# Semantic Segmentation

**Danna Gurari**

University of Colorado Boulder

Fall 2023

# Review

- Last lecture:
  - Single Object Tracking lecture from Dr. Samreen Anjum

- Assignments (Canvas)
  - Reading assignment was due earlier today
  - Next reading assignments due next Monday and Wednesday
  - Project proposal due in one week

- Questions?

# Semantic Segmentation: Today's Topics

- Problem

- Applications

- Datasets

- Evaluation metric

- Computer vision models: fully convolutional networks

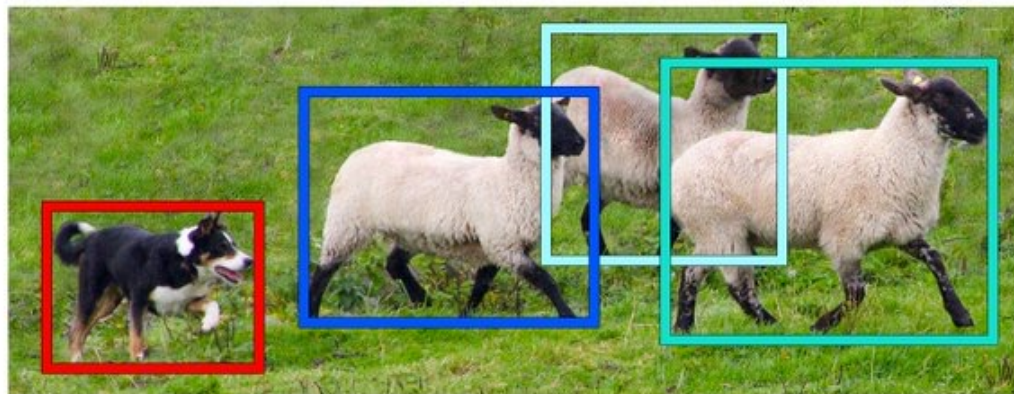- Discussion

# Semantic Segmentation: Today's Topics

- **Problem**

- Applications

- Datasets

- Evaluation metric

- Computer vision models: fully convolutional networks

- Discussion

# Recall: Object Recognition and Detection Tasks



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

**Image Recognition**

**Recognize** categories of interest



**Object Detection**

**Localize** categories of interest

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

# Today's Scope: Localize Pixels for Each Category



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

**Image Recognition**

**Object Detection**

**Semantic Segmentation**

Note: instances of the same category are NOT separated

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

# Today's Scope: Localize Pixels for Each Category



Semantic Segmentation

Instance Segmentation

Separating instances of the same category will be covered in a future lecture

https://ai-pool.com/d/could-you-explain-me-how-instance-segmentation-works
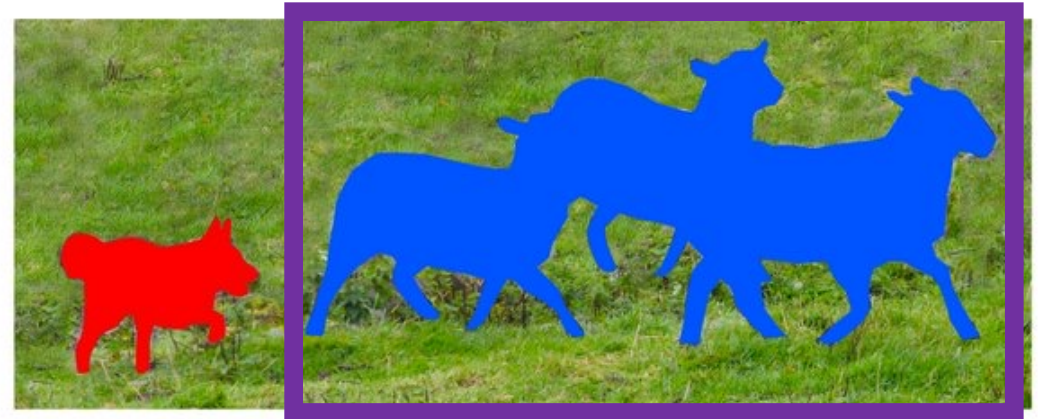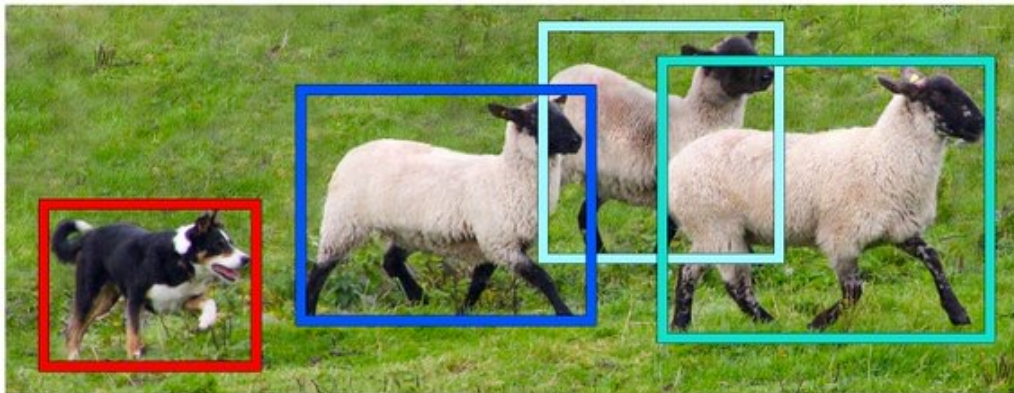
# Challenge: When to Choose Which Task?



