# Trustworthy AI Systems

## -- Image Segmentation

Instructor: Guangjing Wang

guangjingwang@usf.edu

# Last Lecture

- Image classification

- Convolutional neural network

- Some practices for project

# Homework 1: Paper Review

- Paper review is a basic task for a researcher
  - Paper Summary
  - Strengths
  - Weaknesses
  - Questions
  - <span style="color:red">Future Opportunities</span>

When you read a paper, thinking:
  - What is the research problem and motivation?
  - What are the challenges and technical contributions?
  - How is the experimental evaluation?
  - How is the related work, and overall presentation?

# Computer Vision Tasks



Classification — CAT — No spatial extent

Semantic Segmentation — GRASS, CAT, TREE, SKY — No objects, just pixels
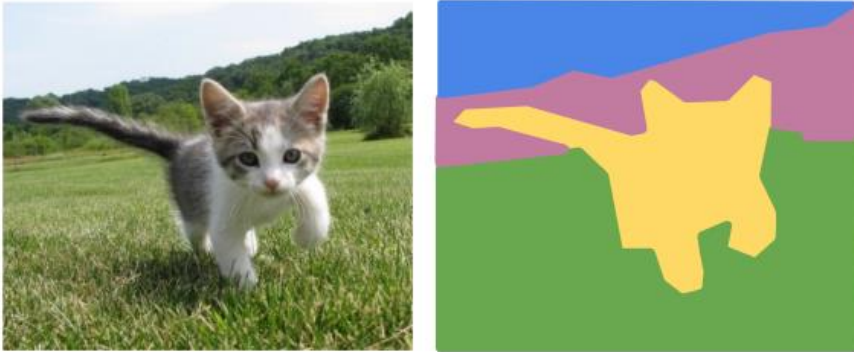
Object Detection — DOG, DOG, CAT

Instance Segmentation — DOG, DOG, CAT

Multiple Object
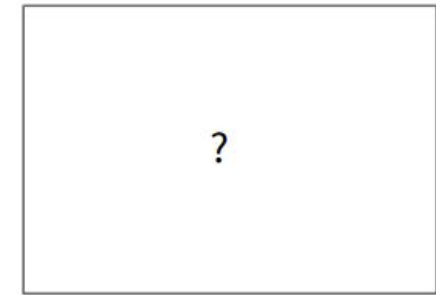
This image is CC0 public domain

# Semantic Segmentation: Problem
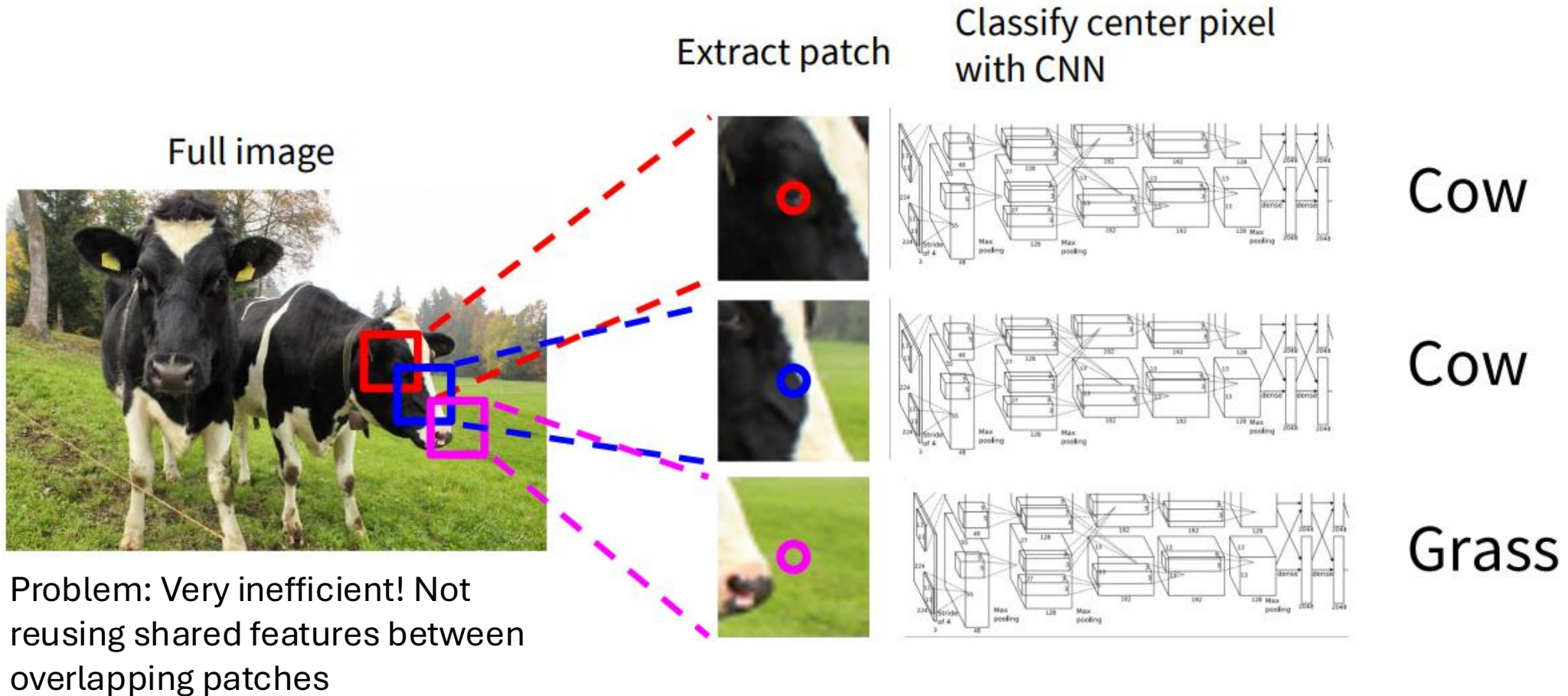


GRASS, CAT, TREE, SKY, ...

Paired training data: for each training image, each pixel is labeled with a semantic category.



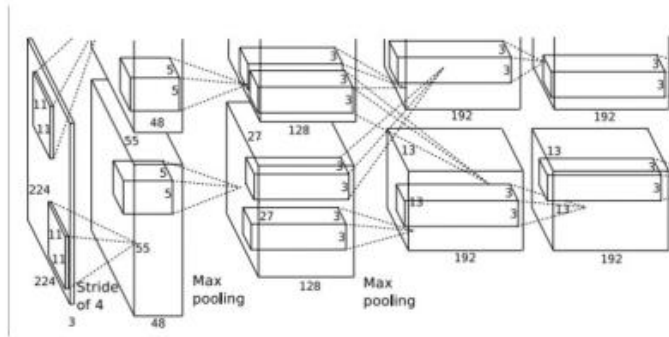At test time, classify each pixel of a new image.

Label each pixel in the image with a category label.

