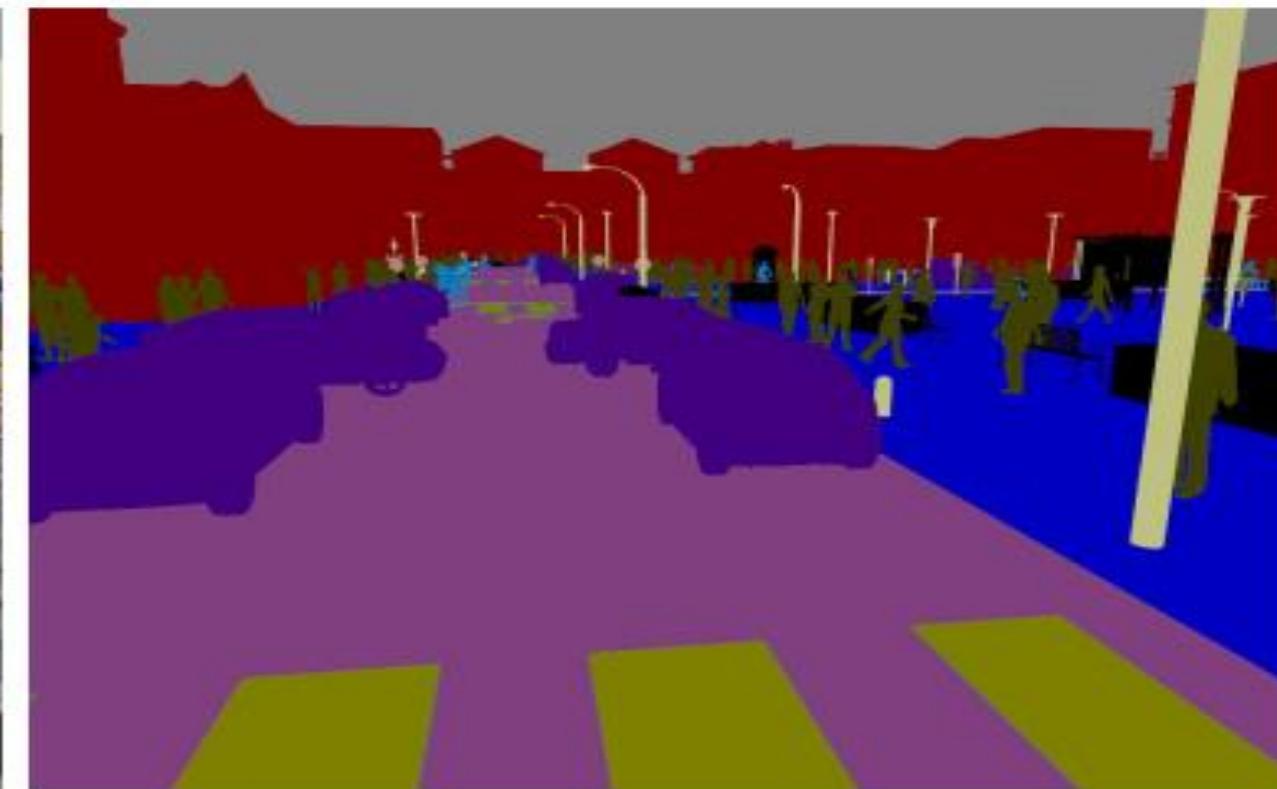


Image Segmentation using Deep Learning



Previous Class: Object bounding box detections

Deep learning based detectors

Two-stage detectors

R-CNN

Fast R-CNN

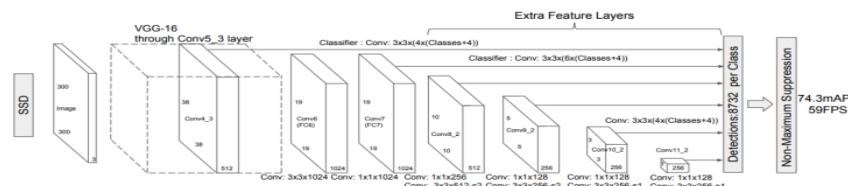
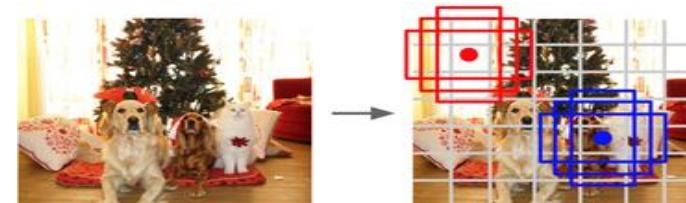
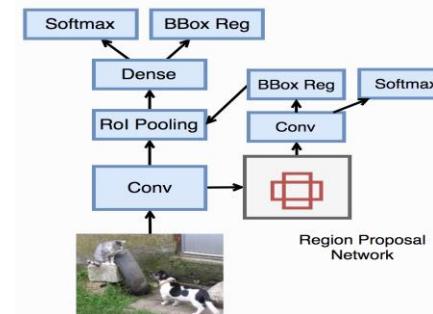
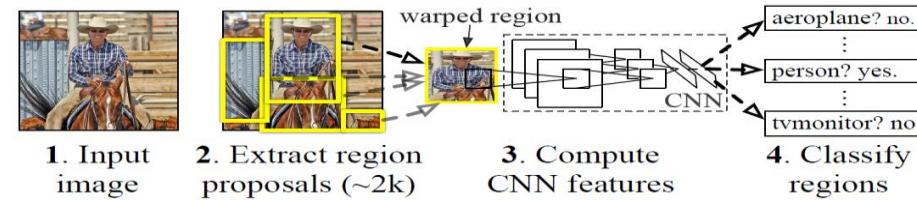
Faster R-CNN

One-stage detectors

YOLO

SSD

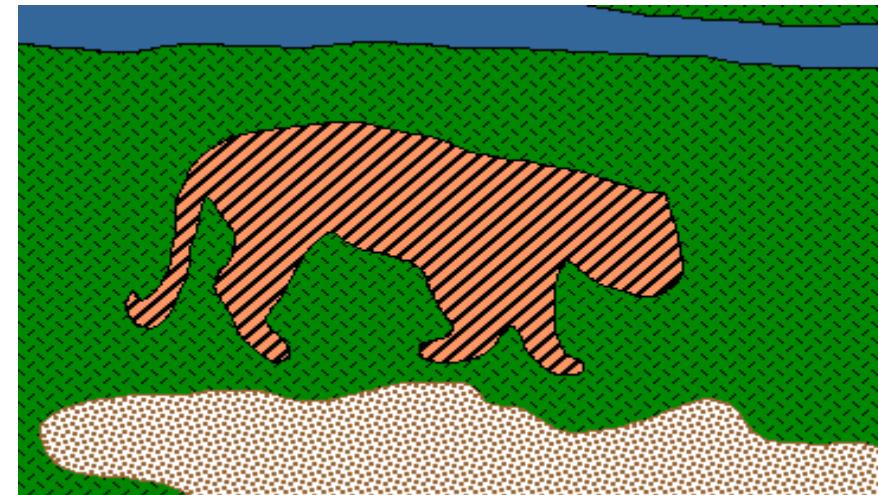
RetinaNet



Today's class

Image Segmentation

Group pixels into meaningful or perceptually similar regions



Outline

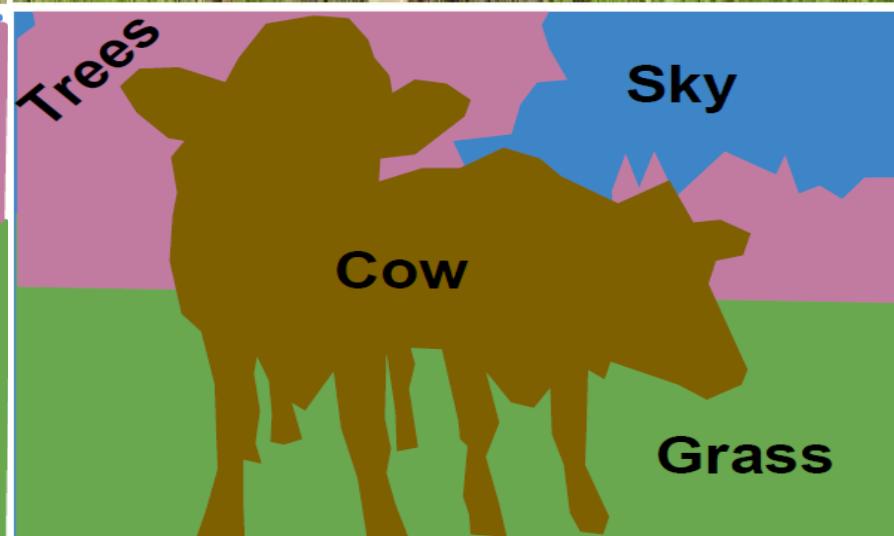
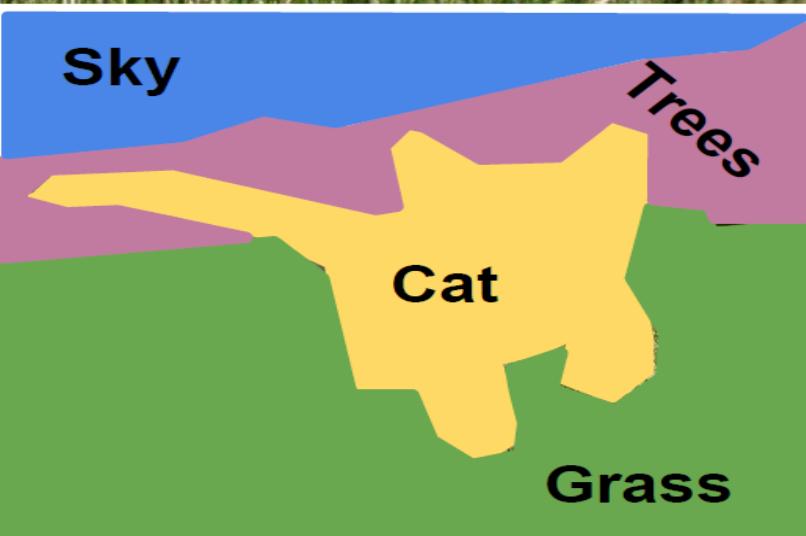
- Early “hacks”
 - Sliding window
 - Fully convolutional networks
- Deep network operations for dense prediction
 - Transposed convolutions
 - Unpooling
 - Dilated convolutions
- Instance segmentation
 - Mask R-CNN
- Other dense prediction problems

