

image detection

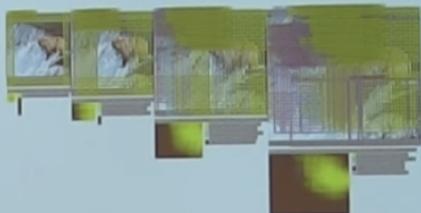
- task is also called image recognition
- Also see YOLO paper notes, R-CNN notes and follow up paper notes.

Sliding Window: Overfeat

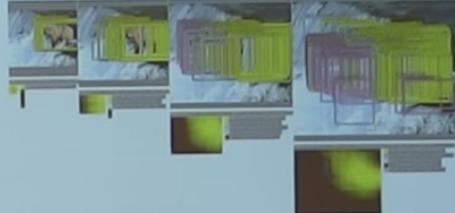
 Comment

In practice use many sliding window locations and multiple scales

Window positions + score maps



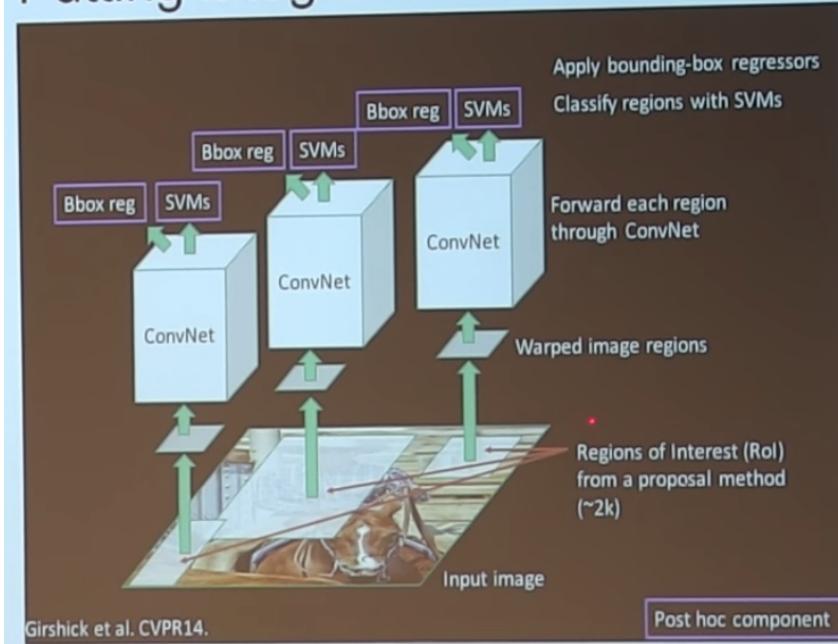
Box regression outputs



Final Predictions



Putting it together: R-CNN

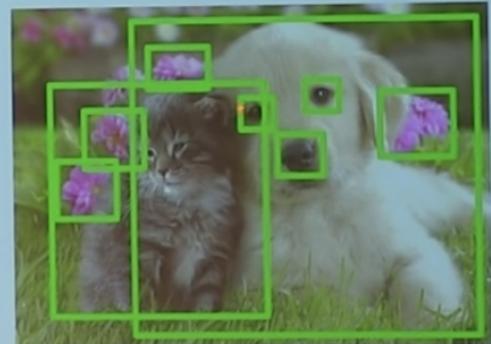


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

Slide credit: Ross Girshick

Region Proposals

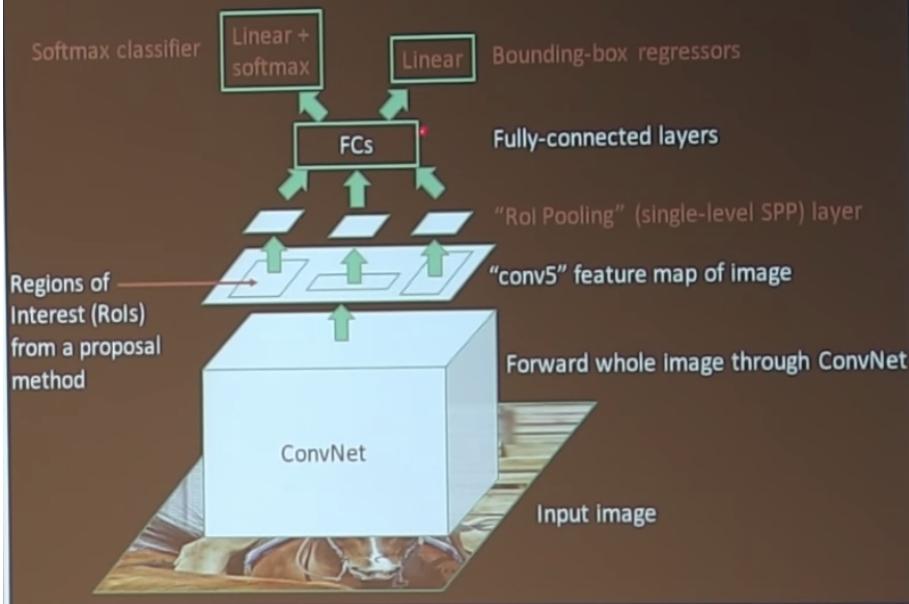
- Find “blobby” image regions that are likely to contain objects
- “Class-agnostic” object detector
- Look for “blob-like” regions



