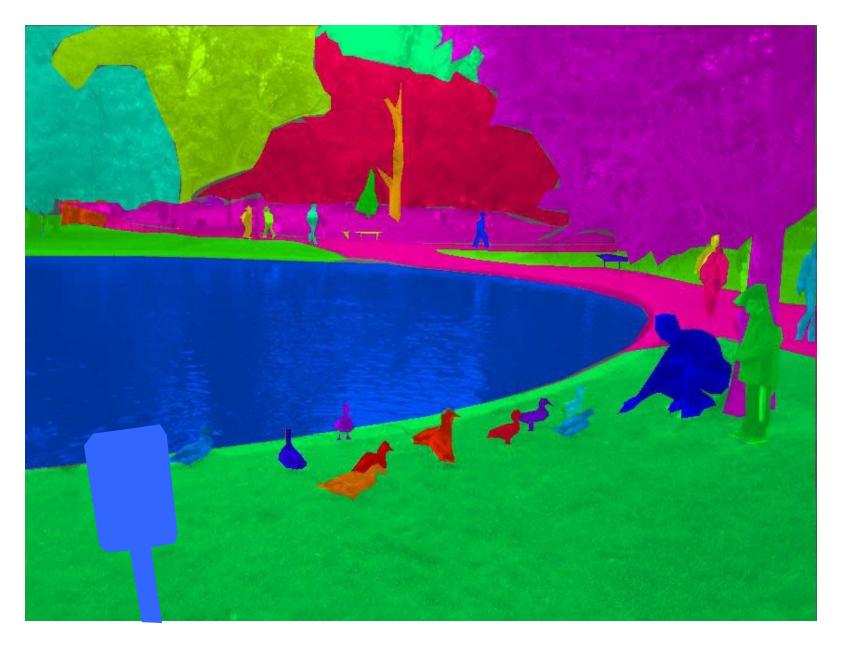
Lecture 11

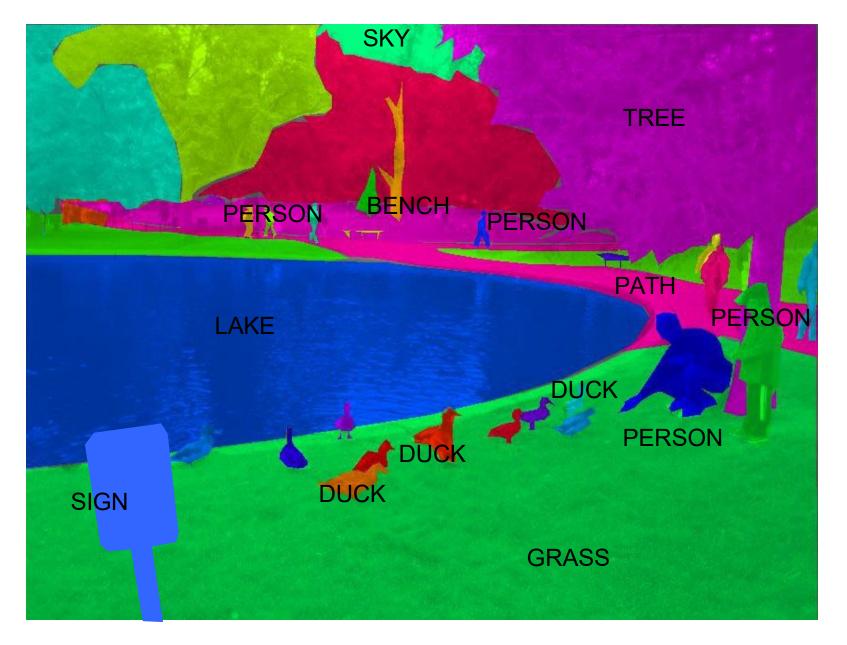
Visual Recognition



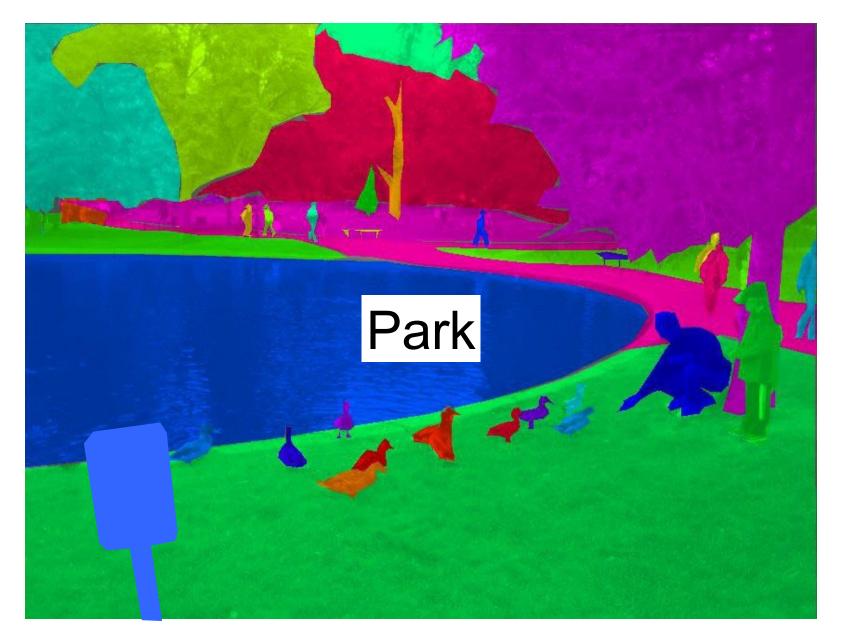




Label each pixel as a category. Each category has a unique color.



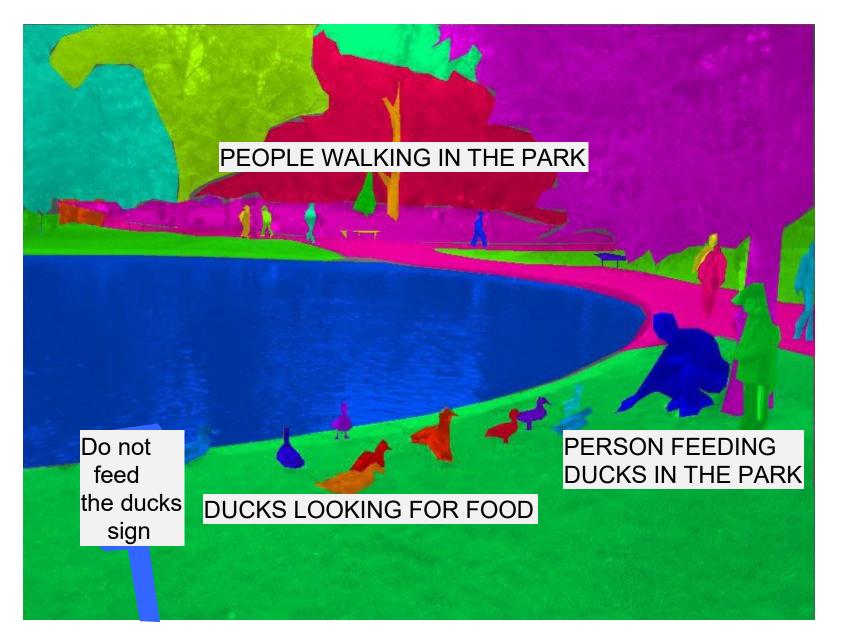
Label each pixel as a category. Each category has a unique color.



Scene-Level Classification: This is a "PARK"



Image Captioning: Describe the image in human language (e.g. English)



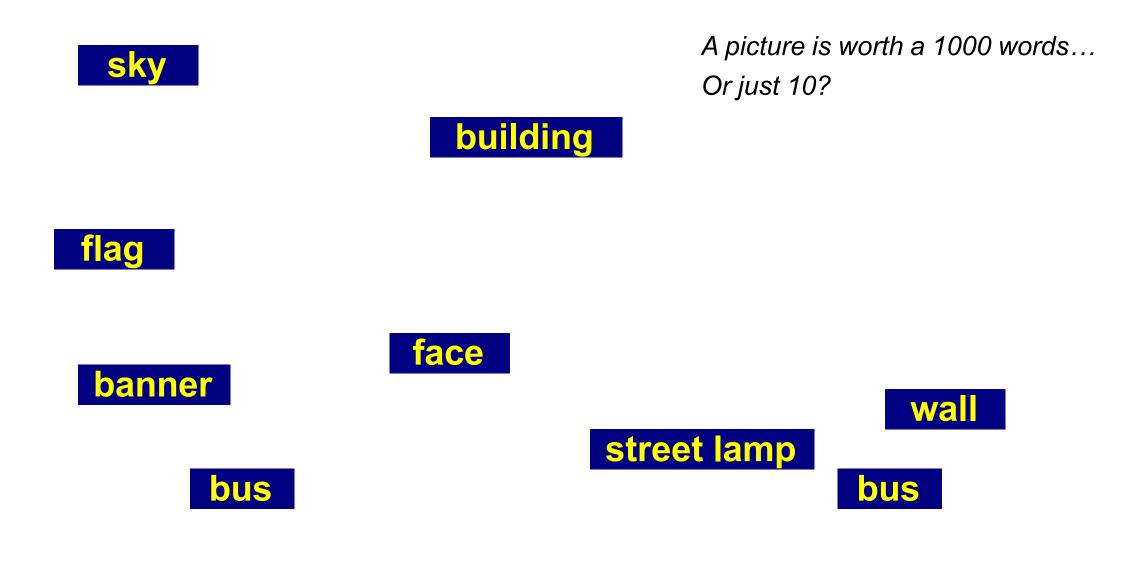
Dense Image Captioning: Describe several parts of the image

Why do we care about recognition?



- The concept of "categories" encapsulates semantic information that humans use when communicating with each other.
- Categories are also linked with what can we do with those objects.