Semantic Segmentation

Label each pixel in the image with a category label

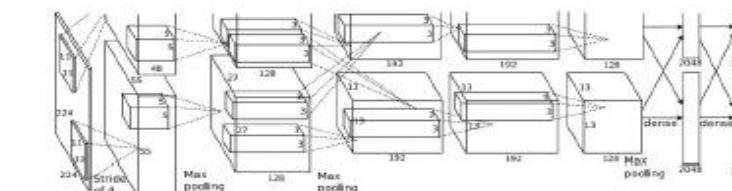
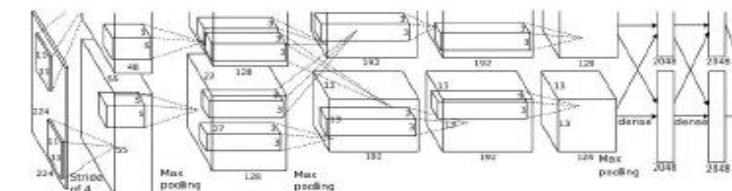
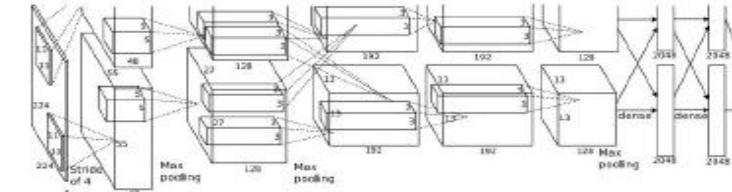
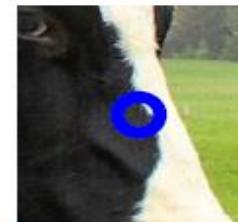
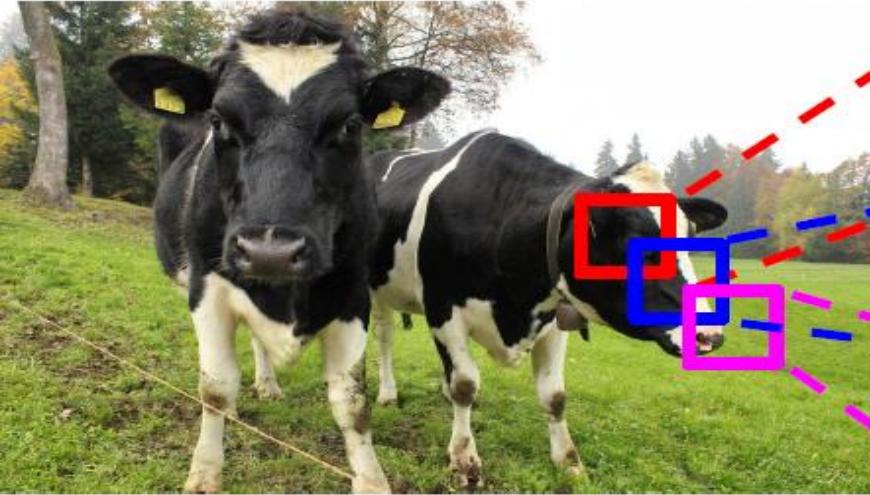
Don't differentiate instances, only care about pixels

These images are CC0 public domain [source1](#) [source2](#)



Semantic Segmentation: Sliding Window

Full image

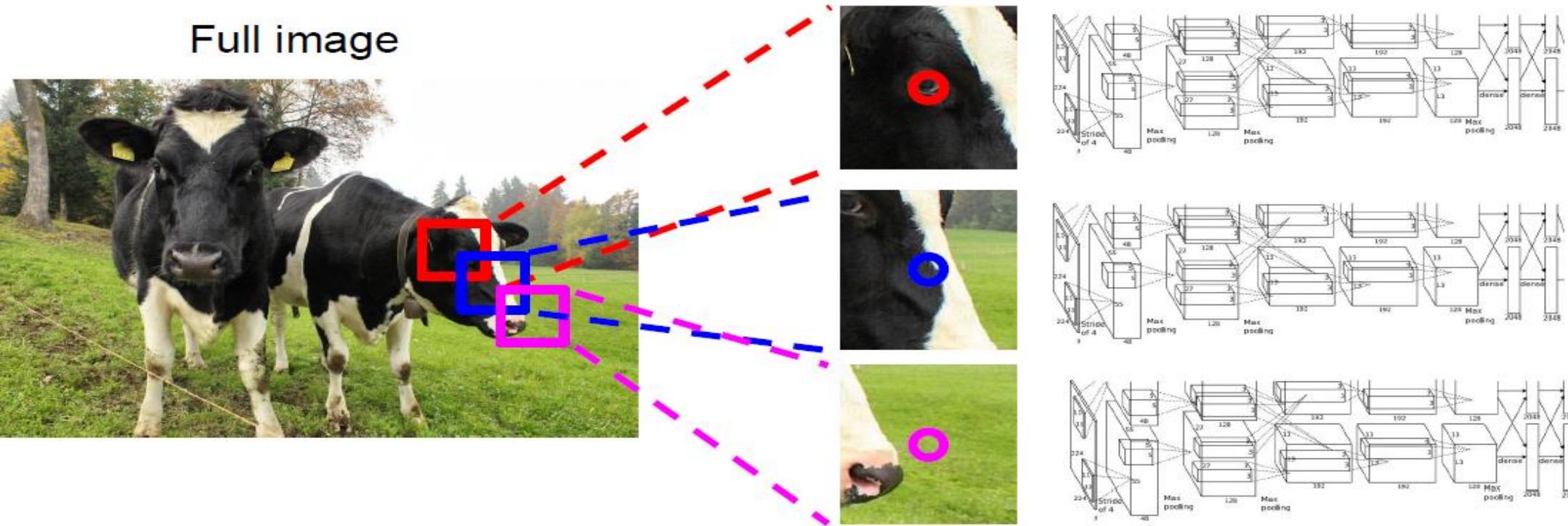


Cow

Cow

Grass

Semantic Segmentation: Sliding Window



Problem:

Very inefficient!

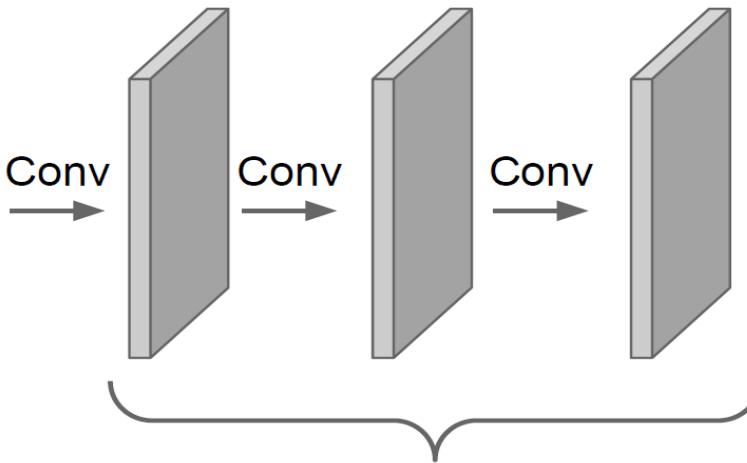
Not reusing shared features between overlapping patches

Semantic Segmentation: Fully Convolutional

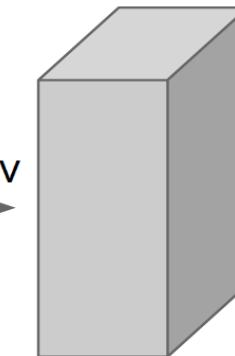
Design a network as a bunch of convolutional layers
to make predictions for pixels all at once!



Input:
 $3 \times H \times W$

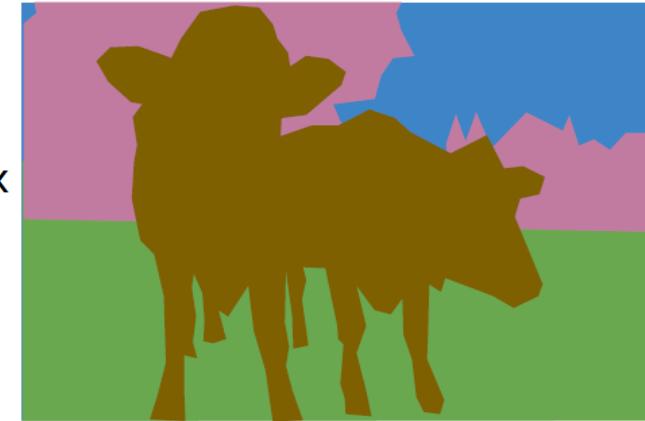


Convolutions:
 $D \times H \times W$



Conv

argmax



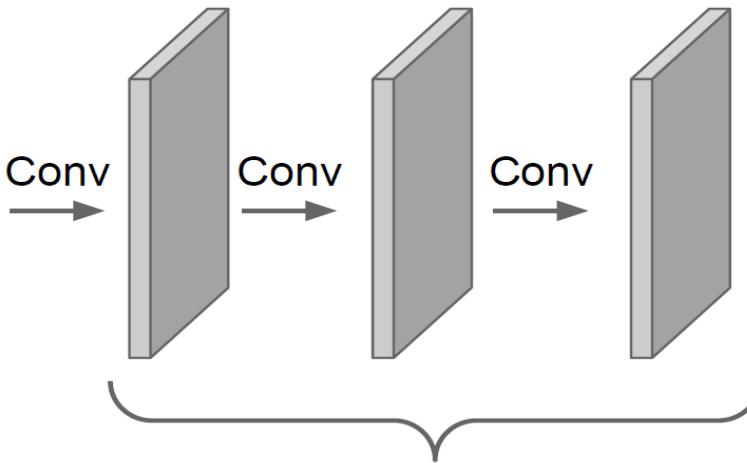
Predictions:
 $H \times W$

Semantic Segmentation: Fully Convolutional

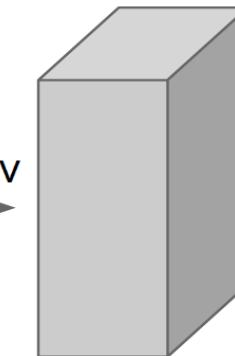
Design a network as a bunch of convolutional layers
to make predictions for pixels all at once!



Input:
 $3 \times H \times W$

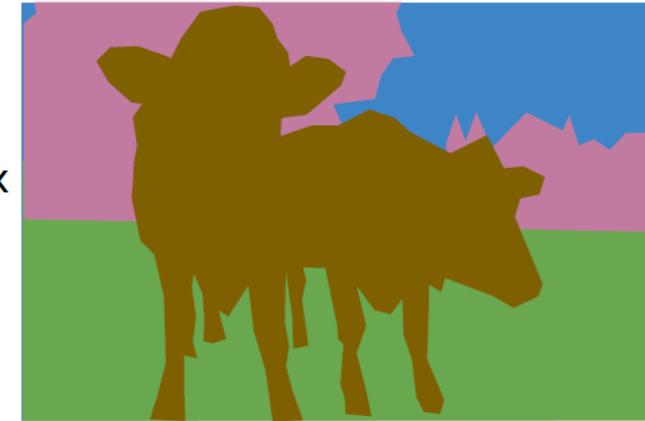


Convolutions:
 $D \times H \times W$



Conv

argmax



Scores:
 $C \times H \times W$

Predictions:
 $H \times W$

Problem:

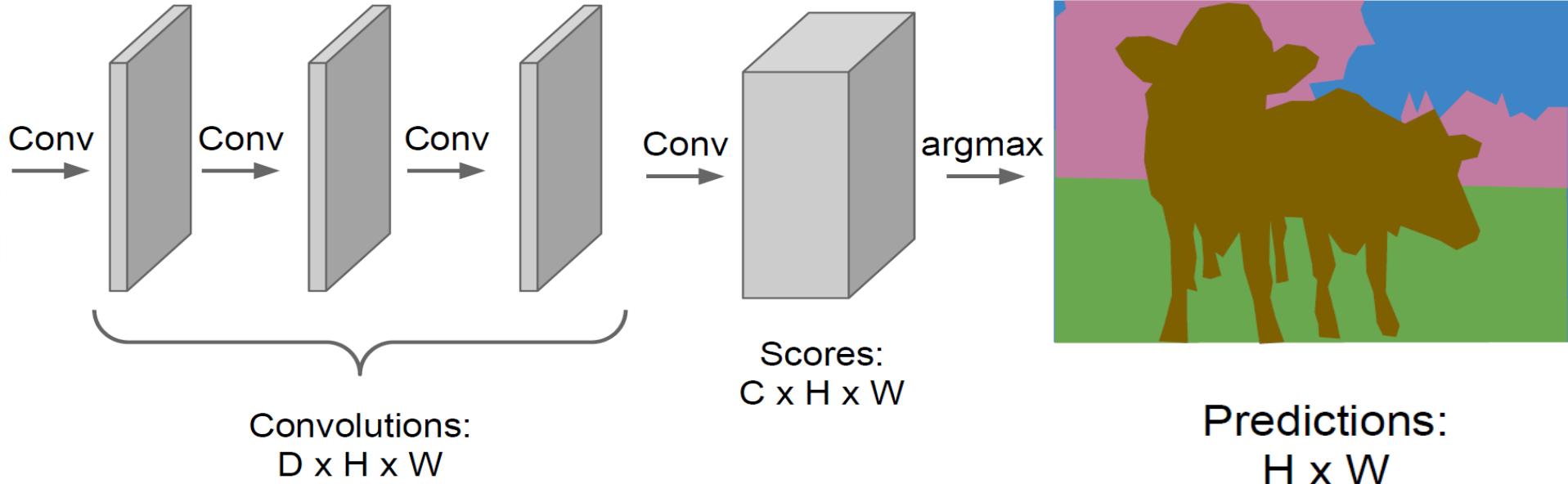
Convolutions at original image resolution will be very
expensive ...

