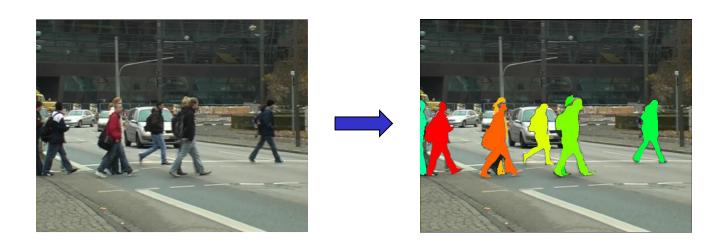
Lecture 9: CNNs in Computer Vision II: Object Localization

Lan Xu SIST, ShanghaiTech Fall, 2021

Outline

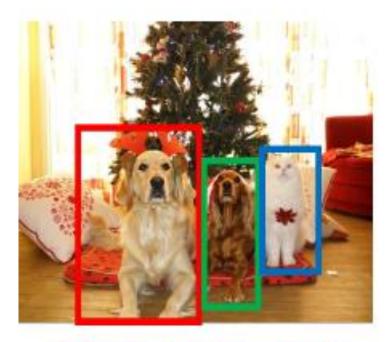


- Object Detection
- Semantic Instance Segmentation

Acknowledgement: Feifei Li et al's cs231n notes

Object Detection

- Problem setup
 - Input: image, object class(es)
 - Output: object instance bounding box locations + object scores



DOG, DOG, CAT

Detection: Why it is so hard?



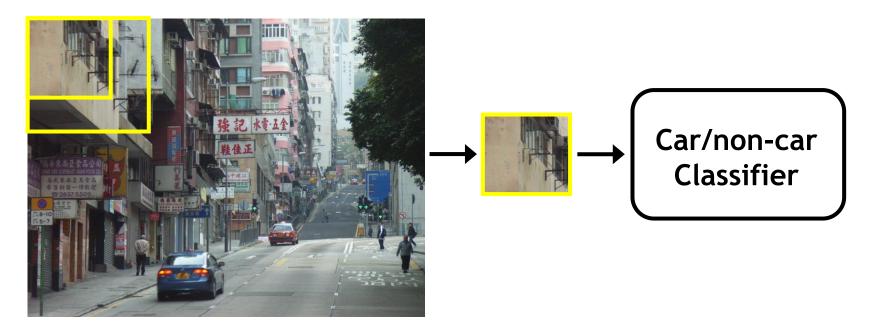




- Realistic scenes are crowded, cluttered, have overlapping objects.
- Object instances have large intra-class variance.

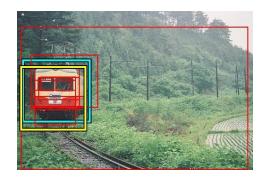
Object Detection

- Problem formulation
 - Object detection as a series of classification problems
 - Step1: Generate object candidate boxes
 - Step2: Scoring the object candidate with a classifier
- Example: sliding window method

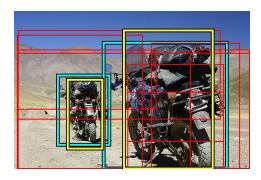


Object Detection

- CNN-based object detection
 - Reducing the search space by focusing on "object proposals" -image regions that are likely to contain objects







- ☐ Classifying each object proposals based on CNNs
- □ Refining the object location afterwards
- Typical method: Fast(er) R-CNN

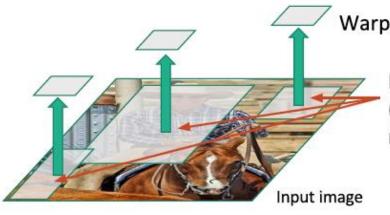


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.



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

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

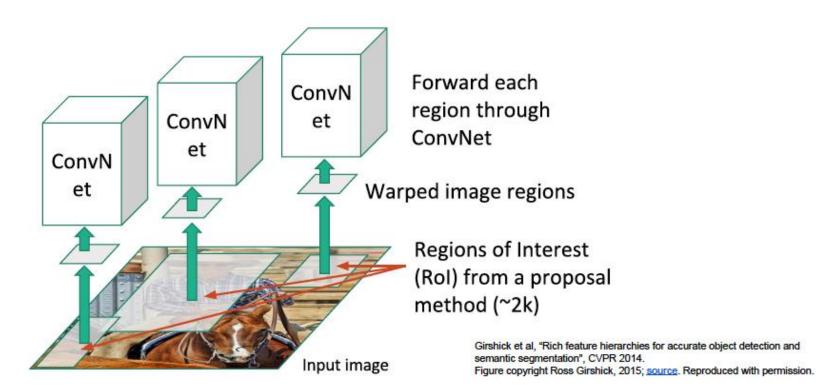


Warped image regions

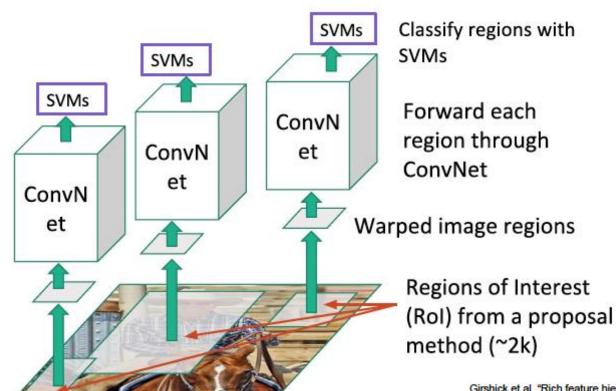
Regions of Interest (RoI) from a proposal method (~2k)