Object categories aren't everything



cars

Object categories aren't everything



How finegrained should categories be?



- Bus
- Yellow & Green Bus
- •
- •
- "Beijing City Transit Bus #17, serial number 43253"

Need more general (useful) information



What can we say the very first time we see this thing?

Functional:

- A large vehicle that may be moving fast
- Will hurt you if you stand in its way.
- However, at specified places, it will allow you to enter it and transport you quickly over large distances.

Communicational:

● bus, autobus, λεωφορείο, ônibus, автобус, 公共汽车, etc.

Source: A. Efros

What makes this challenging?

Visual challenges with categories

Chair

A lot of categories are functional







 Categories are 3D, but images are 2D





World is highly varied





train

Source: A. Efros

Limits to direct perception





Today

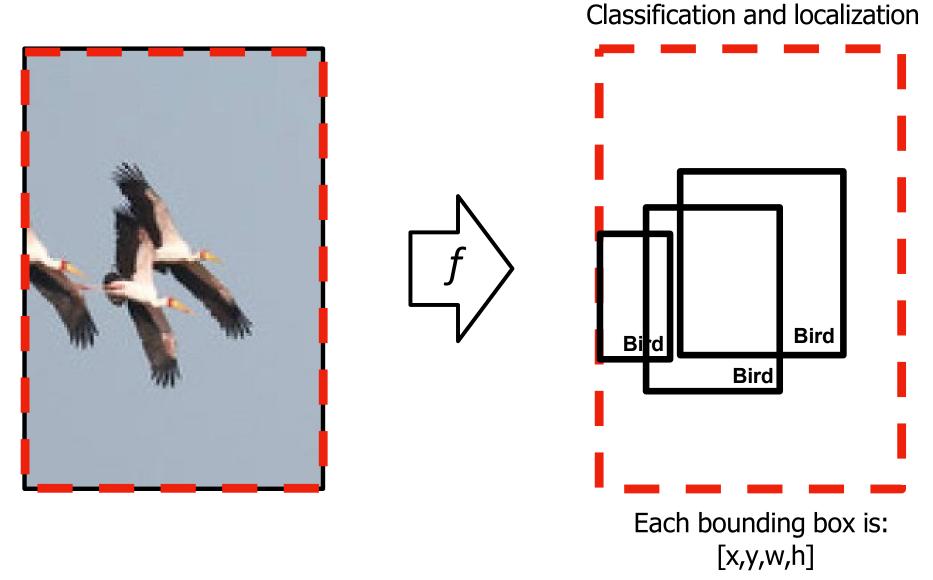
- Introduction to scene understanding
- Object detection models
- Evaluating object detectors
- Future challenges

Image Classification



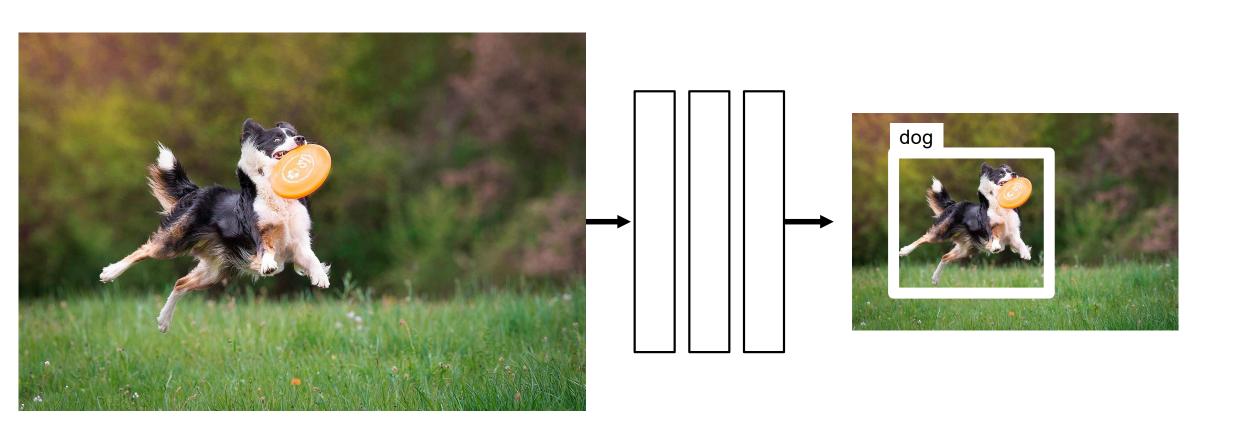


Object detection

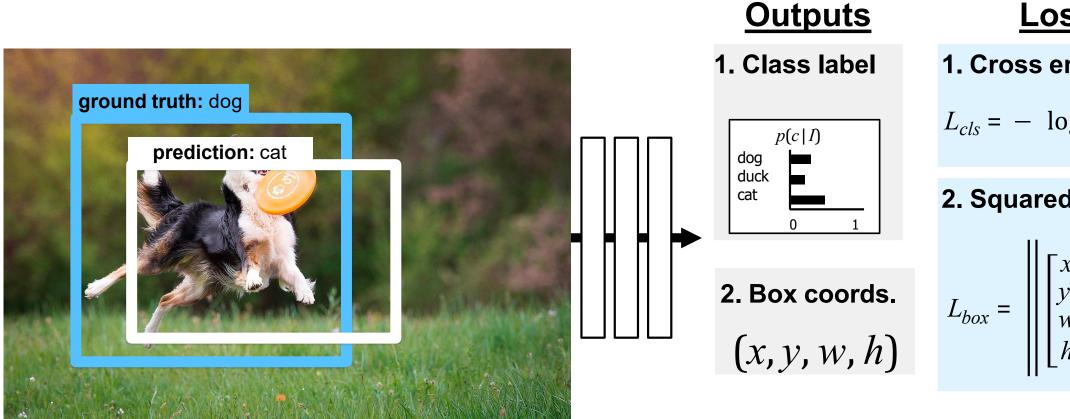


Challenge: unbounded number of detections, possibly multiple detections per pixel

Idea #1: regress bounding box



Idea #1: regress bounding box



Losses

1. Cross entropy loss

$$L_{cls} = -\log(p(y = \log))$$

2. Squared distance

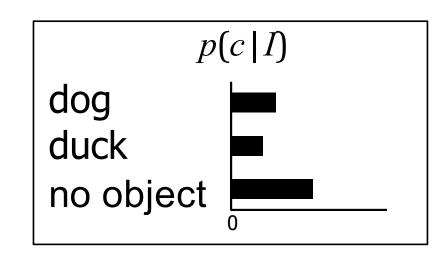
$$L_{box} = \left\| \begin{bmatrix} x \\ y \\ w \\ h \end{bmatrix} - \begin{bmatrix} x_{gt} \\ y_{gt} \\ w_{gt} \\ h_{gt} \end{bmatrix} \right\|^{2}$$

Doesn't scale well to multiple objects.

Idea #2: sliding window



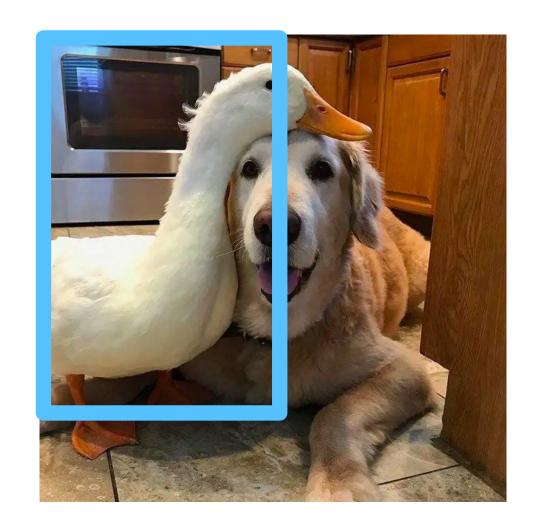




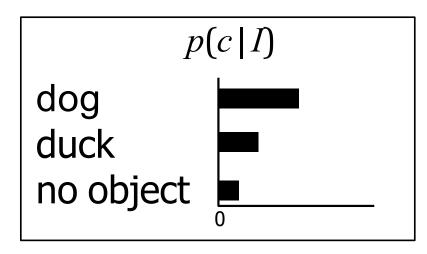
Bounding box (x, y, w, h)

Need multiple scales and aspect ratios

Idea #2: sliding window

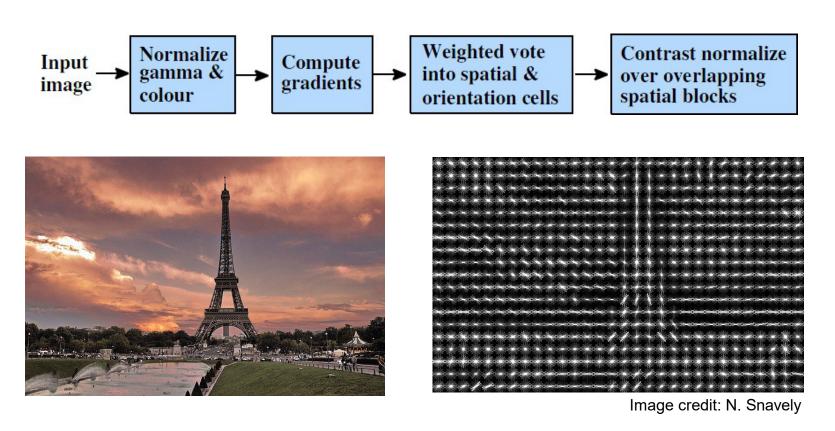






Bounding box (x, y, w, h)

Another Sliding Window Approach: Histograms of oriented gradients (HOG)