Image Recognition

Semantic Segmentation

Object Detection

# Semantic Segmentation: Today's Topics

- Problem

- **Applications**

- Datasets

- Evaluation metric

- Computer vision models: fully convolutional networks
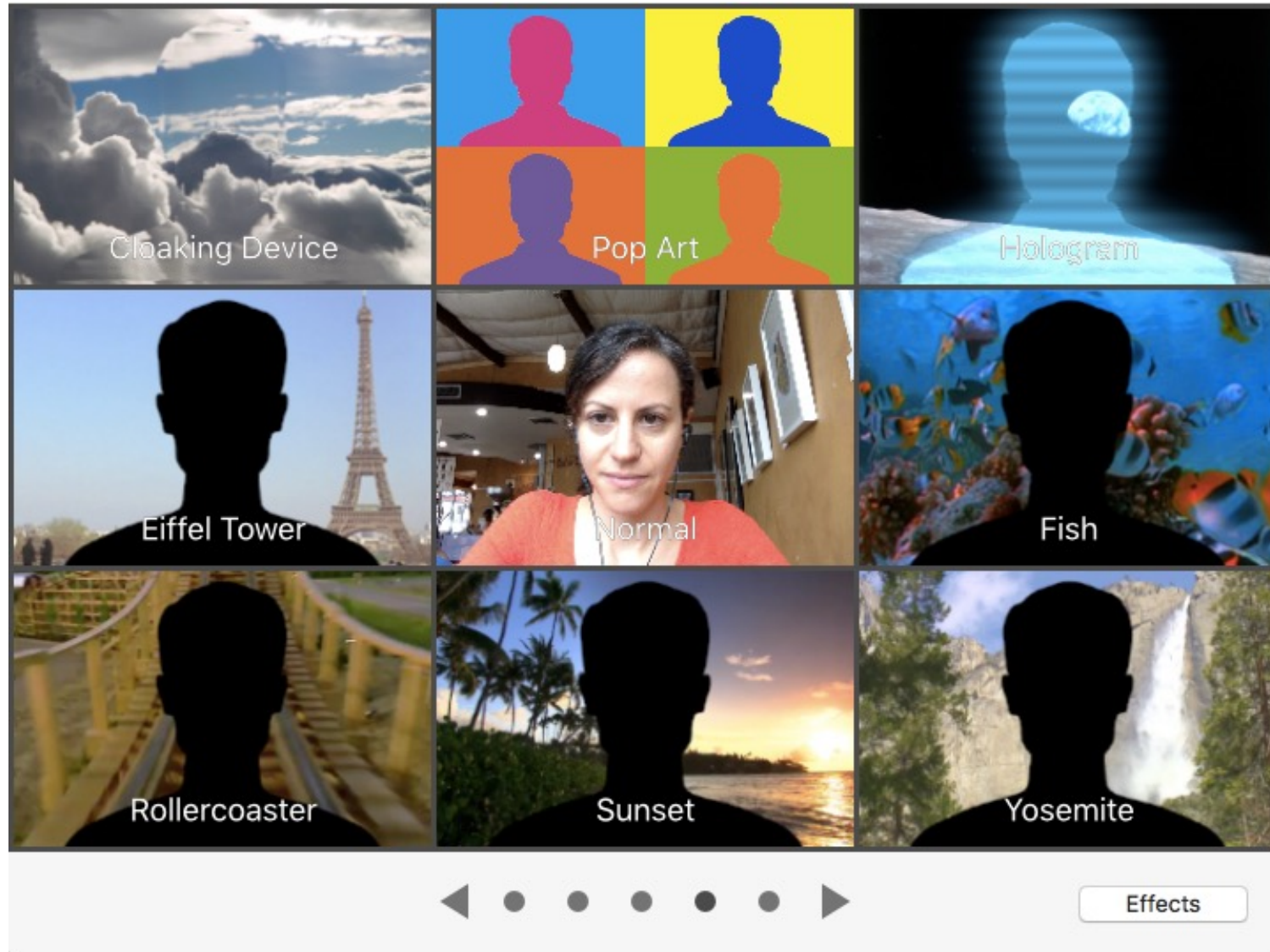
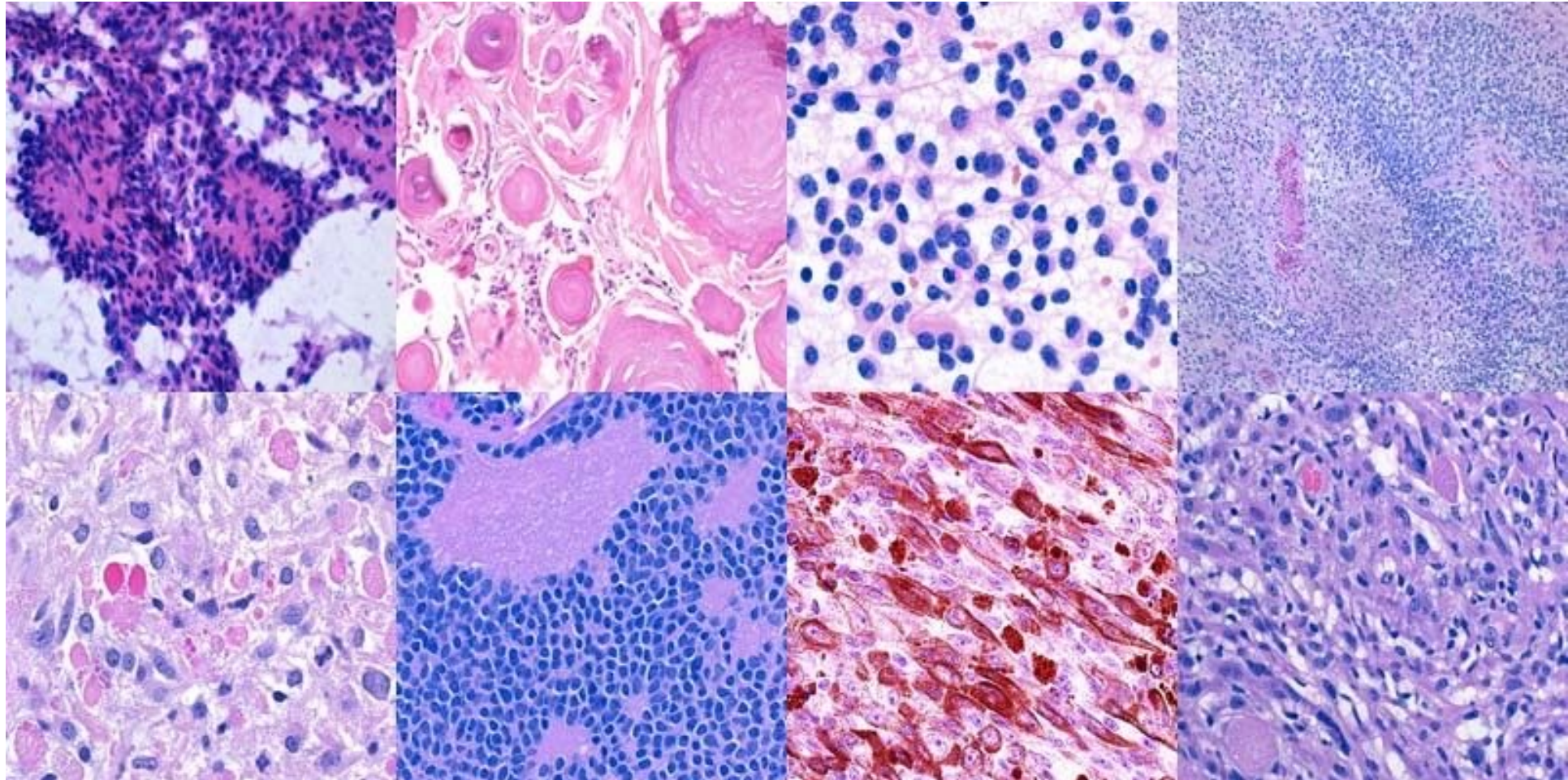- Discussion

