# CNN applications

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

School of ECE
University of Oklahoma

Spring, 2017

## Overview

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

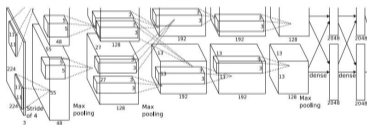# So far: Image Classification



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
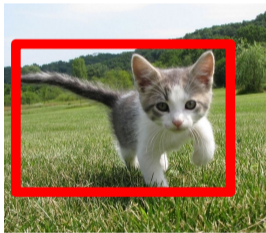Dog: 0.05
Car: 0.01
...

# Other Computer Vision Tasks



| Semantic Segmentation | Classification + Localization | Object Detection | Instance Segmentation |
|---|---|---|---|
| GRASS, CAT, TREE, SKY | CAT | DOG, DOG, CAT | DOG, DOG, CAT |
| No objects, just pixels | Single Object | Multiple Object | |

This image is CC0 public domain

# Semantic Segmentation



**GRASS**, **CAT**,
**TREE**, **SKY**

No objects, just pixels

**CAT**

Single Object

**DOG**, **DOG**, **CAT**

**DOG**, **DOG**, **CAT**

Multiple Object

This image is CC0 public domain

# Semantic Segmentation
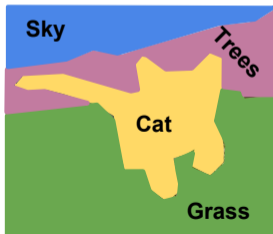


This image is CC0 public domain
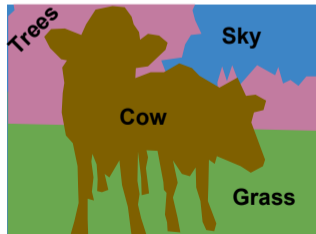
Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

# Semantic Segmentation Idea: Sliding Window



Full image

Extract patch

Classify center pixel with CNN

Cow

Cow

Grass

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: Sliding Window

Full image

Extract patch

Classify center pixel with CNN



Cow

Cow

Grass

Problem: Very inefficient! Not reusing shared features between overlapping patches
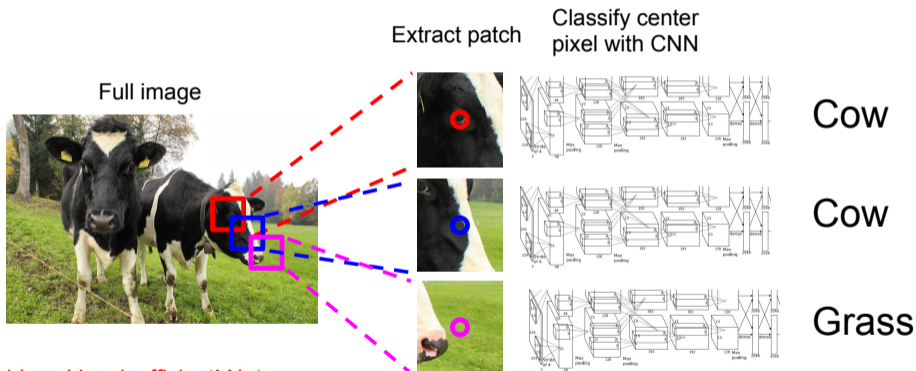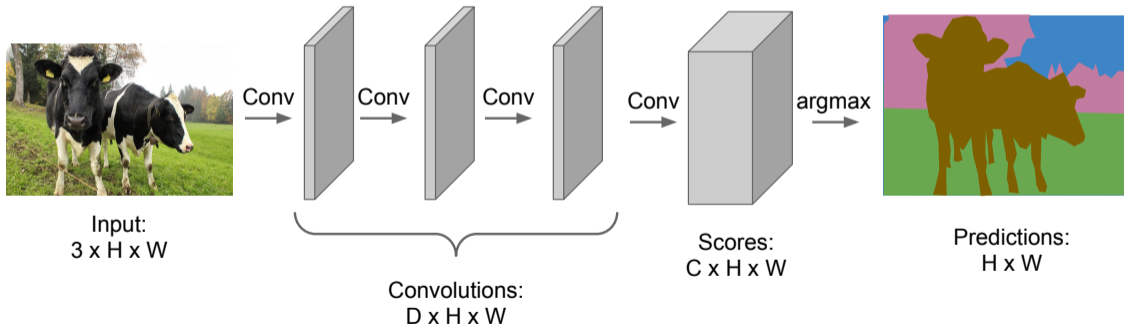
Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Semantic Segmentation Idea: Fully Convolutional

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



Input:
3 x H x W

Conv  Conv  Conv

Convolutions:
D x H x W
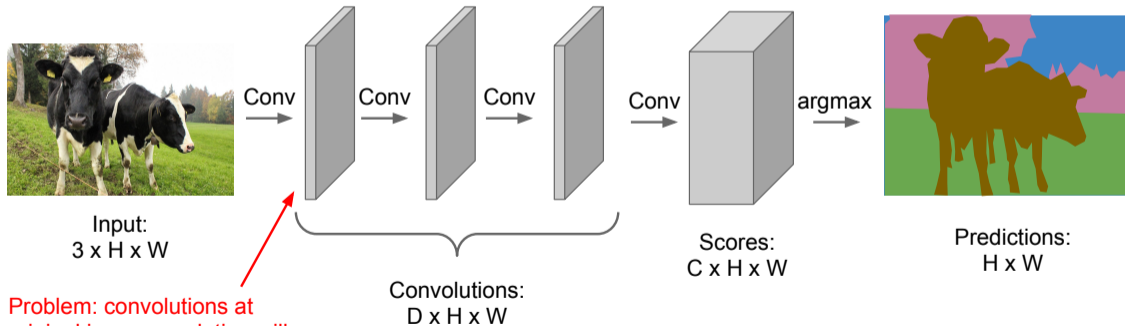
Conv

Scores:
C x H x W

argmax

Predictions:
H x W

# Semantic Segmentation Idea: Fully Convolutional

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



Input:
3 x H x W

Conv

Conv

Conv

Conv

argmax

Convolutions:
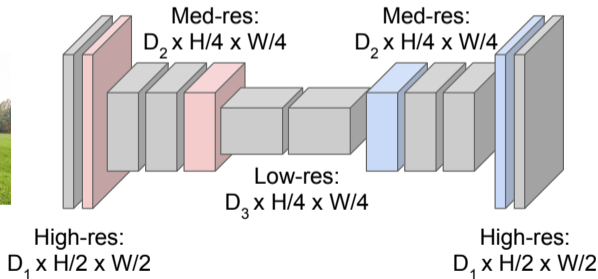D x H x W

Scores:
C x H x W

Predictions:
H x W

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

# Semantic Segmentation Idea: Fully Convolutional

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



Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

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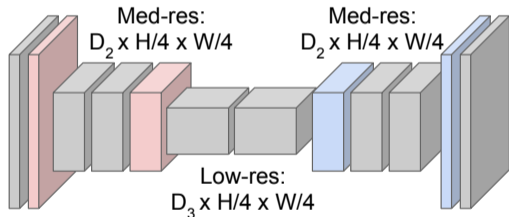Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# Semantic Segmentation Idea: Fully Convolutional