Example: pedestrian detection with HOG

Train a pedestrian template using a linear classifier. Represent each window using HOG

positive training examples



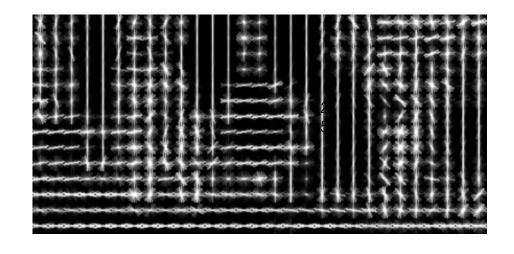
negative training examples



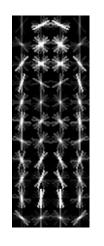
Example: pedestrian detection with HOG

For multi-scale detection, repeat over multiple levels of a HOG pyramid

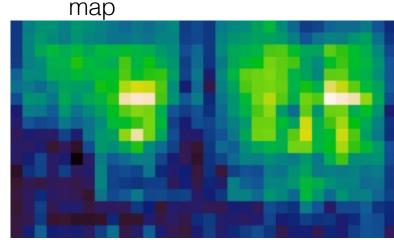
HOG feature map



Template

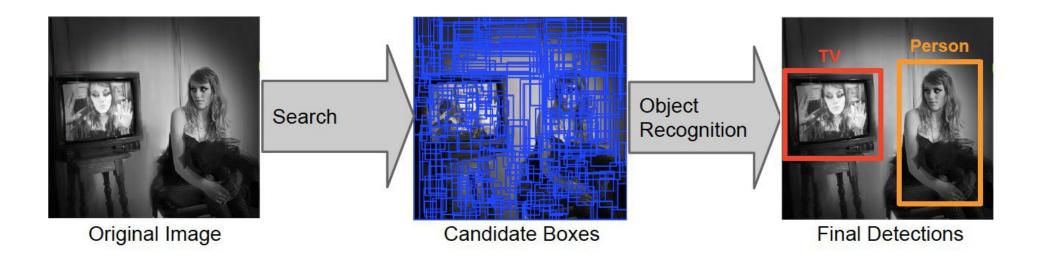


Detector response



Source: S. Lazebnik

Idea #3: selective search



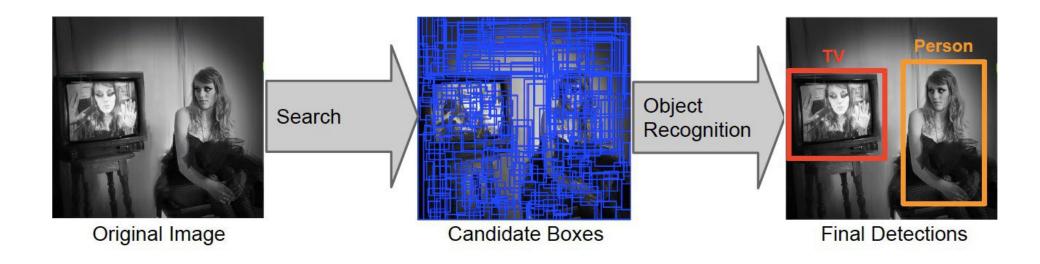
- Problem: evaluating a detector is very expensive
- An image with $m{n}$ pixels has $O(n^2)$ windows
- Only generate and evaluate a few hundred region proposals for regions that are "likely" to be an object of interest.

Selective search



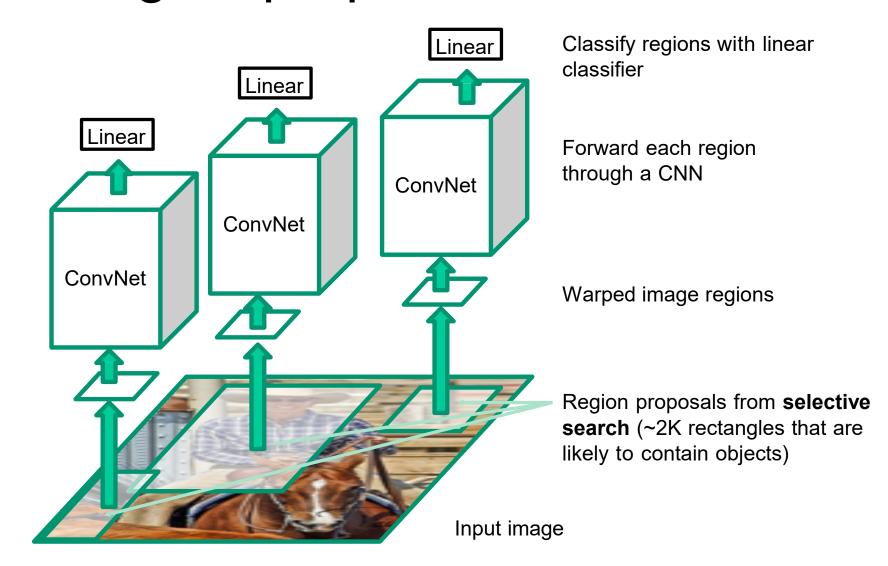
- Example: edge boxes
 [Zitnick & Dollar, 2014]
- Heuristic: detect edges, group them into contours
- Rank each window based on number of contours in window
- These are the only windows our detector will see

Selective search

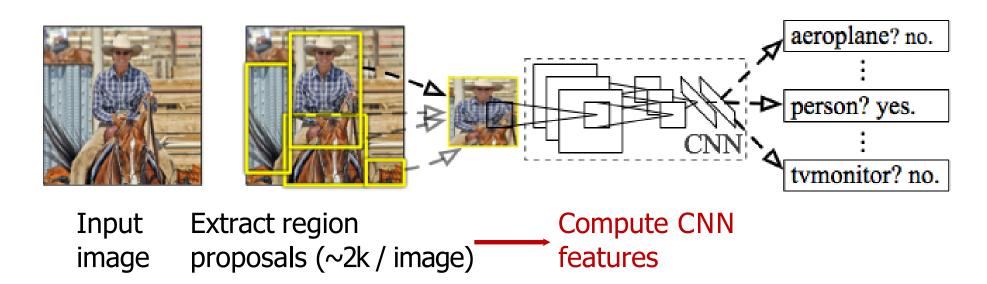


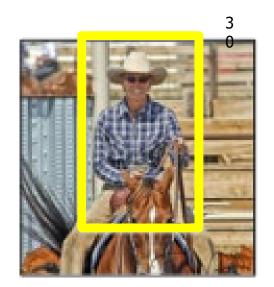
- Problem: evaluating a detector is very expensive
- An image with n pixels has $O(n^2)$ windows
- Only generate and evaluate a few hundred region proposals for regions that are "likely" to be an object of interest.

R-CNN: Region proposals + CNN features



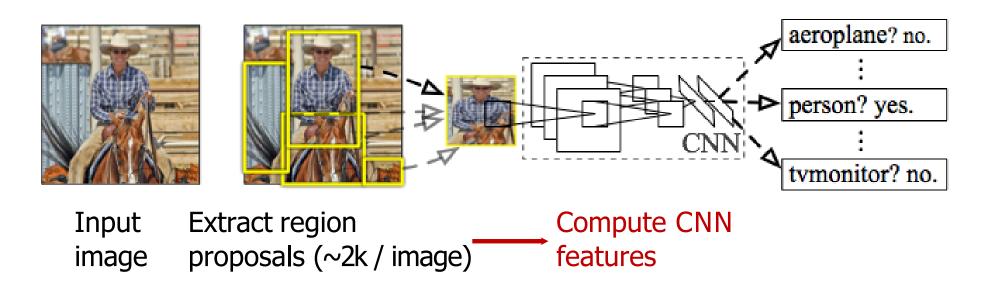
R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014