# Remodeling Inspiration



(a) Target photo

(b) Retextured

Bell et al; SIGGRAPH; 2013

# Rotoscoping (many examples on Wikipedia)



https://www.starnow.co.uk/ahmedmohammed1/photos/4650871/before-and-after-rotoscopinggreen-screening

# Disease Diagnosis; e.g., PathAI

# Face Makeover

# Self-Driving Vehicles



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

Can you think of any other potential applications?

# Semantic Segmentation: Today's Topics

- Problem

- Applications

- **Datasets**

- Evaluation metric

- Computer vision models: fully convolutional networks

- Discussion

# Datasets



| 1945 | 1957 | 1966 | 1983 | 1987 | 1990 | | 2009 | 2017 |
|------|------|------|------|------|------|---|------|------|
| | | | CVPR | ICCV | ECCV | e.g., | VOC | ADE20k |

| | | |
|---|---|---|
| **# Categories:** | 21 | 3,169 |
| **# Images:** | 1112 train/val | 25,210 |

## Trend: build bigger datasets

# VOC

- A subset of images from the VOC detection dataset were used

- Annotation party annually

- Annotation guidelines & real-time assistance – refine detections into segmentations

- Post-hoc correction/feedback about the number and kind of errors made

- Annotations for each of the 20 object classes were merged into class-specific segmentation regions and 1 more class was added for background

Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# VOC: Recall Categories Included (Leaf Nodes)



Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.
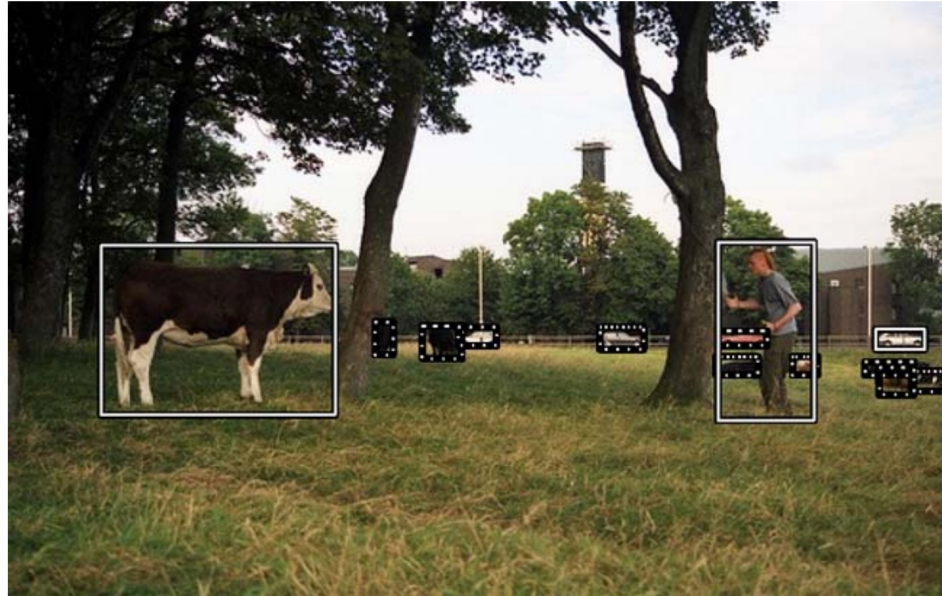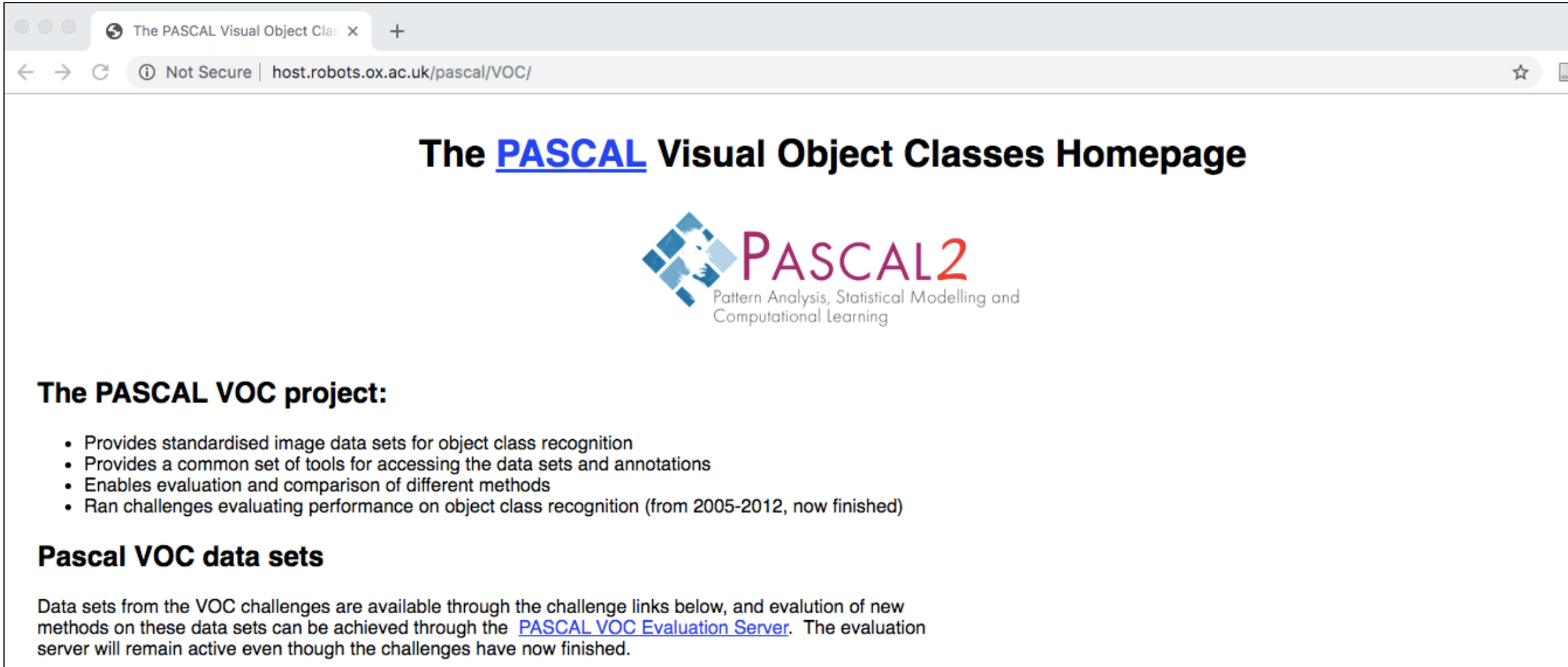
# VOC: Boundary Accuracy Heuristic