# Semantic Segmentation: Sliding Window

Full image

Extract patch

Classify center pixel with CNN

Cow

Cow

Grass

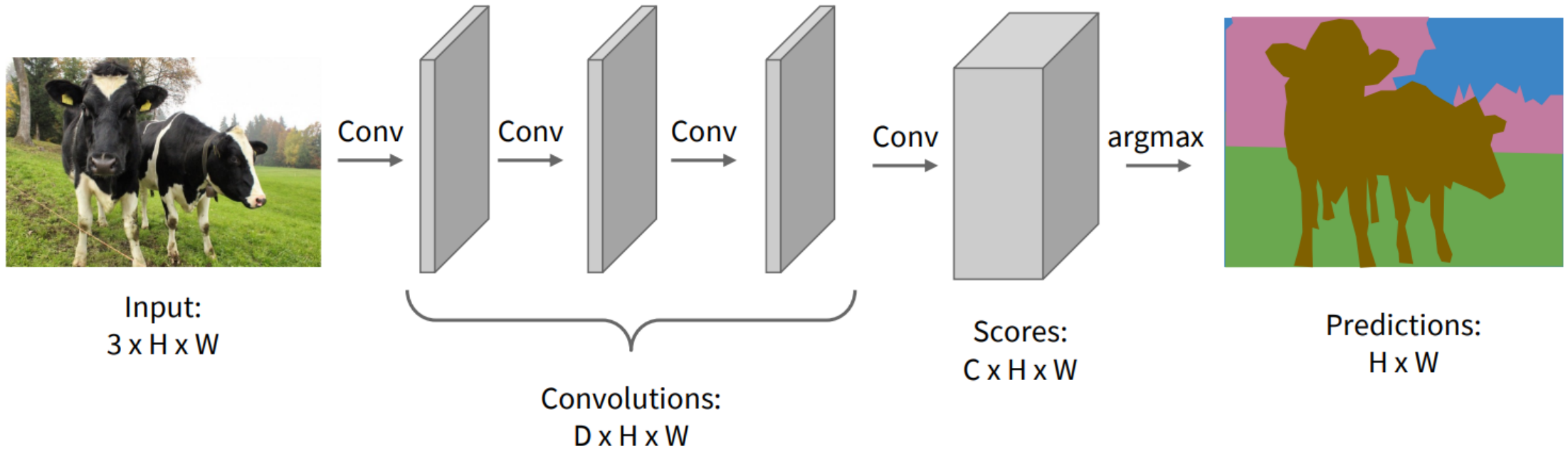Problem: Very inefficient! Not reusing shared features between overlapping patches

# Semantic Segmentation: Convolution (1)



Encode the entire image with conv net, and do semantic segmentation on top

# Semantic Segmentation: Convolution (2)



Input:
3 x H x W

Conv → Conv → Conv → Conv → argmax

Convolutions:
D x H x W

Scores:
C x H x W

Predictions:
H x W

Potential problem?
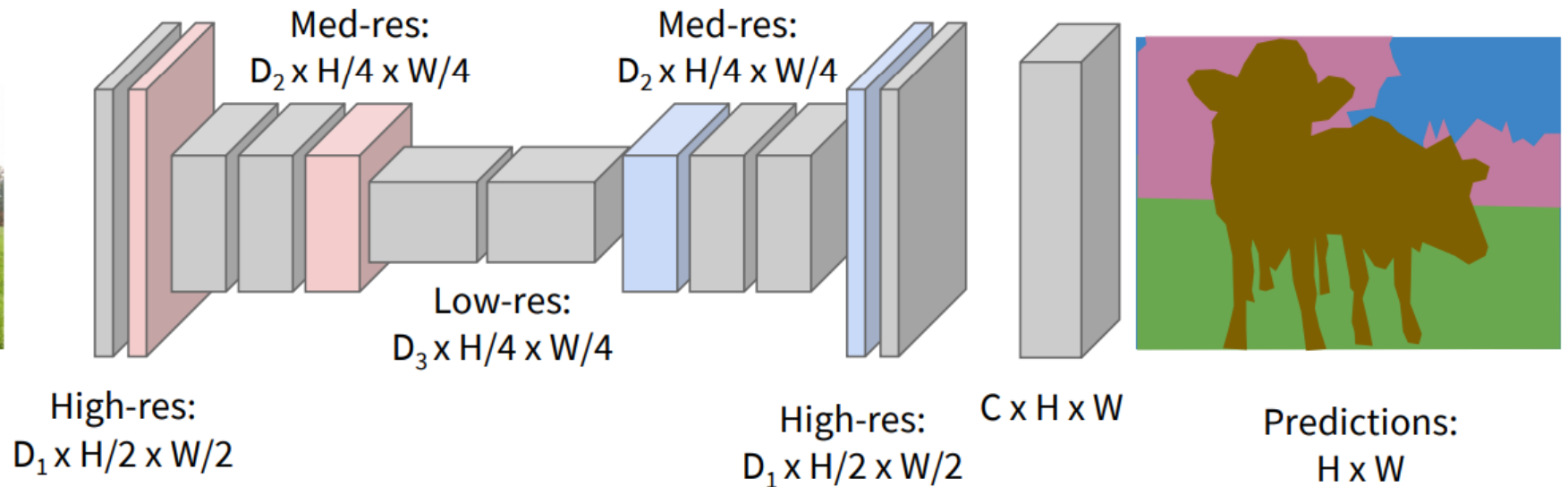
# Semantic Segmentation: Convolution (3)



Downsampling:
Pooling, strided convolution

Design network as a bunch of convolutional layers, with downsampling and upsampling inside the network!

Upsampling:
???

Med-res:
$D_2 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

High-res:
$D_1 \times H/2 \times W/2$

$C \times H \times W$

Predictions:
$H \times W$

# Upsampling

- Non-learnable upsampling
  - Fill the same
  - Fill zeros
  - Remember location then fill
  - You design it…

