

# Deep Learning (for Computer Vision)

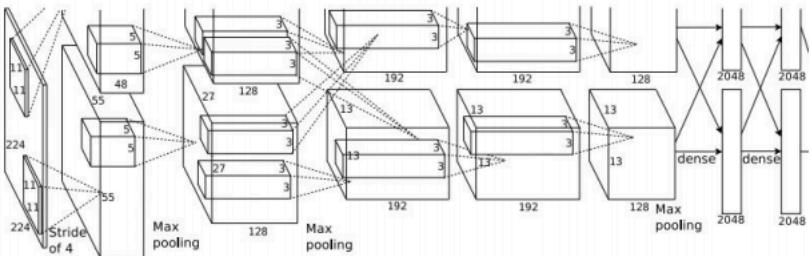
Arjun Jain | 8 April 2017

# Sources

A lot of the material has been shamelessly and gratefully collected from:

- <http://cs231n.stanford.edu/>

# So Far: Image Classification



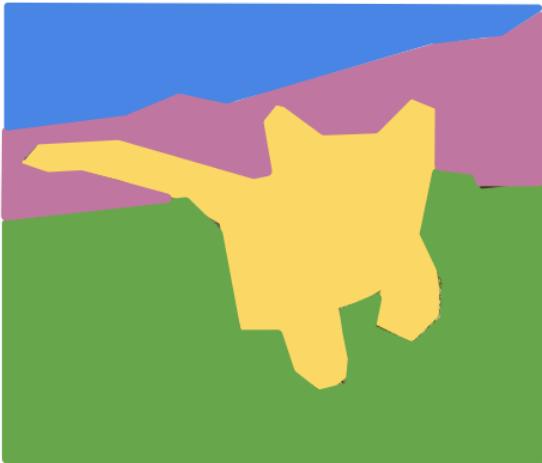
**Vector:**  
4096

**Fully-Connected:**  
4096 to 1000

**Class Scores**  
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

# Other Computer Vision Tasks

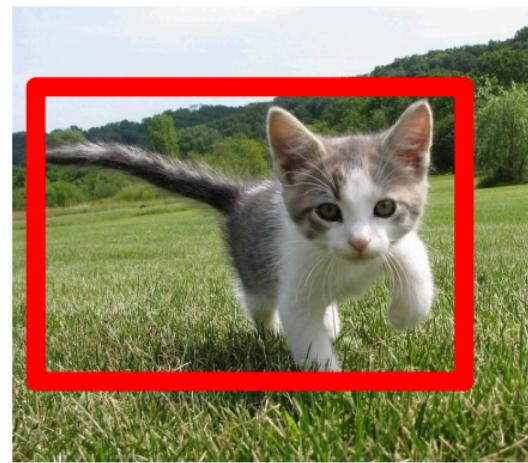
## Semantic Segmentation



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



GRASS, CAT,  
TREE, SKY

No objects, just pixels

## Classification + Localization



CAT

Single Object

## Object Detection



DOG, DOG, CAT

Multiple Object

## Instance Segmentation

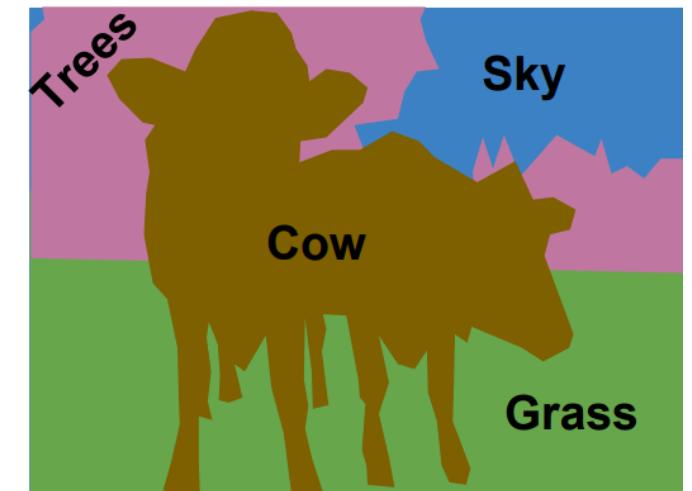
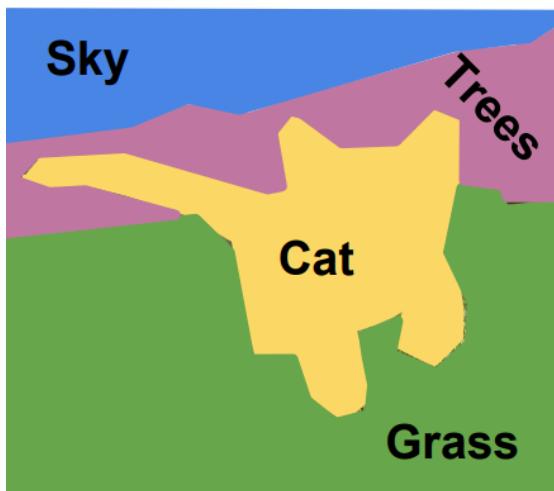


DOG, DOG, CAT

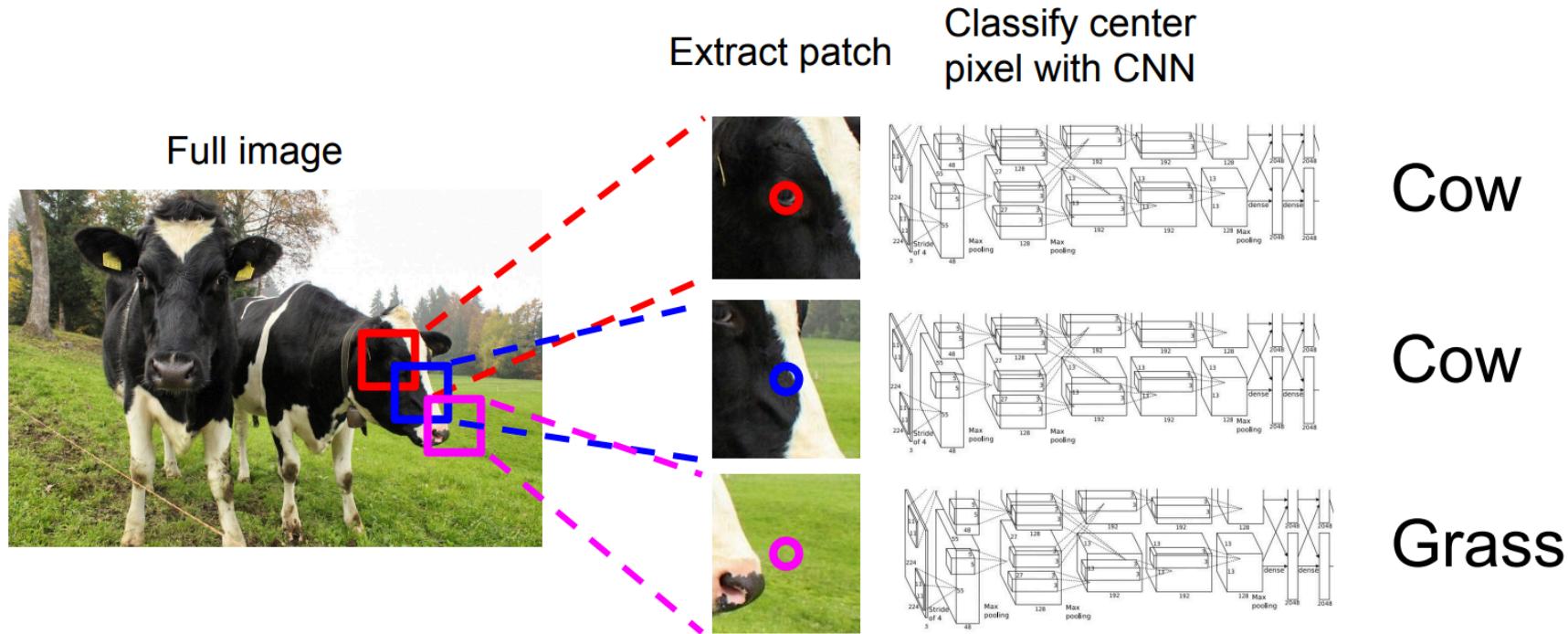
This image is CC0 public domain

# Semantic Segmentation

- Label each pixel in the image with a category label
- Don't differentiate instances, only care about pixels



