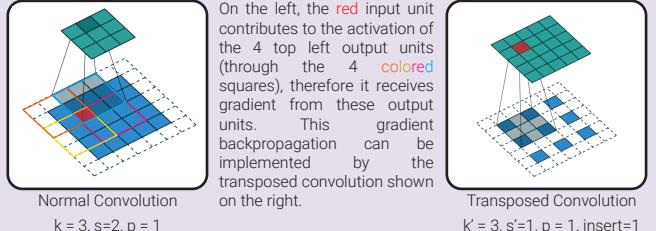
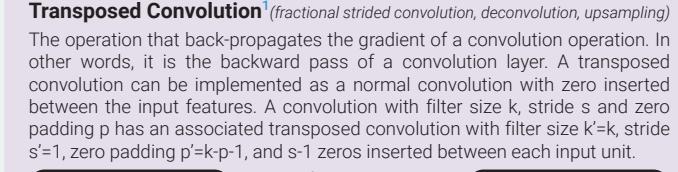
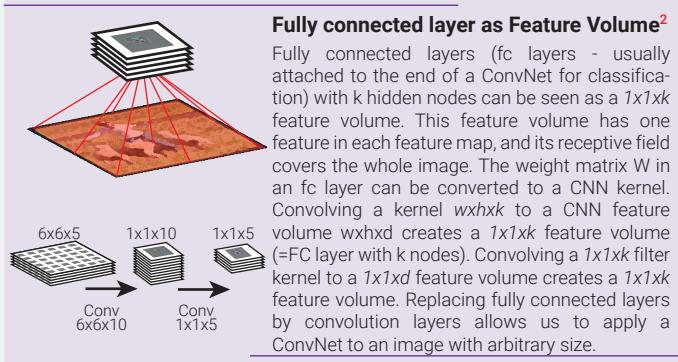
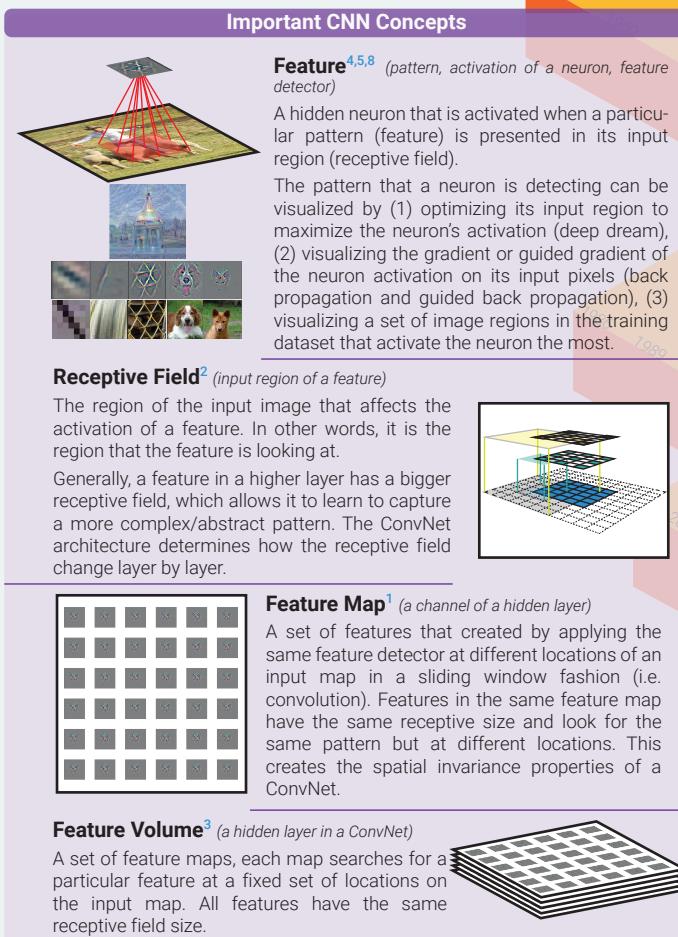
