, y, w, h) **12 numbers**
CAT: (x, y, w, h)



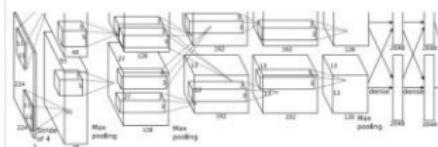
DUCK: (x, y, w, h) **Many numbers!**
DUCK: (x, y, w, h) **Many numbers!**

....

Detecting Multiple Objects: Sliding Window



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

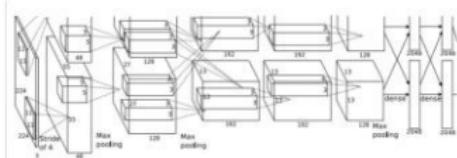


Dog? NO
Cat? NO
Background? YES

Detecting Multiple Objects: Sliding Window (cont.)



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

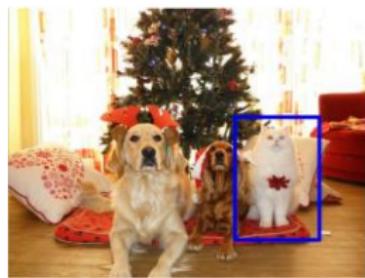


Dog? YES

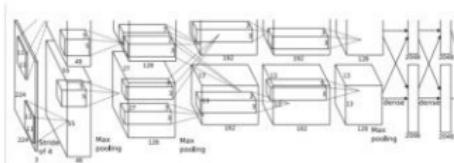
Cat? NO

Background? NO

Detecting Multiple Objects: Sliding Window (cont.)



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



Dog? NO

Cat? YES

Background? NO

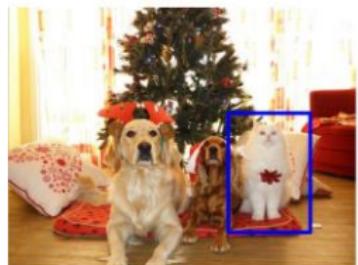
Detecting Multiple Objects: Sliding Window (cont.)



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

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

Detecting Multiple Objects: Sliding Window (cont.)



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

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

Consider a box of size $h \times w$:

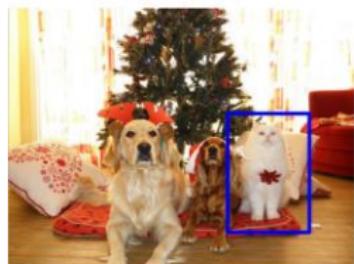
Possible x positions: $W - w + 1$

Possible y positions: $H - h + 1$

Possible positions:

$$(W - w + 1) * (H - h + 1)$$

Detecting Multiple Objects: Sliding Window (cont.)



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

Question: How many possible boxes are there in an image of 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:

$$(W - w + 1) * (H - h + 1)$$

Total possible boxes:

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

$$= \frac{H(H+1)}{2} \frac{W(W+1)}{2}$$

Detecting Multiple Objects: Sliding Window (cont.)



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

800 x 600 image
has ~58M boxes!
No way we can evaluate them all

Question: How many possible boxes are there in an image of 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:

$$(W - w + 1) * (H - h + 1)$$

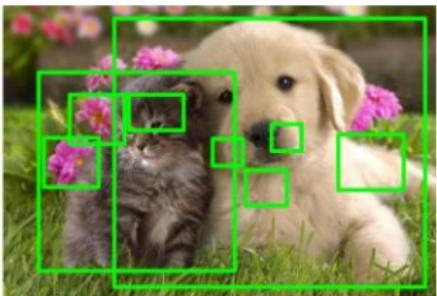
Total possible boxes:

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

$$= \frac{H(H+1)}{2} \frac{W(W+1)}{2}$$

Region Proposal

- ▶ 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 gives 2000 region proposals in a few seconds on CPU





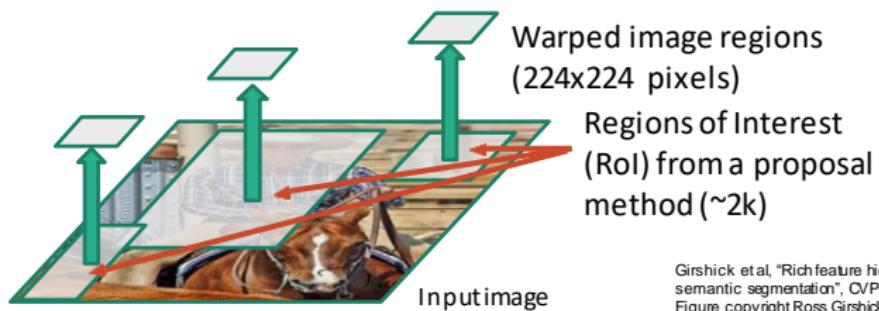
Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



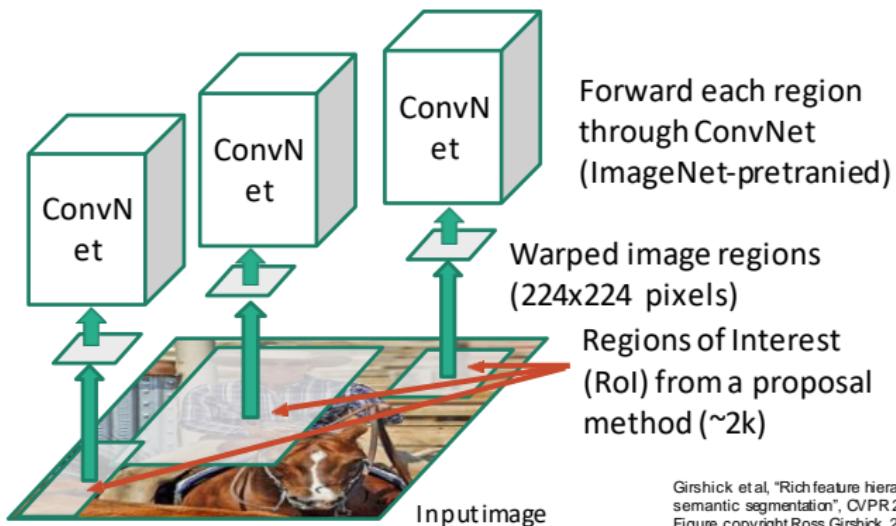
Regions of Interest
(RoI) from a proposal
method (~2k)

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



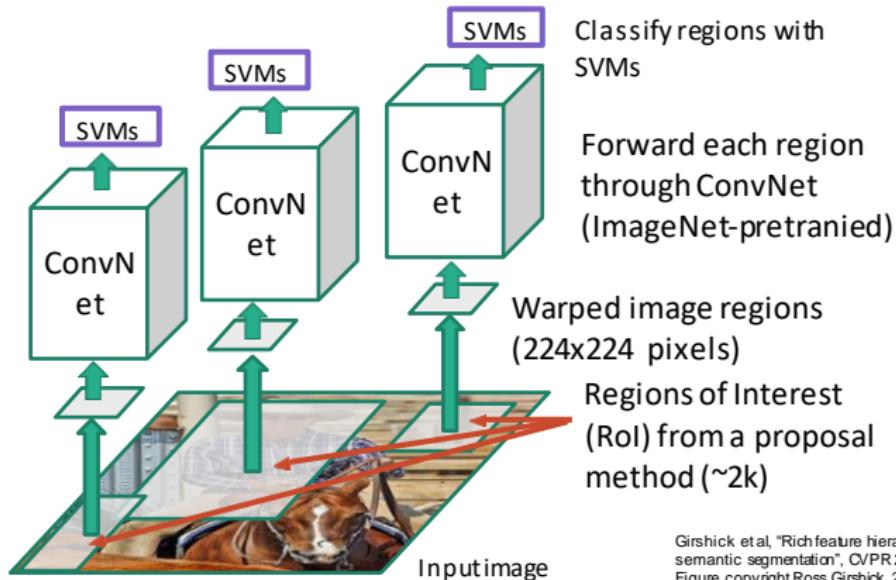
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



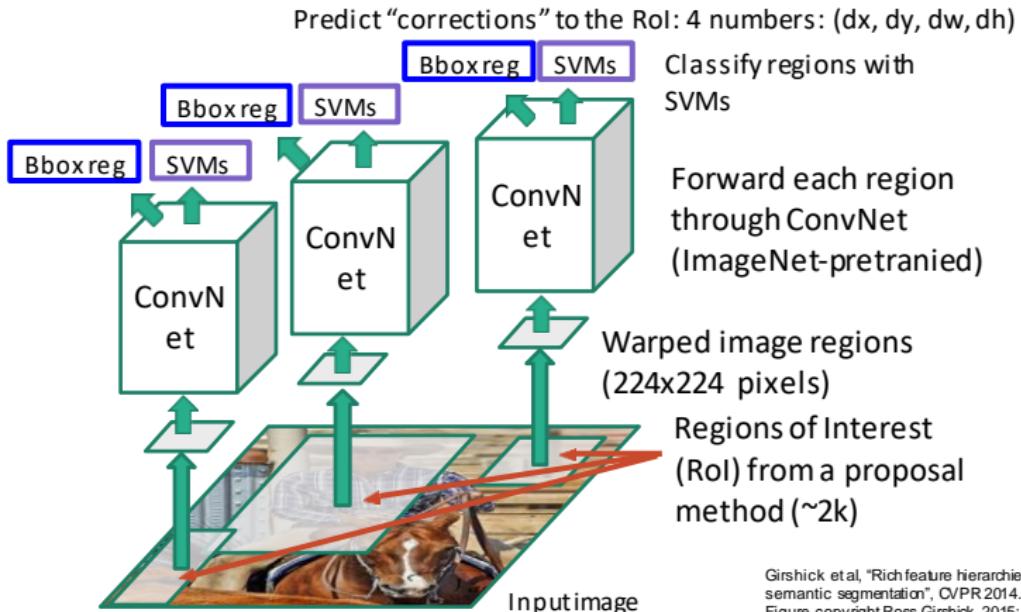
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



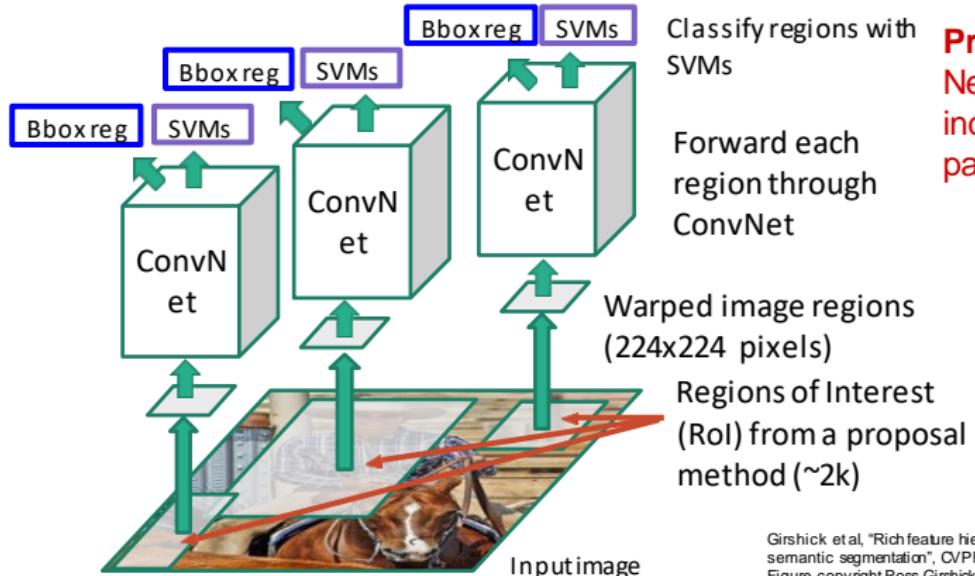
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



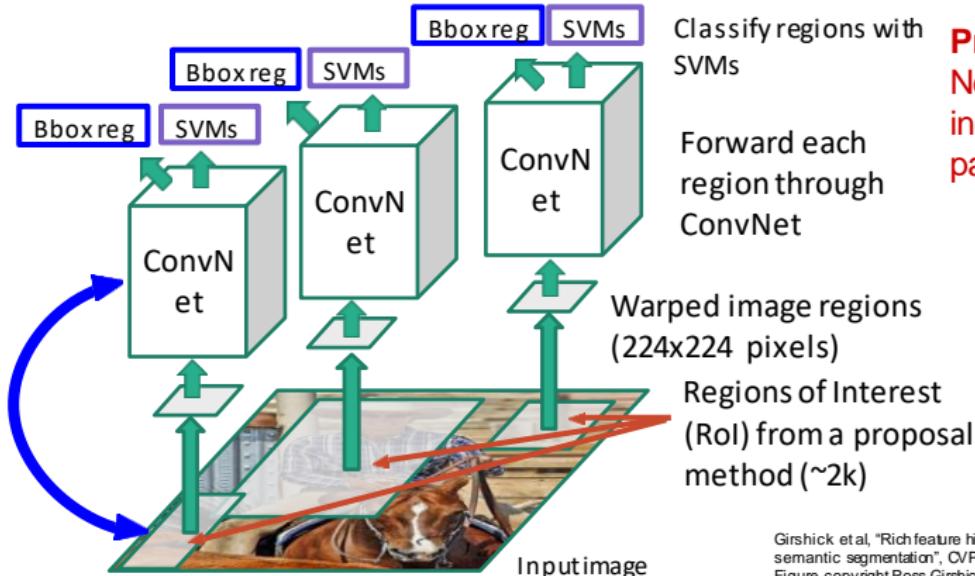
Problem: Very slow!
Need to do ~2k independent forward passes for each image!