"To give high accuracy but to keep the annotation time short enough to provide a large image set, a border area of 5 pixels width was allowed around each object where the pixels were labelled neither object nor background."

Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# VOC: "Difficult" Objects Excluded



Objects that are challenging to recognize are discarded (i.e., dashed regions): flagged for reasons of "small size, illumination, image quality or the need to use significant contextual information… no penalty is incurred for detecting them. The aim of this annotation is to maintain a reasonable level of difficulty…"

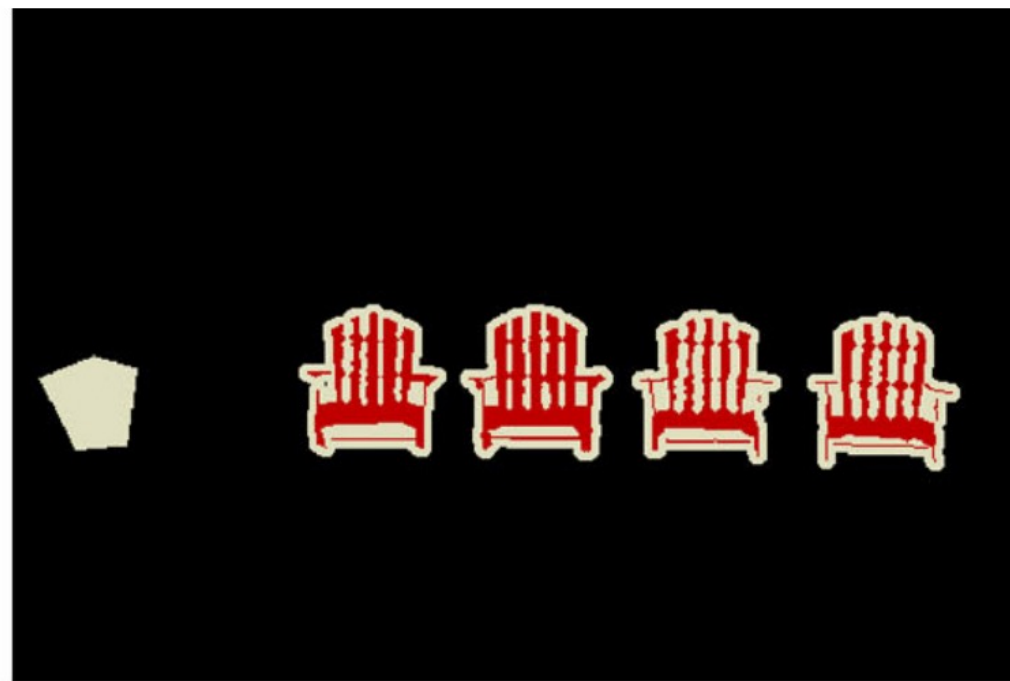Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. IJCV 2010.

# VOC Annual Workshop



The PASCAL Visual Object Cla ×    +

← → C    ⓘ Not Secure | host.robots.ox.ac.uk/pascal/VOC/    ☆

## The **PASCAL** Visual Object Classes Homepage

**PASCAL2**
Pattern Analysis, Statistical Modelling and
Computational Learning

### The PASCAL VOC project:

- Provides standardised image data sets for object class recognition
- Provides a common set of tools for accessing the data sets and annotations
- Enables evaluation and comparison of different methods
- Ran challenges evaluating performance on object class recognition (from 2005-2012, now finished)

### Pascal VOC data sets

Data sets from the VOC challenges are available through the challenge links below, and evalution of new methods on these data sets can be achieved through the PASCAL VOC Evaluation Server.  The evaluation server will remain active even though the challenges have now finished.

