

Image Segmentation

Joe Yeh, M.D.

What is image segmentation?

- Add your first bullet point here
- Add your second bullet point here
- Add your third bullet point here

CNN for Image Segmentation

- Fully convolutional network (<https://arxiv.org/abs/1411.4038v2>)
- U-Net (<https://arxiv.org/abs/1505.04597v1>)
- Mask R-CNN (<https://arxiv.org/abs/1703.06870v3>)

[Computer Science](#) > [Computer Vision and Pattern Recognition](#)

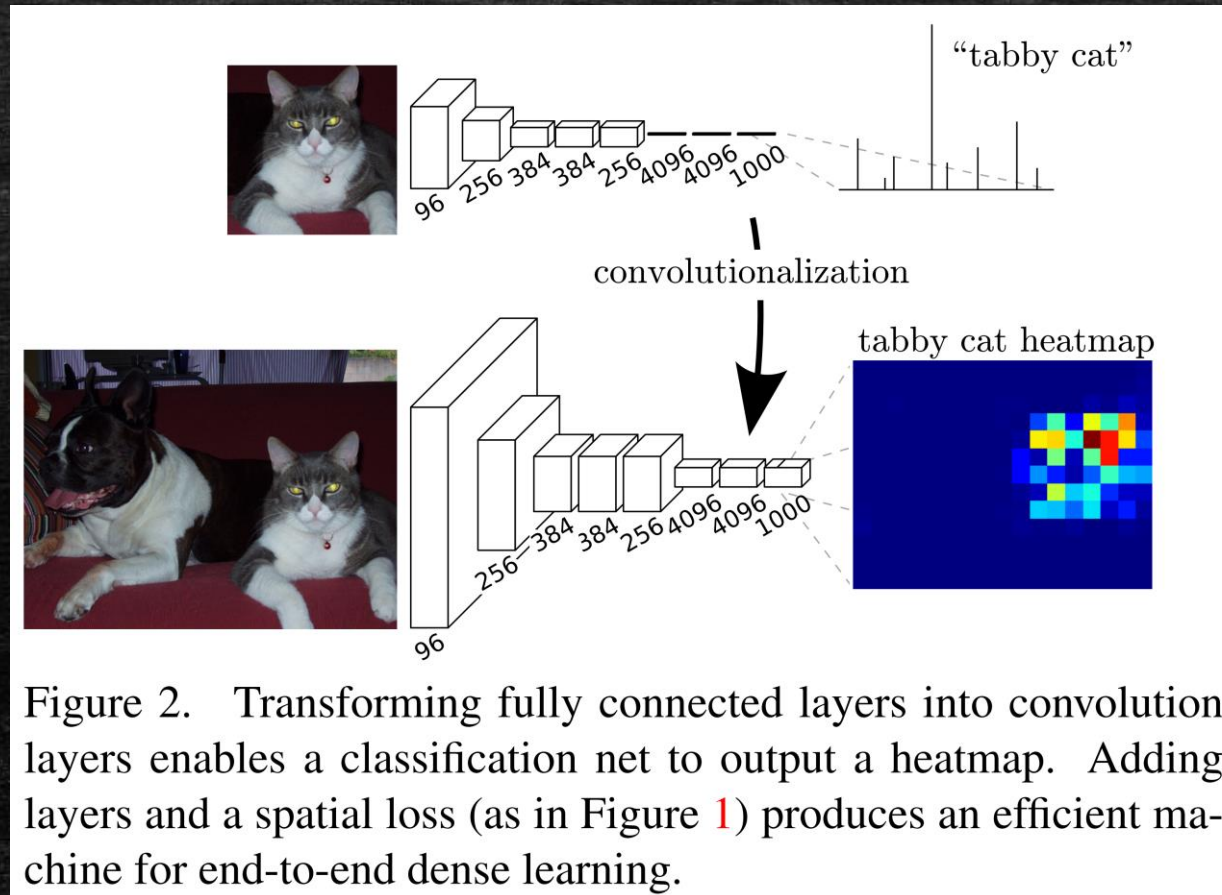
Fully Convolutional Networks for Semantic Segmentation