**Downsampling**:
Pooling, strided convolution

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

**Upsampling**:
???



Input:
$3 \times H \times W$

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

Med-res:
$D_2 \times H/4 \times W/4$

Low-res:
$D_3 \times H/4 \times W/4$

Med-res:
$D_2 \times H/4 \times W/4$

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

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

# In-Network upsampling: "Unpooling"

**Nearest Neighbor**

| 1 | 2 |
|---|---|
| 3 | 4 |

$\longrightarrow$

| 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 |

$\longrightarrow$

| 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"

**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
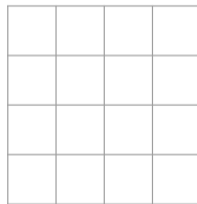upsampling layers

# Learnable Upsampling: Transpose Convolution
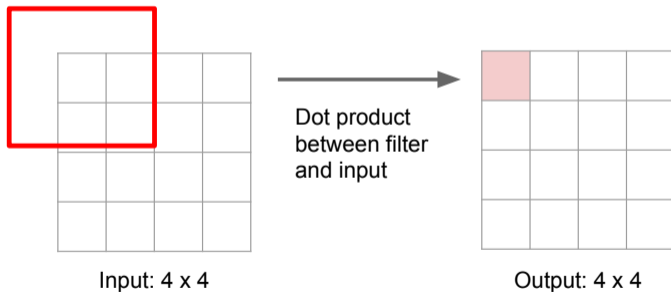
**Recall:** Typical 3 x 3 convolution, stride 1 pad 1
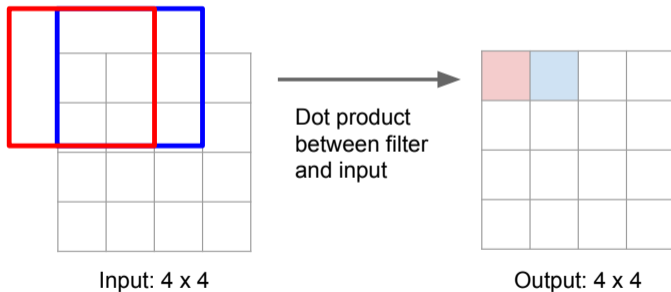


Input: 4 x 4



Output: 4 x 4

# Learnable Upsampling: Transpose Convolution
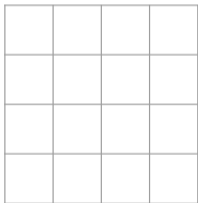
**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



Dot product between filter and input

Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

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



Dot product between filter and input

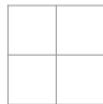Input: 4 x 4

Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1
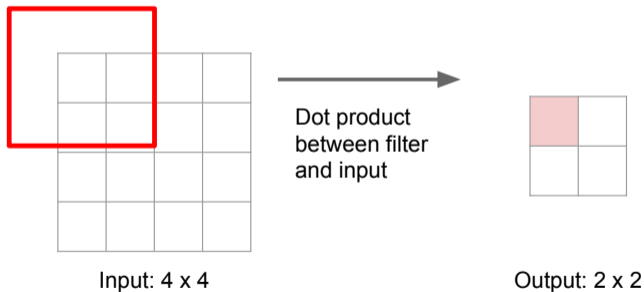
Input: 4 x 4

Output: 2 x 2

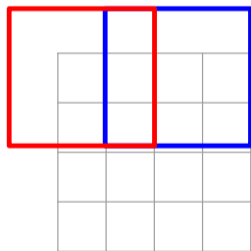# Learnable Upsampling: Transpose Convolution

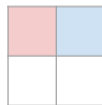**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1



Dot product
between filter
and input

Input: 4 x 4                    Output: 2 x 2

# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, <u>stride 2</u> pad 1



Dot product
between filter
and input
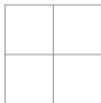
Filter moves 2 pixels in
the input for every one
pixel in the output

Stride gives ratio between
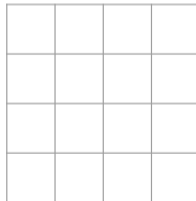movement in input and
output

Input: 4 x 4                    Output: 2 x 2

# 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



Input gives
weight for
filter

Input: 2 x 2                    Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

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

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>
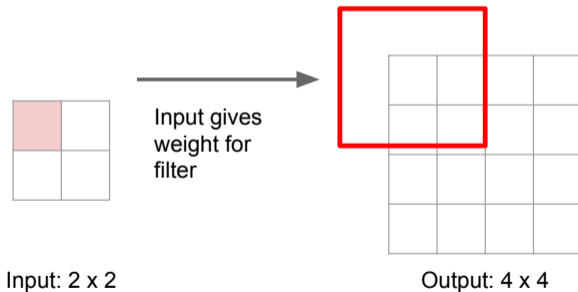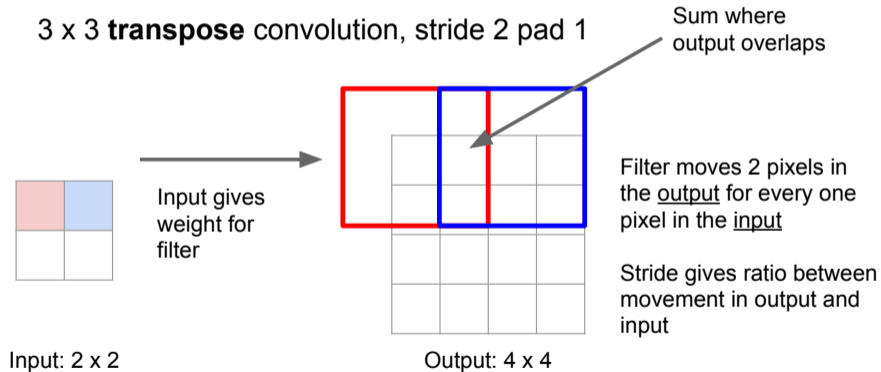
Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