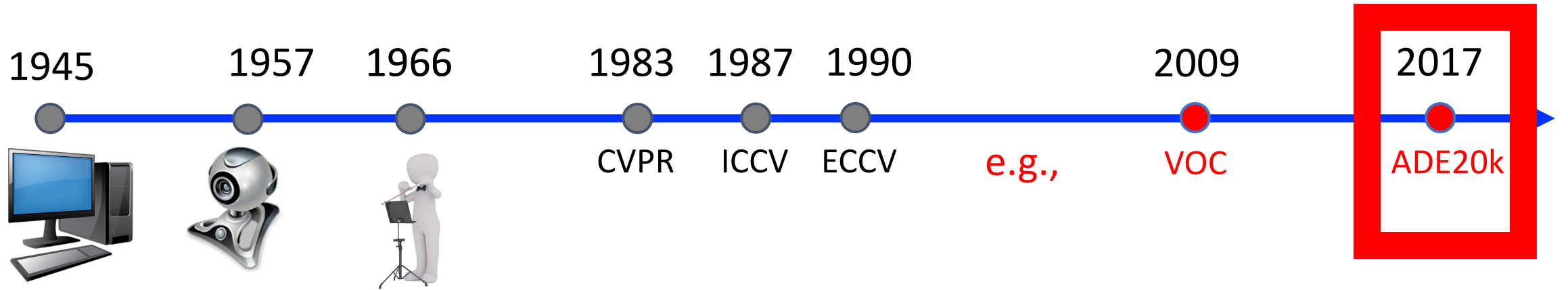http://host.robots.ox.ac.uk/pascal/VOC/

# What is a Limitation of Datasets Built Around Specific Categories (e.g., Objects)?



No knowledge that anything else is in the scene, such as a house, trees or flowers!
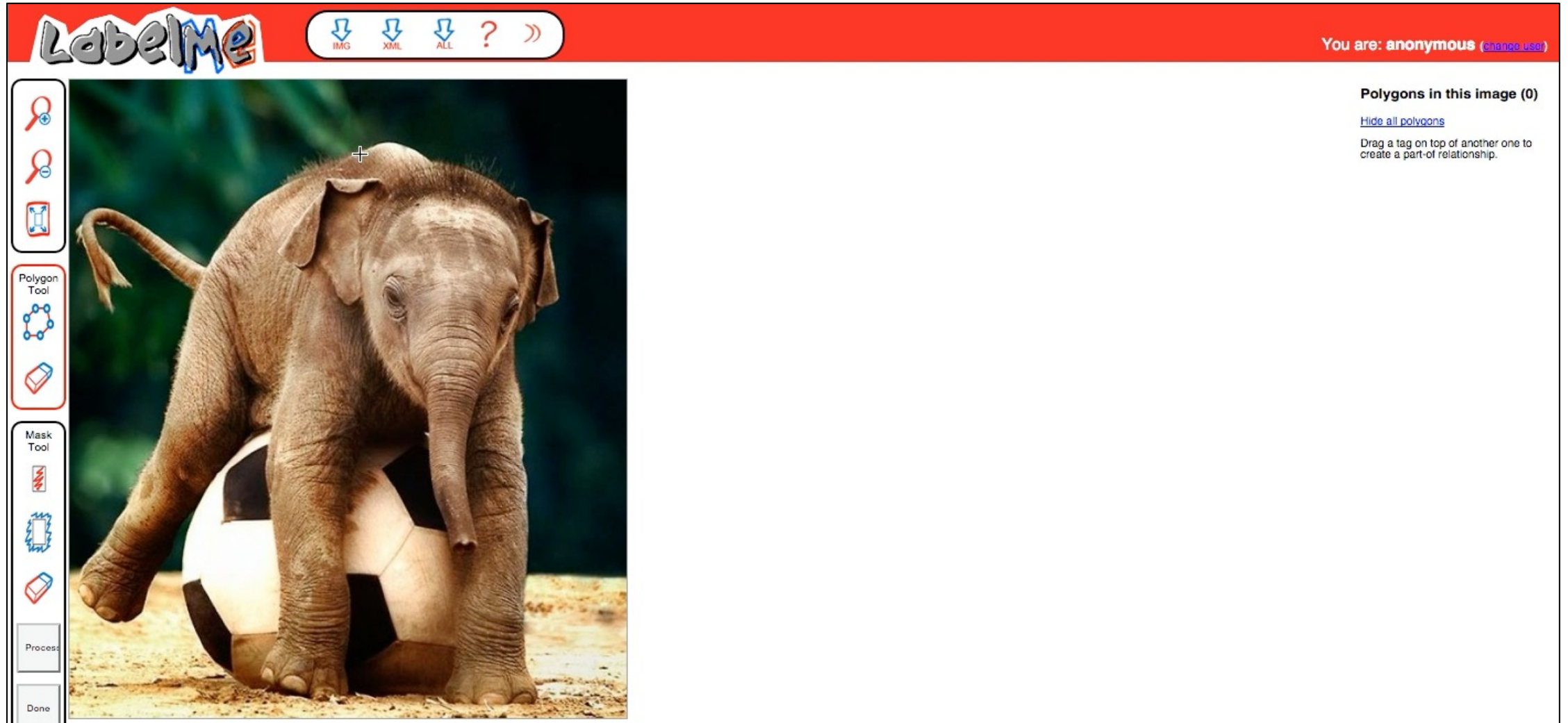
Most pixels are labeled as `background'!

Mark Everingham et al. The PASCAL Visual Object Classes Challenge: A Retrospective. IJCV 2015.

# Datasets

1945　　1957　　1966　　1983　1987　1990　　　　　2009　　2017

CVPR　ICCV　ECCV　　e.g.,　　VOC　　ADE20k

# ADE20K

- 25,210 images collected from existing datasets (SUN, Places, and LabelMe)

- Selected to capture all scene categories defined in SUN

- A single person annotated all images into three types and kept adding new categories as they were observed: (1) objects, (2) object parts, and (3) attributes (e.g., occluded)

# ADE20K: User Annotation Tool

# ADE20K: User Annotation Tool



Bolei Zhou et al. Scene Parsing through ADE20K Dataset. CVPR 2017.

# ADE20K

- Includes:
- "things": objects that can easily be labeled; e.g., person, chair
- "stuff": objects with no clear boundaries; e.g., sky, grass



Image  Ground-truth

Bolei Zhou et al. Scene Parsing through ADE20K Dataset. CVPR 2017.

# Semantic Segmentation: Today's Topics

* Problem

* Applications

* Datasets

* **Evaluation metric**

* Computer vision models: fully convolutional networks

* Discussion

# Evaluation Metric

Ground Truth:

Algorithm:

Evaluation Measure

→ Score

# Recall: IoU Metric

Ground Truth:

Algorithm:

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

Score

# Recall: IoU Metric

Ground Truth:

Algorithm:

?

