香港中文大學
The Chinese University of Hong Kong

# CNN Applications in Computer Vision

ELEG 5491 Tutorial

Xihui Liu

# Table of Contents

- Image Representation & Pre-processing

- Object detection

- Semantic Segmentation

- Instance Segmentation

# Image Representation

- Grayscale image

  – Can be represented by 2D matrices

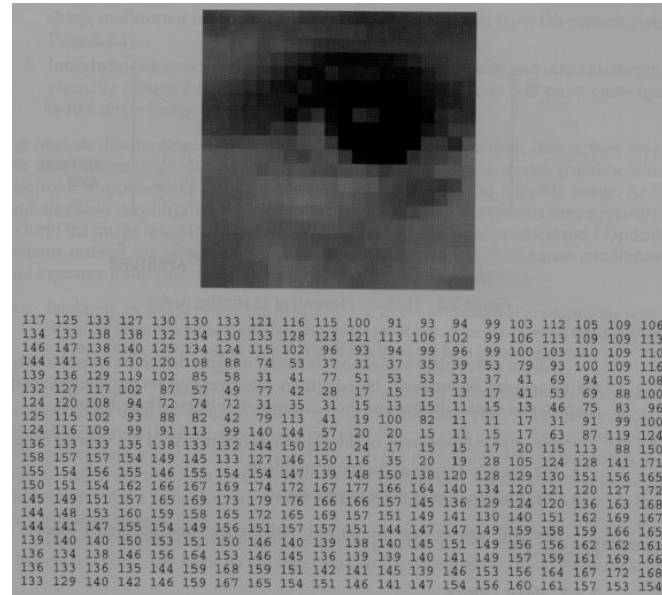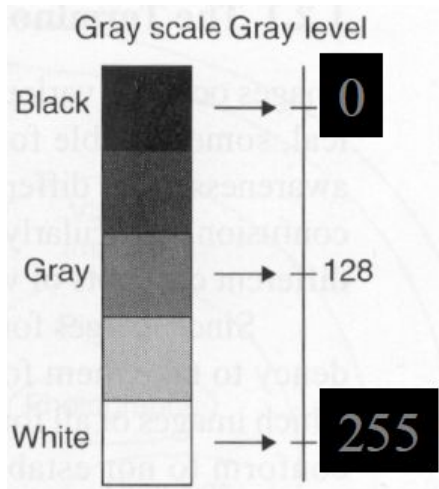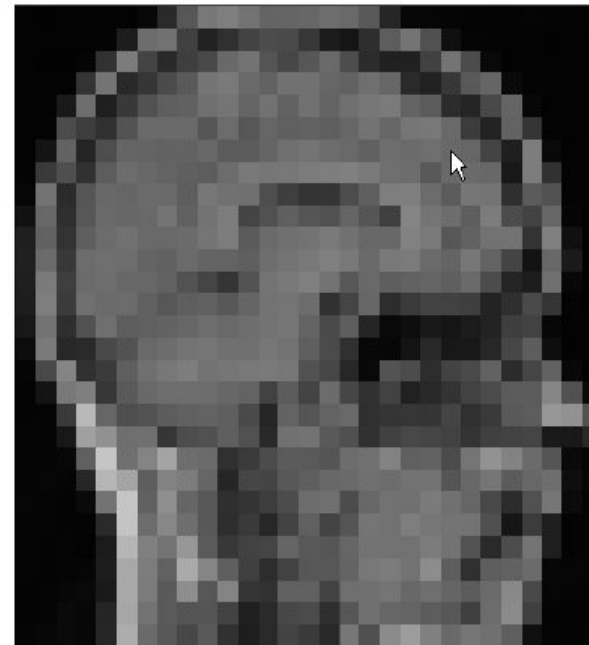  – By default, we use 8 bits per pixel

# Image Representation

- Image is a 2D array of pixels (picture element) with FIXED Number of samples : N x M

N x M = 256 x 256



N x M = 30 x 30
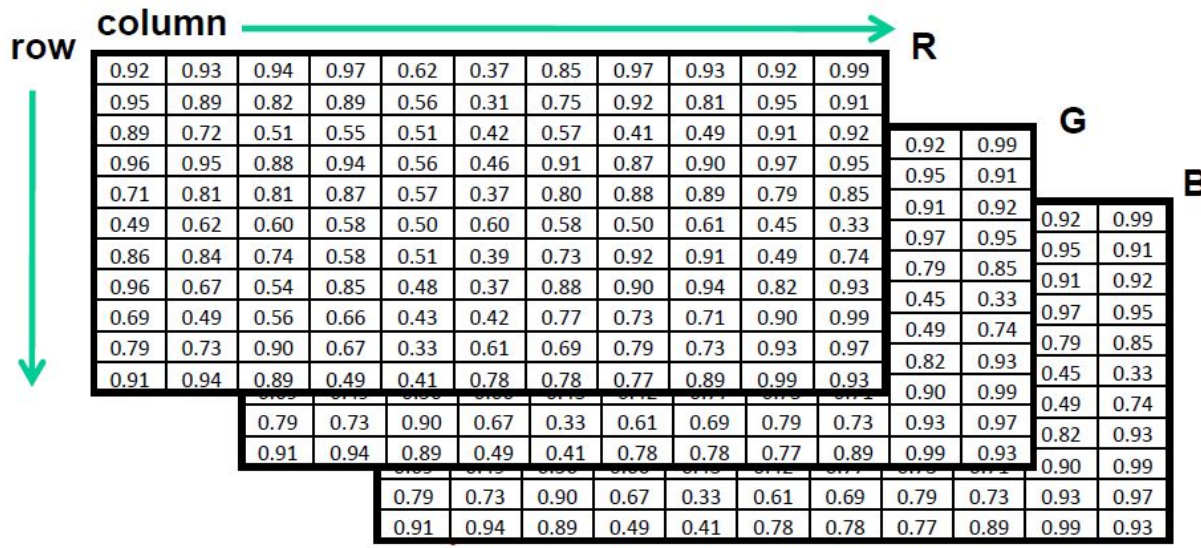
# Color Image Representation