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

Sum where output overlaps

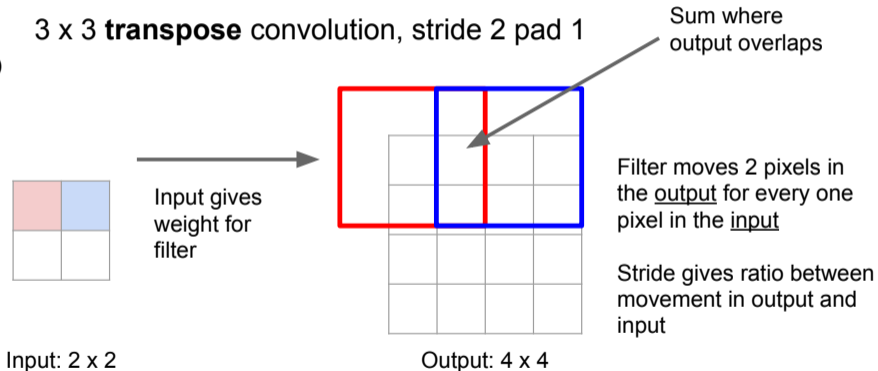Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input
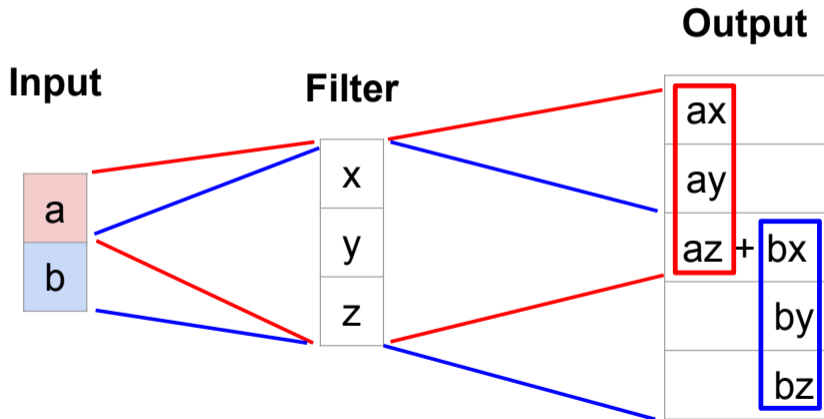
Input: 2 x 2

Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

**Other names:**
-Deconvolution (bad)
-Upconvolution
-Fractionally strided convolution
-Backward strided convolution

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

Sum where output overlaps

Input gives weight for filter

Filter moves 2 pixels in the <u>output</u> for every one pixel in the <u>input</u>

Stride gives ratio between movement in output and input

Input: 2 x 2

Output: 4 x 4

# Transpose Convolution: 1D Example

**Output**

**Input**

**Filter**



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

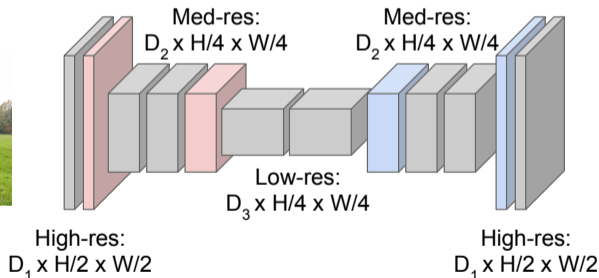Need to crop one pixel from output to make output exactly 2x input

# Semantic Segmentation Idea: Fully Convolutional

**Downsampling**:
Pooling, strided convolution

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

**Upsampling**:
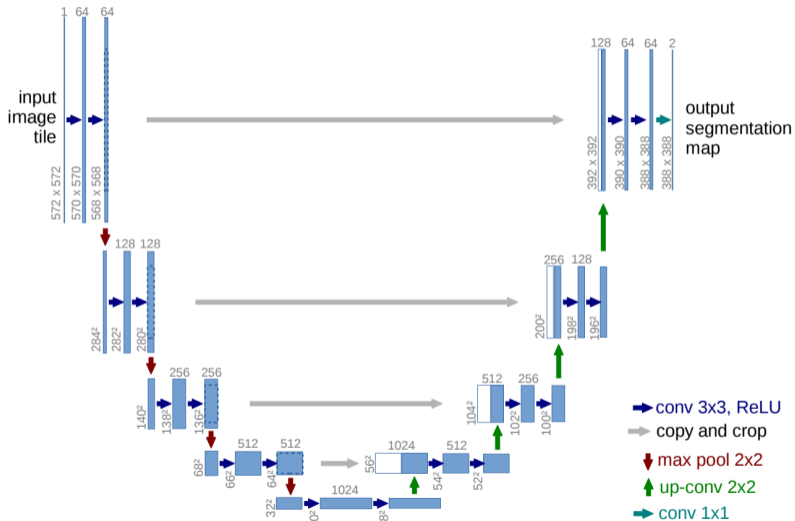Unpooling or strided transpose convolution



Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

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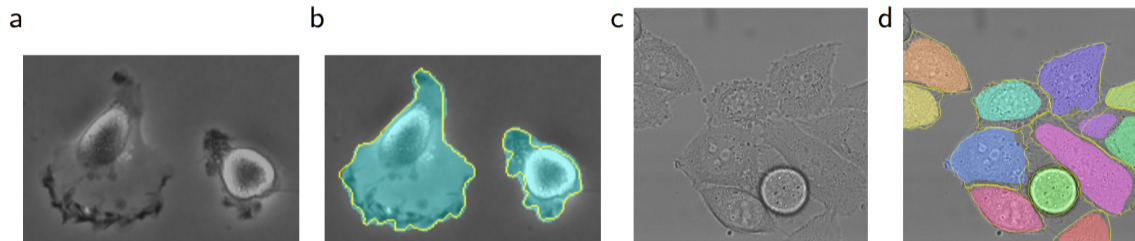
# U-Net

# U-Net



**Fig. 4.** Result on the ISBI cell tracking challenge. (**a**) part of an input image of the "PhC-U373" data set. (**b**) Segmentation result (cyan mask) with manual ground truth (yellow border) (**c**) input image of the "DIC-HeLa" data set. (**d**) Segmentation result (random colored masks) with manual ground truth (yellow border).

## Dice Coefficient

- Dice Coefficient is a similarity measure for two sets.
- Given sets A and B, the Dice Coefficient is defined as:

$$\mathrm{Dice(A, B)} = \frac{2|A \cap B|}{|A| + |B|}$$

- It ranges from 0 (no overlap) to 1 (perfect overlap).