# Recall: IoU Metric

Ground Truth:

Algorithm:

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

# Mean IoU (mIoU)

- Mean IoU score over all categories

# Semantic Segmentation: Today's Topics

- Problem

- Applications

- Datasets

- Evaluation metric

- **Computer vision models: fully convolutional networks**

- Discussion

# Why Fully Convolutional Network?

Named after the proposed technique that excludes fully connected layers:

Jonathon Long, Evan Shelhamer, and Trevor Darrell. "Fully Convolutional Networks for Semantic Segmentation." CVPR 2015.

# Key Novelties of Fully Convolutional Networks

First work for pixelwise prediction to:

1. Train fully convolutional networks end-to-end

2. Use supervised pre-training (recall, R-CNN paper showed this can be a great idea when there is a scarce amount of annotated data)

# Architecture

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



256   384   384   256   4096   4096   21

96

21

pixelwise prediction   segmentation g.t.

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

Source: https://www.jeremyjordan.me/semantic-segmentation/

# Architecture

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



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



Do you recognize this architecture?

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



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

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.

256  384  384  256  4096  4096  21

96

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."

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

# Architecture: Coarse Region Classification (Recall Intuition)



Using VGG16 instead:

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

pixelwise prediction

segmentation g.t.

96

256

21

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



"tabby cat"

convolutionalization

tabby cat heatmap

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"

convolutionalization

tabby cat heatmap

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

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

Challenge: how to decode from coarse region classifications to per pixel classification?



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

# Architecture: Upsampling (Many Approaches)

# Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea**: learn convolutional filters with a fractional sized stride to upsample the coarse image while refining it; e.g., 1/2 stride

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



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

# Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea**: learn convolutional filters with a fractional sized stride to upsample the coarse image while refining it; e.g., 1/2 stride

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



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

# Architecture: Upsampling (Transposed Convolutional Layer)

- **Idea**: learn convolutional filters with a fractional sized stride to upsample the coarse image while refining it; e.g., 1/2 stride

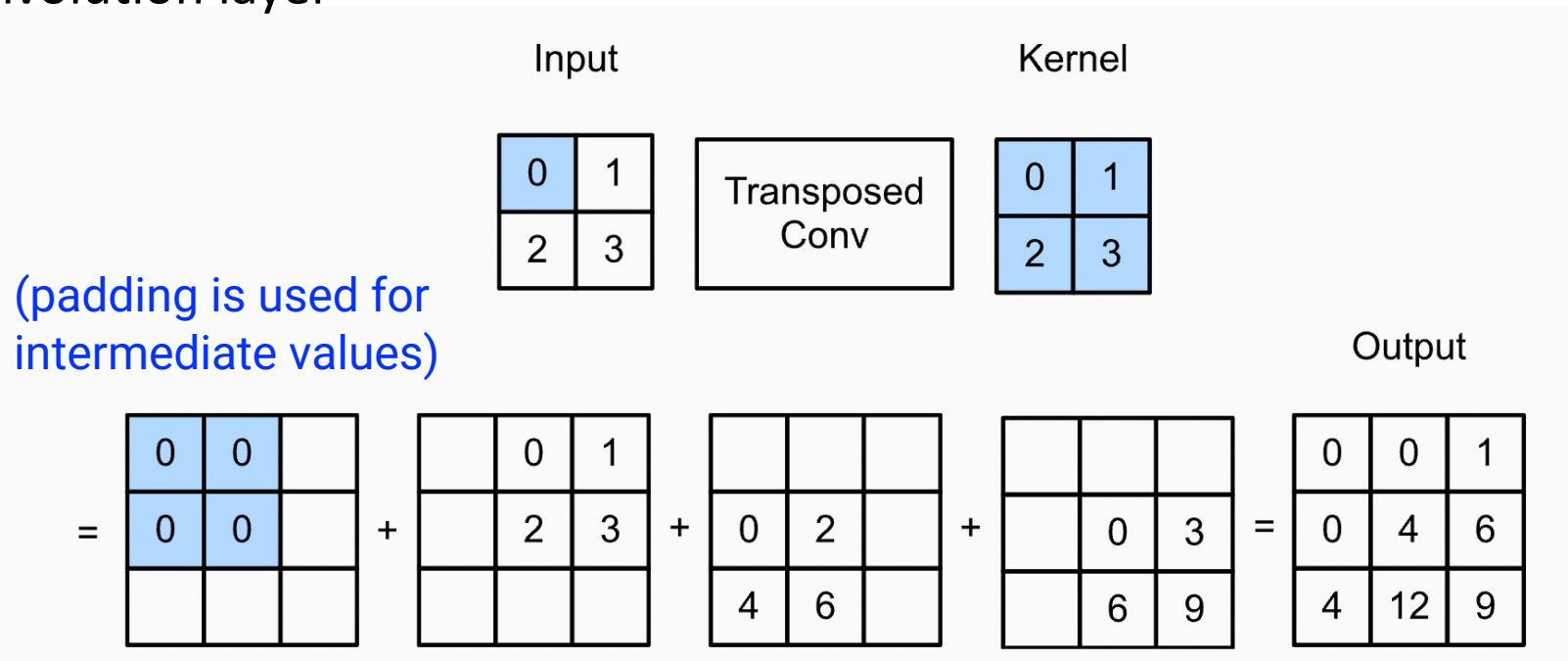- Also called "fractional convolutional layer", "backward convolution", and, incorrectly, "deconvolution layer"



(stride is used to compute intermediate values)

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

# Architecture

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: Results

Ground truth target

Predicted segmentation

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

# Architecture: Update to Use Skip Connections



FCN16: Fuses class predictions of lower-level, more fine-grained features with the predictions at the coarser features

FCN8: Fuses predictions of even lower-level, more fine-grained features with both predictions at the coarser features