- Color image

  – Each pixel is specified by three values, (R, G, B) in the range of [0,255] (8-bit integers)
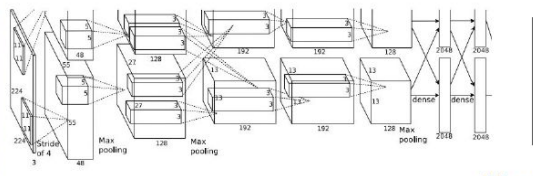


R

G

B

# Color Image Representation

- Color image

  - Color images are stored in a 3 x M x N tensor

  - [0,255] is usually mapped to [0.0,1.0] in PyTorch (a deep learning library)

# CNN Applications in Computer Vision

- Image Classification
  - Given an input image, classify it into a predefined class



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

**Fully-Connected**:
4096 to 1000

**Vector:**
4096

- Other computer vision tasks

Semantic
Segmentation



GRASS, CAT,
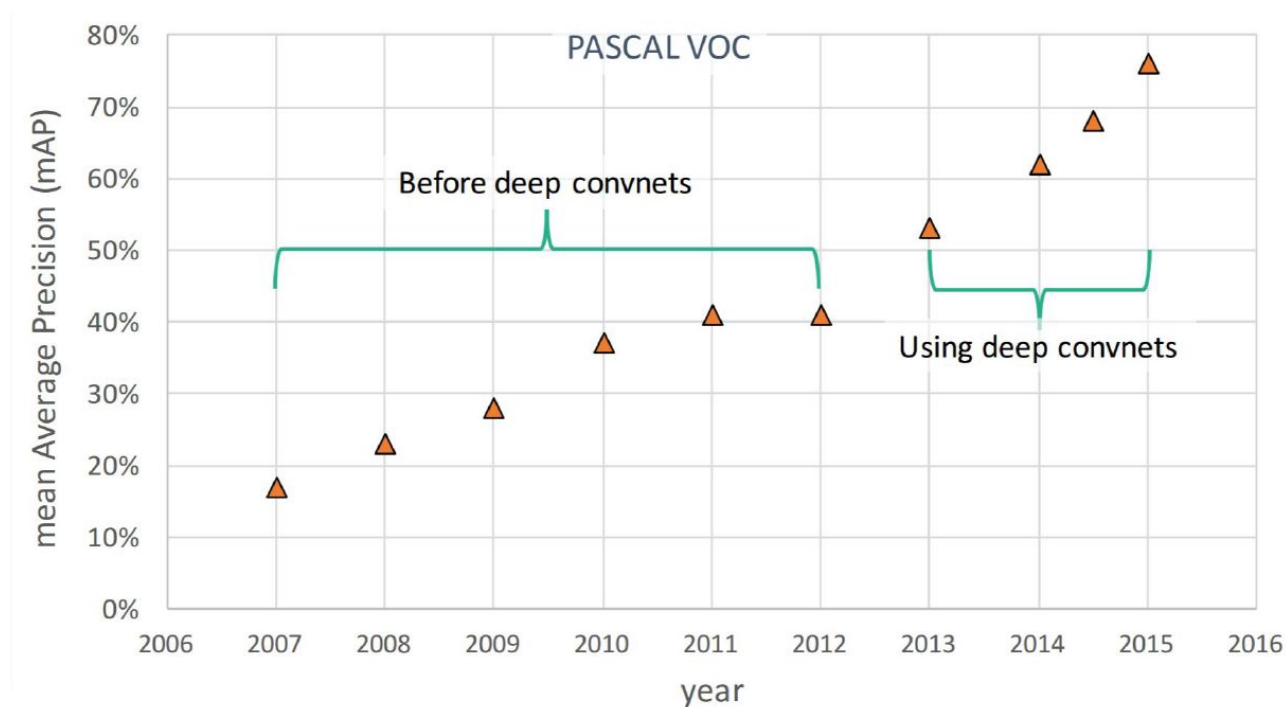TREE, SKY

Object
Detection



DOG, DOG, CAT

7

# Table of Contents

- Image Representation & Pre-processing

- Object detection

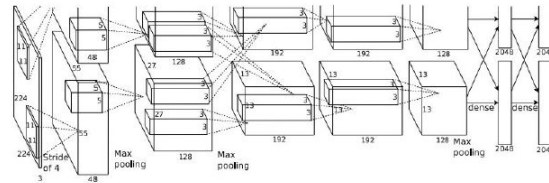- Semantic Segmentation

- Instance Segmentation

# Object Detection: Impact of Deep Learning

- PASCAL VOC is a classical object detection benchmark

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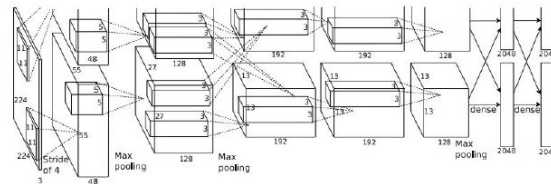
- Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