# Semantic Segmentation Idea: Sliding Window

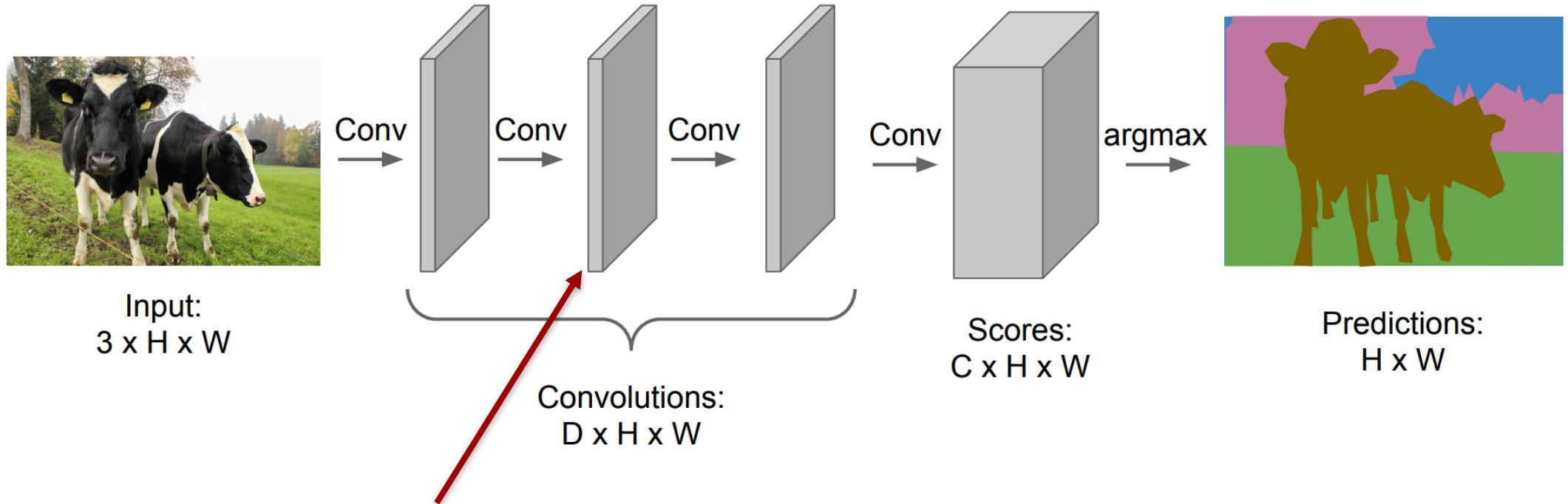


Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013  
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Semantic Segmentation Idea: Fully Convolutional

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



Problem: convolutions at original image  
resolution will be very expensive ...

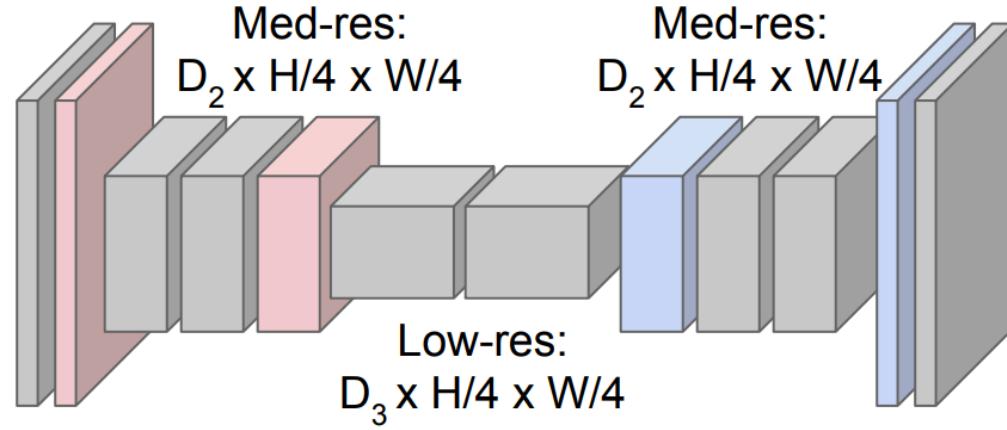
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution

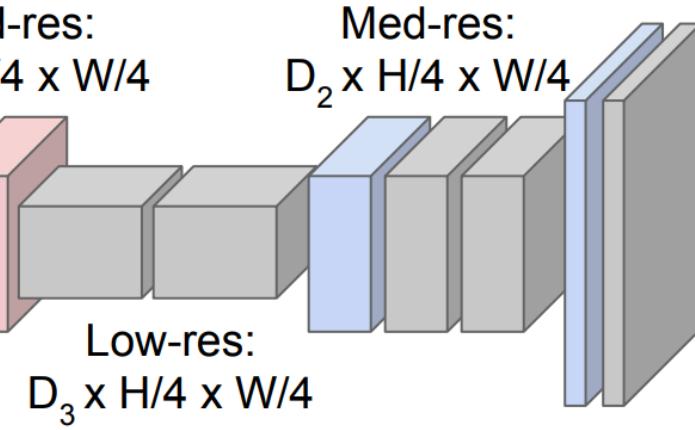


Input:  
 $3 \times H \times W$

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

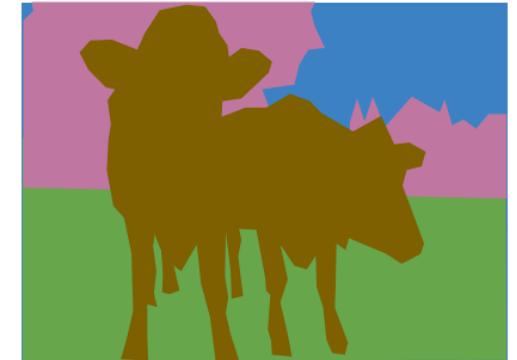


High-res:  
 $D_1 \times H/2 \times W/2$



High-res:  
 $D_1 \times H/2 \times W/2$

**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

# Upsampling: Nearest Neighbor

## 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

# Upsampling: Bed of Nails

**“Bed of Nails”**

1	2
3	4



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

Input:  $2 \times 2$

Output:  $4 \times 4$

# Upsampling: Max Unpooling

## Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Rest of the network

## Max Unpooling

Use positions from pooling layer

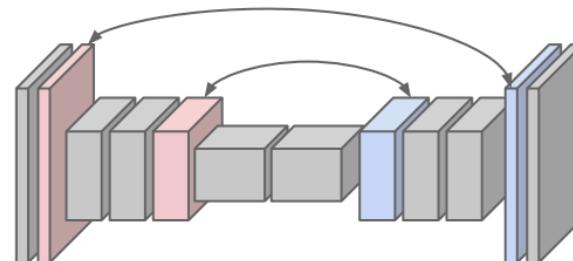
1	2
3	4

Input: 2 x 2

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

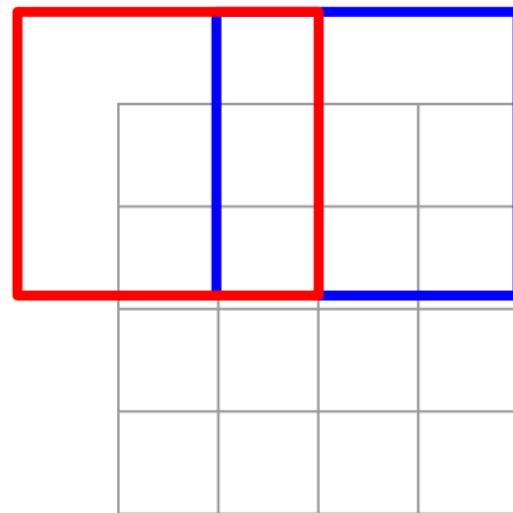
Output: 4 x 4

Corresponding pairs of  
downsampling and  
upsampling layers



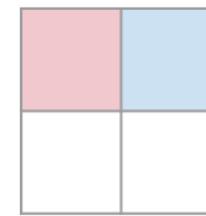
# Upsampling: Transposed Convolution (the backprop of normal convolution with stride > 1)

**Recall:** Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$

Dot product  
between filter  
and input



Output:  $2 \times 2$

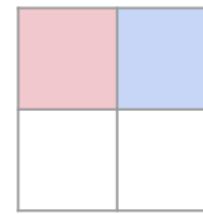
Filter moves 2 pixels in  
the input for every one  
pixel in the output

Stride gives ratio between  
movement in input and  
output

# Upsampling: Transposed Convolution (the backprop of normal convolution with stride > 1)

**Other names:**

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

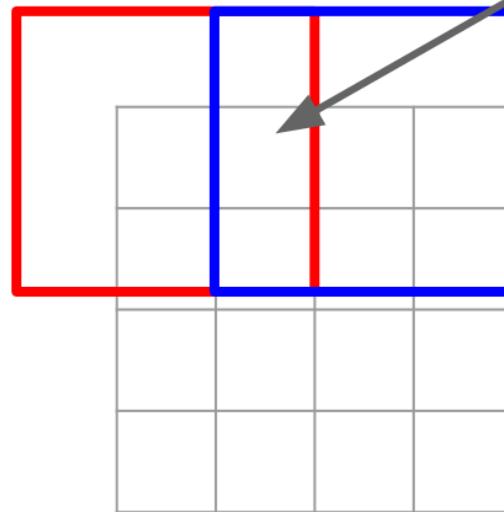


Input: 2 x 2

3 x 3 transpose convolution, stride 2 pad 1



Input gives weight for filter



Output: 4 x 4

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

Stride gives ratio between movement in output and input

Take the scalar value (pink) and then multiply the weights by that scalar and copy the result to the 3x3 region.

