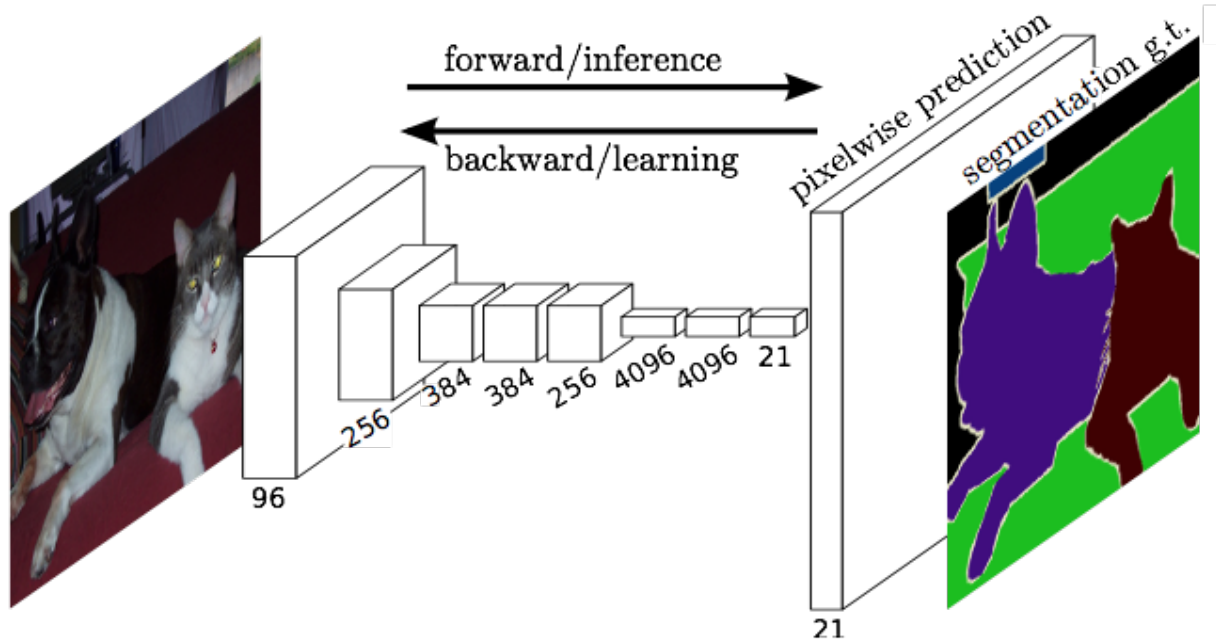


Fully Convolutional Networks

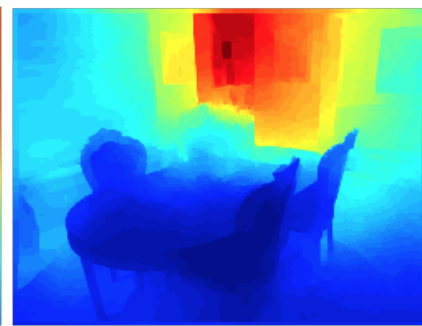
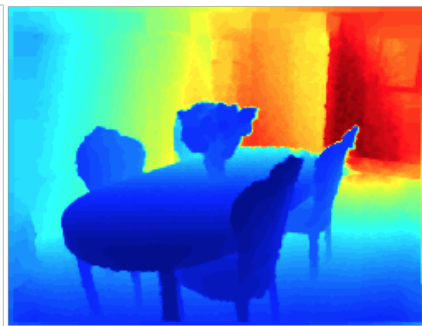


Jon Long and Evan Shelhamer
CVPR15 Caffe Tutorial

pixels in, pixels out

monocular depth estimation (Liu et al. 2015)

