# CSEP 576: Dense Prediction



Jonathan Huang (jonathanhuang@google.com)

University of Washington 26 May 2020

**Google Research**

# Lecture Outline

**Dense Prediction (pixel level prediction)**

- Semantic Segmentation

- Instance Segmentation
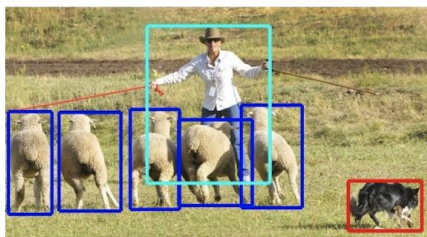
- Panoptic Segmentation

- Keypoint Estimation

We will mainly focus on semantic segmentation as a way to introduce some of technical details behind "dense prediction"

# Problem statement



classify

person, sheep, dog



classify and regress
bounding box per object

**(bounding box)
detection**



classify per pixel

**semantic
segmentation**

# Segmentation Applications



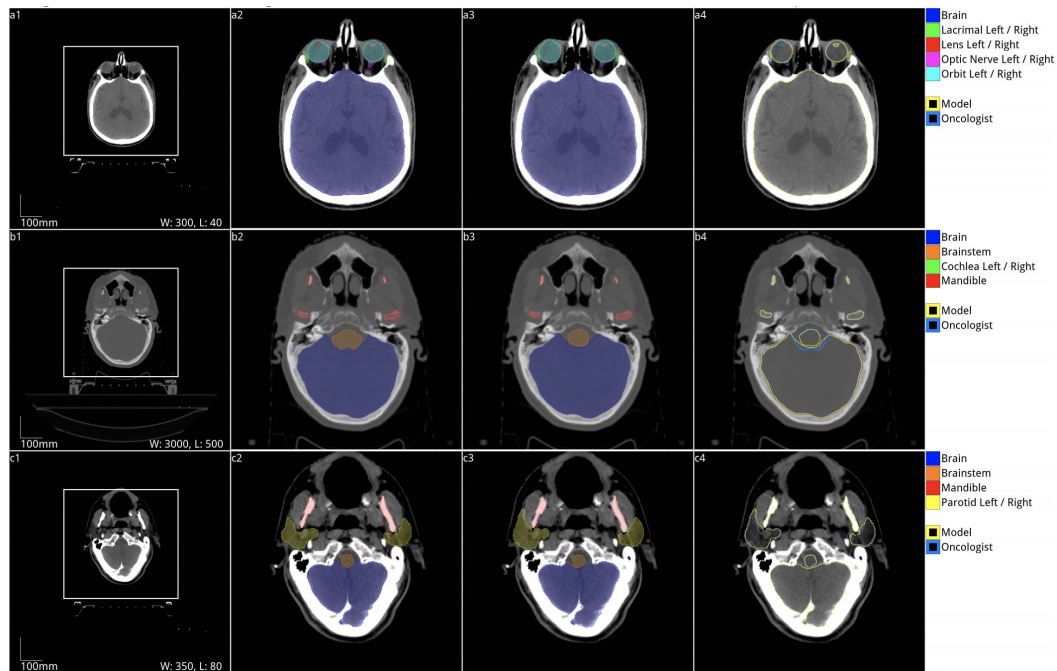Original          Segmentation map          Final

# Segmentation Applications

Water

Forest

Field

Impervious

Developed, Open space

Developed, Low intensity

Developed, Medium intensity

Developed, High intensity

Deciduous Forest

Evergreen Forest

Shrub/Scrub

Cultivated Crops

# Medical Segmentation

# Outline of Semantic Segmentation

- The sliding window connection (again)

- Fully Convolutional models

- How to get high resolution outputs with

    - Atrous convolutions

    - "Upconvolutions"