## Dice Loss

- Dice Loss is derived from the Dice Coefficient and used as a loss function for segmentation tasks.
- The Dice Loss for predicted segmentation P and ground truth segmentation G is defined as:

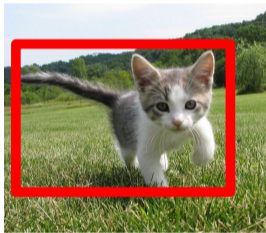$$\text{DiceLoss}(P, G) = 1 - \text{Dice}(P, G)$$

- Lower values of Dice Loss indicate better overlap between predicted and ground truth segmentations.

# Classification + Localization



**GRASS**, **CAT**,
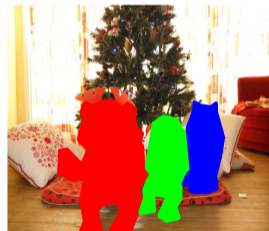**TREE**, **SKY**

No objects, just pixels

**CAT**

Single Object

**DOG**, **DOG**, **CAT**

**DOG**, **DOG**, **CAT**

Multiple Object

This image is CC0 public domain

# Classification + Localization



**Fully Connected**: 4096 to 1000

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

This image is CC0 public domain
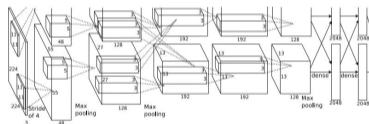
**Vector:** 4096

**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

**Correct label:**
Cat

**Fully Connected**:
4096 to 1000

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

**Softmax Loss**

**Vector:**
4096

**Fully Connected**:
4096 to 4

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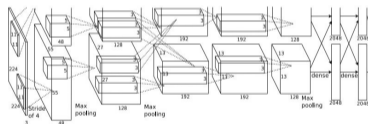(x, y, w, h)

**L2 Loss**

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

Treat localization as a regression problem!

# Classification + Localization



This image is CC0 public domain

**Fully Connected**: 4096 to 1000

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

**Correct label:**
Cat

**Softmax Loss**

Multitask Loss

**+** → **Loss**

**Vector:** 4096

**Fully Connected**: 4096 to 4

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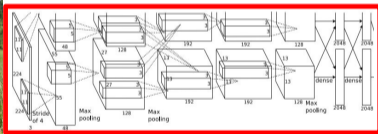(x, y, w, h)

→ **L2 Loss**

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

Treat localization as a regression problem!

# Classification + Localization

This image is CC0 public domain

Often pretrained on ImageNet
(Transfer learning)

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

**Correct label:**
Cat

**Softmax Loss**

**+** ⟶ **Loss**

**L2 Loss**

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

Treat localization as a
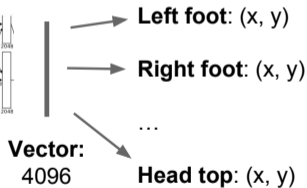regression problem!

# Aside: Human Pose Estimation



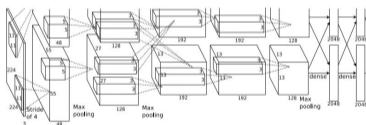This image is licensed under CC-BY 2.0.

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

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

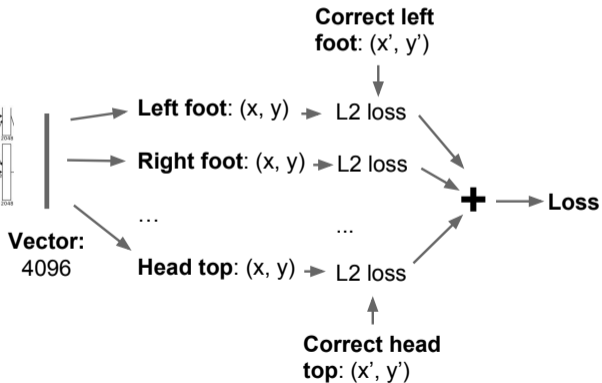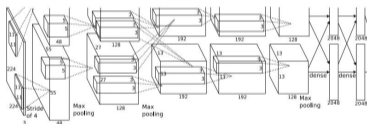# Aside: Human Pose Estimation



**Left foot**: (x, y)

**Right foot**: (x, y)

…

**Vector**: 4096

**Head top**: (x, y)

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

# Aside: Human Pose Estimation



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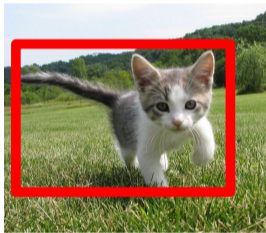
# Object Detection



**GRASS**, **CAT**, **TREE**, **SKY**

No objects, just pixels

**CAT**

Single Object

**DOG**, **DOG**, **CAT**

**DOG**, **DOG**, **CAT**

Multiple Object

# Object Detection: Impact of Deep Learning



mean Average Precision (mAP) vs year

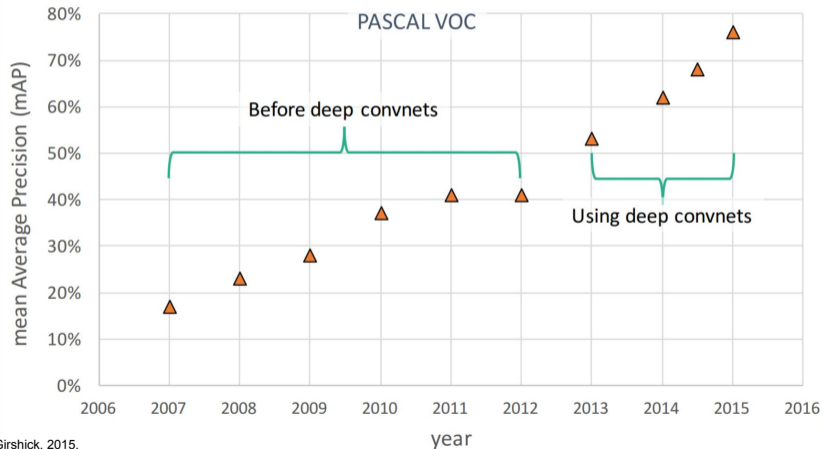PASCAL VOC

Before deep convnets
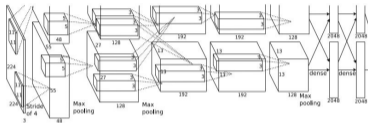
Using deep convnets

Figure copyright Ross Girshick, 2015.
Reproduced with permission.