R-CNN Problems

1. Slow at test-time: need to run full forward pass of CNN for each region proposal
2. SVMs and regressors are post-hoc: CNN features not updated in response to SVMs and regressors
3. Complex multistage training pipeline

Fast R-CNN (test time)



Girshick, "Fast R-CNN", ICCV 2015

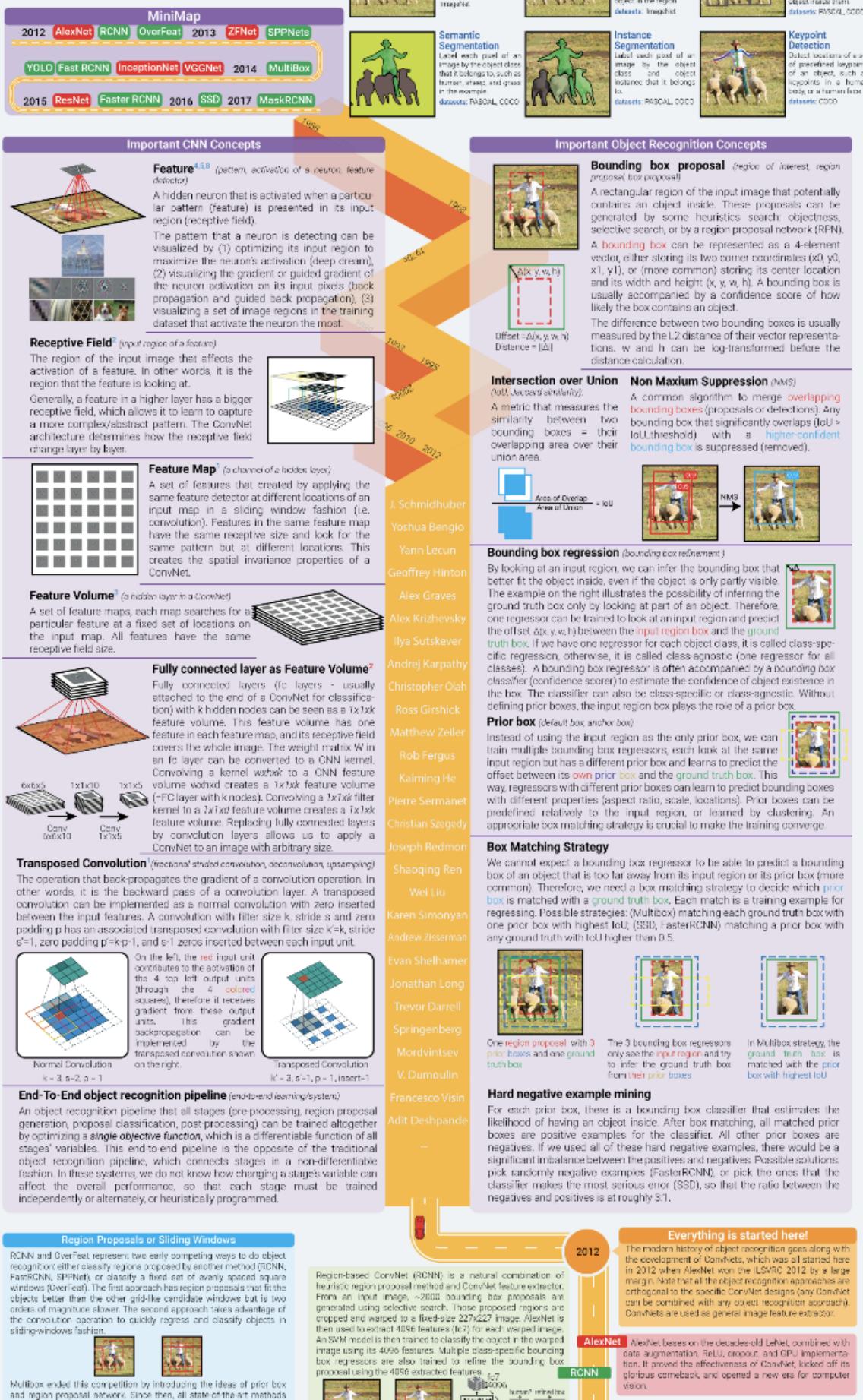
Slide credit: Ross Girshick

- mAP is a common metric for image detection task

Object Detection: Evaluation

- A more detailed explanation of concepts is [here](#) (pdf [here](#))

