

# CNN applications

Samuel Cheng

(Slide credits: Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung)

School of ECE  
University of Oklahoma

Spring, 2017

- We will look into several applications of CNNs besides image recognition
  - Semantic segmentation
  - Object localization
  - Object detection
  - Instance segmentation

# So far: Image Classification



This image is [CC0 public domain](#)

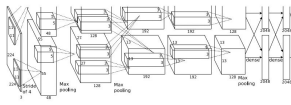


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

**Vector:**  
4096

**Fully-Connected:**  
4096 to 1000

**Class Scores**

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

## Other Computer Vision Tasks

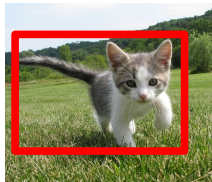
### Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

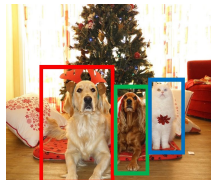
### Classification + Localization



CAT

Single Object

### Object Detection



DOG, DOG, CAT

Multiple Object

### Instance Segmentation

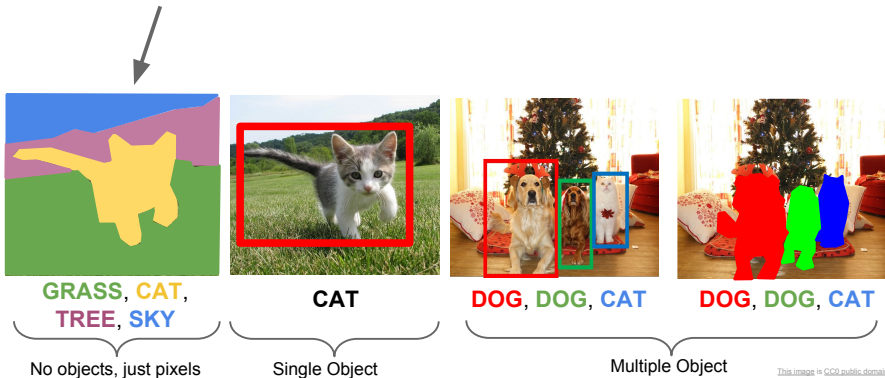


DOG, DOG, CAT

This image is CC0 public domain



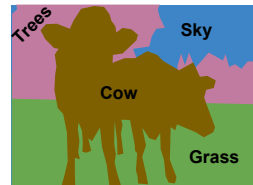
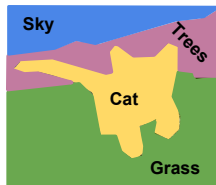
# Semantic Segmentation



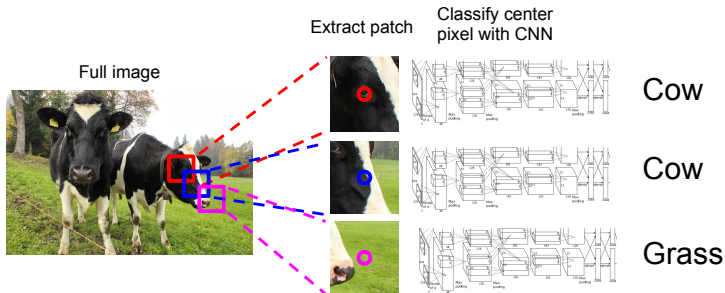
# Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



# Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

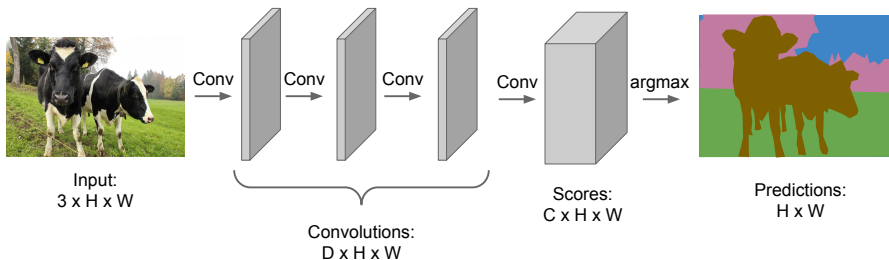
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 20 May 10, 2017



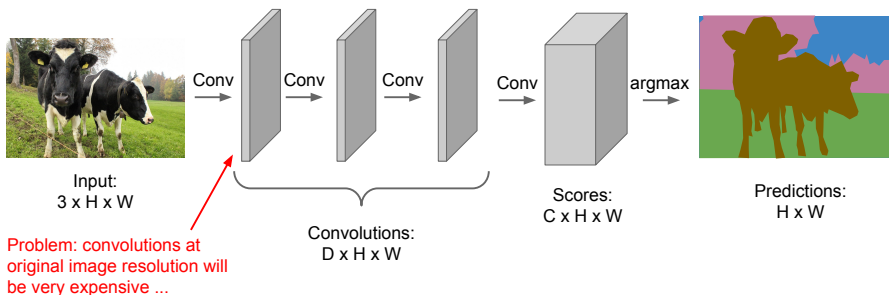
# Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



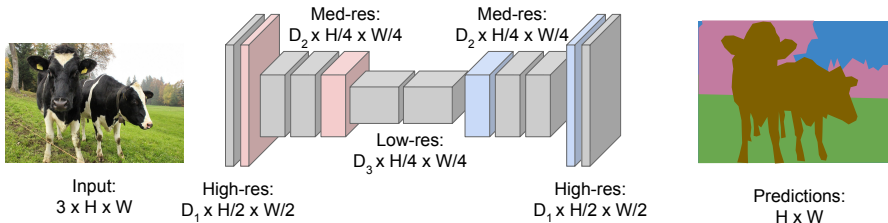
# Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



# Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 24 May 10, 2017

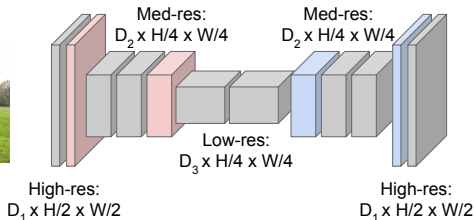
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided  
convolution



