# Object Detection and Semantic Segmentation

**Danna Gurari**

University of Colorado Boulder

Fall 2022

# Review

- Last lecture:
  - Representation learning
  - Pretrained features
  - Fine-tuning
  - Training neural networks: hardware & software
  - Programming tutorial

- Assignments (Canvas)
  - Lab assignment 2 due Wednesday

- Questions?

# Today's Topics

- Problems

- Applications

- PASCAL VOC detection challenge: R-CNNs

- PASCAL VOC semantic segmentation challenge: fully convolutional networks

# Today's Topics

- **Problems**

- Applications

- PASCAL VOC detection challenge: R-CNNs

- PASCAL VOC semantic segmentation challenge: fully convolutional networks

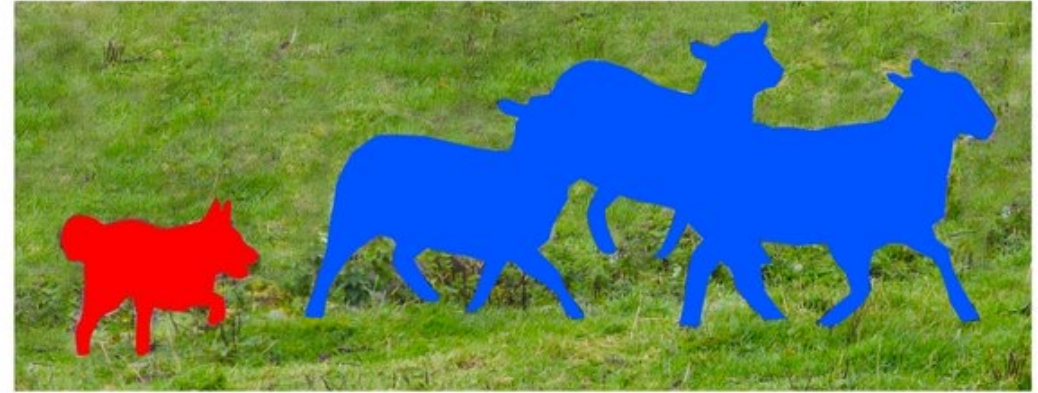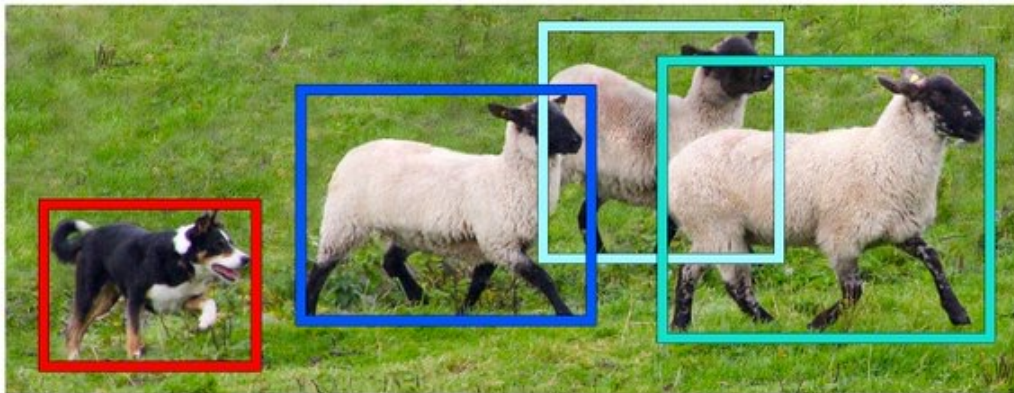# Recall: Image Classification Task

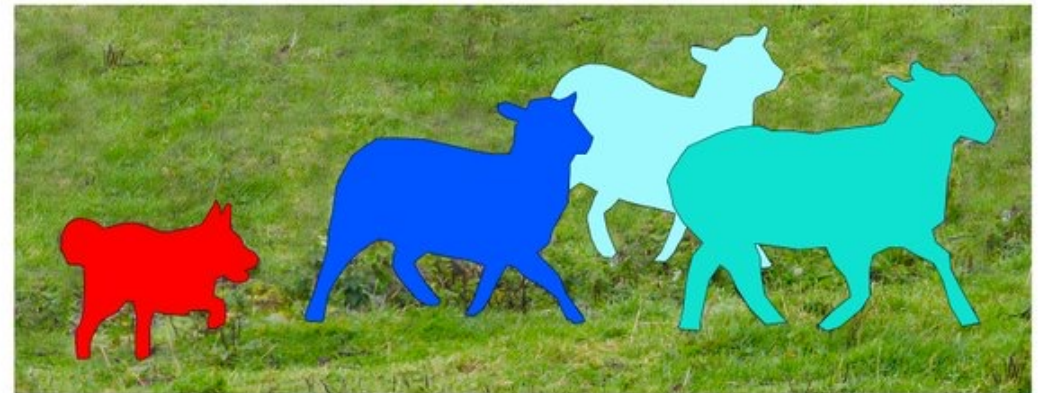# Today's Scope: Localize Content of Interest (Segmentation and Detection)