- Target Assignment

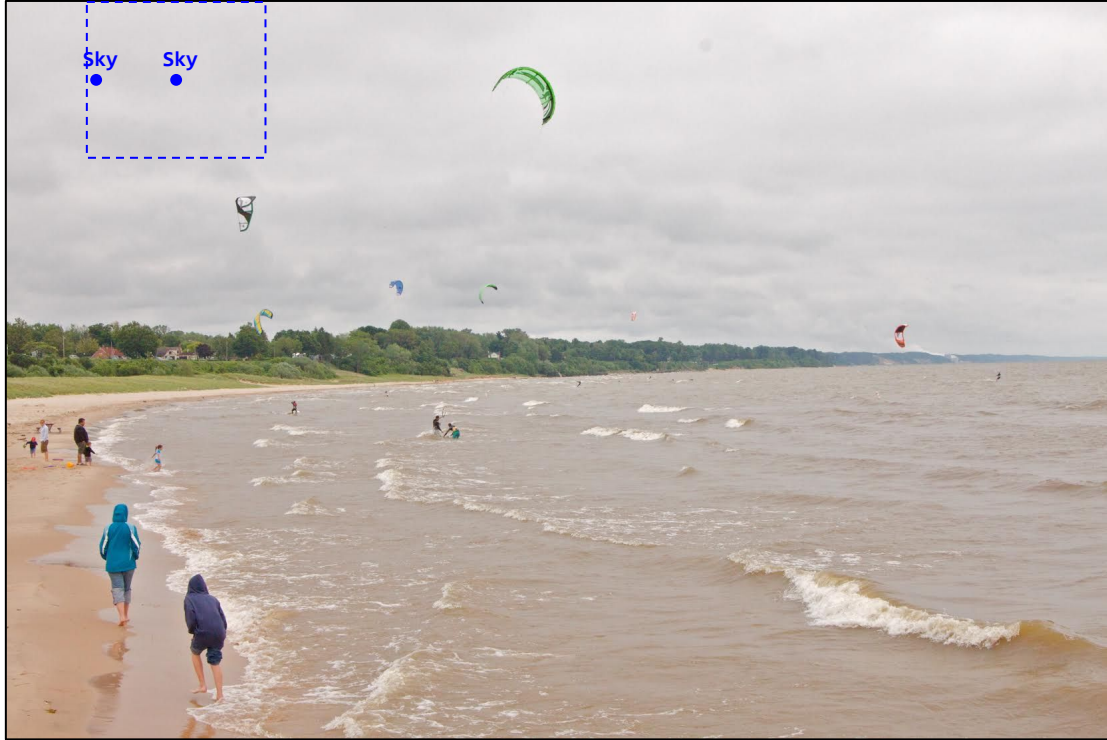- Evaluation of Semantic Segmentation

Relevant for all dense prediction tasks

# "Sliding Window" Segmentation

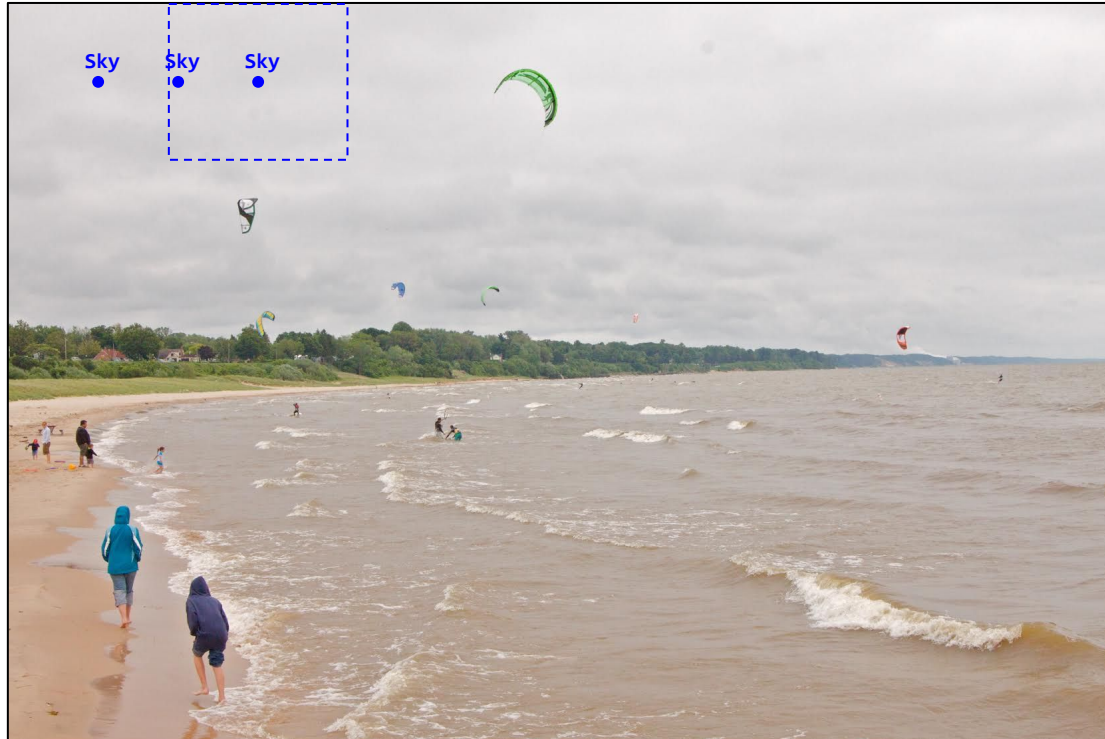

**Same idea as detection**: Extract features from a window around a point; Predict class label for point

# "Sliding Window" Segmentation



**Same idea as detection**: Extract features from a window around a point; Predict class label for point

# "Sliding Window" Segmentation



**Same idea as detection**: Extract features from a window around a point; Predict class label for point

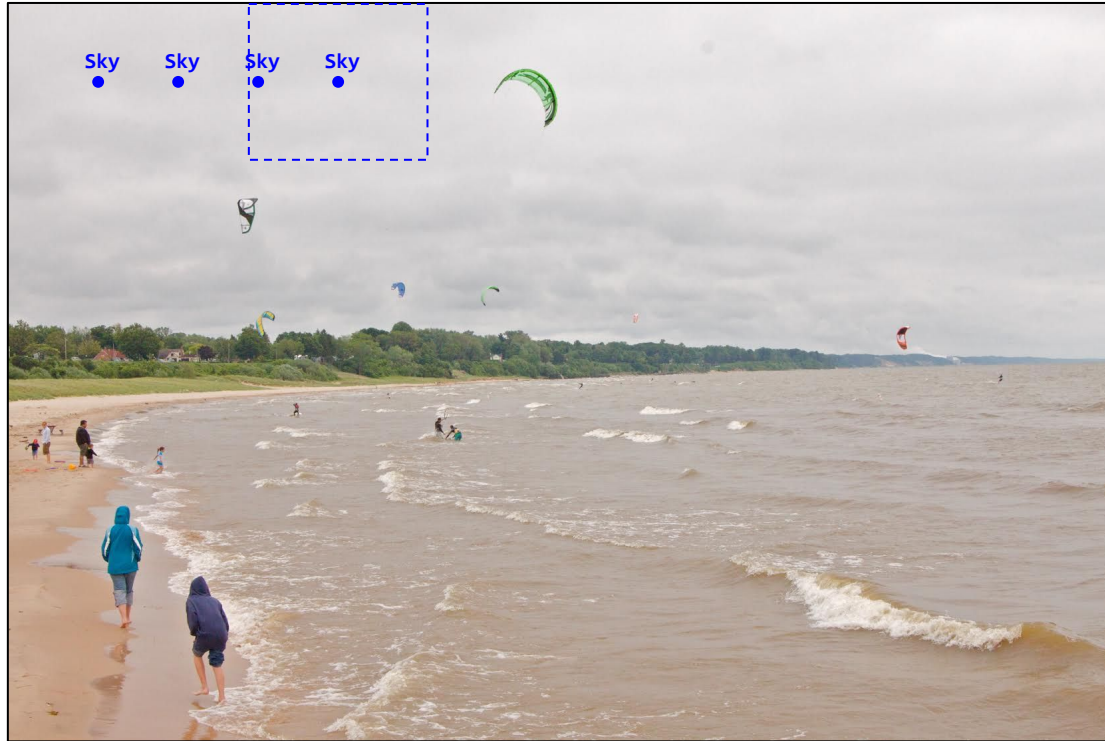# "Sliding Window" Segmentation
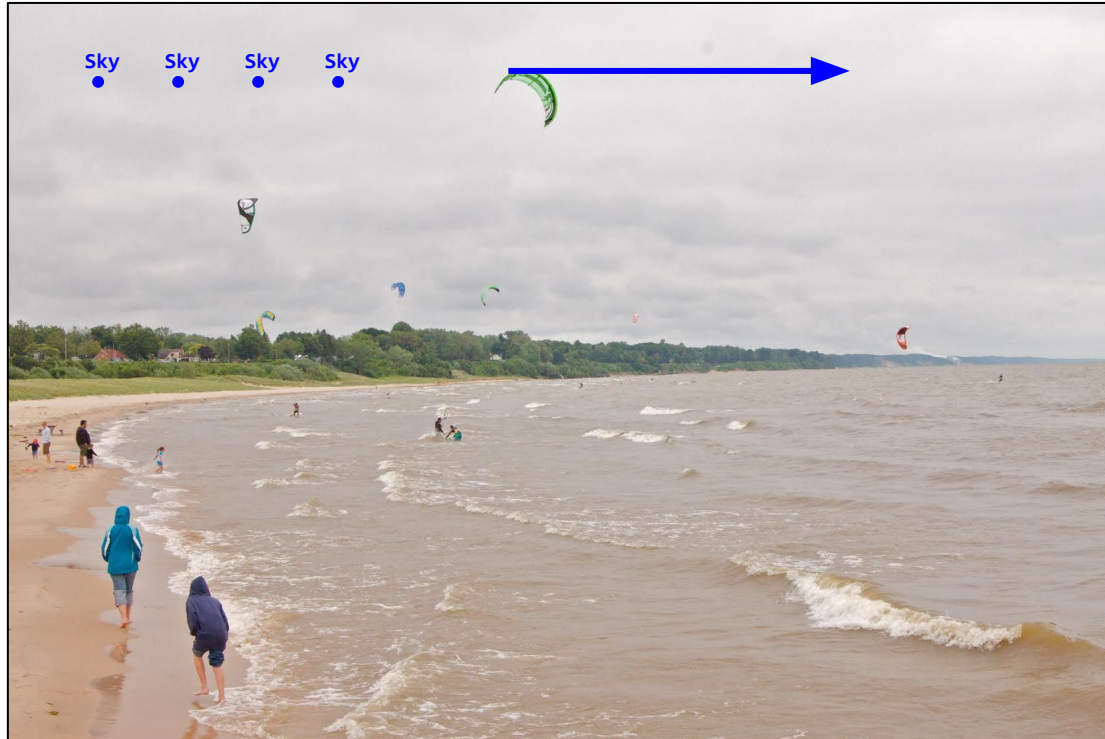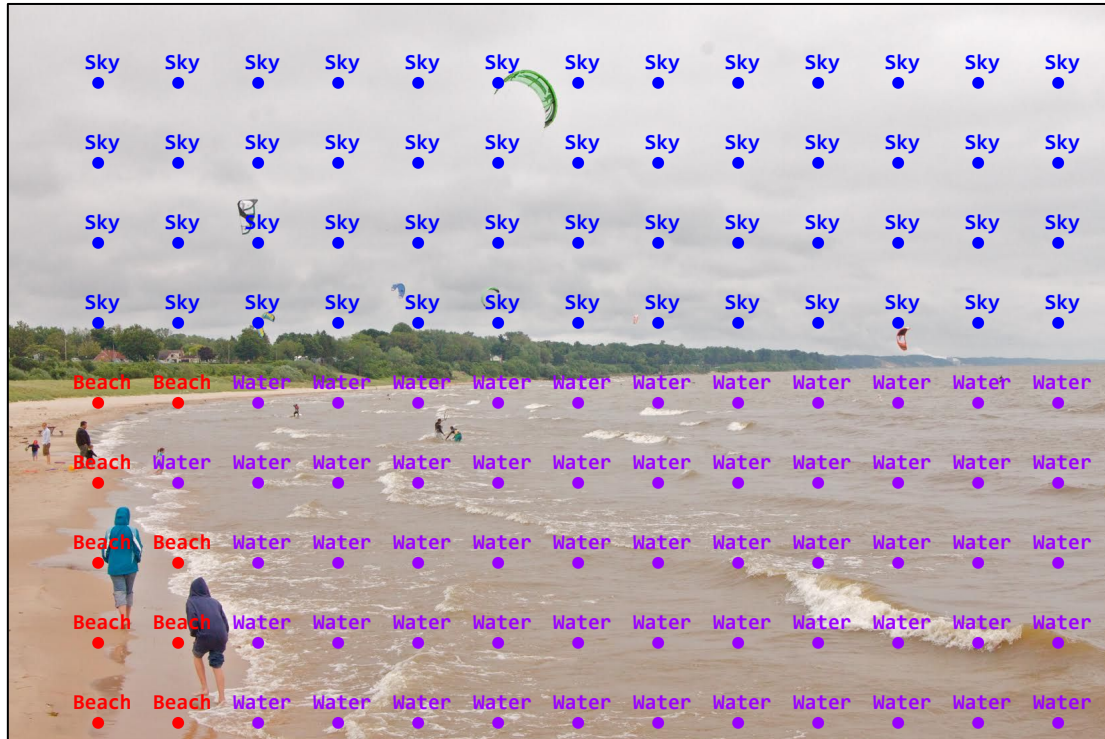