Modern History of Object Recognition Infographic



now has a set of prior boxes (generated based on a set of sliding windows or by clustering ground-truth boxes) from which bounding box regressors are trained to propose regions that better fit the object inside. The new competition is between the direct classification (YOLO, SSD) and refined classification approaches (FasterRCNN, MaskRCNN).

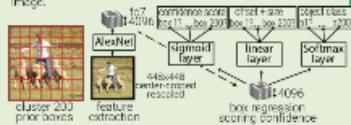
ZFNet is the LSVC 2013 winner, which is basically AlexNet with a minor modification: use 7x7 kernel instead of 11x11 kernel in the first Conv layer to retain more information.

SPPNet (Spatial Pyramid Pooling net) is essentially an enhanced version of RCNN by introducing two important concepts: adaptively-sized pooling (the SPP layer) and computing feature volume only once; in fact, the Fast-RCNN enhanced these ideas to faster RCNN with minor modifications.

SPPNet uses selective search to propose 2000 region proposals per image. It then extracts a common global feature volume from the entire image using ZFNet-Conv5. For each region proposal, SPPNet uses spatial pyramid pooling (SPP) to pool features in that region from the global feature volume to generate its fixed-length representation. This representation is used for training the object classifier and box regressors. Pooling features from a common global feature volume rather than pulling all image crops through a full CNN like RCNN brings two orders of magnitude speed up. Note that although spp position is differentiable, the authors did not do that, so the ZFNet was only trained on ImageNet without fine-tuning.

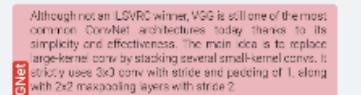
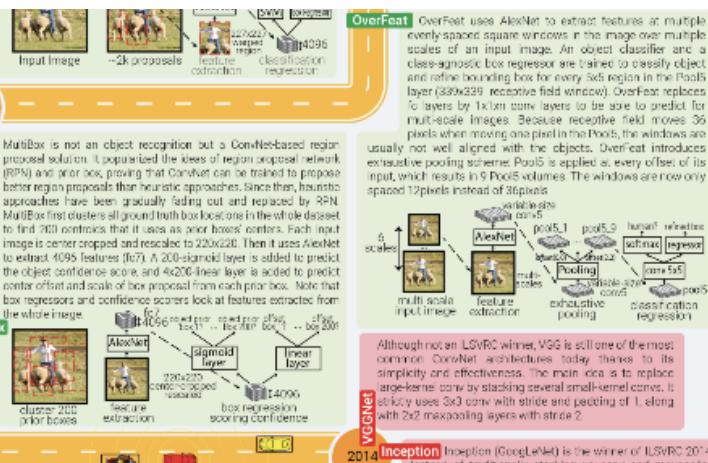


YOLO (You Only Look Once) is a direct development of MultiBox. It turns MultiBox from a region proposal solution to an object recognition method by adding a softmax layer parallel to the box regressor and box classifier layer, to directly predicts the object class. In addition, instead of clustering ground truth box locations to get the prior boxes, YOLO divides the input image into a 7x7 grid where each grid cell is a prior box. The grid cell is also used for box matching: if the center of an object falls into a grid cell, that grid cell is responsible for detecting that object. Like MultiBox, prior box only holds the center location information, not the size, so that box regressor predicts the box size independent with the size of the prior box. Like MultiBox, all the box regressor, conf/cls score, and object classifier look at features extracted from the whole image.

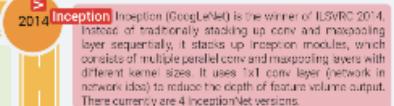


ResNet won the ILSVRC 2015 competition with an unbelievable 30% error rate (human performance is 5-10%), instead of transforming the input representation to output representation. ResNet sequentially stacks residual blocks, each computes the change (residual). It wants to make it as input, and add that to its input to produce its output representation. This is slightly related to boosting.