https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works
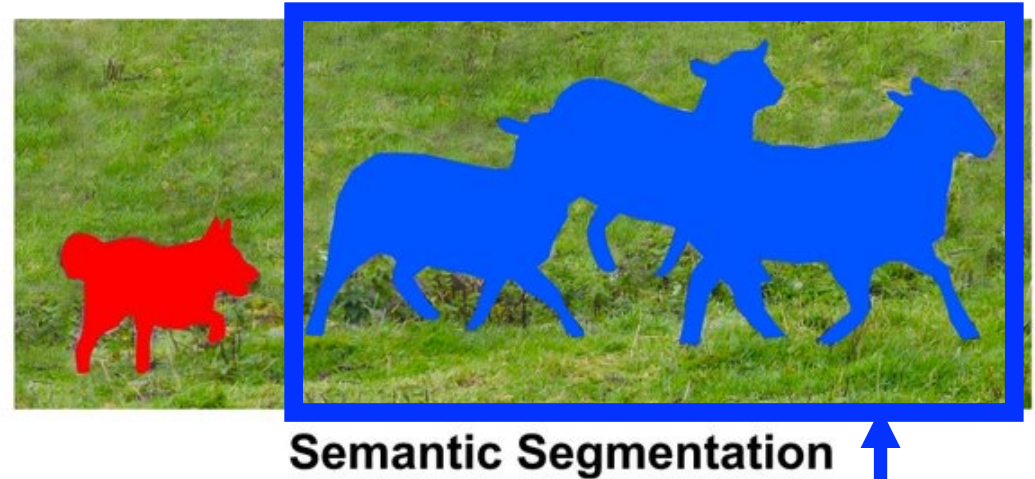
# Today's Scope: Localize Content of Interest (Segmentation and Detection)

Locate all pixels that belong to pre-specified categories



**Semantic Segmentation**

Note: instances of the same class are NOT separated

https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works

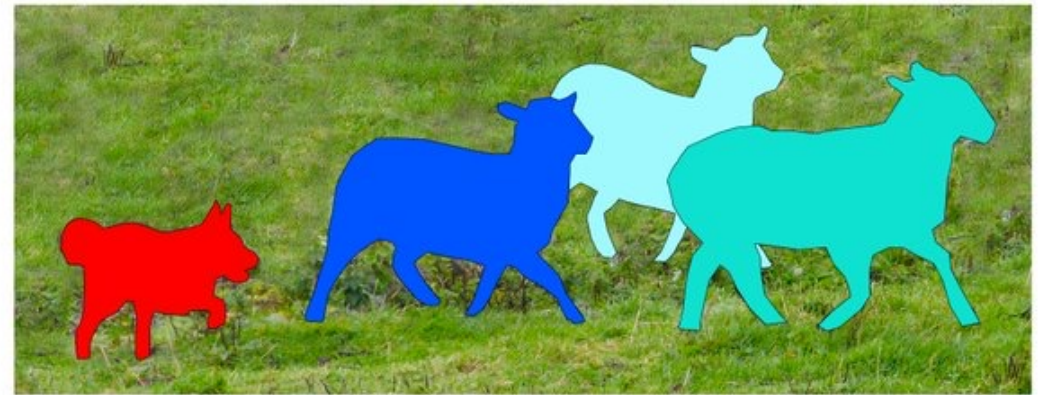# Today's Scope: Localize Content of Interest (Segmentation and Detection)



Object Detection

Use bounding boxes to locate every instance of an object from pre-specified categories

# Today's Scope: Localize Content of Interest (Segmentation and Detection)

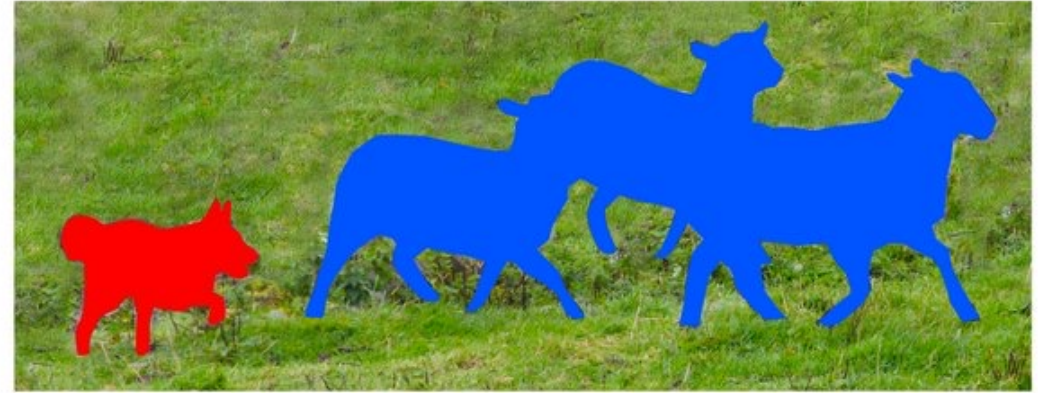Segment every instance of objects from pre-specified categories