**Same idea as detection**: Extract features from a window around a point; Predict class label for point

# "Sliding Window" Segmentation



**Same idea as detection**: Extract features from a window around a point; Predict class label for point

# "Sliding Window" Segmentation



**Same idea as detection**: Extract features from a window around a point; Predict class label for point

# Fully Convolutional Networks - Standard for detection / segmentation / keypoint prediction

**"Fully Convolutional"**: All layers operate on local inputs (e.g. Conv, Pool, ReLU); E.g. no FC layers allowed.

**Properties of FCNs**:

- Operate on input of any size
- Output tensors scale with input size
- Can train with heterogenous resolutions
- Can train and test at different resolutions

# Fully Convolutional Networks - Standard for detection / segmentation / keypoint prediction

**"Fully Convolutional"**: All layers operate on local inputs (e.g. Conv, Pool, ReLU); E.g. no FC layers allowed.
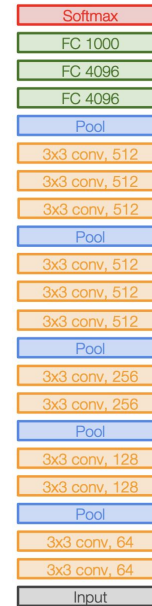
**Properties of FCNs**:

- Operate on input of any size
- Output tensors scale with input size
- Can train with heterogenous resolutions
- Can train and test at different resolutions

| |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

**[7x7x512] "pool5" given 224x224 inputs**

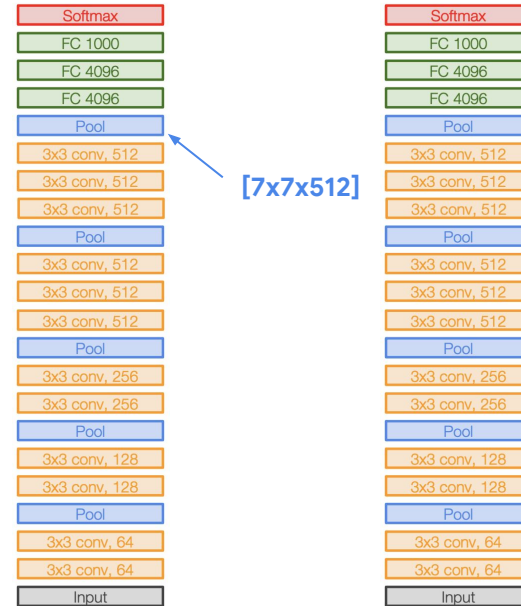**A VGG-16 "non-example" (that is still illustrative)**

VGG trained on 224x224 images

# Fully Convolutional Networks - Standard for detection / segmentation / keypoint prediction

**"Fully Convolutional"**: All layers operate on local inputs (e.g. Conv, Pool, ReLU); E.g. no FC layers allowed.

**Properties of FCNs**:

- Operate on input of any size
- Output tensors scale with input size
- Can train with heterogenous resolutions
- Can train and test at different resolutions

| | |
|---|---|
| Softmax | Softmax |
| FC 1000 | FC 1000 |
| FC 4096 | FC 4096 |
| FC 4096 | FC 4096 |
| Pool | Pool |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| Pool | Pool |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| 3x3 conv, 512 | 3x3 conv, 512 |
| Pool | Pool |
| 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 256 | 3x3 conv, 256 |
| Pool | Pool |
| 3x3 conv, 128 | 3x3 conv, 128 |
| 3x3 conv, 128 | 3x3 conv, 128 |
| Pool | Pool |
| 3x3 conv, 64 | 3x3 conv, 64 |
| 3x3 conv, 64 | 3x3 conv, 64 |
| Input | Input |

[7x7x512]