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









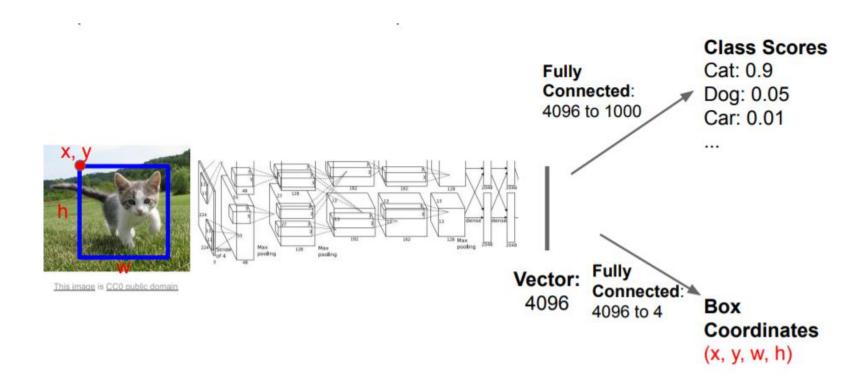
Input image

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

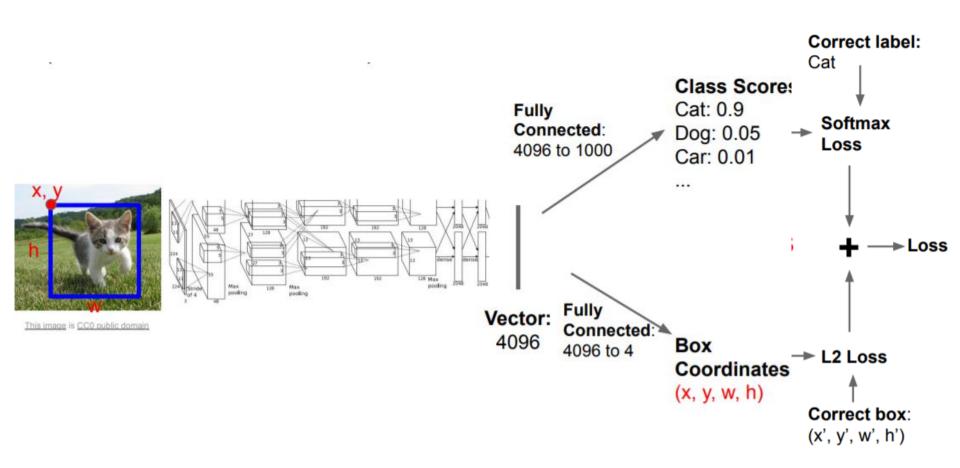
Linear Regression for bounding box offsets Classify regions with SVMs Bbox reg **SVMs SVMs** Bbox reg **SVMs** Bbox reg Forward each ConvN region through ConvN et ConvNet et ConvN Warped image regions et Regions of Interest (RoI) from a proposal method (~2k) Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. Input image



Each proposal



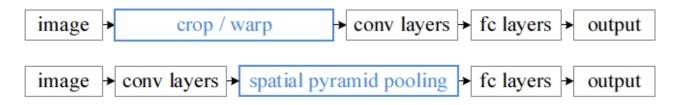
Multi-task loss



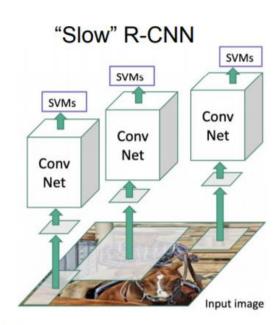


Problems

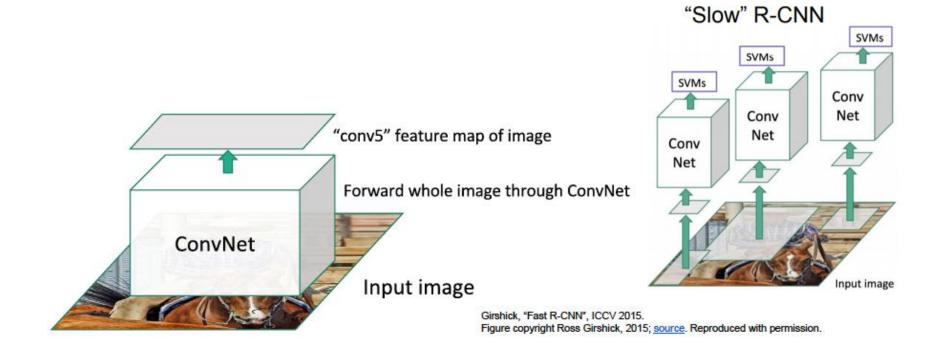
- □ Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Training post-hoc linear SVMs (hinge loss)
 - Training post-hoc bound-box regressions (least squares)
- □ Training is slow (84h), takes a lot of disk space
- □ Inference (detection) is slow
 - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
 - Fixed by SPP-net [He et al. ECCV14]