**Instance Segmentation**

https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works

# Today's Scope: Localize Content of Interest (Segmentation and Detection)



Image Recognition

P 0.6 sheep
P 0.3 dog
P 0.1 cat
P 0.0 horse

Semantic Segmentation

Object Detection

Instance Segmentation

https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works

# Challenge: When to Choose Which Task?



Image Recognition

Semantic Segmentation

Object Detection

Instance Segmentation

https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works

# Today's Topics

- Problems

- **Applications**

- PASCAL VOC detection challenge: R-CNNs

- PASCAL VOC semantic segmentation challenge: fully convolutional networks

# Social Media



Face detection
(e.g., Facebook)

# Banking



Mobile check deposit
(e.g., Bank of America)

# Transportation



License Plate Detection (e.g., AllGoVision)

# Construction Safety



Pedestrian Detection
(e.g., Blaxtair)

# Counting



Counting Fish (e.g., SalmonSoft)

http://www.wecountfish.com/?page_id=143



Business Traffic Analytics

# Remodeling Inspiration



(a) Target photo

(b) Retextured

# Rotoscoping (many examples on Wikipedia)

# Disease Diagnosis; e.g., PathAI

# Face Makeover



Demo: https://www.maybelline.com/virtual-try-on-makeup-tools

# Self-Driving Vehicles



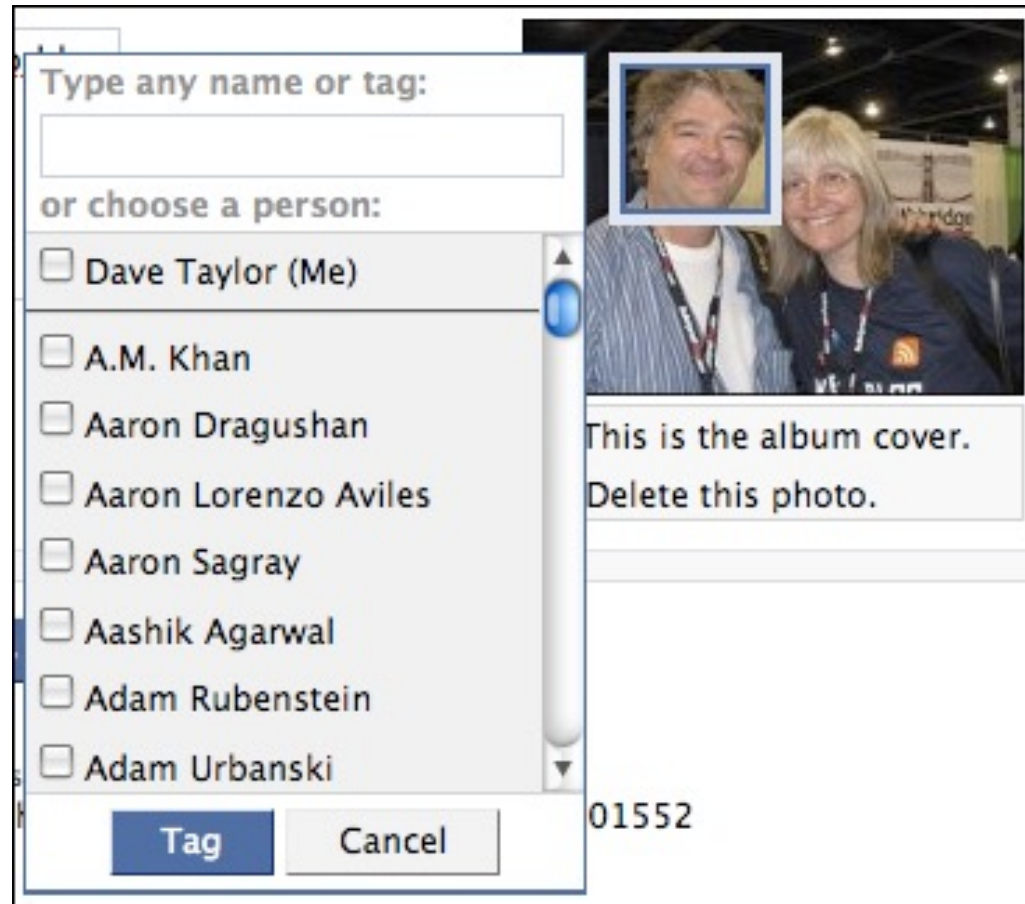Figure Source: https://www.inc.com/kevin-j-ryan/self-driving-cars-powered-by-people-playing-games-mighty-ai.html