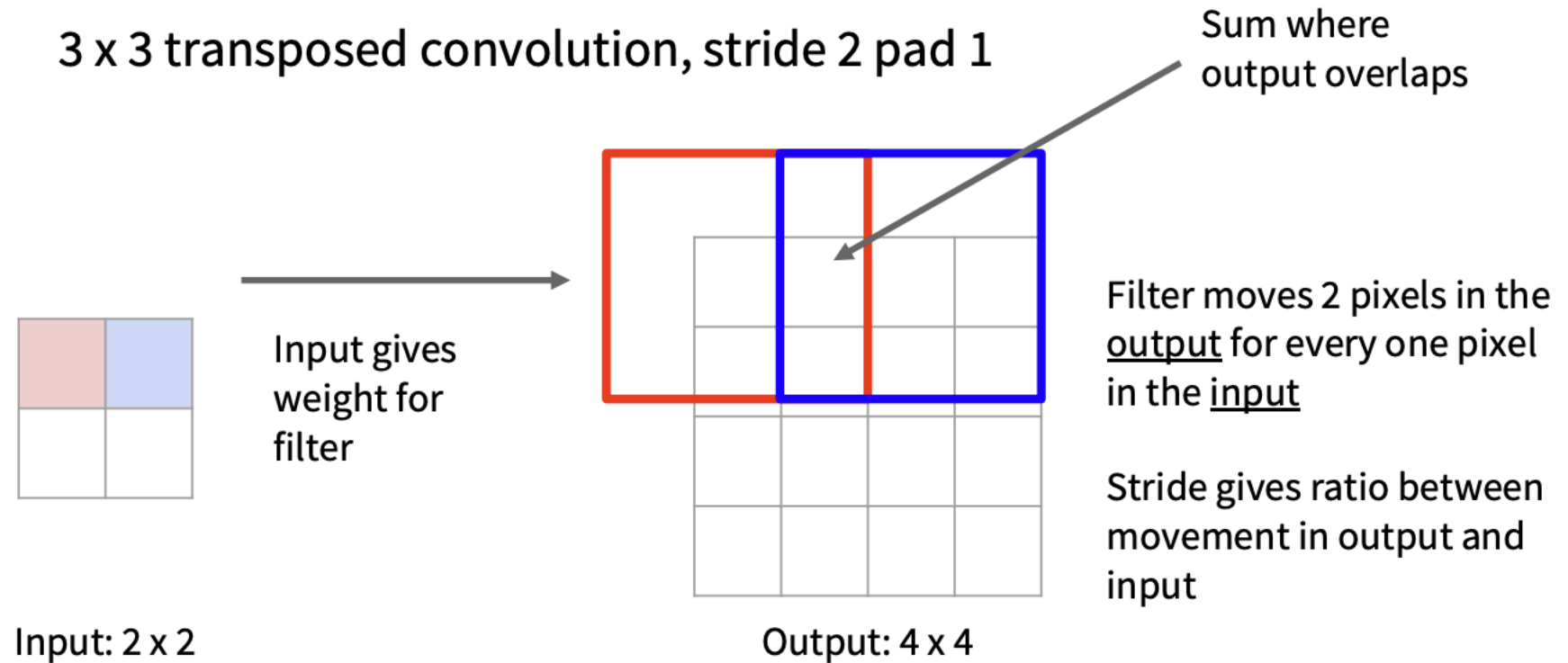- Learnable upsampling
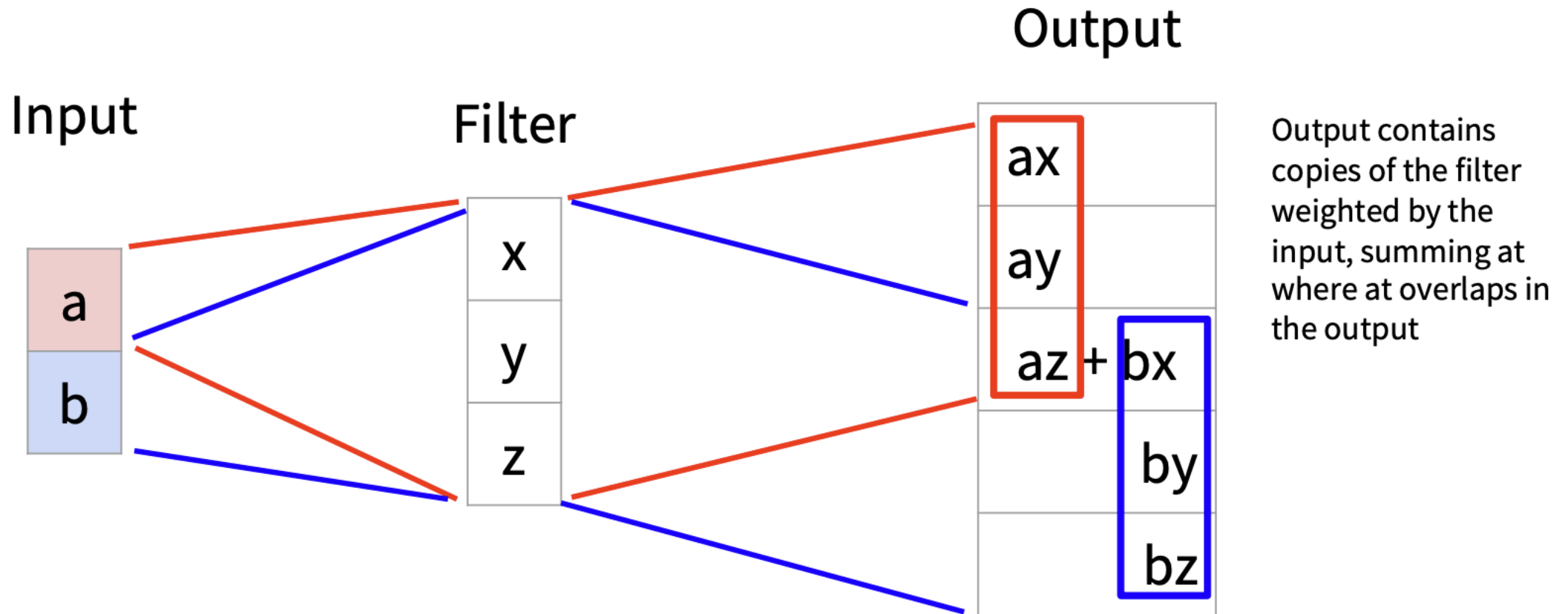  - Transposed convolution



"Bed of Nails"

Input: 2 x 2

| 1 | 2 |
|---|---|
| 3 | 4 |

Output: 4 x 4

| 1 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 3 | 0 | 4 | 0 |
| 0 | 0 | 0 | 0 |

# Upsampling: Transposed Convolution



3 x 3 transposed convolution, stride 2 pad 1

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: 1D Example



Input

| a |
| b |

Filter

| x |
| y |
| z |

Output

| ax |
| ay |
| az + bx |
| by |
| bz |

Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

# Convolution as Matrix Multiplication

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T\vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

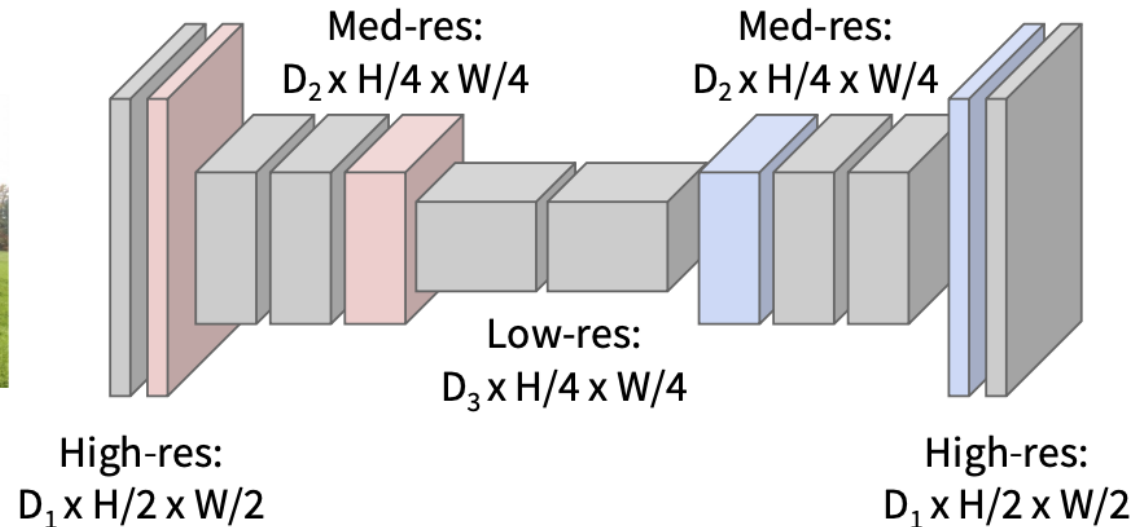Example: 1D transposed conv, kernel size=3, stride=2, padding=0

# Semantic Segmentation: Fully Convolutional



Downsampling:
Pooling, strided convolution

Design network as a bunch of convolutional layers, with downsampling and upsampling inside the network!