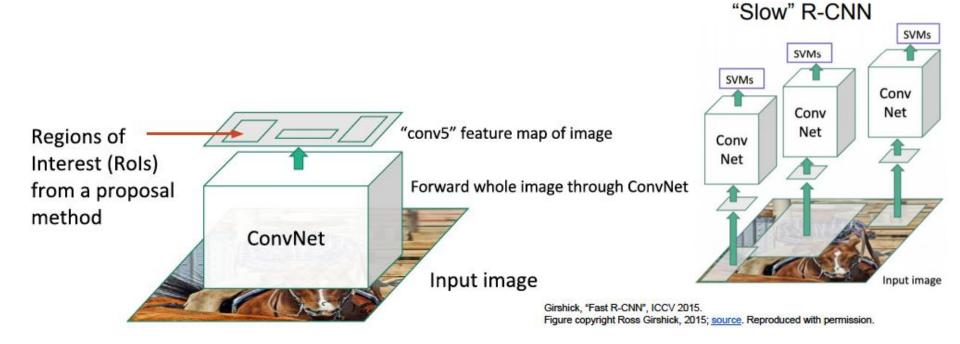


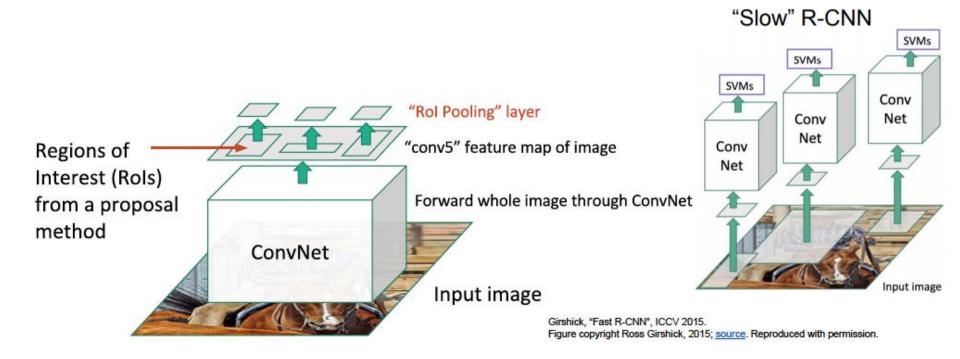


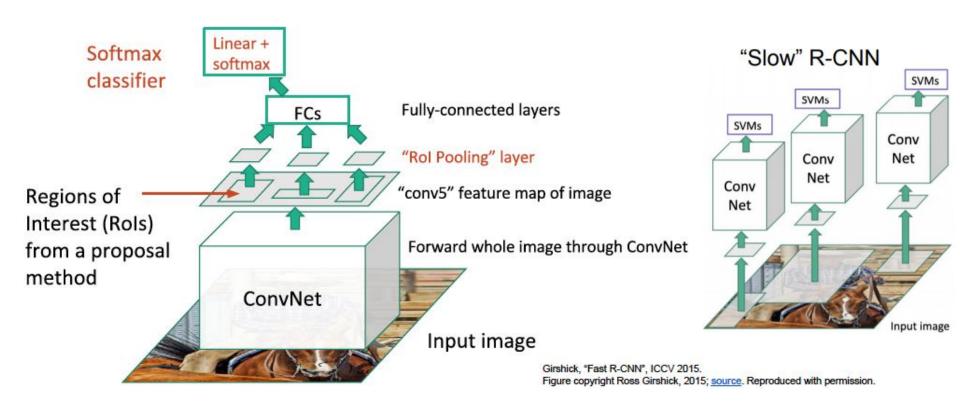


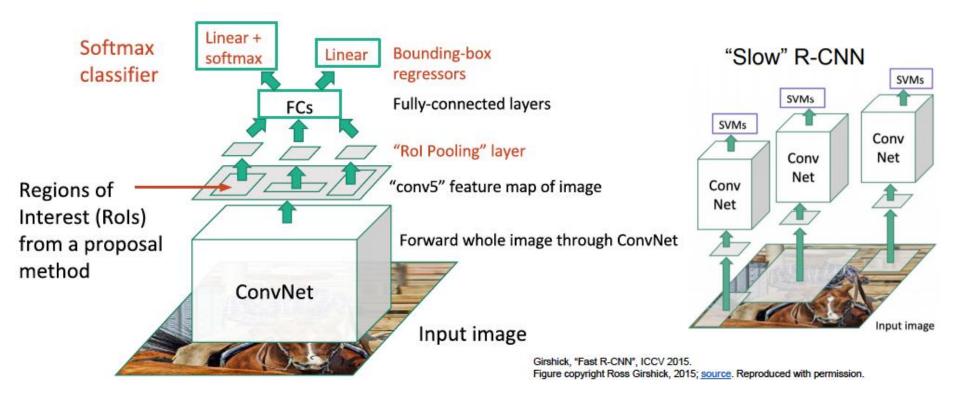
Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.