Semantic Segmentation: Fully Convolutional

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Input:
 $3 \times H \times W$



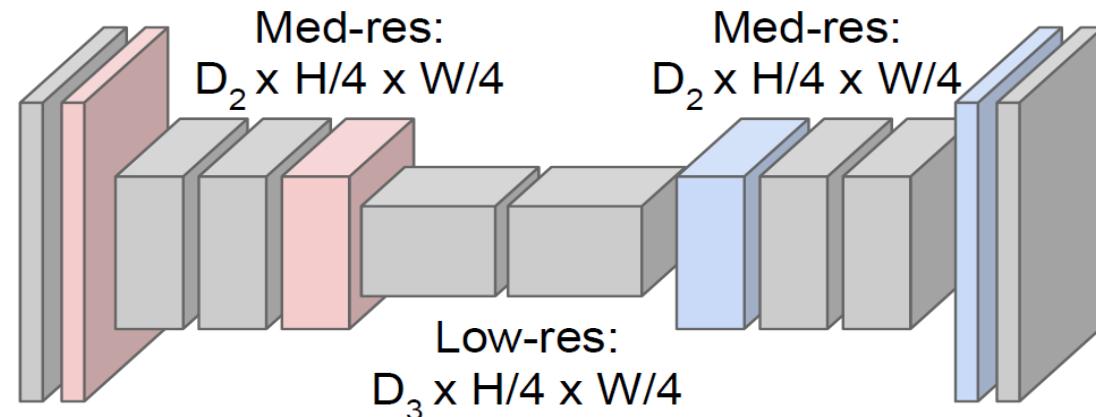
Semantic Segmentation: Fully Convolutional

Downsampling:
Pooling, strided
convolution

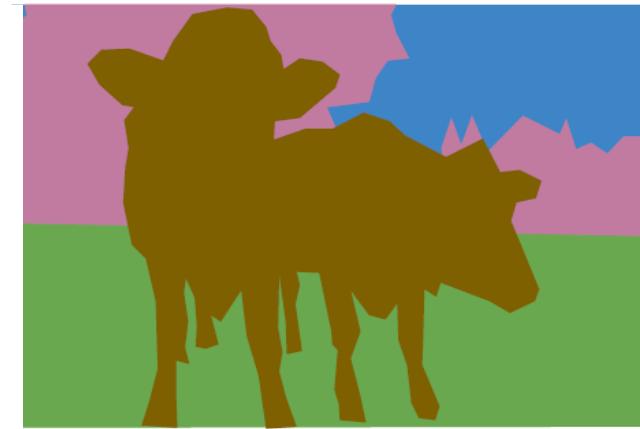


Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Upsampling:
??



Predictions:
 $H \times W$

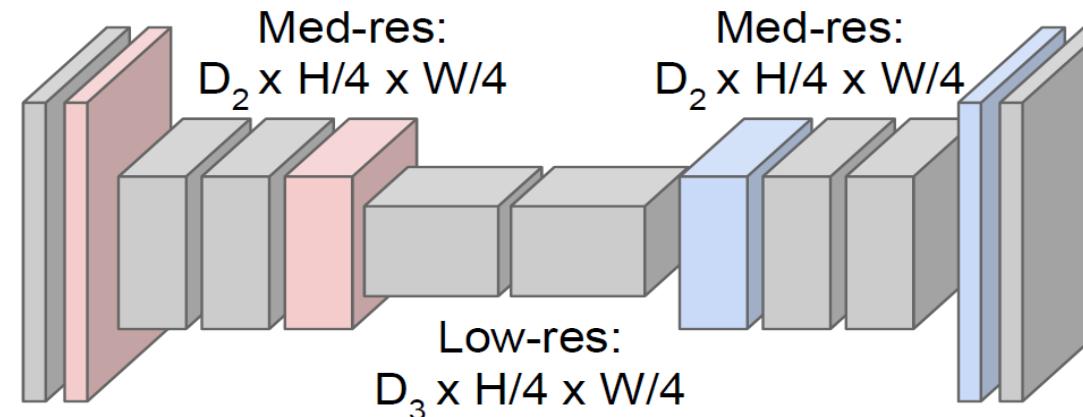
Semantic Segmentation: Fully Convolutional

Downsampling:
Pooling, strided
convolution

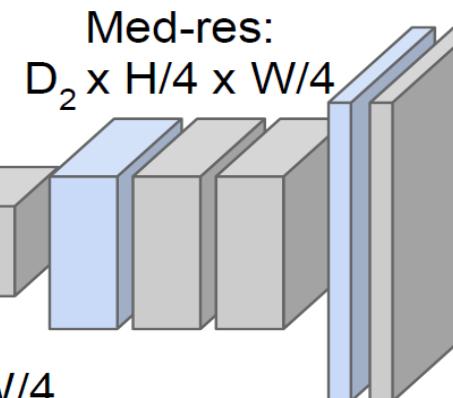


Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



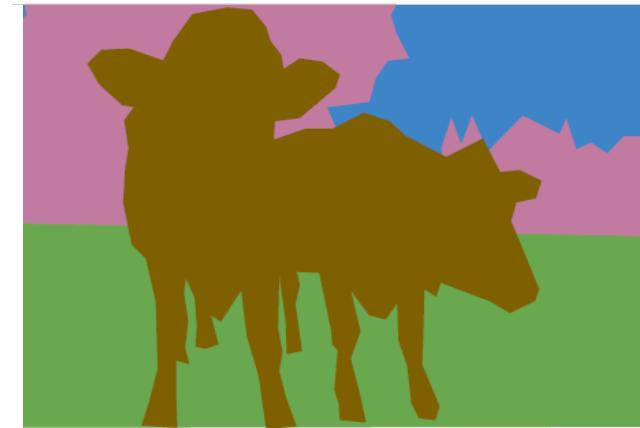
High-res:
 $D_1 \times H/2 \times W/2$



Low-res:
 $D_3 \times H/4 \times W/4$

High-res:
 $D_1 \times H/2 \times W/2$

Upsampling:
Unpooling or strided
transpose convolution



Predictions:
 $H \times W$

Upsampling: Unpooling

Nearest Neighbor

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



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

Input: 2×2

Output: 4×4

Upsampling: Unpooling

Nearest Neighbor

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



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

Output: 4 x 4

Input: 2 x 2

“Bed of Nails”

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



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

Output: 4 x 4

Input: 2 x 2

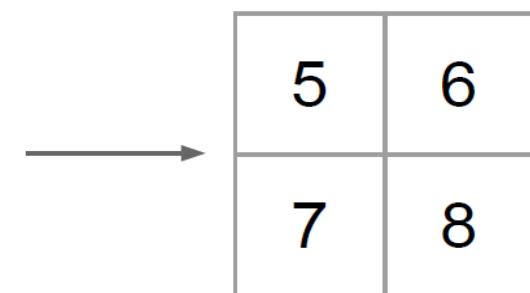
Upsampling: Max Unpooling

Max Pooling

Remember which element was max!

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

Input: 4×4



Output: 2×2

Upsampling: Max Unpooling

Max Pooling

Remember which element was max!

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

Input: 4×4

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

Output: 2×2

Max Unpooling

Use positions from pooling layer

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

Input: 2×2

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

Output: 4×4

Upsampling: Max Unpooling

Max Pooling

Remember which element was max!