[Jonathan Long](#), [Evan Shelhamer](#), [Trevor Darrell](#)

(Submitted on 14 Nov 2014 (v1), last revised 8 Mar 2015 (this version, v2))

Convolutional networks are powerful visual models that yield hierarchies of features. We show that convolutional networks by themselves, trained end-to-end, pixels-to-pixels, exceed the state-of-the-art in semantic segmentation. Our key insight is to build "fully convolutional" networks that take input of arbitrary size and produce correspondingly-sized output with efficient inference and learning. We define and detail the space of fully convolutional networks, explain their application to spatially dense prediction tasks, and draw connections to prior models. We adapt contemporary classification networks (AlexNet, the VGG net, and GoogLeNet) into fully convolutional networks and transfer their learned representations by fine-tuning to the segmentation task. We then define a novel architecture that combines semantic information from a deep, coarse layer with appearance information from a shallow, fine layer to produce accurate and detailed segmentations. Our fully convolutional network achieves state-of-the-art segmentation of PASCAL VOC (20% relative improvement to 62.2% mean IU on 2012), NYUDv2, and SIFT Flow, while inference takes one third of a second for a typical image.

Convolutionalization of classification network



Combining coarse and fine prediction

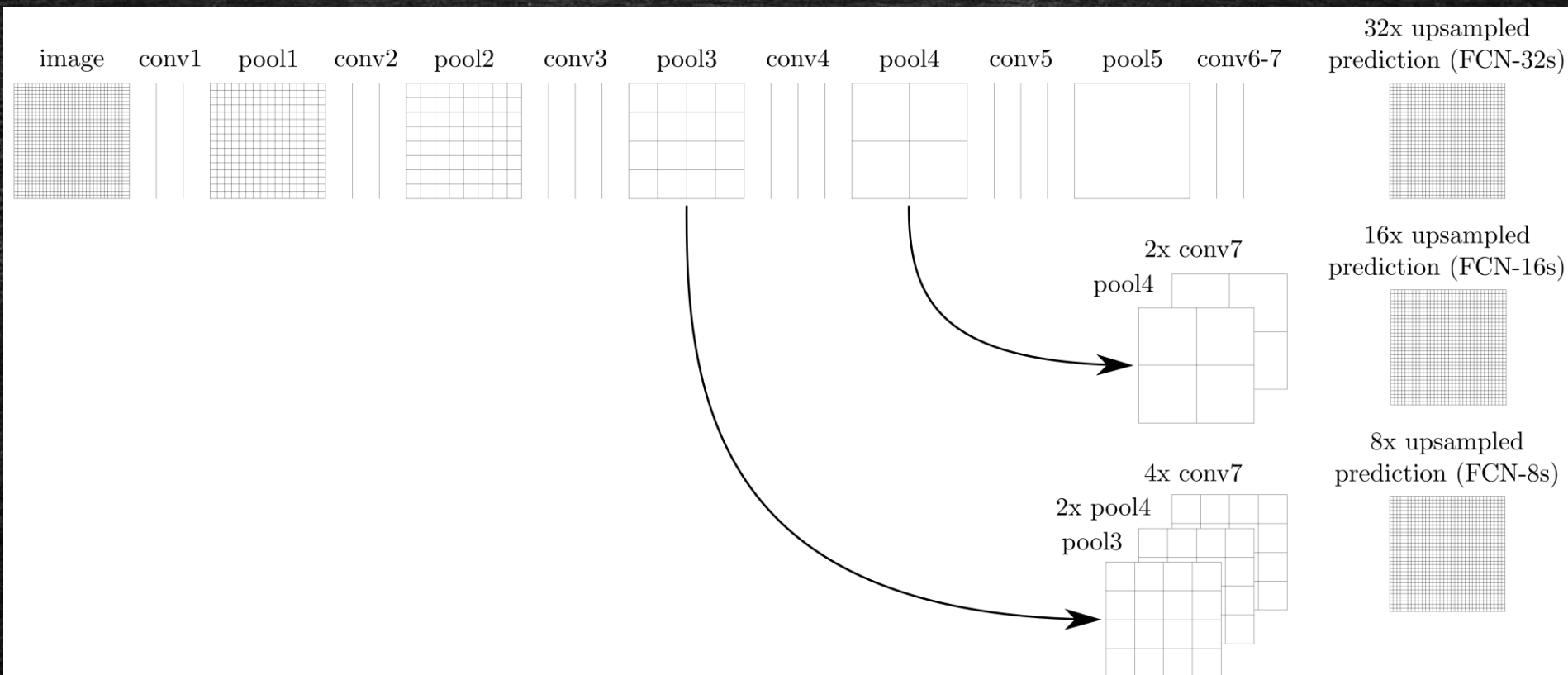


Figure 3. Our DAG nets learn to combine coarse, high layer information with fine, low layer information. Pooling and prediction layers are shown as grids that reveal relative spatial coarseness, while intermediate layers are shown as vertical lines. First row (FCN-32s): Our single-stream net, described in Section 4.1, upsamples stride 32 predictions back to pixels in a single step. Second row (FCN-16s): Combining predictions from both the final layer and the pool4 layer, at stride 16, lets our net predict finer details, while retaining high-level semantic information. Third row (FCN-8s): Additional predictions from pool3, at stride 8, provide further precision.

Segmentation results

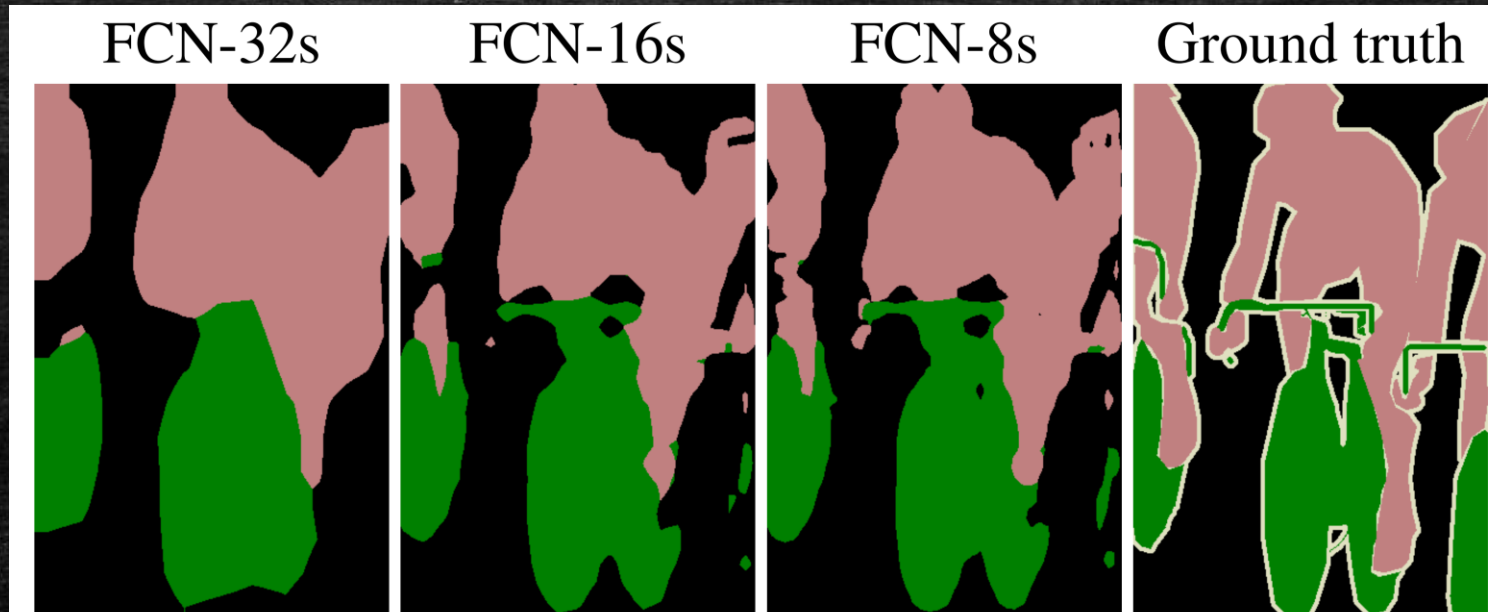


Figure 4. Refining fully convolutional nets by fusing information from layers with different strides improves segmentation detail. The first three images show the output from our 32, 16, and 8 pixel stride nets (see Figure 3).