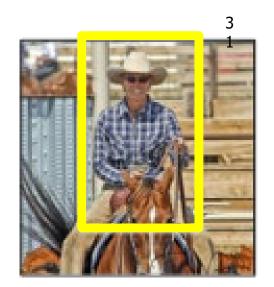






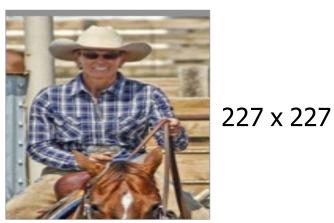
a. Crop



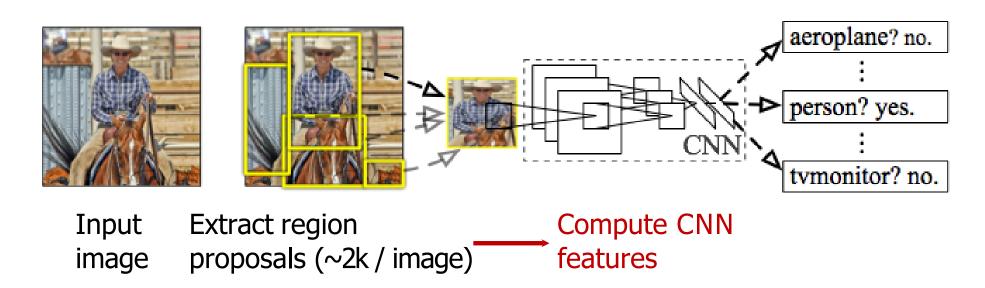




a. Crop



b. Scale



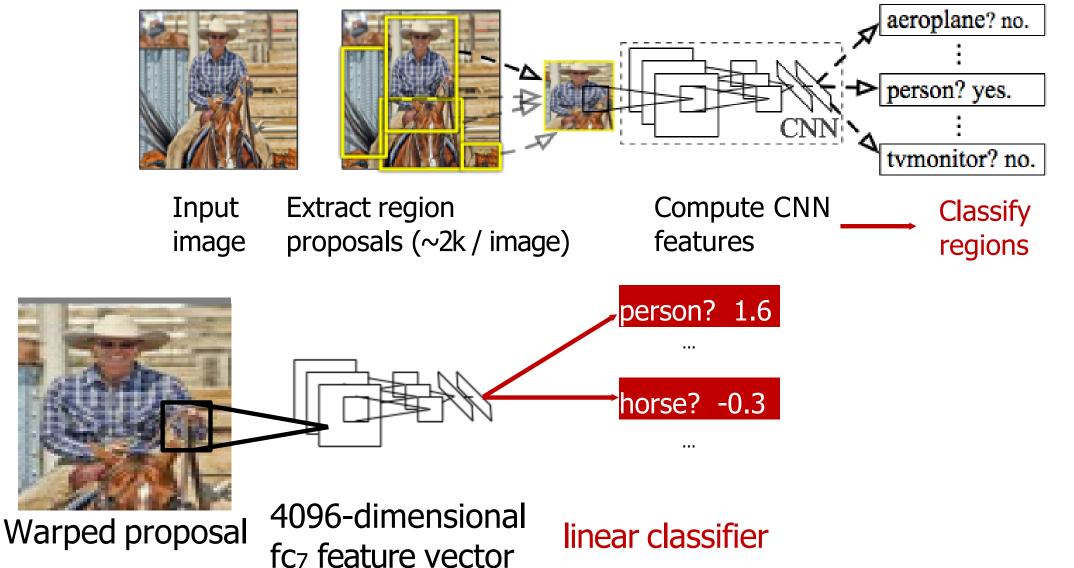


1. Crop



b. Scale

c. Forward propagate Output: "fc7" features

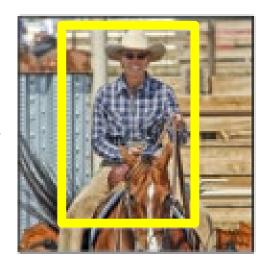


Proposal refinement



Linear regression

on CNN features

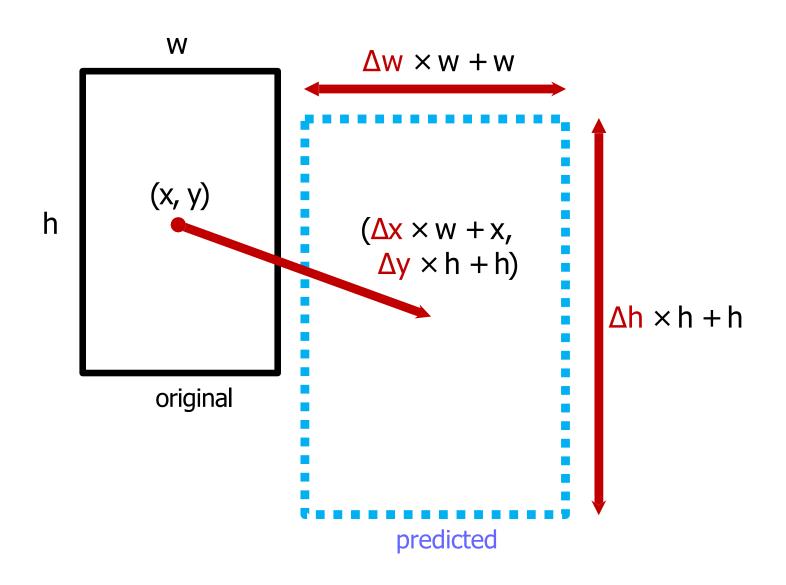


Original proposal

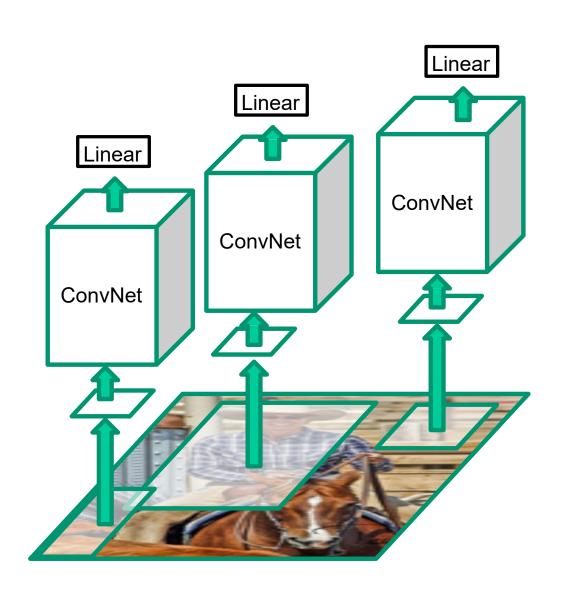
Predicted object bounding box

Bounding-box regression

Bounding-box regression

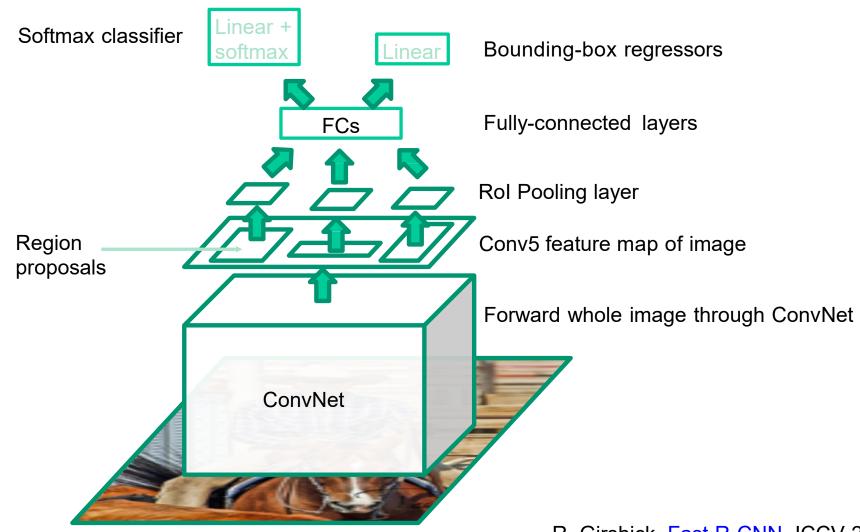


Problems with R-CNN



- 1. Slow! Have to run CNN per window
- 2. Hand-crafted mechanism for region proposal might be suboptimal.

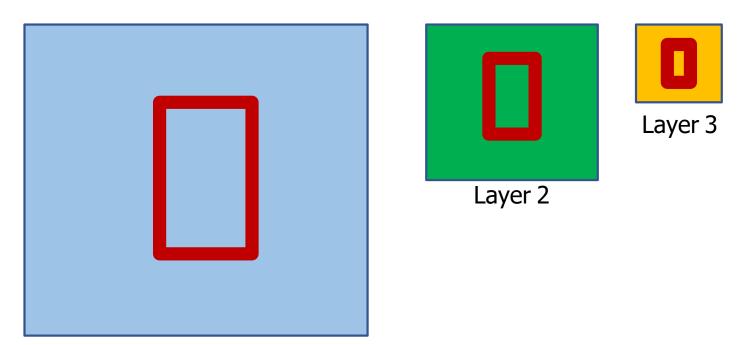
"Fast" R-CNN: reuse features between proposals



Source: R. Girshick

R. Girshick, Fast R-CNN, ICCV 2015

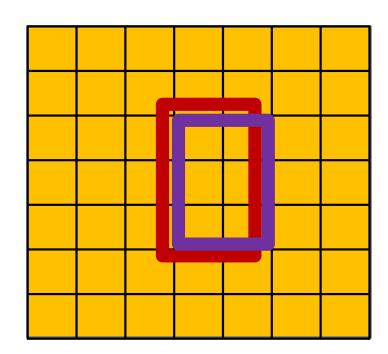
- How do we crop from a feature map?
- Step 1: Resize boxes to account for subsampling