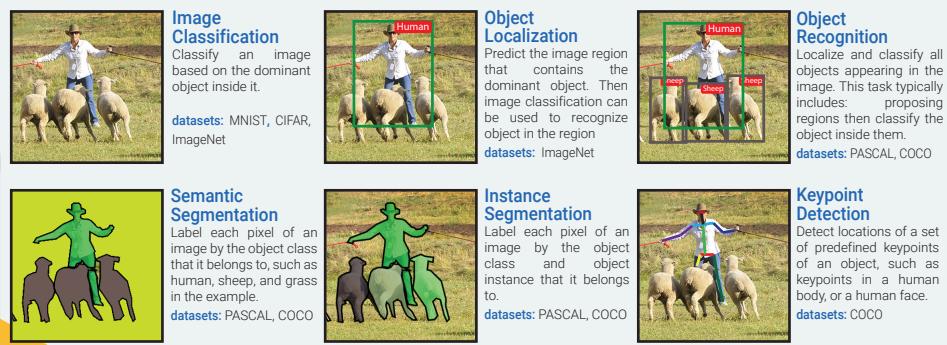
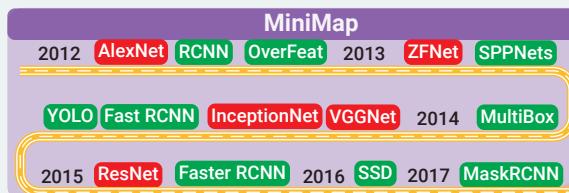
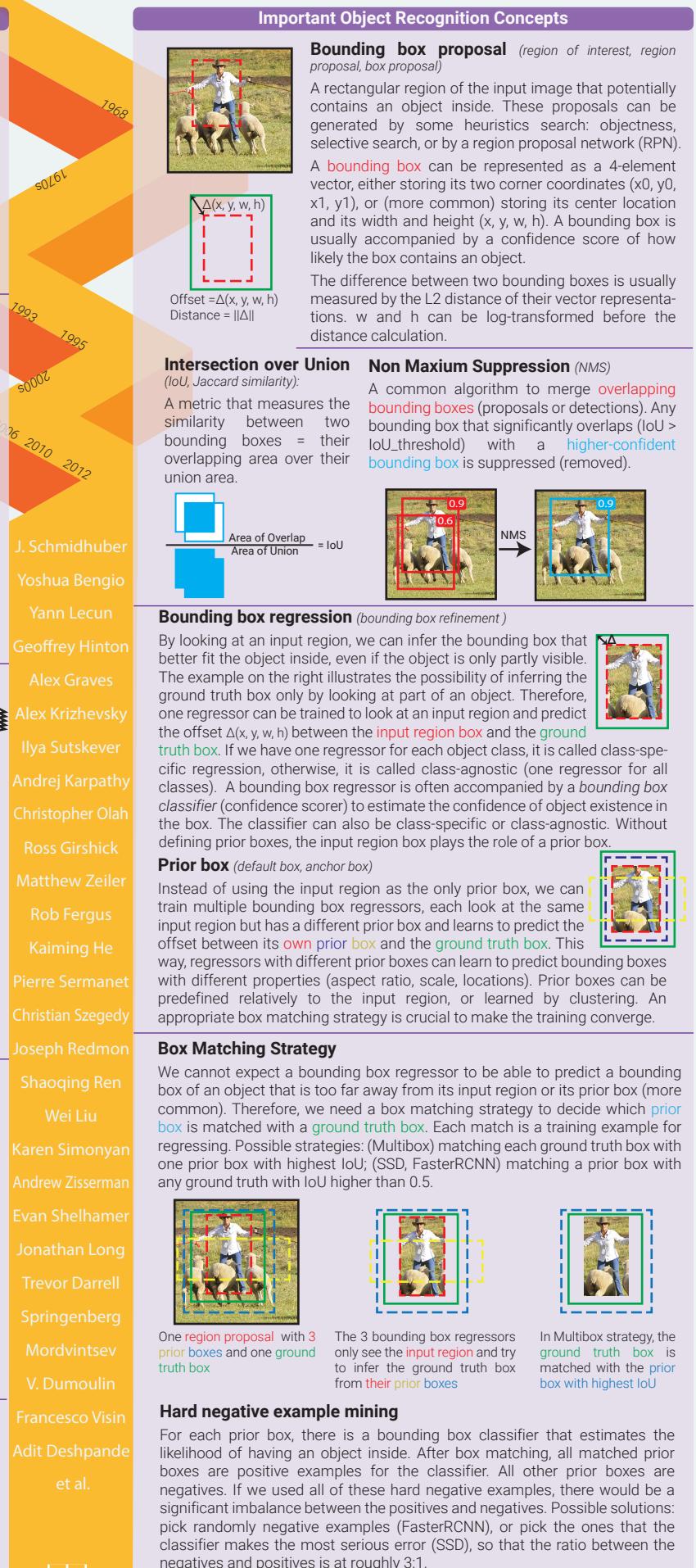