# Object Detection as Classification: Sliding Window
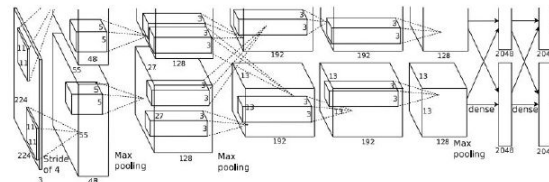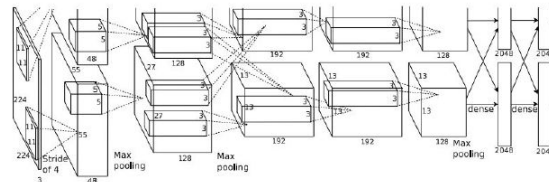
- Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

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

# Region Proposals

- Find plausible 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

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

# R-CNN: Problems

- Ad hoc training objectives
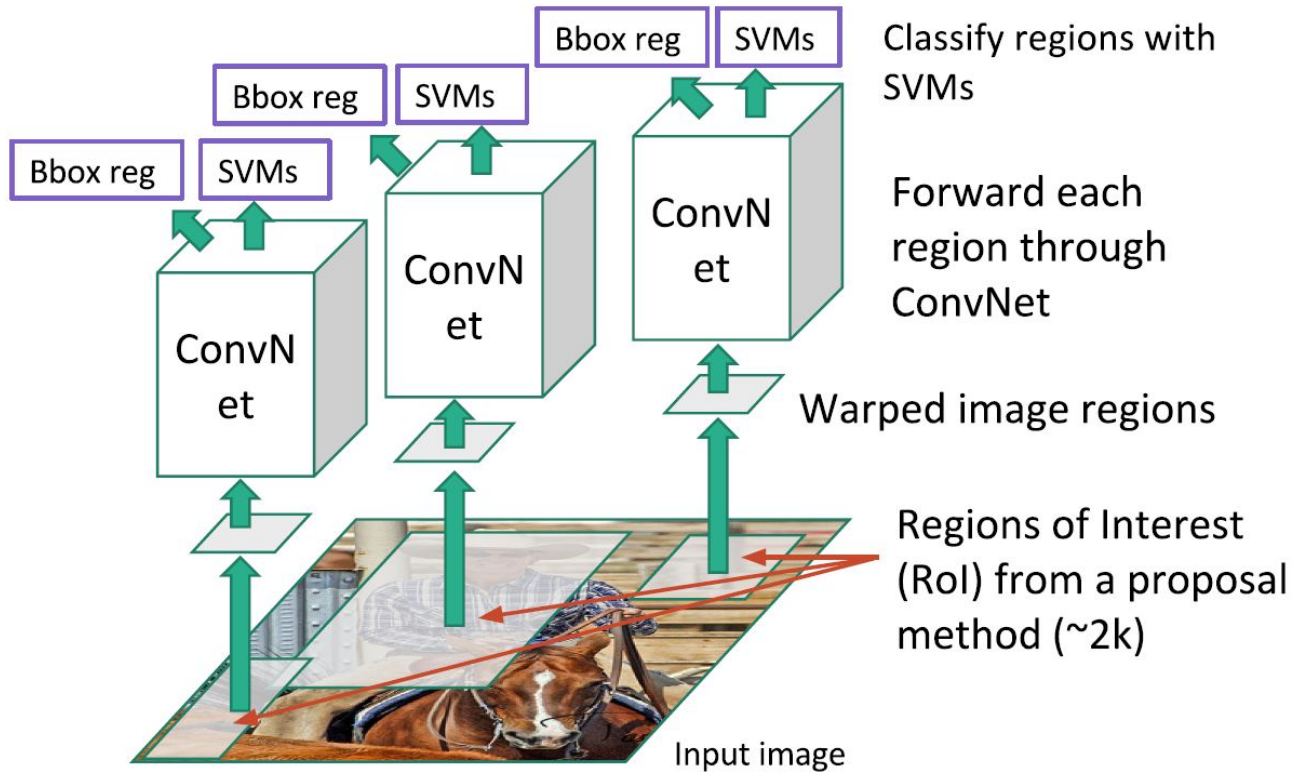
  - Fine-tune network with softmax classifier (log loss)

  - 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.

# 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

"conv5" feature map of image

ConvNet

Forward whole image through ConvNet

Input image

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

# Fast R-CNN: ROI Pooling



Project proposal onto features

Divide projected proposal into 7x7 grid, max-pool within each cell

Fully-connected layers

CNN

Hi-res input image:
3 x 640 x 480
with region
proposal

Hi-res conv features:
512 x 20 x 15;

Projected region
proposal is e.g.
512 x 18 x 8
(varies per proposal)