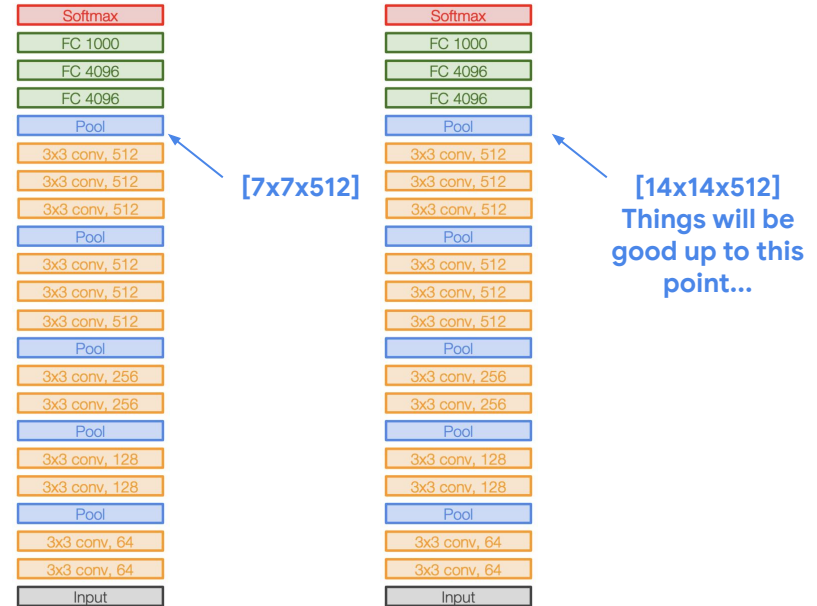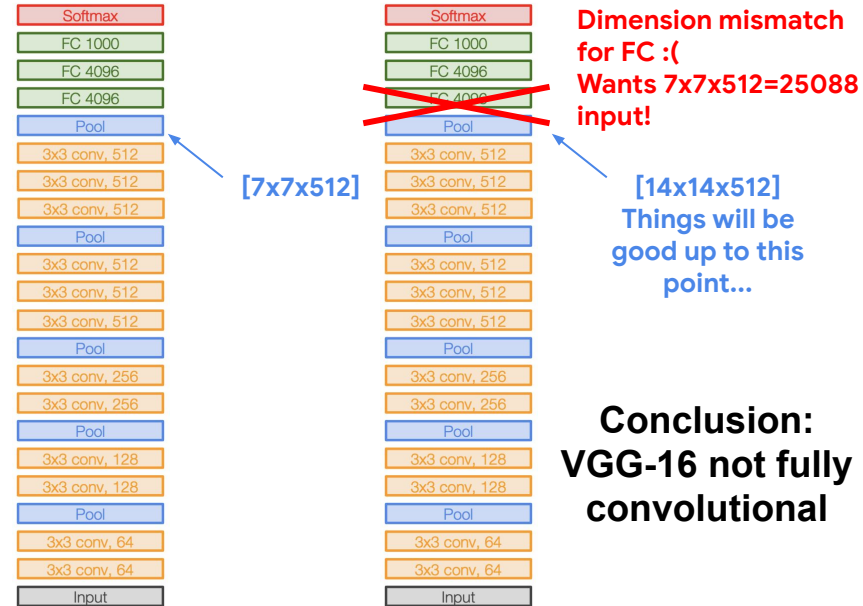VGG trained on 224x224 images

*What if we try running inference 448x448 image?*

# Fully Convolutional Networks - Standard for detection / segmentation / keypoint prediction

**"Fully Convolutional"**: All layers operate on local inputs (e.g. Conv, Pool, ReLU); E.g. no FC layers allowed.
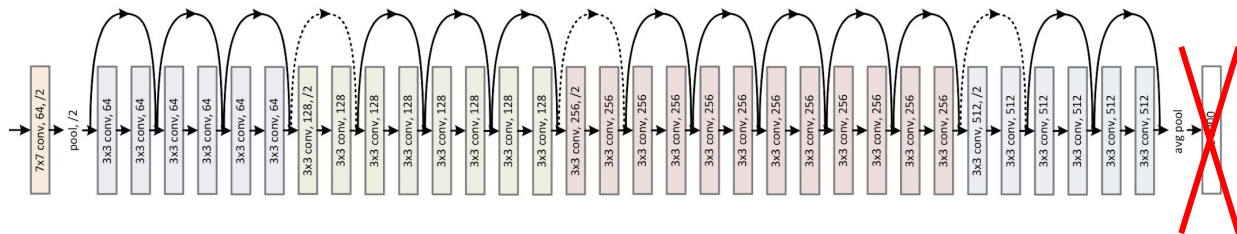
**Properties of FCNs:**

- Operate on input of any size
- Output tensors scale with input size
- Can train with heterogenous resolutions
- Can train and test at different resolutions

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512          **[7x7x512]**
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

VGG trained on 224x224 images

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512          **[14x14x512]**
3x3 conv, 512          **Things will be good up to this point...**
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

What if we try running inference 448x448 image?

# Fully Convolutional Networks - Standard for detection / segmentation / keypoint prediction

**"Fully Convolutional"**: All layers operate on local inputs (e.g. Conv, Pool, ReLU); E.g. no FC layers allowed.

**Properties of FCNs:**

- Operate on input of any size
- Output tensors scale with input size
- Can train with heterogenous resolutions
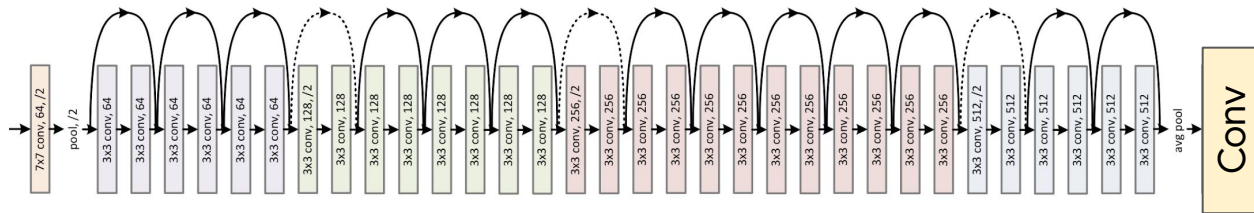- Can train and test at different resolutions

| Softmax |
|---|
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

[7x7x512]

VGG trained on 224x224 images

| Softmax |
|---|
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

**Dimension mismatch for FC :( Wants 7x7x512=25088 input!**

[14x14x512] Things will be good up to this point...

**Conclusion: VGG-16 not fully convolutional**

What if we try running inference 448x448 image?

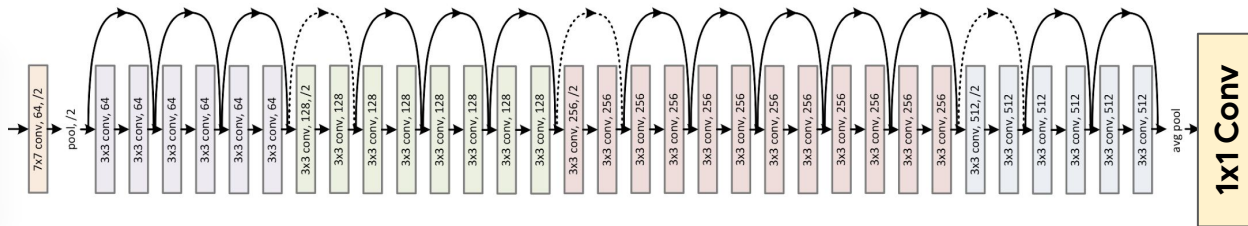# Ways to get an FCN (from an existing non-FCN)



Option 1: Chop off FC (and pooling layers) at top
(and possibly add new convs)



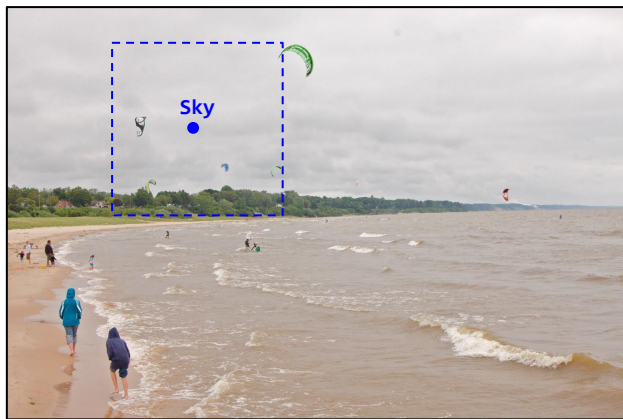Option 2: Convert FC layers to "Equivalent" Convs