# 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
DUCK: (x, y, w, h)  numbers!
….

# 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

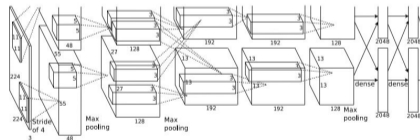Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



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
Uijlings 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

# 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.

# R-CNN



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

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.

# R-CNN



Warped image regions

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

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.

# R-CNN



Forward each region through ConvNet

Warped image regions

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

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.

# R-CNN



SVMs — Classify regions with SVMs

Forward each region through ConvNet

Warped image regions

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

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.

# 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)

Input image

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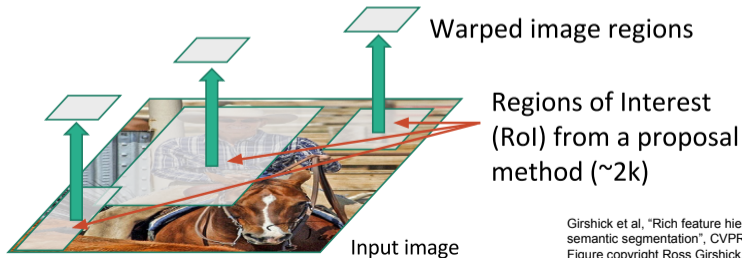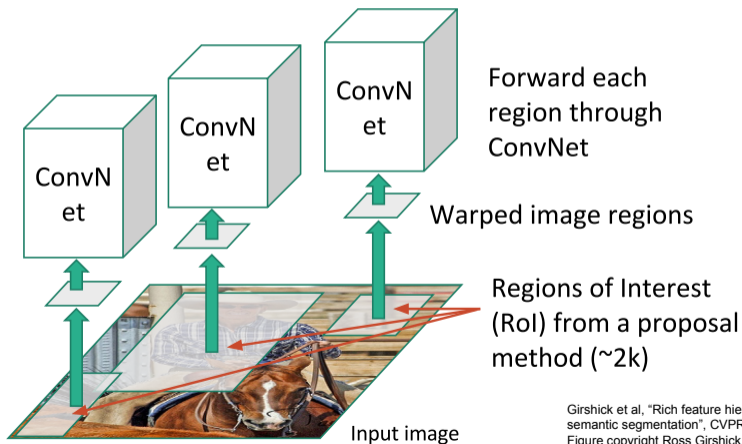Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# 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]
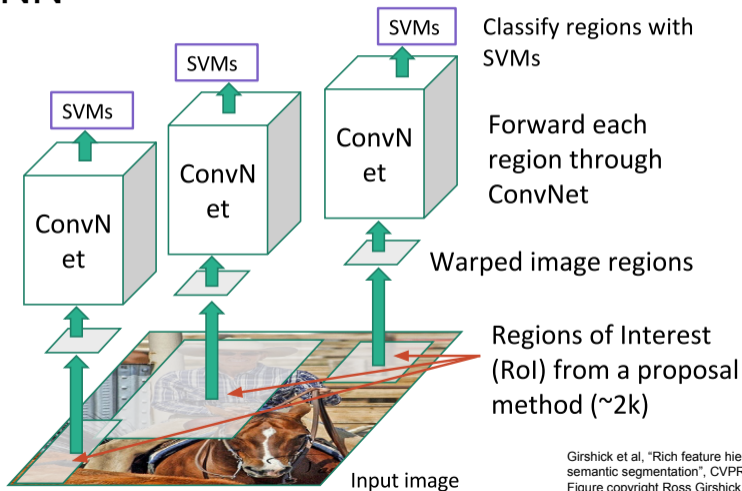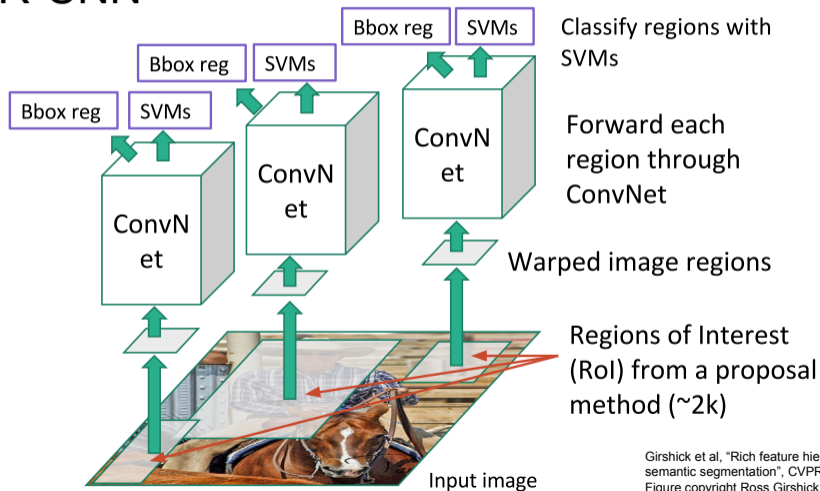  - Fixed by SPP-net [He et al. ECCV14]



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.

# Fast R-CNN



"conv5" feature map of image

Forward whole image through ConvNet

ConvNet

Input image

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

# Fast R-CNN



Regions of Interest (RoIs) from a proposal method

"conv5" feature map of image

Forward whole image through ConvNet

ConvNet

Input image

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

# Fast R-CNN



Regions of Interest (RoIs) from a proposal method

"RoI Pooling" layer

"conv5" feature map of image

Forward whole image through ConvNet

ConvNet

Input image

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

# Fast R-CNN



Softmax classifier

Linear + softmax

FCs

Fully-connected layers

"RoI Pooling" layer

"conv5" feature map of image

Regions of Interest (RoIs) from a proposal method
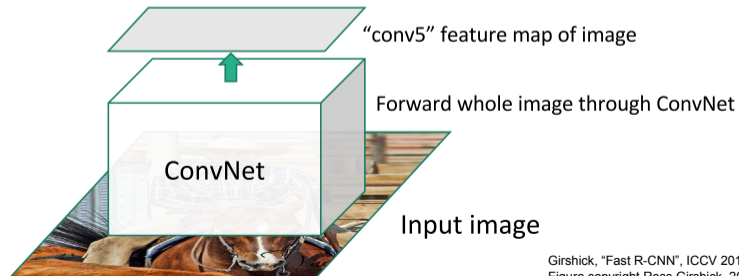
Forward whole image through ConvNet