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

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

“Slow” R-CNN

Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



Problem: Very slow!
Need to do ~2k independent forward passes for each image!

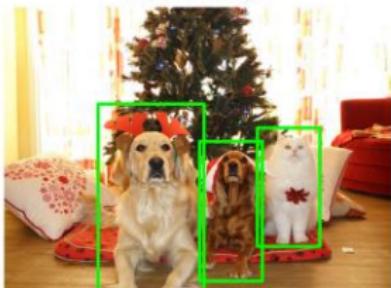
Idea: Pass the image through convnet before cropping! Crop the conv feature instead!

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

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

"Slow" R-CNN Training

Input Image

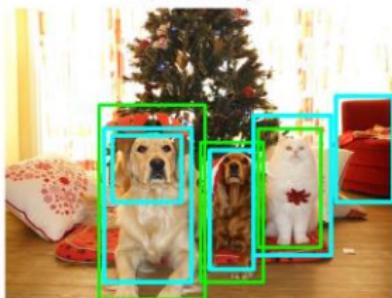


Ground-Truth boxes

[This image is CC0 public domain](#)

"Slow" R-CNN Training

Input Image



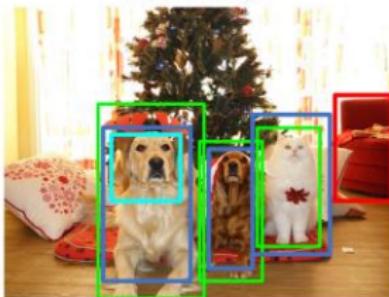
Ground-Truth boxes

Region Proposals

[This image is CC0 public domain](#)

"Slow" R-CNN Training

Input Image

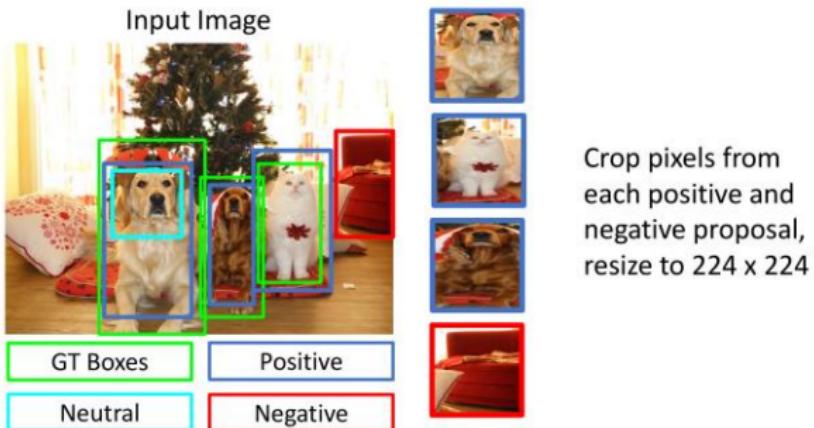


Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes

GT Boxes	Positive
Neutral	Negative

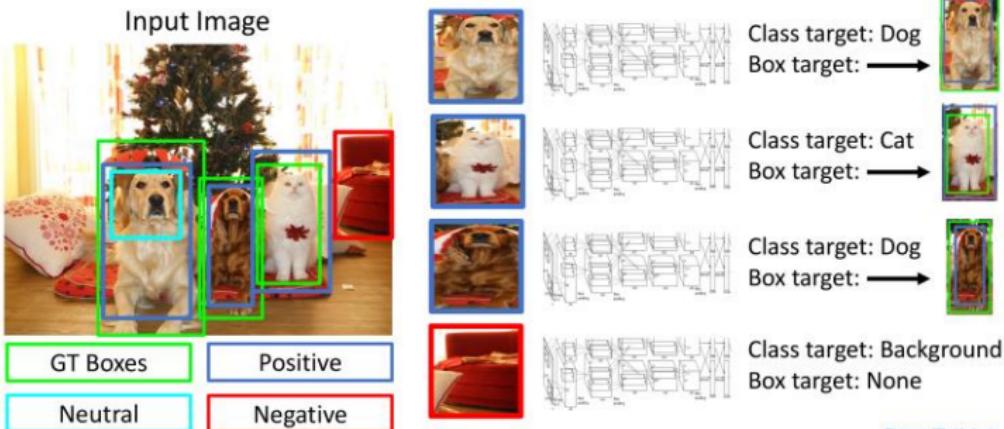
[This image is CC0 public domain](#)

"Slow" R-CNN Training



[This image is CC0 public domain](#)

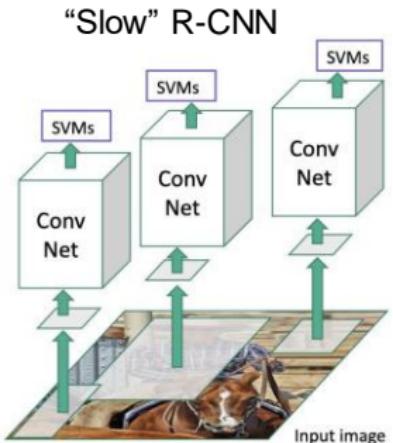
"Slow" R-CNN Training



Fast R-CNN

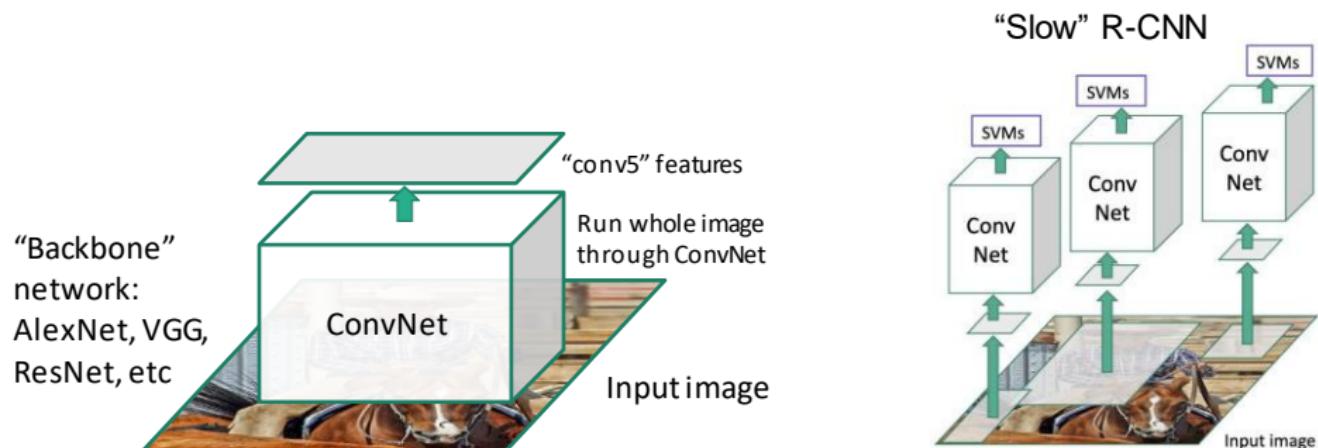


Input image

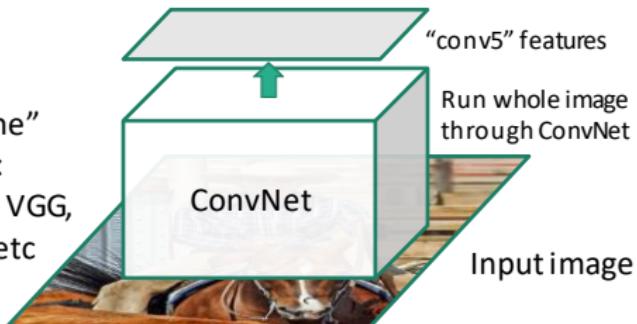


Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN

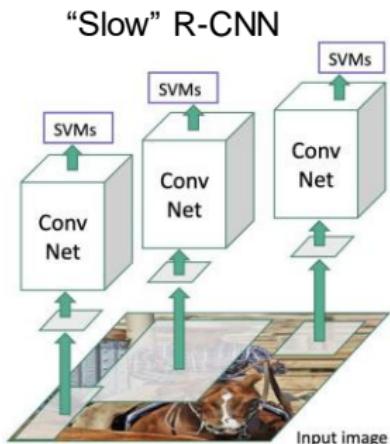
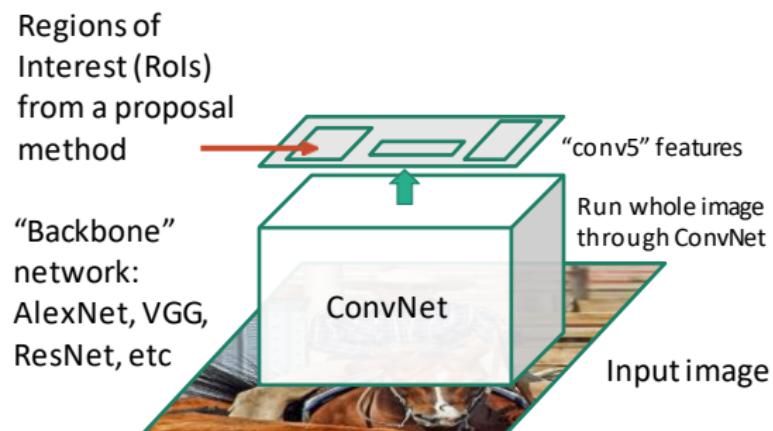


"Backbone"
network:
AlexNet, VGG,
ResNet, etc



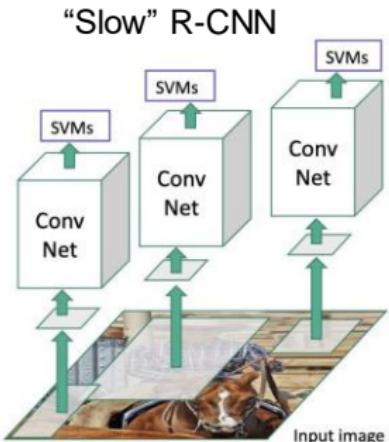
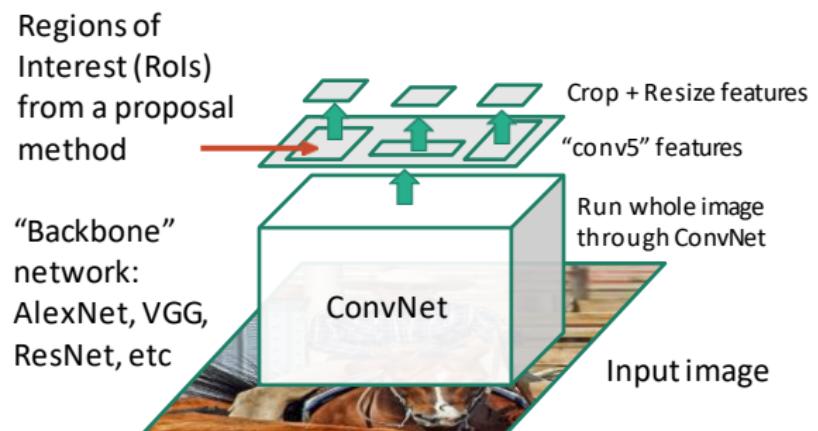
Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN



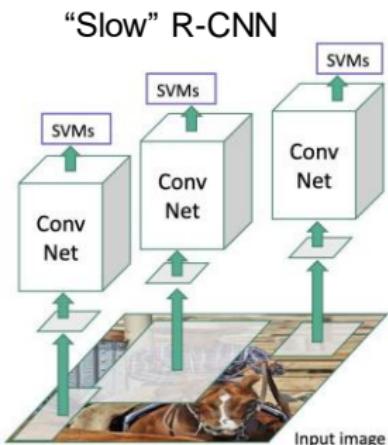
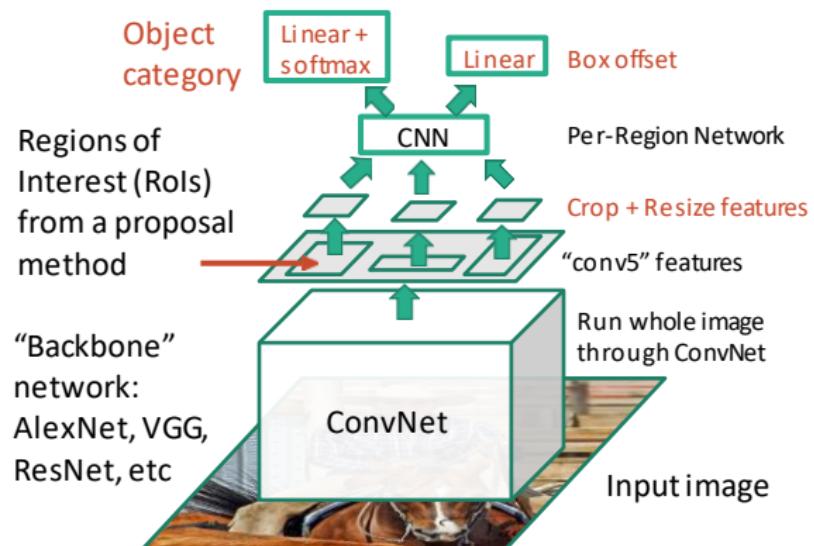
Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN