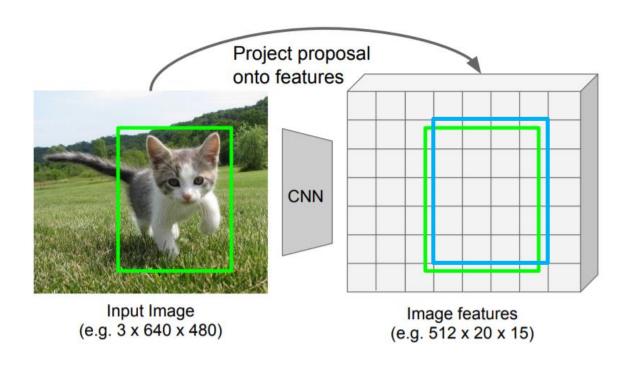






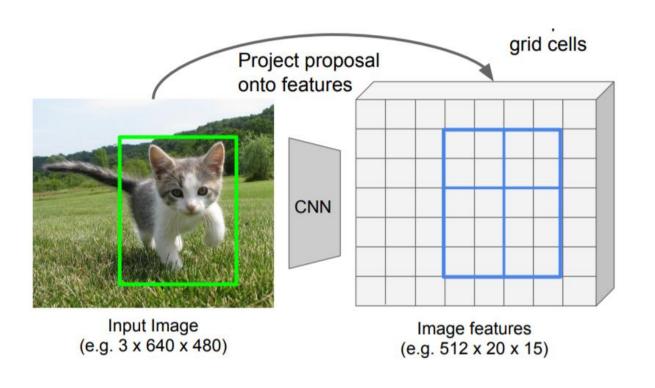


Rol Pooling



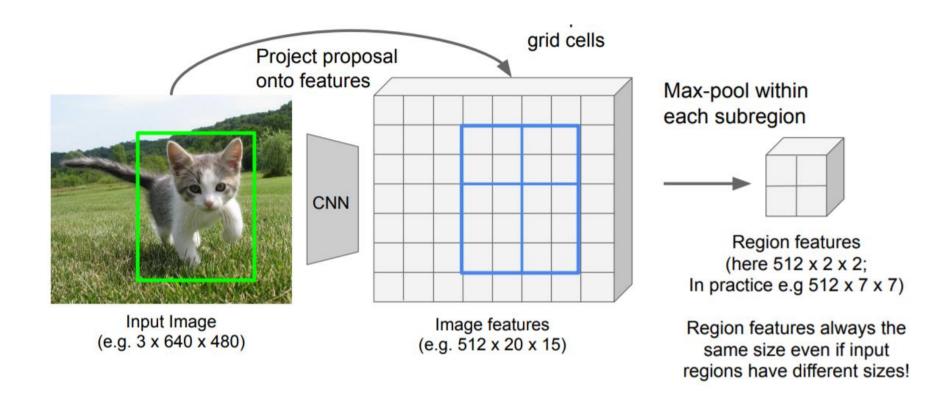
Snap to grid cells

Rol Pooling

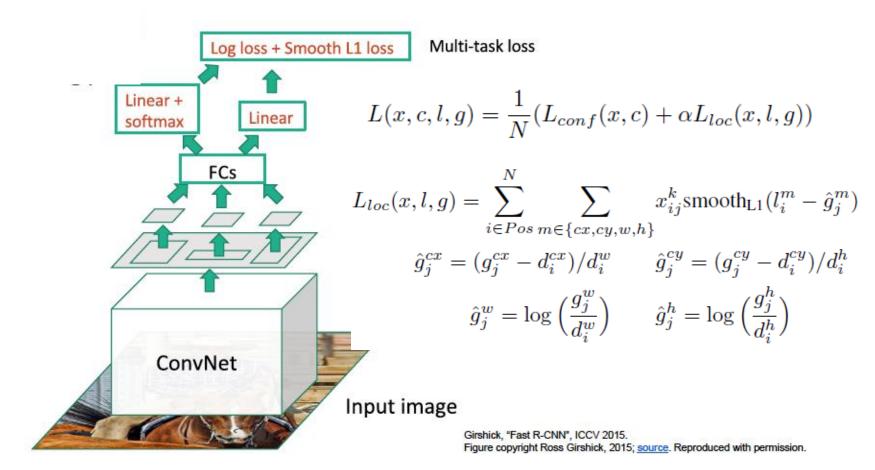


Divide into 2x2 grid of (roughly) equal subregions

Rol Pooling

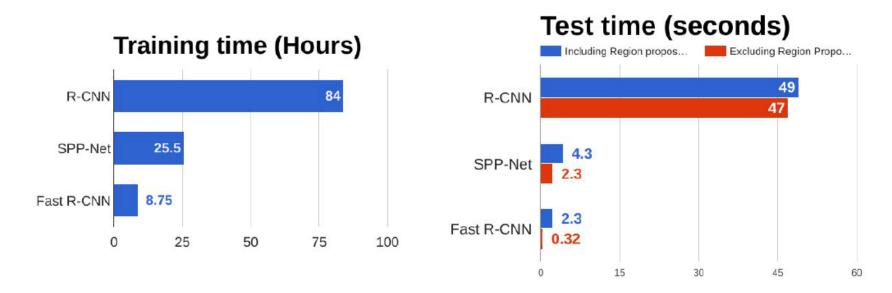


Deep network training



Speed Comparison

R-CNN vs SPP 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

loss

- Region proposal network
 - Make CNN do proposal generation

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

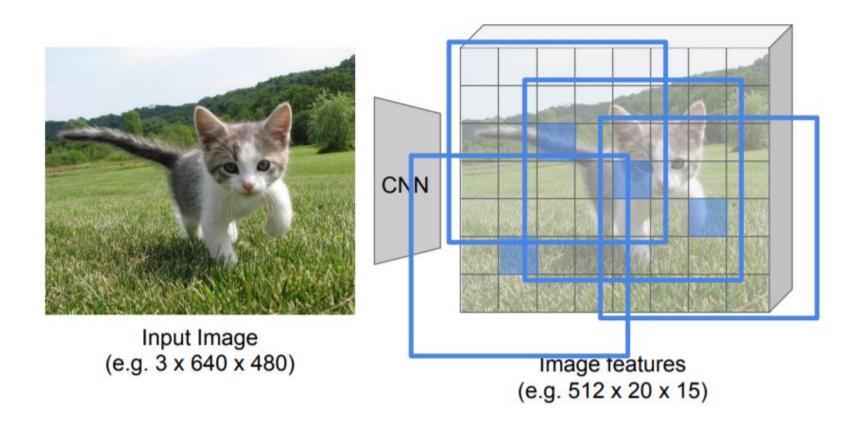
Jointly train with 4 losses:

- RPN classify object / not object
- RPN regress box coordinates
- Final classification score (object classes)
- Final box coordinates

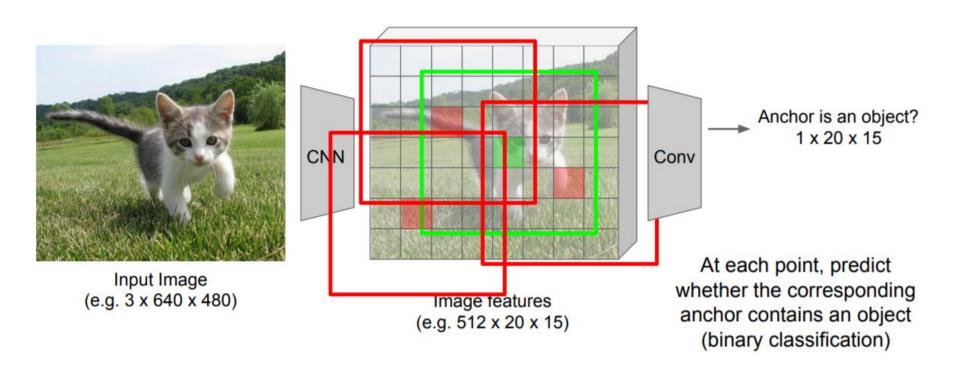
Classification Bounding-box regression loss Classification Bounding-box Rol pooling regression loss proposals Region Proposal Network feature map CNN