RoI conv features:
512 x 7 x 7
for region proposal

Fully-connected layers expect
low-res conv features:
512 x 7 x 7

Girshick et al, "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 |

0   25   50   75   100

**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 |

0   15

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

He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick et al, "Fast R-CNN", ICCV 2015.
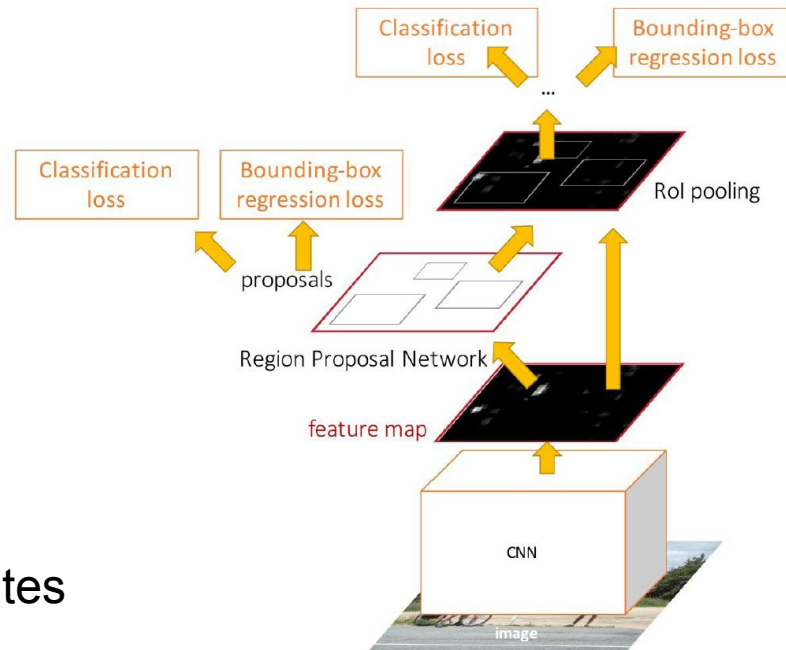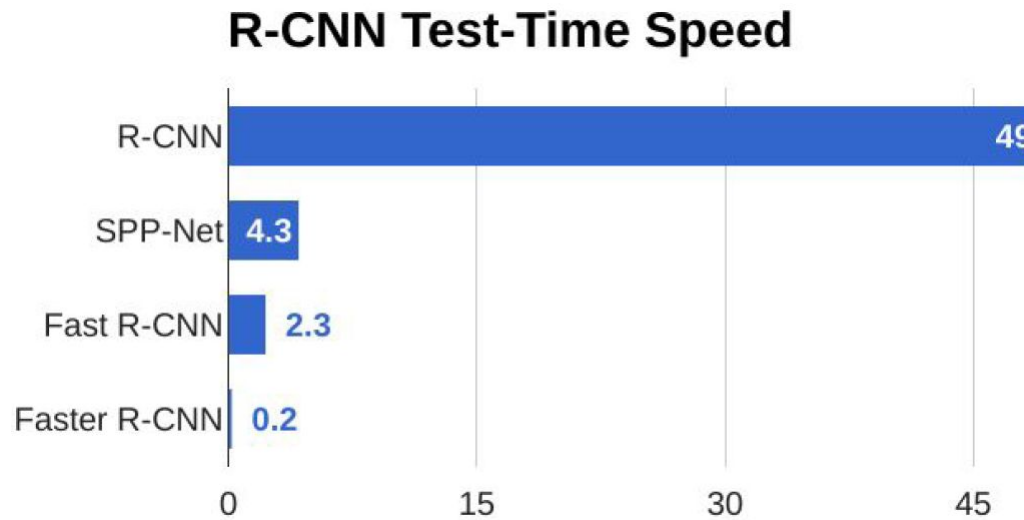
# Faster R-CNN

- Make CNN do proposals!

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

- Jointly train with 4 losses:

  – RPN classify object / not object

  – RPN regress box coordinates

  – Final classification score (object classes)

  – Final box coordinates

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

# Faster R-CNN



**R-CNN Test-Time Speed**

| Model | Speed |
|---|---|
| R-CNN | 49 |
| SPP-Net | 4.3 |
| Fast R-CNN | 2.3 |
| Faster R-CNN | 0.2 |

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

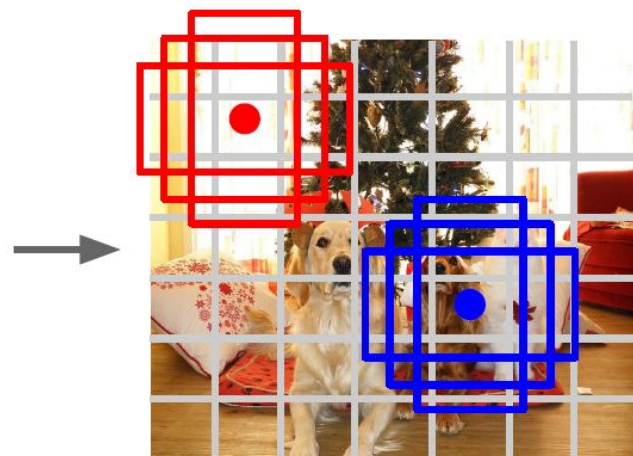# One-stage Methods without Proposals: YOLO / SSD

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

# 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
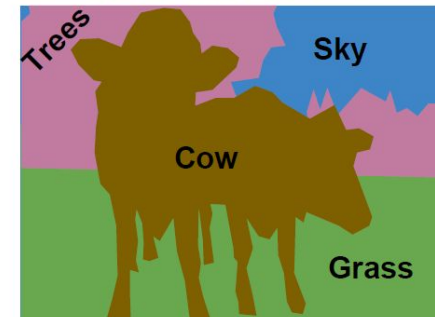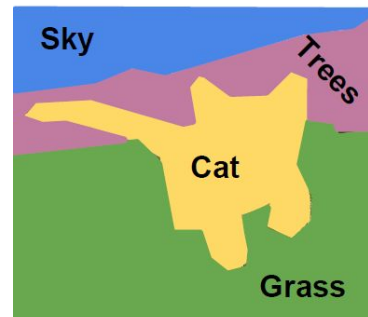
# Table of Contents