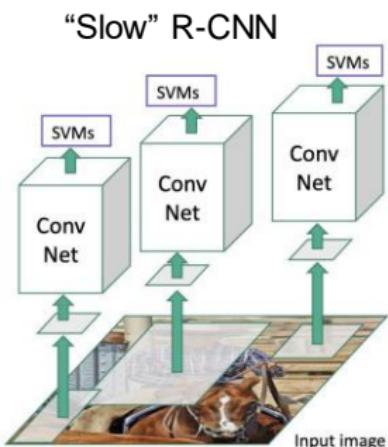
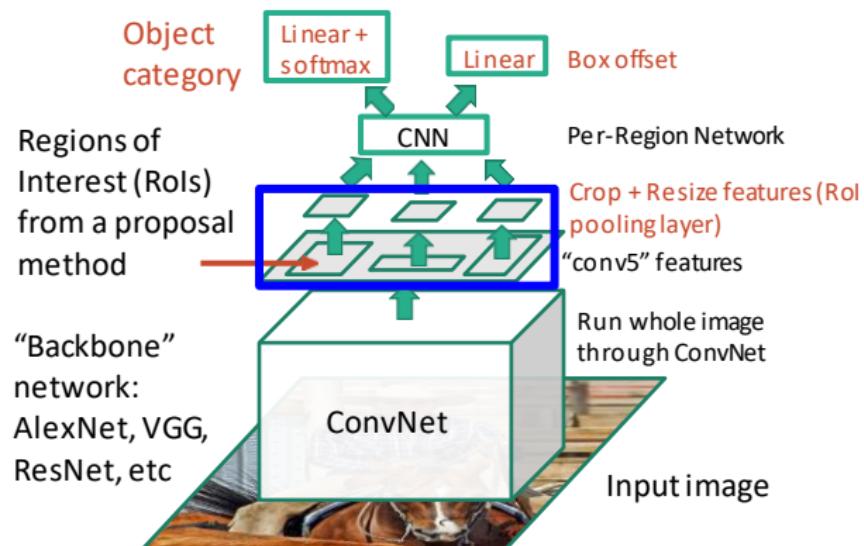
Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN



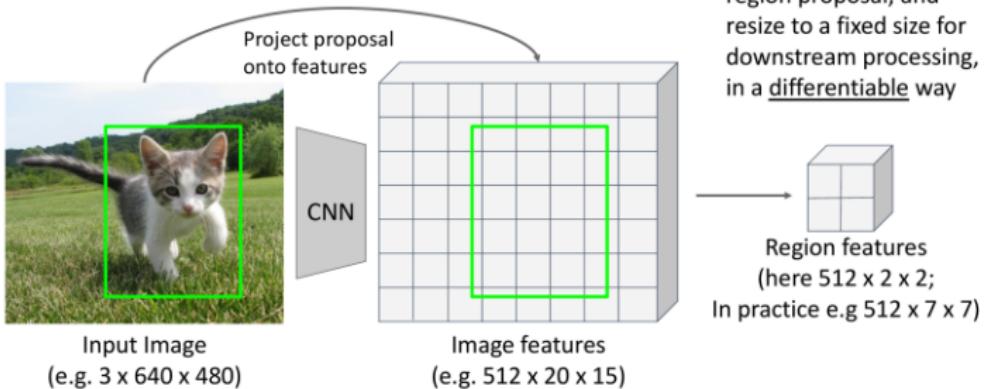
Girshick, “Fast R-CNN”, ICCV 2015. Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN



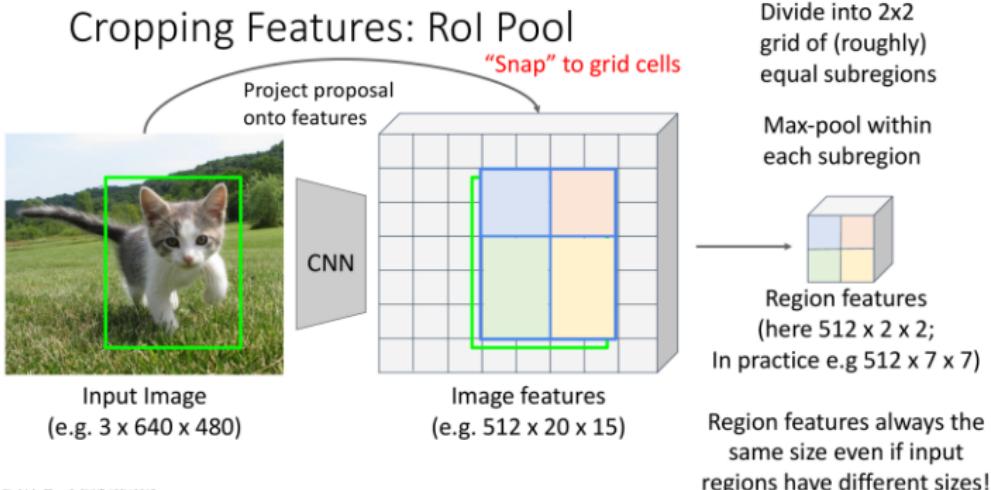
Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source Reproduced with permission.

Cropping Features



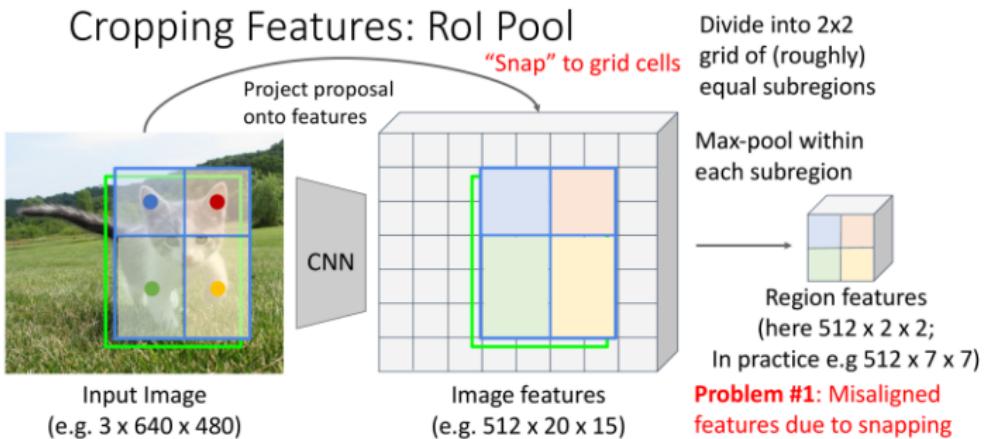
Girshick, "Fast R-CNN", ICCV 2015.

RoI pooling

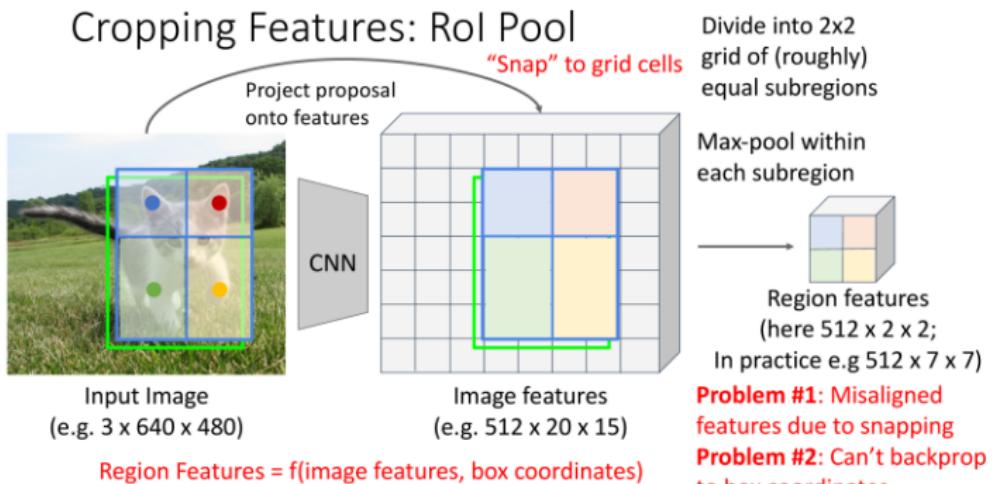


Girshick, "Fast R-CNN", ICCV 2015.

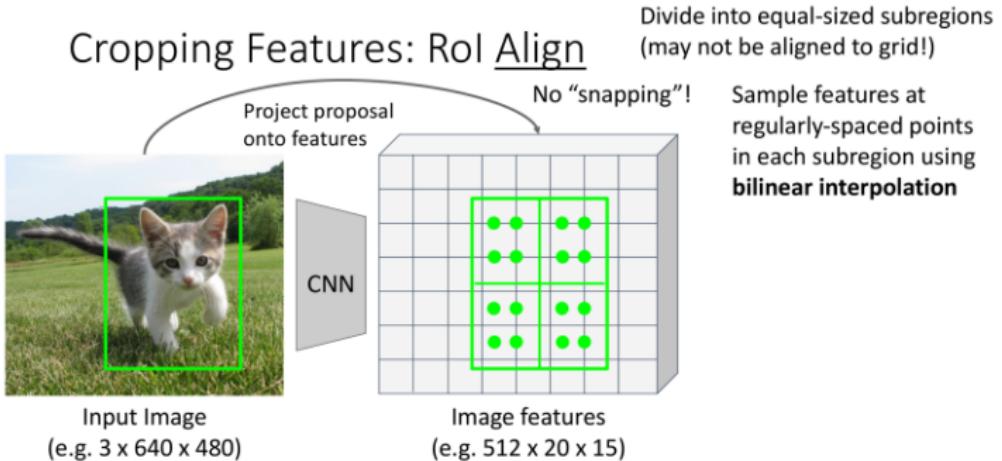
RoI pooling



RoI pooling



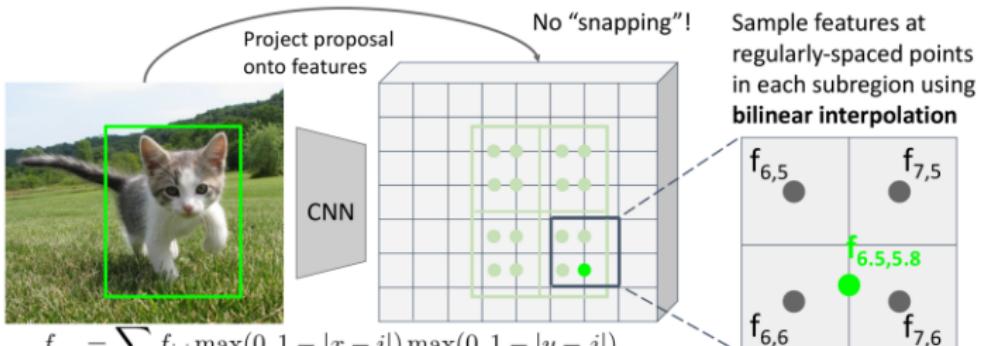
Cropping Features: RoI Align



He et al., "Mask R-CNN", ICCV 2017

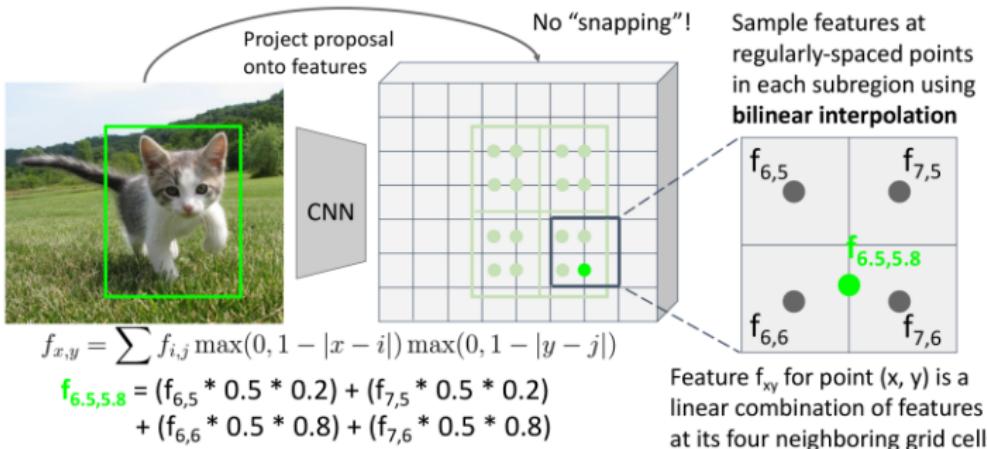
Cropping Features: RoI Align

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

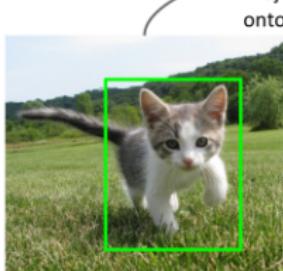


Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



Cropping Features: RoI Align



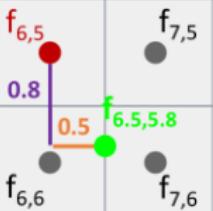
Project proposal onto features

CNN

No “snapping”!

