

Neural Networks 3/3

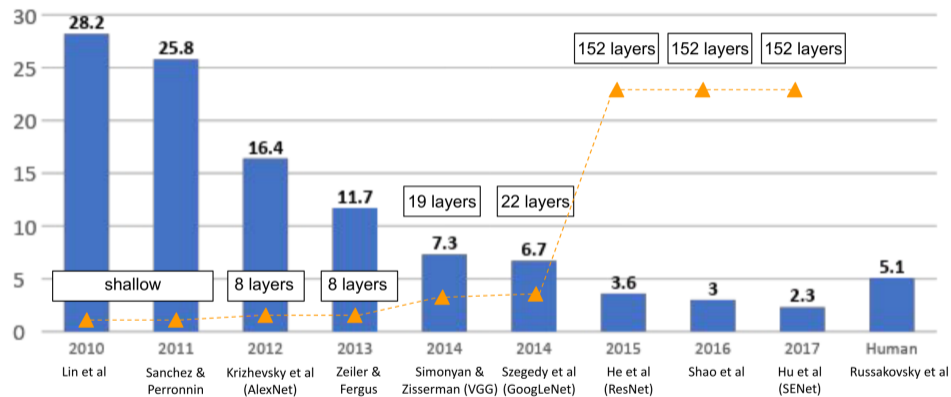
Lecture 12

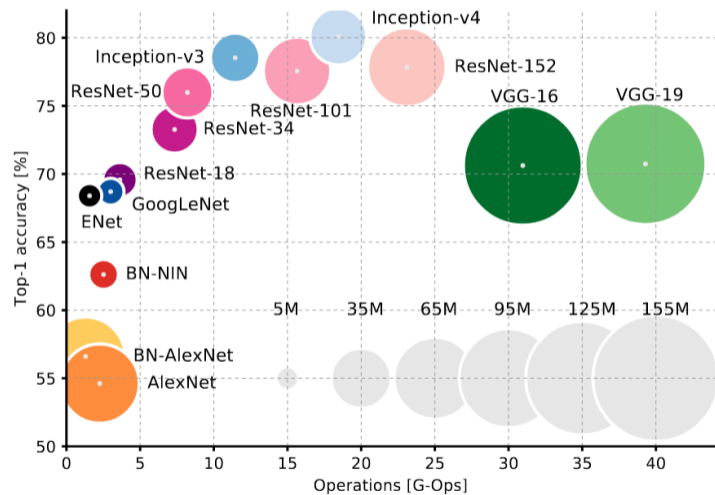
Computer Vision for Geosciences

June 4, 2021



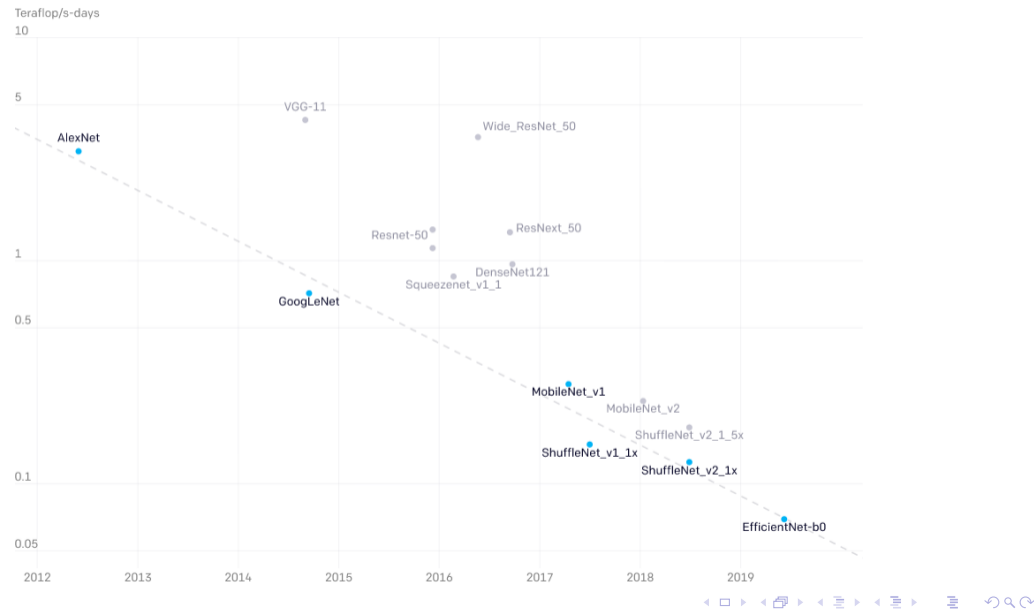
- Image from Stanford CS231n Lecture 9, Fei-Fei Li
http://cs231n.stanford.edu/slides/2021/lecture_9.pdf





- Image from An Analysis of Deep Neural Network Models for Practical Applications, Canziani et al, 2017

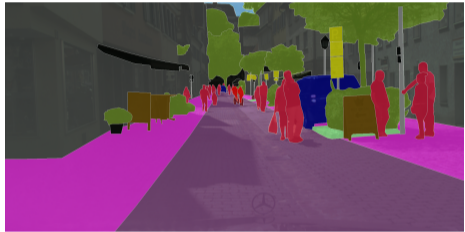
Efficiency



- Total amount of compute in teraflops/s-days used to train to AlexNet level performance. Lowest compute points at any given time shown in blue, all points measured shown in gray.
- Image from <https://openai.com/blog/ai-and-efficiency/>

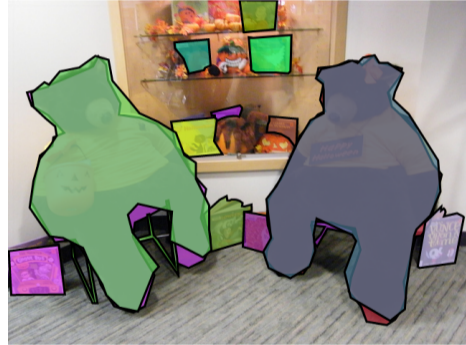


- Semantic Segmentation is the task of classifying every pixel of an image with an object class.
- Often including a background class.