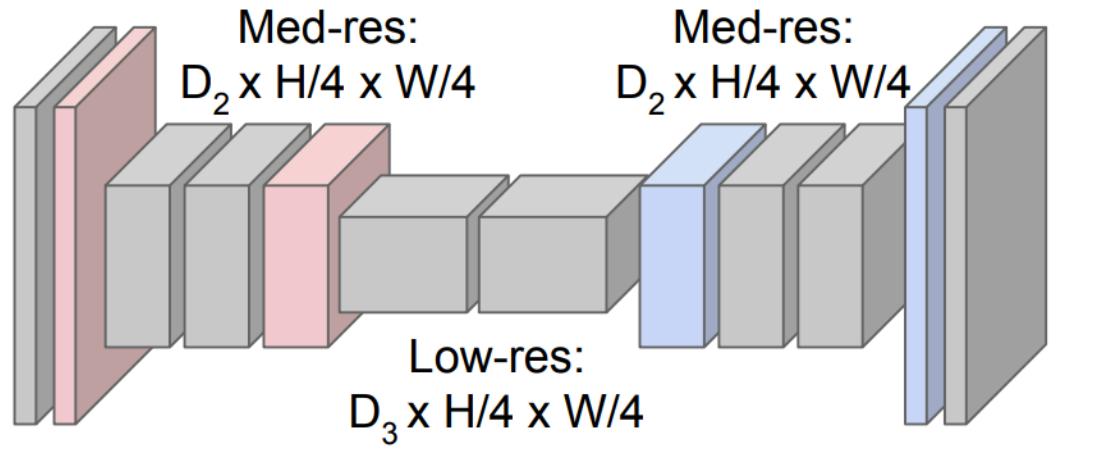
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution

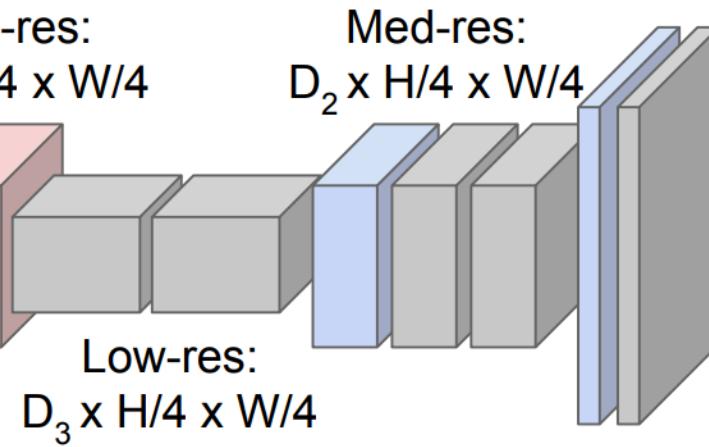


Input:  
 $3 \times H \times W$

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



High-res:  
 $D_1 \times H/2 \times W/2$



High-res:  
 $D_1 \times H/2 \times W/2$

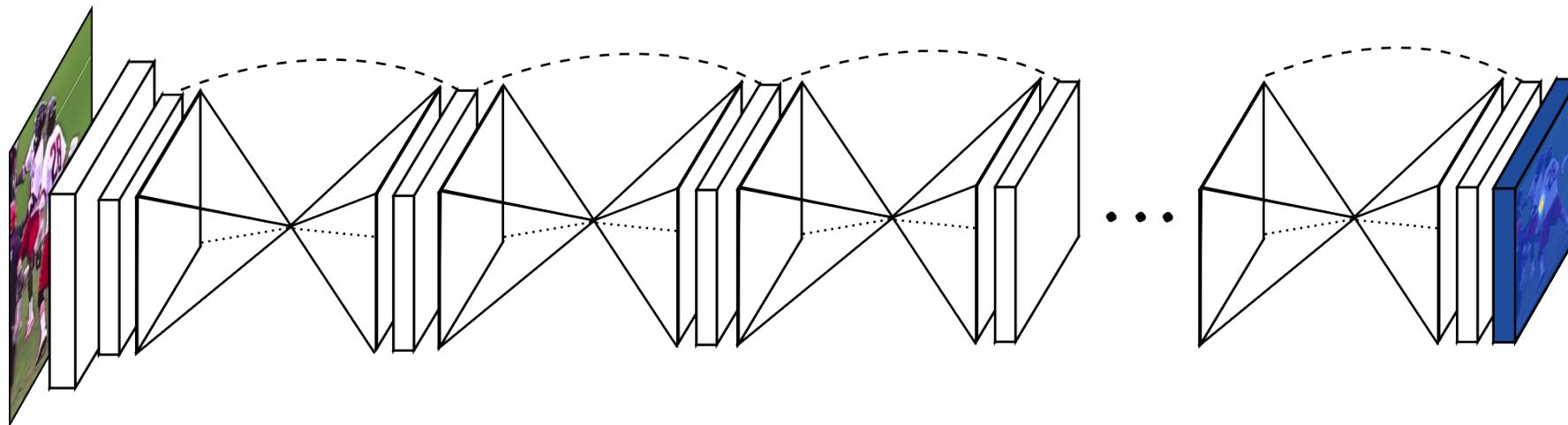
**Upsampling:**  
Unpooling or strided transpose convolution



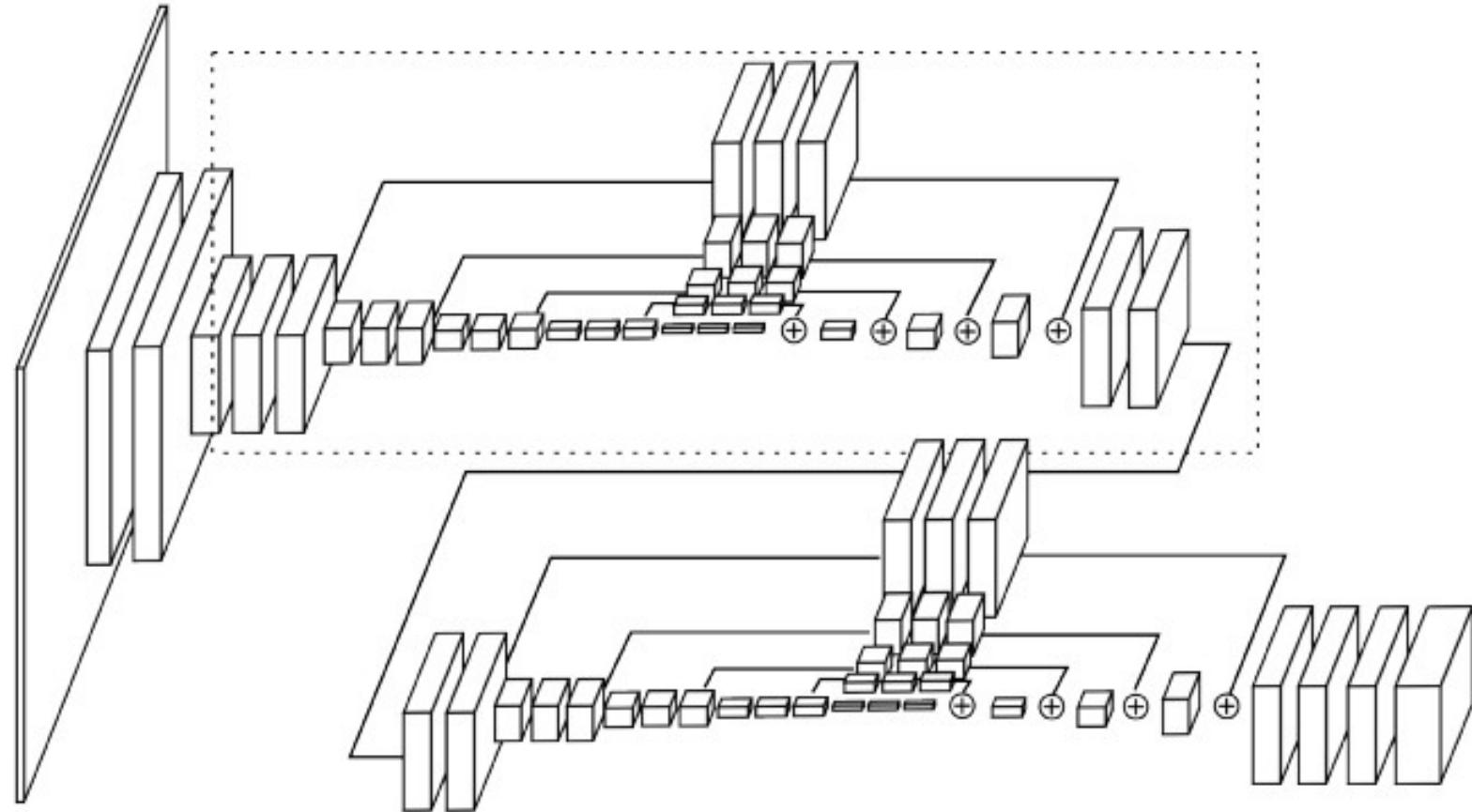
Predictions:  
 $H \times W$

# Aside: Stacked Hour Glass Networks

- Popular architecture for human pose estimation
- Progressive pooling followed by progressive up sampling, creating an hourglass form.
- Input are images showing a person's body.
- Outputs are K heatmaps