# Architecture: Results



Ground truth target | FCN-32s | FCN-16s | FCN-8s

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

# 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
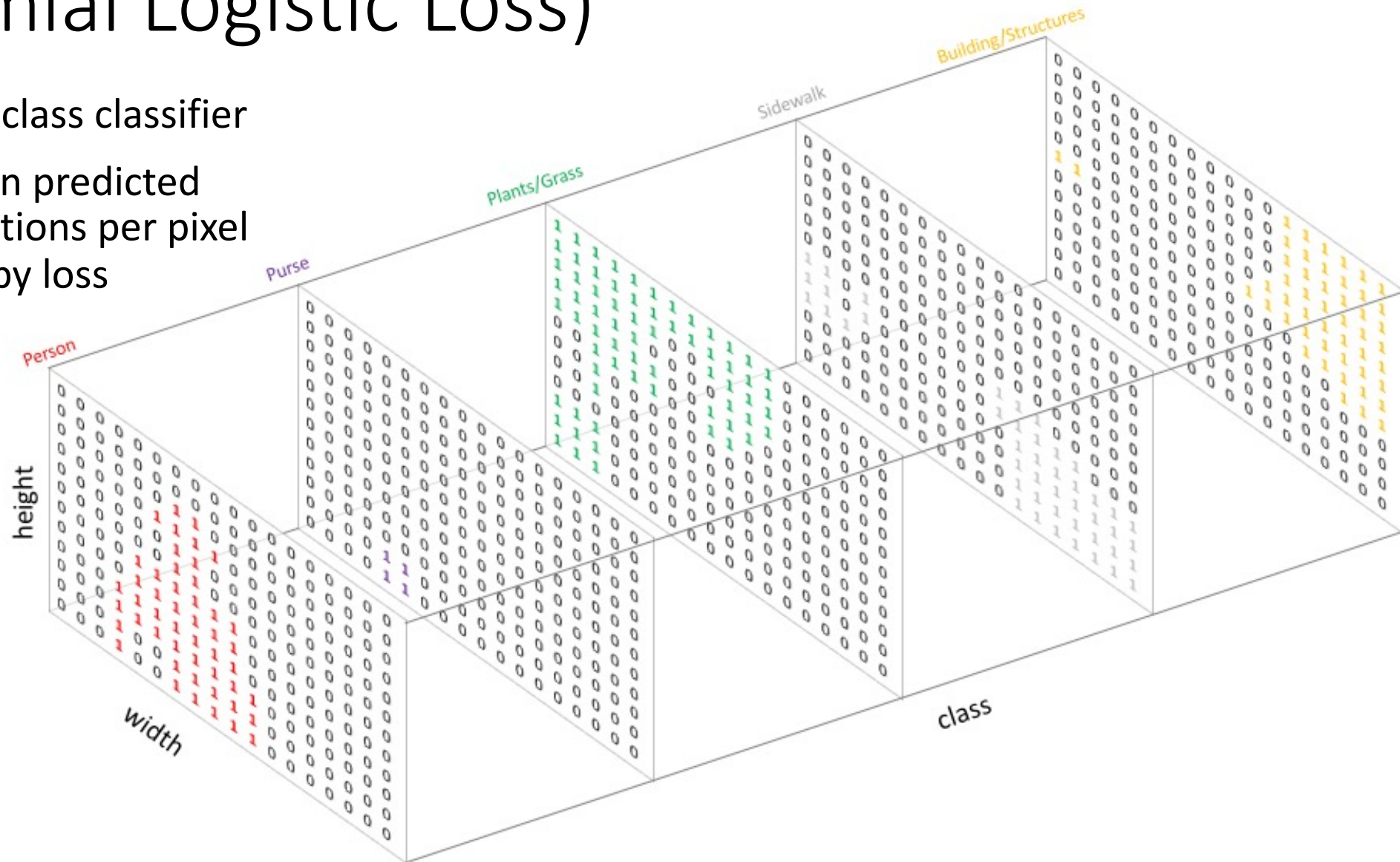  4. Update each parameter using calculated gradients

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
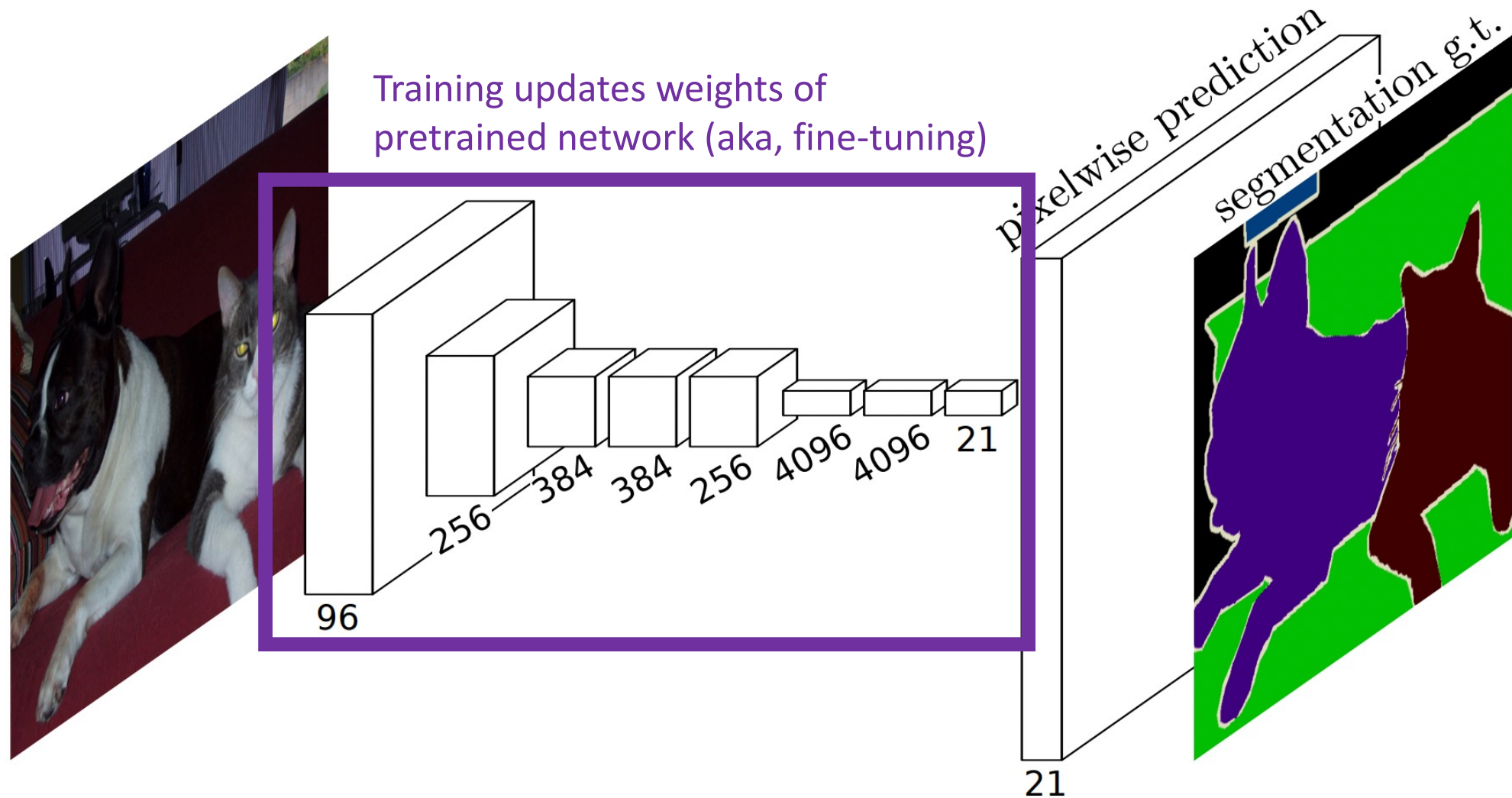
# Training: Cross Entropy Loss (Multinomial Logistic Loss)