- ▶ 30 classes
- ▶ 5000 annotated images with fine annotation
- ▶ 20000 annotated images with coarse annotations

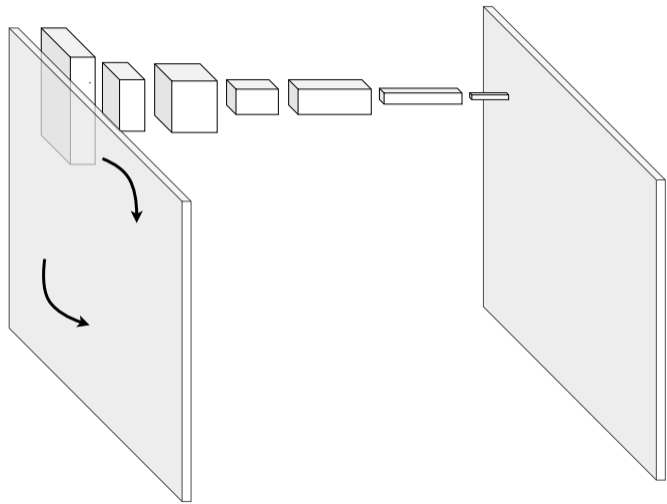
-
-
-



- ▶ 1.5 million object instances
- ▶ 80 object categories
- ▶ 91 stuff categories
- ▶ 330K images (>200K labeled)

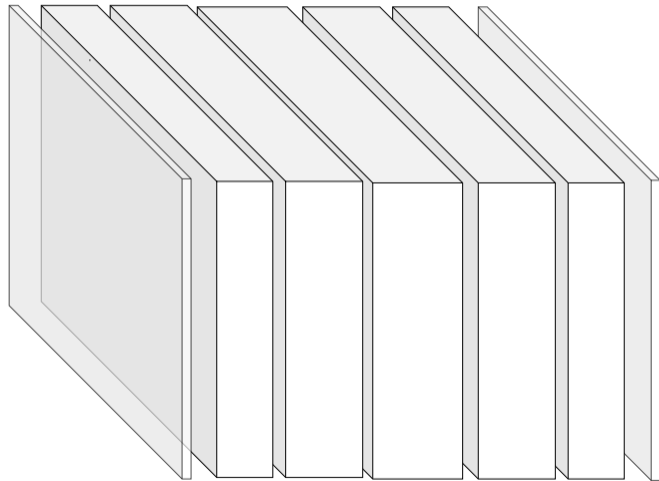
-
-
-

Semantic Segmentation: sliding window?



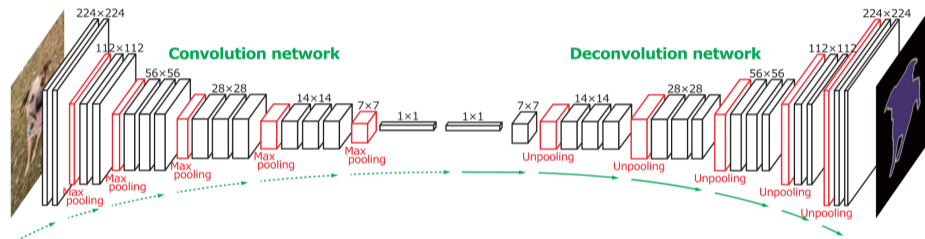
-
-
-

Semantic Segmentation: without downsampling?

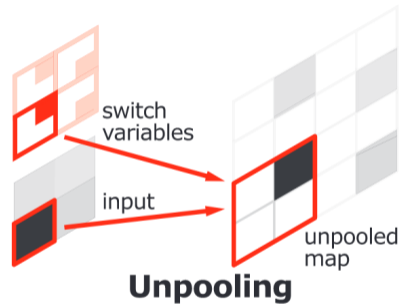
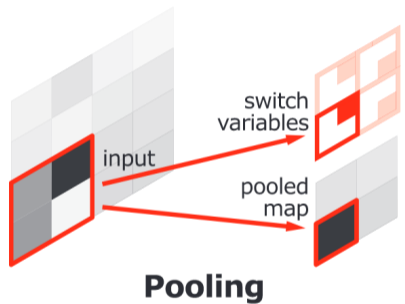