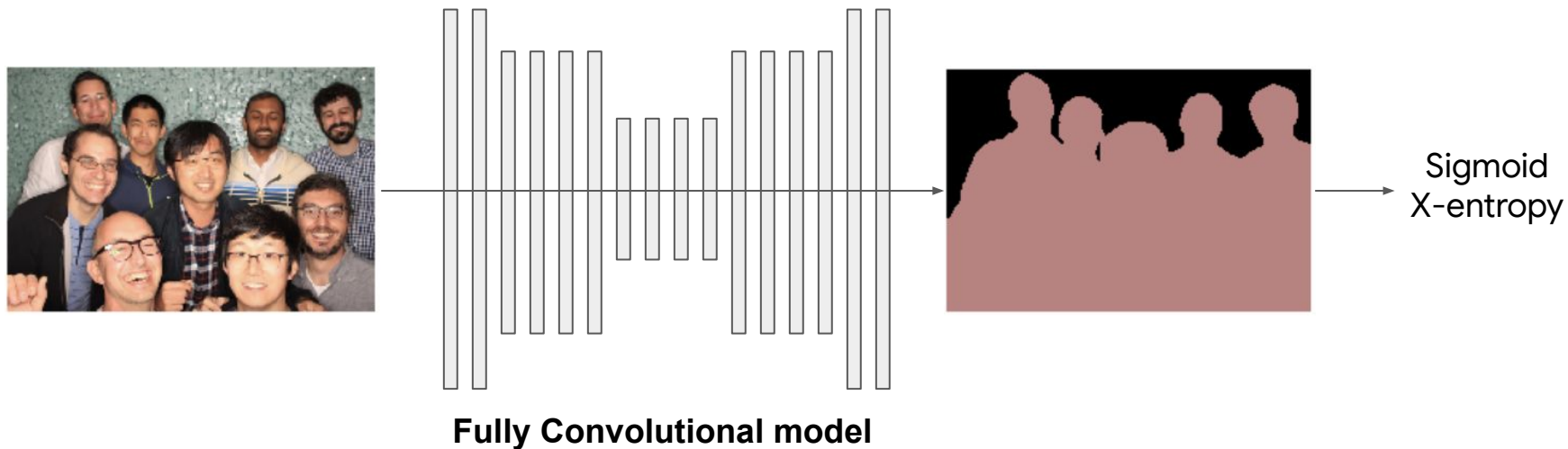# Ways to get an FCN (from an existing non-FCN)

**224x224**



Convert top FC layer to Conv layer that takes full extent of input (in this case, FC is 1x1 with 1000 output channels)

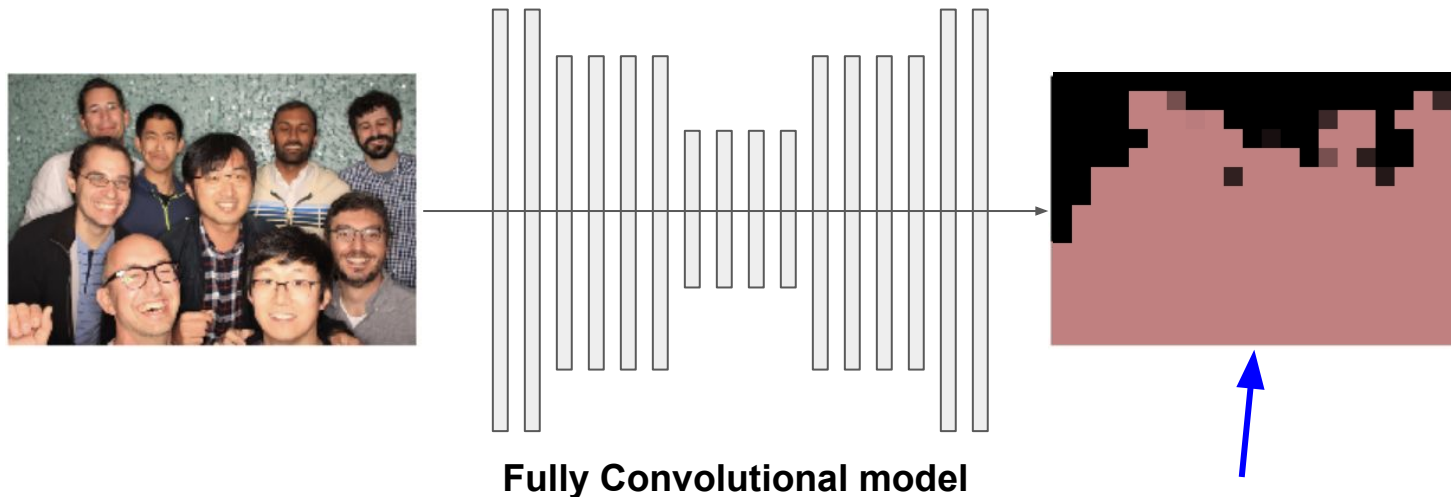Note: w/o the avg pool, we'd convert the FC to a 7x7 conv with 1000 output channels

Now can run network on much larger image (even after training!)

Sky

# Typical Semantic Segmentation model



**Fully Convolutional model**

Sigmoid
X-entropy

- Run image through FCN
- Train with per-pixel sigmoid X-entropy

# Typical Semantic Segmentation model



**Fully Convolutional model**

- Run image through FCN
- Train with per-pixel sigmoid X-entropy

**But**: if we directly convert typical classification model (e.g. VGG) to FCN, we'll get something like this :(
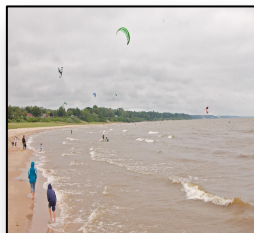
# Typical CNN Output sizes are too small



**Stride 2 ops (5 ops)**

640x640

20x20

Total Network Stride = 2^5 = 32;
Output size = (640/32) x (640x32) = 20x20
Too small!! :(

- Network stride = product of layer strides (for single path network)
  - For typical ImageNet networks (e.g. AlexNet, VGG, Resnet) stride prior to FC layers is 32
- For segmentation we typically want smaller network stride (e.g. 2, 4 or 8)
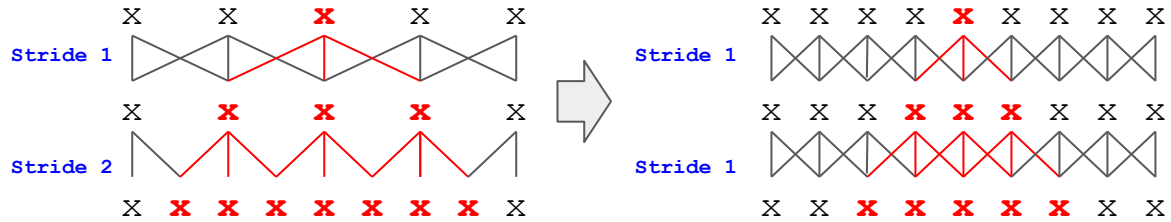
# How to get high resolution outputs (e.g. w/stride < 32)

- **Use fewer stride 2 convolutions**

- Use "upconvolution" operators

# Approach 1: Just don't downsample that many times



7x7 conv, 64, /2 · pool, /2 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 128, /2 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 256, /2 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 512, /2 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512 · avg pool

**1x1 Conv**

**80x80**

**Make these stride 2 ops
stride 1 instead**

Resulting network stride: 8

# Replace stride 2 convolutions with stride 1



**Problem**: Doing this directly can significantly reduce receptive field size...

Chen, Papandreou, et al, 2015

# Some Receptive Field arithmetic

How big is our receptive field?