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

Input: 4×4

Corresponding pairs of
downsampling and
upsampling layers

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

Output: 2×2

Max Unpooling

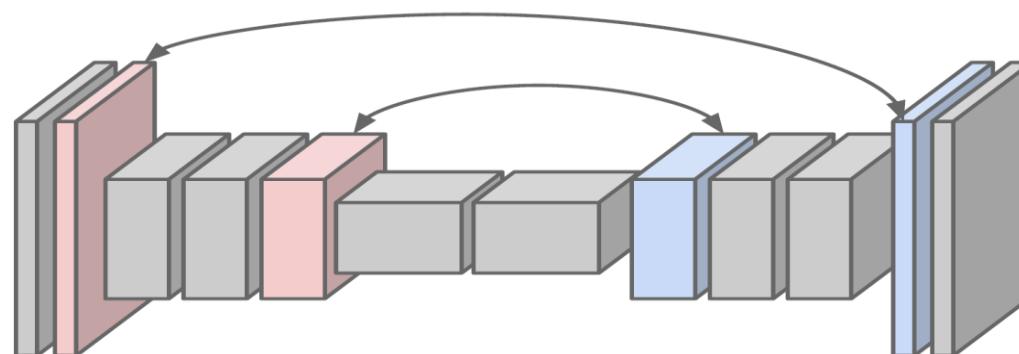
Use positions from pooling layer

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

Input: 2×2

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

Output: 4×4



Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

| | |
|--|--|
| | |
| | |

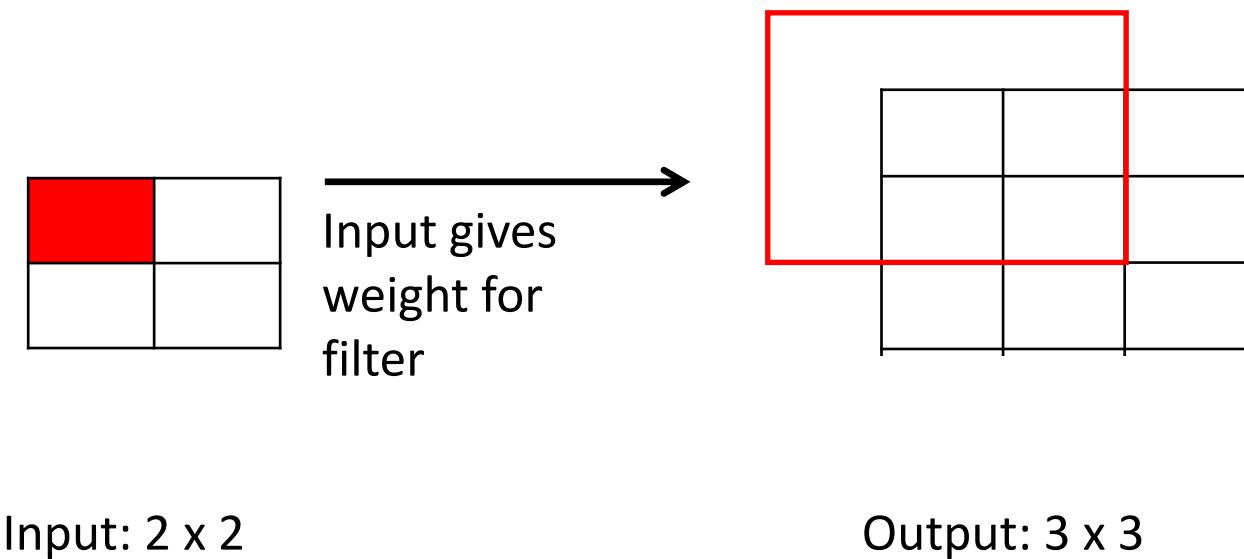
Input: 2 x 2

| | | |
|--|--|--|
| | | |
| | | |
| | | |

Output: 3 x 3

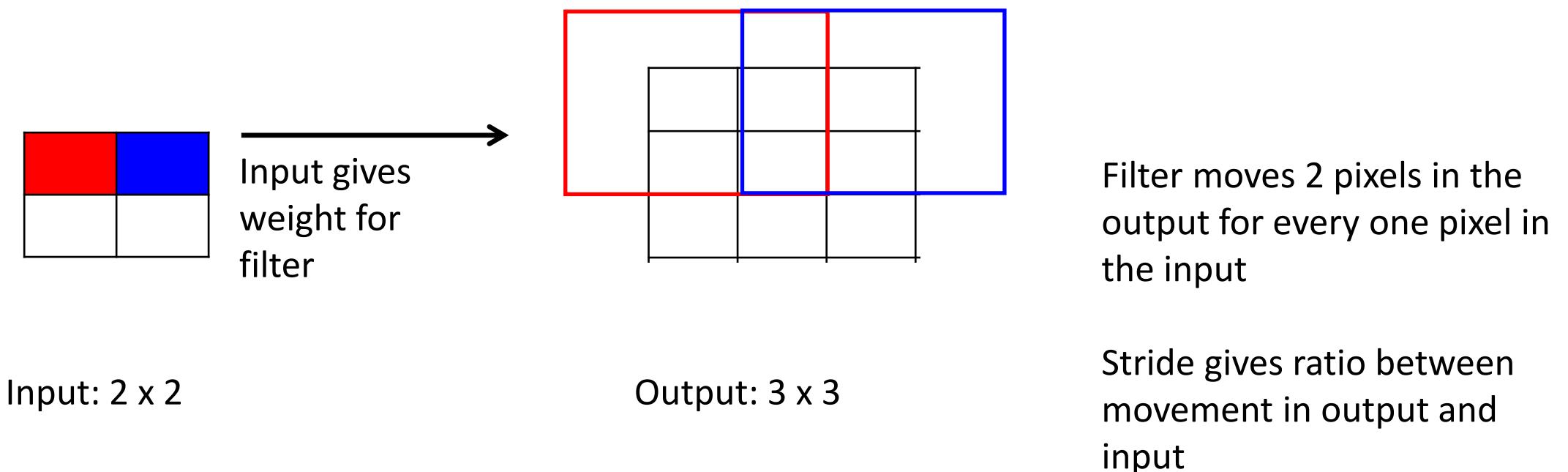
Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



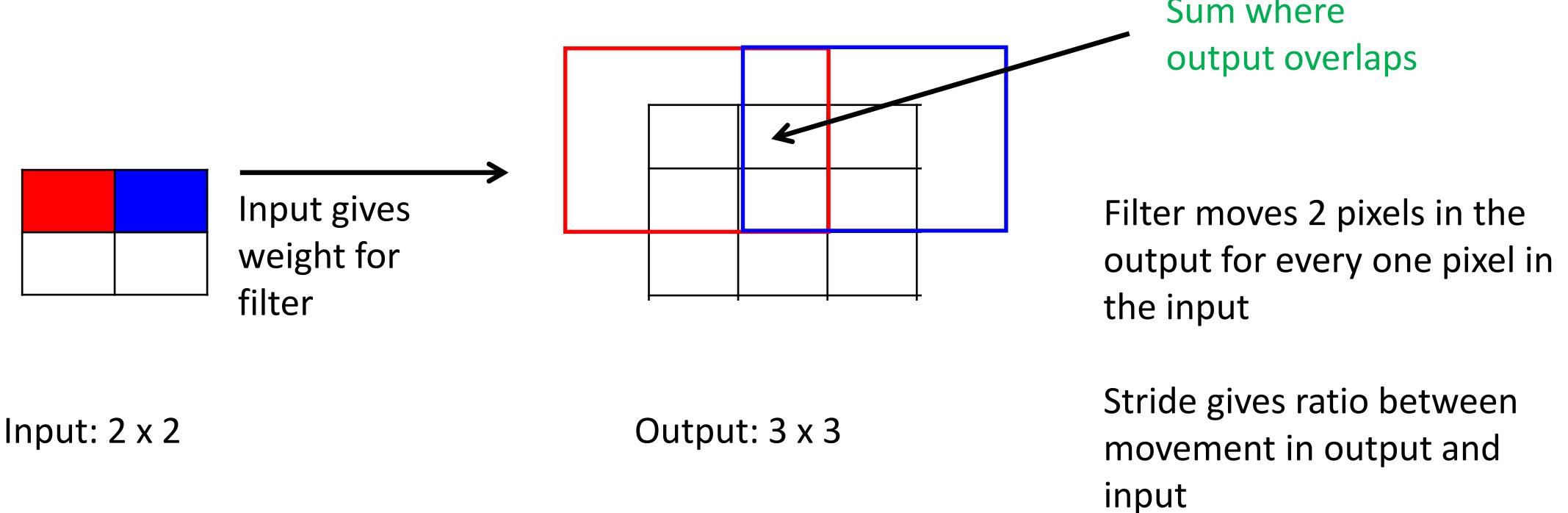
Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



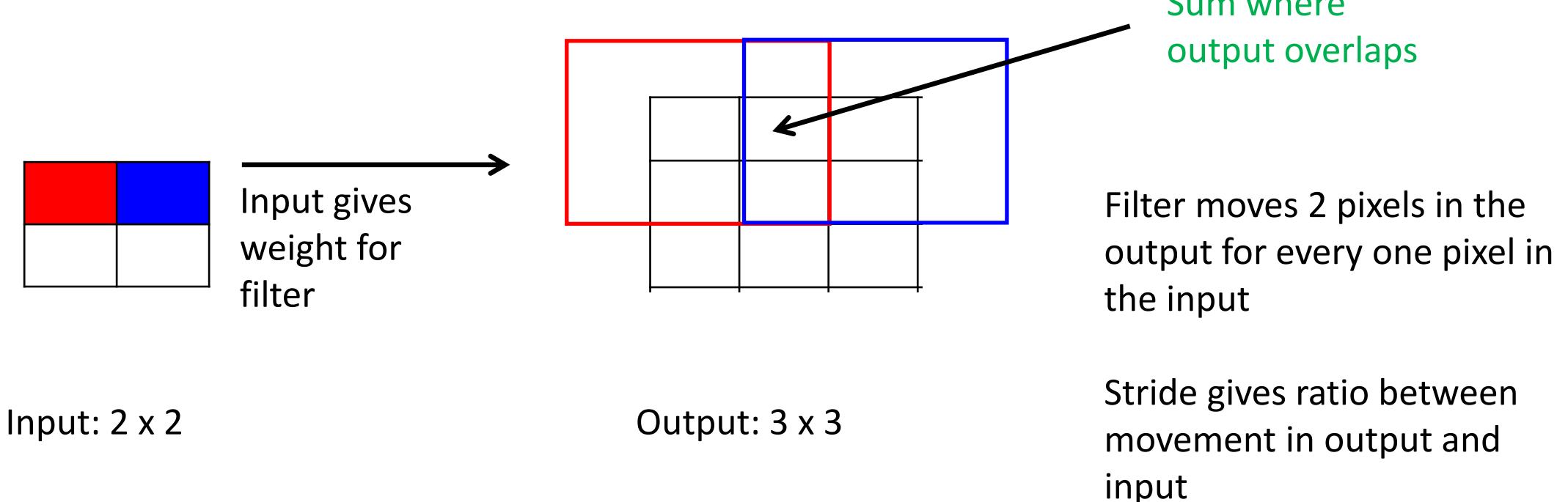
Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