ConvNet

Input image

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

# Fast R-CNN



Softmax classifier

Linear + softmax

Linear

Bounding-box regressors

FCs

Fully-connected layers

"RoI Pooling" layer

"conv5" feature map of image

Regions of Interest (RoIs) from a proposal method

Forward whole image through ConvNet

ConvNet

Input image

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

# Fast R-CNN (Training)



Log loss + Smooth L1 loss

Multi-task loss

Linear + softmax

Linear

FCs

ConvNet

Input image

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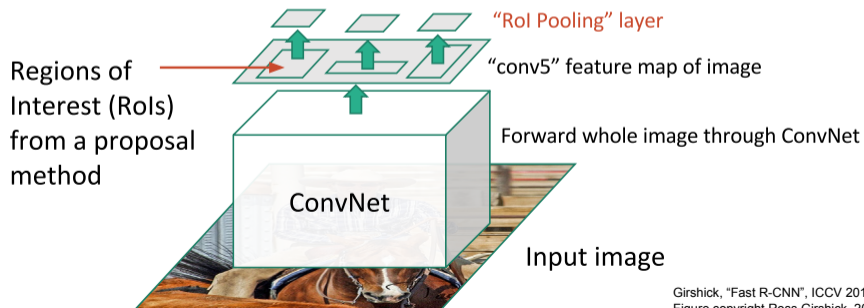Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# Fast R-CNN (Training)



Log loss + Smooth L1 loss — Multi-task loss
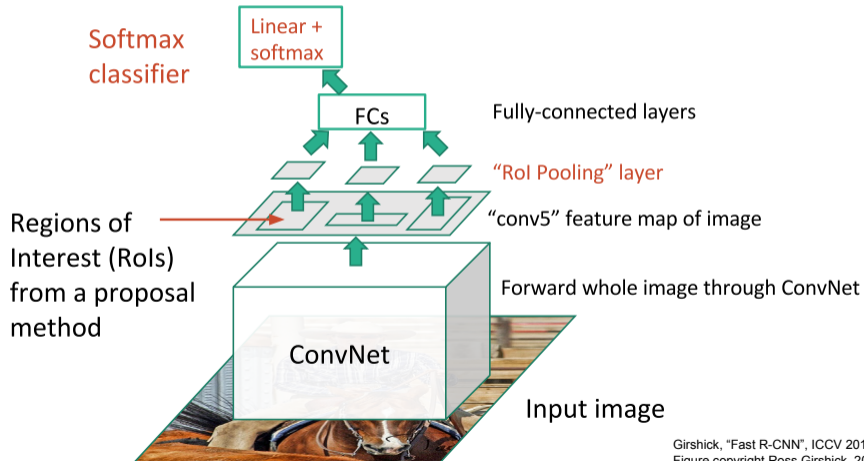
Linear + softmax

Linear

FCs

ConvNet

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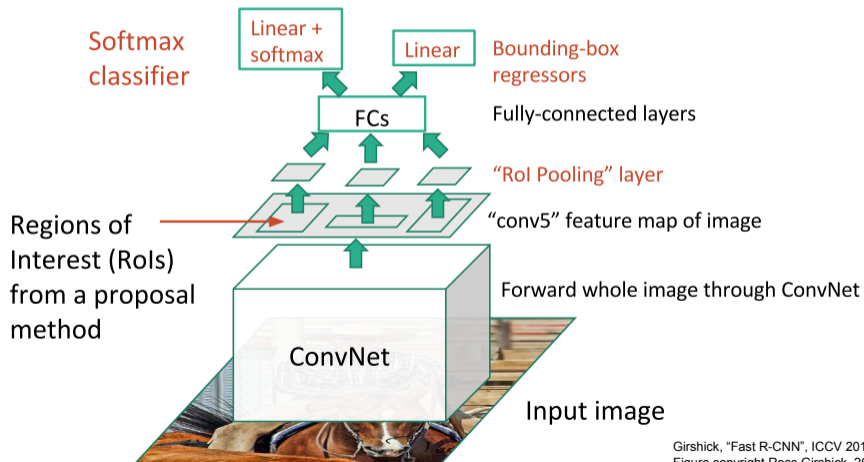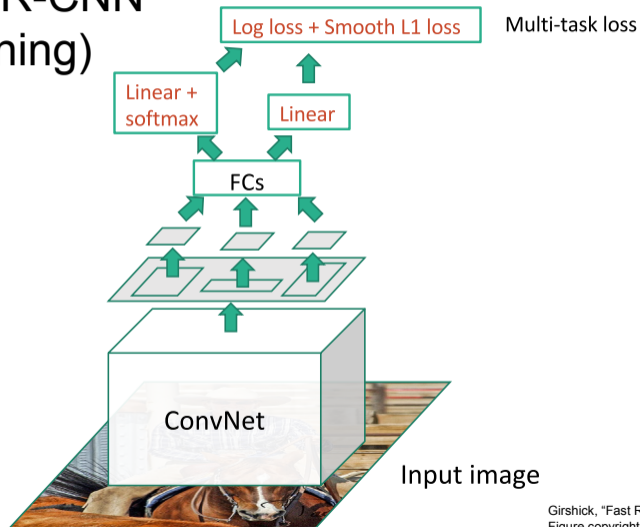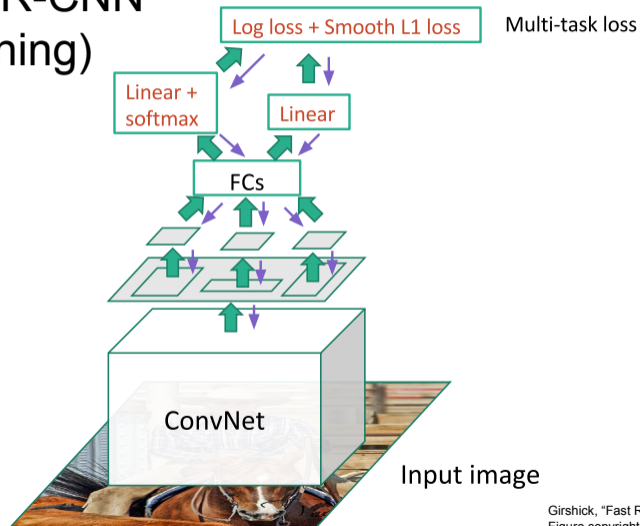
Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; source. Reproduced with permission.