Kernel size at layer l

Product of strides up to layer l

Receptive field size

$$r_0 = \sum_{l=1}^{L} \left( (k_l - 1) \prod_{i=1}^{l-1} s_i \right) + 1$$

Sum over network layers

**Resnet-{34,50}**

| # layers | stride @ layer |
|----------|----------------|
| 1 | 1 |
| 1 | 2 |
| 3 | 4 |
| 4 | 8 |
| 6 | 16 |
| 3 | 32 |

RO = 1 + (3-1) * ( 1 * 1
+ 1 * 2
+ 3 * 4
+ 4 * 8
+ 6 * 16
+ 3 * 32)

= **479**

**Resnet-{34,50}**
*after converting last 2 stride 2 layers to stride 1*

| # layers | stride @ layer |
|----------|----------------|
| 1 | 1 |
| 1 | 2 |
| 3 | 4 |
| 4 | 8 |
| 6 | 8 |
| 3 | 8 |

RO = 1 + (3-1) * ( 1 * 1
+ 1 * 2
+ 3 * 4
+ 4 * 8
+ 6 * 8
+ 3 * 8)

= **239**

Receptive field area reduced 4x :(

https://distill.pub/2019/computing-receptive-fields/

# Replace stride 2 convolutions with stride 1



Convolution with atrous rate=1,
(i.e., ordinary convolution)

Convolution with atrous rate=2

**Problem**: Doing this directly can reduce receptive field size...

**Solution**: Use dilated/*atrous* convolution (convolution with holes, *en français*) to compensate at the second layer.

Chen, Papandreou, et al, 2015

# Stringing atrous through multiple layers

*Compensation needs to happen at all higher layers*



Use convolution with atrous rate=2 at both layers above to maintain receptive field size

Stride 1
Stride 1
Stride 2

Receptive Field Size = 11

Stride 1
Stride 1
Stride 1

Receptive Field Size = 11

# Atrous Cost/Benefit

- Quadrupled memory
- Quadrupled theoretical FLOPS
- Same # parameters

Only in affected layers, and due to larger inputs (Atrous Conv itself is not more expensive than ordinary Conv)

- High resolution outputs
- Large receptive field
- *Can initialize model from ImageNet w/o retraining*

# Case Study (2015): DeepLab-LargeFOV Architecture

Start with VGG; Remove last two pools; Use Atrous Convs in higher layers



**Original VGG-16**

| | |
|---|---|
| Softmax | |
| FC 1000 | |
| FC 4096 | |
| FC 4096 | |
| Pool | Stride 32 |
| 3x3 conv, 512 | |
| 3x3 conv, 512 | Stride 16 |
| 3x3 conv, 512 | |
| Pool | |
| 3x3 conv, 512 | |
| 3x3 conv, 512 | Stride 8 |
| 3x3 conv, 512 | |
| Pool | |
| 3x3 conv, 256 | Stride 4 |
| 3x3 conv, 256 | |
| Pool | |
| 3x3 conv, 128 | Stride 2 |
| 3x3 conv, 128 | |
| Pool | |
| 3x3 conv, 64 | Stride 1 |
| 3x3 conv, 64 | |
| Input | |

**DeepLab-LargeFOV**

| | |
|---|---|
| Conv 1x1 1024 | |
| Conv 3x3 1024 rate 12 | |
| Pool | |
| 3x3 conv, 512 | |
| 3x3 conv, 512 | Stride 8, 3x3 convs w/atrous rate 2 |
| 3x3 conv, 512 | |
| Pool | |
| 3x3 conv, 512 | |
| 3x3 conv, 512 | Stride 8 |
| 3x3 conv, 512 | |
| Pool | |
| 3x3 conv, 256 | Stride 4 |
| 3x3 conv, 256 | |
| Pool | |
| 3x3 conv, 128 | Stride 2 |
| 3x3 conv, 128 | |
| Pool | |
| 3x3 conv, 64 | Stride 1 |
| 3x3 conv, 64 | |
| Input | |

Semantic Image Segmentation with Deep Convolutional Nets and Fully Connected CRFs by Chen et al
DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs by Chen et al

# DeepLab results (Pascal VOC dataset)



**VGG based
DeepLab**

**Resnet-101 based
DeepLab**

# How to get high resolution outputs (e.g. w/stride < 32)

- Use fewer stride 2 convolutions

- **Use "upconvolution" operators**

# "Upconvolution" operators

- Resize + Conv

- Fractional / Sub-pixel Convolution

- Transpose Convolution

- Convolution + "Periodic Reshuffling"

- Unpool (not super common)



To reduce spatial resolution, use
**Convolution w/stride 2**



To *increase* spatial resolution, use **???**

# Resize + Conv



2x NN or
bilinear resize

Conv

(often merging with
lower level
features)

# Fractionally Strided / Subpixel Convolution



"Bed of nails"

**Fractional indices w/half stride**

Conv

# Convolution + "Periodic reshuffling"



Is the deconvolution layer the same as a convolutional layer? by Shi et al

# Transpose Conv



Input

Zero paddings

Output

1-D Convolution (stride 2)

We can always write (ordinary) convolution as a matrix multiplication

# Transpose Conv



1-D Convolution (stride 2)

1-D Transpose Convolution (stride 2)

Interesting fact: Swapping forwards and backwards passes of Conv op
will give Transpose Conv op

# Transpose Conv (2-d example)

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

***Crop* borders**

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Crop borders**

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

**Input**

***Crop* borders**

| 0 | 2 |
|---|---|
| 2 | 0 |
| 0 | 0 |

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

*Crop* **borders**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)



*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)



**Filter/"stamp"**

**Transpose Conv (stride 2)**

Input

Output

*Crop* borders

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)



**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

**Input**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

*Crop* borders

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Input**

**Transpose Conv (stride 2)**

***Crop* borders**

| 0 | 0 | 1 | 0 | 2 | 0 | 3 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 2 | 0 | 3 | 0 |
| 1 | 0 | 6 | 0 | 8 | 0 | 6 |
| 0 | 4 | 0 | 5 | 0 | 6 | 0 |
| 4 0 | 0 0 | 5 7 | 0 | 6 | 0 | 0 |
| 0 | 7 | 0 | | | | |
| 7 | 0 | 0 | | | | |

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

**Crop borders**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | **8** | 9 |

**Transpose Conv (stride 2)**

**Input**

| 0 | 0 | 1 | 0 | 2 | 0 | 3 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 2 | 0 | 3 | 0 |
| 1 | 0 | 6 | 0 | 8 | 0 | 6 |
| 0 | 4 | 0 | 5 | 0 | 6 | 0 |
| 4 | 0 | 12 | 0 | 0 | 8 | 0 | 0 |
| 0 | 7 | 0 | 0 | 8 | 0 | 0 |
| 7 | 0 | 8 | 0 | 0 | | |

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)



**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

*Crop* borders

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Crop borders**

**Filter/"stamp"**

| 0 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)**

| 0 | 0 | 1 | 0 | 2 | 0 | 3 |
|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 2 | 0 | 3 | 0 |
| 1 | 0 | 6 | 0 | 8 | 0 | 6 |
| 0 | 4 | 0 | 5 | 0 | 6 | 0 |
| 4 | 0 | 12 | 0 | 14 | 0 | 9 |
| 0 | 7 | 0 | 8 | 0 | 9 | 0 |
| 7 | 0 | 8 | 0 | 9 | 0 | 0 |

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Transpose Conv (2-d example)

**Filter/"stamp"**

| | | |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |

| | | |
|---|---|---|
| 1 | 2 | 3 |
| 4 | 5 | 6 |
| 7 | 8 | 9 |

**Transpose Conv (stride 2)** →

| | | | | |
|---|---|---|---|---|
| 1 | 0 | 2 | 0 | 3 |
| 0 | 6 | 0 | 8 | 0 |
| 4 | 0 | 5 | 0 | 6 |
| 0 | 12 | 0 | 14 | 0 |
| 7 | 0 | 8 | 0 | 9 |

**Input**

**Output**

*Think of "stamping" filter across the output image*

# Which one should I use??

- Fractional / Sub-pixel Convolution
- Transpose Convolution
- Convolution + "Periodic Reshuffling"

> Representationally Equivalent!