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

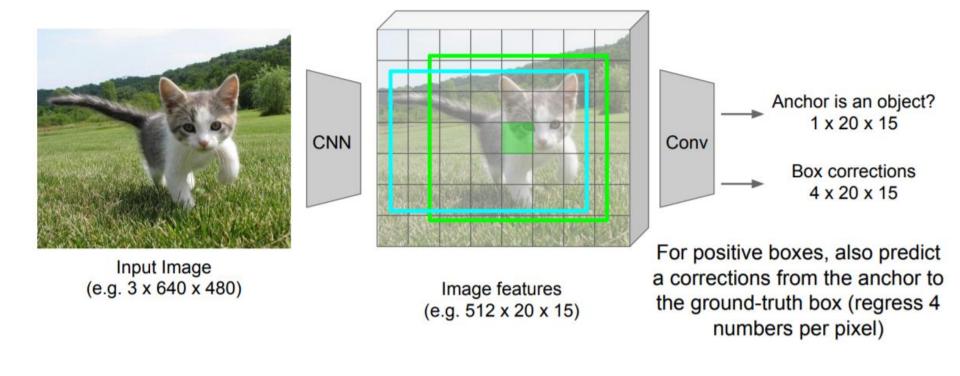
Imagine an anchor box of fixed size at each point in the feature map



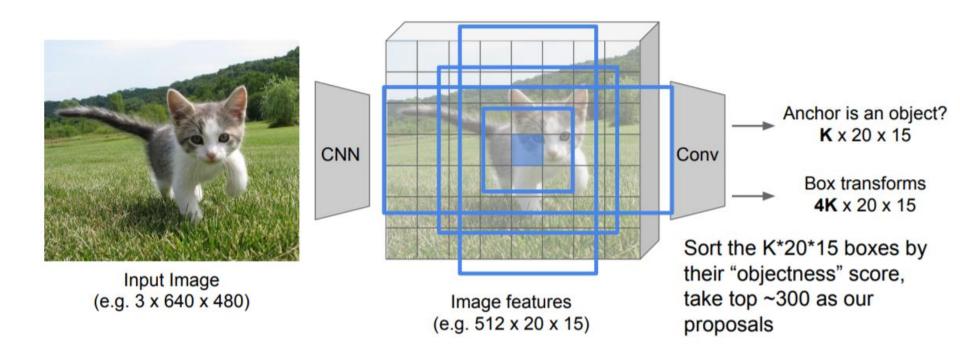
Imagine an anchor box of fixed size at each point in the feature map



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



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



Classification

loss

- Region proposal network
 - ☐ Make CNN do proposal generation

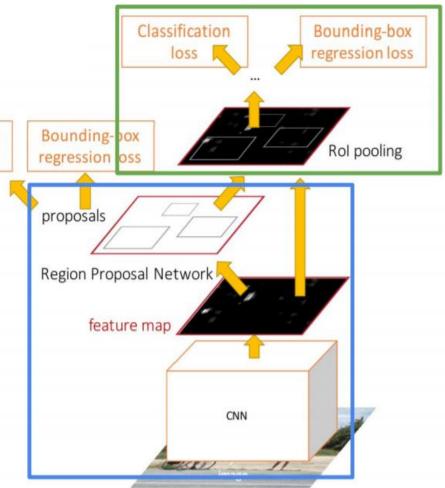
Faster R-CNN is a Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

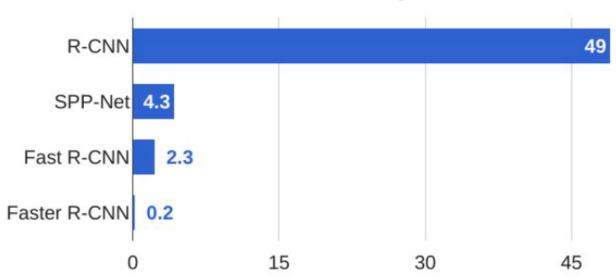
- Crop features: Rol pool / align
- Predict object class
- Prediction bbox offset