Input:  
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



**Upsampling:**  
???



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 25 May 10, 2017



# In-Network upsampling: “Unpooling”

**Nearest Neighbor**

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

**“Bed of Nails”**

1	2
3	4

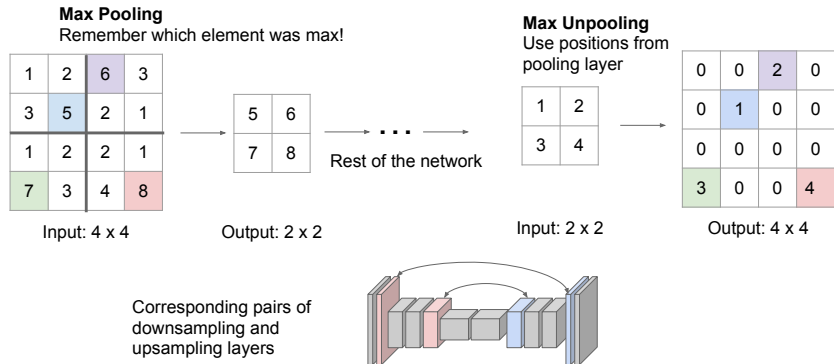


1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Input: 2 x 2

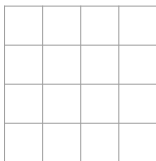
Output: 4 x 4

# In-Network upsampling: “Max Unpooling”

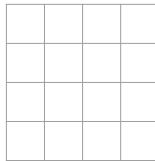


# Learnable Upsampling: Transpose Convolution

**Recall:** Typical 3 x 3 convolution, stride 1 pad 1



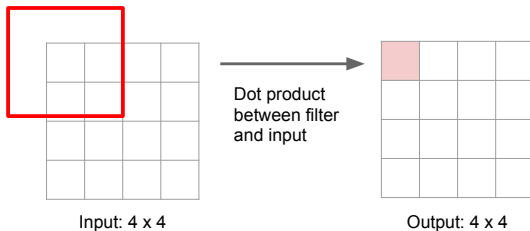
Input: 4 x 4



Output: 4 x 4

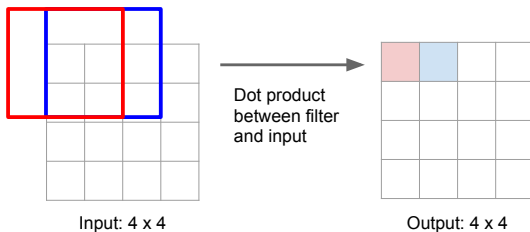
# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



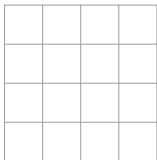
# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



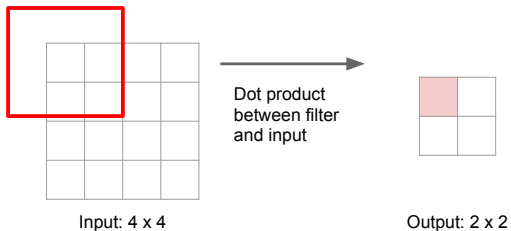
Input: 4 x 4



Output: 2 x 2

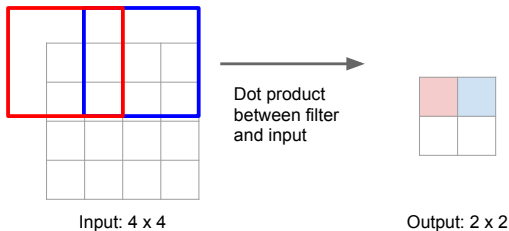
# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

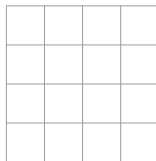


# Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1



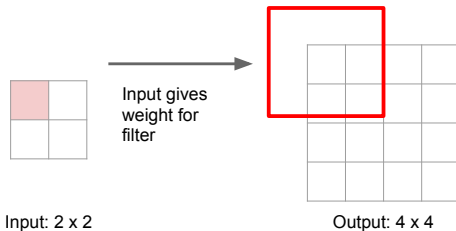
Input: 2 x 2



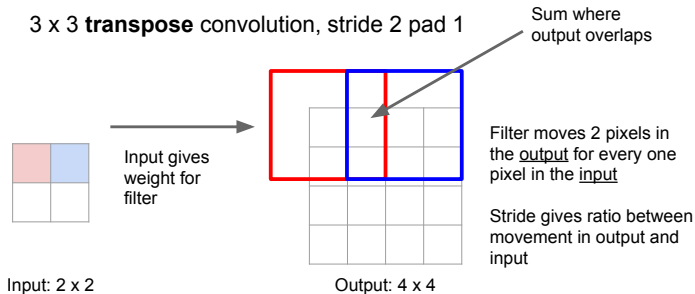
Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

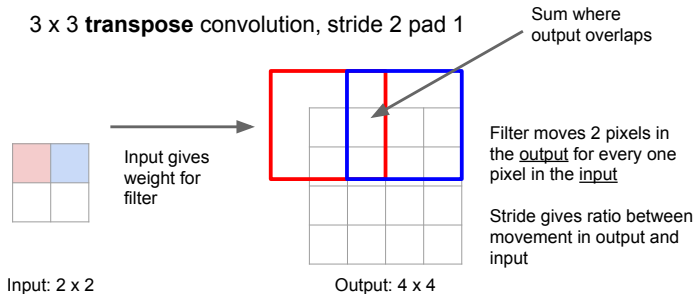
3 x 3 **transpose** convolution, stride 2 pad 1



# Learnable Upsampling: Transpose Convolution



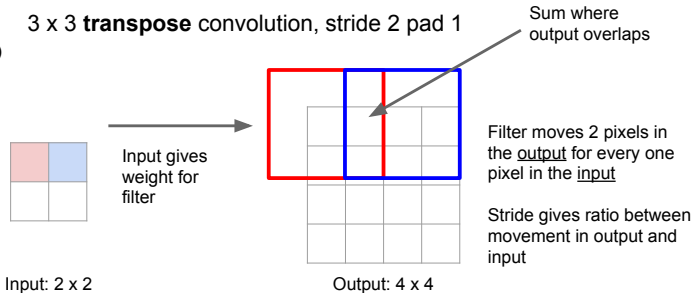
# Learnable Upsampling: Transpose Convolution