- Resize + Conv

> Slightly less expressive



Resize + Conv equivalent to Bed-of-Nails followed by an "all ones" 2x2 Conv then ordinary Conv

# Checkerboard artifacts

Transpose Convolutions "want" to generate checkerboards



Deconv in last two layers.
Other layers use resize-convolution.
*Artifacts of frequency 2 and 4.*

Deconv only in last layer.
Other layers use resize-convolution.
*Artifacts of frequency 2.*

All layers use resize-convolution.
*No artifacts.*

Resize + Conv less expressive than Transpose Conv, but less susceptible to checkerboard artifacts

https://distill.pub/2016/deconv-checkerboard/

# Case Study (2015): FCN



32x upsampled prediction (FCN-32s)

image    pool1    pool2    pool3    pool4    pool5

**VGG-based FCN (stride 32)**

Fully Convolutional Networks for Semantic Segmentation by Long et al.

# Case Study (2015): FCN



Transpose Conv

32x upsampled prediction (FCN-32s)

2x upsampled prediction

16x upsampled prediction (FCN-16s)

image   pool1   pool2   pool3   pool4   pool5

pool4 prediction

1x1 Conv

*"Make local predictions that respect global structure"*

**VGG-based FCN (stride 16)**

Fully Convolutional Networks for Semantic Segmentation by Long et al.

# Case Study (2015): FCN



VGG-based FCN (stride 8)

Fully Convolutional Networks for Semantic Segmentation by Long et al.

# Case Study (2015): FCN

# Case Study (2019): FPN (revisited)

# Case Study(2018) DeepLabV3+



Atrous Spatial Pyramid

(Bilinear) Resize + Conv

Spatial Pyramid Pooling

0.5x

0.5x

0.5x

0.5x

4x

4x

Image

Prediction

# Outline of Semantic Segmentation

- The sliding window connection (again)

- Fully Convolutional models

- How to get high resolution outputs with

  - Atrous convolutions

  - "Upconvolutions"

- **Target Assignment**

- Evaluation of Semantic Segmentation

Relevant for all dense prediction tasks

# Target Assignment / Alignment



**Subsample groundtruth**

**Loss**

or...

**Loss**

GT

**Upsample predictions**

Image

# Target Assignment / Alignment

**A reasonable desideratum**: *groundtruth target for a particular logit should be sampled at center of that prediction's receptive field*

- *Getting this right requires thinking about padding, specific resizing algorithm*



**"Valid padded" network**

Assign "Water"

**"Same padded" network**

Assign "Beach"

# Recap

We want

- High output resolution

- Large receptive fields

- "Alignment" between receptive fields and targets

# Outline of Semantic Segmentation

- The sliding window connection (again)

- Fully Convolutional models

- How to get high resolution outputs with

  - Atrous convolutions

  - "Upconvolutions"

- Target Assignment

- **Evaluation of Semantic Segmentation**

Relevant for all dense prediction tasks

# How to evaluate a segmentation model: Per-Pixel Accuracy

*Problem with per-pixel accuracy --- not fair to small/thin classes*



Categories: **Water**, **Land**

# How to evaluate a segmentation model: Per-Pixel Accuracy

*Problem with per-pixel accuracy --- not fair to small/thin classes*



*Setting every pixel to "**Land**" is >90% Accuracy*

Categories: **Water**, **Land**

# How to evaluate a segmentation model: "Mask IOU"



$$IOU = \frac{Intersection}{Union}$$

- Masks are disjoint if and only if IOU=0
- Masks are identical if and only if IOU=1

# How to evaluate a semantic segmentation model



| Sky | Building | Pole | Road | Sidewalk | Vegetation | Sign | Fence | Car | Pedestrian | Cyclist | Void |



**Image**



**Groundtruth**



**Prediction**

**Mean IOU** = Mean(IOU(groundtruth_c, predicted_c)
for c in {Sky, Building, Pole, ...})

# Lecture Outline

**Dense Prediction (pixel level prediction)**

- Semantic Segmentation

- **Instance Segmentation**

- **Panoptic Segmentation**

- **Keypoint Estimation**

# Semantic vs Instance Segmentation: Don't get confused!



classify

classify and regress
bounding box per object

**(bounding box)
detection**

classify per pixel

**semantic
segmentation**

classify per pixel per object

**instance
segmentation**

Figure from Lin et al 2014

# Mask R-CNN



**Boxes first paradigm**:

1. Run detector (Faster R-CNN)
2. Produce segmentation relative to each predicted box

*Mask R-CNN combines both steps into an end-to-end trainable model*

Mask R-CNN by He et al, 2017

# Mask R-CNN Training



**Faster R-CNN Stage 1**

7x7 conv, 64, /2 · pool, /2 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 64 · 3x3 conv, 128, /2 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 128 · 3x3 conv, 256, /2 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256 · 3x3 conv, 256

**Block 1**   **Block 2**   **Block 3**

RPN Classification Loss

RPN Localization Loss

**Faster R-CNN Stage 2**

Crop/Resize (a.k.a. ROIAlign) and Stack

3x3 conv, 512, /2 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512 · 3x3 conv, 512

**7x7x512**

**Block 4**

Object Classification Loss

Object Localization Loss

Note: exact dimensions in this figure are a bit off as this figure is based on "basic residual unit"

**Transpose Conv**

**"Mask Head"**

14x14x512 · 14x14xK

**Mask Loss**
Per-pixel Sigmoid XEnt

# Mask R-CNN Inference



**Faster R-CNN Stage 1**

Block 1  Block 2  Block 3

**Faster R-CNN Stage 2**

Crop/Resize (a.k.a. ROIAlign) and Stack

7x7x512

Block 4

**Mask R-CNN Stage 3**

Crop/Resize (a.k.a. ROIAlign) and Stack

Block 4

Transpose Conv

14x14x512  14x14xK

*Apply Mask head to top-100 scoring boxes (for speed)*

# Evaluation for Instance Segmentation

- We care about the same things as object detection
  - E.g. Precision, Recall, Average Precision (AP), mean Average Precision (mAP)

$$\text{Box IOU} = \frac{\text{Intersection}}{\text{Union}}$$

$$\text{Mask IOU} = \frac{\text{Intersection}}{\text{Union}}$$

*But... with Mask IOU instead of Box IOU*

# Stuff vs Things



*Semantic segmentation makes more sense*

**Stuff**

**Things**

*Instance segmentation makes more sense*

Figure from Caesar et al 2018

# Handle both stuff *and* things: Panoptic Segmentation



(a) image

(b) semantic segmentation

(c) instance segmentation

(d) panoptic segmentation

- Assign (category, instance id) pair to each pixel in image.

- Instance label ignored for "stuff" categories.