- Image Representation & Pre-processing

- Object detection

- Semantic Segmentation

- Instance Segmentation

# Semantic Segmentation

- Classical Computer Vision problem

- Label each pixel in the image with a class label

- Does not differentiate instance, only care about pixels

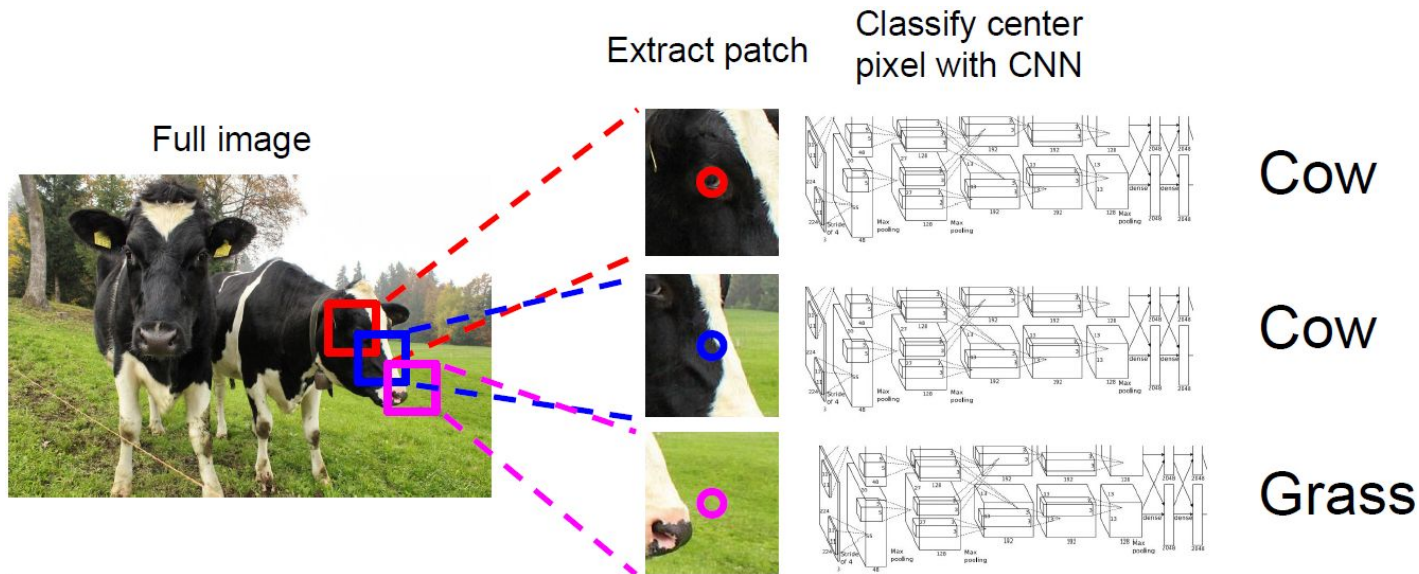# Some Public Semantic Segmentation Datasets



Pascal Visual Object Classes
20 Classes
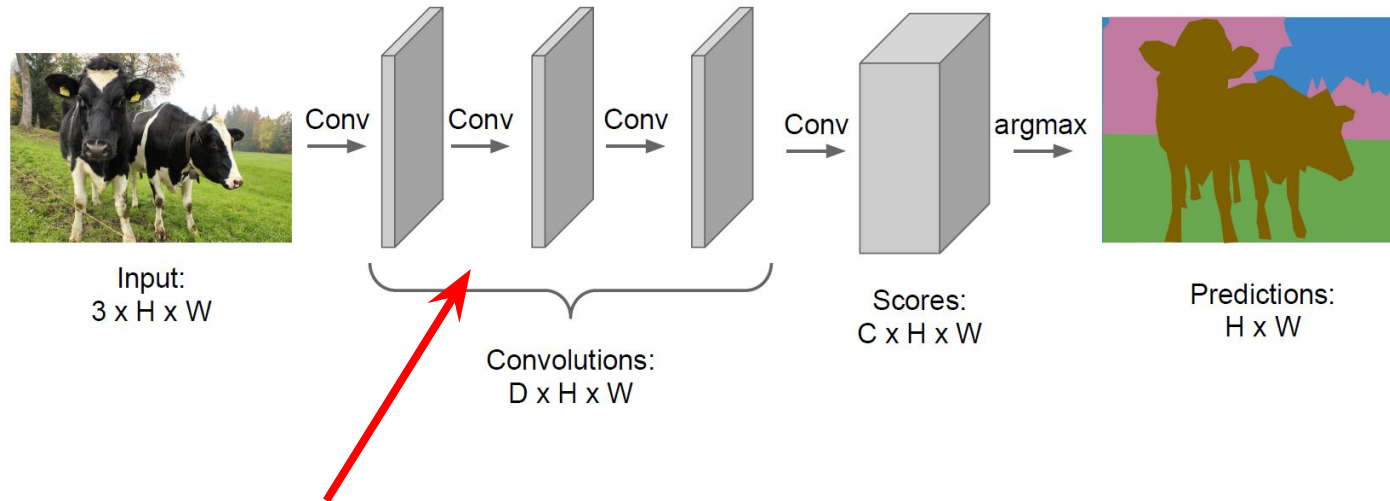~ 5.000 images