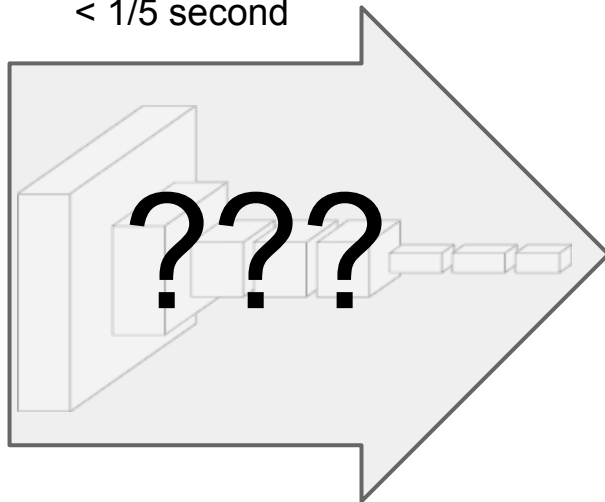
semantic segmentation



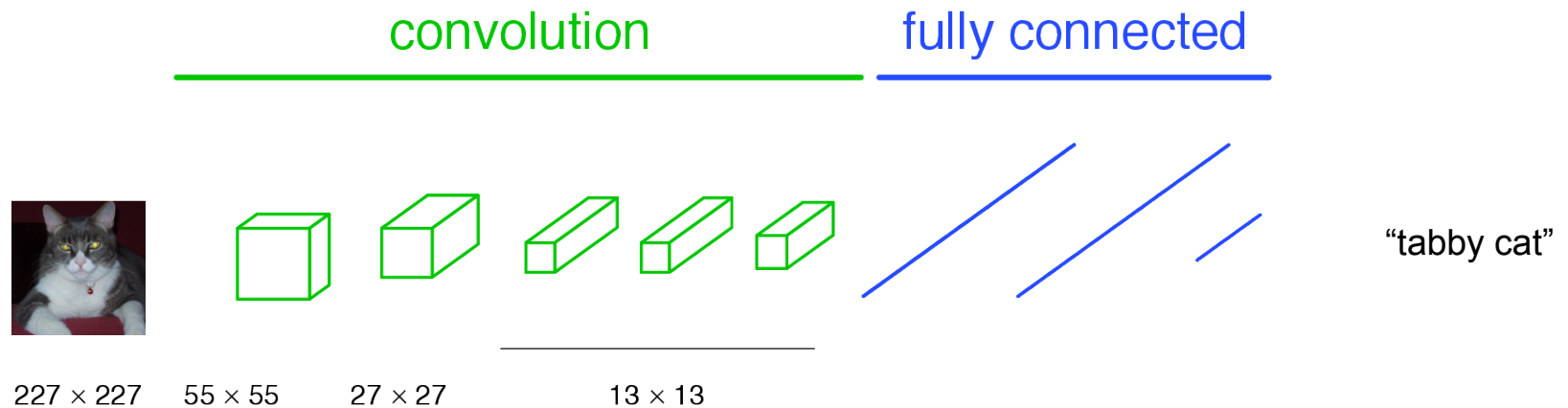
boundary prediction (Xie & Tu 2015)