# Today's Topics

- Problems

- Applications

- **PASCAL VOC detection challenge: R-CNNs**

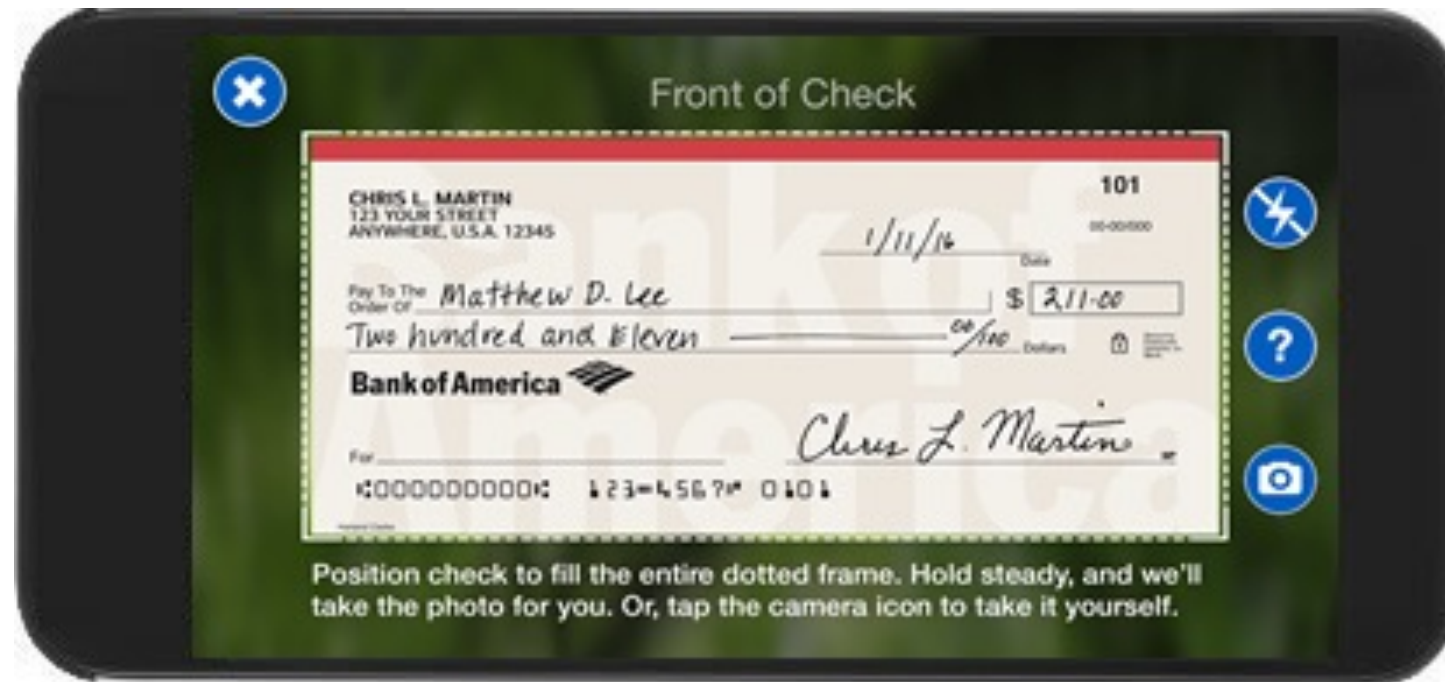- PASCAL VOC semantic segmentation challenge: fully convolutional networks

# VOC Challenge

- Goal: locate all instances of 20 object categories with BBs

- Dataset: 11,530 images collected from Flickr and annotated by annotators at University of Leeds



https://cv-tricks.com/artificial-intelligence/object-detection-using-deep-learning-for-advanced-users-part-1/

Dataset location: http://host.robots.ox.ac.uk/pascal/VOC/index.html
Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# VOC Challenge: Evaluation Metric (IoU)

Ground Truth:

Algorithm:

$$\frac{|A \cap B|}{|A \cup B|}$$

Score

# VOC Challenge: Evaluation Metric (IoU)

Ground Truth:

Algorithm:

$$\frac{28}{47}$$

(60%)

Then, threshold:

e.g., 50% or greater means correct detection!

# VOC Challenge: Evaluation Metric (mAP)

- For each object class (e.g., cat, dog, …), compute:
  - Precision: fraction of correct detections from all detections using 0.5 IoU threshold



Ground truth

Algorithm BB + its Confidence

[Russakovsky et al; IJCV 2015]

*https://jonathan-hui.medium.com/map-mean-average-precision-for-object-detection-45c121a31173*

# VOC Challenge: Evaluation Metric (mAP)

- For each object class (e.g., cat, dog, …), compute:
  - Precision: fraction of correct detections from all detections using 0.5 IoU threshold



Ground truth

AP: ? ? ? ?

[Russakovsky et al; IJCV 2015]
https://jonathan-hui.medium.com/map-mean-average-precision-for-object-detection-45c121a31173

# VOC Challenge: Evaluation Metric (mAP)

- For each object class (e.g., cat, dog, ...), compute:
  - Precision: fraction of correct detections from all detections using 0.5 IoU threshold
- Then, compute mean precision across all classes



Ground truth

AP: 0.0  0.5  1.0  0.3

[Russakovsky et al; IJCV 2015]
https://jonathan-hui.medium.com/map-mean-average-precision-for-object-detection-45c121a31173

# Evaluation Metric (mAP):
# Why "Mean" and Why "Average"

- More generally, for each object class (e.g., cat, dog, ...) :
  - AP: compute area under a precision-recall curve, created by varying IoU threshold

- Then, compute mean AP across all classes

# Naïve Solution: Sliding Window Approach

Person?