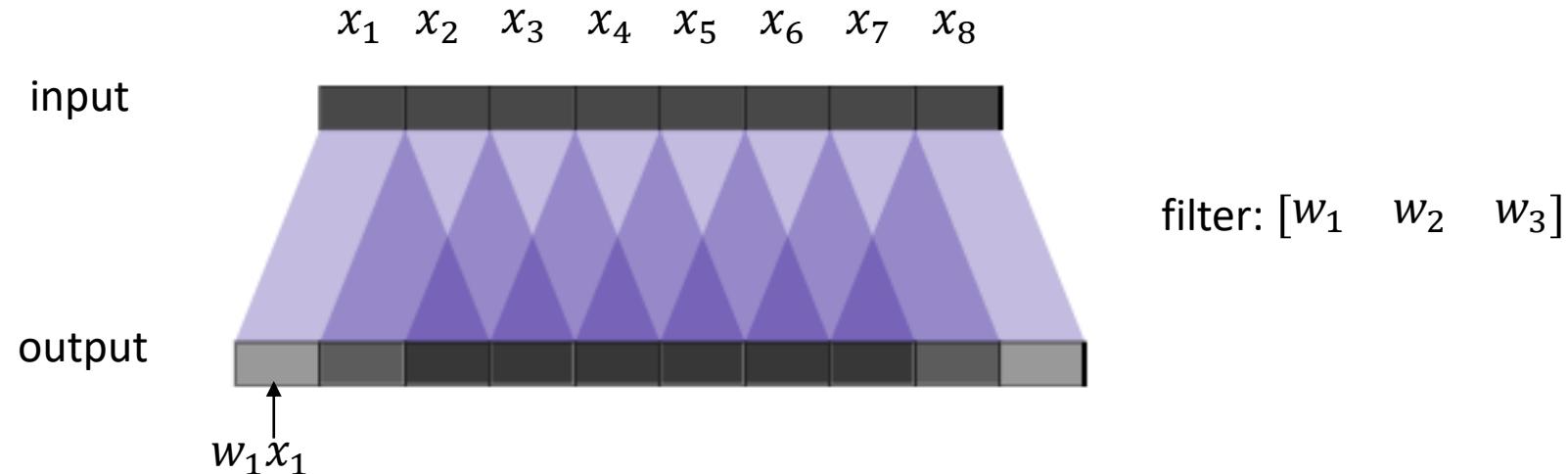


Actual output: 5*5

Need to crop 3*3 centrally

Upsampling in a deep network

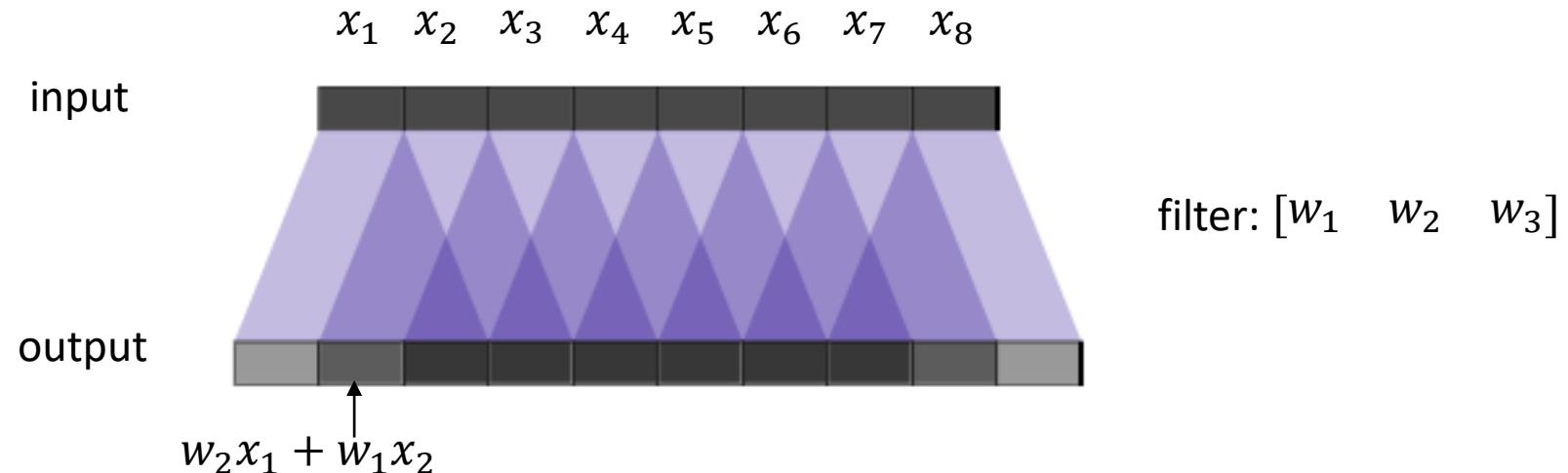
- 1D example



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

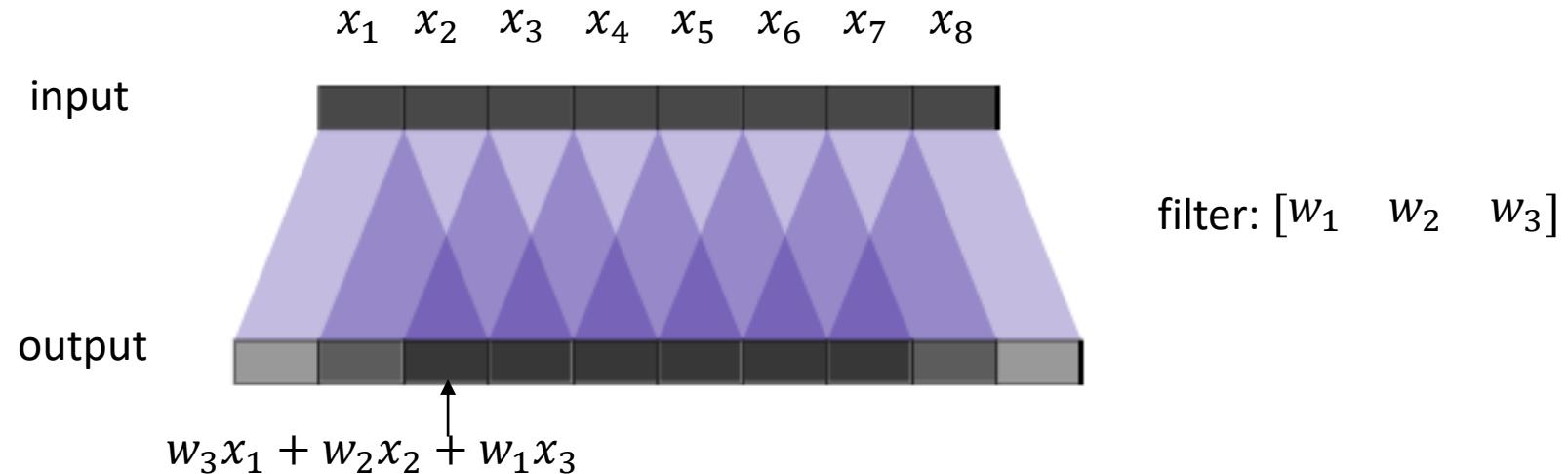
- 1D example



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

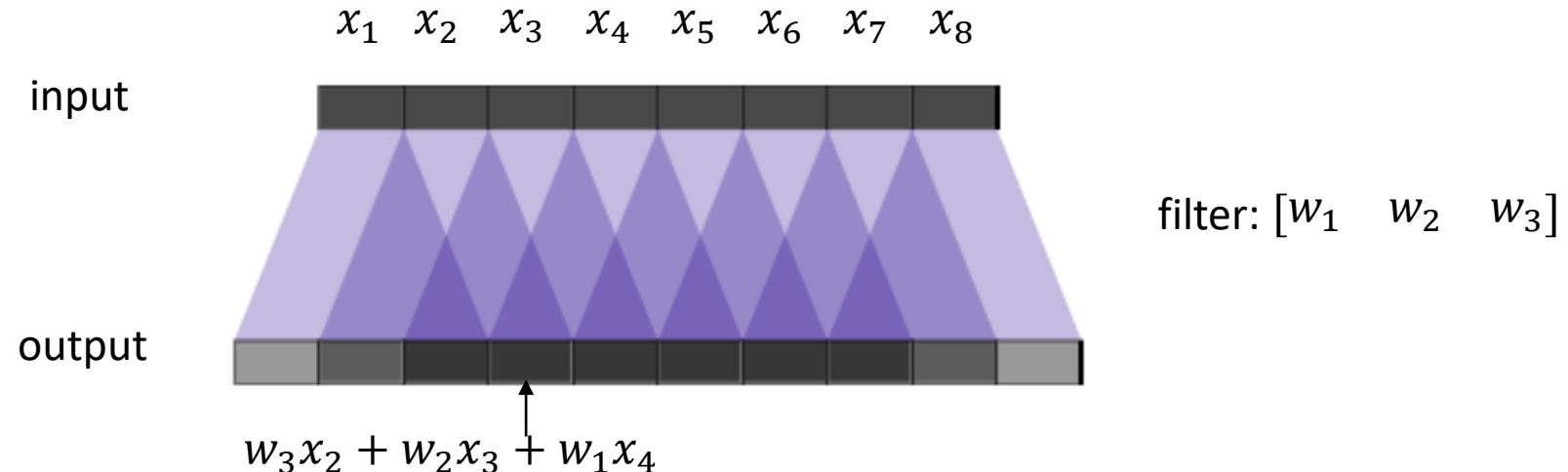
- 1D example



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

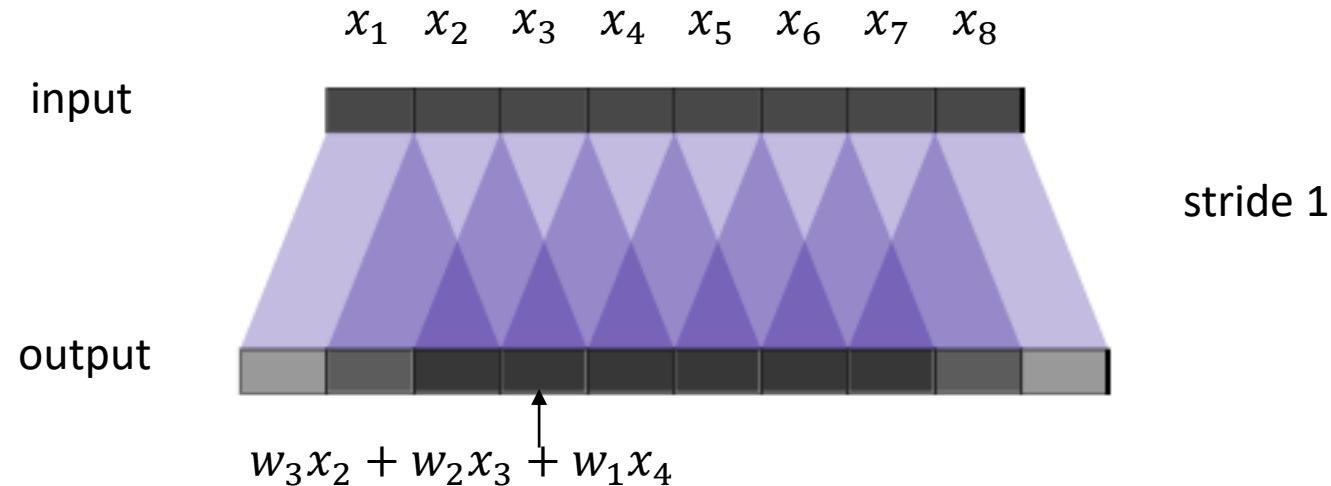
- 1D example



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

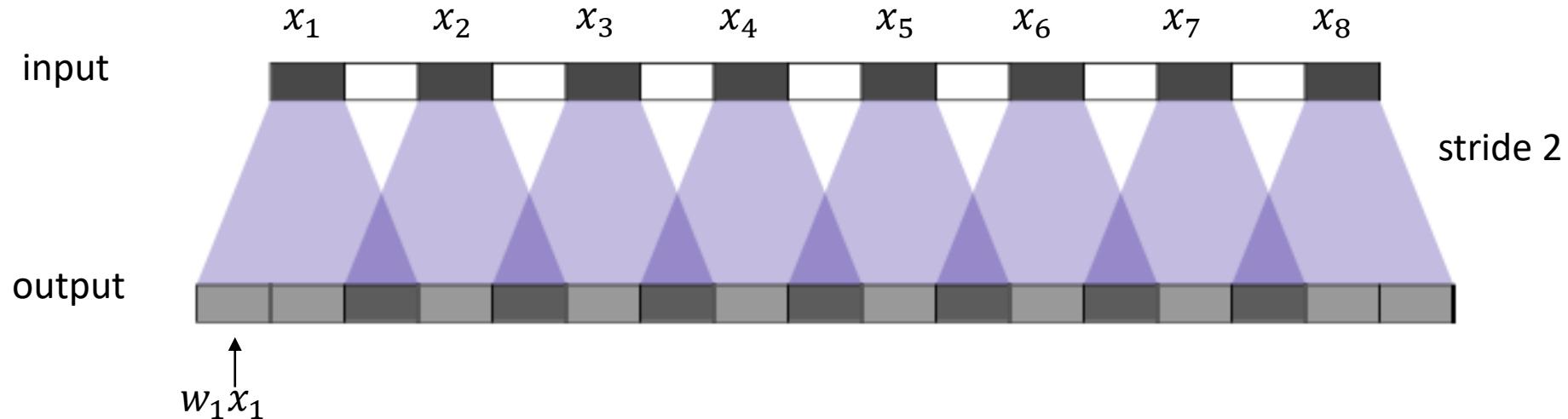
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

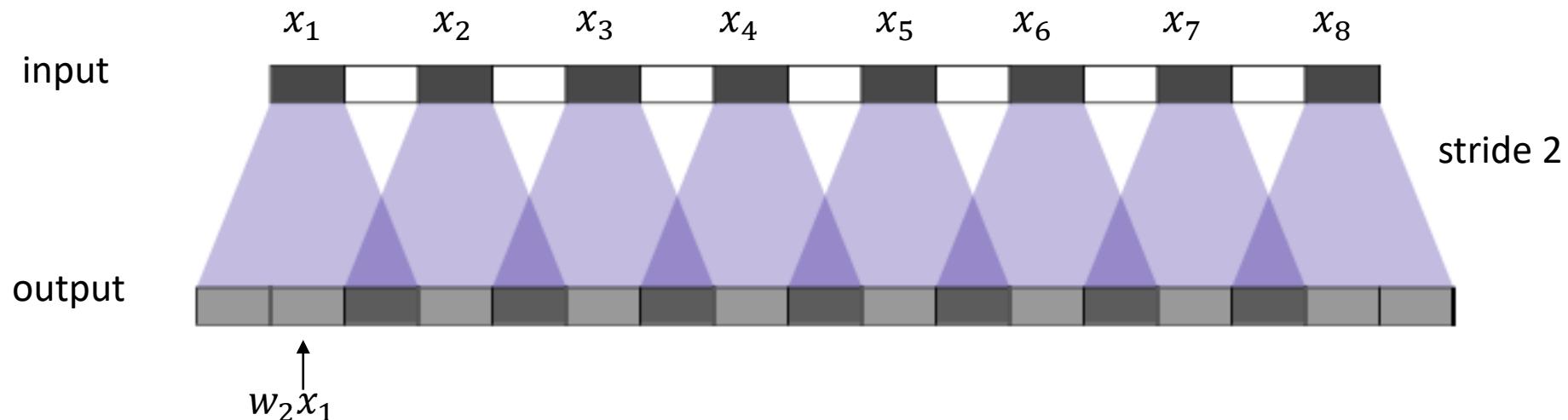
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