-
-
-

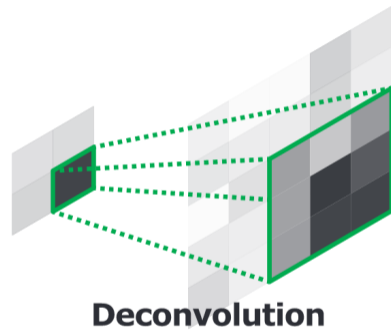
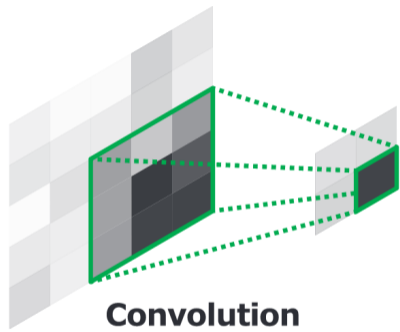
Encoder-Decoder-Architecture



-
-
- Image from Learning Deconvolution Network for Semantic Segmentation, Noh et al, ICCV 2015

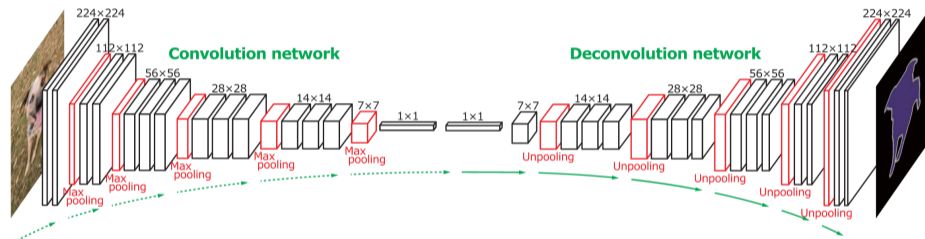


- Nearest Nighbour
- Bed of Nails
- Max unpooling
- Image from Learning Deconvolution Network for Semantic Segmentation, Noh et al, ICCV 2015

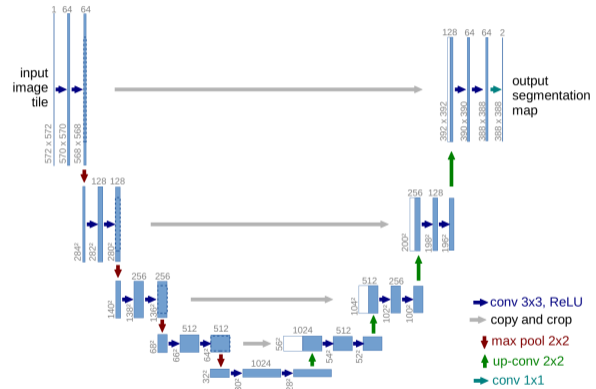


- Transpose convolution, deconvolution
- stride 2, pad 1, the other way
- Image from Learning Deconvolution Network for Semantic Segmentation, Noh et al, ICCV 2015

Encoder-Decoder-Architecture

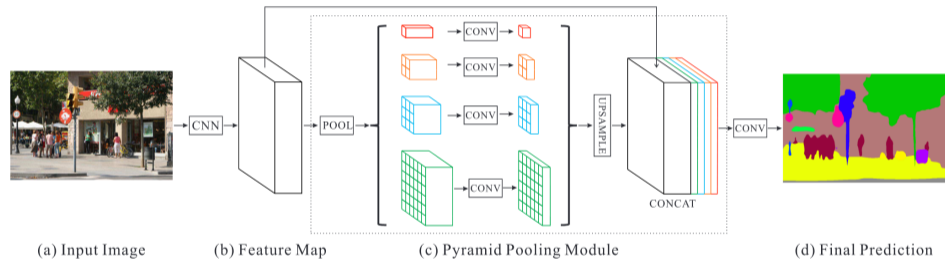


-
-
- Image from Learning Deconvolution Network for Semantic Segmentation, Noh et al, ICCV 2015

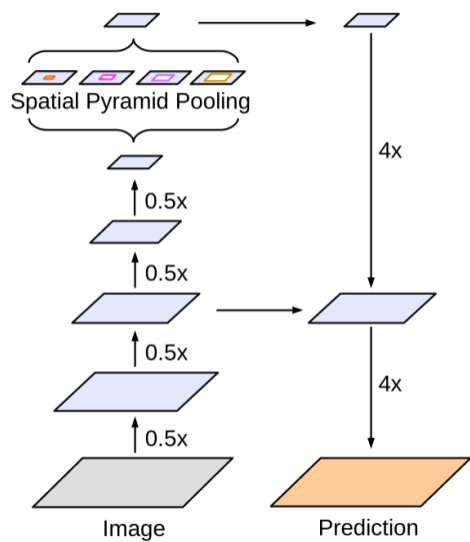


-
- Image from U-Net: Convolutional Networks for Biomedical Image Segmentation, Ronnenberger et al, MICCAI 2015
- SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, Badrinarayanan et al, TPAMI 2017

