# Lecture 16: Detection + Segmentation

Justin Johnson

Lecture 16 - 1

### Reminder: A4

#### A4 due Friday October 30, 11:59pm

A4 covers:

- PyTorch autograd
- Residual networks
- Recurrent neural networks
- Attention
- Feature visualization
- Style transfer
- Adversarial examples

## Last Time: Computer Vision Tasks



### Last Time: Object Detection



### **Object Detection: Impact of Deep Learning**



Reproduced with permission.

Justin Johnson

Lecture 16 - 5

### Last Time: Object Detection Methods

**"Slow" R-CNN**: Run CNN independently for each region



# **Fast R-CNN**: Apply differentiable cropping to shared image features



#### Justin Johnson

#### Lecture 16 - 6

## Recap: Slow R-CNN Training

# **"Slow" R-CNN**: Run CNN independently for each region



# **Fast R-CNN**: Apply differentiable cropping to shared image features



#### Justin Johnson

#### Lecture 16 - 7

Input Image



#### Ground-Truth boxes

This image is CC0 public domain

#### Justin Johnson

#### Lecture 16 - 8

Input Image



Ground-Truth boxes

**Region Proposals** 

This image is CC0 public domain

Justin Johnson

Lecture 16 - 9

Input Image



Categorize each region proposal as positive, negative, or neutral based on overlap with ground-truth boxes

This image is CC0 public domain

#### Justin Johnson

#### Lecture 16 - 10

Input Image











Crop pixels from each positive and negative proposal, resize to 224 x 224

This image is CC0 public domain

#### Justin Johnson

#### Lecture 16 - 11



Run each region through CNN. For positive boxes predict class and box offset; for negative boxes just predict background class

Input Image







Class target: Dog Box target: ------





Class target: Cat Box target: ——





Class target: Dog Box target: ------



Class target: Background Box target: None

This image is CC0 public domain

#### Justin Johnson

#### Lecture 16 - 12



#### Justin Johnson

Lecture 16 - 13

# Faster R-CNN: Learnable Region Proposals

loss

Jointly train with 4 losses:

- **RPN classification**: anchor box is object / not an object
- **RPN regression**: predict transform 2. from anchor box to proposal box
- **Object classification**: classify 3. proposals as background / object class
- **Object regression**: predict transform 4. from proposal box to object box

Anchor -> Region Proposal -> Object Box

(Stage 1) (Stage 2) Training each stage looks a lot like Fast R-CNN:



#### Justin Johnson

Lecture 16 - 14

Faster R-CNN Training: RPN Training Class target: Obj Input Image Box target: — Image Features Backbone Class target: Obj CNN Box target: -Class target: Obj Box target: -----RPN gives lots of **anchors** which we classify as pos / neg / neutral **Class target: Background** Positive **GT** Boxes by matching with ground-truth Box target: None Neutral Negative

Lecture 16 - 15

October 28, 2020

This image is CC0 public domain

RPN predicts Object / Background for each anchor, as well as regresses from anchor to object box

Crop features for each proposal, use them to predict class and box targets per region

# Faster R-CNN Training: Stage 2

Input Image





Now proposals come from RPN rather than selective search, but otherwise this works the same as Fast R-CNN training

#### Image Features





**Class target: Background** Box target: None

#### Justin Johnson

#### Lecture 16 - 16

## Recap: Fast R-CNN Feature Cropping

**"Slow" R-CNN**: Run CNN independently for each region



# **Fast R-CNN**: Apply differentiable cropping to shared image features



#### Justin Johnson

#### Lecture 16 - 17



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

Justin Johnson

Lecture 16 - 18



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

Justin Johnson

#### Lecture 16 - 19



#### Justin Johnson

#### Lecture 16 - 20



#### Justin Johnson

Lecture 16 - 21

# Cropping Features: Rol <u>Align</u>



Divide into equal-sized subregions (may not be aligned to grid!)