- 1 Dumoulin, Vincent, and Francesco Visin. "A guide to convolution arithmetic for deep learning." [Conv](#)
- 2 The Hien Bang Ha. "A guide to receptive field arithmetic for CNN" [ReceptiveField](#)
- 3 Karpathy, Andrej. "Cs231n: Convolutional neural networks for visual recognition" [DetailSummary](#)
- 4 Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." [DataProc](#)
- 5 Mordvintsev, Alexander, Christopher Olah, and Mike Tyka. "Inceptionism: Going deeper into neural networks" [DeepVision](#)
- 6 Adit Deshpande. "The 9 Deep Learning Papers You Need To Know About" [Summary](#)
- 7 Shelhamer, Evan, Jonathan Long, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." [LocNet](#)
- 8 Springerberg, Jost Tobias, et al. "Striving for simplicity: The all convolutional net." [QuickerBio-Proc](#)
- 9 Dhruv Pathak. "A Brief History of CNNs in Image Segmentation: From R-CNN to Mask R-CNN" [Summary](#)

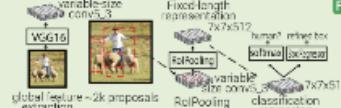


Although not an LSVC 2013 winner, VGG is still one of the most common ConvNet architectures today, thanks to its simplicity and effectiveness. The input idea is to replace layers conv by stacking several small-kernel convs. It also uses 3x3 conv with stride one and padding of 1, along with 2x2 maxpooling layers with stride 2.



Direct Classification or Refined Classification
These are the two competing approaches for now. Direct classification simultaneously regresses prior box and classifies object directly from the same input region, while the refined classification approach first regresses the prior box for a refined bounding box, and then pools the features of the refined box from a common feature volume and classifies object by those features. The former is faster but less accurate since the features it uses to classify are not extracted exactly from the refined prior box region.

Faster RCNN is essentially SPPNet with trainable feature extraction network and RoIPooling in replacement of the SPP layer. RoIPooling (region of interest pooling) is simply a special case of SPP where here only one pyramid level is used. RoIPooling generates a fixed 7x7 feature volume for each ROI (region proposal) by dividing the ROI feature volume into a 7x7 grid of sub-windows and then maxpooling the values from each sub-window.



Faster RCNN is Faster RCNN with heuristic region proposal replaced by region proposal network (RPN) inspired by MultiBox. In Faster RCNN RPN is a small ConvNet (3x3 conv + 1x1 conv + 1x1 conv) looking at the conv5_3 global feature volume in the sliding window fashion. Each sliding window has 9 prior boxes that relate to its receptive field (3 scales x 3 aspect ratios). RPN does bounding box regression and box confidence scoring for each prior box. The whole pipeline is trainable by combining the loss of box regression, box confidence scoring, and object classification into one common global objective function. Note that here RPN only looks at a small input region and prior boxes hold both the prior location and the box size, which are different from the MultiBox one RPN does.



References

- 1 Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." [AlexNet](#)
- 2 Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." [ZFNet](#)
- 3 Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition" [VGGNet](#)
- 4 Szegedy, Christian, et al. "Going deeper with convolutions" [Inception](#)
- 5 He, Kaiming, et al. "Deep residual learning for image recognition." [ResNet](#)
- 6 Girshick, Ross, et al. "Rich feature hierarchies for accurate object detection and semantic segmentation." [RCNN](#)
- 7 Sermanet, Pierre, et al. "Overfeat: Integrated recognition, localization and detection using convolutional networks." [OverFeat](#)
- 8 He, Kaiming, et al. "Spatial pyramid pooling in deep convolutional networks for visual recognition." [SPPNet](#)
- 9 Szepegy, Christian, et al. "Scalable, high-quality object detection." [MultiBox](#)
- 10 Girshick, Ross. "Fast-rcnn." [FasterRCNN](#)
- 11 Redmon, Joseph, et al. "You only look once: Unified real-time object detection." [YOLO](#)
- 12 Ren, Shaoting, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." [FasterRCNN](#)
- 13 Liu, Wei, et al. "SSD: Single shot multibox detector." [SSD](#)
- 14 He, Kaiming, et al. "Mask R-CNN." [MaskRCNN](#)

Nikasa

<https://nikasa1889.github.io/>
<https://medium.com/@nikasa1889/>



Mask RCNN extends Faster RCNN for instance segmentation by adding a branch for predicting class-specific object mask in parallel with each prior box instead of just scoring the object confidence (similar to YOLO). It improves the diversity of prior boxes realizations by running the RPN on multiple conv layers at different depth levels.