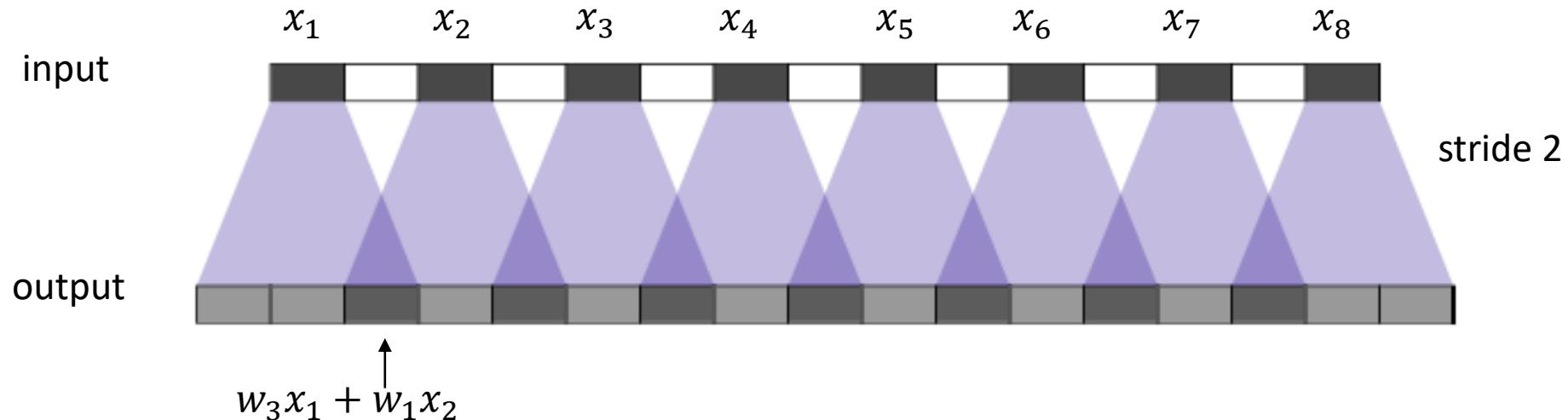
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

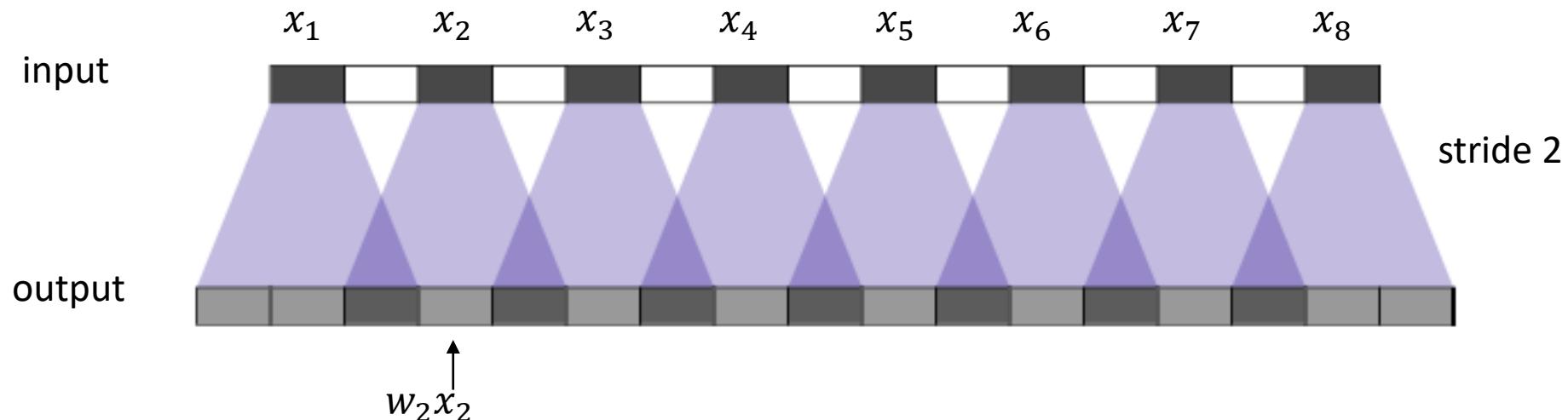
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

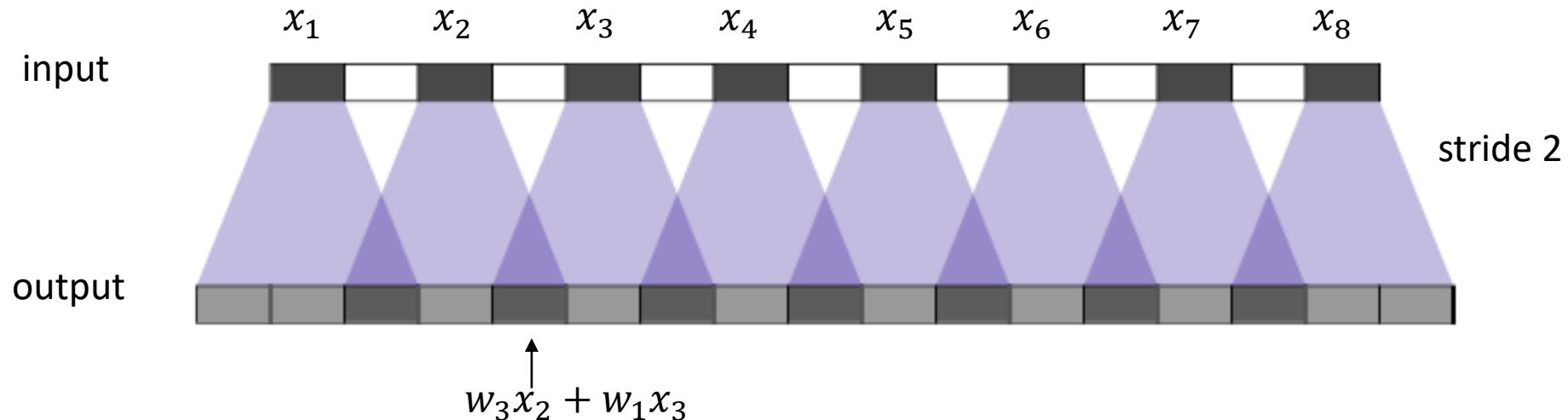
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

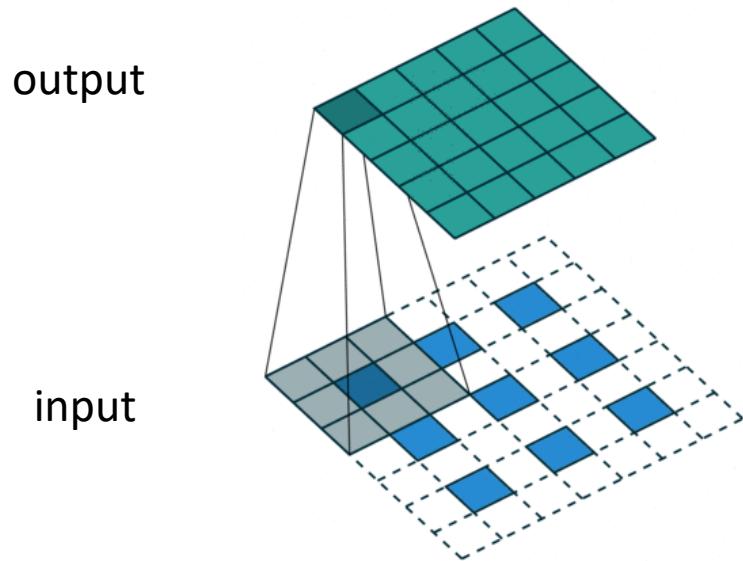
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



Animation: <https://distill.pub/2016/deconv-checkerboard/>

Upsampling in a deep network

- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1
 - For stride 2, dilate the input by inserting rows and columns of zeros between adjacent entries
 - Sometimes called convolution with *fractional input stride* 1/2

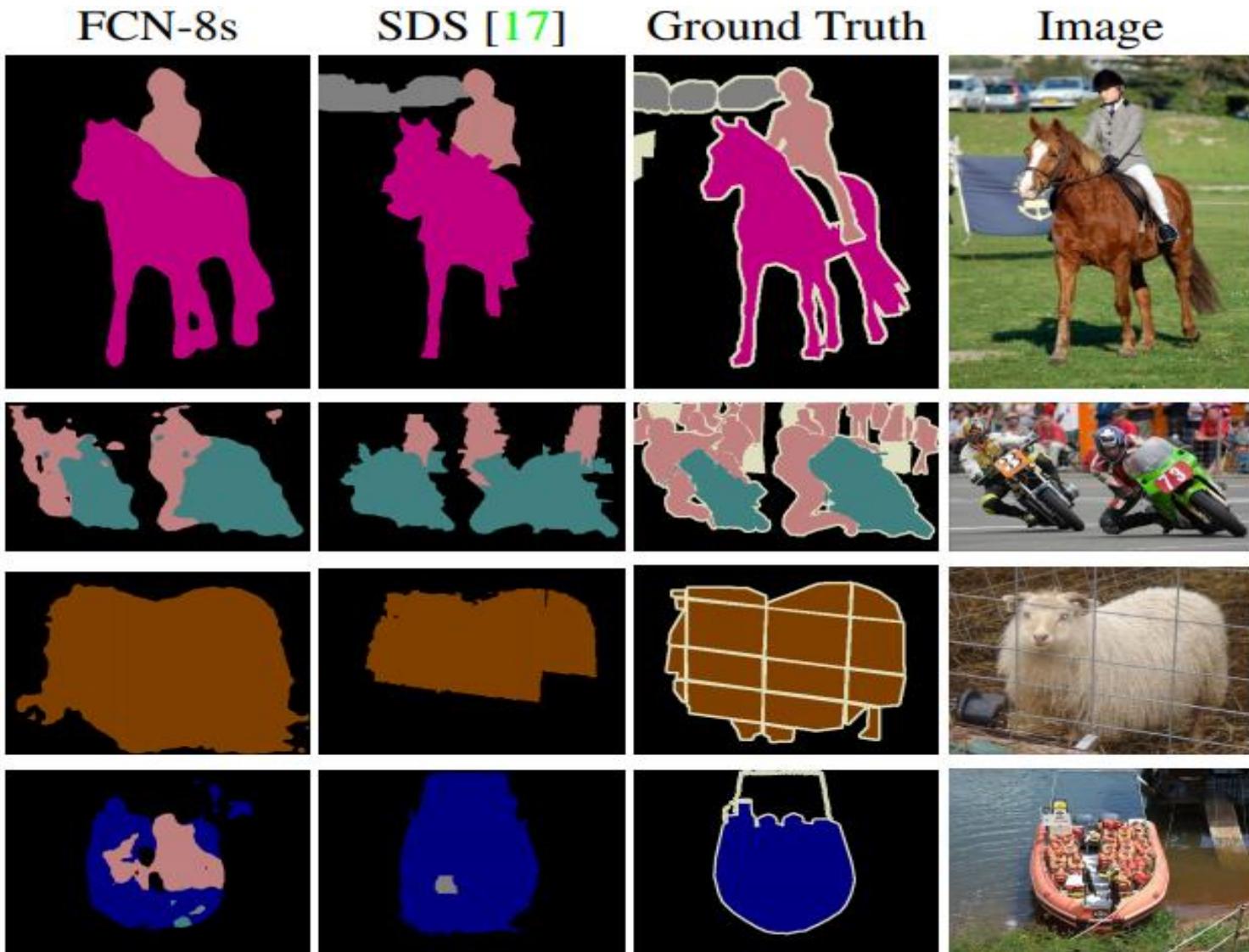


Q: What 3x3 filter would correspond to bilinear upsampling?