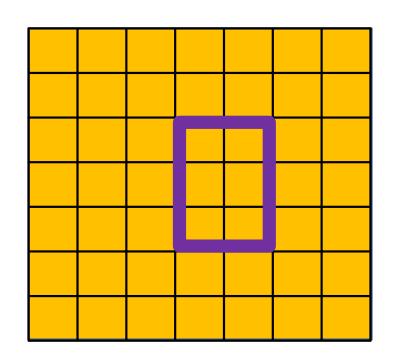
Layer 1

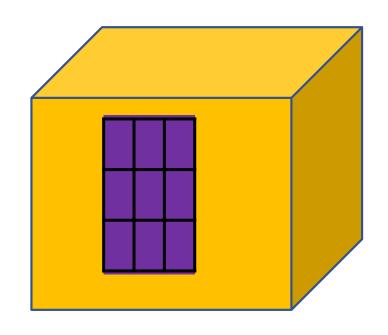
Source: B. Hariharan

- How do we crop from a feature map?
- Step 2: Snap to feature map grid

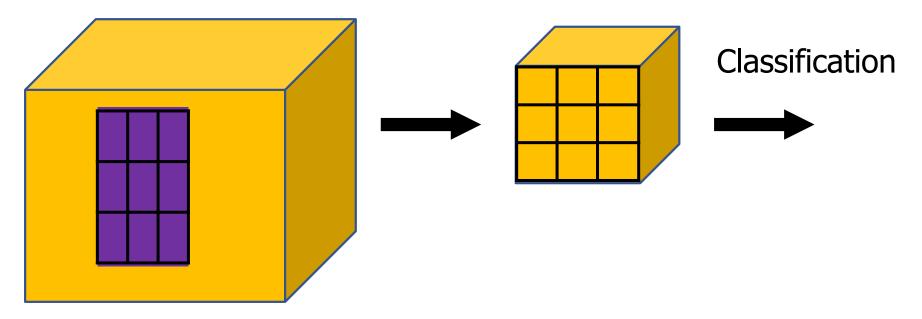


- How do we crop from a feature map?
- Step 3: Overlay a new grid of fixed size



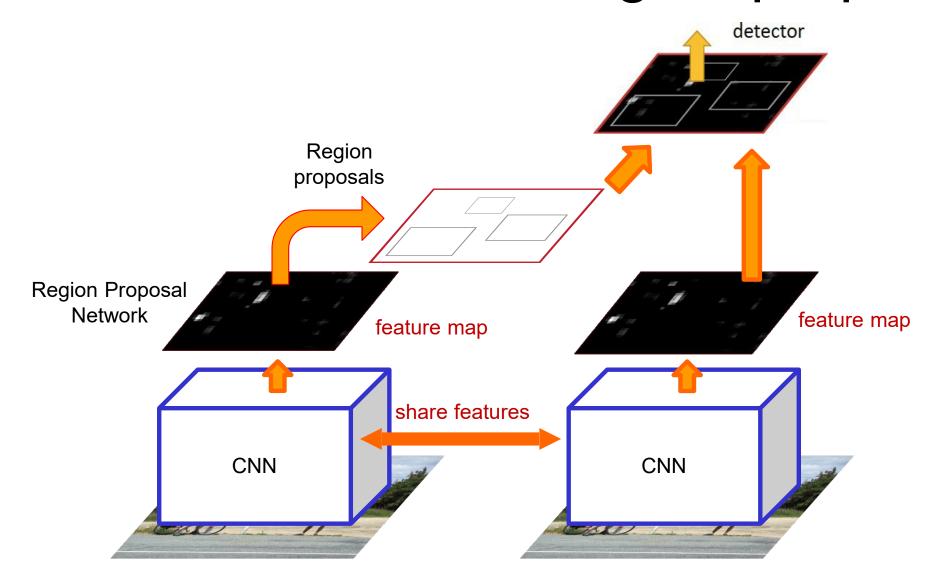


- How do we crop from a feature map?
- Step 4: Take max in each cell
- Can improve with bilinear sampling

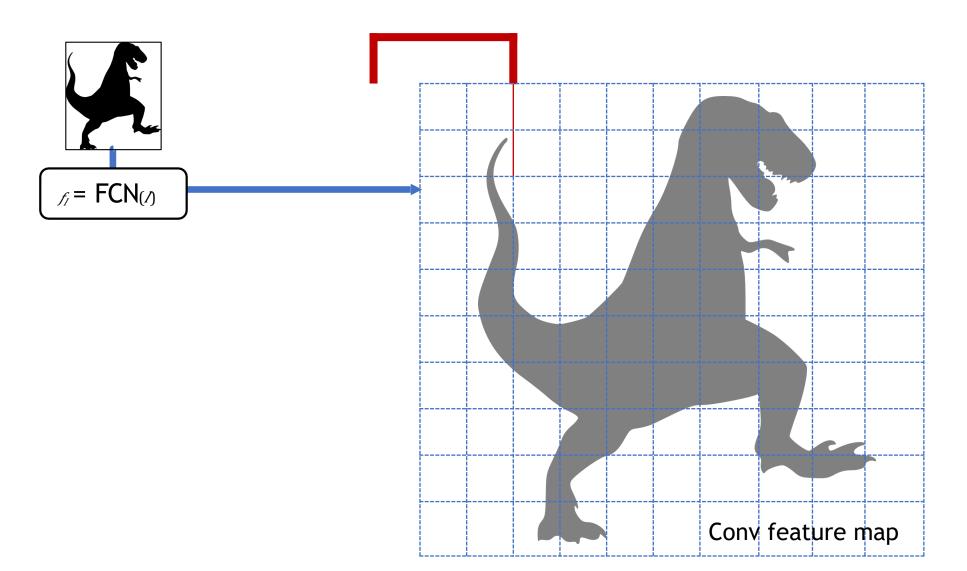


See more here: https://deepsense.ai/region-of-interest-pooling-explained/

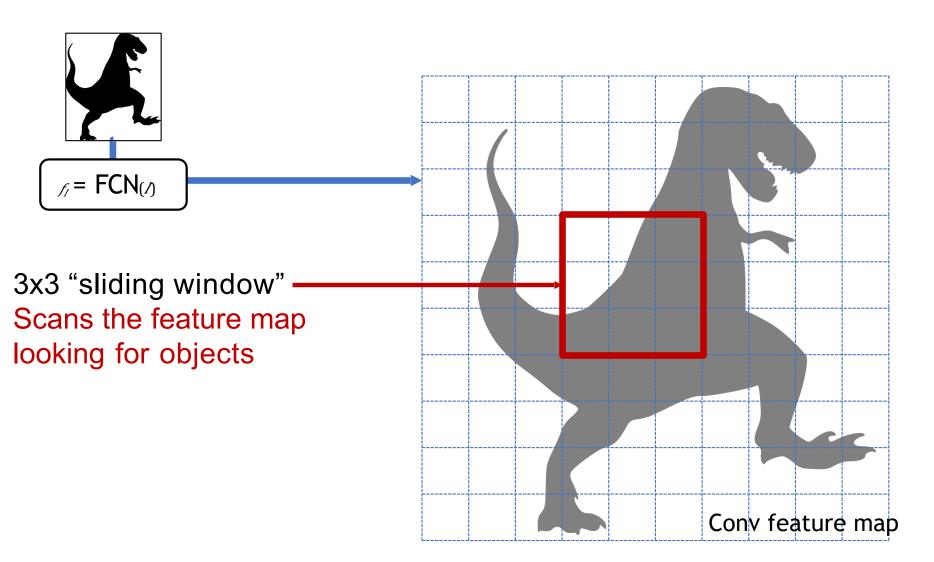
"Faster" R-CNN: learn region proposals



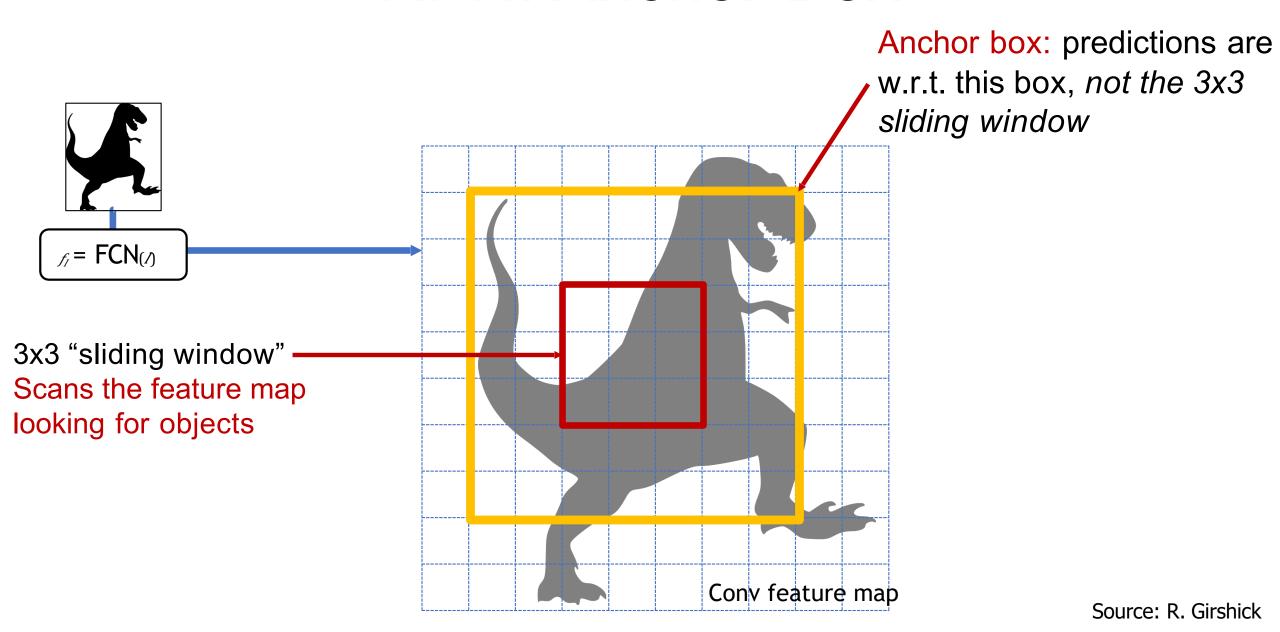
RPN: Region Proposal Network



RPN: Region Proposal Network



RPN: Anchor Box

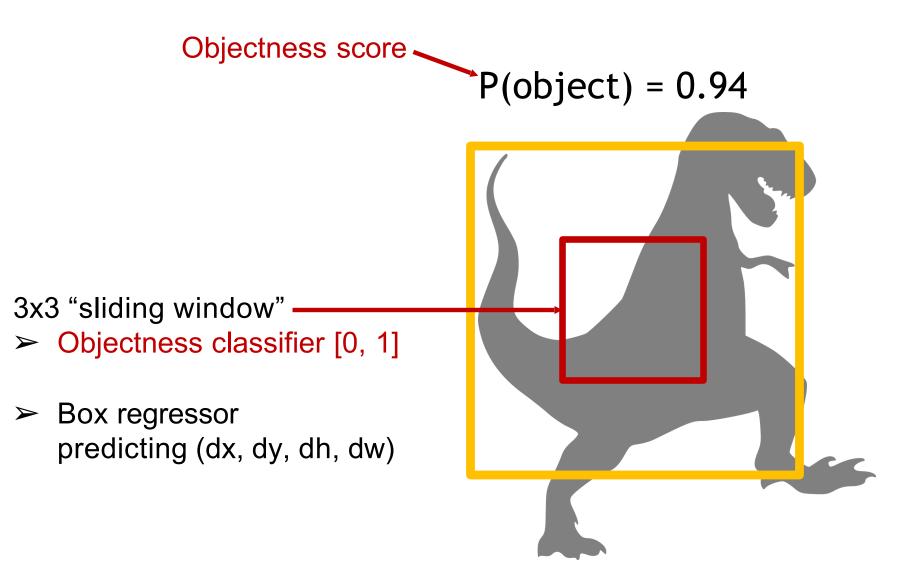


RPN: Anchor Box

Anchor box: predictions are w.r.t. this box, *not the 3x3* sliding window $f_I = FCN_{(I)}$ 3x3 "sliding window" ➤ Objectness classifier [0, 1] Box regressor predicting (dx, dy, dh, dw) Conv feature map