Person?

Person?

Person?

Person?

Person?

Person?

Person?

Person?



Image Source: https://yourboulder.com/boulder-neighborhood-downtown/

# Naïve Solution: Sliding Window Approach

Car?

Car?

Car?

Car?

Car?

Car?

Car?

Car?

Car?



Image Source: https://yourboulder.com/boulder-neighborhood-downtown/

# Naïve Solution: Sliding Window Approach

- Sliding window approach: must test different locations at…
  - Different scales
  - Different aspect ratios (e.g., for person vs car or car viewed at different angles)

- Number of regions to test? (e.g., 1920 x 1080 image)
  - Easily can explode to hundreds of thousands or millions of windows

- Key limitation
  - Very slow!

# Historical Context: R-CNN Methods

# R-CNN



- First CNN to outperform hand-crafted features on detection challenges
- Named after technique: **R**egion proposals with **CNN** features

# R-CNN



Apply bounding-box regressors

Classify regions with SVMs

Forward each region through ConvNet

Warped image regions

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

Input image

Locate "object"-like regions using objectness methods
- Considerably fewer regions than sliding window approach
- Regions likely contain objects of interest (i.e., high recall)

# R-CNN

# Describe Each Region with Fixed-length Vector

Given relatively little amount of training data, devise good feature by fine-tuning pre-trained model

1) Replace final layer of AlexNet (trained on ImageNet) with # of categories in detection dataset

2) Train for image classification (use max IoU class, if IoU >= 0.5)

How many classes should be predicted?

Input: 227 x 227 x 3 image

C1

C2

C3

C4

C5

FC6

FC7

FC8

5

5

3

3

3

3

3

3

13

13

384

13

13

384

13

13

256

27

27

256

55

55

96

4096

4096

1000

# Describe Each Region with Fixed-length Vector



Image Source: https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454

# Describe Each Region with Fixed-length Vector

Challenge: how to resize a proposed region to the required size?



Image Source: https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454

# Describe Each Region with Fixed-length Vector

Region anisotropically scaled to fit the required resolution



Image Source: https://www.researchgate.net/figure/Architecture-of-Alexnet-From-left-to-right-input-to-output-five-convolutional-layers_fig2_312303454

# Describe Each Region with Fixed-length Vector



Input: 227 x 227 x 3 image

# R-CNN



1. SVM classifier trained to use a region's CNN feature to assign a category from pre-defined set

2. Regressor trained to refine each region's position, width, and height

Figure Source: https://www.analyticsvidhya.com/blog/2018/10/a-step-by-step-introduction-to-the-basic-object-detection-algorithms-part-1/

# R-CNN: Region Refinement

Original region proposal with center $(p_x, p_y)$, width $(p_w)$, and height $(p_h)$ is refined using model parameters $(d_x, d_y, d_w, d_h)$

# Algorithm Training: Linear Regression Model

- **Aim**: learn transformation from region proposal to ground truth

- **Input**: original region location; BB described by a center ($p_x$, $p_y$), width ($p_w$), and height ($p_h$)

- **Output**: learns four refinement functions: $d_x$, $d_y$, $d_w$, $d_y$

- Loss function for learning: SSE

$$\sum_{i \in \{x,y,w,h\}} (t_i - d_i(\mathbf{p}))^2$$

True location

Predicted location

# R-CNN Limitations



- **Slow training procedure**
  - Must train three models

- **Slow at test time**
  (~1 minute per image)

# Fast R-CNN: Single Stage Training (rather than 3)



For each region, assign it to a class and refine it

Extract feature description per proposed region with section of feature map corresponding to region

# Fast R-CNN Training: Multi-task Loss

Objective function sums classification and localization losses for each region proposal

# Fast R-CNN Training: Multi-task Loss

Objective function sums classification and localization losses for each region proposal

# Fast R-CNN Training: Classification Loss (Recall Cross Entropy Loss, aka Log Loss)



Range of negative log-likelihood

Greater penalty when predicted probability of true class is confidently wrong

Lesser penalty otherwise

$$-\log \frac{\exp\left(w_k \cdot x + b_k\right)}{\sum_{j=1}^{K} \exp\left(w_j \cdot x + b_j\right)}$$

Figure source: https://ljvmiranda921.github.io/notebook/2017/08/13/softmax-and-the-negative-log-likelihood/

# Fast R-CNN Training: Multi-task Loss

Objective function sums classification and localization losses for each region proposal

# Fast R-CNN Training: Measure Localization Loss

$$\mathcal{L}_{\text{box}}(t^u, v) = \sum_{i \in \{x,y,w,h\}} L_1^{\text{smooth}}(t_i^u - v_i) \longrightarrow L_1^{\text{smooth}}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