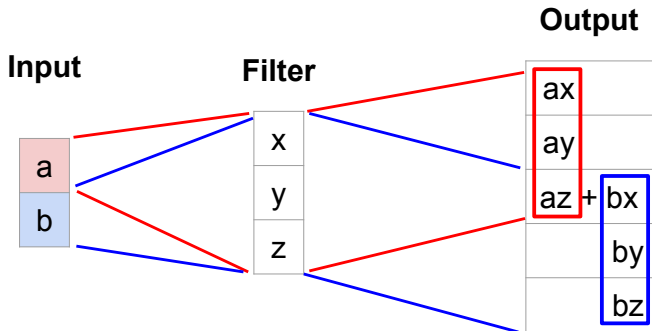
# Learnable Upsampling: Transpose Convolution

## Other names:

- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution



# Transpose Convolution: 1D Example



Output contains copies of the filter weighted by the input, summing at where it overlaps in the output

Need to crop one pixel from output to make output exactly 2x input

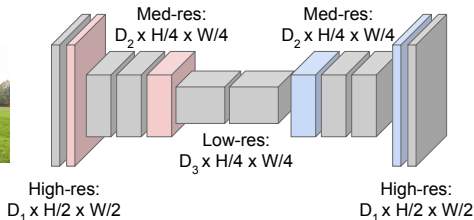
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided  
convolution



Input:  
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



**Upsampling:**  
Unpooling or strided  
transpose convolution



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

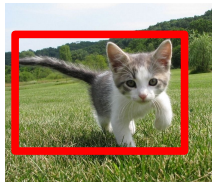
Lecture 11 - 44 May 10, 2017

# Classification + Localization



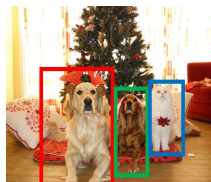
GRASS, CAT,  
TREE, SKY

No objects, just pixels



CAT

Single Object



DOG, DOG, CAT

Multiple Object



DOG, DOG, CAT

This image is CC0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

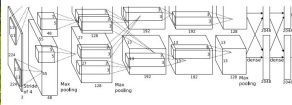
Lecture 11 - 45 May 10, 2017



# Classification + Localization



This image is CC0 public domain



**Fully  
Connected:**  
4096 to 1000

**Class Scores**

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

**Vector:**  
4096

**Fully  
Connected:**  
4096 to 4

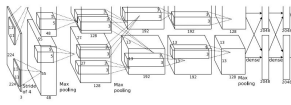
**Box  
Coordinates**  
(x, y, w, h)

Treat localization as a  
regression problem!

# Classification + Localization



This image is CC0 public domain



Fully  
Connected:  
4096 to 1000

**Class Scores**

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

Correct label:  
Cat

Softmax  
Loss

Vector:  
4096

Fully  
Connected:  
4096 to 4

**Box  
Coordinates**  
(x, y, w, h)

L2 Loss

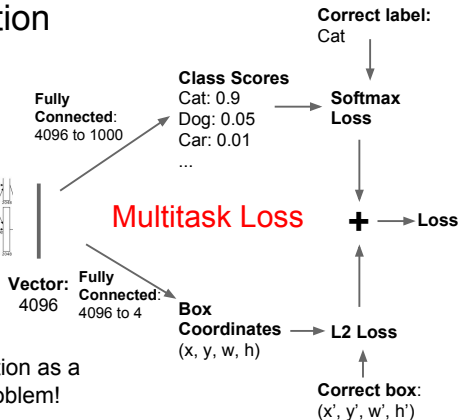
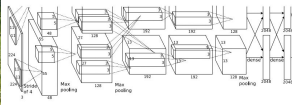
Correct box:  
(x', y', w', h')

Treat localization as a  
regression problem!

# Classification + Localization



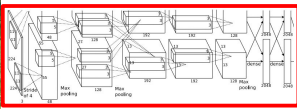
This image is CC0 public domain



# Classification + Localization



This image is CC0 public domain



Often pretrained on ImageNet  
(Transfer learning)

Vector:  
4096

Fully  
Connected:  
4096 to 1000

**Class Scores**

Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

Fully  
Connected:  
4096 to 4

**Box  
Coordinates**  
(x, y, w, h)

**Correct label:**  
Cat

**Softmax  
Loss**

+

**Loss**

**L2 Loss**

**Correct box:**  
(x', y', w', h')

Treat localization as a  
regression problem!

## Aside: Human Pose Estimation



Represent pose as a set of 14 joint positions:

- Left / right foot
- Left / right knee
- Left / right hip
- Left / right shoulder
- Left / right elbow
- Left / right hand
- Neck
- Head top

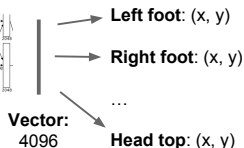
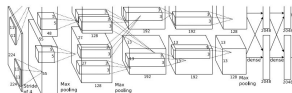
This image is licensed under CC-BY 2.0

Johnson and Everingham, "Clustered Pose and Nonlinear Appearance Models for Human Pose Estimation", BMVC 2010