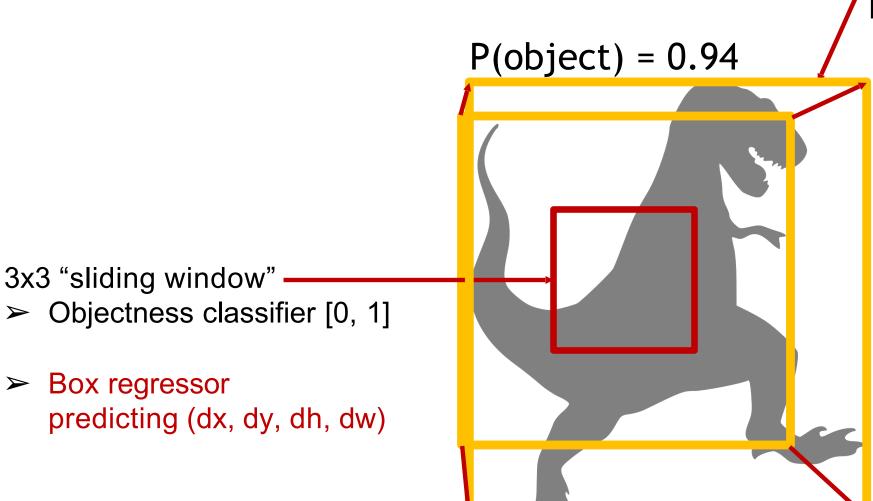
Source: R. Girshick

RPN: Prediction (on object)



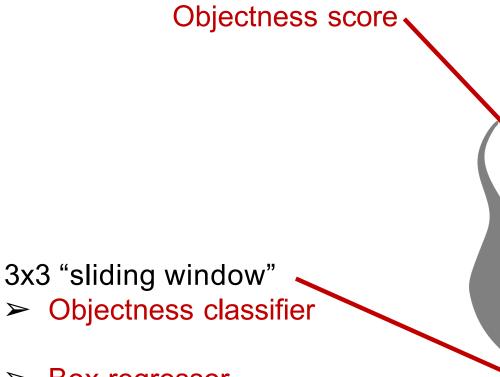
RPN: Prediction (on object)

Anchor box: transformed by box regressor



RPN: Prediction (off object)

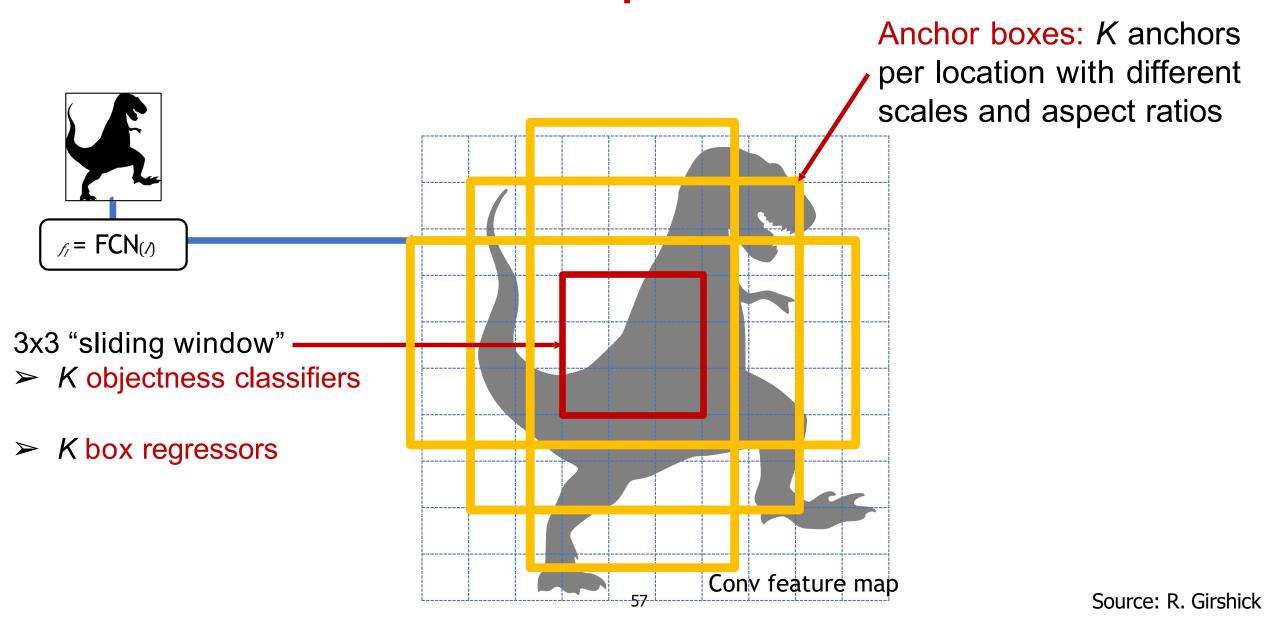
Anchor box: transformed by box regressor



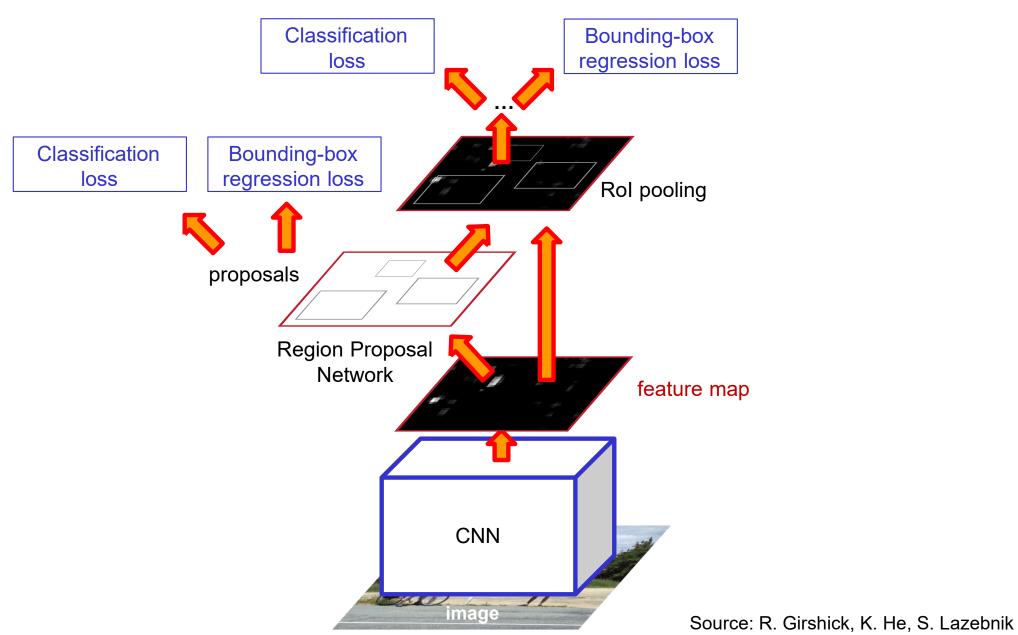
P(object) = 0.02

Box regressor predicting (dx, dy, dh, dw)

RPN: Multiple Anchors



One network, four losses



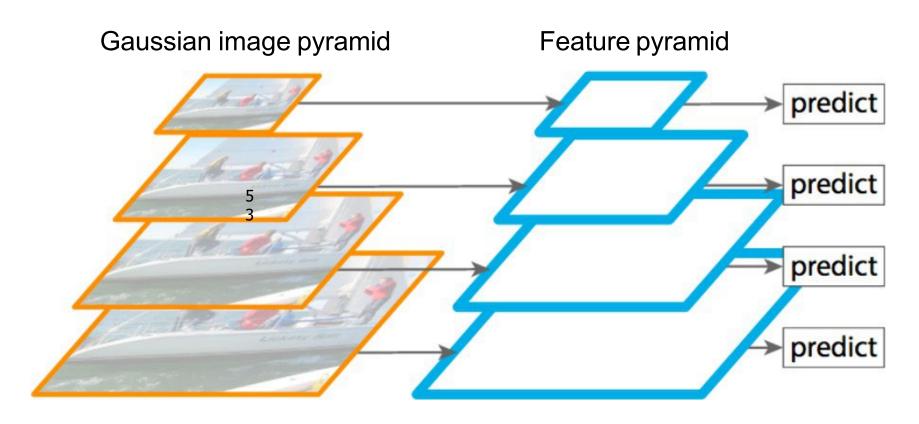
Faster R-CNN results

system	time	07 data	07+12 data
R-CNN	~50s	66.0	-
Fast R-CNN	~2s	66.9	70.0
Faster R-CNN	198ms	69.9	73.2

detection mAP on PASCAL VOC 2007, with VGG-16 pre-trained on ImageNet

Source: S. Lazebnik

How do we deal with scale? Idea #1: Gaussian pyramid

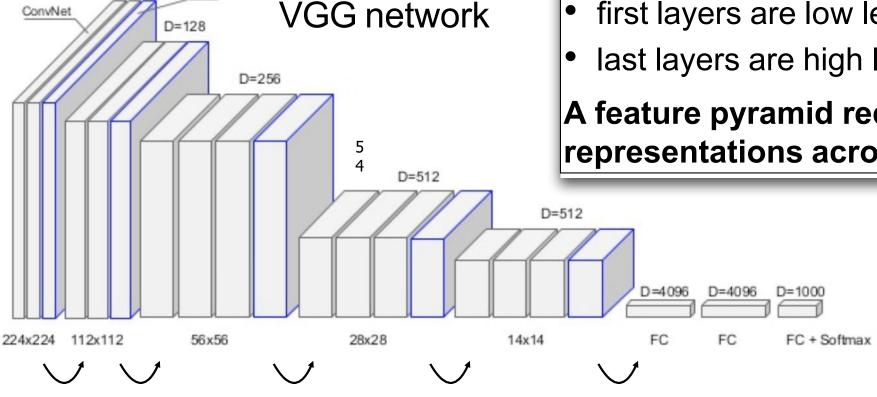


[Lin et al., "Feature Pyramid Networks for Object Detection", 2017]