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

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

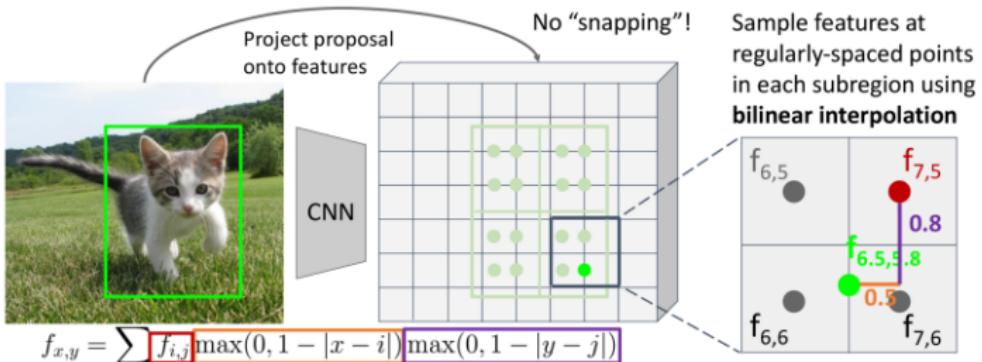


$$f_{x,y} = \sum [f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)]$$

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

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

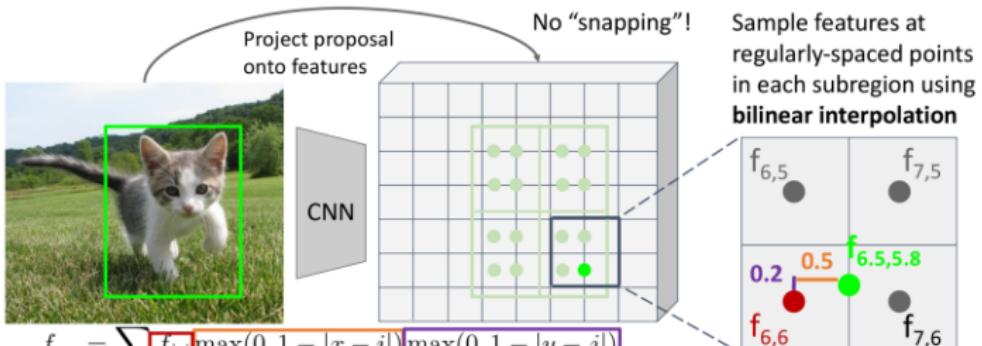
Cropping Features: RoI Align



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

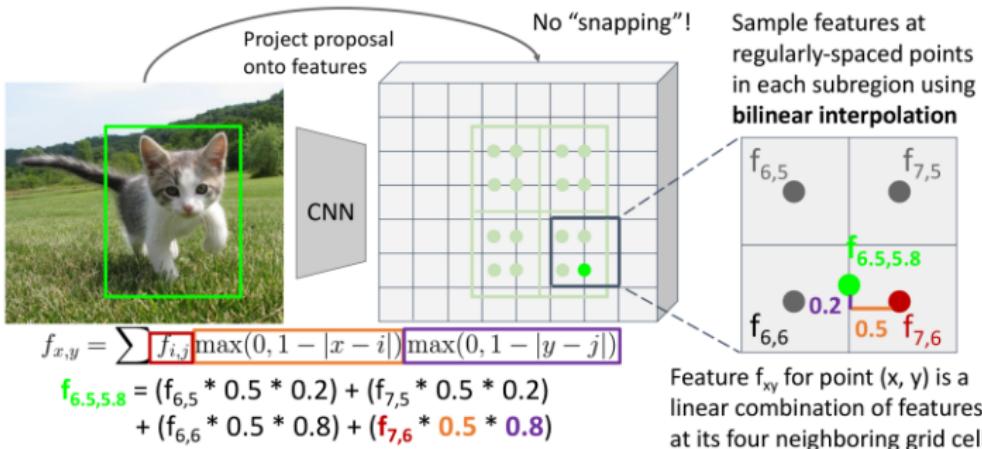
Cropping Features: RoI Align

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



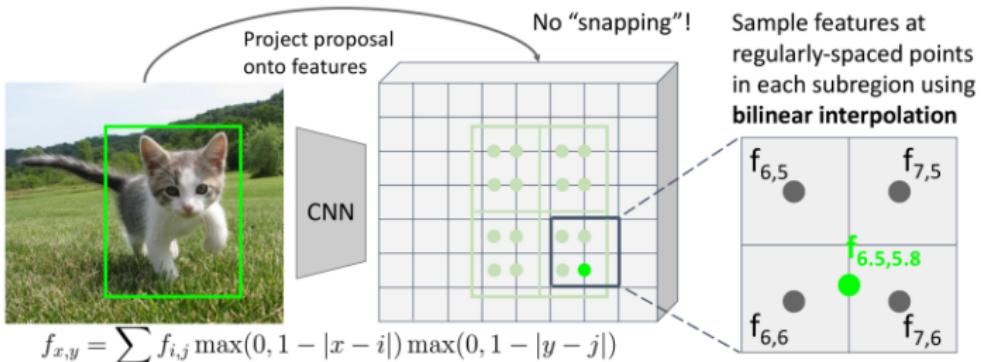
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

Cropping Features: RoI Align



Cropping Features: RoI Align

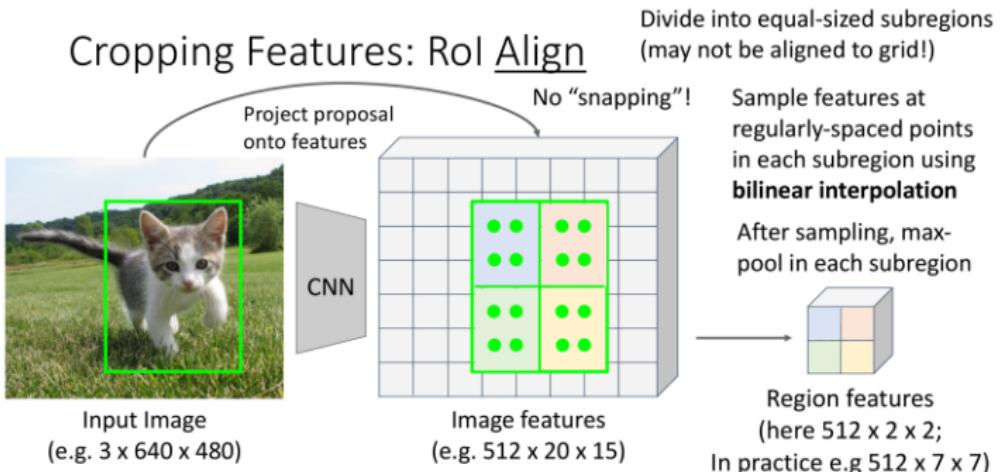
Divide into equal-sized subregions
(may not be aligned to grid!)



This is differentiable! Upstream gradient for sampled feature will flow backward into each of the four nearest-neighbor gridpoints

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

Cropping Features: RoI Align



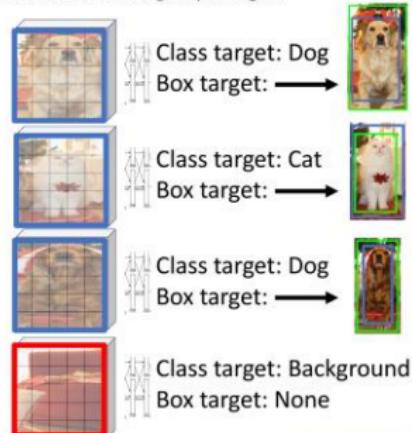
Output features now aligned to input box! And we can backprop to box coordinates!

Fast R-CNN Training



Image Features

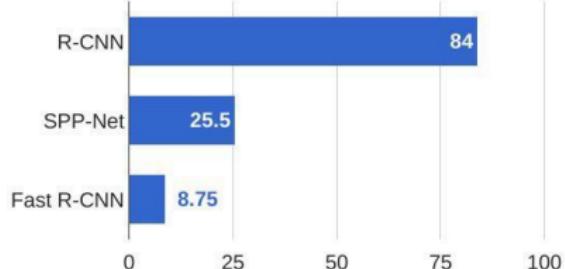
Crop features for each region, use them to predict class and box targets per region



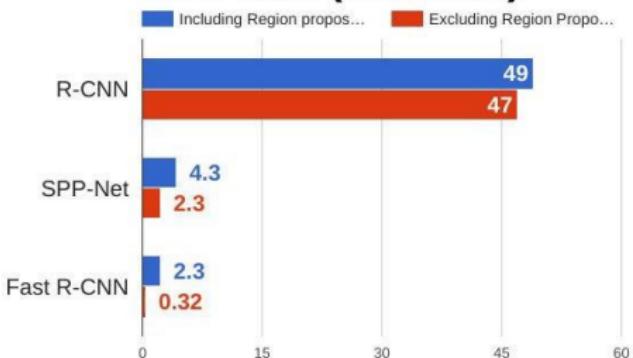
This image is CC0 public domain

R-CNN vs Fast R-CNN

Training time (Hours)



Test time (seconds)

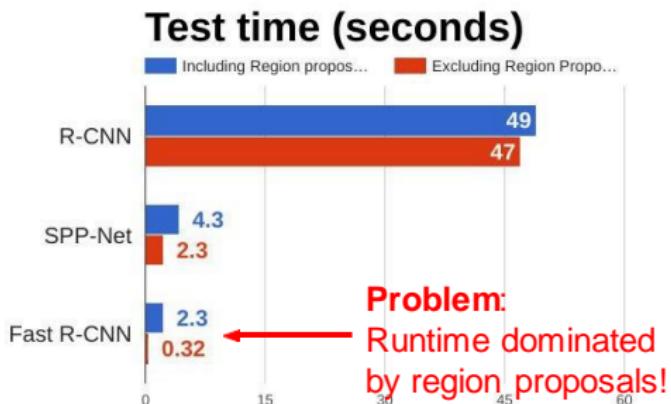


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

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

Girshick, "Fast R-CNN", ICCV 2015

R-CNN vs Fast R-CNN



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

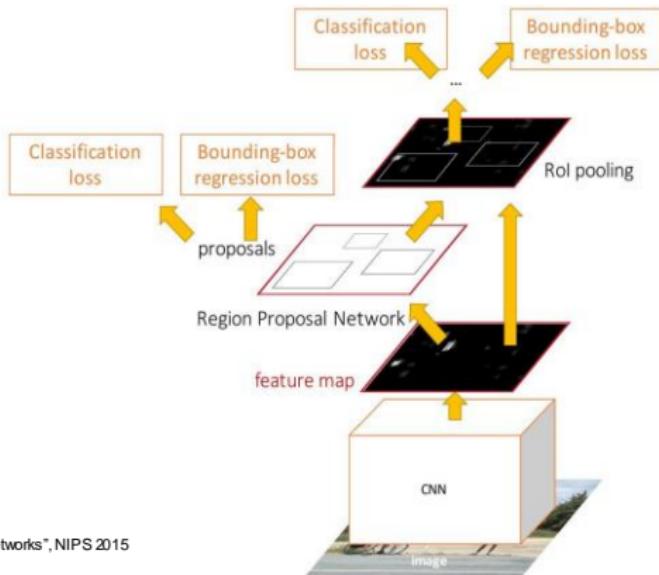
He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014

Girshick, "Fast R-CNN", ICCV 2015

Make CNN do 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



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

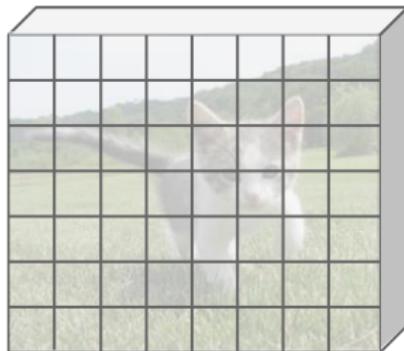
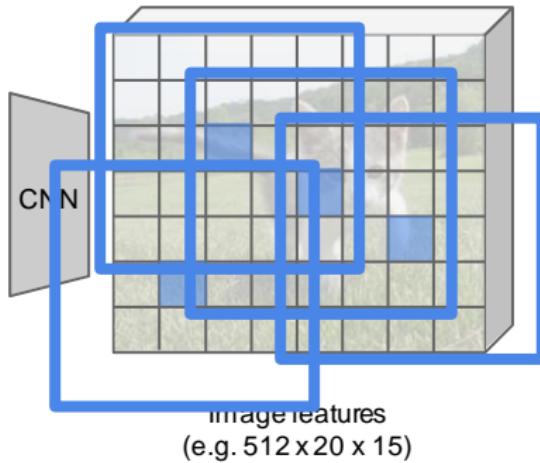


Image features
(e.g. $512 \times 20 \times 15$)