Upsampling:
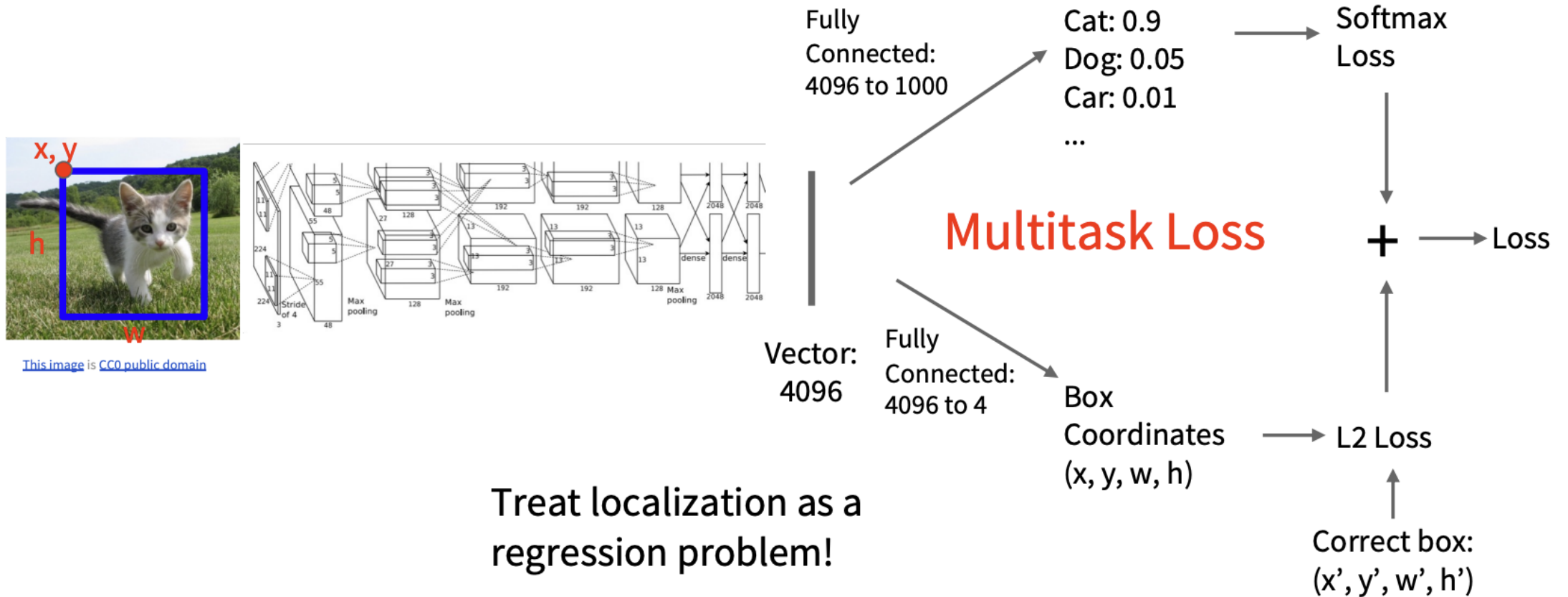Unpooling or strided transposed convolution

Input:
$3 \times H \times W$

High-res:
$D_1 \times H/2 \times W/2$

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

High-res:
$D_1 \times H/2 \times W/2$

Predictions:
$H \times W$

# Take a break



https://www.youtube.com/watch?v=JIPbilHxFbI

# Object Detection: Classification + Localization



Fully Connected: 4096 to 1000

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

Softmax Loss

**Multitask Loss**

+ → Loss

Vector: 4096

Fully Connected: 4096 to 4

Box Coordinates (x, y, w, h)

L2 Loss

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

Treat localization as a regression problem!

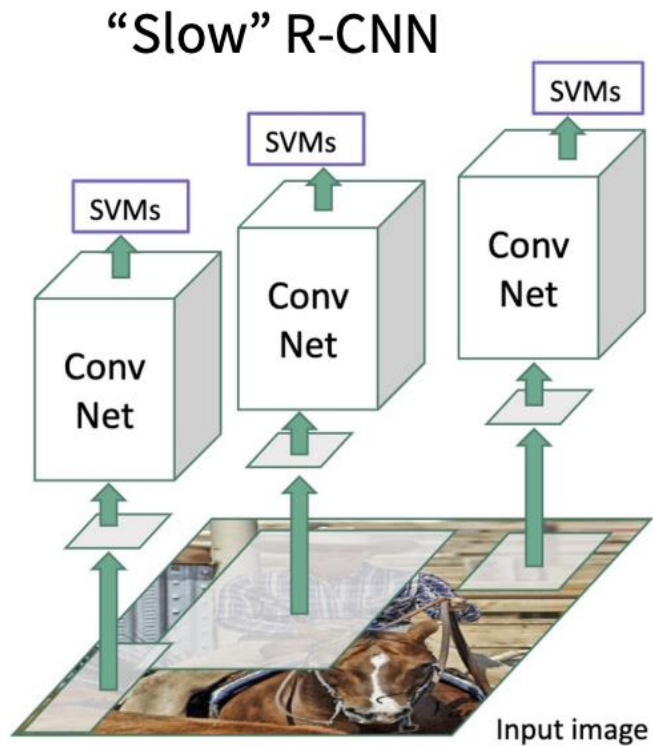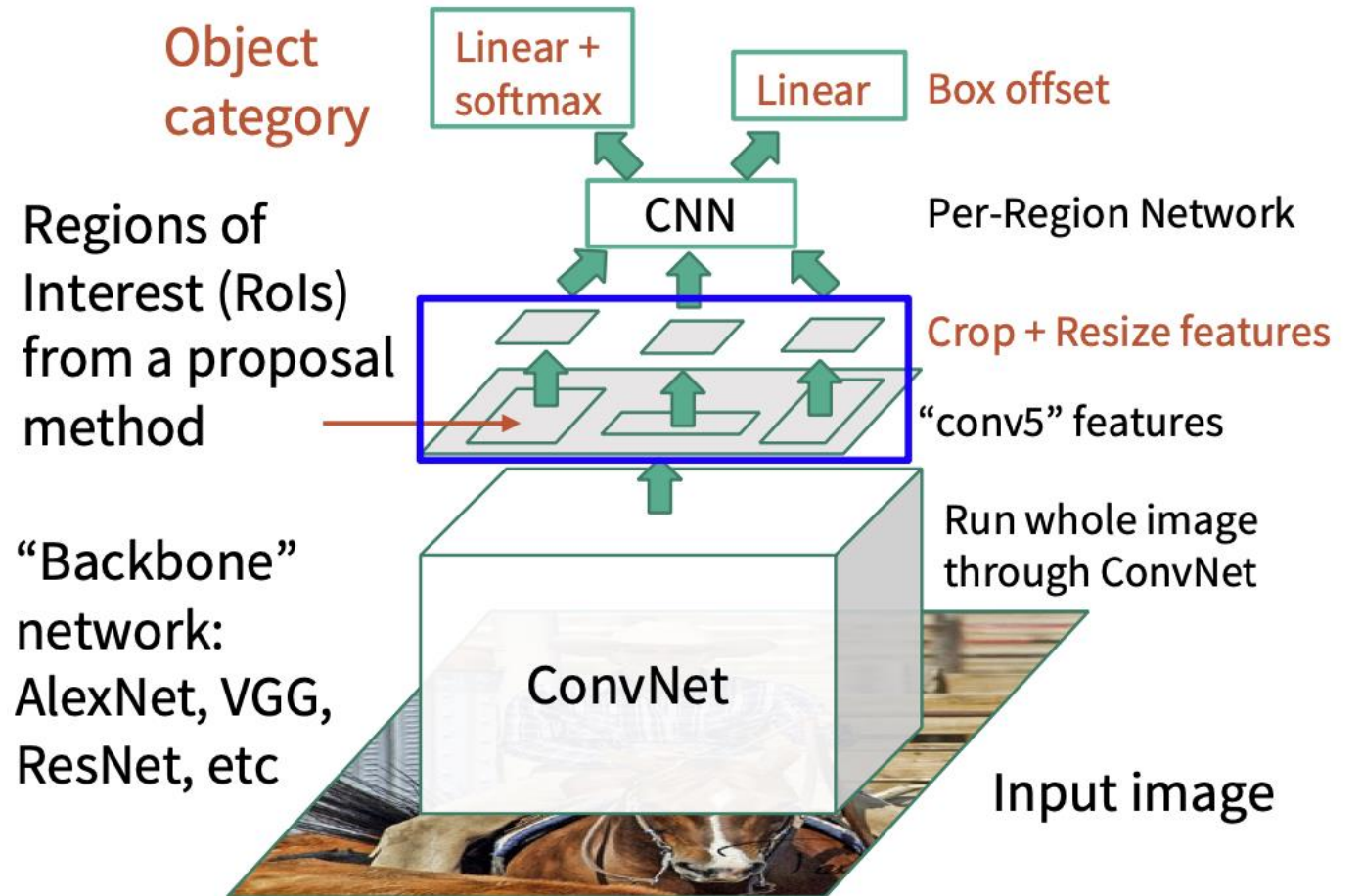This image is CC0 public domain