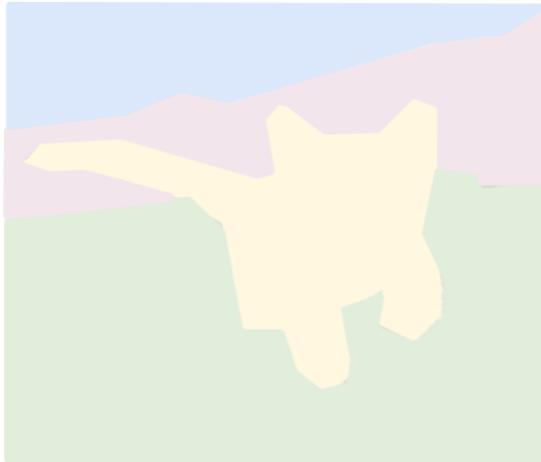


# Aside: Stacked Hour Glass Networks



# Classification + Localization

Semantic Segmentation



GRASS, CAT,  
TREE, SKY

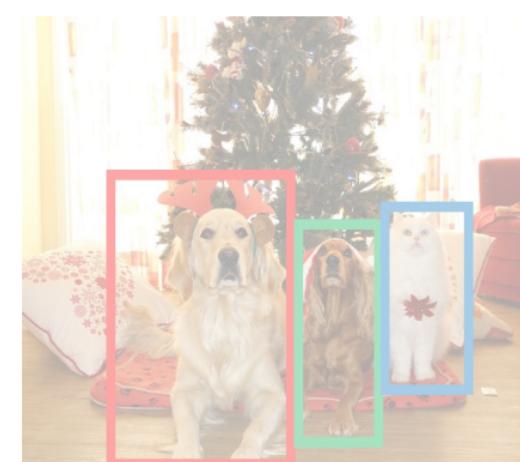
No objects, just pixels

Classification  
+ Localization



CAT  
Single Object

Object  
Detection



DOG, DOG, CAT

Multiple Object

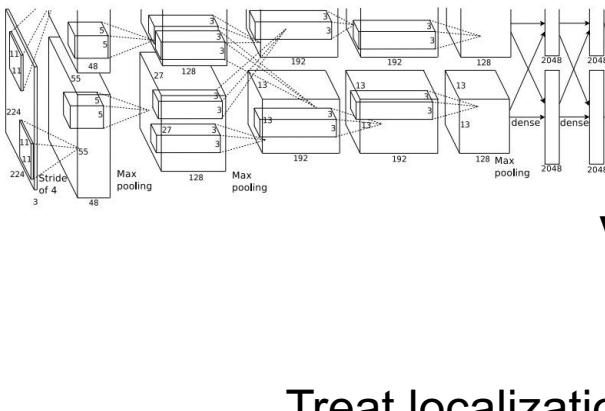
Instance  
Segmentation



DOG, DOG, CAT

This image is CC0 public domain

# Classification + Localization

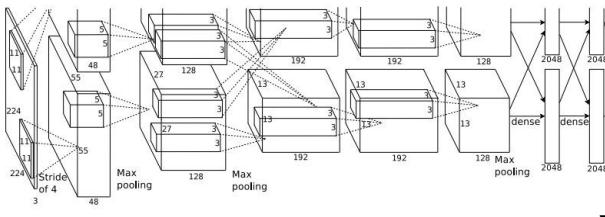


**Fully Connected:** 4096 to 1000  
**Class Scores**  
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...

**Vector:** 4096  
**Fully Connected:** 4096 to 4  
**Box Coordinates**  
( $x, y, w, h$ )

Treat localization as a regression problem!

# Classification + Localization



Treat localization as a  
regression problem!

Fully  
Connected:  
4096 to 1000

Vector:  
4096  
Fully  
Connected:  
4096 to 4

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

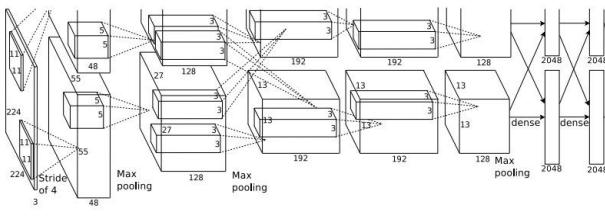
Box  
Coordinates —> L2 Loss  
( $x, y, w, h$ )

Correct label:  
Cat

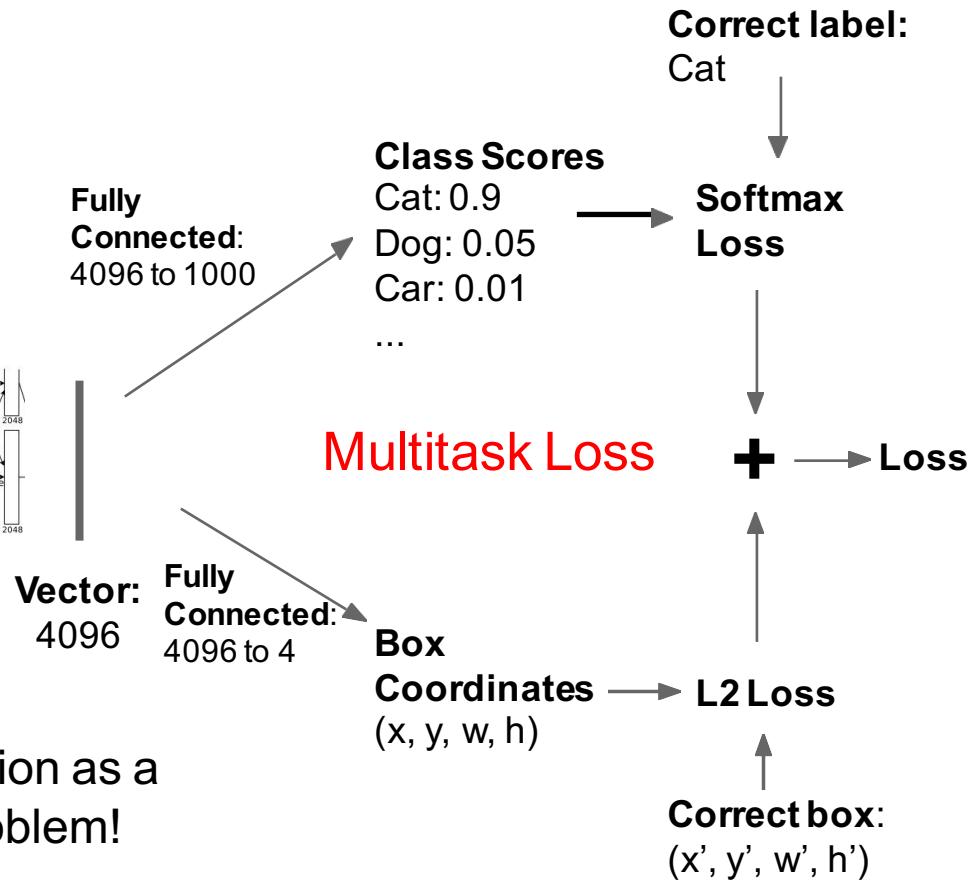
Softmax  
Loss

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

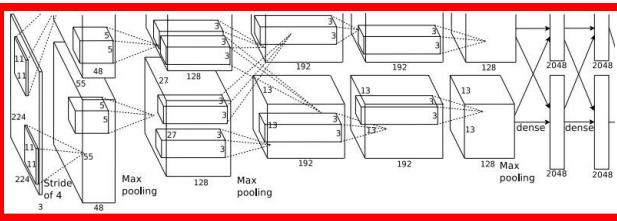
# Classification + Localization



Treat localization as a  
regression problem!



# Classification + Localization

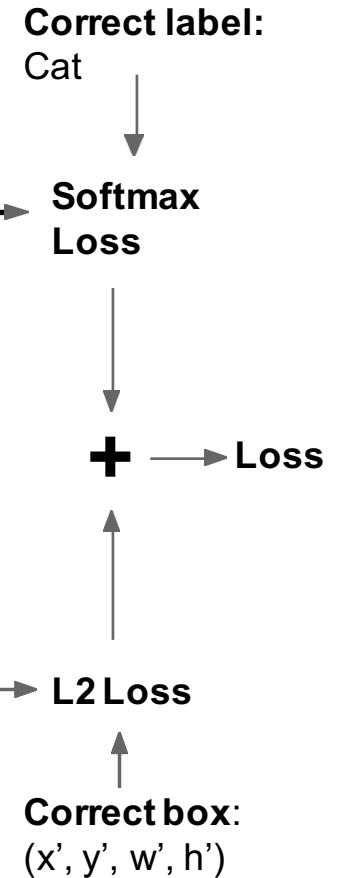


Often pretrained on ImageNet  
(Transfer learning)

Treat localization as a  
regression problem!