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 50 May 10, 2017

## Aside: Human Pose Estimation

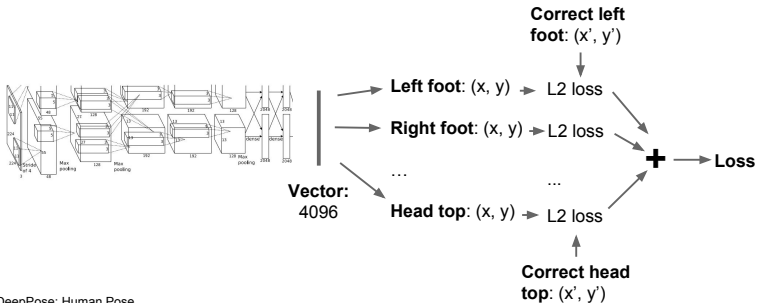


Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 51 May 10, 2017

## Aside: Human Pose Estimation

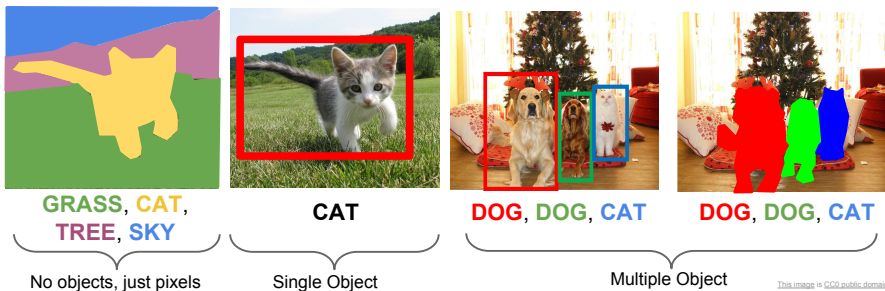


Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 52 May 10, 2017

## Object Detection



Fei-Fei Li &amp; Justin Johnson &amp; Serena Yeung

Lecture 11 - 53 May 10, 2017



# Object Detection: Impact of Deep Learning

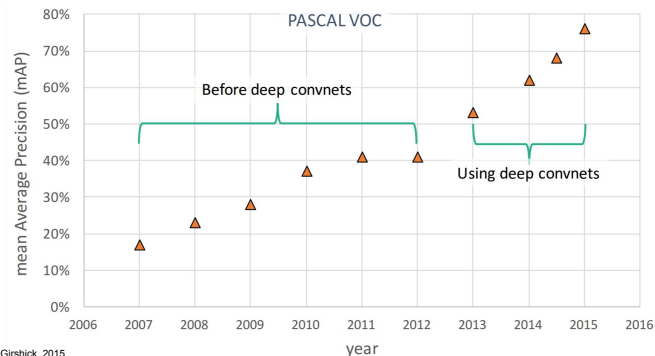


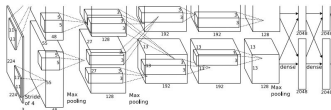
Figure copyright Ross Girshick, 2015.  
Reproduced with permission.





# Object Detection as Classification: Sliding Window

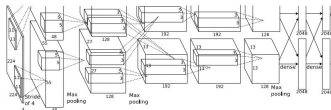
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
 Cat? NO  
 Background? YES

# Object Detection as Classification: Sliding Window

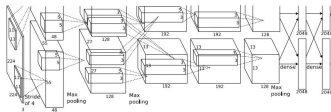
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection as Classification: Sliding Window

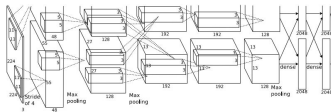
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection as Classification: Sliding Window

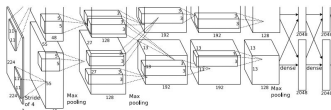
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? YES  
Background? NO

# Object Detection as Classification: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



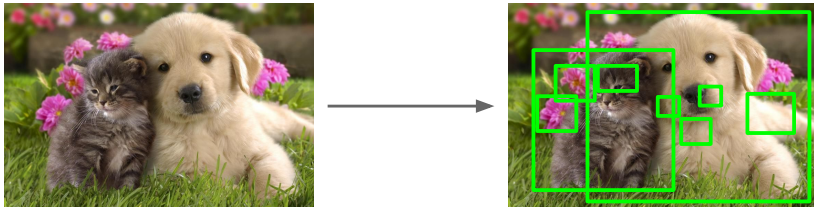
Dog? NO  
Cat? YES  
Background? NO

Problem: Need to apply CNN to huge number of locations and scales, very computationally expensive!



## Region Proposals

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



Alexe et al, "Measuring the objectness of image windows", TPAMI 2012

Uijings et al, "Selective Search for Object Recognition", IJCV 2013

Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014

Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 62 May 10, 2017

# R-CNN



Input image

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 63 May 10, 2017

# R-CNN



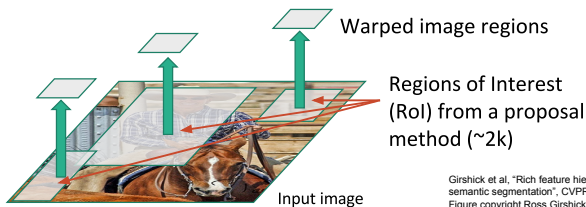
Regions of Interest  
(RoI) from a proposal  
method (~2k)

Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 64 May 10, 2017

## R-CNN

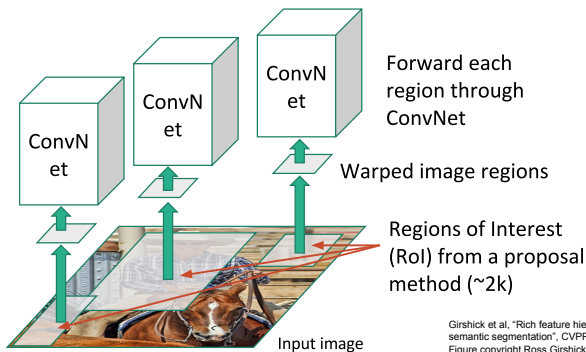


Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

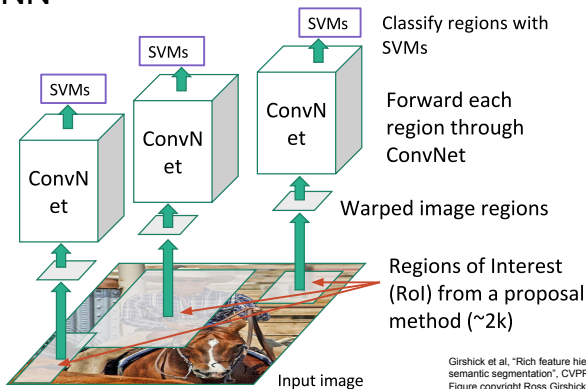
Lecture 11 - 65 May 10, 2017

## R-CNN



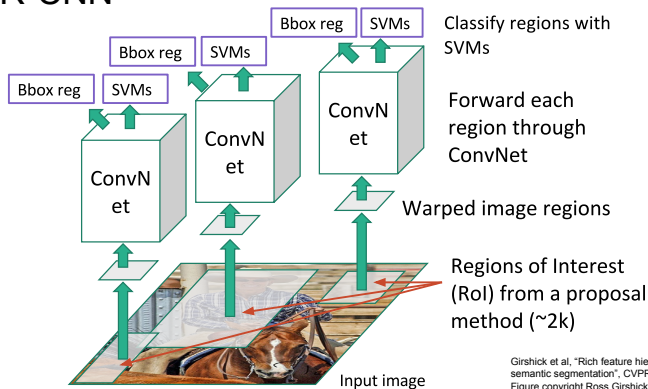
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

## R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

## R-CNN



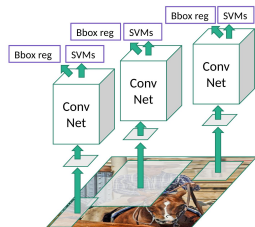
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 68 May 10, 2017

## R-CNN: Problems

- Ad hoc training objectives
  - Fine-tune network with softmax classifier (log loss)
  - Train post-hoc linear SVMs (hinge loss)
  - Train post-hoc bounding-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 et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.  
 Slide copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



# Fast R-CNN



Input image

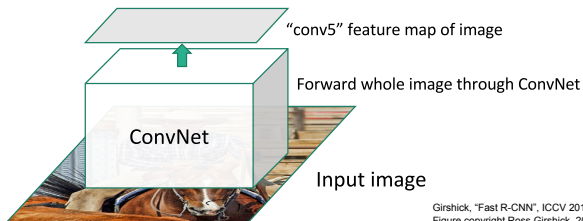
Girshick, "Fast R-CNN", ICCV 2015.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

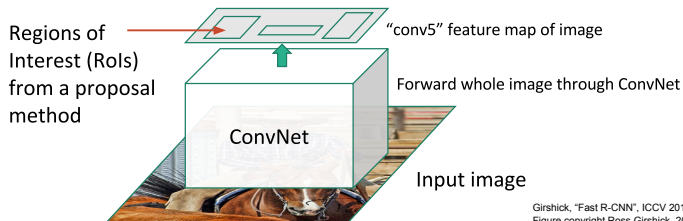
Lecture 11 - 70 May 10, 2017

# Fast R-CNN



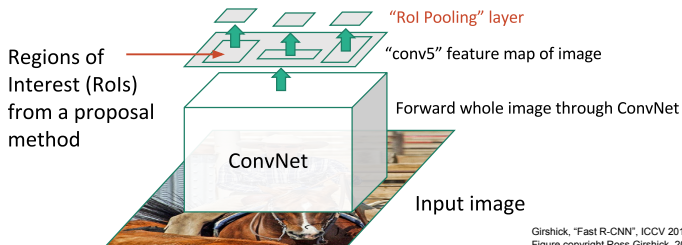
Girshick, “Fast R-CNN”, ICCV 2015.  
 Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



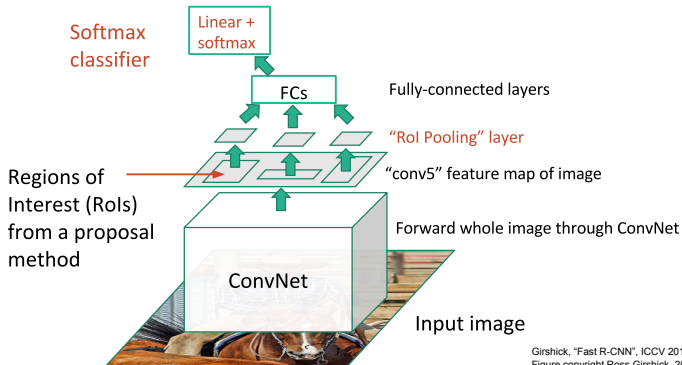
Girshick, “Fast R-CNN”, ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



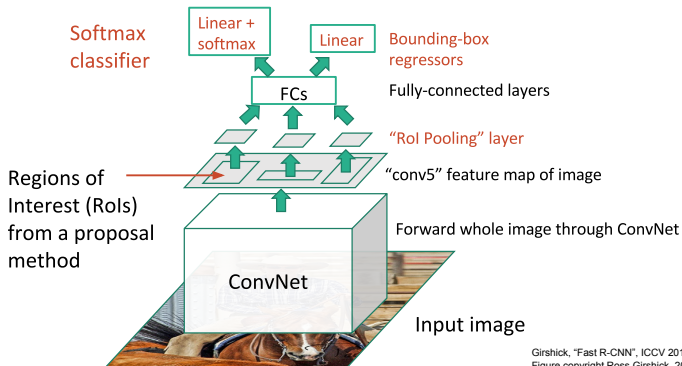
Girshick, “Fast R-CNN”, ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



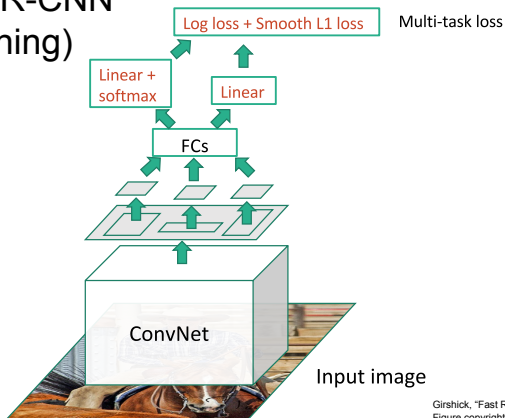
Girshick, "Fast R-CNN", ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015.  
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