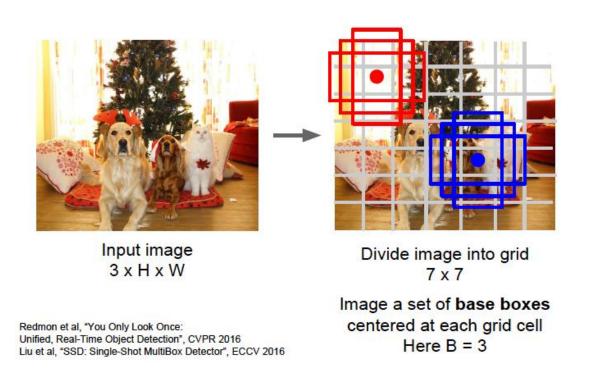
Speed comparison





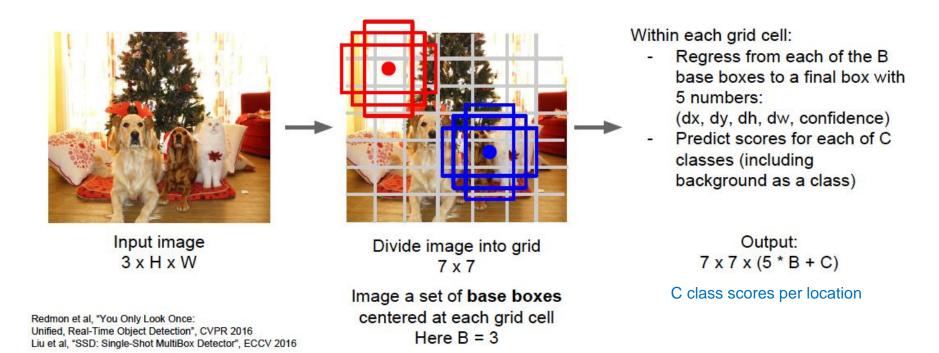
Object Detection without Proposals

- YOLO / SSD
- Alternative formulation: regression task



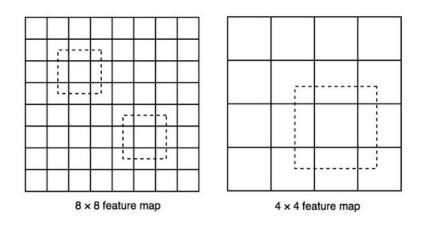
Object Detection without Proposals

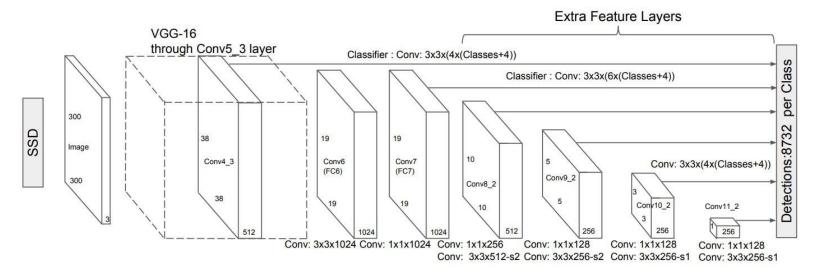
- YOLO / SSD
- Alternative formulation: regression task



Object Detection without Proposals

SSD: multi-scale feature maps





.

Object Detection Summary

Many configurations

Backbone
Network
VGG16
ResNet-101
Inception V2
Inception V3
Inception
ResNet
MobileNet

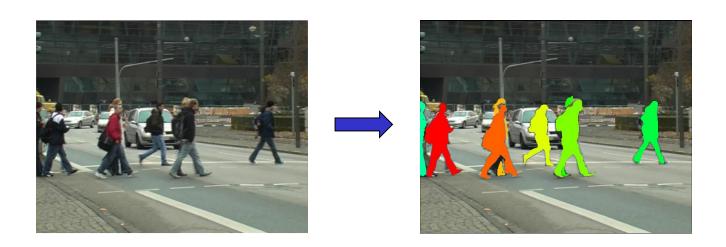
"Meta-Architecture"
Two-stage: Faster R-CNN
Single-stage: YOLO / SSD
Hybrid: R-FCN
Image Size # Region Proposals

Takeaways Faster R-CNN is slower but more accurate SSD is much faster but not as accurate Bigger / Deeper

backbones work better

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017 Zou et al, "Object Detection in 20 Years: A Survey", arXiv 2019

Outline



- Object Detection
- Semantic Instance Segmentation

Acknowledgement: Feifei Li et al's cs231n notes

Object Instance Segmentation

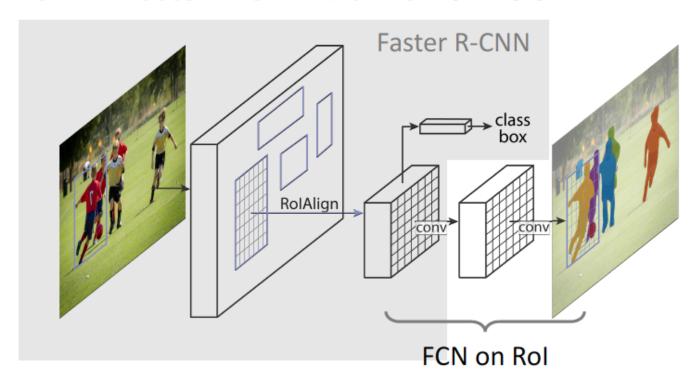
- Problem setup
 - Input: image, object class(es)
 - Output: object instance masks + object scores



DOG, DOG, CAT



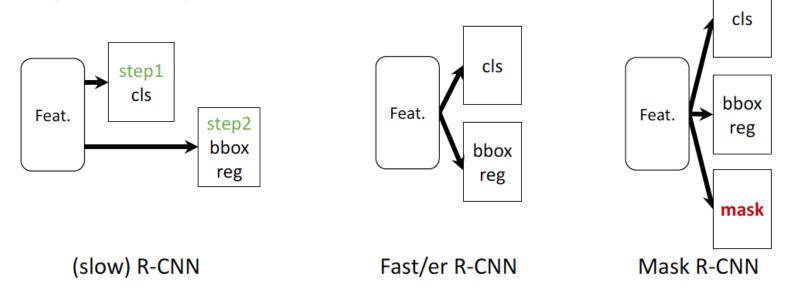
- Problem formulation
 - Mask R-CNN = Faster R-CNN with FCN on Rols





Parallel heads

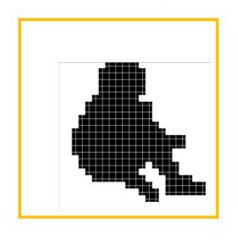
• Easy, fast to implement and train

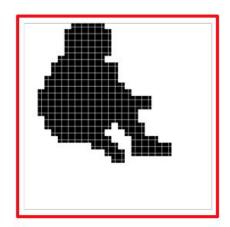


Fully-Convolution on Rol



target masks on Rols



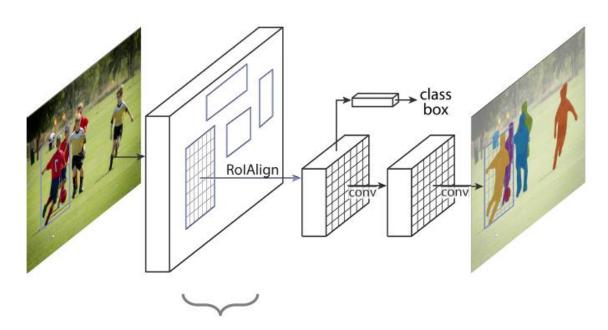


Translation of object in Rol => Same translation of mask in Rol

- Equivariant to small translation of Rols
- More robust to Rol's localization imperfection