Region Proposal Network



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

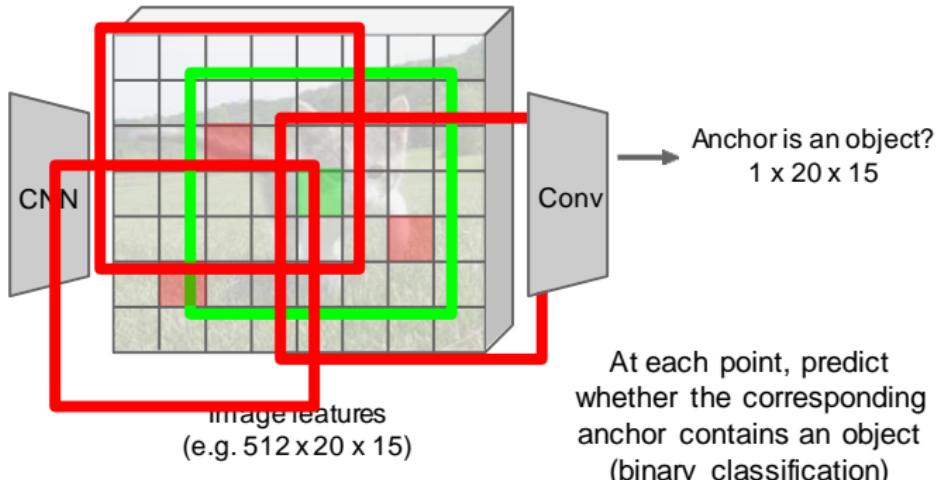


Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network



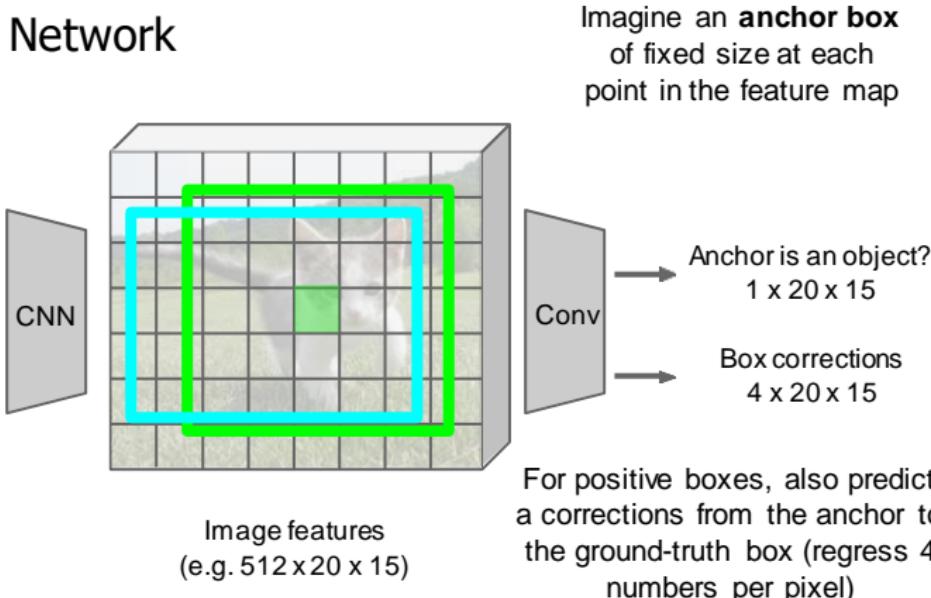
Input Image
(e.g. $3 \times 640 \times 480$)



Region Proposal Network



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

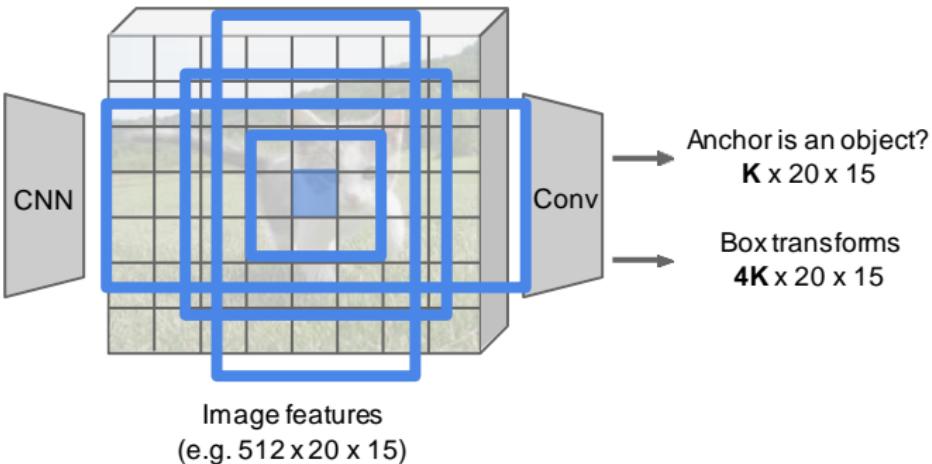


Region Proposal Network

In practice use K different anchor boxes of different size / scale at each point



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

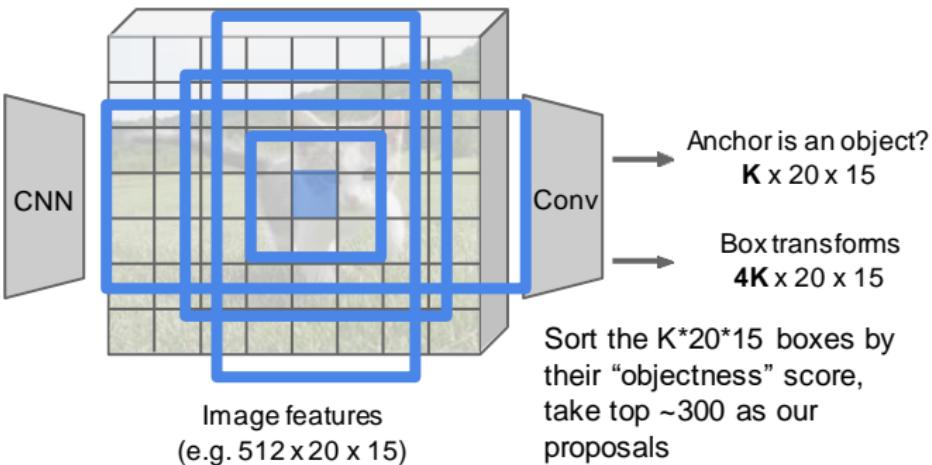


Region Proposal Network

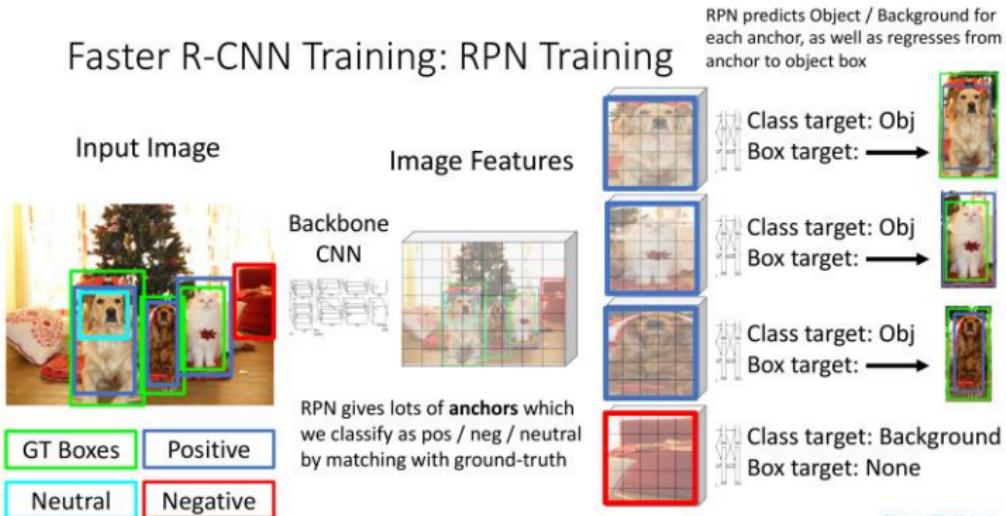
In practice use K different anchor boxes of different size / scale at each point



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

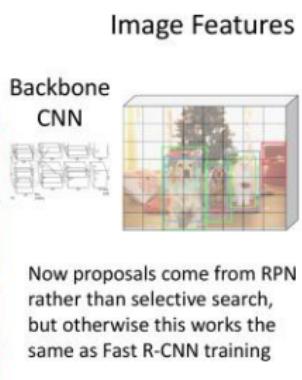


Faster R-CNN Training: RPN Training

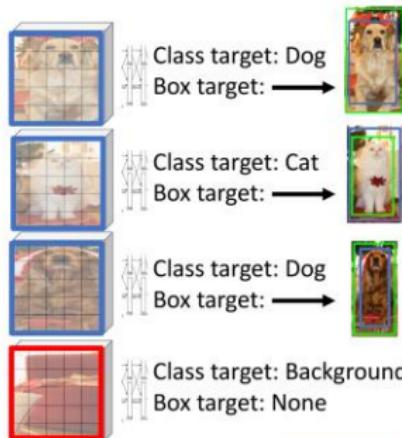


[This image is CC0 public domain](#)

Faster R-CNN Training: Stage 2



Crop features for each proposal, use them to predict class and box targets per region

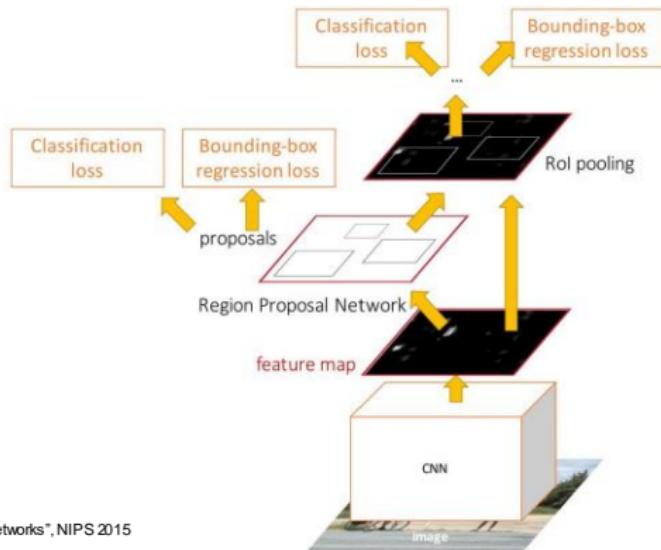


This image is CC0 public domain

Make CNN do proposals!

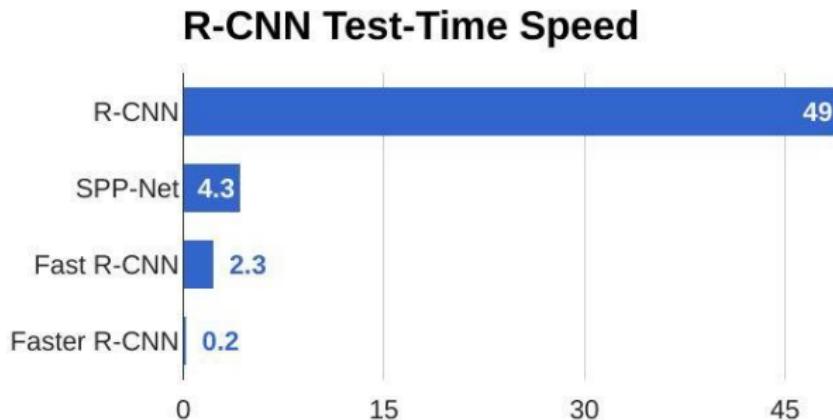
Jointly train with 4 losses:

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



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

Make CNN do proposals!



Make CNN do proposals!

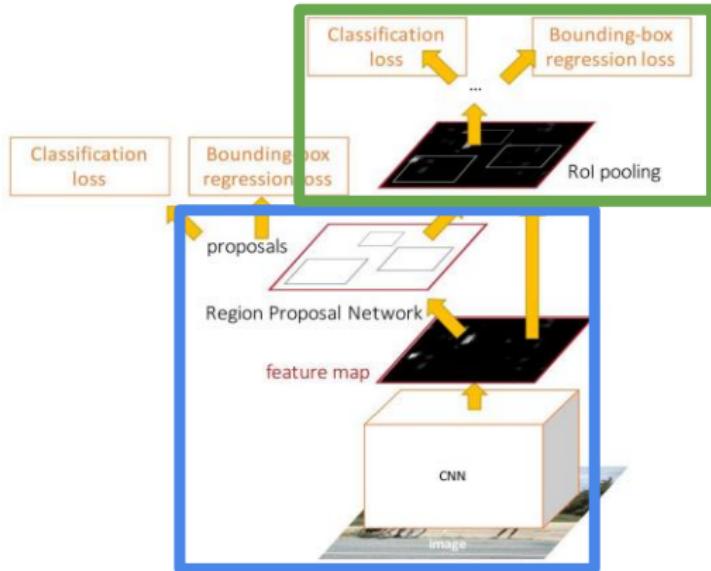
Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

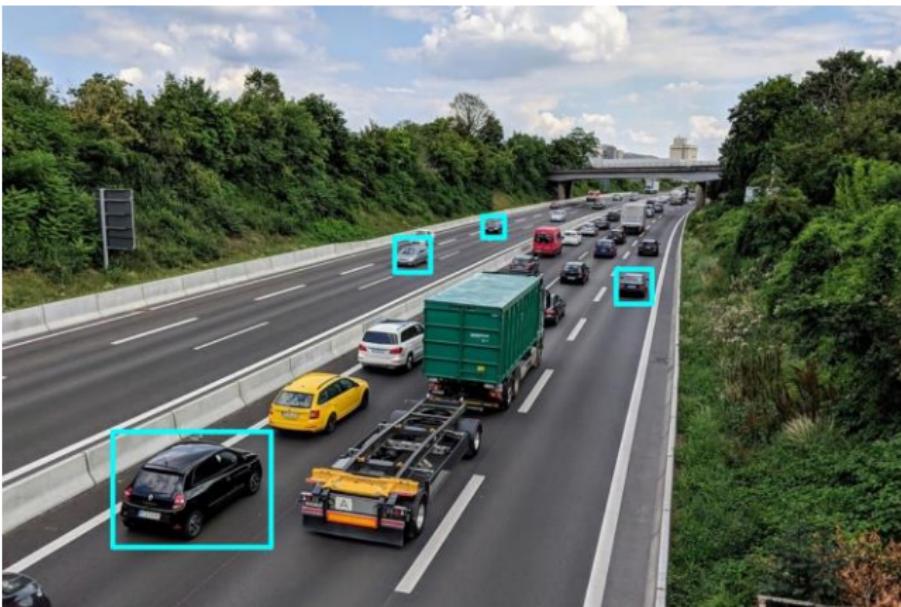
Second stage: Run once per region

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



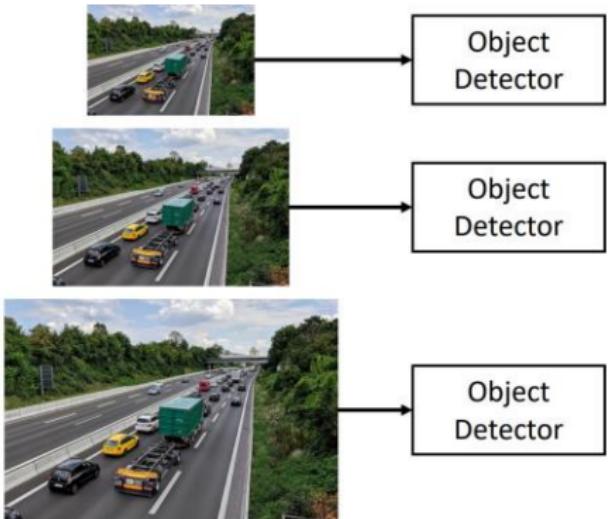
Dealing with Scale

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



Dealing with Scale: Image Pyramid

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.