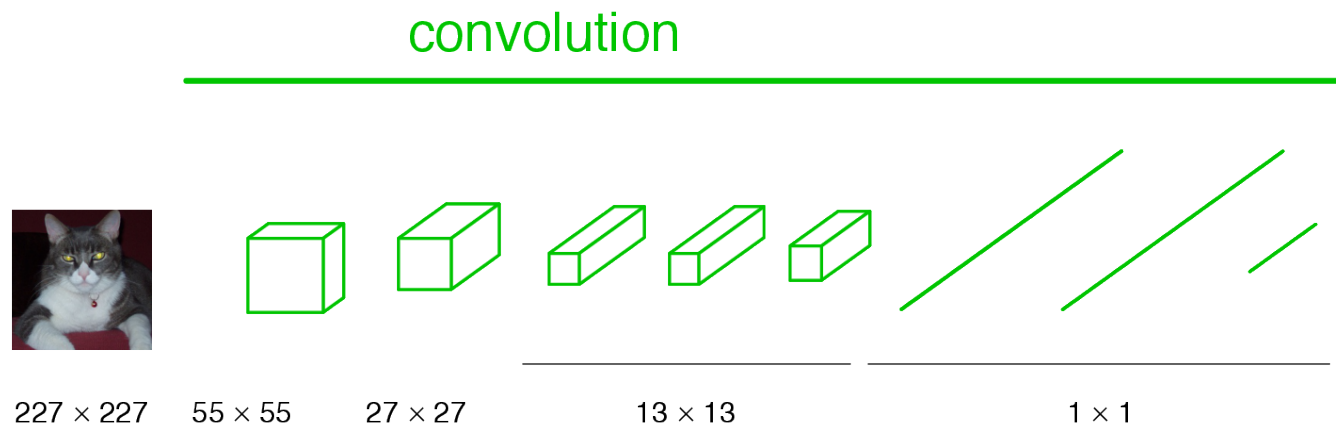
< 1/5 second



a classification network



becoming fully convolutional



becoming fully convolutional

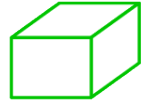
convolution



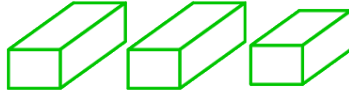
$H \times W$



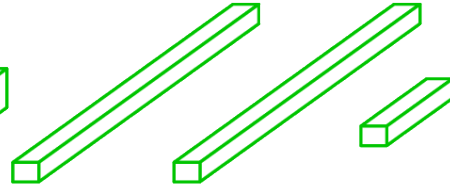
$H/4 \times W/4$



$H/8 \times W/8$



$H/16 \times W/16$



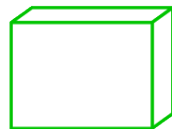
$H/32 \times W/32$

upsampling output

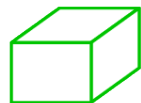
convolution



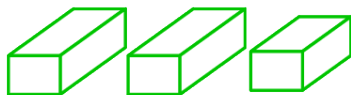
$H \times W$



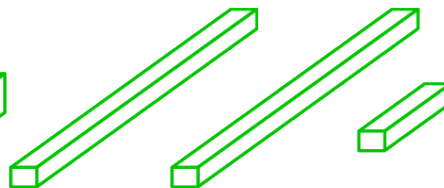
$H/4 \times W/4$



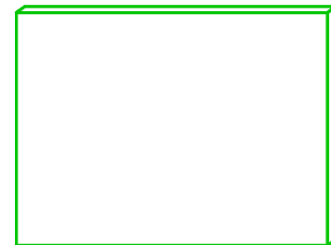
$H/8 \times W/8$



$H/16 \times W/16$