Image and features pyramids

Each pooling reduces the resolution by a factor of 2

D=64

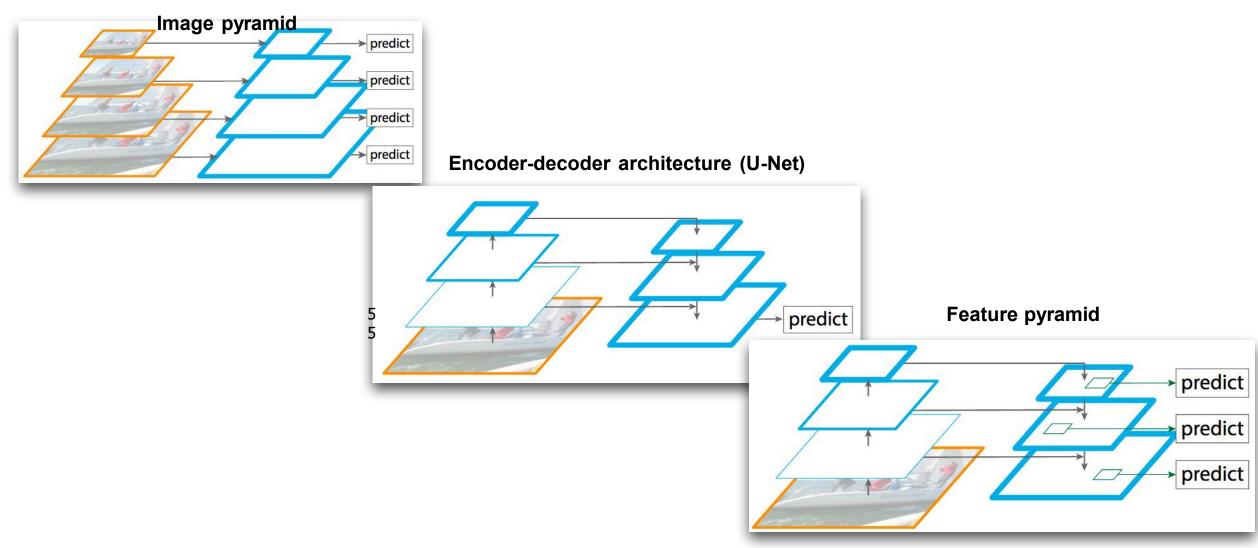


CNN architectures build:

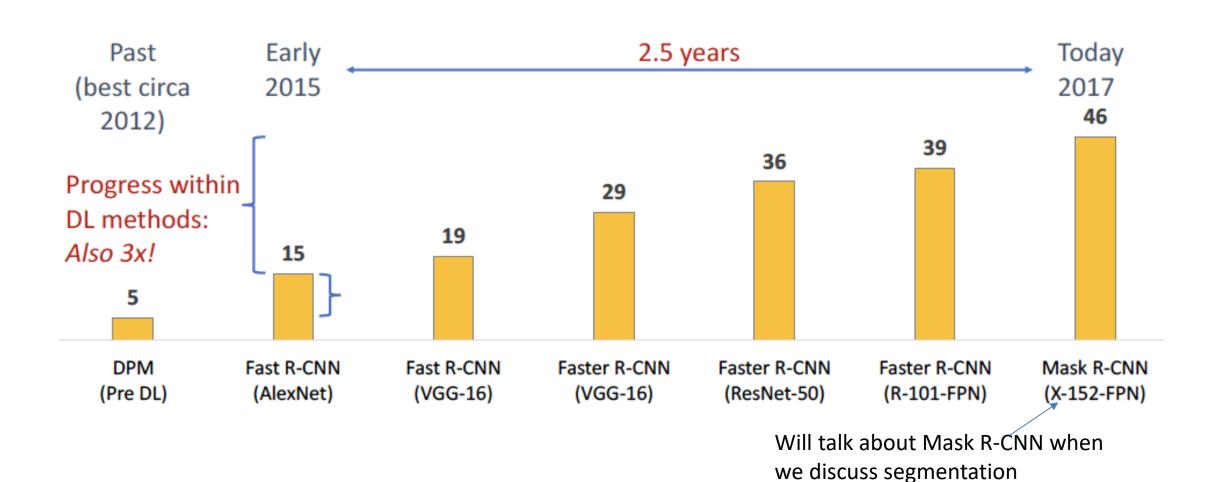
- Multiscale feature hierarchies, but
- each layer builds a different representation
- first layers are low level, while
- last layers are high level.

A feature pyramid requires a uniform representations across scales.

Idea #2: Feature pyramid network

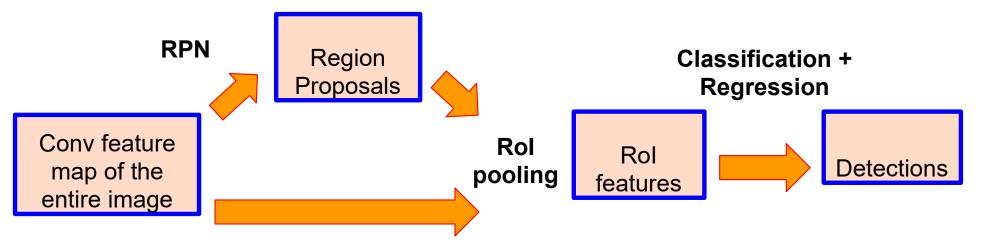


Summary: Progress through Better Detector Design and Better Backbones



Streamlined detection architectures

 The Faster R-CNN pipeline separates proposal generation and region classification:



Is it possible do detection in one shot?

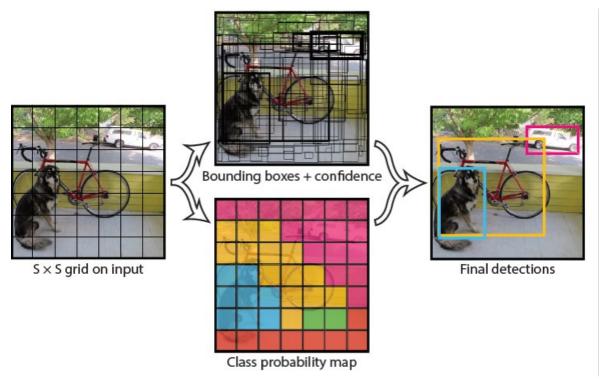


64

Source: S. Lazebnik

Single-stage object detector

- Divide the image into a coarse grid using a fully convolutional net
- Directly predict class label, confidence, and a few candidate boxes for each grid cell.



YOLO detector

- Take convolutional feature maps at 7x7 resolution
- Predict, at each location, a score for each class and 2 bounding boxes (with confidence)

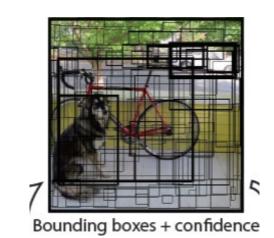
E.g. for 20 classes, output is 7x7x30 (30 = 20 + 2*(4+1))

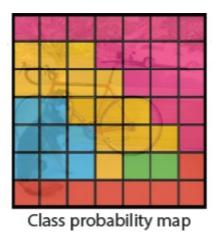
7x speedup over Faster R-CNN (45-155 FPS vs. 7-18 FPS)

but less accurate (e.g. 65% vs. 72 mAP%)

Extension:

use anchor boxes in last layer to try a few possible aspect ratios

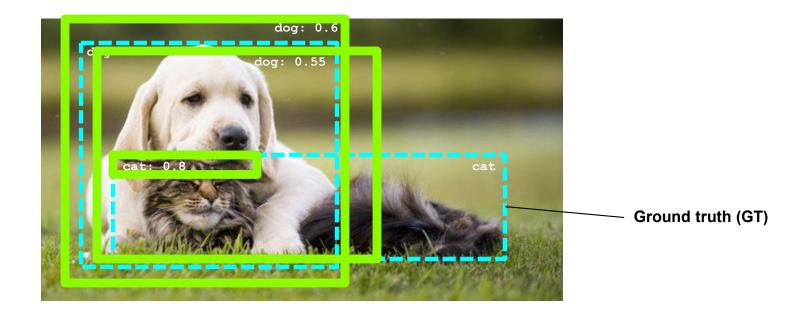


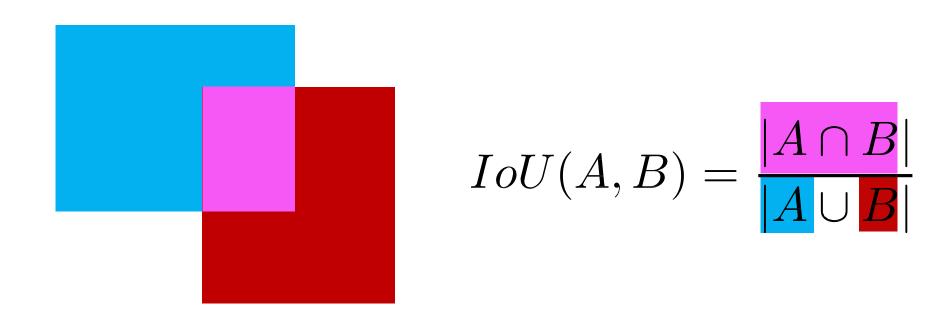


- Introduction to scene understanding
- Object detection models
- Evaluating object detectors
- Future challenges

- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
 Intersection over union (IoU):