Sample features at regularly-spaced points in each subregion using **bilinear interpolation** 

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

Justin Johnson

Lecture 16 - 22

**Project proposal** 

**CNN** 

 $f_{x,y} = \sum f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$ 

 $j \in \{ |y| - 1, \dots, [y] + 1 \}$ 

 $^{i,j} \quad i \in \{ |x| - 1, \dots, \lceil x \rceil + 1 \}$ 

onto features

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Feature  $f_{xy}$  for point (x, y) is a linear combination of features at its four neighboring grid cells:

Justin Johnson

Lecture 16 - 23

No "snapping"!



Justin Johnson

Lecture 16 - 24

**Project proposal** 

**CNN** 

 $f_{x,y} = \sum f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$ 

+  $(f_{6.6} * 0.5 * 0.8) + (f_{7.6} * 0.5 * 0.8)$ 

 $\mathbf{f}_{6.5.5.8} = (\mathbf{f}_{6,5} * \mathbf{0.5} * \mathbf{0.2}) + (\mathbf{f}_{7,5} * \mathbf{0.5} * \mathbf{0.2})$ 

onto features

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Feature f<sub>xy</sub> for point (x, y) is a linear combination of features at its four neighboring grid cells:

Justin Johnson

Lecture 16 - 25

No "snapping"!



#### Justin Johnson

#### Lecture 16 - 26

**Project proposal** 

**CNN** 

 $f_{x,y} = \sum f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$ 

+  $(f_{6.6} * 0.5 * 0.8) + (f_{7.6} * 0.5 * 0.8)$ 

 $\mathbf{f}_{6.5,5.8} = (\mathbf{f}_{6,5} * 0.5 * 0.2) + (\mathbf{f}_{7,5} * 0.5 * 0.2)$ 

onto features

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Feature  $f_{xy}$  for point (x, y) is a linear combination of features at its four neighboring grid cells:

#### Justin Johnson

#### Lecture 16 - 27

No "snapping"!



#### Justin Johnson

Lecture 16 - 28

**Project proposal** 

**CNN** 

 $f_{x,y} = \sum f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|)$ 

This is differentiable! Upstream gradient for sampled feature will

flow backward into each of the four nearest-neighbor gridpoints

onto features

Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Feature f<sub>xy</sub> for point (x, y) is a linear combination of features at its four neighboring grid cells:

#### Justin Johnson

#### Lecture 16 - 29

No "snapping"!

Justin Johnson



Divide into equal-sized subregions (may not be aligned to grid!)

> Sample features at regularly-spaced points in each subregion using **bilinear interpolation**

After sampling, maxpool in each subregion



Region features (here 512 x 2 x 2; In practice e.g 512 x 7 x 7)

Output features now aligned to input box! And we can backprop to box coordinates!

Lecture 16 - 30

### **Object Detection Methods**

Both of these rely on anchor boxes. Can we do detection without anchors?

**"Slow" R-CNN**: Run CNN independently for each region



Fast R-CNN: Apply differentiable cropping to shared image features



**Faster R-CNN**: Compute proposals with CNN **Single-Stage**: Fully convolutional detector





#### Justin Johnson

#### Lecture 16 - 31

### Detection without Anchors: CornerNet



# Represent bounding boxes by pairs of corners

Law and Deng, "CornerNet: Detecting Objects as Paired Keypoints", ECCV 2018

Justin Johnson



Upper left corners Heatmap: C x H x W Embeddings: D x H x W

Matching corners are trained to have similar embeddings



Lower right corners Heatmap: C x H x W Embeddings: D x H x W

#### October 28, 2020

Lecture 16 - 32

**Backbone CNN** 

## Detection without Anchors: Transformers



# (+) No RPN!(+) No anchors!(+) No NMS!

# (-) Slow and brittle to train(-) Worse than Faster R-CNN for small objects

Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

Justin Johnson

Lecture 16 - 33

## Computer Vision Tasks: Object Detection

Classification

Semantic Segmentation

### Object Detection

Instance Segmentation



### Computer Vision Tasks: Semantic Segmentation



### Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

Trees Sky Sky Cow Cat Grass Grass

Justin Johnson

Lecture 16 - 36

October 28, 2020

This image is CC0 public domain
### 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

#### Justin Johnson

#### Lecture 16 - 37

### 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

#### Justin Johnson

#### Lecture 16 - 38

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



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



Input:**Problem #1**: Effective receptive3 x H x Wfield size is linear in number of<br/>conv layers: With L 3x3 conv<br/>layers, receptive field is 1+2L