# R-CNN vs SPP vs Fast R-CNN



**Training time (Hours)**

| | |
|---|---|
| R-CNN | 84 |
| SPP-Net | 25.5 |
| Fast R-CNN | 8.75 |

**Test time (seconds)**

Including Region propos...     Excluding Region Propo...

| | Including Region proposal | Excluding Region Proposal |
|---|---|---|
| R-CNN | 49 | 47 |
| SPP-Net | 4.3 | 2.3 |
| Fast R-CNN | 2.3 | 0.32 |

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

# R-CNN vs SPP vs Fast R-CNN



**Training time (Hours)**

| | |
|---|---|
| R-CNN | 84 |
| SPP-Net | 25.5 |
| Fast R-CNN | 8.75 |

**Test time (seconds)**

Including Region proposals | Excluding Region Proposals

| | Including | Excluding |
|---|---|---|
| R-CNN | 49 | 47 |
| SPP-Net | 4.3 | 2.3 |
| Fast R-CNN | 2.3 | 0.32 |

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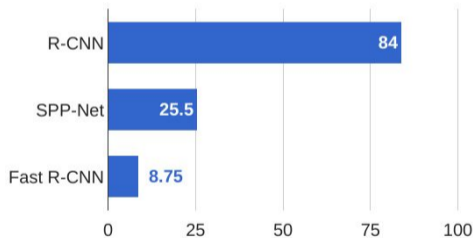by region proposals!

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
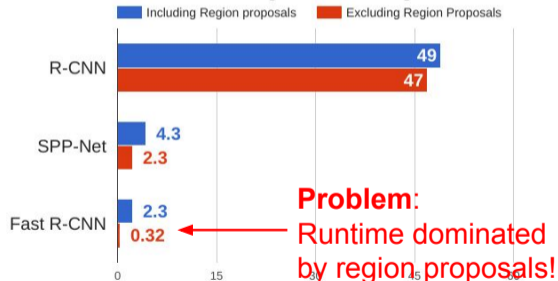Girshick, "Fast R-CNN", ICCV 2015

# 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



Classification loss

Bounding-box regression loss

RoI pooling

Classification loss

Bounding-box regression loss

proposals

Region Proposal Network

feature map

CNN

image

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

# Fast**er** R-CNN:
Make CNN do proposals!



**R-CNN Test-Time Speed**

# Still many computations are not shared for RCNN-like methods

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



Average Voting

Score Maps

Fully connected
layers are replaced by average pooling

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



$$k = 3$$

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



image and RoI

position-sensitive score maps

position-sensitive RoI-pool

vote → yes

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



image and RoI

position-sensitive score maps

position-sensitive RoI-pool

vote → no

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

Table 1: Methodologies of *region-based* detectors using **ResNet-101** [9].

| | R-CNN [7] | Faster R-CNN [19, 9] | R-FCN [ours] |
|---|---|---|---|
| depth of shared convolutional subnetwork | 0 | 91 | 101 |
| depth of RoI-wise subnetwork | 101 | 10 | **0** |

# Feature pyramid network (FPN)



(a) Featurized image pyramid

(b) Single feature map

(c) Pyramidal feature hierarchy

(d) Feature Pyramid Network

(e) Similar Structure with (d)

- a) hand-engineered features
- c) Multiscale prediction (e.g. ssd)
- e) U-Net
- b) Alexnet-like
- d) Feature pyramid network

# Feature pyramid network (FPN)

# Detection without Proposals: YOLO / SSD



Input image
3 x H x W

Divide image into grid
7 x 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 x 7 x (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 x H x W

Divide image into grid
7 x 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 x 7 x (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

# Focal loss



$$\text{CE}(p_t) = -\log(p_t)$$
$$\text{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

Figure 1. We propose a novel loss we term the *Focal Loss* that adds a factor $(1 - p_t)^\gamma$ to the standard cross entropy criterion. Setting $\gamma > 0$ reduces the relative loss for well-classified examples ($p_t > .5$), putting more focus on hard, misclassified examples. As our experiments will demonstrate, the proposed focal loss enables training highly accurate dense object detectors in the presence of vast numbers of easy background examples.

# RetinaNet



Figure 3. The one-stage **RetinaNet** network architecture uses a Feature Pyramid Network (FPN) [20] backbone on top of a feedforward ResNet architecture [16] (a) to generate a rich, multi-scale convolutional feature pyramid (b). To this backbone RetinaNet attaches two subnetworks, one for classifying anchor boxes (c) and one for regressing from anchor boxes to ground-truth object boxes (d). The network design is intentionally simple, which enables this work to focus on a novel focal loss function that eliminates the accuracy gap between our one-stage detector and state-of-the-art two-stage detectors like Faster R-CNN with FPN [20] while running at faster speeds.

# RetinaNet

| | backbone | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|
| *Two-stage methods* | | | | | | | |
| Faster R-CNN+++ [16] | ResNet-101-C4 | 34.9 | 55.7 | 37.4 | 15.6 | 38.7 | 50.9 |
| Faster R-CNN w FPN [20] | ResNet-101-FPN | 36.2 | 59.1 | 39.0 | 18.2 | 39.0 | 48.2 |
| Faster R-CNN by G-RMI [17] | Inception-ResNet-v2 [34] | 34.7 | 55.5 | 36.7 | 13.5 | 38.1 | 52.0 |
| Faster R-CNN w TDM [32] | Inception-ResNet-v2-TDM | 36.8 | 57.7 | 39.2 | 16.2 | 39.8 | **52.1** |
| *One-stage methods* | | | | | | | |
| YOLOv2 [27] | DarkNet-19 [27] | 21.6 | 44.0 | 19.2 | 5.0 | 22.4 | 35.5 |
| SSD513 [22, 9] | ResNet-101-SSD | 31.2 | 50.4 | 33.3 | 10.2 | 34.5 | 49.8 |
| DSSD513 [9] | ResNet-101-DSSD | 33.2 | 53.3 | 35.2 | 13.0 | 35.4 | 51.1 |
| **RetinaNet** (ours) | ResNet-101-FPN | 39.1 | 59.1 | 42.3 | 21.8 | 42.7 | 50.2 |
| **RetinaNet** (ours) | ResNeXt-101-FPN | **40.8** | **61.1** | **44.1** | **24.1** | **44.2** | 51.2 |

# Precision and recall

- Precision and recall are important metrics to evaluate classification models.
- They are particularly useful when the dataset is imbalanced.
  - i.e., one class has significantly more samples than another class
- These metrics give a better understanding of model performance compared to accuracy.

## Confusion Matrix

- A confusion matrix is a table that helps to visualize the performance of a classification model.
- It shows the actual and predicted classes.
- The confusion matrix consists of four elements: True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN).