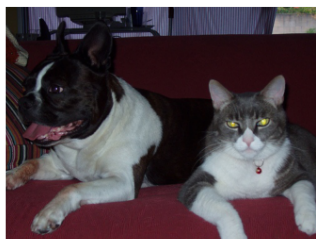
$H/32 \times W/32$



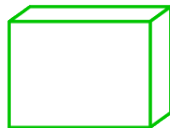
$H \times W$

end-to-end, pixels-to-pixels network

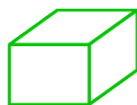
convolution



$H \times W$



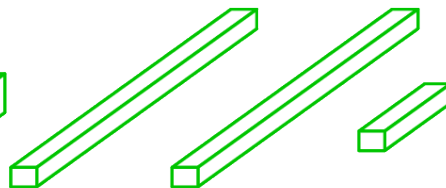
$H/4 \times W/4$



$H/8 \times W/8$



$H/16 \times W/16$

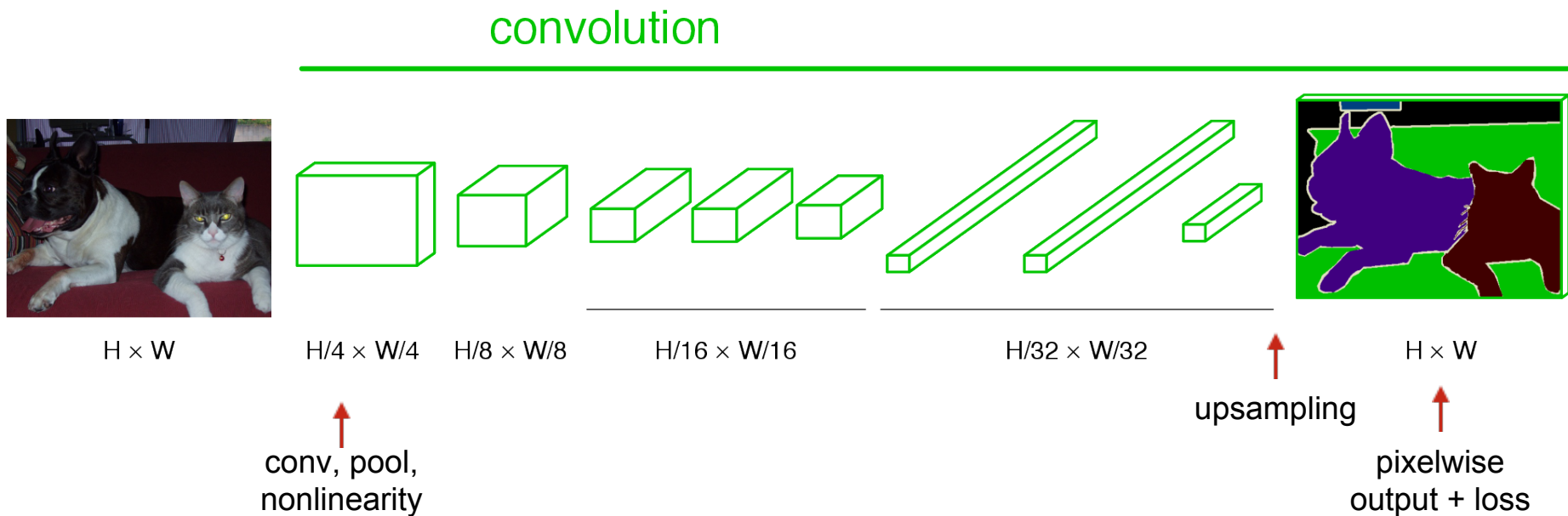


$H/32 \times W/32$



$H \times W$

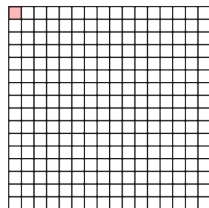
end-to-end, pixels-to-pixels network



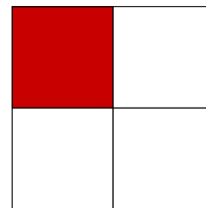
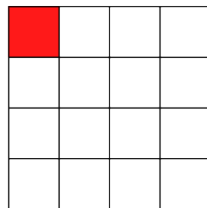
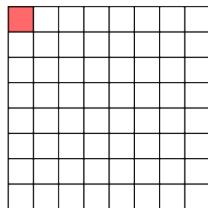
spectrum of deep features

combine *where* (local, shallow) with *what* (global, deep)