Fully-Convolution on Rol

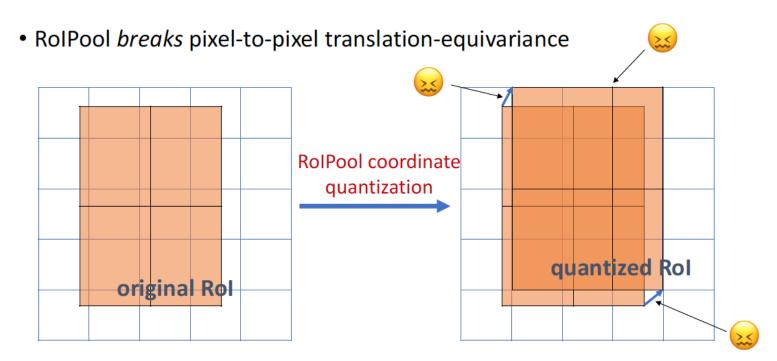


3. RolAlign:

3a. maintain translation-equivariance before/after Rol

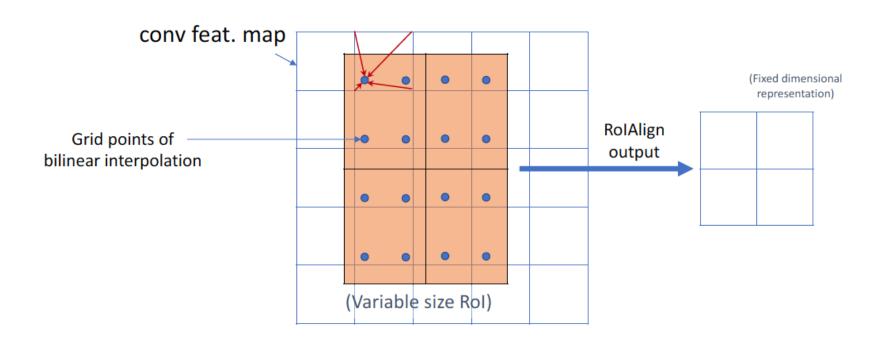


- Fully-Convolution on Rol
 - □ RolAlign vs RolPool



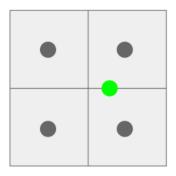


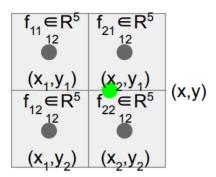
- Fully-Convolution on Rol
 - □ RolAlign





- Fully-Convolution on Rol
 - □ RolAlign
 - Sample at regular points in each subregion using bilinear interpolation



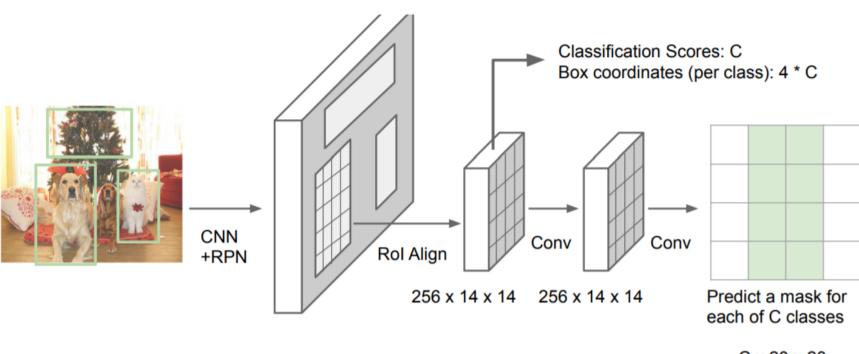


 \Box Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

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



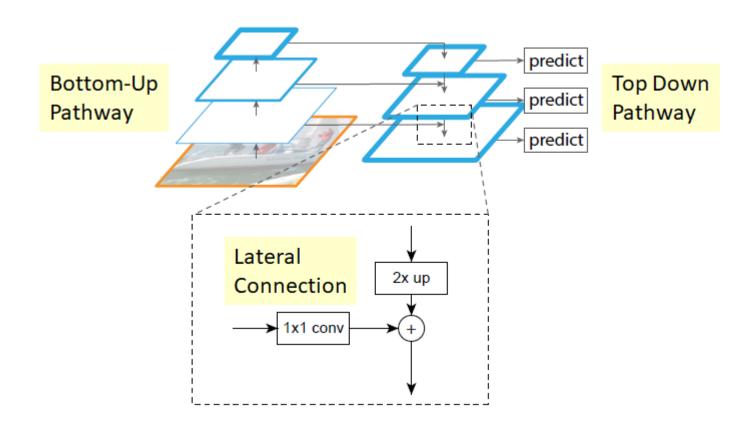
Overall architecture



C x 28 x 28

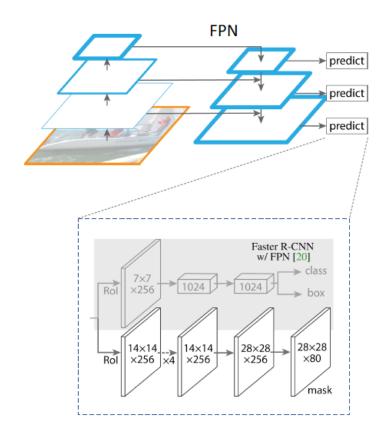


- Multiscale representation
 - □ Feature Pyramid Network (FPN)