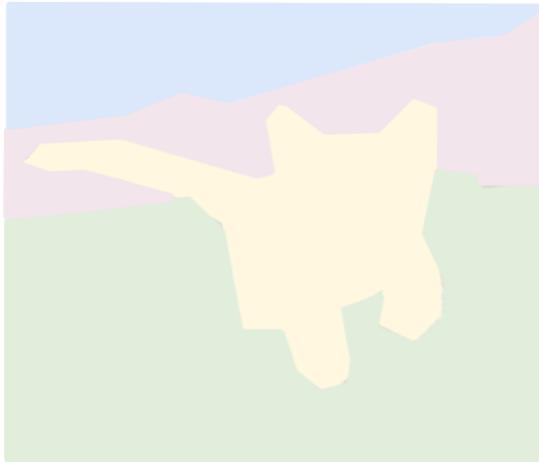
Fully Connected: 4096 to 1000  
  
Vector: 4096

Class Scores  
Cat: 0.9  
Dog: 0.05  
Car: 0.01  
...  
  
Fully Connected: 4096 to 4  
  
Box Coordinates  $\rightarrow$  L2 Loss  
( $x, y, w, h$ )



# Object Detection

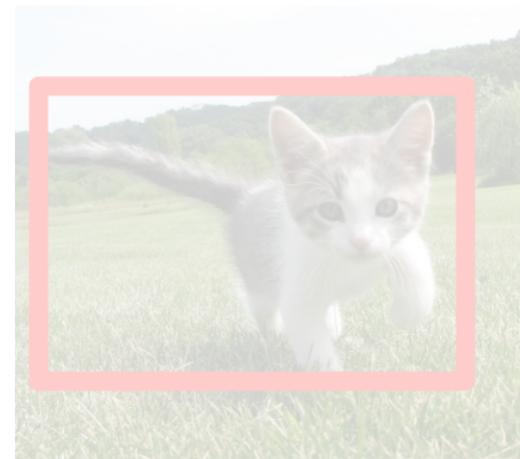
Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

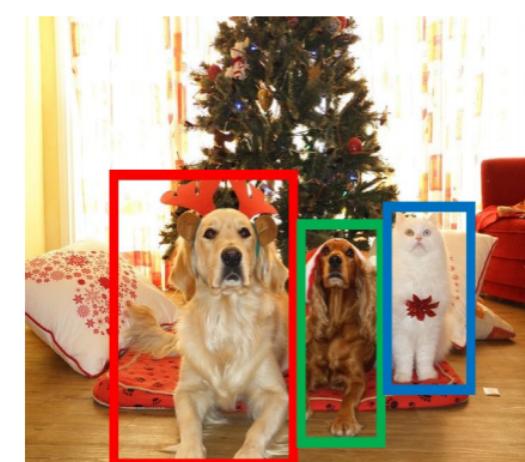
Classification  
+ Localization



CAT

Single Object

Object  
Detection



DOG, DOG, CAT

Multiple Object

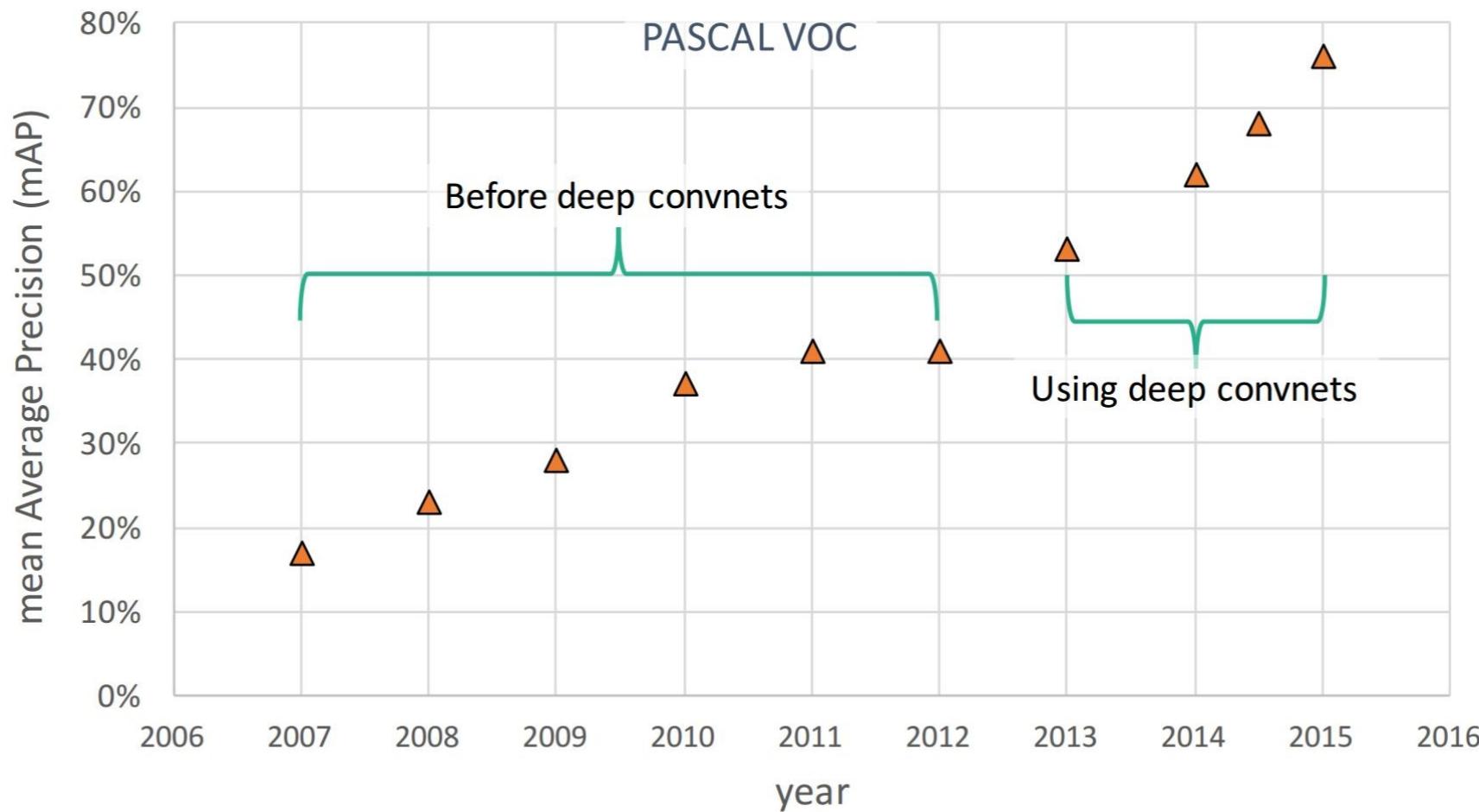
Instance  
Segmentation



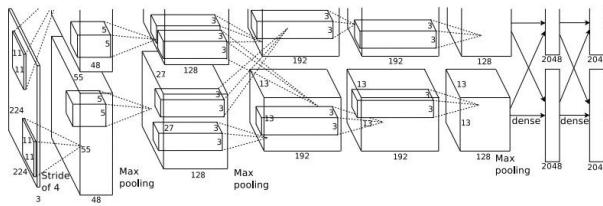
DOG, DOG, CAT

This image is CC0 public domain

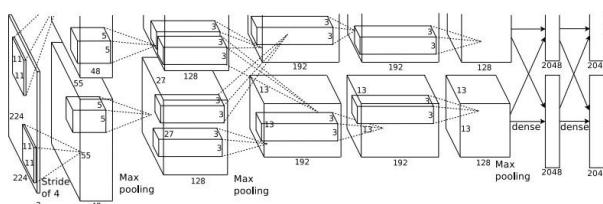
# Object Detection: Impact of ConvNets



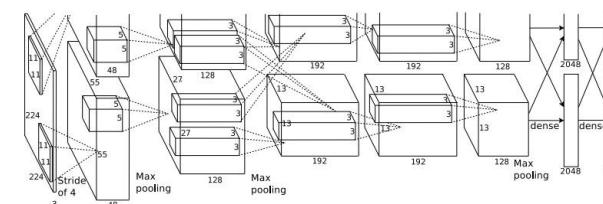
# Object Detection as Regression?



CAT: (x, y, w, h)



DOG: (x, y, w, h)  
DOG: (x, y, w, h)  
CAT: (x, y, w, h)

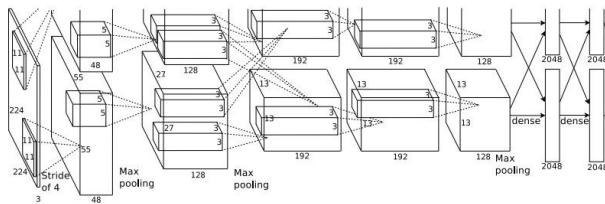


DUCK: (x, y, w, h)  
DUCK: (x, y, w, h)

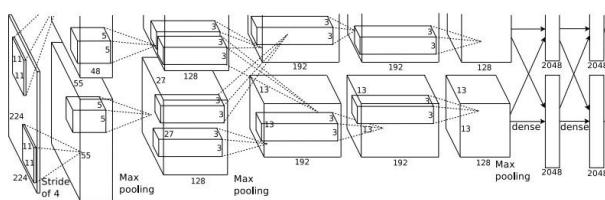
...