Kirillov et al, 2018

# Measuring Panoptic Quality



Ground Truth      Prediction

|  | mAP | mIOU |
|---|---|---|
| **Things** | Standard for thing categories (instance segmentation) | Does not account for False Positives/Negatives |
| **Stuff** | Stuff segments typically do not come with a score needed to compute mAP | Standard for stuff categories (semantic segmentation) |

Figure by Kirillov et al 2018

# Measuring Panoptic Quality



Ground Truth

Prediction

$$\text{TP}_l = \{(\ \boxed{\phantom{x}}\ ,\ \boxed{\phantom{x}}\ ),\ (\ \boxed{\phantom{x}}\ ,\ \boxed{\phantom{x}}\ )\}$$

$$\text{FP}_l = \{\ \boxed{\phantom{x}}\ \}$$

$$\text{FN}_l = \{\ \boxed{\phantom{x}}\ \}$$

**Match Groundtruth and Predicted segments _if IOU>50%_**

Figure by Kirillov et al 2018

# Measuring Panoptic Quality



Ground Truth                     Prediction

$$\mathrm{TP}_l = \{(\ \blacksquare\ ,\ \blacksquare\ ),\ (\ \blacksquare\ ,\ \blacksquare\ )\}$$

$$\mathrm{FP}_l = \{\ \blacksquare\ \}$$

$$\mathrm{FN}_l = \{\ \blacksquare\ \}$$

$$\mathrm{PSQ}_l = \frac{\mathrm{IoU}(\ \blacksquare\ ,\ \blacksquare\ ) + \mathrm{IoU}(\ \blacksquare\ ,\ \blacksquare\ )}{|\mathrm{TP}_l| + |\mathrm{FP}_l| + |\mathrm{FN}_l|}$$

Figure by Kirillov et al 2018

# Measuring Panoptic Quality



Ground Truth

Prediction

**More generally:**
$$\text{PQ} = \frac{\sum_{(p,g) \in TP} \text{IoU}(p,g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$$

# Measuring Panoptic Quality



Ground Truth

Prediction

**More generally:** $\text{PQ} = \dfrac{\sum_{(p,g) \in TP} \text{IoU}(p,g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$    ( = $F_1$-score * mIOU )

*Another common detection metric*

Figure by Kirillov et al 2018

# Keypoint Detection



Slide courtesy of George Papandreou, Tyler Zhu

# "Top-down" approach: Mask R-CNN



Predict heatmap for each pose keypoint

He et al, 2017

# "Bottom up" approach: Predict keypoint positions (Step 1)



Image credit: DeeperCut paper

# "Bottom up" approach: Group keypoints (Step 2)



vs.

# Example "bottom up" method: PersonLab



Papandreou, Zhu et al, 2018

# "Bottom up" vs "Top down"

## Performance on COCO keypoints task

| | $AP$ | $AP^{.50}$ | $AP^{.75}$ | $AP^M$ | $AP^L$ | $AR$ | $AR^{.50}$ | $AR^{.75}$ | $AR^M$ | $AR^L$ |
|---|---|---|---|---|---|---|---|---|---|---|
| Bottom-up methods: | | | | | | | | | | |
| CMU-Pose [32] (+refine) | 0.618 | 0.849 | 0.675 | 0.571 | 0.682 | 0.665 | 0.872 | 0.718 | 0.606 | 0.746 |
| Assoc. Embed. [2] (multi-scale) | 0.630 | 0.857 | 0.689 | 0.580 | 0.704 | - | - | - | - | - |
| Assoc. Embed. [2] (mscale, refine) | 0.655 | 0.879 | 0.777 | 0.690 | 0.752 | 0.758 | 0.912 | 0.819 | 0.714 | 0.820 |
| Top-down methods: | | | | | | | | | | |
| Mask-RCNN [34] | 0.631 | 0.873 | 0.687 | 0.578 | 0.714 | 0.697 | 0.916 | 0.749 | 0.637 | 0.778 |
| G-RMI *COCO-only* [33] | 0.649 | 0.855 | 0.713 | 0.623 | 0.700 | 0.697 | 0.887 | 0.755 | 0.644 | 0.771 |
| PersonLab (ours): | | | | | | | | | | |
| ResNet101 (single-scale) | 0.655 | 0.871 | 0.714 | 0.613 | 0.715 | 0.701 | 0.897 | 0.757 | 0.650 | 0.771 |
| ResNet152 (single-scale) | **0.665** | 0.880 | 0.726 | 0.624 | 0.723 | 0.710 | 0.903 | 0.766 | 0.661 | 0.777 |
| ResNet101 (multi-scale) | 0.678 | 0.886 | 0.744 | 0.630 | 0.748 | 0.745 | 0.922 | 0.804 | 0.686 | 0.825 |
| ResNet152 (multi-scale) | **0.687** | 0.890 | 0.754 | 0.641 | 0.755 | 0.754 | 0.927 | 0.812 | 0.697 | 0.830 |

Papandreou, Zhu et al, 2018

# Another example "bottom up" method: "Objects as Points"



Predict heatmap for each pose keypoint

Predict heatmap for object center

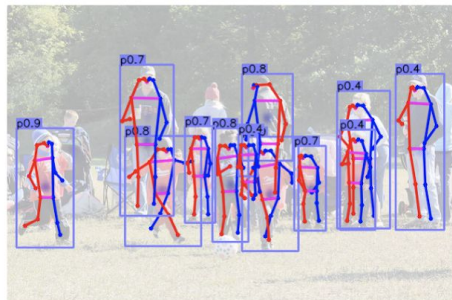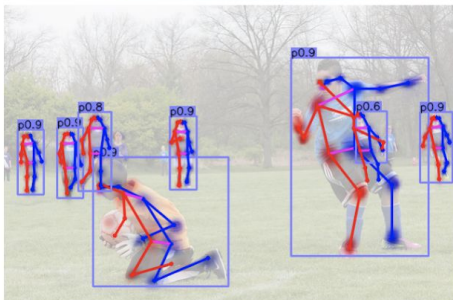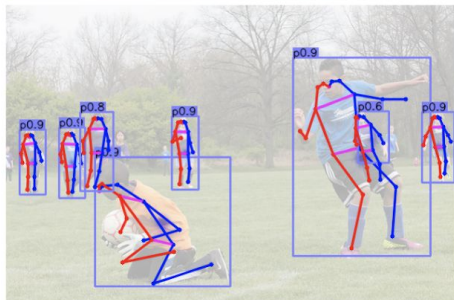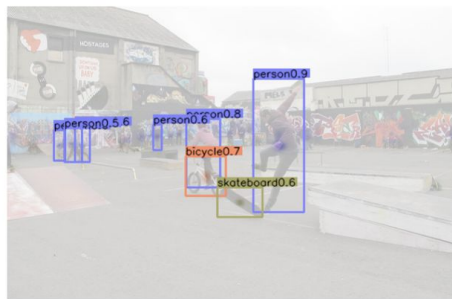Predict offset to each pose keypoint

Predict object height/width

Zhou et al, 2019

Zhou et al, 2019

# New kid on the block: "Anchor-free" object detection



Zhou et al, 2019

# Wrap up

**Dense Prediction (pixel level prediction)**

- Semantic Segmentation

- Instance Segmentation

- Panoptic Segmentation

- Keypoint Estimation



**Technology!**

room

person

Wants to weigh himself

Standing on

scale

person

Stepping on

Watching and laughing

person

Wants to play a prank

**Stepping on a scale adds weight and ups the reading.**

© Pete Souza / The White House

ImageNet: Where have we been? Where are we going? by Fei Fei Li and Jia Deng