# Precision

---

**Definition (Precision)**

Precision is the ratio of correctly predicted positive instances to the total predicted positive instances. It is also known as Positive Predictive Value (PPV).

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

# Recall

## Definition (Recall)

Recall is the ratio of correctly predicted positive instances to the total actual positive instances. It is also known as Sensitivity, Hit Rate, or True Positive Rate (TPR).

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

## Mean Average Precision (mAP)

- mAP is a widely used evaluation metric for object detection tasks.
- It measures both precision (how many predicted objects are actually objects) and recall (how many objects are detected by the model).
- Average precision (AP) is computed for each class and then averaged to obtain mAP.

# Intersection over Union (IoU)

- IoU is a measure of the overlap between the predicted bounding box and the ground truth bounding box.
- IoU ranges from 0 (no overlap) to 1 (perfect overlap).
- A higher IoU threshold requires tighter overlap between predicted and ground truth bounding boxes.

# mAP@0.5:0.95

- mAP@0.5:0.95 evaluates the model's performance across a range of IoU thresholds.
- It computes the AP at IoU thresholds from 0.5 to 0.95 with a step of 0.05.
- The final mAP@0.5:0.95 is the average of the AP values computed at each IoU threshold.
- This metric provides a better understanding of the model's performance at various levels of bounding box overlap.

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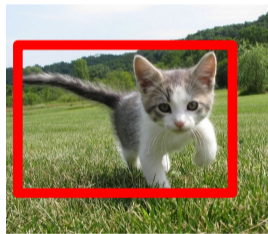MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

# Instance Segmentation



**GRASS**, **CAT**, **TREE**, **SKY**
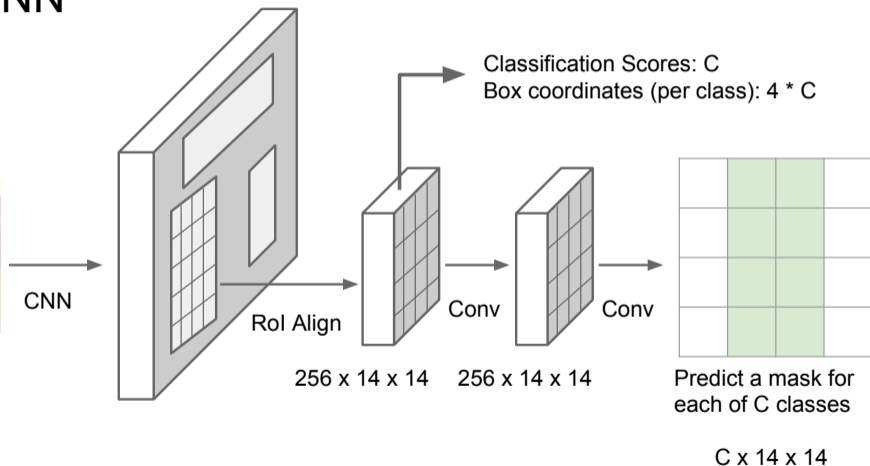
No objects, just pixels

**CAT**

Single Object

**DOG**, **DOG**, **CAT**

**DOG**, **DOG**, **CAT**

Multiple Object

This image is CC0 public domain

# Mask R-CNN



Classification Scores: C
Box coordinates (per class): 4 * C

CNN

RoI Align

256 x 14 x 14

Conv

256 x 14 x 14

Conv

Predict a mask for
each of C classes

C x 14 x 14

He et al, "Mask R-CNN", arXiv 2017
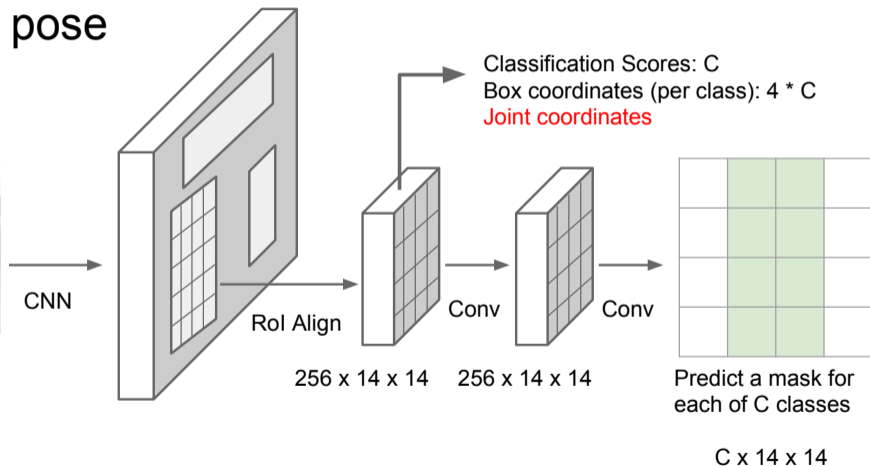
# Mask R-CNN: Very Good Results!



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

# Mask R-CNN
# Also does pose



Classification Scores: C
Box coordinates (per class): 4 * C
Joint coordinates

CNN

RoI Align

256 x 14 x 14

Conv

256 x 14 x 14

Conv

Predict a mask for
each of C classes

C x 14 x 14

He et al, "Mask R-CNN", arXiv 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.
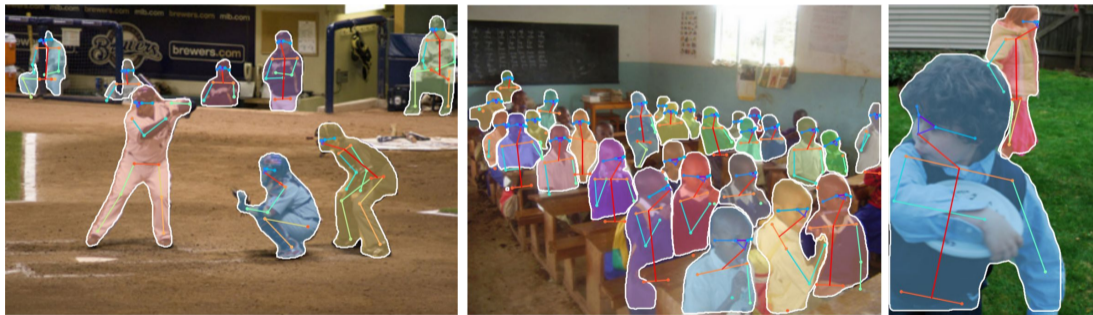Reproduced with permission.