# Object Detection

- What if there are multiple objects?
  - Apply a CNN to many different crops of the image, CNN classifies each crop as object or background
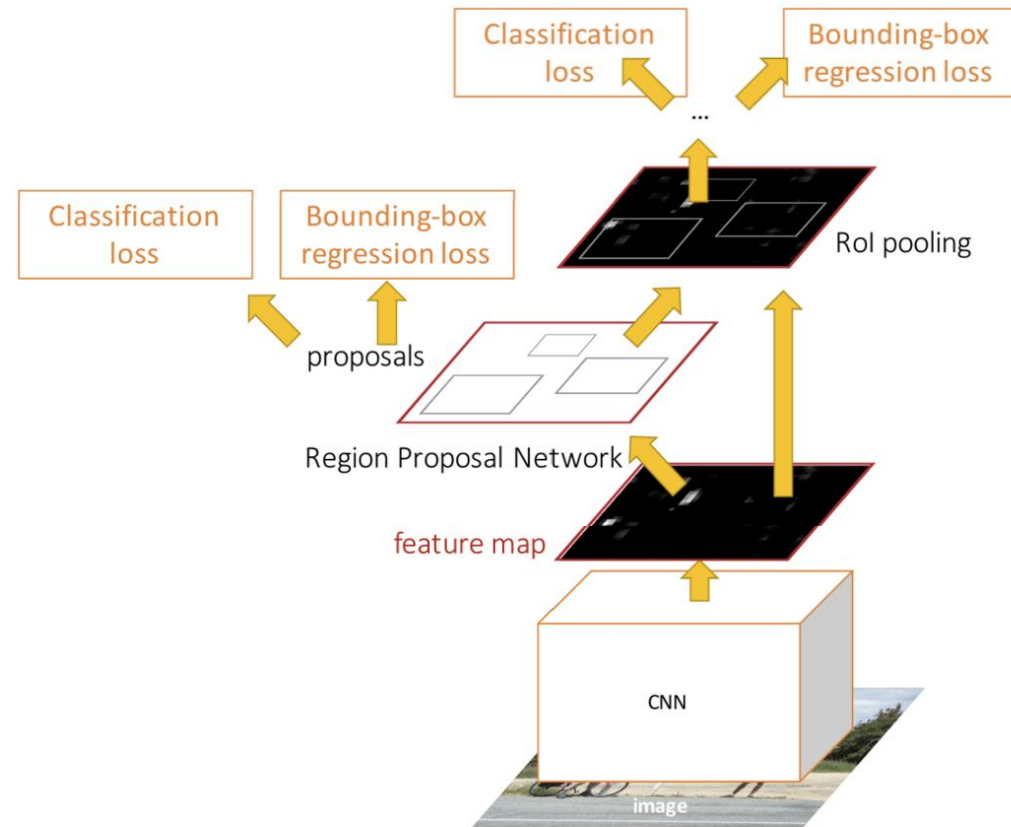
# R-CNN and Fast R-CNN