# Modern History of Object Recognition Infographics



### End-To-End object recognition pipeline (end-to-end learning/system)

An object recognition pipeline that all stages (pre-processing, region proposal generation, proposal classification, post-processing) can be trained altogether by optimizing a **single objective function**, which is a differentiable function of all stages' variables. This end-to-end pipeline is the opposite of the traditional object recognition pipeline, which connects stages in a non-differentiable fashion. In these systems, we do not know how changing a stage's variable can affect the overall performance, so that each stage must be trained independently or alternately, or heuristically programmed.



## Region Proposals or Sliding Windows

RCNN and OverFeat represent two early competing ways to do object recognition: either classify regions proposed by another method (RCNN, FastRCNN, SPPNet), or classify a fixed set of evenly spaced square windows (OverFeat). The first approach has region proposals that fit the objects better than the other grid-like candidate windows but is two orders of magnitude slower. The second approach takes advantage of the convolution operation to quickly regress and classify objects in sliding-windows fashion.

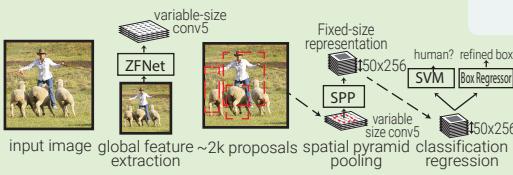


Multibox ended this competition by introducing the ideas of prior box and region proposal network. Since then, all state-of-the-art methods now have a set of prior boxes (generated based on a set of sliding windows or by clustering ground-truth boxes) from which bounding box regressors are trained to propose regions that better fit the object inside. The new competition is between the *direct classification* (YOLO, SSD) and *refined classification* approaches (FasterRCNN, MaskRCNN).

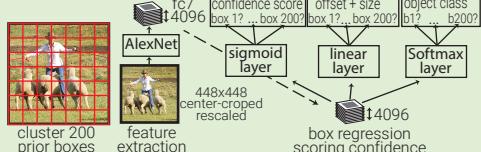
ZFNet is the ILSVRC 2013 winner, which is basically AlexNet with a minor modification: use 7x7 kernel instead of 11x11 kernel in the first Conv layer to retain more information.

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 embraced these ideas to fasten 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 pushing all image crops through a full CNN like RCNN brings two orders of magnitude speed up. Note that although SPP operation is differentiable, the authors did not do that, so the ZFNet was only trained on ImageNet without finetuning.



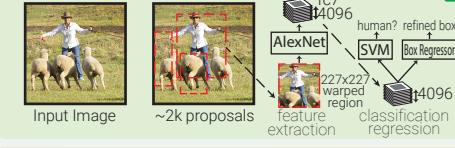
YOLO (You Only Look Once) is a direct development of MultiBox. It turns MultiBox from a region proposal solution to an object recognition method by adding a softmax layer, parallel to the box regressor and box classifier layer, to directly predicts the object class. In addition, instead of clustering 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, confidence scorer, and object classifier look at features extracted from the whole image.



ResNet won the ILSVRC 2015 competition with an unbelievable 3.6% 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 to its input, and add that to its input to produce its output representation. This is slightly related to boosting.



Region-based ConvNet (RCNN) is a natural combination of heuristic region proposal method and ConvNet feature extractor. From an input image, ~2000 bounding box proposals are generated using selective search. Those proposed regions are cropped and warped to a fixed-size 227x227 image. AlexNet is then used to extract 4096 features (fc7) for each warped image. An SVM model is then trained to classify the object in the warped image using its 4096 features. Multiple class-specific bounding box regressors are also trained to refine the bounding box proposal using the 4096 extracted features.



AlexNet

AlexNet bases on the decades-old LeNet, combined with data augmentation, ReLU, dropout, and GPU implementation. It proved the effectiveness of ConvNet, kicked off its glorious comeback, and opened a new era for computer vision.

RCNN

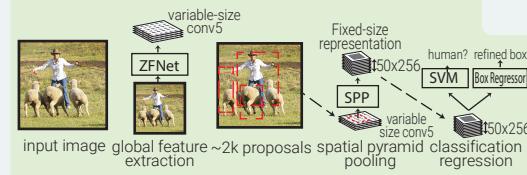
OverFeat

**Everything is started here!**

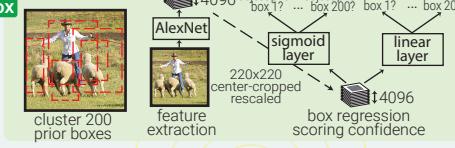
The modern history of object recognition goes along with the development of ConvNets, which was all started here in 2012 when AlexNet won the ILSVRC 2012 by a large margin. Note that all the object recognition approaches are orthogonal to the specific ConvNet designs (any ConvNet can be combined with any object recognition approach). ConvNets are used as general image feature extractor.