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

# Architecture: Algorithm Training
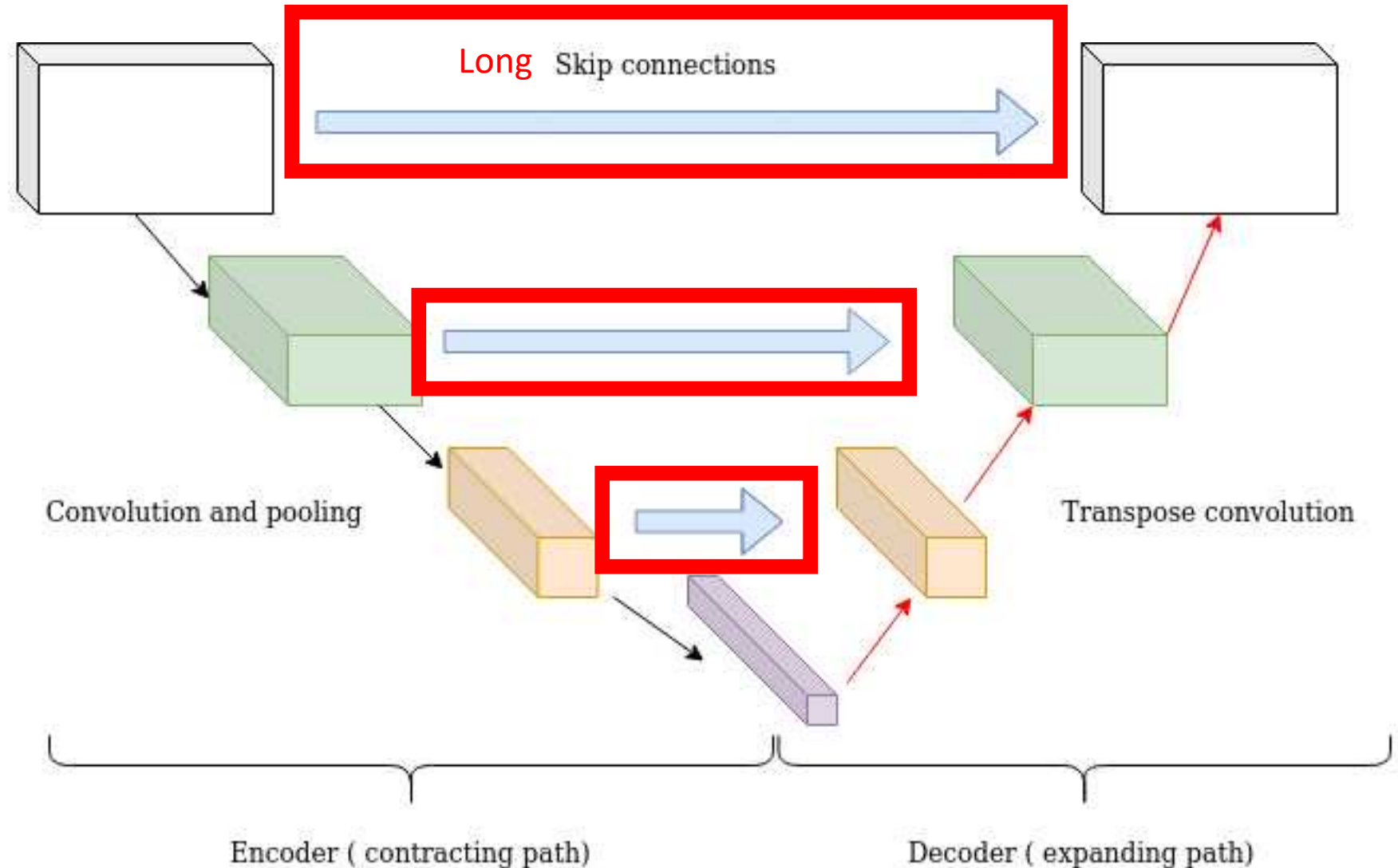


Training updates weights of pretrained network (aka, fine-tuning)

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.

# 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.

# Improved Architecture: U-Net



Passes information lost in the encoder to the decoder from **each downsampling layer** in the encoder to its corresponding upsampling layer in the decoder, while also keeping the computation low.

Long Skip connections

Convolution and pooling

Transpose convolution

Encoder ( contracting path )

Decoder ( expanding path )

Image Source: https://theaisummer.com/skip-connections/

# U-Net



Ronneberger, Fischer, and Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. MICCAI 2015.

# Semantic Segmentation: Today's Topics

- Problem

- Applications

- Datasets

- Evaluation metric

- Computer vision models: fully convolutional networks

- **Discussion**

# Semantic Segmentation: Today's Topics

- Problem

- Applications

- Datasets

- Evaluation metric

- Computer vision models: fully convolutional networks

- Discussion