Pyramid Pooling

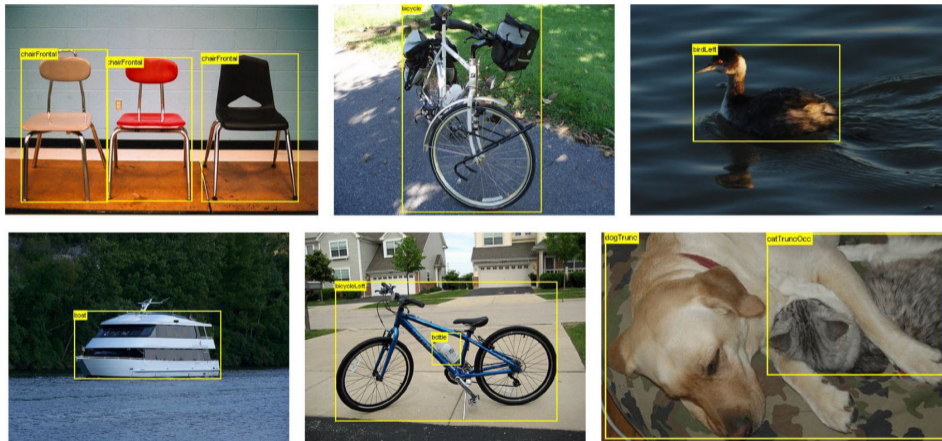


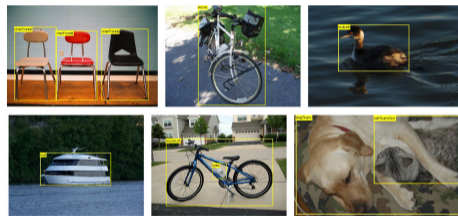
-
-
- Image from Pyramid Scene Parsing Network, Zhao et al, CVPR 2017



-
-
- Image from Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation, Chen et al, ECCV 2018

- Image from The PASCAL Visual Object Classes Challenge: A Retrospective, Everingham et al, IJCV 2014





- ▶ 20 classes
- ▶ 11k annotated images
- ▶ 27k annotated objects

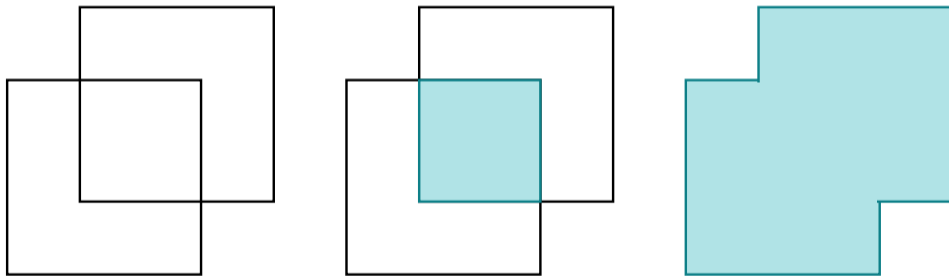
- Pascal VOC (DPM 33.6%)

-
-

- Default threshold was 0.5 for a long time but is now often higher.

Detection is correct if

$$\textit{intersection/union} > \textit{threshold}$$

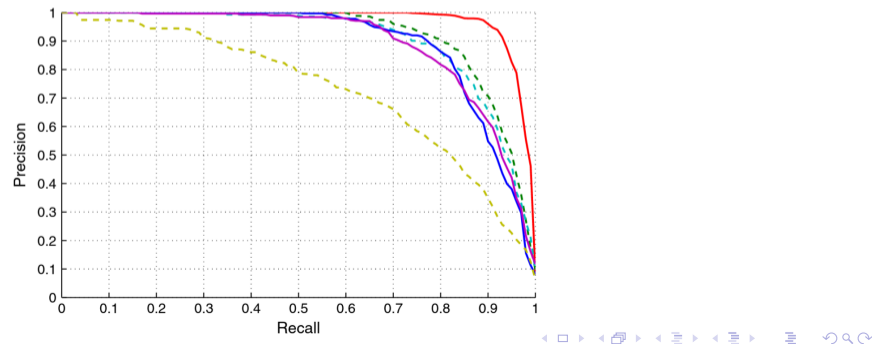


$$\textit{precision} = \#(\textit{correct detections}) / \#(\textit{all objects})$$

$$\textit{recall} = \#(\textit{correct detections}) / \#(\textit{all detections})$$

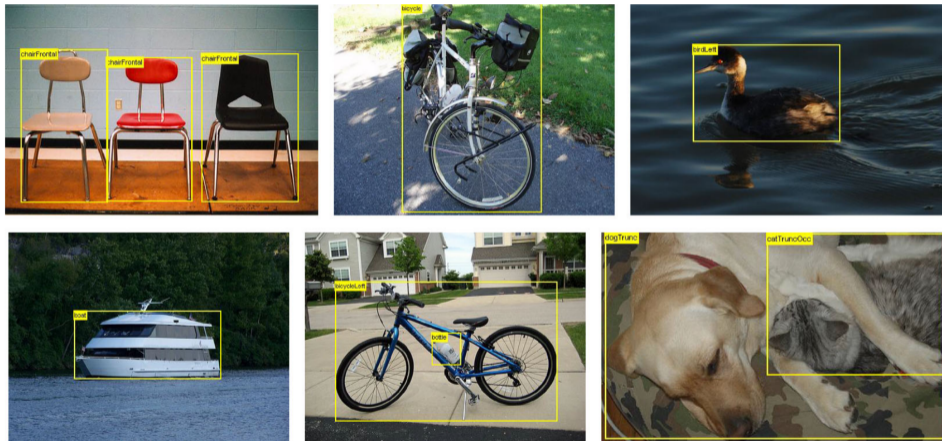
Average Precision: area under PR curve for specific class

mean Average Precision: AP averaged over all classes



-
-
- Image from The PASCAL Visual Object Classes Challenge: A Retrospective, Everingham et al, IJCV 2014

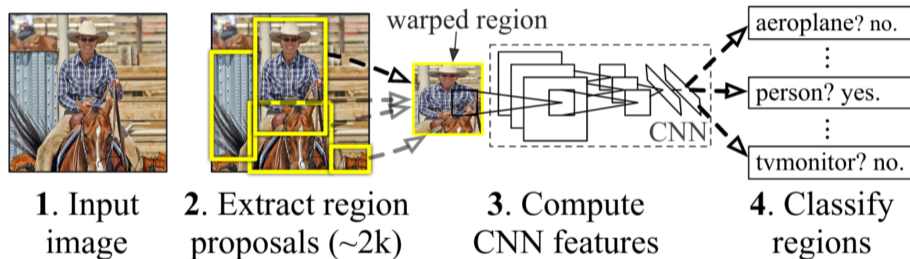
Object Detection: output dimensionality?