True location for true class "u"

Predicted location for class u

Less sensitive to outliers than SSE



Image Source: https://lilianweng.github.io/lil-log/2017/12/31/object-recognition-for-dummies-part-3.html#bounding-box-regression

# Fast R-CNN: Limitation



Still requires slow, initial step of generating region proposals

Figure Source: https://www.analyticsvidhya.com/blog/2018/10/a-step-by-step-introduction-to-the-basic-object-detection-algorithms-part-1/

# Faster R-CNN

Adds finding region proposals to network so that all parts of model are learned in end-to-end fashion

Convolutional layers are shared for region proposal and detection



Ren Shaoqing Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Neurips 2015.

# Faster R-CNN: Region Proposal Network

Probability of object/not object

Parameters to refine anchor box to match GT box (center, width, and height)

2k scores

4k coordinates

k anchor boxes

cls layer

reg layer

256-d

intermediate layer

sliding window

conv feature map

**Based on convolution, so uses sliding window**

- At each sliding window position, region proposals are predicted with respect to an anchor point (i.e., center of sliding window position)
- At each anchor point, k = 9 anchors are used to represent 3 scales and 3 aspect ratios

Ren Shaoqing Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Neurips 2015.

# Faster R-CNN: Region Proposal Network

At training, loss for each region proposal is sum of classification and localization losses



**Based on convolution, so uses sliding window**

- At each sliding window position, region proposals are predicted with respect to an anchor point (i.e., center of sliding window position)
- At each anchor point, k = 9 anchors are used to represent 3 scales and 3 aspect ratios

Ren Shaoqing Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Neurips 2015.

# Faster R-CNN Training

1. Train RPN

2. Train Fast R-CNN using proposals from pretrained RPN

3. Fine-tune layers unique to RPN

4. Fine-tune the fully connected layers of Fast R-CNN



Ren Shaoqing Ren et al. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Neurips 2015.

Historical Context: In 2017, Mask R-CNN Introduced for Instance Segmentation

1847 — Gradient descent

1945 — First programmable machine

1950 — Turing test

1956 — AI
1957 — Perceptron
1959 — Machine learning

1980 — Neocognitron

1986 — Neural networks with effective learning strategy

1989 — Backpropagation for CNNs

1998 — MNIST, LeNet

2009 — ImageNet

2012 — AlexNet

2014-5 — R-CNN, Fast R-CNN, Faster R-CNN

# Today's Topics

- Problems

- Applications

- PASCAL VOC detection challenge: R-CNNs

- **PASCAL VOC semantic segmentation challenge: fully convolutional networks**

# VOC Challenge

- Goal: locate all pixels belonging to 20 categories (e.g., person, cat, bus, mortorbike, potted plant, bottle) plus background

- Dataset: 11,530 images collected from Flickr and annotated by annotators at University of Leeds



Dataset location: http://host.robots.ox.ac.uk/pascal/VOC/index.html
Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# VOC Challenge: Evaluation Metric (IoU)



Ground Truth:

Algorithm:

$$\frac{|A \cap B|}{|A \cup B|}$$

Score

# VOC Challenge: Evaluation Metric (IoU)

Ground Truth:

Algorithm:

?

# VOC Challenge: Evaluation Metric (IoU)

**Mean IoU**: IoU between predicted and ground-truth pixels, averaged over all 21 categories

Ground Truth:

Algorithm:

$$\frac{19}{27}$$

# Architecture

For each image pixel, the probability of each class is predicted



256   384   384   256   4096   4096   21

96

pixelwise prediction

segmentation g.t.

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Output Layer

- e.g., assume a 5-class classifier

# Architecture: Output Layer

- e.g., assume a 5-class classifier; output 1-hot encoding collapsed into single mask image



0: Background/Unknown
1: Person
2: Purse
3: Plants/Grass
4: Sidewalk
5: Building/Structures

# Architecture

Input: RGB image of ANY size

Output: Image of same size as input

How many classes are there?
- 21
Why 21?
- 20 object classes plus background

pixelwise prediction

segmentation g.t.

256   384   384   256   4096   4096   21

96

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture



Do you recognize this architecture?

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture



Can use your favorite
pretrained ImageNet classifier;
AlexNet, VGG, GoogleNet

96 256 384 384 256 4096 4096 21

pixelwise prediction

segmentation g.t.

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture

To make the architecture fully convolutional, fully connected layers are converted to convolutional layers.

In the absence of fully connected layers, there are no constraints on the number of input nodes (and so any input image size can be supported).



pixelwise prediction

segmentation g.t.

96   256   384   384   256   4096   4096   21

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture



Another result of this change is that, unlike for classification, a class can be assigned to each "coarse region."

pixelwise prediction

segmentation g.t.