U-Net

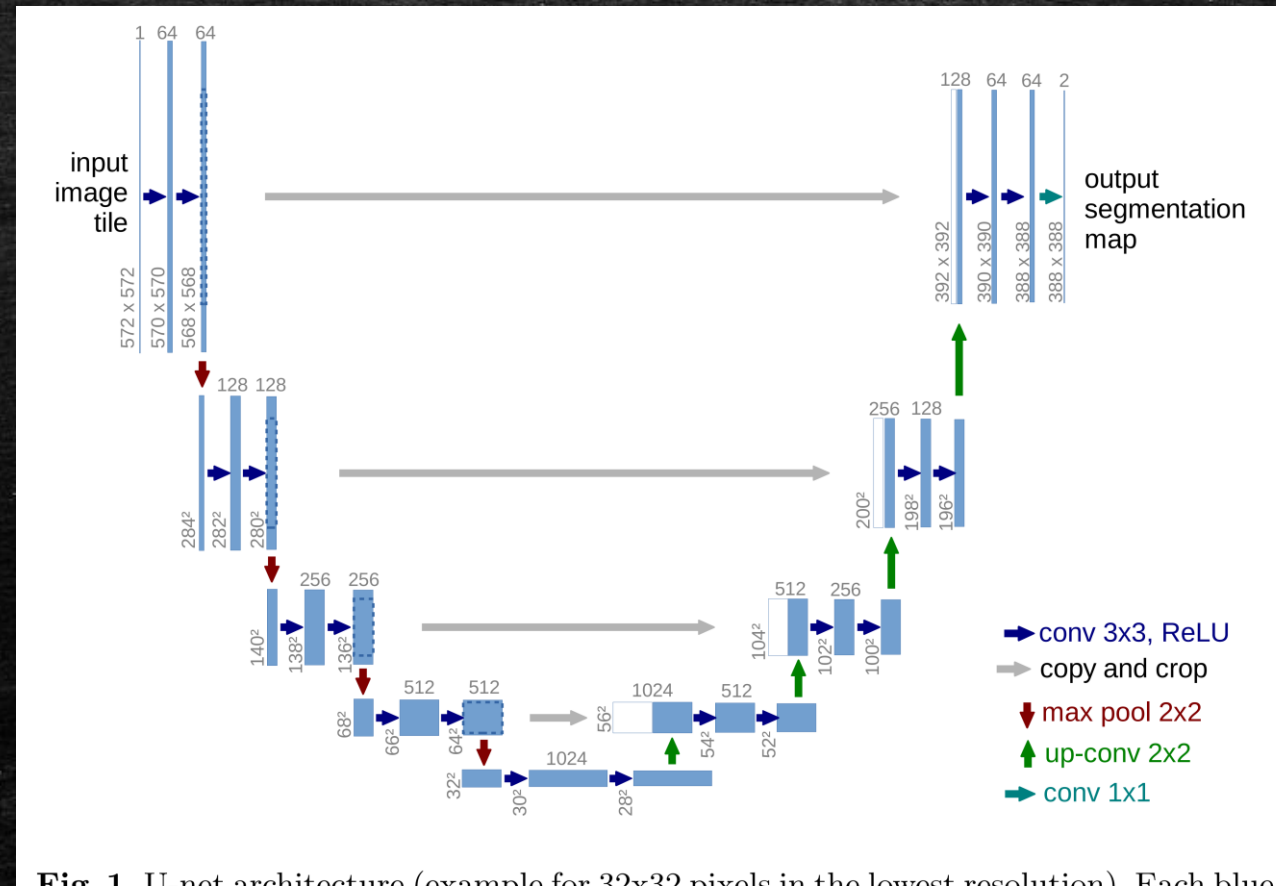


Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.