# Object Detection as Regression?

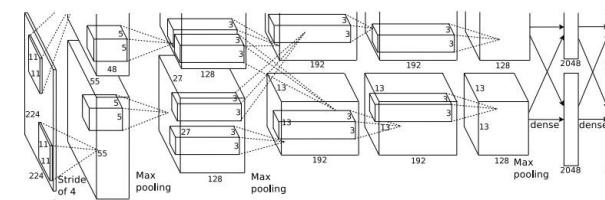
Each image needs a different number of outputs!



CAT:  $(x, y, w, h)$  4 numbers



DOG:  $(x, y, w, h)$   
DOG:  $(x, y, w, h)$  16 numbers  
CAT:  $(x, y, w, h)$

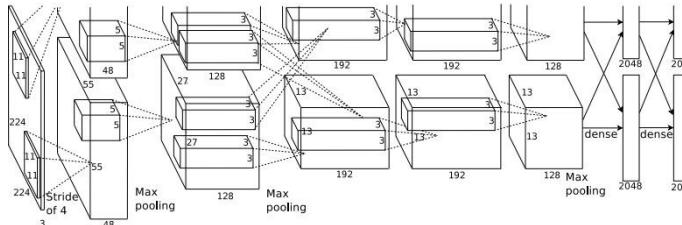


DUCK:  $(x, y, w, h)$  Many numbers!  
DUCK:  $(x, y, w, h)$

...

# Object Detection as Classification: Sliding Window

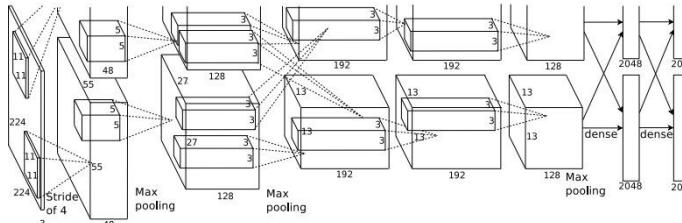
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO  
Cat? NO  
Background? YES

# Object Detection as Classification: Sliding Window

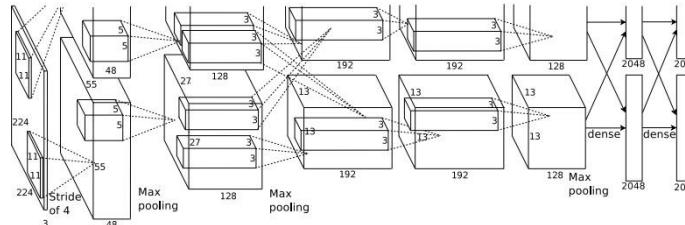
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES  
Cat? NO  
Background? NO

# Object Detection as Classification: Sliding Window

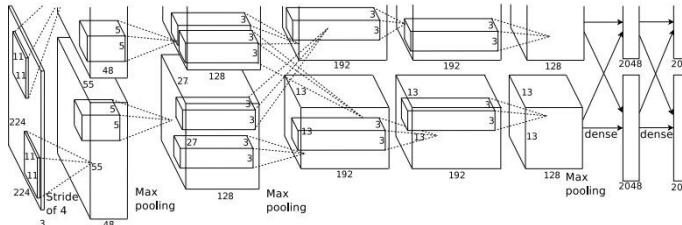
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



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

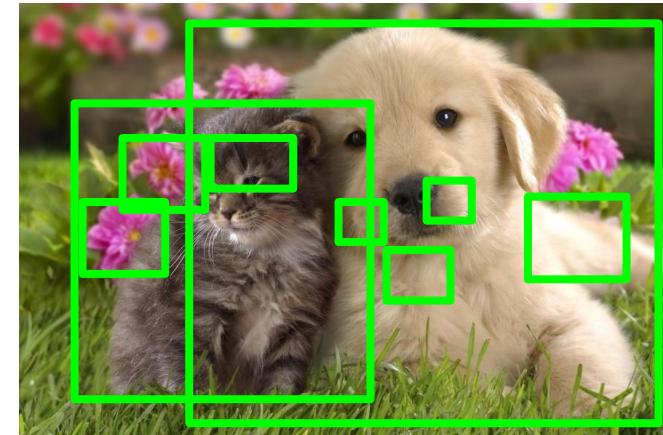


Dog? NO  
Cat? YES  
Background? NO

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 – non DL based algorithm
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU



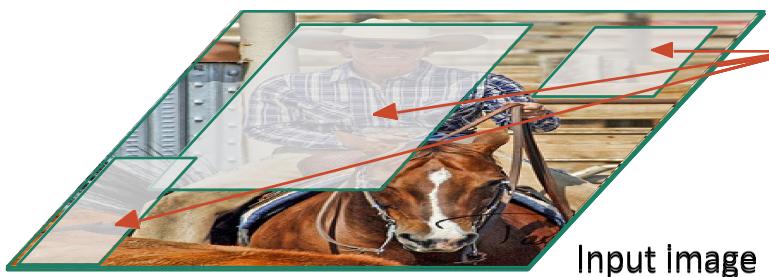
# R-CNN



Input image

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

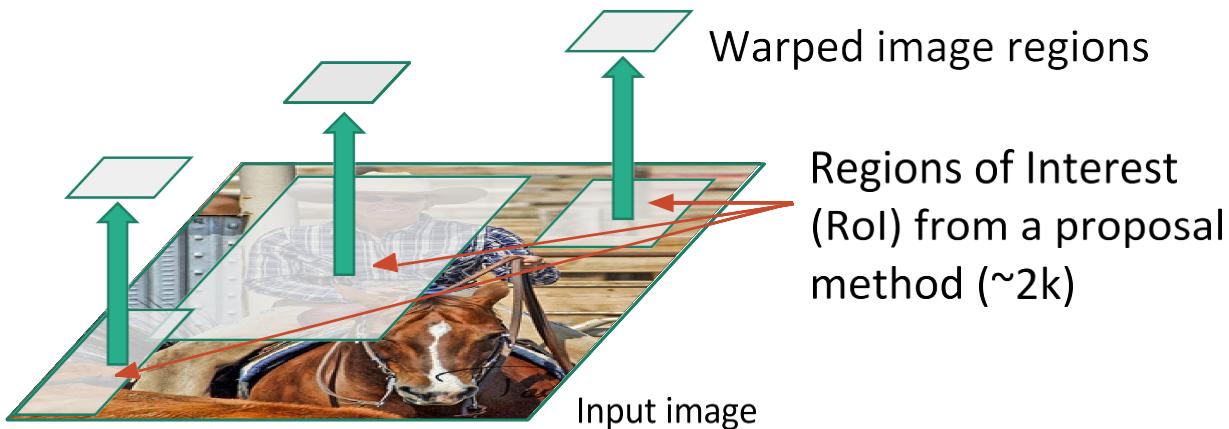
# R-CNN



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

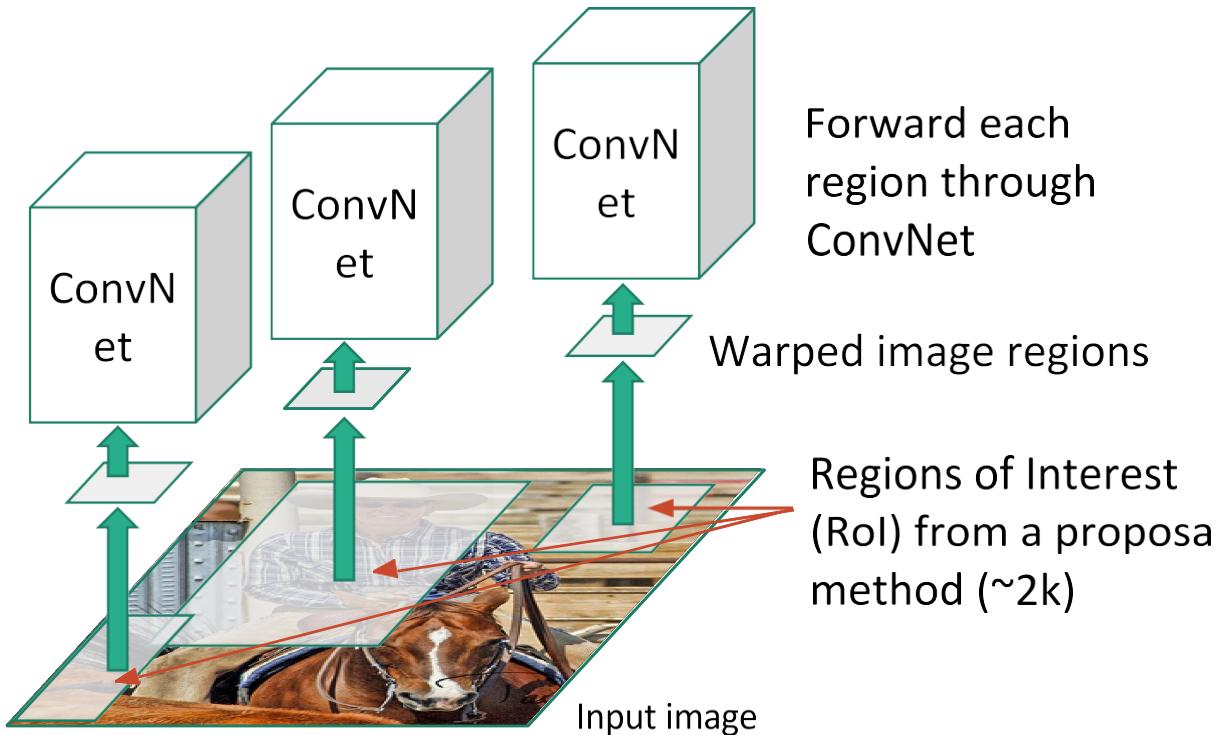
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