"Slow" R-CNN

SVMs

SVMs

SVMs

ConvNet

ConvNet

ConvNet

Input image

extracts around 2000
bottom-up region proposals,

Object category

Linear + softmax

Linear    Box offset

CNN    Per-Region Network

Regions of Interest (RoIs) from a proposal method

Crop + Resize features

"conv5" features

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

ConvNet

Run whole image through ConvNet

Input image

# Faster R-CNN: Make CNN Do Proposals
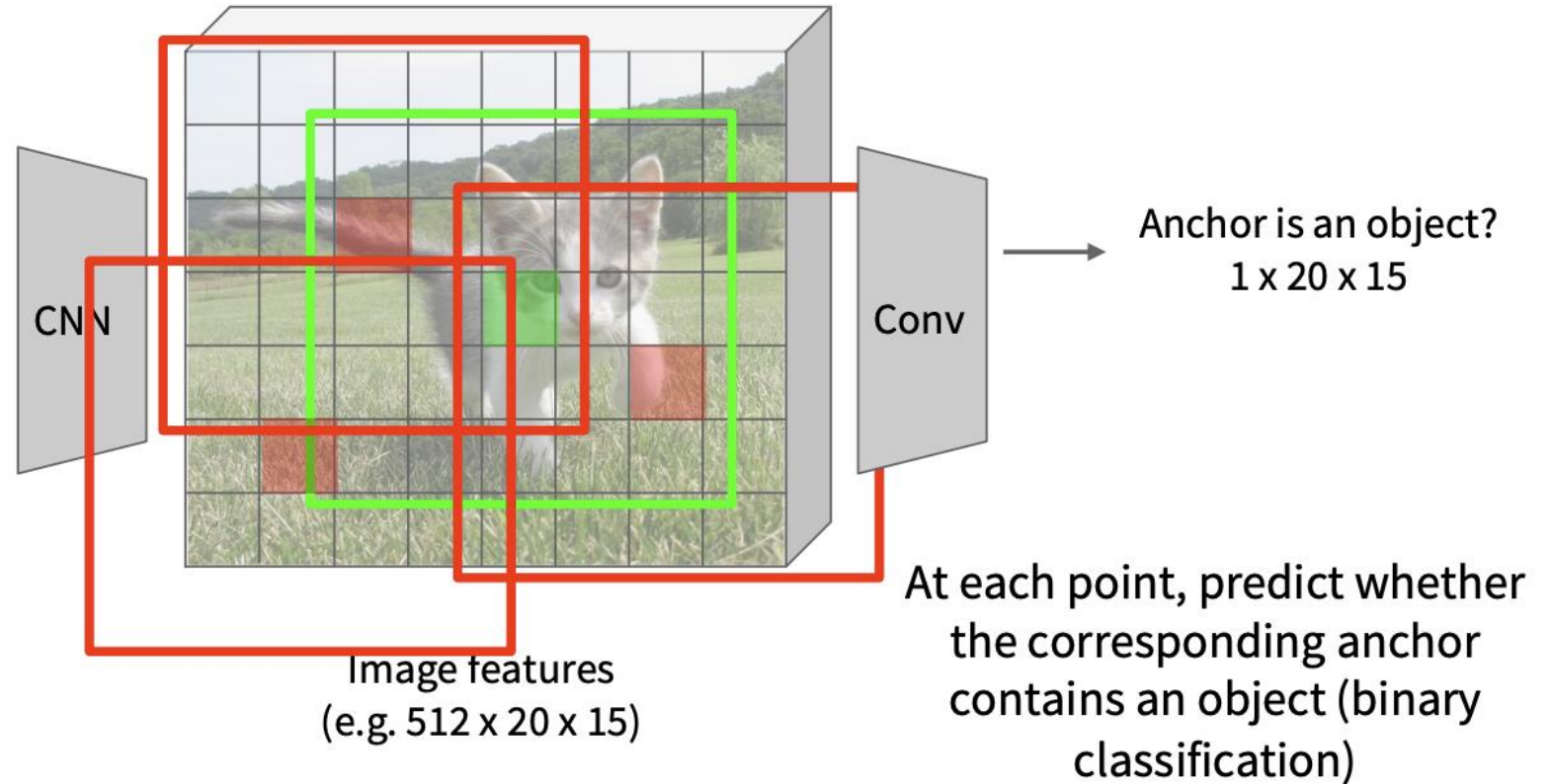
- Insert Region Proposal Network (RPN) to predict proposals from features

# Region Proposal Network (1)



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

CNN

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

Conv

Anchor is an object?
1 x 20 x 15

At each point, predict whether the corresponding anchor contains an object (binary classification)

# Region Proposal Network (2)



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

CNN

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

Conv

Anchor is an object?
1 x 20 x 15

Box corrections
4 x 20 x 15

For positive boxes, also predict a corrections from the anchor to the ground-truth box (regress 4 numbers per pixel)

# Faster R-CNN: Two Stages
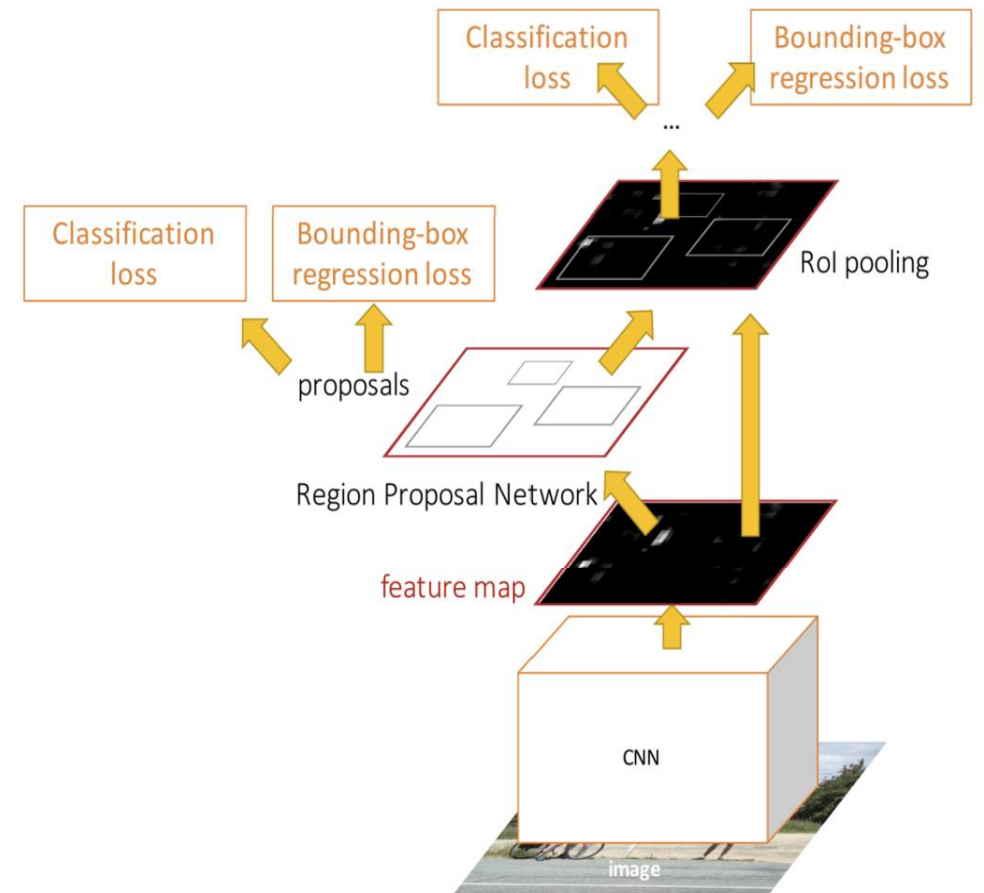
Jointly train with 4 losses:
- RPN classify object / not object
- RPN regress box coordinates
- Final classification score (object classes)
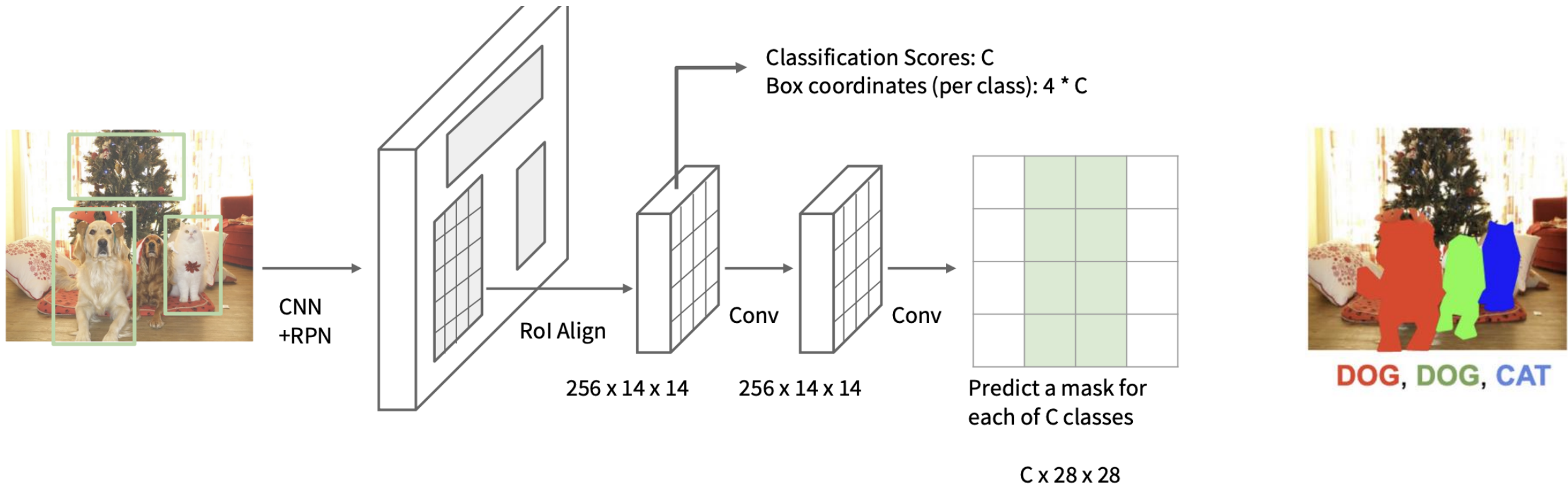- Final box coordinates