# Fast R-CNN (Training)



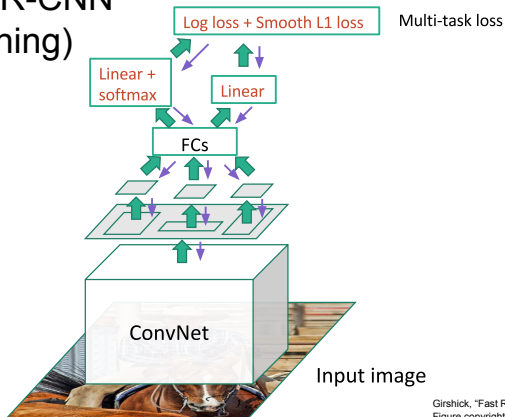
Girshick, "Fast R-CNN", ICCV 2015.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 76 May 10, 2017

# Fast R-CNN (Training)



Girshick, "Fast R-CNN", ICCV 2015.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

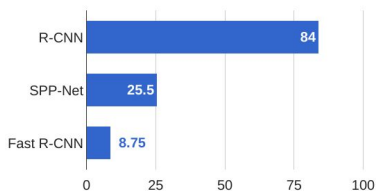
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 77 May 10, 2017

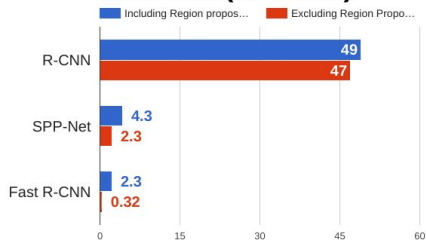


# R-CNN vs SPP vs Fast R-CNN

## Training time (Hours)

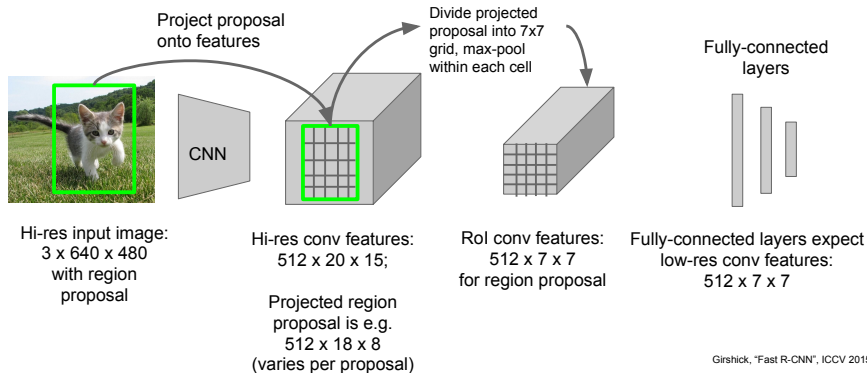


## Test time (seconds)



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

# Faster R-CNN: RoI Pooling

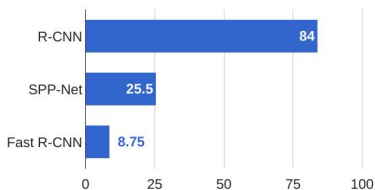


Fei-Fei Li & Justin Johnson & Serena Yeung

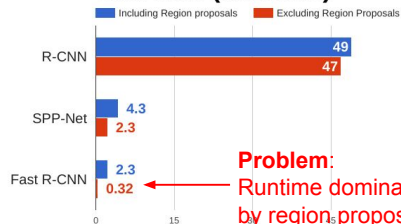
Lecture 11 - 78 May 10, 2017

## R-CNN vs SPP vs Fast R-CNN

Training time (Hours)



Test time (seconds)



**Problem:**  
Runtime dominated  
by region proposals!

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 80 May 10, 2017

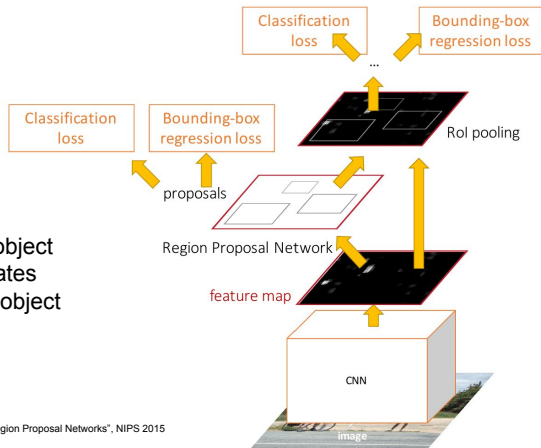
## Faster R-CNN:

Make CNN do proposals!

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

Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



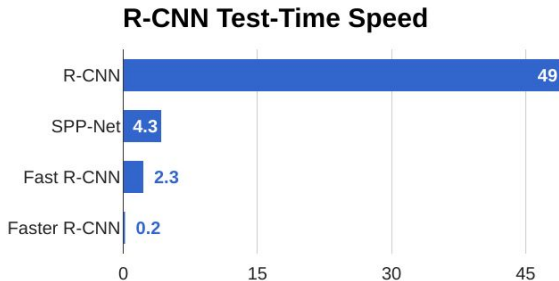
Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015  
Figure copyright 2015, Ross Girshick; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

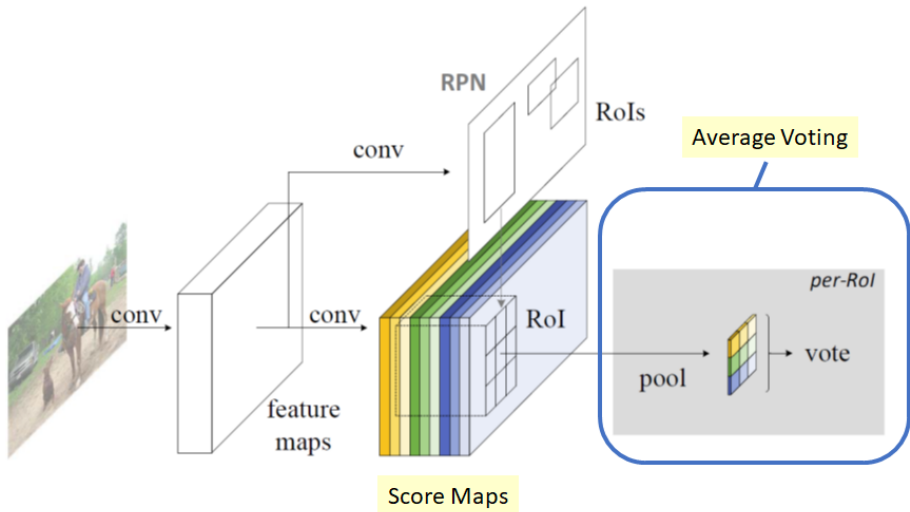
Lecture 11 - 81 May 10, 2017

## Faster R-CNN:

Make CNN do proposals!