Microsoft COCO
80 Classes
~ 300.000 images

# Semantic Segmentation Idea: Sliding Window



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

# Semantic Segmentation Idea: Fully Convolutional

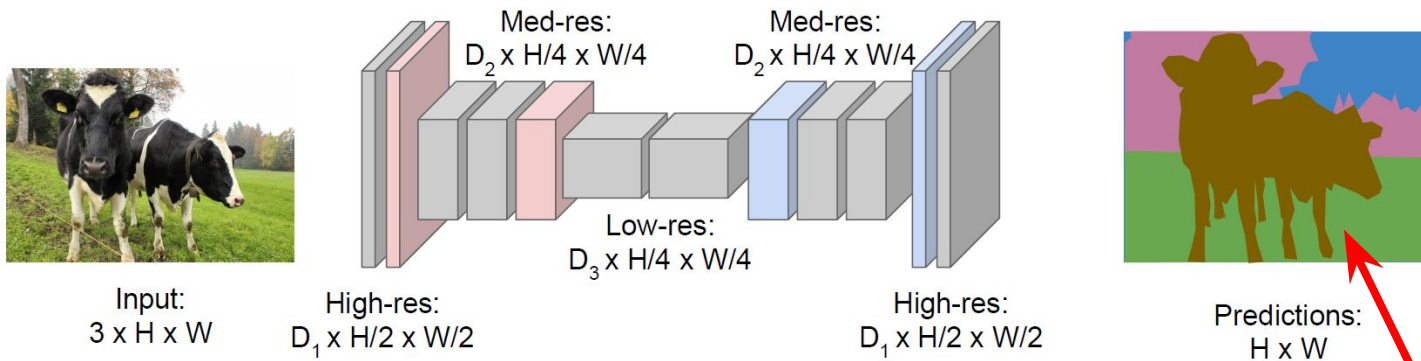Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



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

**Downsampling:**
Pooling, strided convolution

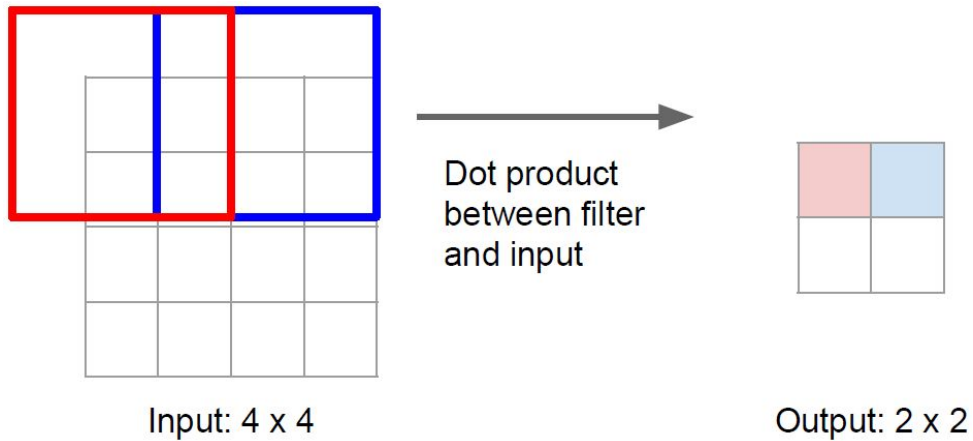Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Upsampling:**
???