image



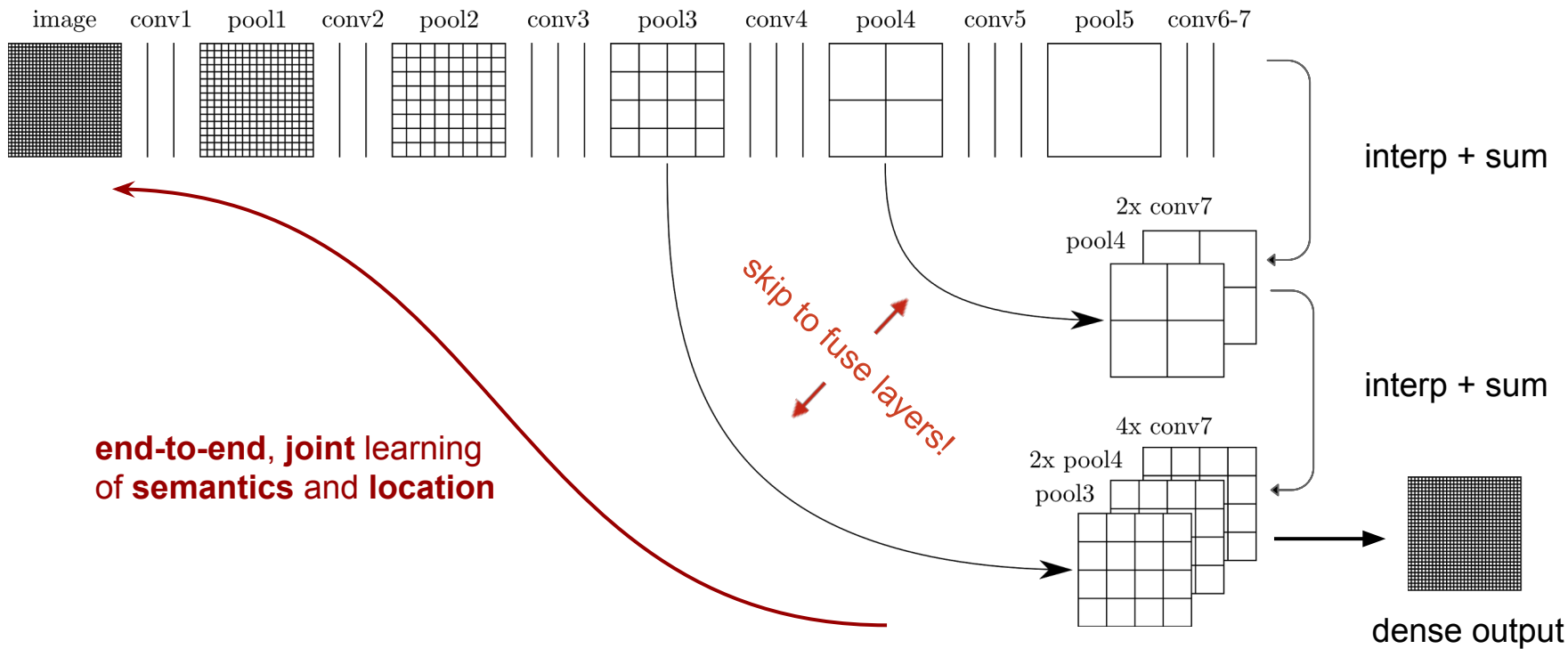
intermediate layers



fuse features into **deep jet**

(cf. Hariharan et al. CVPR15 “hypercolumn”)

skip layers



skip layer refinement

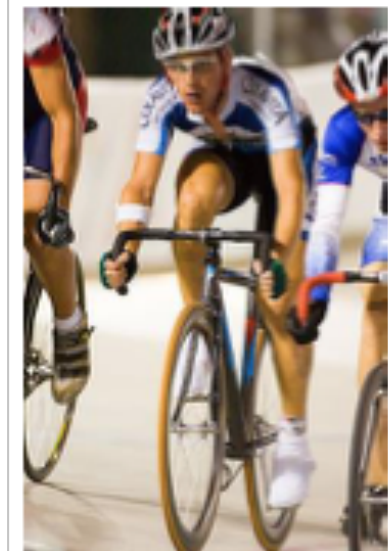
input image

stride 32

stride 16

stride 8

ground truth



no skips



1 skip

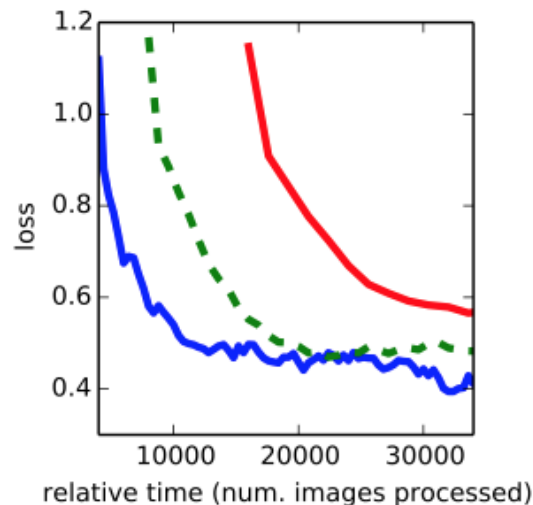
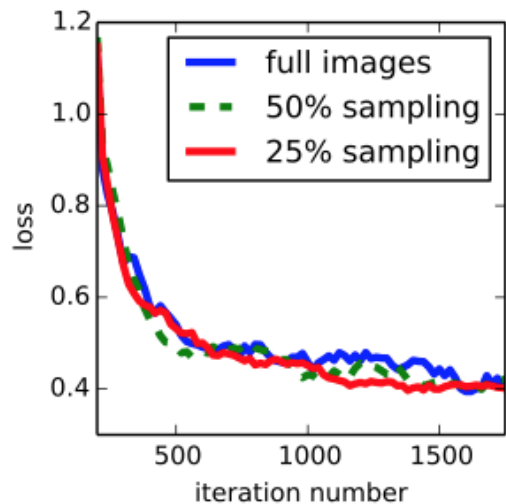