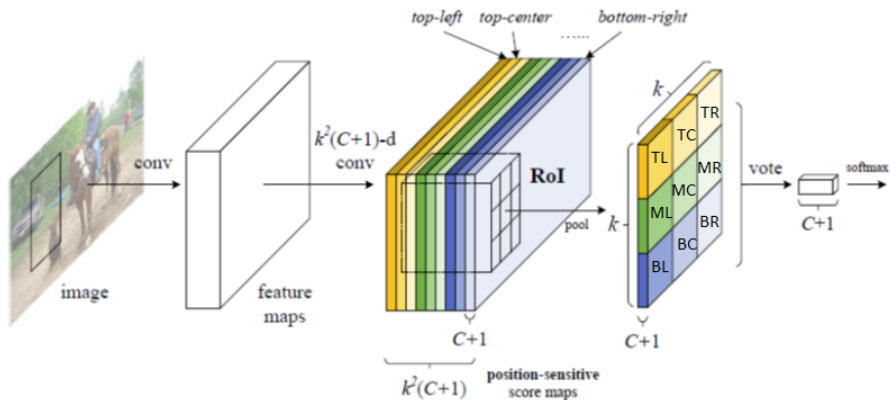


# Region-based fully convolutional network (R-FCN)

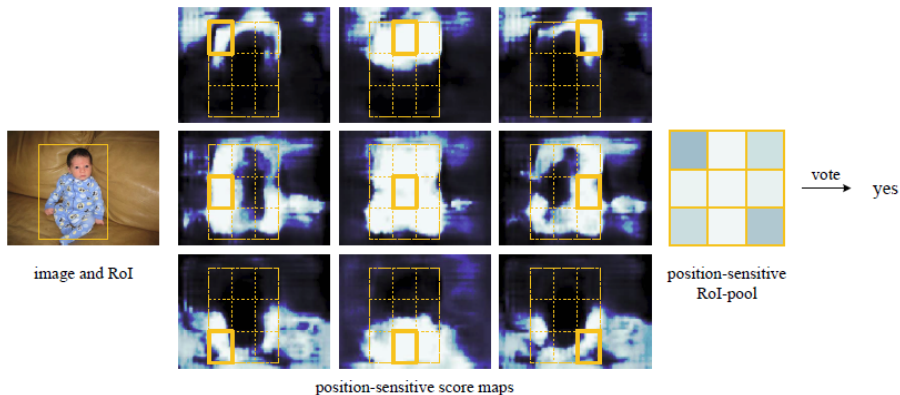


Fully connected layers are replaced by average pooling

## Region-based fully convolutional network (R-FCN)

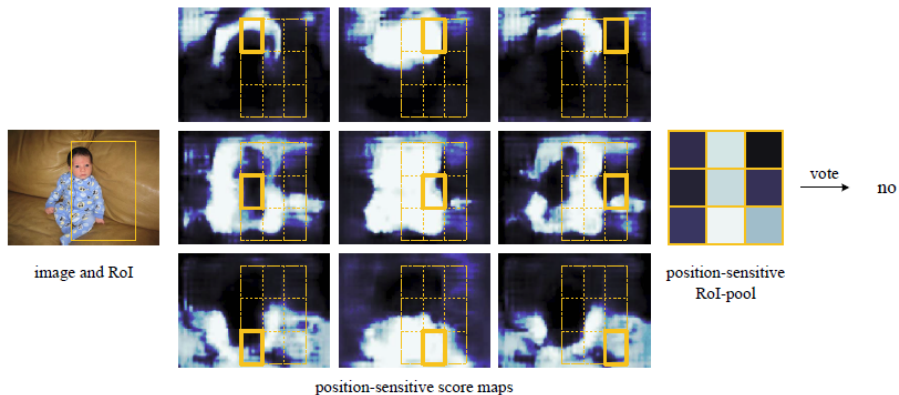


## Region-based fully convolutional network (R-FCN)

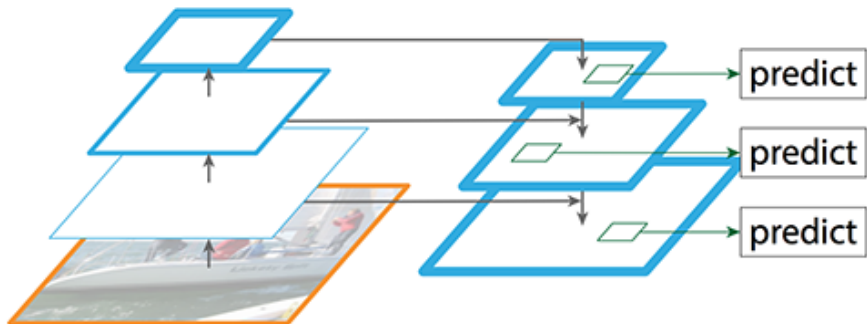




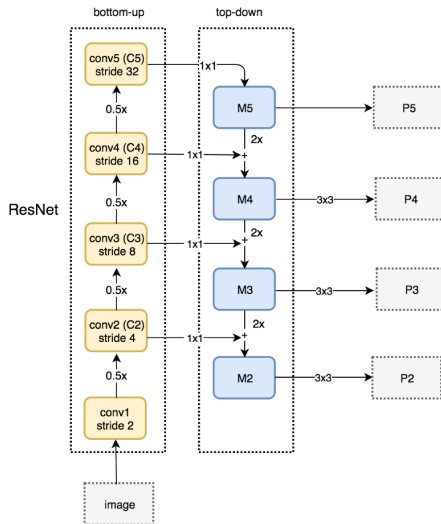
## Region-based fully convolutional network (R-FCN)