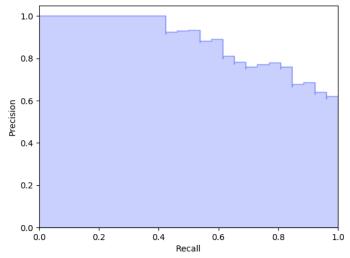
Area(GT \cap Det) / Area(GT \cup Det) > 0.5





Intersection over union (also known as Jaccard similarity)

- For each class, plot Precision-Recall curve and compute Average Precision (area under the curve)
- Take mean of AP over classes to get mAP



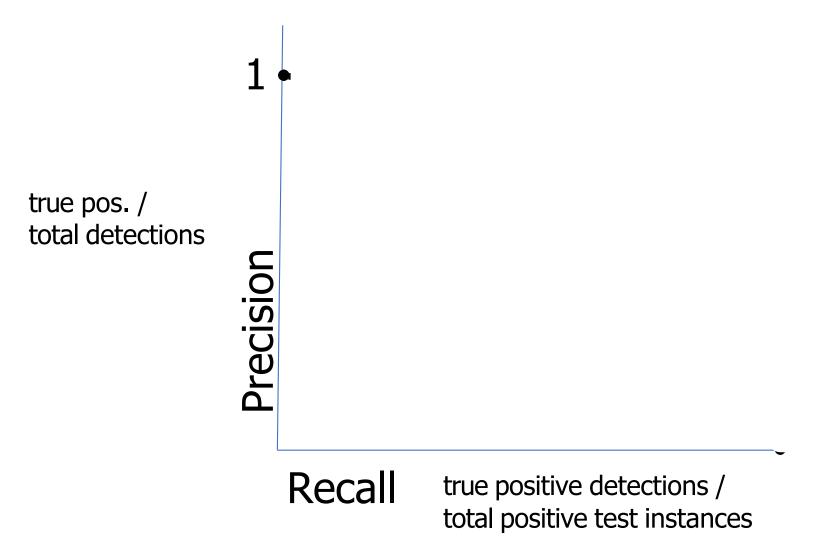
Precision:

true positive detections / total detections

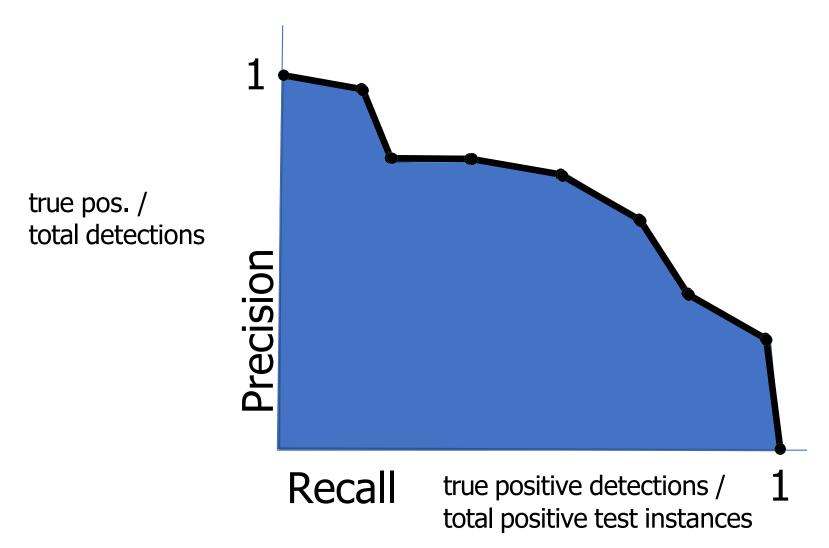
Recall:

true positive detections / total positive test instances

Average precision



Average precision



Non-maximum suppression



- Subtlety: we predict a bounding box for every sliding window. Which ones should we keep?
- Keep only "peaks" in detector response.
- Discard low-prob boxes near high-prob ones
- Often use a simple greedy algorithm

Non-maximum suppression

Greedy algorithm, run on each class independently

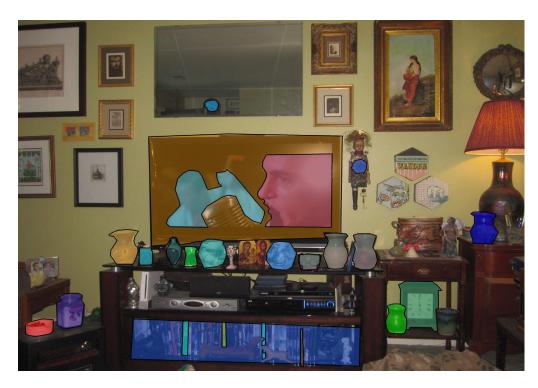
```
let be the set of all bounding boxes let be the set of detections we'll keep, D=\varnothing while A\neq\varnothing: remove the box with highest probability from if doesn't significantly overlap with an existing box in D=D\cup\{x\} return
```

- Introduction to scene understanding
- Object detection models
- Evaluating object detectors
- Challenges

Handle the long tail of the distribution

Person, dog, table, ... Frequency Kale, colander, Himalayan Salt, birdfeeder, humidifier, ... Object categories

Handle the "long tail" of the distribution



From COCO (80 categories) [Lin et al., 2014]



LVIS dataset (1000+ categories) "Few shot" (e.g. < 20 examples) [Gupta et al., 2019]

OWL-ViT: Open-vocabulary object detection model

[Submitted on 12 May 2022 (v1), last revised 20 Jul 2022 (this version, v2)]

Simple Open-Vocabulary Object Detection with Vision Transformers

Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, Neil Houlsby

Combining simple architectures with large-scale pre-training has led to massive improvements in image classification. For object detection, pre-training and scaling approaches are less well established, especially in the long-tailed and open-vocabulary setting, where training data is relatively scarce. In this paper, we propose a strong recipe for transferring image-text models to open-vocabulary object detection. We use a standard Vision Transformer architecture with minimal modifications, contrastive image-text pre-training, and end-to-end detection fine-tuning. Our analysis of the scaling properties of this setup shows that increasing image-level pre-training and model size yield consistent improvements on the downstream detection task. We provide the adaptation strategies and regularizations needed to attain very strong performance on zero-shot text-conditioned and one-shot image-conditioned object detection. Code and models are available on GitHub.

Comments: ECCV 2022 camera-ready version

OWL-ViT: Open-vocabulary object detection model

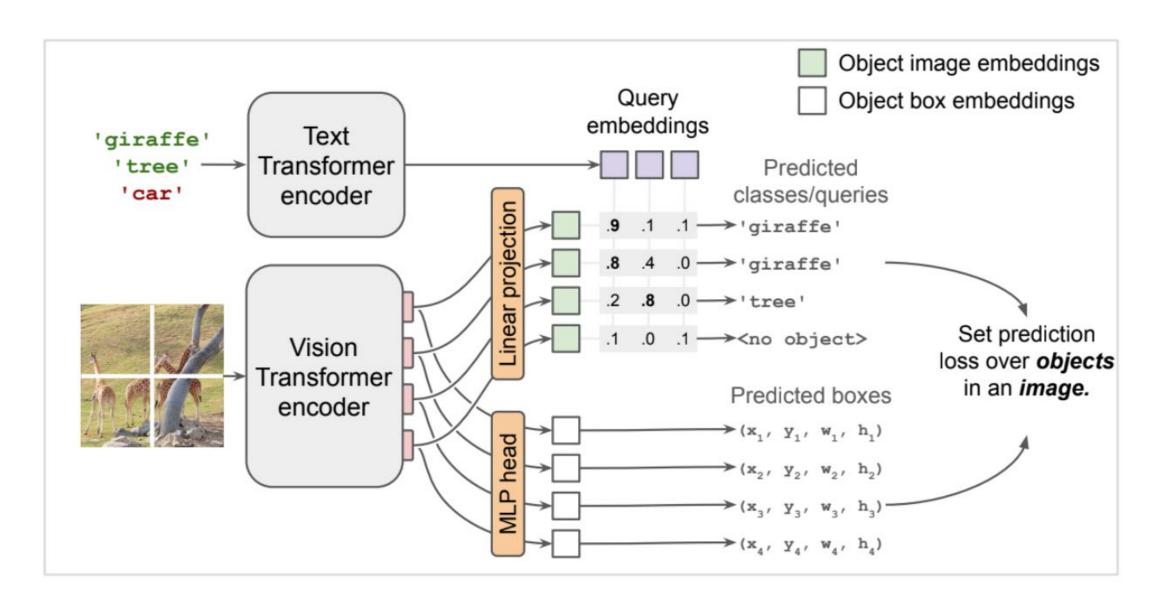


Image-level contrastive pre-training Transfer to open-vocabulary detection Object image embeddings Query Object box embeddings Text embeddings embedding Text Text 'bird 'giraffe' Transformer sitting Transformer Predicted on a tree 'car' encoder encoder classes/queries Set prediction Token .1 → <no object> Vision Vision Contrastive loss over objects pooling loss over ransformer Transformer in an image. Predicted boxes images in a encoder encoder projection batch. embedding

Fig. 1. Overview of our method. Left: We first pre-train an image and text encoder contrastively using image-text pairs, similar to CLIP [33], ALIGN [19], and LiT [44]. Right: We then transfer the pre-trained encoders to open-vocabulary object detection by removing token pooling and attaching light-weight object classification and localization heads directly to the image encoder output tokens. To achieve open-vocabulary detection, query strings are embedded with the text encoder and used for classification. The model is fine-tuned on standard detection datasets. At inference time, we can use text-derived embeddings for open-vocabulary detection, or image-derived embeddings for few-shot image-conditioned detection.

OWL-ViT Demo

https://colab.research.google.com/github/huggingface/notebooks/blob/main/examples/zeroshot object detection with owlvit.ipynb

https://colab.research.google.com/drive/1evZkcsq4FTreFxGV6JXDmymnYcq WK43n