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
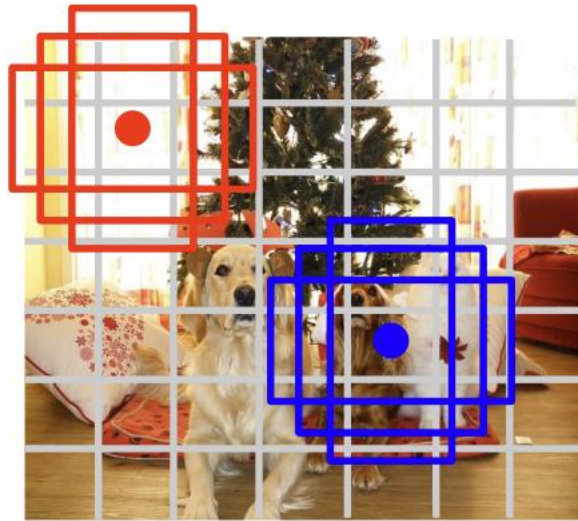- Prediction bbox offset

**CIS6930 Trustworthy AI Systems**

# Instance Segmentation: Mask R-CNN



CNN +RPN

RoI Align

256 x 14 x 14

Conv

256 x 14 x 14

Conv

Classification Scores: C
Box coordinates (per class): 4 * C

Predict a mask for each of C classes

C x 28 x 28

DOG, DOG, CAT

# Yolo: Single Stage Object Detector



Input image
3 x H x W

Divide image into grid
7 x 7
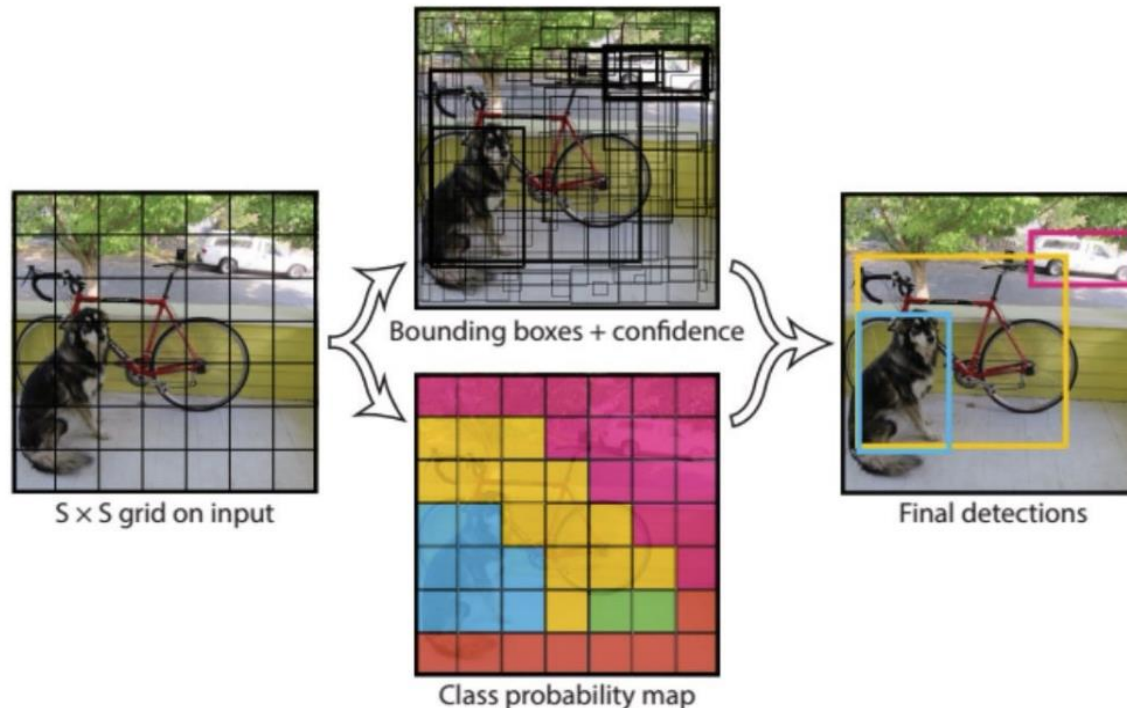
Image a set of base
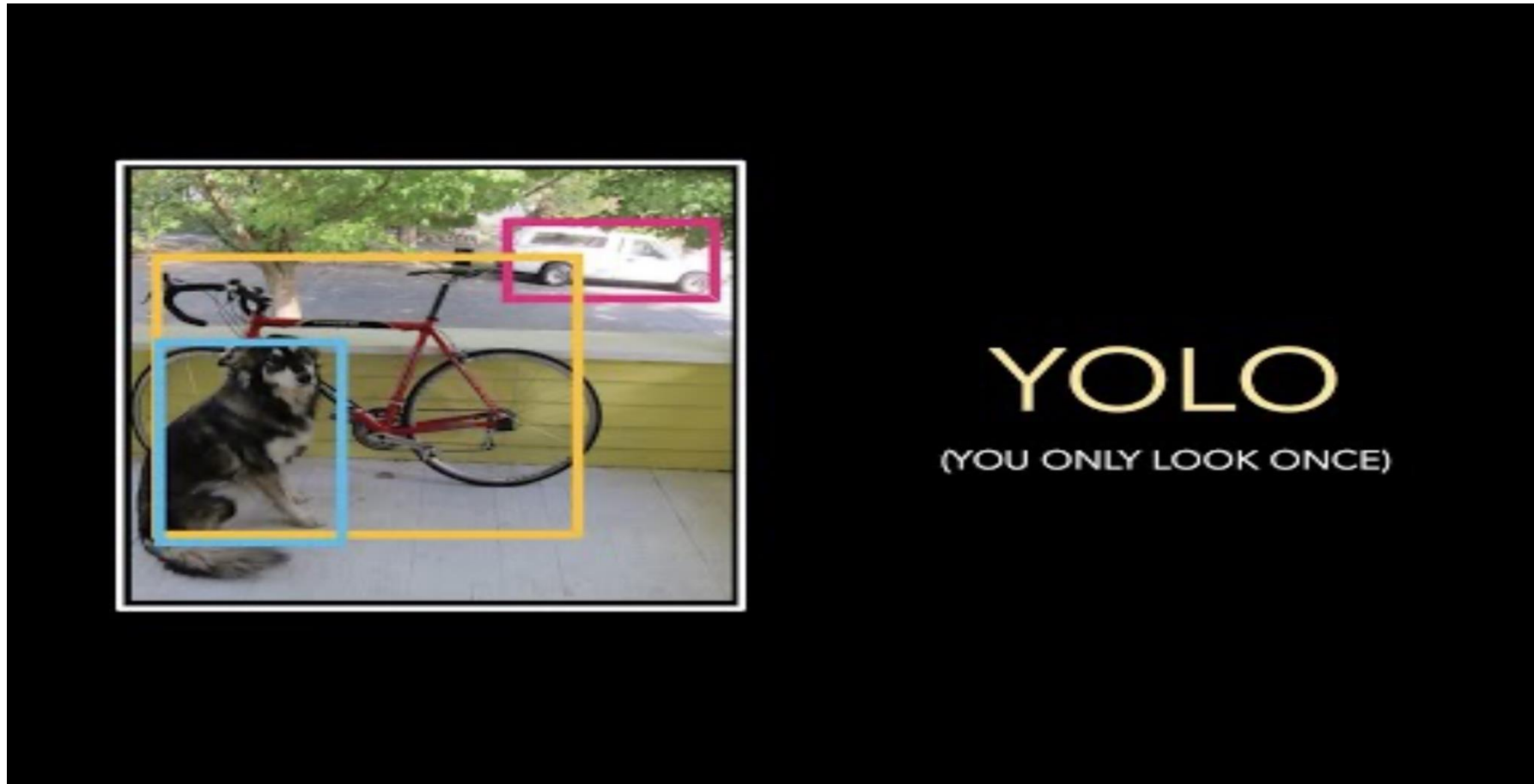boxes centered at each
grid cell Here B = 3