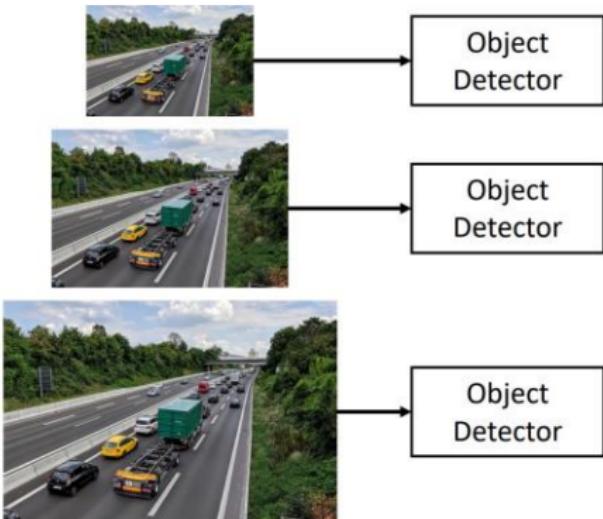
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

Dealing with Scale: Image Pyramid (cont.)

Classic idea: build an *image pyramid* by resizing the image to different scales, then process each image scale independently.

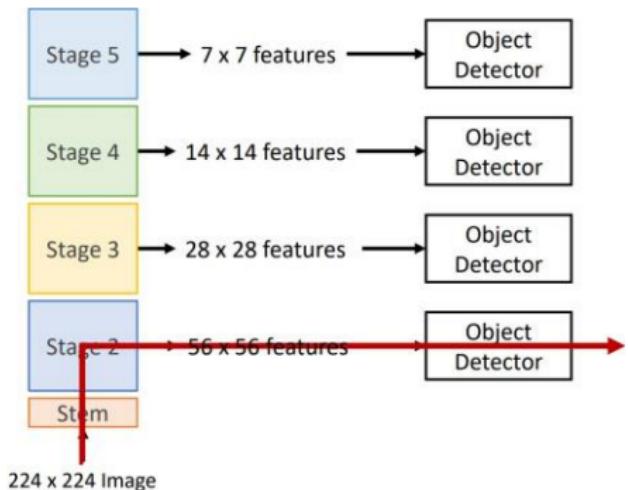
Problem: Expensive! Don't share any computation between scales

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



Dealing with Scale: Image Pyramid

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level



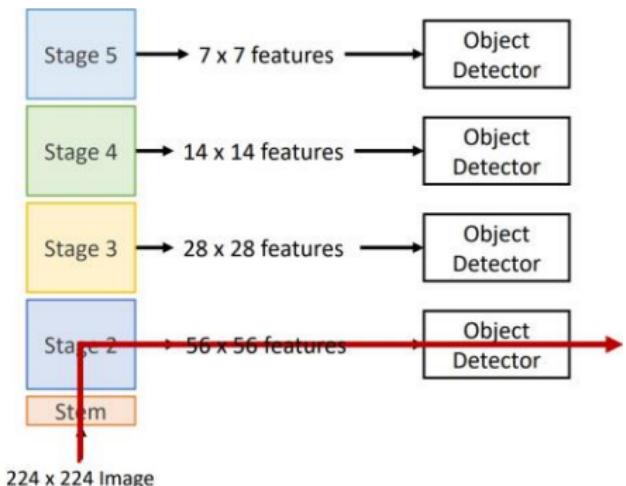
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

Dealing with Scale: Image Pyramid (cont.)

CNNs have multiple *stages* that operate at different resolutions. Attach an independent detector to the features at each level

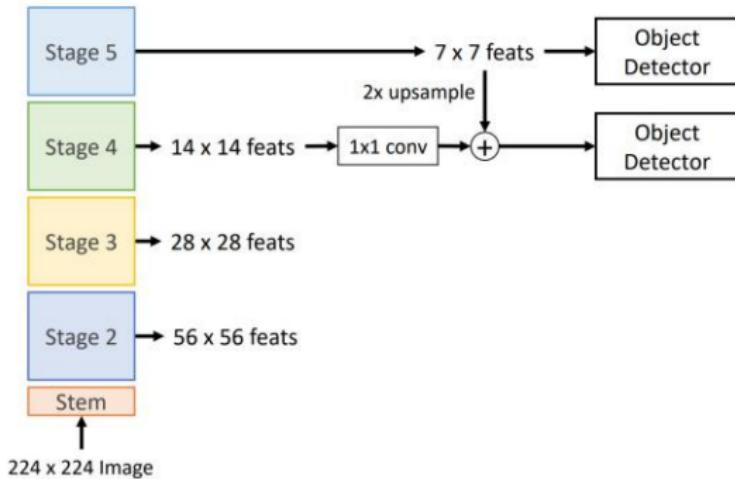
Problem: detector on early features doesn't make use of the entire backbone; doesn't get access to high-level features

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



Dealing with Scale: Feature Pyramid Network

Add *top down* connections that feed information from high level features back down to lower level features



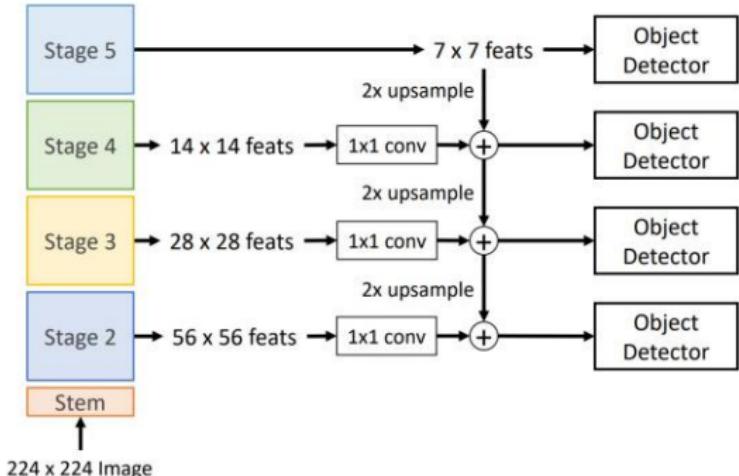
Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

Dealing with Scale: Feature Pyramid Network (cont.)

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

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017

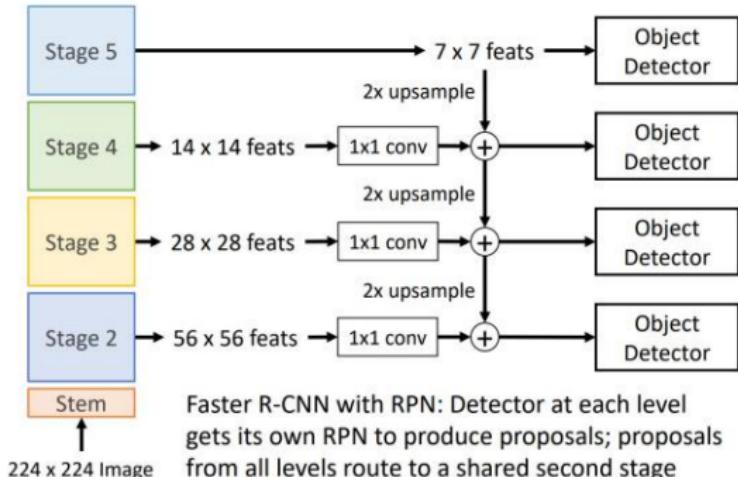


Dealing with Scale: Feature Pyramid Network (cont.)

Add top down connections that feed information from high level features back down to lower level features

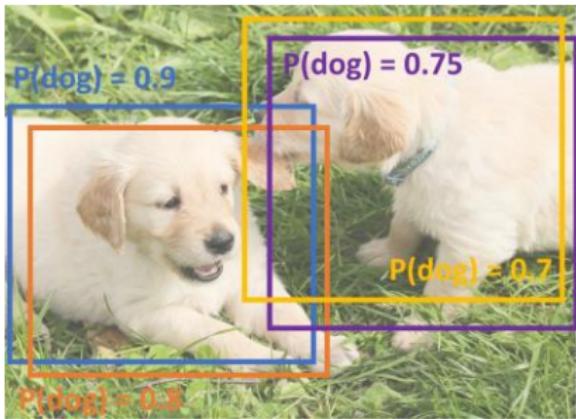
Efficient multiscale features where all levels benefit from the whole backbone! Widely used in practice

Lin et al, "Feature Pyramid Networks for Object Detection", ICCV 2017



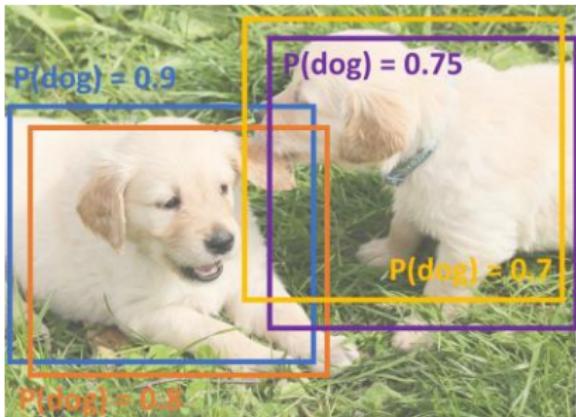
Overlapping Boxes: Non-Max Suppression (NMS)

- ▶ **Problem:** Object detectors often output many overlapping detections



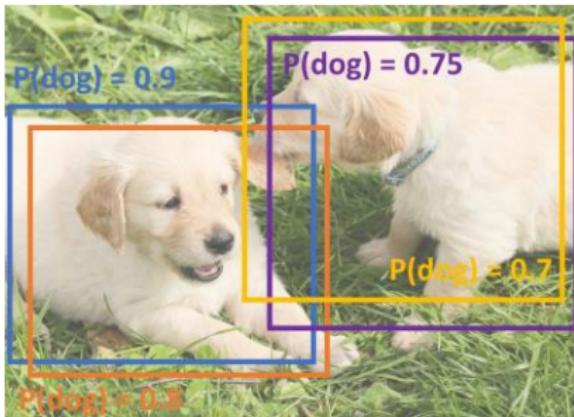
Overlapping Boxes: Non-Max Suppression (NMS)

- ▶ **Problem:** Object detectors often output many overlapping detections
- ▶ **Solution:** Post-process raw detections using Non-Max Suppression (NMS)



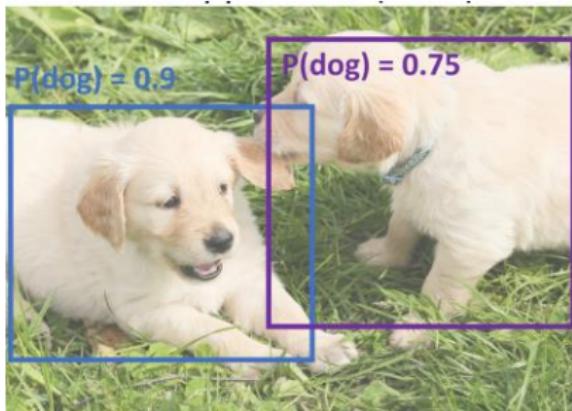
Overlapping Boxes: Non-Max Suppression (NMS)

- ▶ **Problem:** Object detectors often output many overlapping detections
- ▶ **Solution:** Post-process raw detections using Non-Max Suppression (NMS)
 1. Select next highest-scoring box
 2. Eliminate lower-scoring boxes
 3. with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
 4. If any boxes remain, GOTO 1



Overlapping Boxes: Non-Max Suppression (NMS)

- ▶ **Problem:** Object detectors often output many overlapping detections
- ▶ **Solution:** Post-process raw detections using Non-Max Suppression (NMS)
 1. Select next highest-scoring box
 2. Eliminate lower-scoring boxes
 3. with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
 4. If any boxes remain, GOTO 1





- ▶ **Problem:** Object detectors often output many overlapping detections
- ▶ **Solution:** Post-process raw detections using Non-Max Suppression (NMS)
 1. Select next highest-scoring box
 2. Eliminate lower-scoring boxes
 3. with $\text{IoU} > \text{threshold}$ (e.g. 0.7)
 4. If any boxes remain, GOTO 1
- ▶ **Problem:** NMS may eliminate "good" boxes when objects are highly overlapping. . . no good solution =(



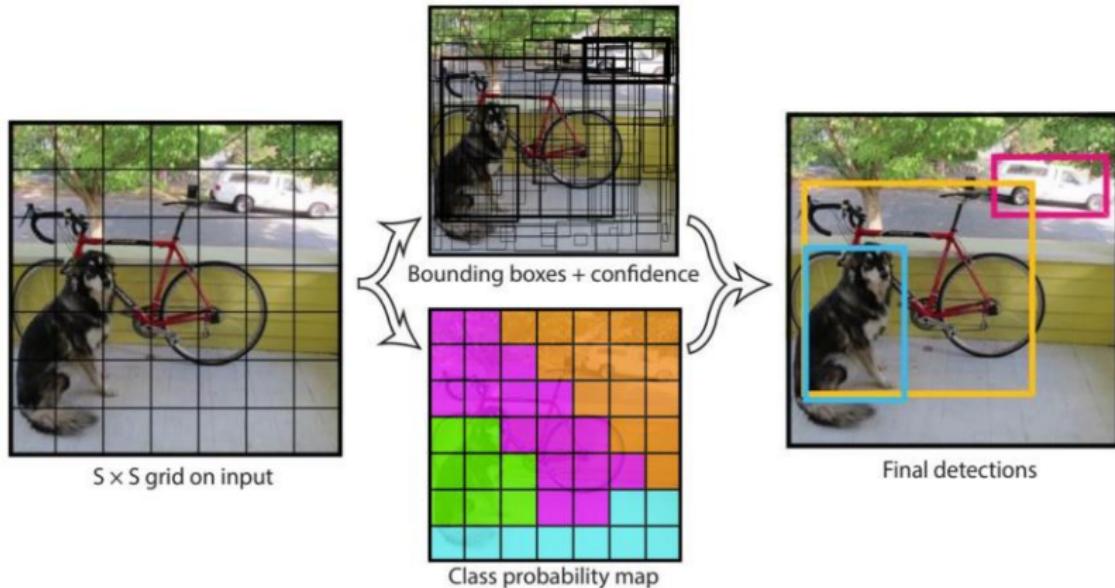
Single Shot Object Detection



Single Shot:
SSD, YOLO ...