# R-CNN



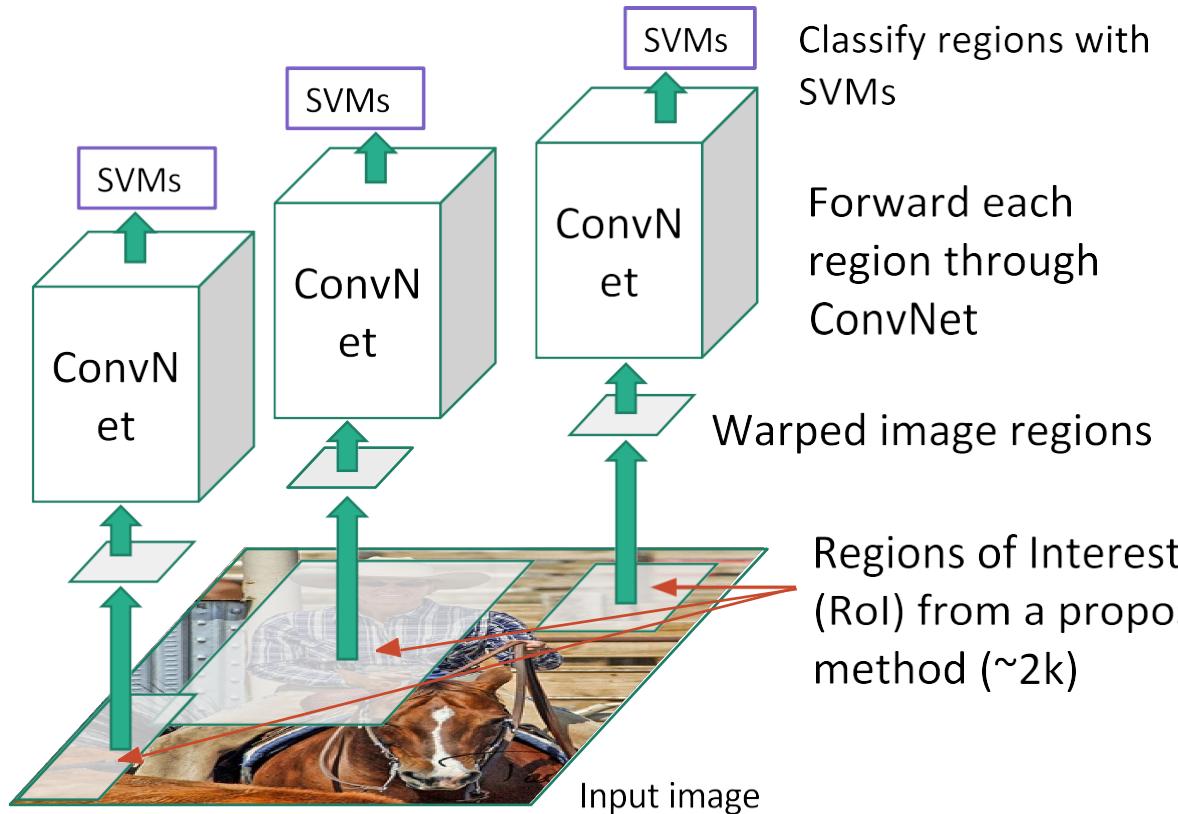
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

# R-CNN



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

# R-CNN



Linear Regression for bounding box offsets

Classify regions with  
SVMs

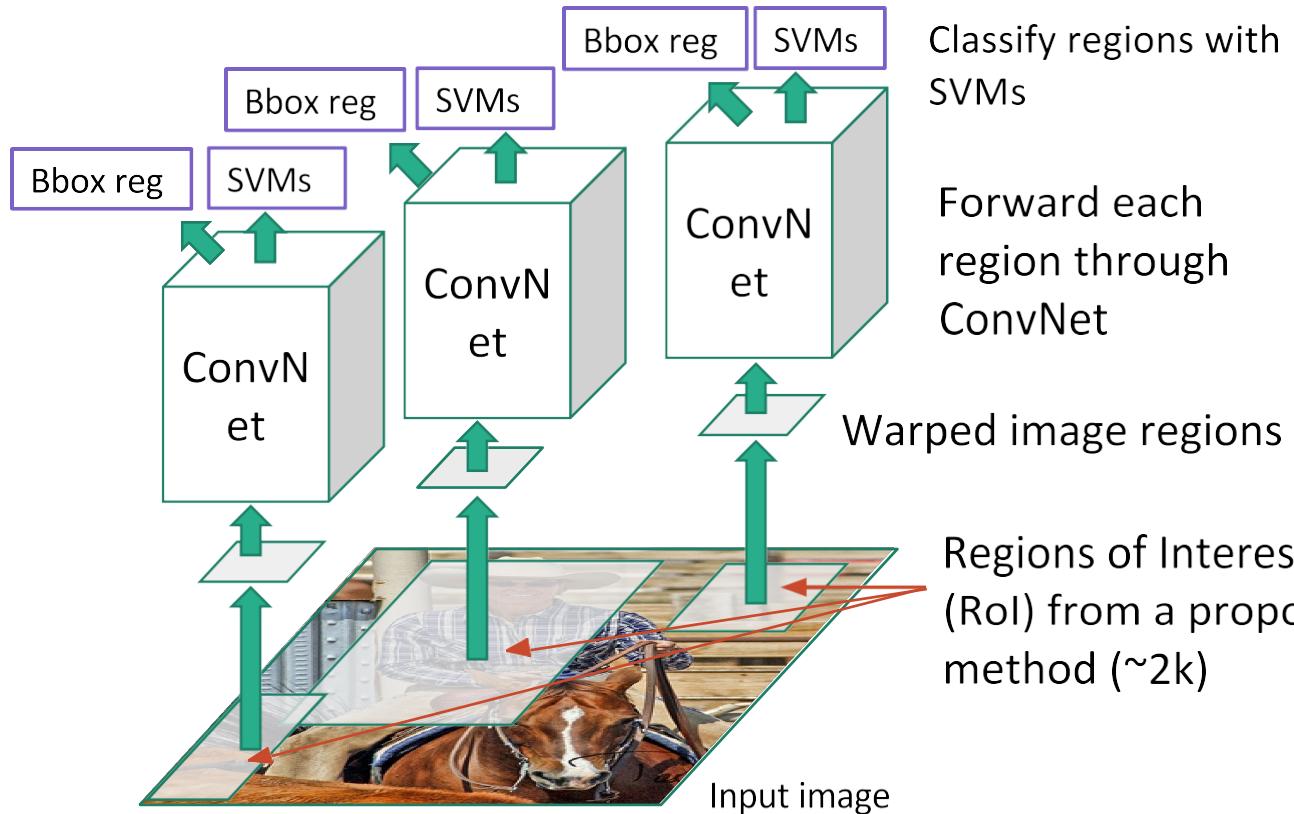
Forward each  
region through  
ConvNet

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.