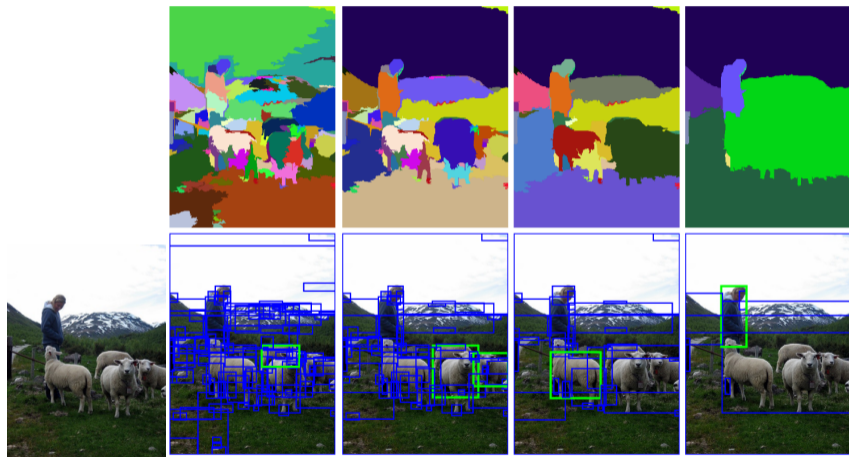
- How would the head of this network look like?
- Image from The PASCAL Visual Object Classes Challenge: A Retrospective, Everingham et al, IJCV 2014

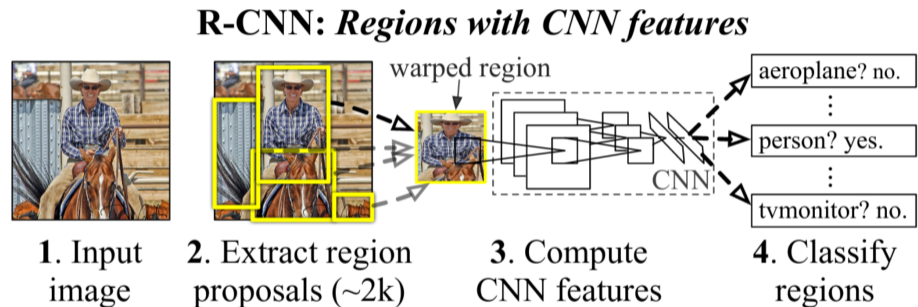
- Same author as DPM.
- Sliding window as in DPM. But NN much slower as SVM, therefore they used region proposals (2k).
- Image from Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al, CVPR 2014

R-CNN: *Regions with CNN features*

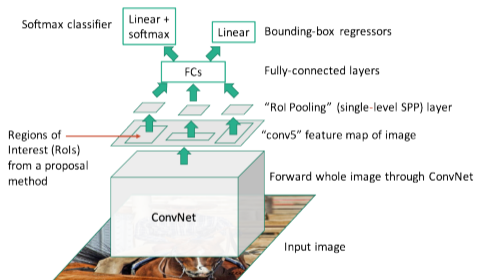
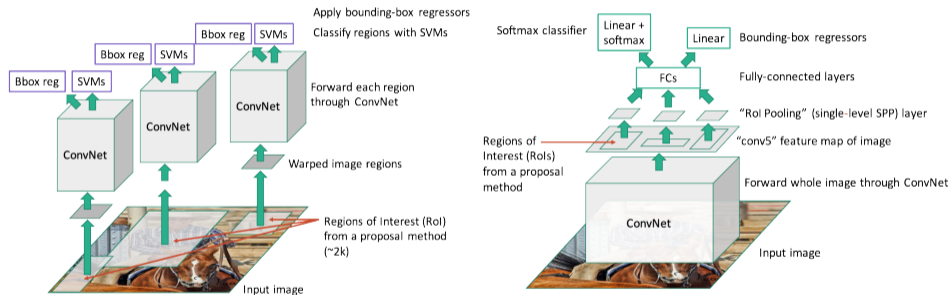


- Image from Selective Search for Object Recognition, Uijlings et al, IJCV 2013

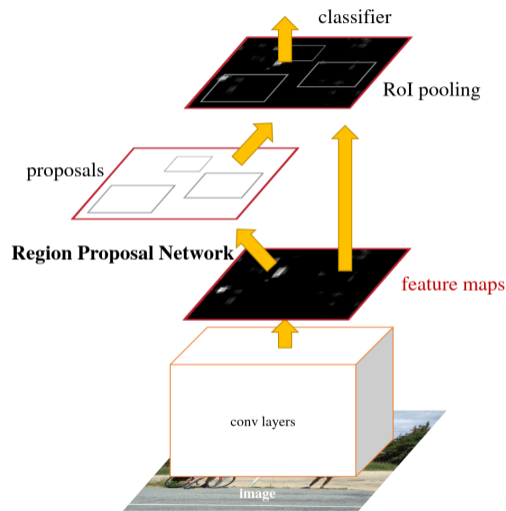




- Network also needs to predict bounding box parameters (size and offset from patch center).
- Non maximum suppression in prediction space.
- Often some high level reasoning (coherence in object relations).
- mAP for Pascal VOC improved to 53% with AlexNet as ConvNet and 62% with VGG (from 33% DPM)
- Image from Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al, CVPR 2014