Fast
High false rate

YOLO - Overview



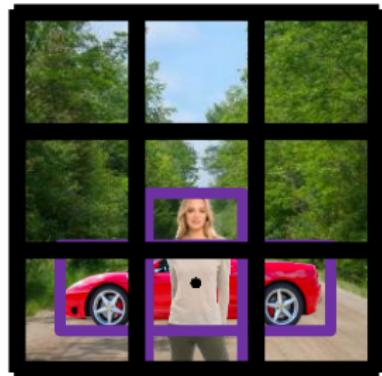
YOLO algorithm



Labels for training
For each grid cell:

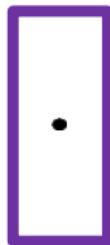
[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

YOLO algorithm



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



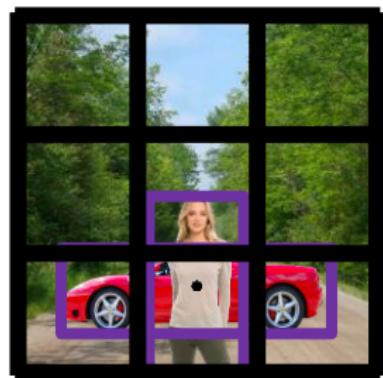
Anchor box 2:



[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

Andrew Ng

YOLO algorithm



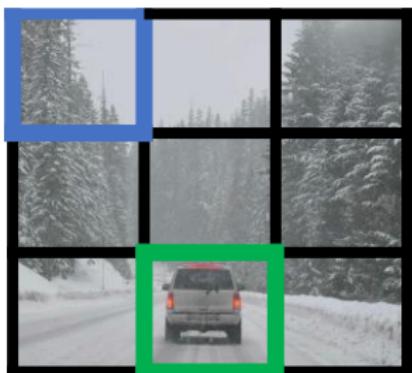
$$\mathbf{y} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1: Anchor box 2:



Andrew Ng

YOLO algorithm



- 1 - pedestrian
- 2 - car
- 3 - motorcycle

$y =$

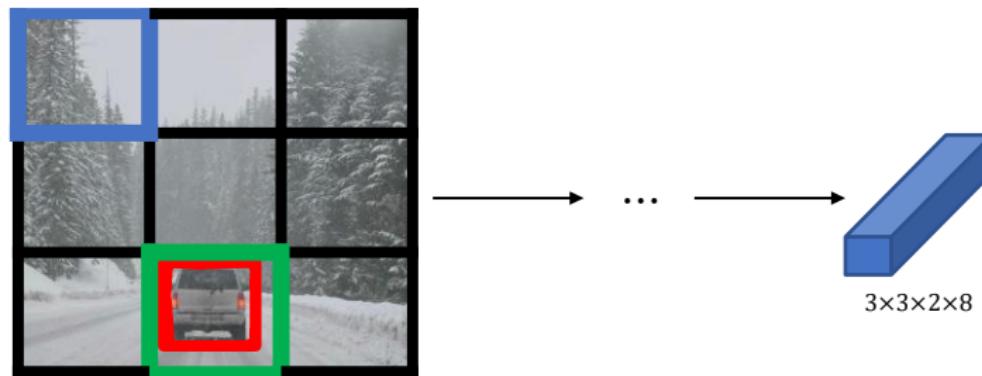
y is $3 \times 3 \times 2 \times 8$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

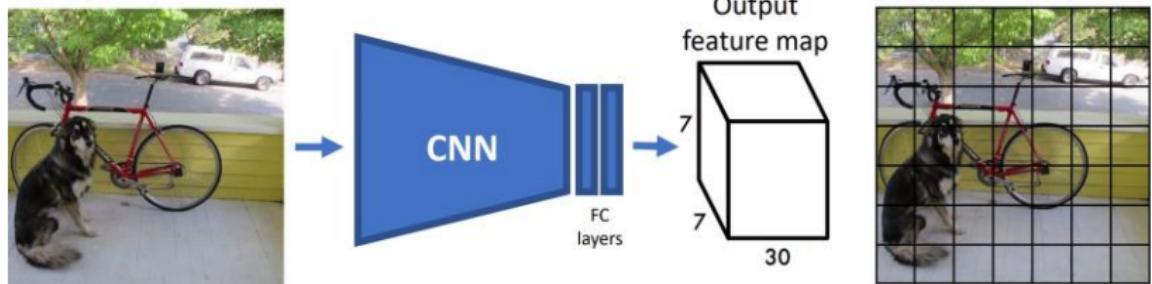
Andrew Ng

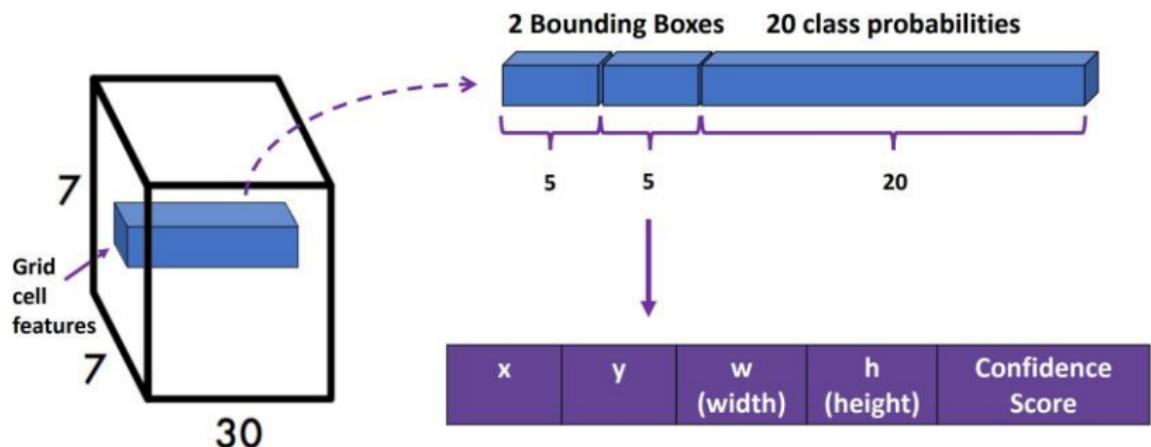
YOLO algorithm

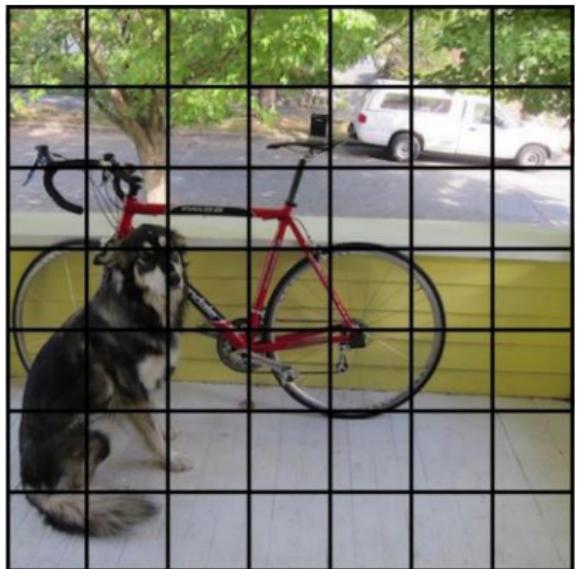


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

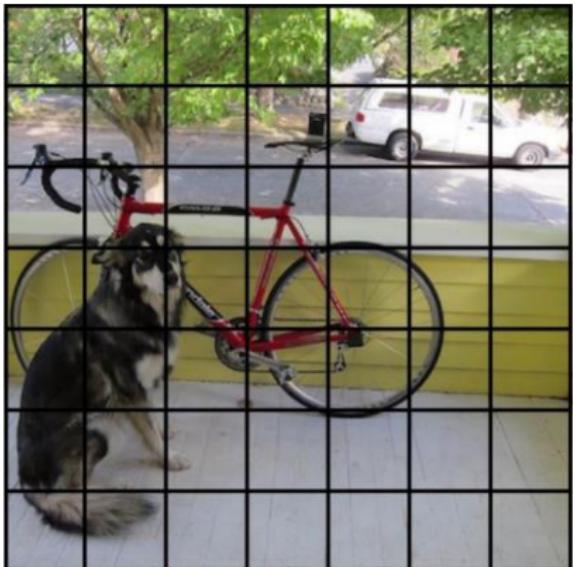
Andrew Ng





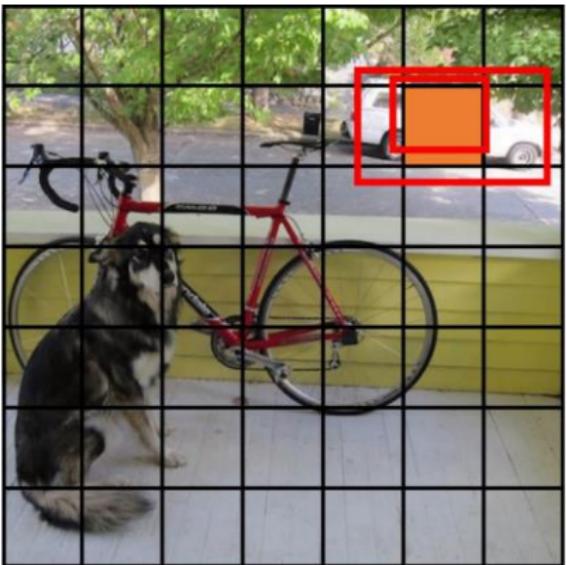


Each cell predicts



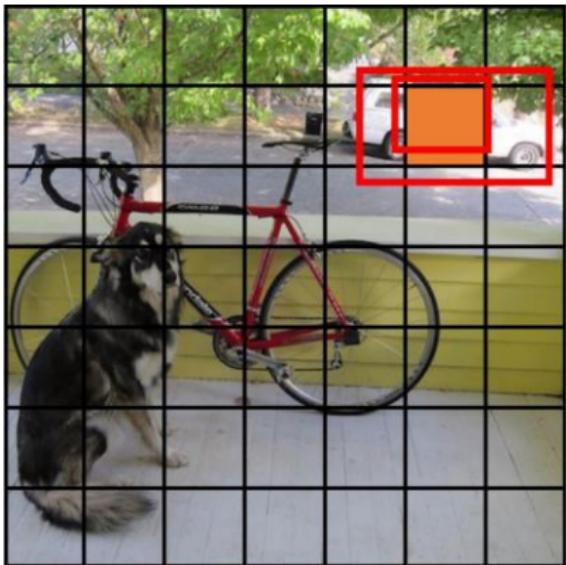
Each cell predicts

- ▶ $B = 2$ bounding boxes
 $(x, y, w, h) + \text{confidence score}$



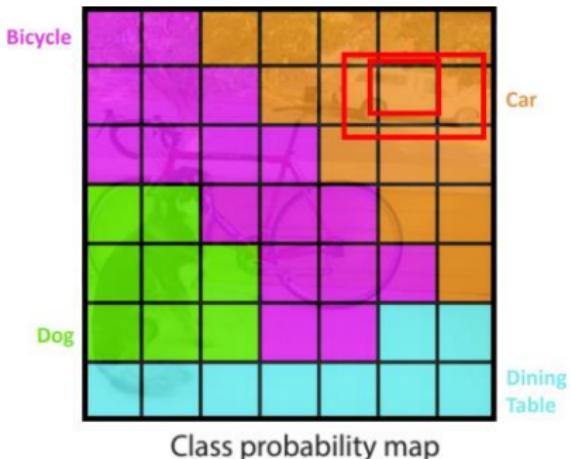
Each cell predicts

- ▶ $B = 2$ bounding boxes
 (x, y, w, h) + confidence score
- ▶ $C = 20$ class probabilities