2 skips



training + testing

- train full image at a time *without patch sampling*
- reshape network to take input of any size
- forward time is ~150ms for 500 x 500 x 21 output

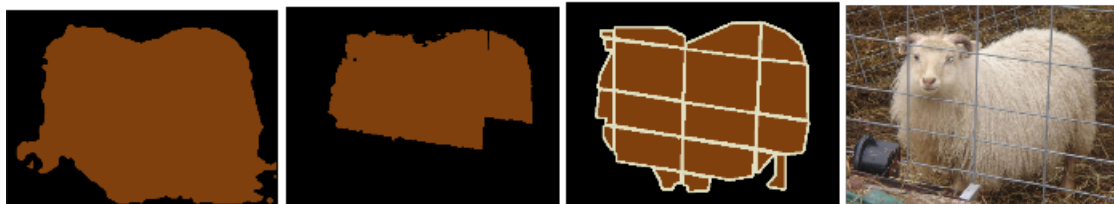
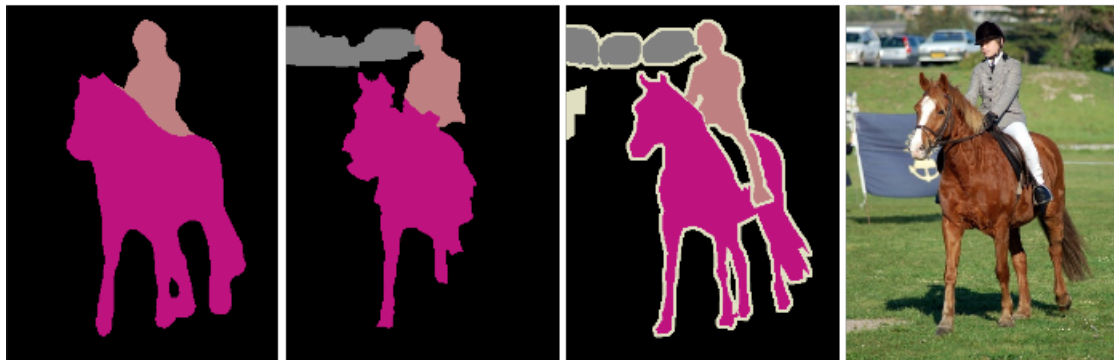


FCN

SDS*

Truth

Input



Relative to prior state-of-the-art SDS:

- 20% improvement for mean IoU
- 286× faster

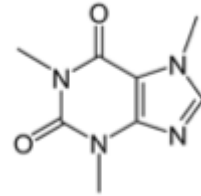
*Simultaneous Detection and Segmentation
Hariharan et al. ECCV14

models + code

fully convolutional networks are fast, end-to-end models for pixelwise problems

- **code** in Caffe branch (merged soon)
- **models** for PASCAL VOC, NYUDv2, SIFT Flow, PASCAL-Context in Model Zoo

fcn.berkeleyvision.org



caffe.berkeleyvision.org



github.com/BVLC/caffe

models

- **PASCAL VOC** standard for object segmentation
- **NYUDv2** multi-modal rgb + depth scene segmentation
- **SIFT Flow** multi-task for semantic + geometric segmentation
- **PASCAL-Context** object + scene segmentation