Multibox ended this competition by introducing the ideas of prior box and region proposal network. Since then, all state-of-the-art methods now have 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).

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 pushing all image crops through a full CNN like RCNN brings two orders of magnitude speed up. Note that although SPP operation is differentiable, the authors did not do that, so the ZFNet was only trained on ImageNet without finetuning.



Multibox is not an object recognition but a ConvNet-based region proposal solution. It popularized the ideas of region proposal network (RPN) and prior box, proving that ConvNet can be trained to propose better region proposals than heuristic approaches. Since then, heuristic approaches have been gradually fading out and replaced by RPN. Multibox first clusters all ground truth box locations in the whole dataset to find 200 centroids that it uses as prior boxes' centers. Each input image is center cropped and rescaled to 220x220. Then it uses AlexNet to extract 4096 features (fc7). A 200-sigmoid layer is added to predict the object confidence score, and 4x200-linear layer is added to predict center offset and scale of box proposal from each prior box. Note that box regressors and confidence scorers look at features extracted from the whole image.



Multibox

VGGNet

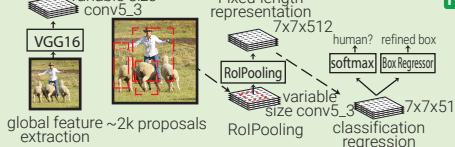
Although not an ILSVRC winner, VGG is still one of the most common ConvNet architectures today thanks to its simplicity and effectiveness. The main idea is to replace large-kernel conv by stacking several small-kernel convs. It strictly uses 3x3 conv with stride and padding of 1, along with 2x2 maxpooling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

Inception (GoogLeNet) is the winner of ILSVRC 2014. Instead of traditionally stacking up conv and maxpooling layer sequentially, it stacks up Inception modules, which consists of multiple parallel conv and maxpooling layers with different kernel sizes. It uses 1x1 conv layer (network in network idea) to reduce the depth of feature volume output. There currently are 4 InceptionNet versions.

**Direct Classification or Refined Classification**

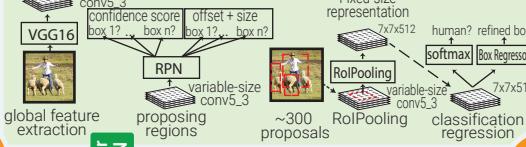
These are the two competing approaches for now. Direct classification simultaneously regresses prior box and classifies object directly from the same input region, while the refined classification approach first regresses the prior box for a refined bounding box, and then pools the features of the refined box from a common feature volume and classifies object by these 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.

Fast 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 max-pooling the values from each sub-window.



FastRCNN

Faster RCNN is Fast 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 relative 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 center location and the box size, which are different from the MultiBox and YOLO design.



References

- 1 Dumoulin, Vincent, and Francesco Visin. "A guide to convolution arithmetic for deep learning." [Conv](#)
- 2 The-Hien Dang-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." [ZFNet](#)
- 5 Mordvintsev, Alexander, Christopher Olah, and Mike Tyka. "Inceptionism: Going deeper into neural networks." [DeepDream](#)
- 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." [Deconv](#)
- 8 Springenberg, Jost Tobias, et al. "Striving for simplicity: The all convolutional net." [GuidedBackProp](#)
- 9 Dhruv Parthasarathy "A Brief History of CNNs in Image Segmentation: From R-CNN to Mask R-CNN" [Summary](#)
- 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](#)
- 1 Girshick, Ross, et al. "Rich feature hierarchies for accurate object detection and semantic segmentation." [RCNN](#)
- 2 Sermanet, Pierre, et al. "Overfeat: Integrated recognition, localization and detection using convolutional networks." [OverFeat](#)
- 3 He, Kaiming, et al. "Spatial pyramid pooling in deep convolutional networks for visual recognition." [SPPNet](#)
- 4 Szegedy, Christian, et al. "Scalable, high-quality object detection." [MultiBox](#)
- 5 Girshick, Ross. "Fast r-cnn." [FastRCNN](#)
- 6 Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." [YOLO](#)
- 7 Ren, Shaoqing, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." [FasterRCNN](#)
- 8 Liu, Wei, et al. "SSD: Single shot multibox detector." [SSD](#)
- 9 He, Kaiming, et al. "Mask R-CNN." [MaskRCNN](#)