Each cell predicts

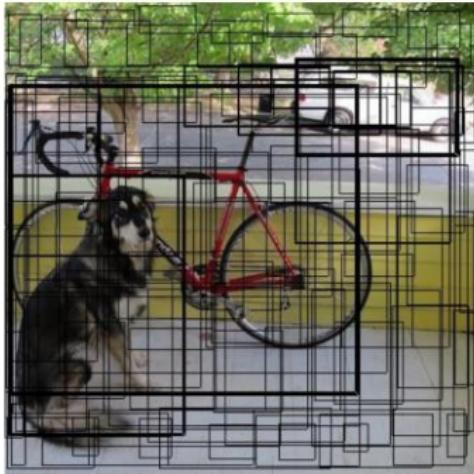
- ▶ $B = 2$ bounding boxes
- ▶ $(x, y, w, h) + \text{confidence score}$
- ▶ $C = 20$ class probabilities



Each cell predicts

- ▶ $B = 2$ bounding boxes
 $(x, y, w, h) + \text{confidence score}$
- ▶ $C = 20$ class probabilities

SxSxB Bounding-Boxes ($S=7, B=2 \rightarrow 96 \text{ Bboxes}$)

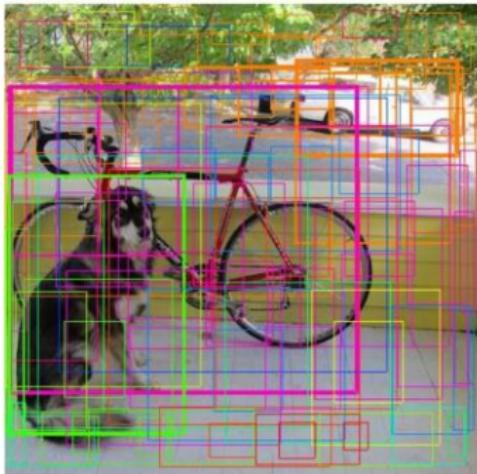


$S \times S$ grid on input

Each cell predicts

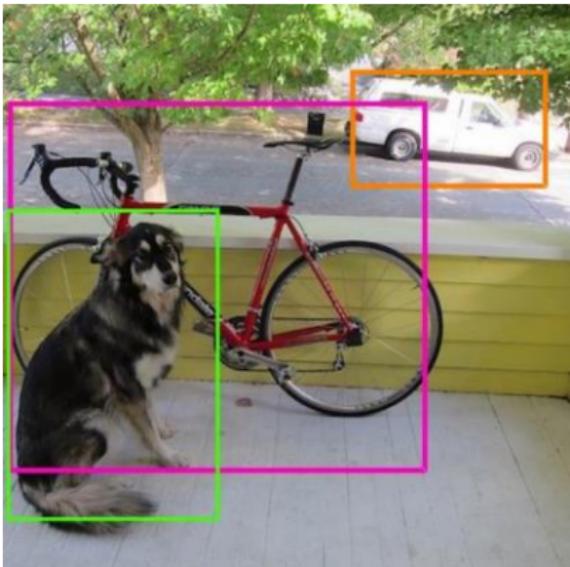
- ▶ $B = 2$ bounding boxes
 $(x, y, w, h) + \text{confidence score}$
- ▶ $C = 20$ class probabilities

SxSxB Bounding-Boxes ($S=7, B=2 \rightarrow 96$ Bboxes)



Each cell predicts

- ▶ $B = 2$ bounding boxes
 $(x, y, w, h) + \text{confidence score}$
- ▶ $C = 20$ class probabilities
- ▶ Apply Non-Maximum Suppression



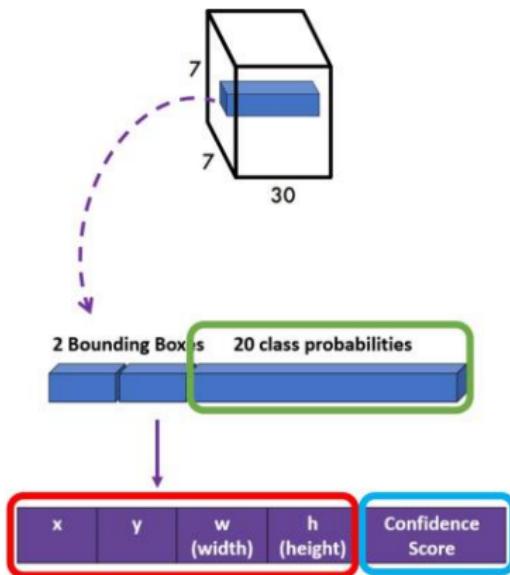
YOLO - Loss Function

YOLO – Loss function

$$\mathcal{L} = \mathcal{L}_{Localization\ Loss}$$

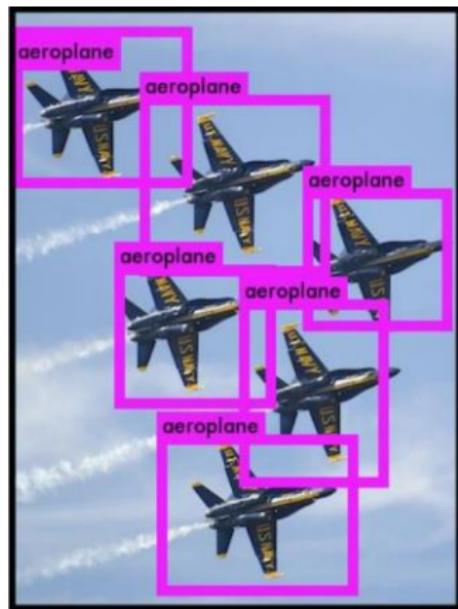
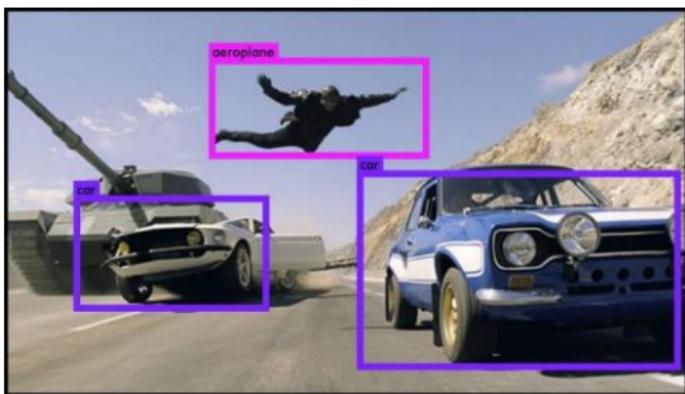
$$+ \mathcal{L}_{Confidence\ Loss}$$

$$+ \mathcal{L}_{Classification\ Loss}$$

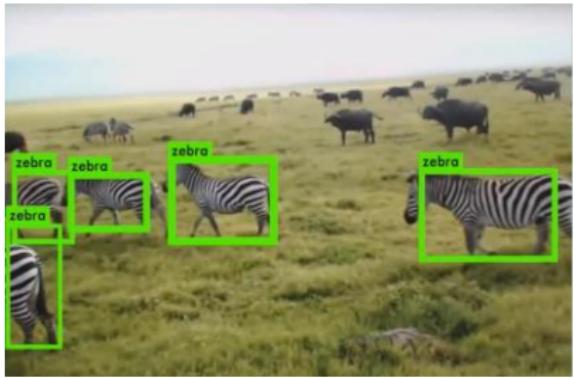


YOLO - Benefits

- ▶ Fast. Good for real-time processing
- ▶ End-to-end training

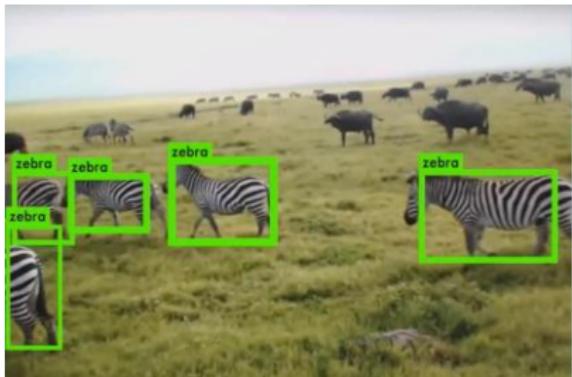


YOLO - Limitations



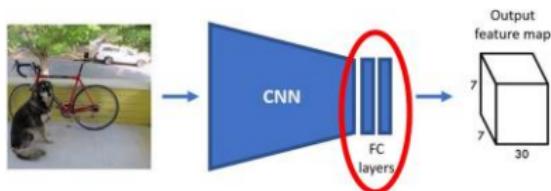
YOLO - Limitations

- ▶ Difficult to detect small objects
- ▶ Coarse predictions



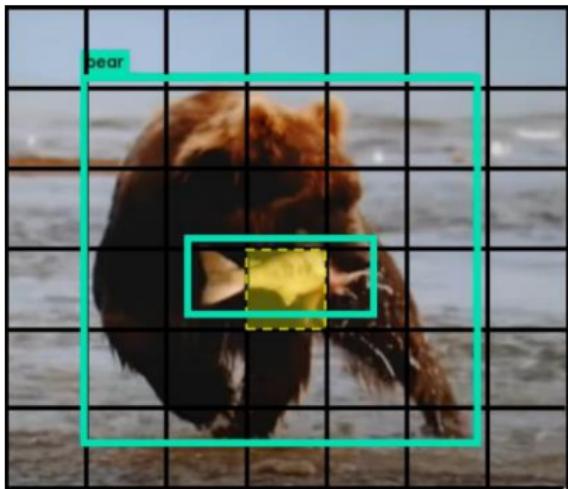
YOLO - Limitations

- ▶ Difficult to detect small objects
- ▶ Coarse predictions
- ▶ Fixed input size



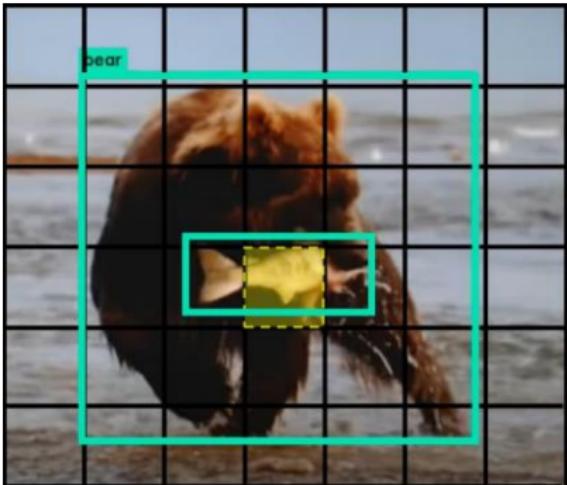
YOLO - Limitations

- ▶ Difficult to detect small objects
 - ▶ Coarse predictions
 - ▶ Fixed input size
 - ▶ A grid cell can predict only one class

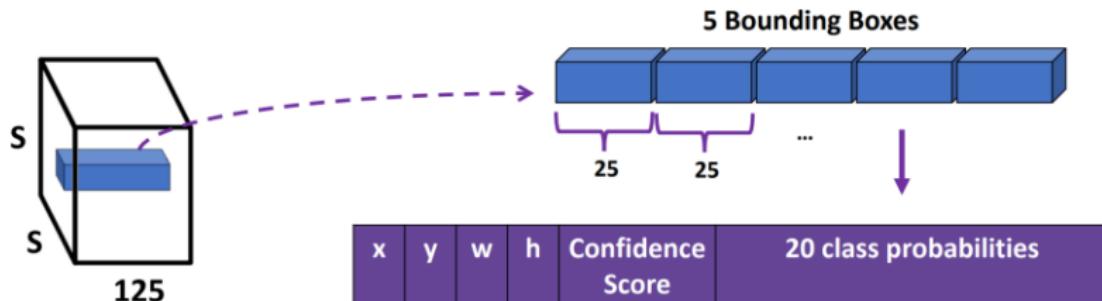


YOLO - Limitations

- ▶ Difficult to detect small objects
 - ▶ Coarse predictions
 - ▶ Fixed input size
 - ▶ A grid cell can predict only one class
-
- ▶ **Solutions:**
 - Remove fc layers!
 - Predict class per bbox (not per cell)



- ▶ Removed fully connected layers
- ▶ A grid cell predicts class probabilities for each box



► YOLOv3

- J. Redmon, A. Farhadi. Yolov3: An incremental improvement, 2018

► YOLOv4

- A. Bochkovskiy, C. Wang, H. Liao. Yolov4: Optimal speed and accuracy of object detection (Feb. 2020)

► YOLOv5

- YOLOv5 by ultralytics (June 2020)

► PP-YOLO

- X. Long, K. Deng, G. Wang, Y. Zhang, Q. Dang, Y. Gao, H. Shen, J. Ren, S. Han, E. Ding, S. Wen. Pp-yolo: An effective and efficient implementation of object detector (June 2020)

► PP-YOLOv2 (2021)

- J. X. Huang, X. Wang, W. Lv, X. Bai, X. Long, K. Deng, Q. Dang, S. Han, Q. Liu, X. Hu, D. Yu, Y. Ma, O. Yoshie. PP-YOLOv2: A Practical Object Detector (2021)

These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV [Deep Learning for Computer Vision: Fundamentals and Applications](#)
- ▶ Justin Johnson, UMich EECS 498.008/598.008: [Deep Learning for Computer Vision](#)
- ▶ [cs231n@stanford](#)
- ▶ [EECS 498.008 / 598.008](#)
- ▶ [CNN DeepLearning.AI](#)