- Multiscale representation
 - □ Feature Pyramid Network (FPN)



Results



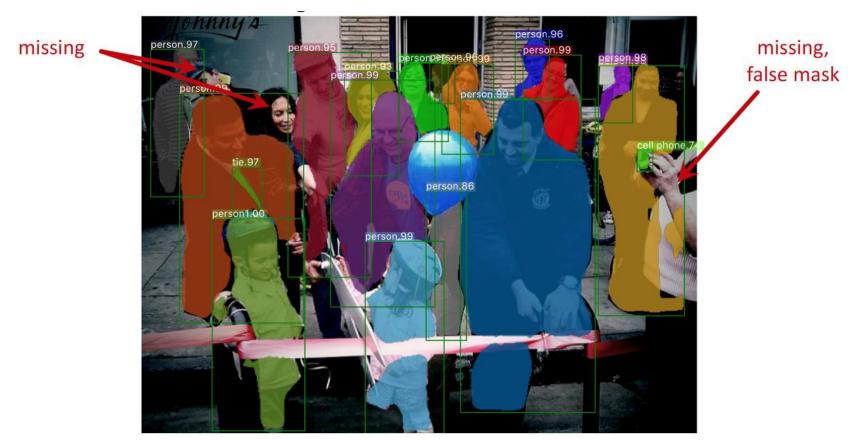
Mask R-CNN results on COCO

Results



Mask R-CNN results on COCO

Results



Mask R-CNN results on COCO

Pose results



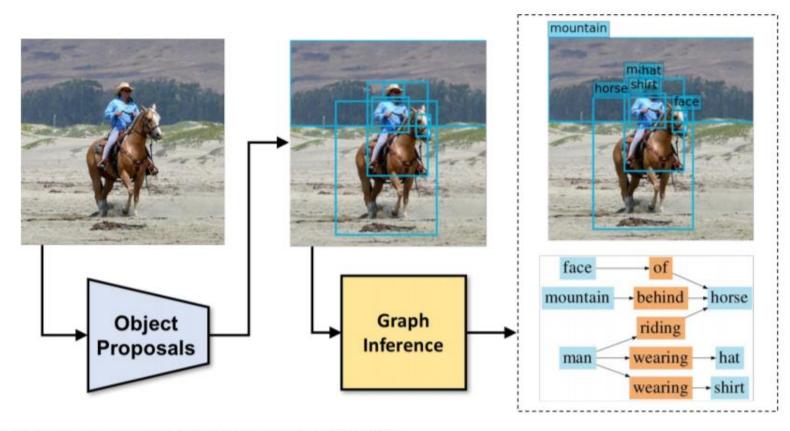


Open Source Implementations

- TensorFlow Detection API: https://github.com/tensorflow/models/tree/master/researc h/object_detection
 - □ Faster RCNN, SSD, RFCN, Mask R-CNN, ...
- Detectron2 (PyTorch) https://github.com/facebookresearch/detectron2
 - Mask R-CNN, RetinaNet, Faster R-CNN, RPN, Fast R-CNN, R-FCN, ...
- Finetune on your own dataset with pre-trained models

Other tasks

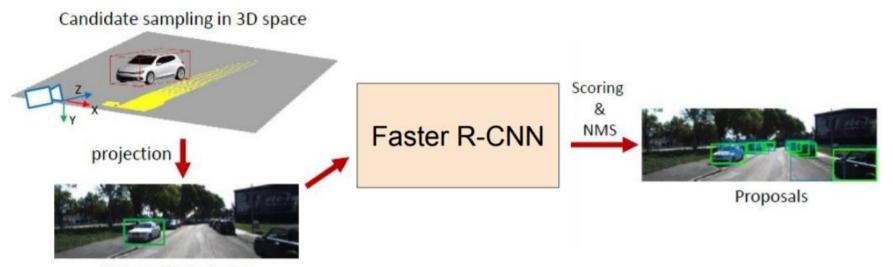
Objects + Relationships = Scene Graphs



Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017

Other tasks

3D Object Detection: Single view



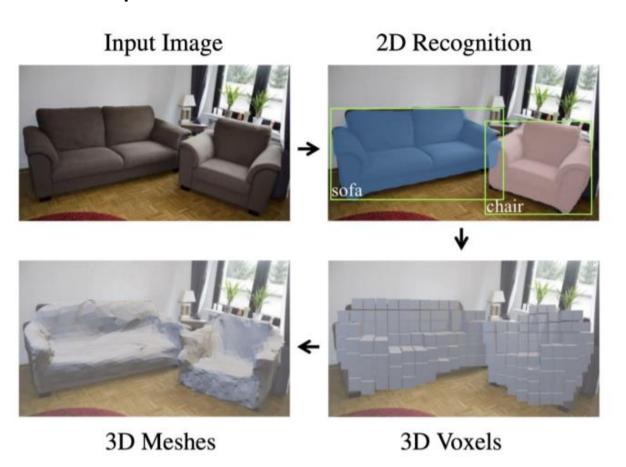
2D candidate boxes

- Same idea as Faster RCNN, but proposals are in 3D
- 3D bounding box proposal, regress 3D box parameters + class score

Chen, Xiaozhi, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. "Monocular 3d object detection for autonomous driving." CVPR 2016.

Other tasks

3D Shape Prediction: Mesh R-CNN



Gkioxari et al., Mesh RCNN, ICCV 2019



Summary

- CNNs in computer vision
 - Localization, Detection
 - Instance segmentation
- Next time:
 - □ Understanding CNNs
 - Limitations of CNNs