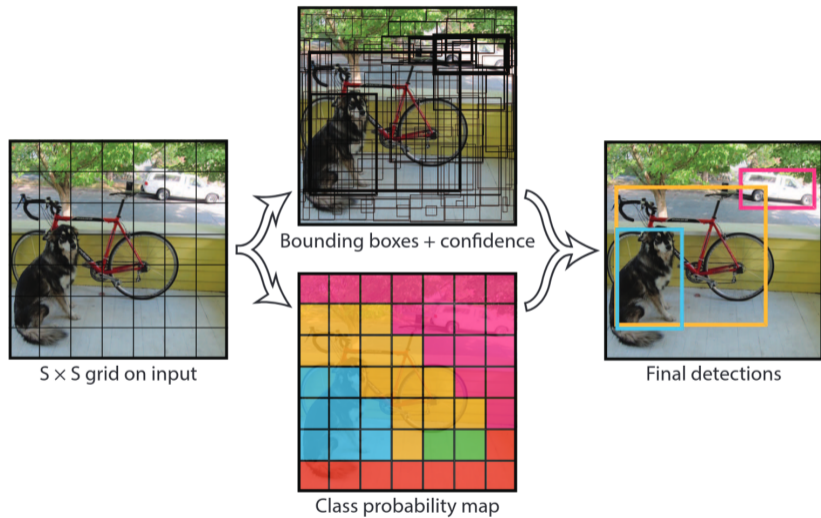


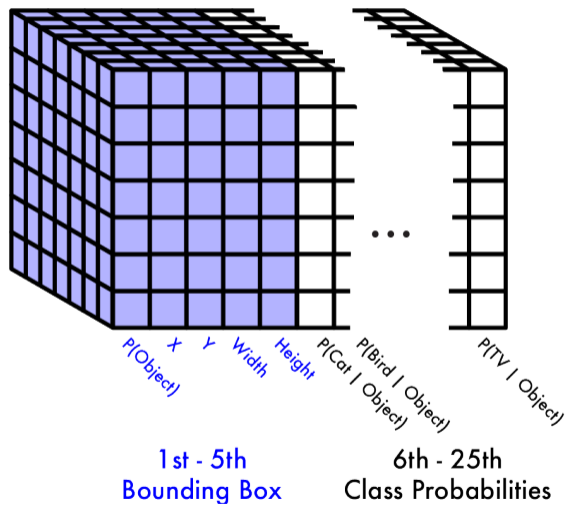
-
- Image from Talk at ICCV 2015 by Ross Girshick
<https://dl.dropboxusercontent.com/s/vlyrkgd8nz8gy5l/fast-rcnn.pdf?dl=0>



- Region proposal is now the expensive step in Fast-RNN
- Solution: Do region proposal in feature map.

- Image from You Only Look Once: Unified, Real-Time Object Detection, Redmon et al, CVPR 2016





- Newer versions of YOLO have multiple detections per cell for different object sizes.
- Image from Ancient Secrets of Computer Vision Lecture 18, Joseph Redmon

- weighted loss, binary and multi-class cross entropy, MSE
- What would happen without conditional probability?

$$\mathcal{L} = \alpha_1 \mathcal{L}_{\text{localization}} + \alpha_2 \mathcal{L}_{\text{object confidence}} + \alpha_3 \mathcal{L}_{\text{classification}}$$

$\mathcal{L}_{\text{localization}}$: root mean squared error

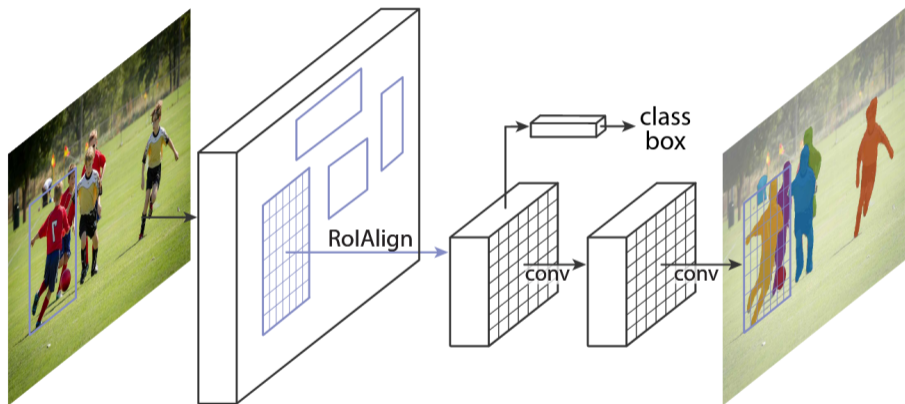
$\mathcal{L}_{\text{object confidence}}$: binary cross entropy