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Justin Johnson

Lecture 16 - 40

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



Input:Problem #1: Effective receptive3 x H x Wfield size is linear in number of<br/>conv layers: With L 3x3 conv<br/>layers, receptive field is 1+2L

Problem #2: Convolution on high res images is expensive! Recall ResNet stem aggressively downsamples

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Justin Johnson

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

#### Justin Johnson

#### Lecture 16 - 42

**Downsampling:** Pooling, strided convolution



Input: 3 x H x W Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

**Upsampling**: ???





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

#### Justin Johnson

#### Lecture 16 - 43

In-Network Upsampling: "Unpooling"

### **Bed of Nails**



C x 2 x 2 C x 4 x 4

### Justin Johnson

#### Lecture 16 - 44

In-Network Upsampling: "Unpooling"

**Bed of Nails** 

**Nearest Neighbor** 



### In-Network Upsampling: Bilinear Interpolation



Input: C x 2 x 2 Output: C x 4 x 4

 $f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$ Use two closest neighbors in x and y to construct linear approximations  $j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$ 

#### Justin Johnson

#### Lecture 16 - 46

# In-Network Upsampling: Bicubic Interpolation



Input: C x 2 x 2

Output: C x 4 x 4

Use **three** closest neighbors in x and y to construct **cubic** approximations (This is how we normally resize images!)

Justin Johnson

In-Network Upsampling: "Max Unpooling"

Max Pooling: Remember which position had the max

Max Unpooling: Place into remembered positions





Pair each downsampling layer with an upsampling layer

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Justin Johnson

Lecture 16 - 48

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

Input: 4 x 4

Output: 4 x 4

#### Justin Johnson

Lecture 16 - 49

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



Input: 4 x 4

Output: 4 x 4

Justin Johnson

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



Input: 4 x 4

Output: 4 x 4

Justin Johnson

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



Input: 4 x 4

Output: 2 x 2

Justin Johnson

Lecture 16 - 52

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



### Input: 4 x 4

Output: 2 x 2

#### Justin Johnson

Lecture 16 - 53

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



### Input: 4 x 4

Output: 2 x 2

#### Justin Johnson

Lecture 16 - 54

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



Convolution with stride > 1 is "Learnable Downsampling" Can we use stride < 1 for "Learnable Upsampling"?

Dot product between input and filter



Input: 4 x 4

Output: 2 x 2

#### Justin Johnson

Lecture 16 - 55

3 x 3 convolution transpose, stride 2



### Justin Johnson

3 x 3 convolution transpose, stride 2



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

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

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

#### Justin Johnson

### Learnable Upsampling: Transposed Convolution Sum where

3 x 3 convolution transpose, stride 2

This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output



Output: 4 x 4

Weight filter by input value and copy to output

Input: 2 x 2

### Transposed Convolution: 1D example

Filter Output Input ах ay Х a az**+**bx Y b by Ζ bz

Output has copies of filter weighted by input

Stride 2: Move 2 pixels output for each pixel in input

Sum at overlaps

Justin Johnson

Lecture 16 - 61

### Transposed Convolution: 1D example



This has many names:

- Deconvolution (bad)!
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution
- <u>Transposed Convolution</u> (best name)

#### Justin Johnson

#### Lecture 16 - 62

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

#### Justin Johnson

We can express convolution in terms of a matrix multiplication

**Transposed convolution** multiplies by the transpose of the same matrix:

$$\vec{x} * \vec{a} = X\vec{a} \qquad \qquad \vec{x} *^{T} \vec{a} = X^{T}\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 \\ 0 & x & y & z & 0 \\ 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix} \qquad \begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

Example: 1D conv, kernel When stride=1, transposed conv is just a size=3, stride=1, padding=1 regular conv (with different padding rules)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