# Feature pyramid network (FPN)



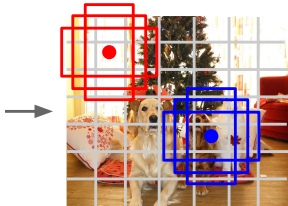
# Feature pyramid network (FPN)



# Detection without Proposals: YOLO / SSD



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
 centered at each grid cell  
 Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
 $(dx, dy, dh, dw, \text{confidence})$
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

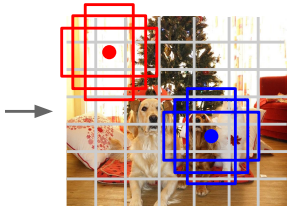
Redmon et al, "You Only Look Once:  
 Unified, Real-Time Object Detection", CVPR 2016  
 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

# Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network! →



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

Image a set of **base boxes**  
centered at each grid cell  
Here  $B = 3$

Within each grid cell:

- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  
( $dx, dy, dh, dw, confidence$ )
- Predict scores for each of  $C$  classes (including background as a class)

Output:  
 $7 \times 7 \times (5 * B + C)$

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 84 May 10, 2017

## Object Detection: Lots of variables ...

### Base Network

VGG16

ResNet-101

Inception V2

Inception V3

Inception

ResNet

MobileNet

### Object Detection architecture

Faster R-CNN

R-FCN

SSD

### Image Size # Region Proposals

...

### Takeaways

Faster R-CNN is  
slower but more  
accurate

SSD is much  
faster but not as  
accurate

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016

Inception-V2: Ioffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2016

Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016

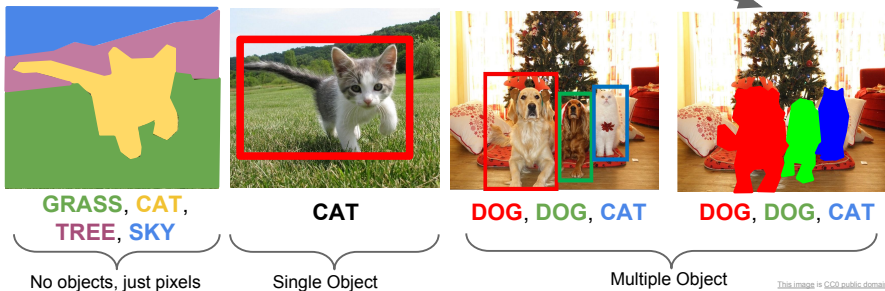
Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016

MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 85 May 10, 2017

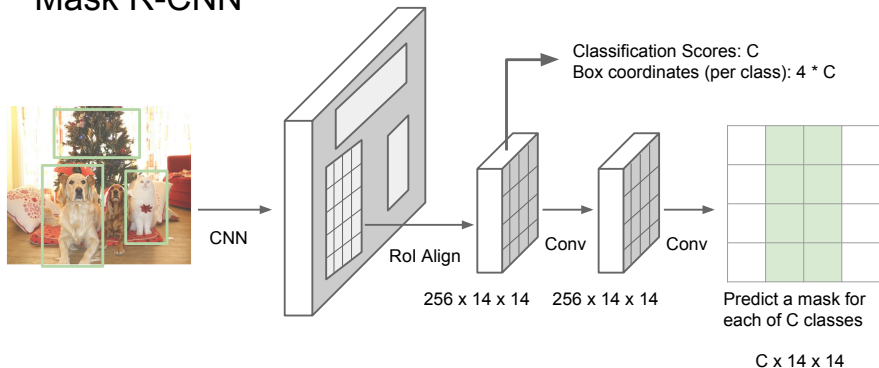
# Instance Segmentation



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 89 May 10, 2017

## Mask R-CNN



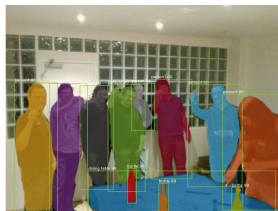
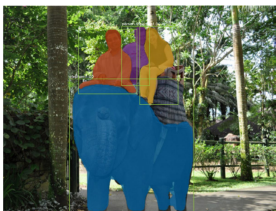
He et al, "Mask R-CNN", arXiv 2017

Fei-Fei Li &amp; Justin Johnson &amp; Serena Yeung

Lecture 11 - 90 May 10, 2017



# Mask R-CNN: Very Good Results!



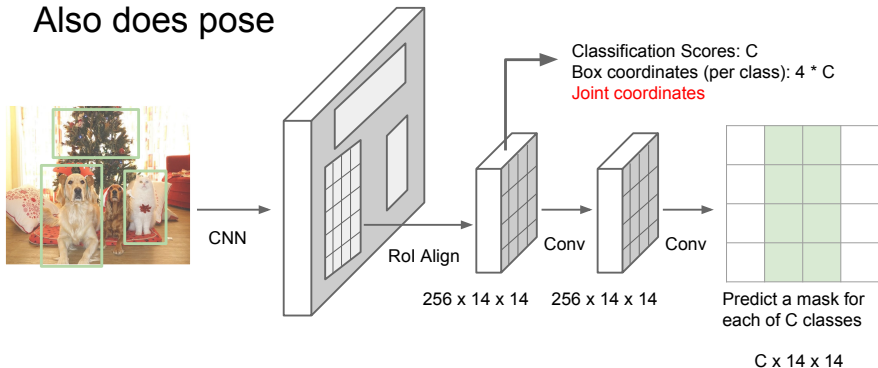
He et al, "Mask R-CNN", arXiv 2017  
 Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.  
 Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 91 May 10, 2017

# Mask R-CNN

## Also does pose



He et al, "Mask R-CNN", arXiv 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 92 May 10, 2017

# Mask R-CNN

## Also does pose



He et al, "Mask R-CNN", arXiv 2017  
 Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.  
 Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 93 May 10, 2017