$\mathcal{L}_{\text{classification}}$: multi – class cross entropy

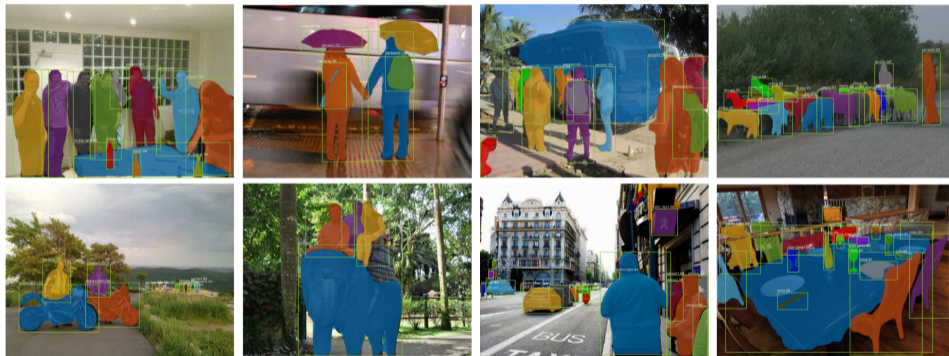
Why not both? Instance Segmentation

-
-



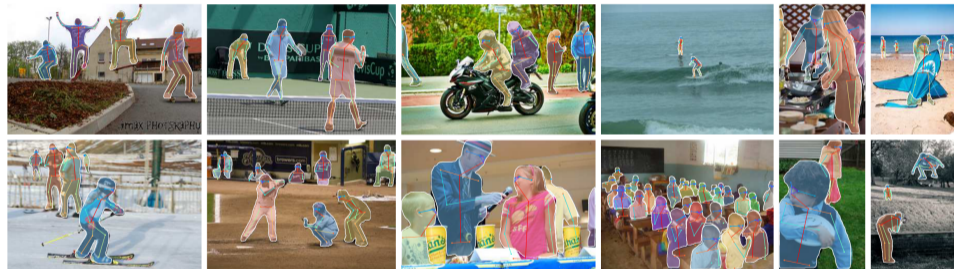


-
-
- Image from Mask R-CNN, He et al, ICCV 2017

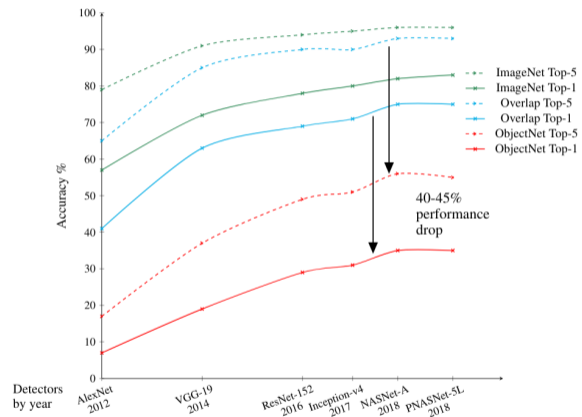
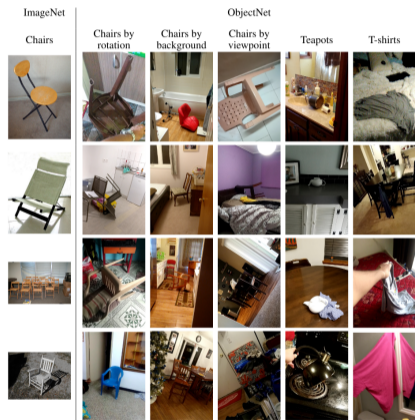


-
-
- Image from Mask R-CNN, He et al, ICCV 2017

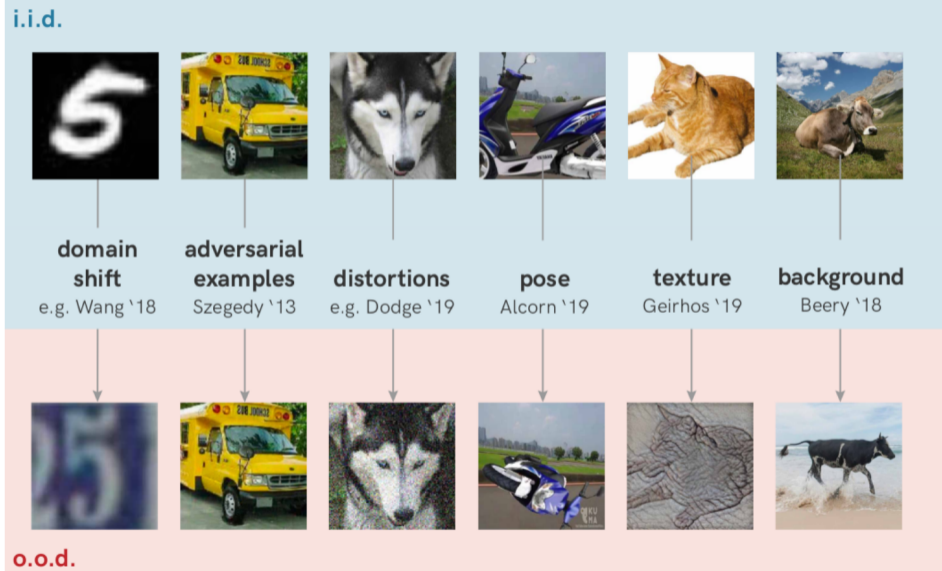
-
-
- Image from Mask R-CNN, He et al, ICCV 2017