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

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

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

Input:
$3 \times H \times W$

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

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

Predictions:
$H \times W$

Apply cross-entropy loss at every pixel of the predicted label map

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

# Convolution Layer

Typical 3 x 3 convolution, stride 2 pad 1



Dot product between filter and input

Input: 4 x 4

Output: 2 x 2

# "Deconvolution" Layer for Upsampling

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

Sum where output overlaps

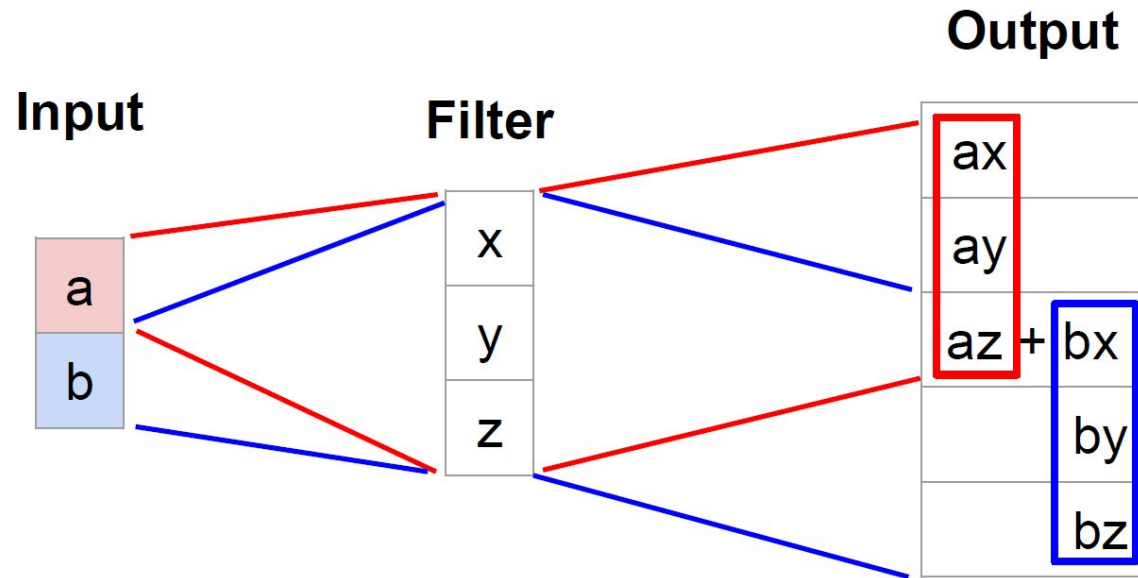Input gives weight for filter

Input: 2 x 2

Output: 4 x 4

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

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

Stride gives ratio between movement in output and input

31

# Transpose Convolution: 1D Example



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

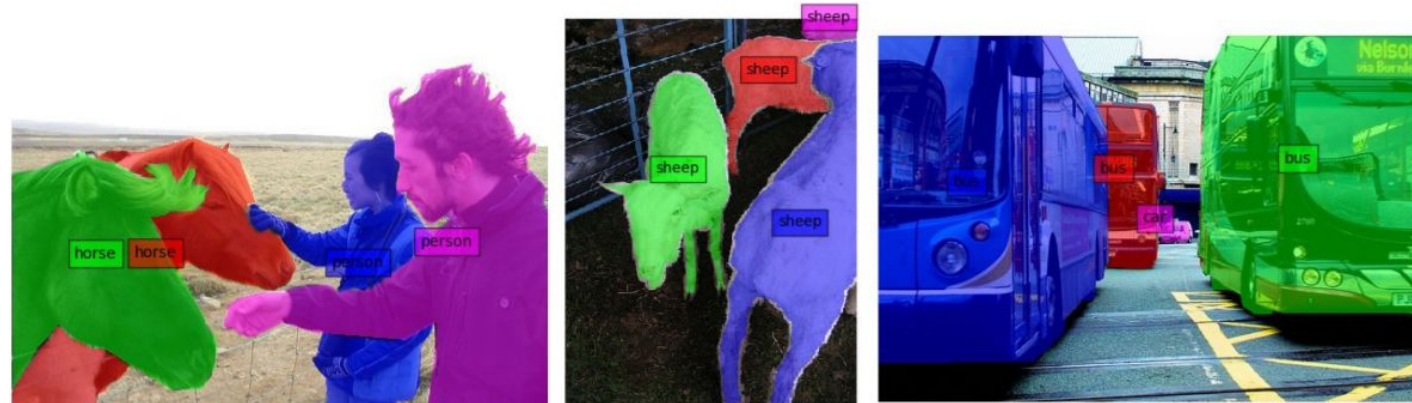Need to crop one pixel from output to make output exactly 2x input

# Table of Contents

- Image Representation & Pre-processing

- Object detection
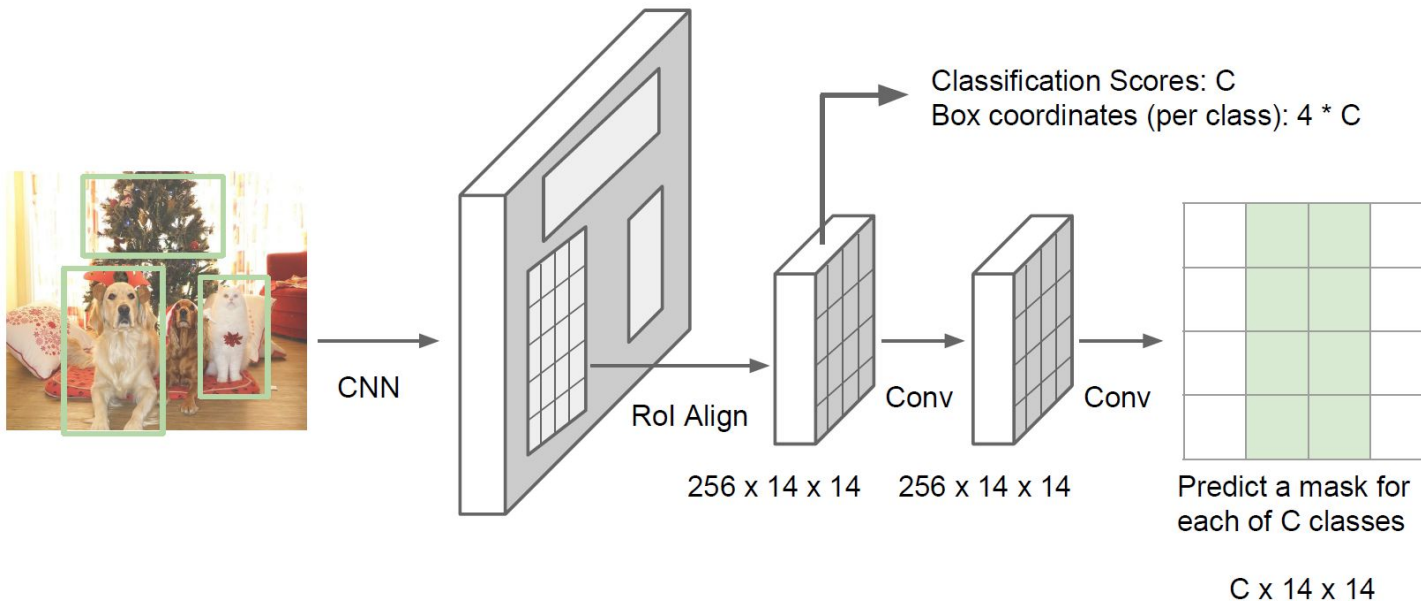
- Semantic Segmentation