inference

eval.py

Raw

```
1 import numpy as np
2 from PIL import Image
3
4 import caffe
5
6 # load image, switch to BGR, subtract mean, and make dims C x H x W for Caffe
7 im = Image.open('pascal/VOC2010/JPEGImages/2007_000129.jpg')
8 in_ = np.array(im, dtype=np.float32)
9 in_ = in_[::-1]
10 in_ -= np.array((104.00698793, 116.66876762, 122.67891434))
11 in_ = in_.transpose((2,0,1))
12
13 # load net
14 net = caffe.Net('deploy.prototxt', 'fcn-32s-pascalcontext.caffemodel', caffe.TEST)
15 # shape for input (data blob is N x C x H x W), set data
16 net.blobs['data'].reshape(1, *in_.shape)
17 net.blobs['data'].data[...] = in_
18 # run net and take argmax for prediction
19 net.forward()
20 out = net.blobs['score'].data[0].argmax(axis=0)
```

[inference script \(gist\)](#)

solving

```
31 # base net -- follow the editing model parameters example to make
32 # a fully convolutional VGG16 net.
33 # http://nbviewer.ipython.org/github/BVLC/caffe/blob/master/examples/net_surgery.ipynb
34 base_weights = 'vgg16fc.caffemodel'
35
36 # init
37 caffe.set_mode_gpu()
38 caffe.set_device(0)
39
40 solver = caffe.SGDSolver('solver.prototxt')
41
42 # do net surgery to set the deconvolution weights for bilinear interpolation
43 interp_layers = [k for k in solver.net.params.keys() if 'up' in k]
44 interp_surgery(solver.net, interp_layers)
45
46 # copy base weights for fine-tuning
47 solver.net.copy_from(base_weights)
48
49 # solve straight through -- a better approach is to define a solving loop to
50 # 1. take SGD steps
51 # 2. score the model by the test net `solver.test_nets[0]`
52 # 3. repeat until satisfied
53 solver.step(80000)
```

[solving script \(gist\)](#)

Reshape

- Decide shape on-the-fly in C++ / Python / MATLAB
- DataLayer automatically reshapes
for batch size == 1
- Essentially free
(only reallocates when necessary)

Helpful Layers

- Losses can take spatial predictions + truths
- Deconvolution / “backward convolution”
can compute interpolation
- Crop: maps coordinates between layers

FCN for Pose Estimation

Georgia Gkioxari
UC Berkeley

FCN for Pose Estimation

Input data:

Image



FCN for Pose Estimation

Input data:

Image



Keypoints



FCN for Pose Estimation

Input data:

Image



Keypoints



Define an area around the keypoint as its positive neighborhood with radius r .

FCN for Pose Estimation

Input data:

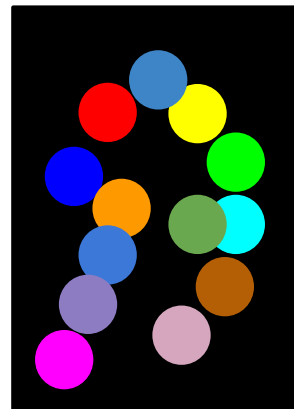
Image



Keypoints



Labels



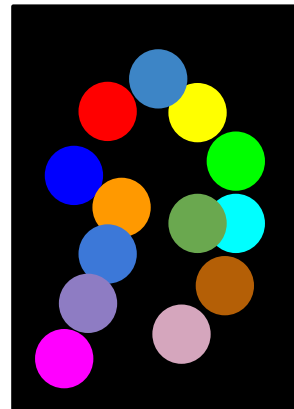
FCN for Pose Estimation

Input data:

Image



Labels

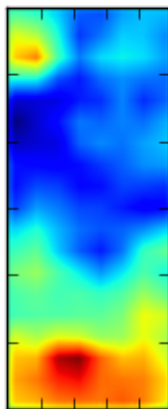


Heat Map Predictions from FCN

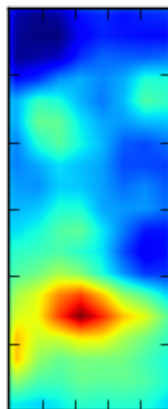
Test Image



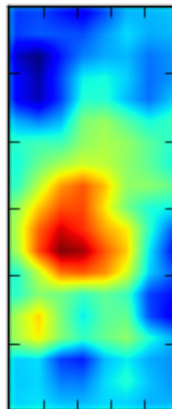
Right Ankle



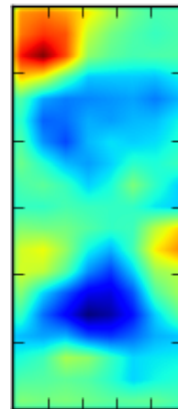
Right Knee



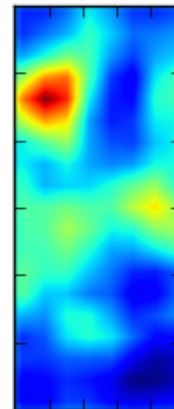
Right Hip



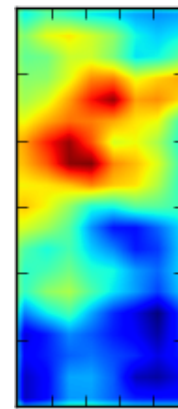
Right Wrist



Right Elbow



Right Shoulder

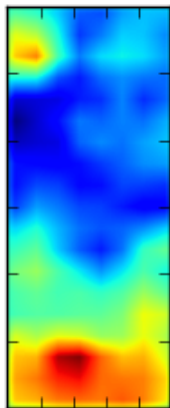


Heat Map Predictions from FCN

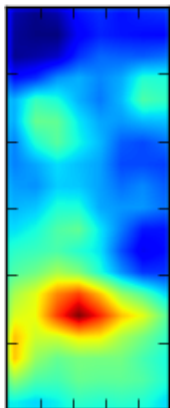
Test Image



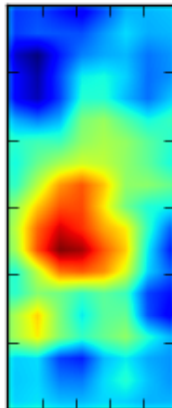
Right Ankle



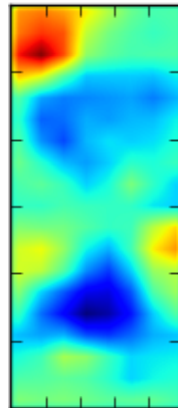
Right Knee



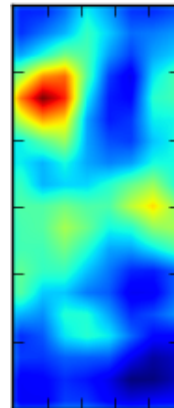
Right Hip



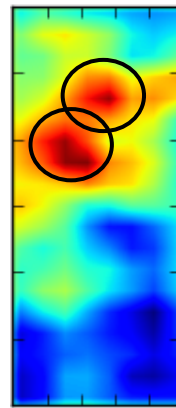
Right Wrist



Right Elbow



Right Shoulder



Two modes because there are two Right Shoulders in the image!

Heat Maps to Keypoints

PCK @ 0.2	LSP test set
Ankle	56.5
Knee	60.0
Hip	56.6
Wrist	62.9
Elbow	71.8
Shoulder	78.8
Head	93.6

FCN baseline PCK == ~69%

State-of-the-art == ~72%

Details

Architecture:

- FCN - 32 stride. **No** data augmentation.
- radius = $0.1 * \text{im.shape}[0]$ (**no** cross validation)

Runtime on a K40:

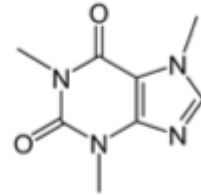
- 0.7 sec/iteration for training (15hrs for 80K iterations)
- 0.25 sec/image for inference for all keypoints

conclusion

fully convolutional networks are fast, end-to-end models for pixelwise problems

- **code** in Caffe branch (merged soon)
- **models** for PASCAL VOC, NYUDv2, SIFT Flow, PASCAL-Context in Model Zoo

fcn.berkeleyvision.org



caffe.berkeleyvision.org



github.com/BVLC/caffe