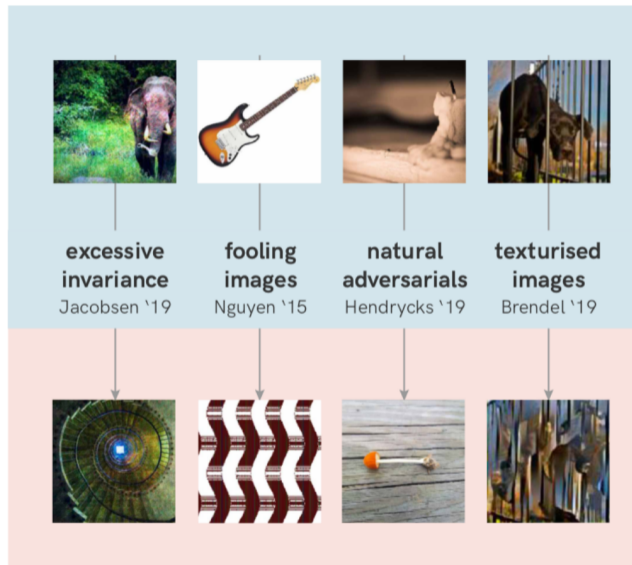
Generalization



- Depth adds complexity in training.
-
- Image from ObjectNet: A large-scale bias-controlled dataset for pushing the limits of object recognition models, Barbu et al, NeurIPS 2019

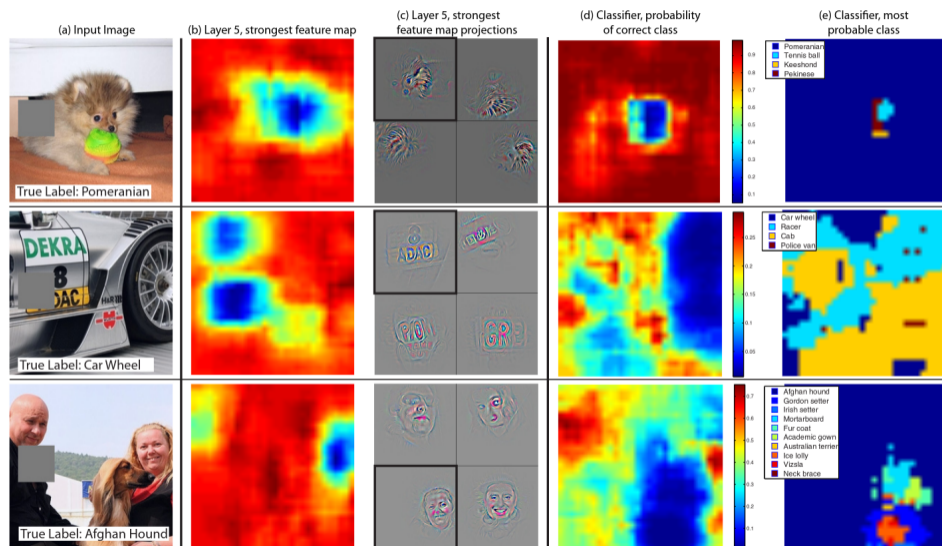


-
-
- Image from Shortcut Learning in Deep Neural Networks, Geirhos et al, Nature Machine Intelligence 2020



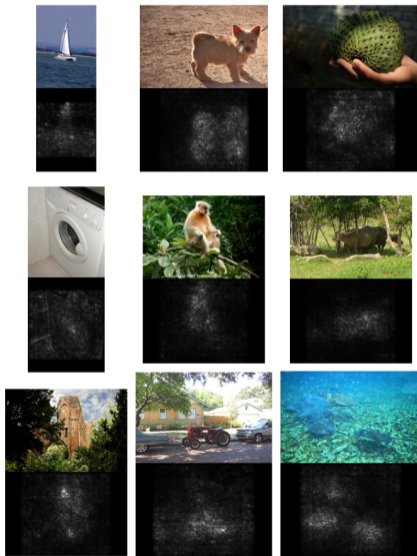
-
-
- Image from Shortcut Learning in Deep Neural Networks, Geirhos et al, Nature Machine Intelligence 2020

Investigate decisions: partial occlusion



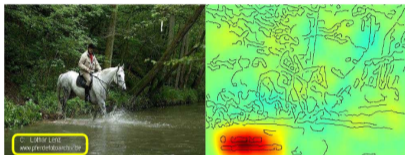
-
-
- Image from Visualizing and Understanding Convolutional Networks, Zeiler & Fergus, ECCV 2014

Investigate decisions: image gradient



-
-
- Image from Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al, 2013

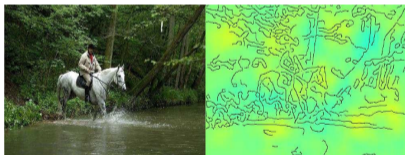
Horse-picture from Pascal VOC data set



Source tag present



Classified as horse

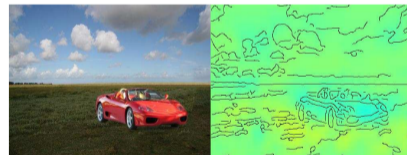


No source tag present



Not classified as horse

Artificial picture of a car



- Explain the output, not the local variation.
- Image from Unmasking Clever Hans Predictors and Assessing What Machines Really Learn, Lapuschkin et al, Nature Communications 2019