- **Instance Segmentation**

# Instance Segmentation

- Not only to segment each pixel but differentiate different instances of the same class

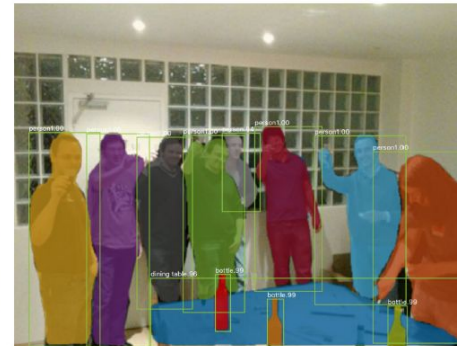- Idea: combining object detection and semantic segmentation for instance segmentation
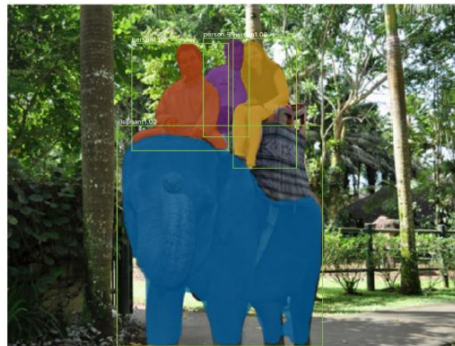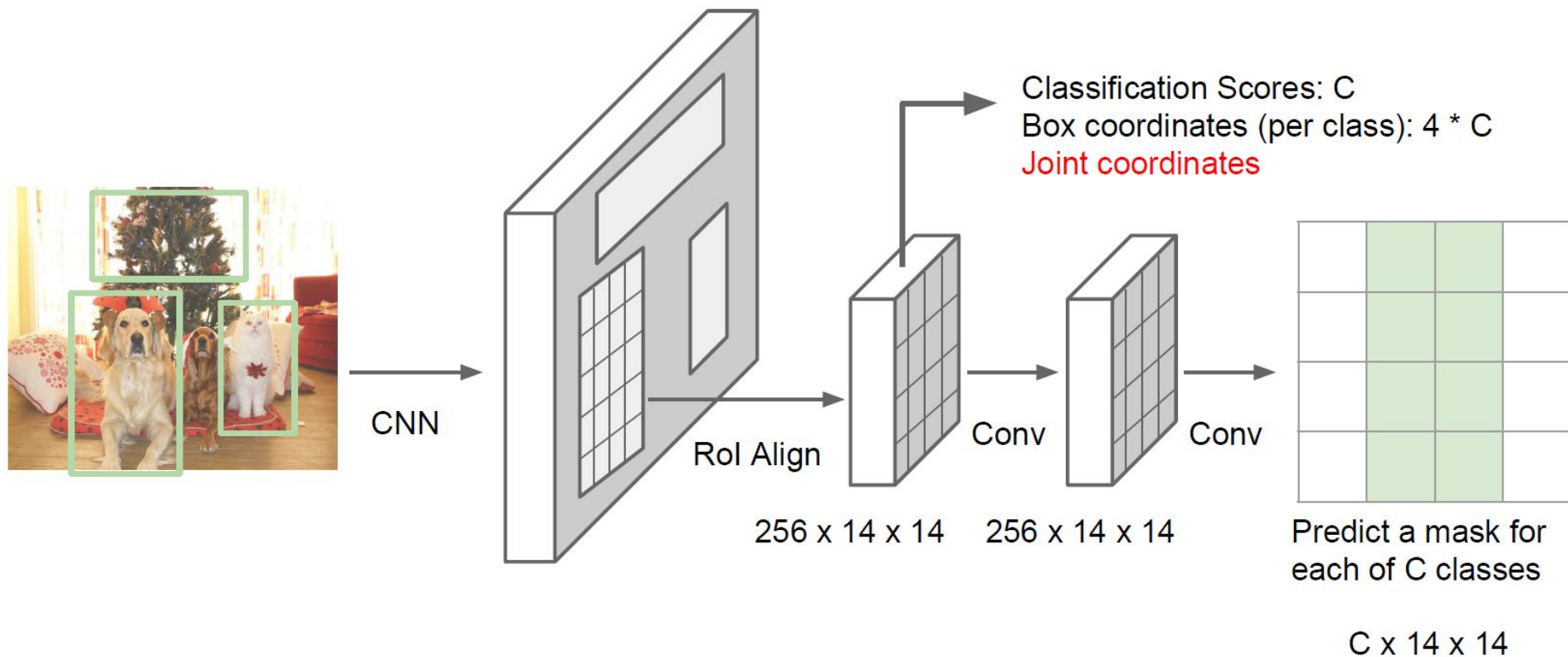
# Mask R-CNN

- **Idea:** combining object detection and semantic segmentation for instance segmentation



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", ICCV 2017

# Mask R-CNN: Very Good Results



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

# Mask R-CNN: Also Can Estimate Human Poses



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", ICCV 2017

# Mask R-CNN: Also Can Estimate Human Poses



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

# Thanks!

ELEG 5491 Tutorial
Xihui Liu