96   256   384   384   256   4096   4096   21

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Coarse Region Classification (Recall Intuition)



Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Coarse Region Classification (Recall Intuition)



Each line represents a convolutional layer

Using VGG16 instead:

image  conv1  pool1  conv2  pool2  conv3  pool3  conv4  pool4  conv5  pool5  conv6-7

Grids reflect relative spatial coarseness at each layer

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Coarse Region Classification (Recall Intuition)

Stacking many convolutional layers leads to learning patterns in increasingly **larger regions of the input (e.g., pixel) space.**

# Architecture: Fully vs Convolution Layers



Each slice indicates the likelihood each pixel in the coarse region belongs to the class identified by the filter

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Fully vs Convolution Layers



"tabby cat"

96 256 384 384 256 4096 4096 1000

convolutionalization

tabby cat heatmap

96 256 384 384 256 4096 4096 1000

If convolutionizing ImageNet trained classifiers, how many classes would be predicted for each coarse region?

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Coarse Region Classification



Locates 20 object classes plus background for VOC

pixelwise prediction

segmentation g.t.

96 256 384 384 256 4096 4096 21

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Upsampling (Many Approaches)

**Nearest Neighbor**

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

→

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

Input: 2 x 2          Output: 4 x 4

**"Bed of Nails"**

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

→

| 1 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 3 | 0 | 4 | 0 |
| 0 | 0 | 0 | 0 |

Input: 2 x 2          Output: 4 x 4

**Max Pooling**
Remember which element was max!

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

→

| 5 | 6 |
|---|---|
| 7 | 8 |

→ • • • →  Rest of the network

Input: 4 x 4          Output: 2 x 2

**Max Unpooling**
Use positions from pooling layer

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

→

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 4 |

Input: 2 x 2          Output: 4 x 4

# Architecture: Upsampling (Transposed Convolutional Layer)

- Also called "fractional convolutional layer", "backward convolution", and, incorrectly, "deconvolution layer"

- Idea: learn filters with a fractional sized stride to upsample the coarse image while refining it based on the filter values; e.g.,



https://www.machinecurve.com/index.php/2019/09/29/understanding-transposed-convolutions/#the-goal-reconstructing-the-original-input

# Architecture: Upsampling (Transposed Convolutional Layer)

- Also called "fractional convolutional layer", "backward convolution", and, incorrectly, "deconvolution layer"

- Idea: learn filters with a fractional sized stride to upsample the coarse image while refining it based on the filter values; e.g.,



https://d2l.ai/chapter_computer-vision/transposed-conv.html

# Architecture: Upsampling (Transposed Convolutional Layer)

- Also called "fractional convolutional layer", "backward convolution", and, incorrectly, "deconvolution layer"

- Idea: learn filters with a fractional sized stride to upsample the coarse image while refining it based on the filter values; e.g.,



(stride is used to compute intermediate values)

https://d2l.ai/chapter_computer-vision/transposed-conv.html

# Architecture



Next challenge: how to decode a **highly detailed** per pixel classification from the coarse region classifications?

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Results

Ground truth target

Predicted segmentation

Figure source: https://www.jeremyjordan.me/semantic-segmentation/

# Architecture: Update to Use Skip Connections



32x upsampled prediction (FCN-32s)

2x upsampled prediction

16x upsampled prediction (FCN-16s)

image  pool1  pool2  pool3  pool4  pool5

pool4 prediction

Trained ~1 more day to update the FCN-32 model

FCN16: Sums predictions of lower-level, more fine-grained features (pool4) with the predictions at the coarser features

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Results



Skip connections support capturing finer-grained details while retaining the correct semantic information!

Figure source: https://www.jeremyjordan.me/semantic-segmentation/

# Architecture: Upsampling + Skip Connections

Seems complicated... why not instead preserve the image size and solve for per-pixel classification?
- would result in unreasonable computational burden due to many model parameters



Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Architecture: Encoder Decoder Architecture



For efficiency, the image is encoded (downsampled) into a lower-resolution feature map that effectively discriminates between classes...

Then, the feature map is decoded (upsampled) into a full-resolution segmentation map.

pixelwise prediction

segmentation

96    256    384    384    256    4096    4096    21

21

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Training: Took 3 days on 1 GPU



- Repeat until stopping criterion met:
  1. **Forward pass**: propagate training data through model to make prediction
  2. Quantify the dissatisfaction with a model's results on the training data
  3. **Backward pass**: using predicted output, calculate gradients backward to assign blame to each model parameter
  4. Update each parameter using calculated gradients

Figure from: Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, Jeffrey Mark Siskind; Automatic Differentiation in Machine Learning: a Survey; 2018

# Training: How Neural Networks Learn

- Repeat until stopping criterion met:
  1. **Forward pass**: propagate training data through model to make prediction
  2. Quantify the dissatisfaction with a model's results on the training data
  3. **Backward pass**: using predicted output, calculate gradients backward to assign blame to each model parameter
  4. Update each parameter using calculated gradients

Sum across all pixels the distance between predicted and true distributions using cross entropy loss

Sum of gradients for all pixels (acts like a minibatch)