Within each grid cell:
- Regress from each of the B base boxes to a final box with 5 numbers:
  (dx, dy, dh, dw, confidence)

- Predict scores for each of C classes (including background as a class)

- Looks a lot like RPN, but category-specific!

- Output: 7x7x(5*B+C)

# Yolo: Non-Max Suppression

- If IoU(P1, P2) > Threshold: P = argmax(C(p1), C(p2))
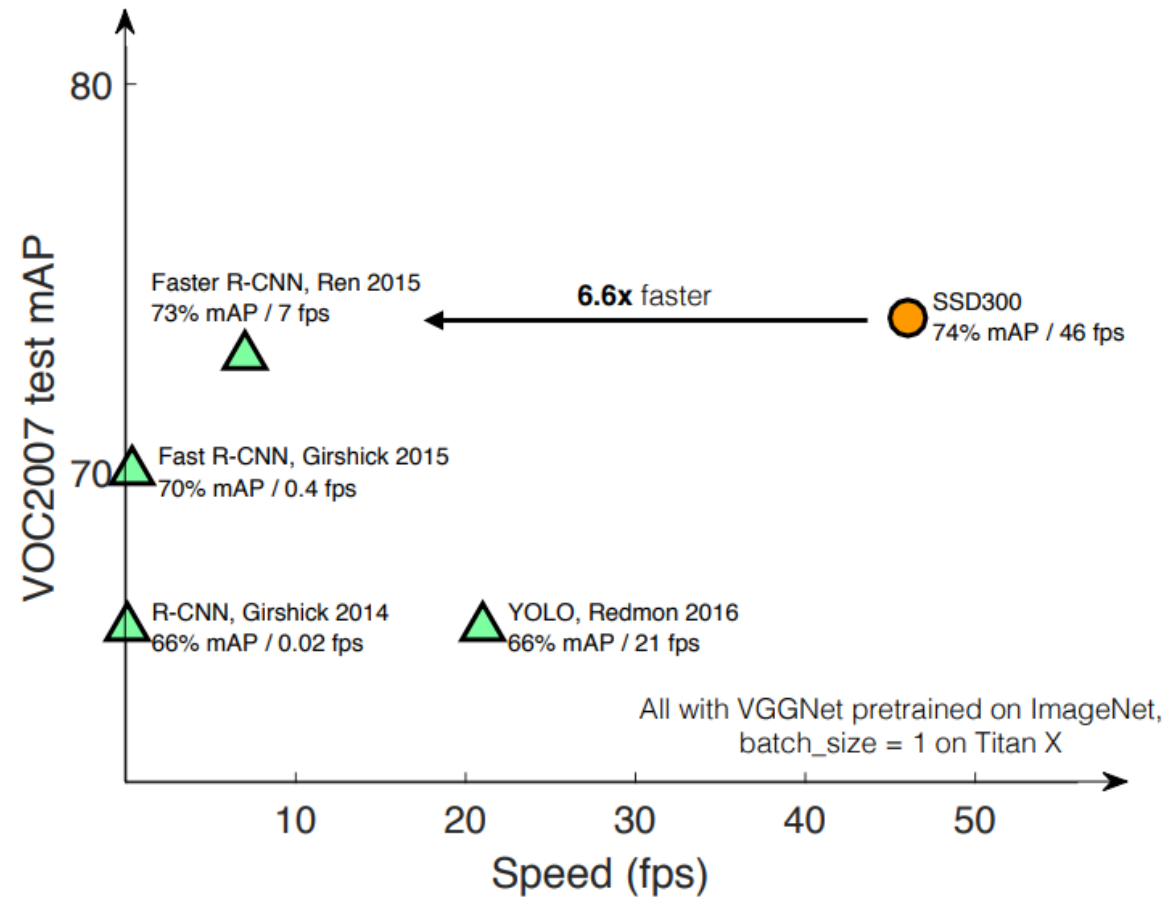  - Eliminating bounding boxes that have a high overlap with the box that has the highest confidence score



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

# YOLO: Model as a Regression Problem



https://youtu.be/svn9-xV7wjk

# Single-shot VS Two-shot Detector



https://www.cs.unc.edu/~wliu/papers/ssd_eccv2016_slide.pdf

# Object Detection: Evaluation Metrics

- Intersection over Union (IoU)
  - Predicted bounding box (A) and ground truth bounding box (B)

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- Average Precision (AP)
  - The precision-recall curve that is created by varying the detection threshold.
  - mean Average Precision (mAP), which calculates AP for each class and then take the average

# Midterm Project Group

- Please find your team member (1-3 members in a group)
- Sign your group in Canvas
- Random sign-up will be executed on Sep. 5$^{th}$.

# References

- https://cs231n.stanford.edu/slides/2024/lecture_9.pdf

- https://encord.com/blog/yolo-object-detection-guide/

- https://github.com/ultralytics/ultralytics

- https://github.com/facebookresearch/detectron2