# R-CNN



Linear Regression for bounding box offsets

Classify regions with  
SVMs

Forward each  
region through  
ConvNet

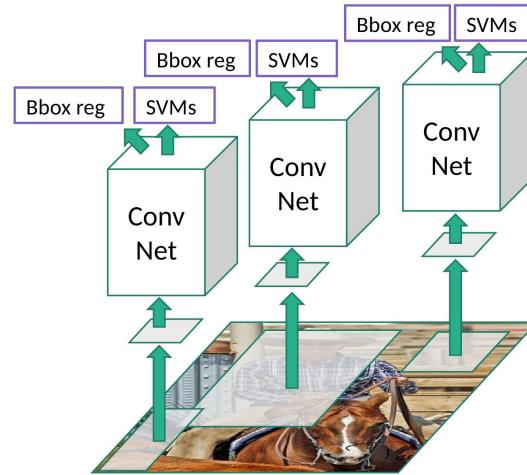
Warped image regions

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

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

# 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.

# R-CNN: Solutions?

## Base Network

VGG16  
ResNet-101  
Inception V2  
Inception V3  
Inception  
ResNet  
MobileNet

## Object Detection architecture

Faster R-CNN  
R-FCN  
SSD  
YOLO  
**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

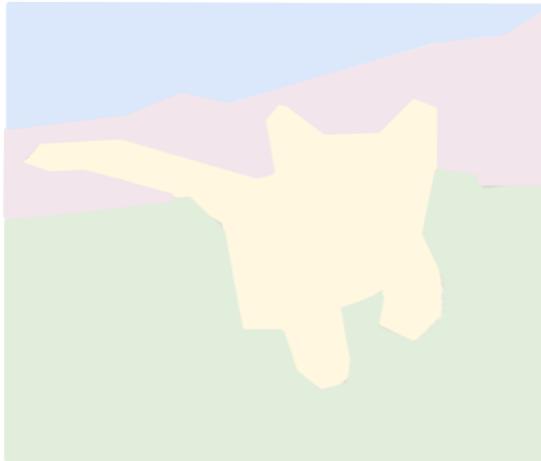
Inception-V2: Ioffe and Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, ICML 2015

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

# Instance Segmentation

Semantic  
Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

Classification  
+ Localization



CAT

Single Object

Object  
Detection



DOG, DOG, CAT

Multiple Object

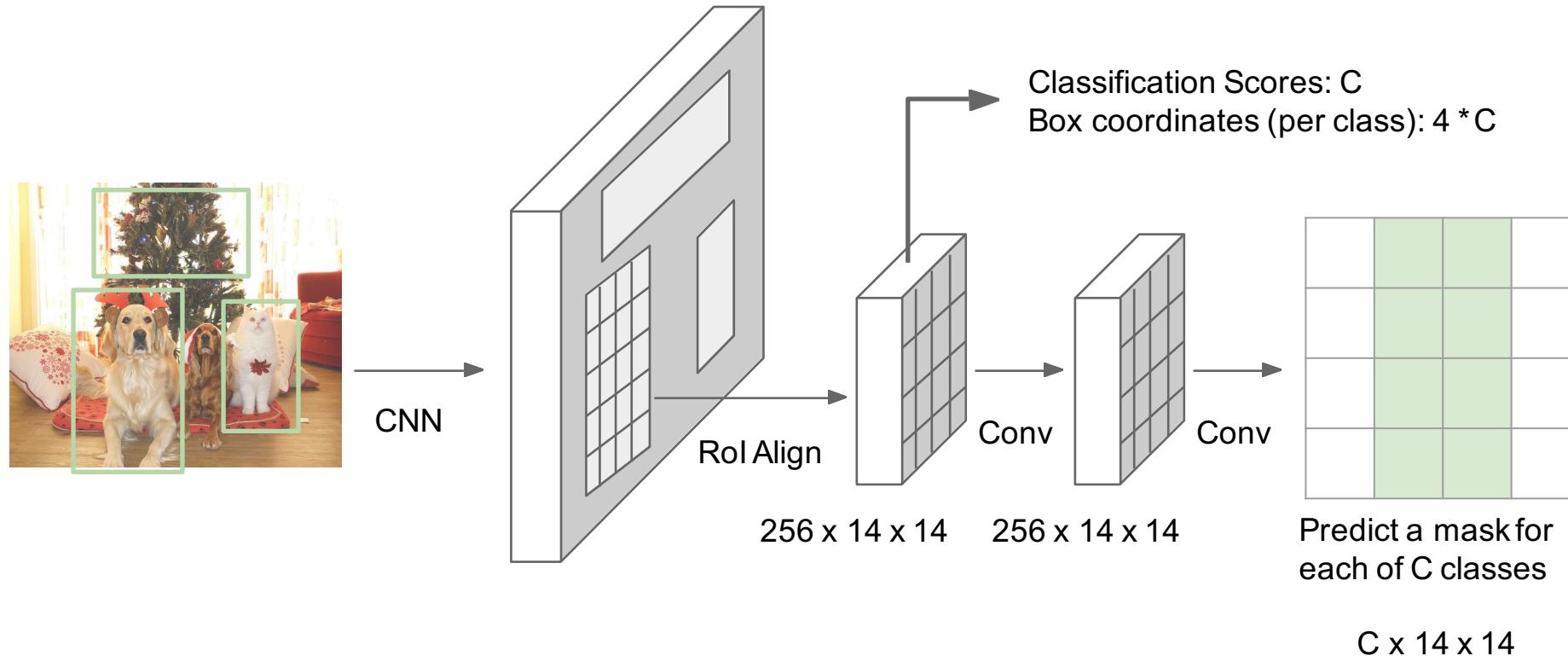
Instance  
Segmentation



DOG, DOG, CAT

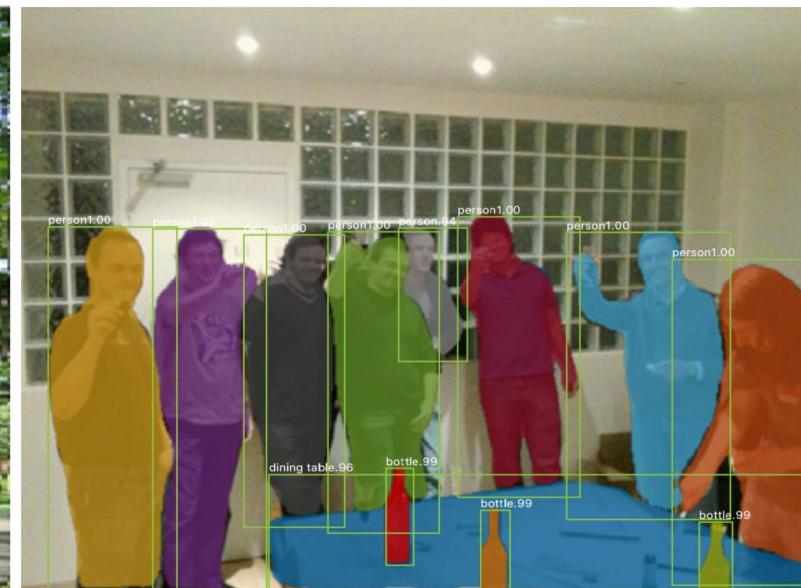
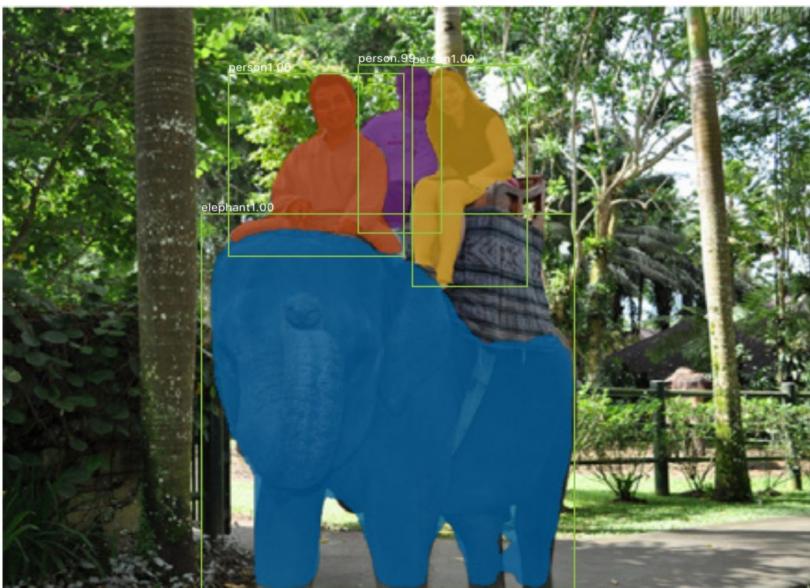
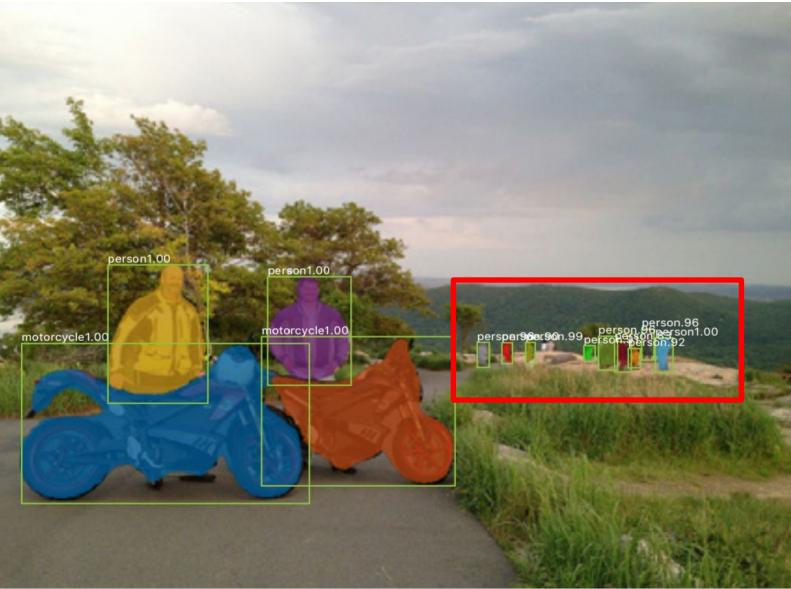
This image is CC0 public domain

# Mask R-CNN



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

# Mask R-CNN: Results



Thank you!