Figure from: Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, Jeffrey Mark Siskind; Automatic Differentiation in Machine Learning: a Survey; 2018
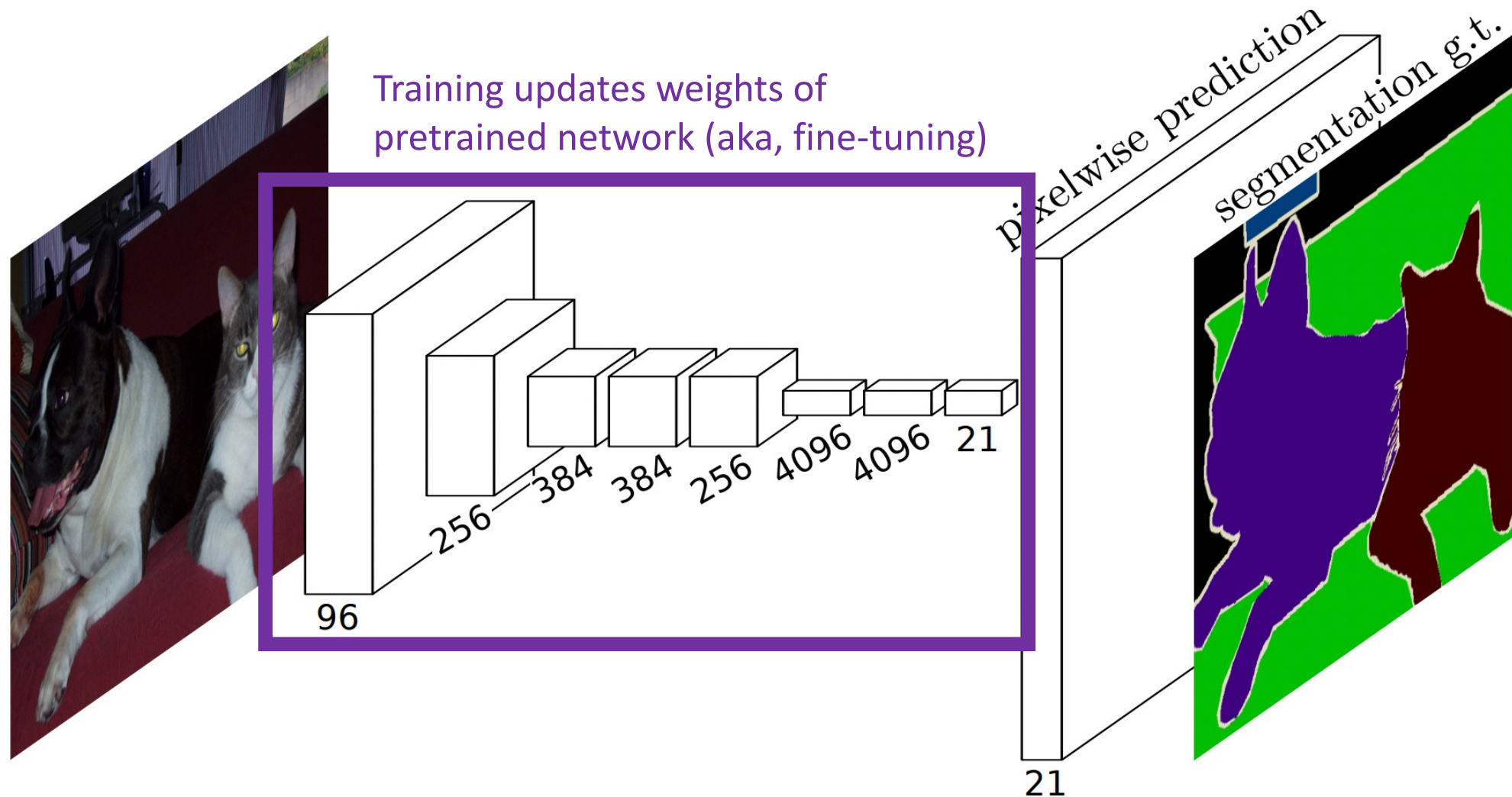
# Training: Cross Entropy Loss (Multinomial Logistic Loss)

- e.g., assume a 5-class classifier
- Distance between predicted and true distributions per pixel with cross entropy loss

# Architecture: Algorithm Training



Training updates weights of pretrained network (aka, fine-tuning)

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Results

| | mean IU VOC2011 test | mean IU VOC2012 test | inference time |
|---|---|---|---|
| R-CNN [12] | 47.9 | - | - |
| SDS [16] | 52.6 | 51.6 | $\sim$ 50 s |
| FCN-8s | **62.7** | **62.2** | $\sim$ **175 ms** |

Compared to existing methods, produces better results at a faster speed!

Long, Shelhamer, and Darrell. Fully Convolutional Networks for Semantic Segmentation. CVPR 2015.

# Today's Topics

- Problems

- Applications

- PASCAL VOC detection challenge: R-CNNs

- PASCAL VOC semantic segmentation challenge: fully convolutional networks