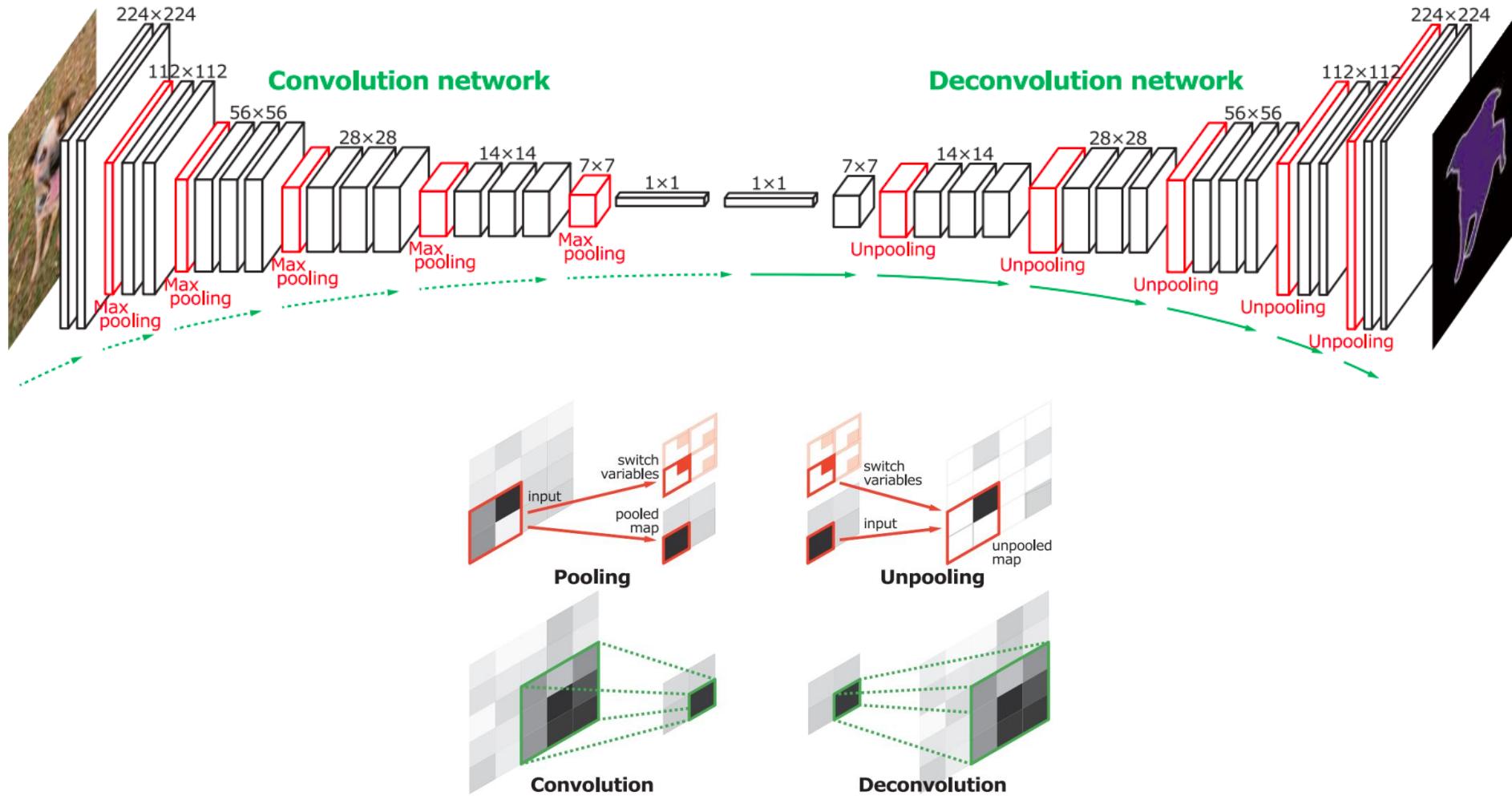
| | | |
|---------------|---------------|---------------|
| $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{4}$ |
| $\frac{1}{2}$ | 1 | $\frac{1}{2}$ |
| $\frac{1}{4}$ | $\frac{1}{2}$ | $\frac{1}{4}$ |

“Fully convolutional (FCN)” results

- Takes advantage of pre-training from classification
- Applied to objects and scenes (NYUD v2)
- But feature pooling reduces spatial sensitivity and resolution



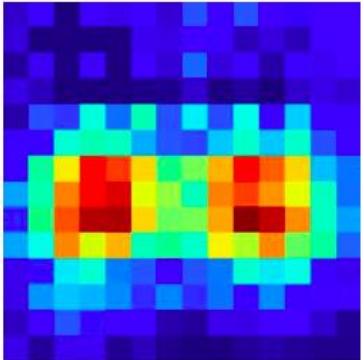
DeconvNet



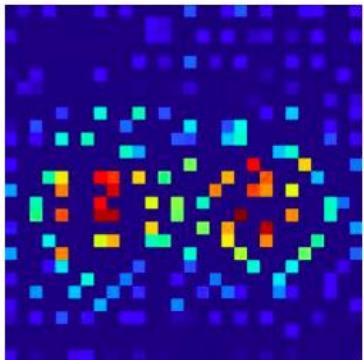
DeconvNet



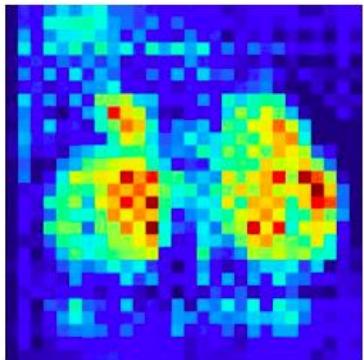
Original image



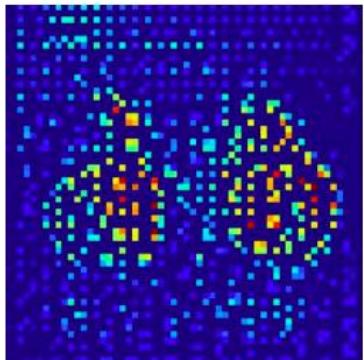
14x14 deconv



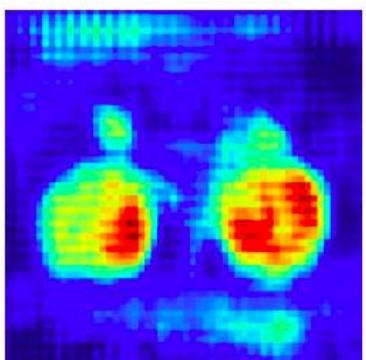
28x28 unpooling



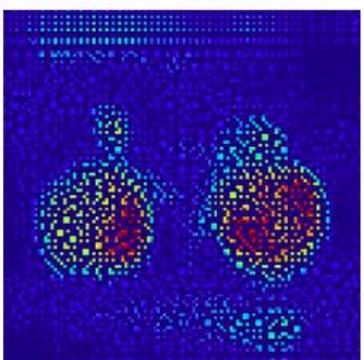
28x28 deconv



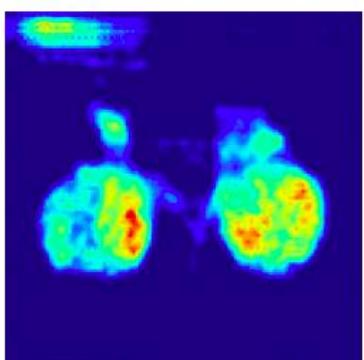
54x54 unpooling



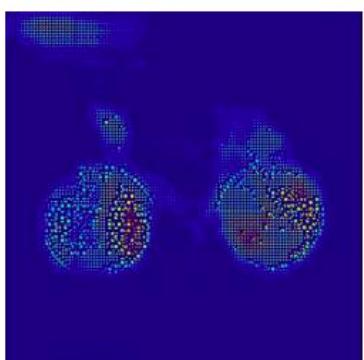
54x54 deconv



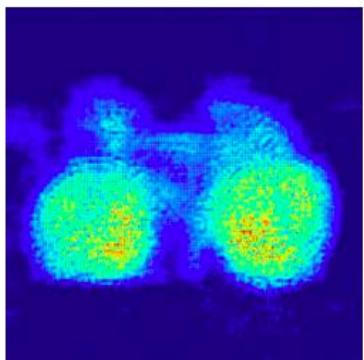
112x112 unpooling



112x112 deconv



224x224 unpooling

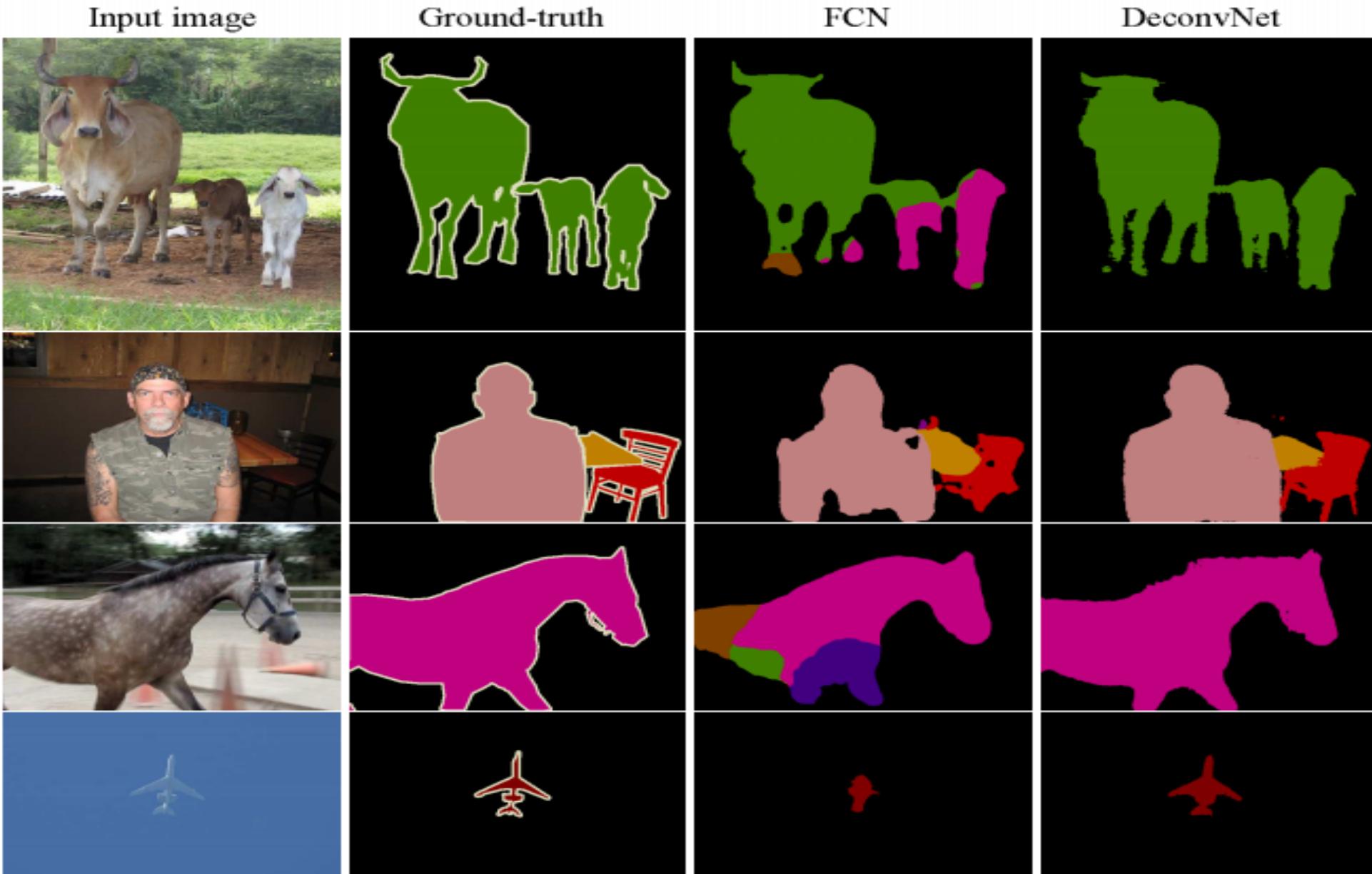


224x224 deconv

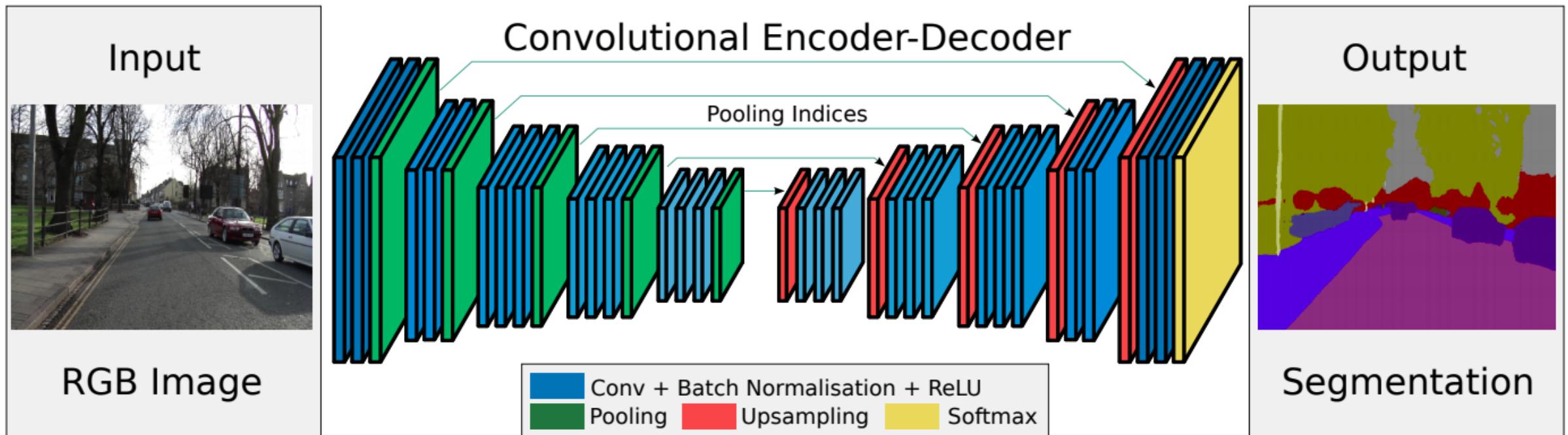
DeconvNet results

| PASCAL VOC 2012 | mIoU |
|-------------------------------|------|
| FCN-8 | 62.2 |
| DeconvNet | 69.6 |
| Ensemble of DeconvNet and FCN | 71.7 |