**Transposed convolution** multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

Justin Johnson

Lecture 16 - 65

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

 $\vec{x} *^T \vec{a} = X^T \vec{a}$ 

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

When stride>1, transposed convolution cannot be expressed as normal conv

#### Justin Johnson

**Downsampling:** Pooling, strided convolution



Input: 3 x H x W Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



**Upsampling**: linterpolation, transposed conv



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

### Loss function: Per-Pixel cross-entropy

#### Justin Johnson

Lecture 16 - 67

Computer Vision Tasks

**Object Detection:** Detects individual object instances, but only gives box



**Semantic Segmentation**: Gives perpixel labels, but merges instances



#### Justin Johnson

Lecture 16 - 68

# Things and Stuff

- Things: Object categories that can be separated into object instances (e.g. cats, cars, person)
- Stuff: Object categories that cannot be separated into instances (e.g. sky, grass, water, trees)



# Computer Vision Tasks

**Object Detection:** Detects individual object instances, but only gives box (Only things!)



Semantic Segmentation: Gives perpixel labels, but merges instances (Both things and stuff)



#### Justin Johnson

Lecture 16 - 70

# Computer Vision Tasks: Instance Segmentation



# Computer Vision Tasks: Instance Segmentation

### **Instance Segmentation**:

Detect all objects in the image, and identify the pixels that belong to each object (Only things!)




## Computer Vision Tasks: Instance Segmentation

### **Instance Segmentation**:

Detect all objects in the image, and identify the pixels that belong to each object (Only things!)

Approach: Perform object detection, then predict a segmentation mask for each object!



Lecture 16 - 73



#### Justin Johnson

Lecture 16 - 74



Justin Johnson

Lecture 16 - 75



C x 28 x 28

#### October 28, 2020

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

Justin Johnson

#### Lecture 16 - 76



#### Justin Johnson

Lecture 16 - 77



#### Justin Johnson

#### Lecture 16 - 78



#### Justin Johnson

Lecture 16 - 79



#### Justin Johnson

Lecture 16 - 80

### Mask R-CNN: Very Good Results!



#### Justin Johnson

#### Lecture 16 - 81

### **Beyond Instance Segmentation**

**Instance Segmentation**: Separate object instances, but only things



Semantic Segmentation: Identify both things and stuff, but doesn't separate instances



#### Justin Johnson

Lecture 16 - 82

### Beyond Instance Segmentation: Panoptic Segmentation

Label all pixels in the image (both things and stuff)

For "thing" categories also separate into instances

Kirillov et al, "Panoptic Segmentation", CVPR 2019 Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019



Justin Johnson

Lecture 16 - 83

### Beyond Instance Segmentation: Panoptic Segmentation



Kirillov et al, "Panoptic Feature Pyramid Networks", CVPR 2019

Justin Johnson

Lecture 16 - 84

### Beyond Instance Segmentation: Human Keypoints

Represent the pose of a human by locating a set of **keypoints** 

e.g. 17 keypoints:

- Nose
- Left / Right eye
- Left / Right ear
- Left / Right shoulder
- Left / Right elbow
- Left / Right wrist
- Left / Right hip
- Left / Right knee
- Left / Right ankle



#### Person image is CC0 public domain

#### Lecture 16 - 85





Lecture 16 - 87



Lecture 16 - 88



Ground-truth has one "pixel" turned on per keypoint. Train with softmax loss

Justin Johnson

Lecture 16 - 89

### Joint Instance Segmentation and Pose Estimation



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

#### Justin Johnson

Lecture 16 - 90



Justin Johnson

Lecture 16 - 91



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

Justin Johnson

Lecture 16 - 92

### **Dense Captioning**



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

#### Justin Johnson

#### Lecture 16 - 93

### **Dense Captioning**



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

#### Justin Johnson

#### Lecture 16 - 94



Gkioxari, Malik, and Johnson, "Mesh R-CNN", ICCV 2019

Justin Johnson

Lecture 16 - 95

### 3D Shape Prediction: Mask R-CNN + Mesh Head

#### Mask R-CNN: 2D Image -> 2D shapes



Mesh R-CNN: 2D Image -> **3D** shapes

chair

# bookcase chair More details next time!

Gkioxari, Malik, and Johnson, "Mesh R-CNN", ICCV 2019

He, Gkioxari, Dollár, and Girshick, "Mask R-CNN", ICCV 2017

#### Justin Johnson

#### Lecture 16 - 96

### Summary: Many Computer Vision Tasks!







## Next Time: 3D Vision

Justin Johnson

Lecture 16 - 98