“DeconNet” results



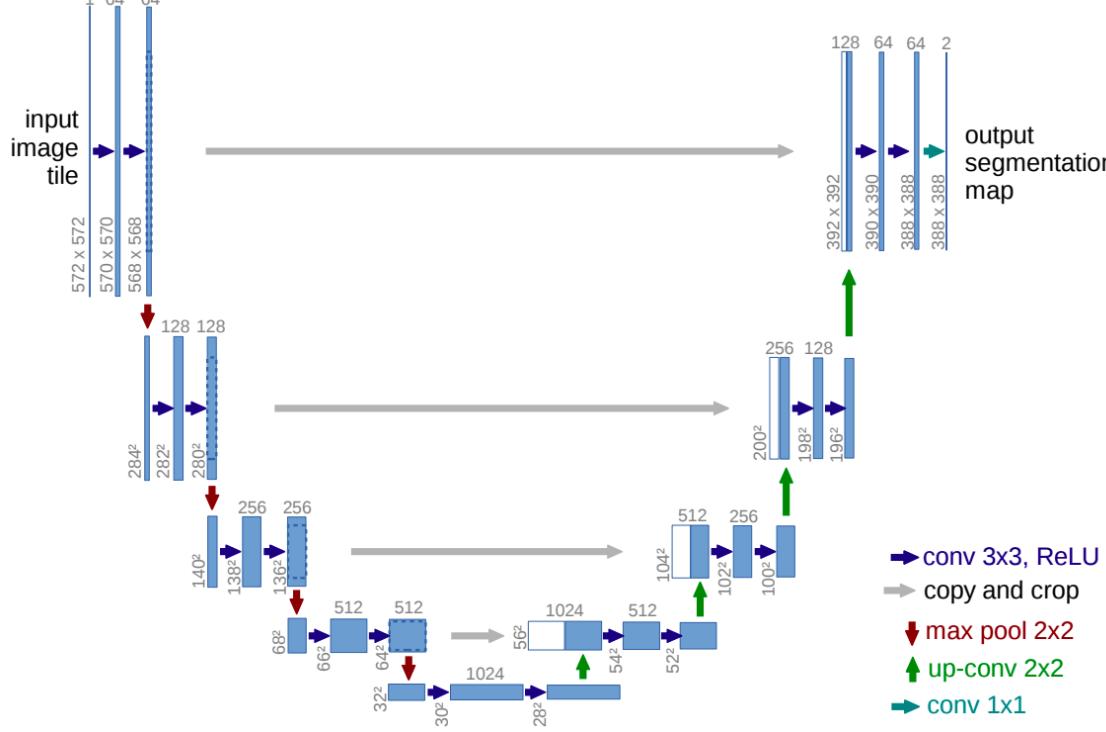
Similar architecture: SegNet



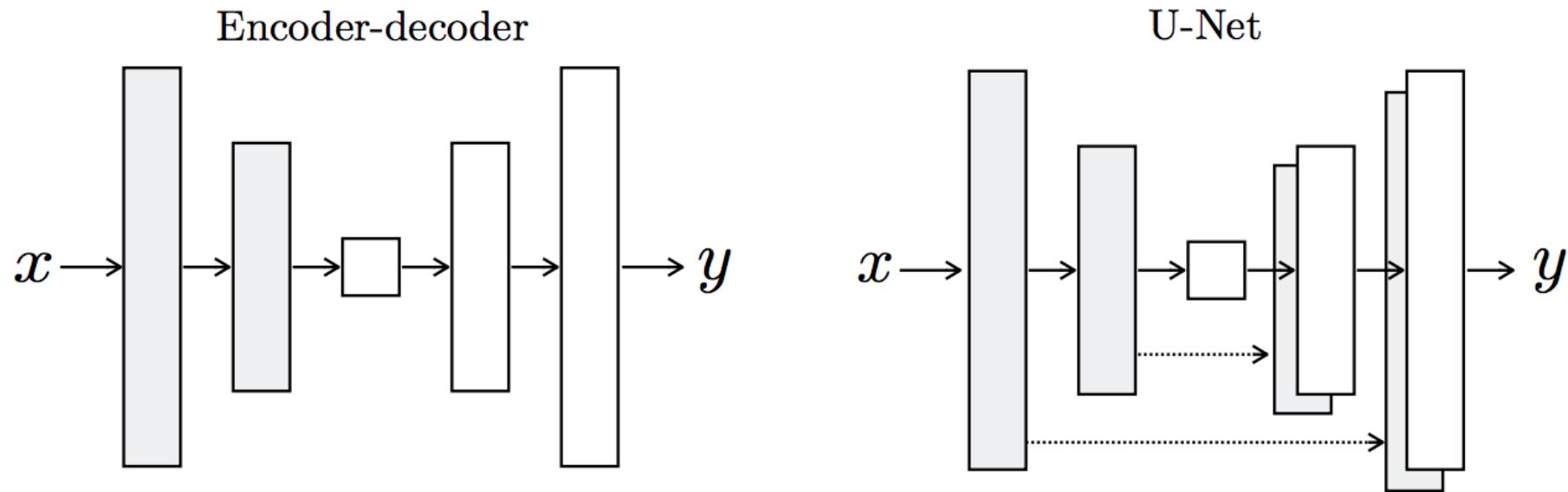
Drop the FC layers,
get better results

U-Net

- Like FCN, fuse upsampled higher-level feature maps with higher-resolution, lower-level feature maps
- Unlike FCN, fuse by concatenation, predict at the end



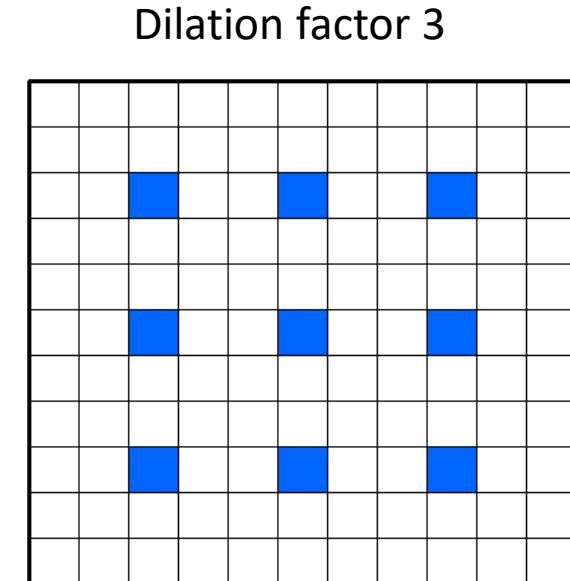
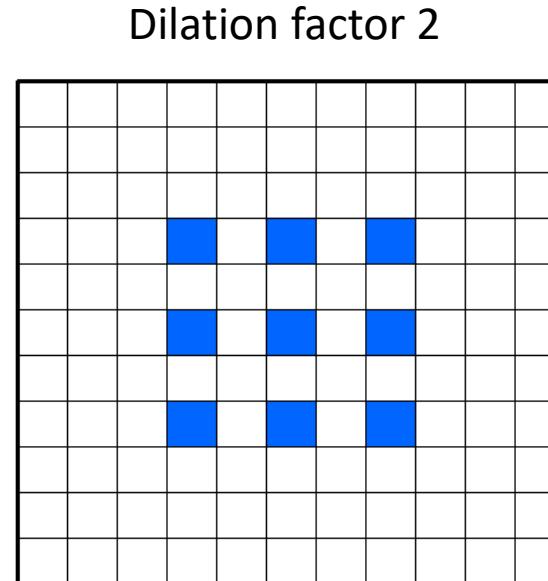
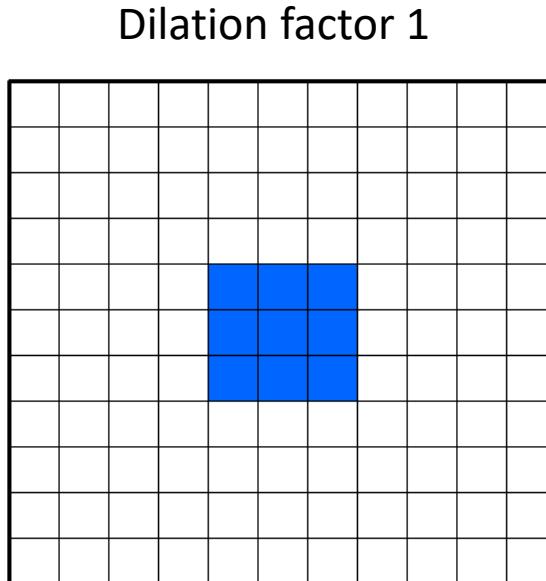
Summary of upsampling architectures



[Figure source](#)

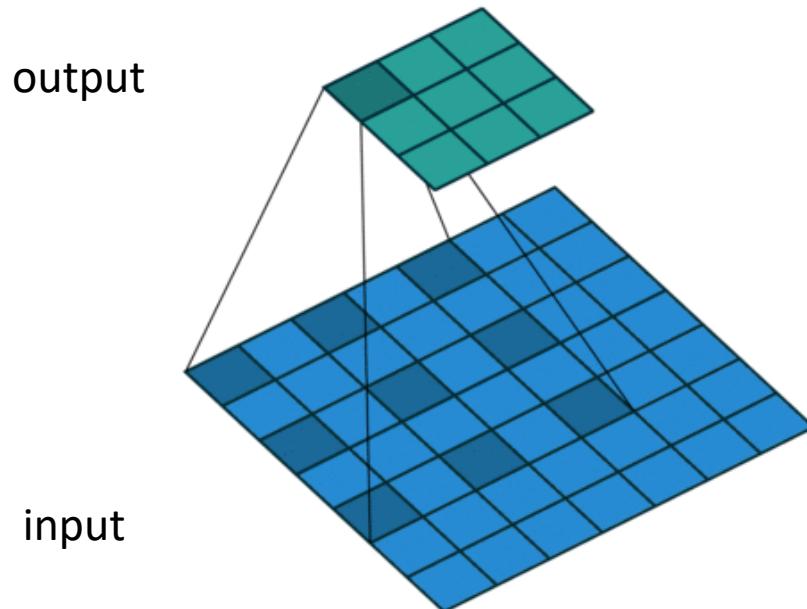
Dilated convolutions

- Idea: instead of reducing spatial resolution of feature maps, use a large sparse filter
 - Also known as *à trous* convolution



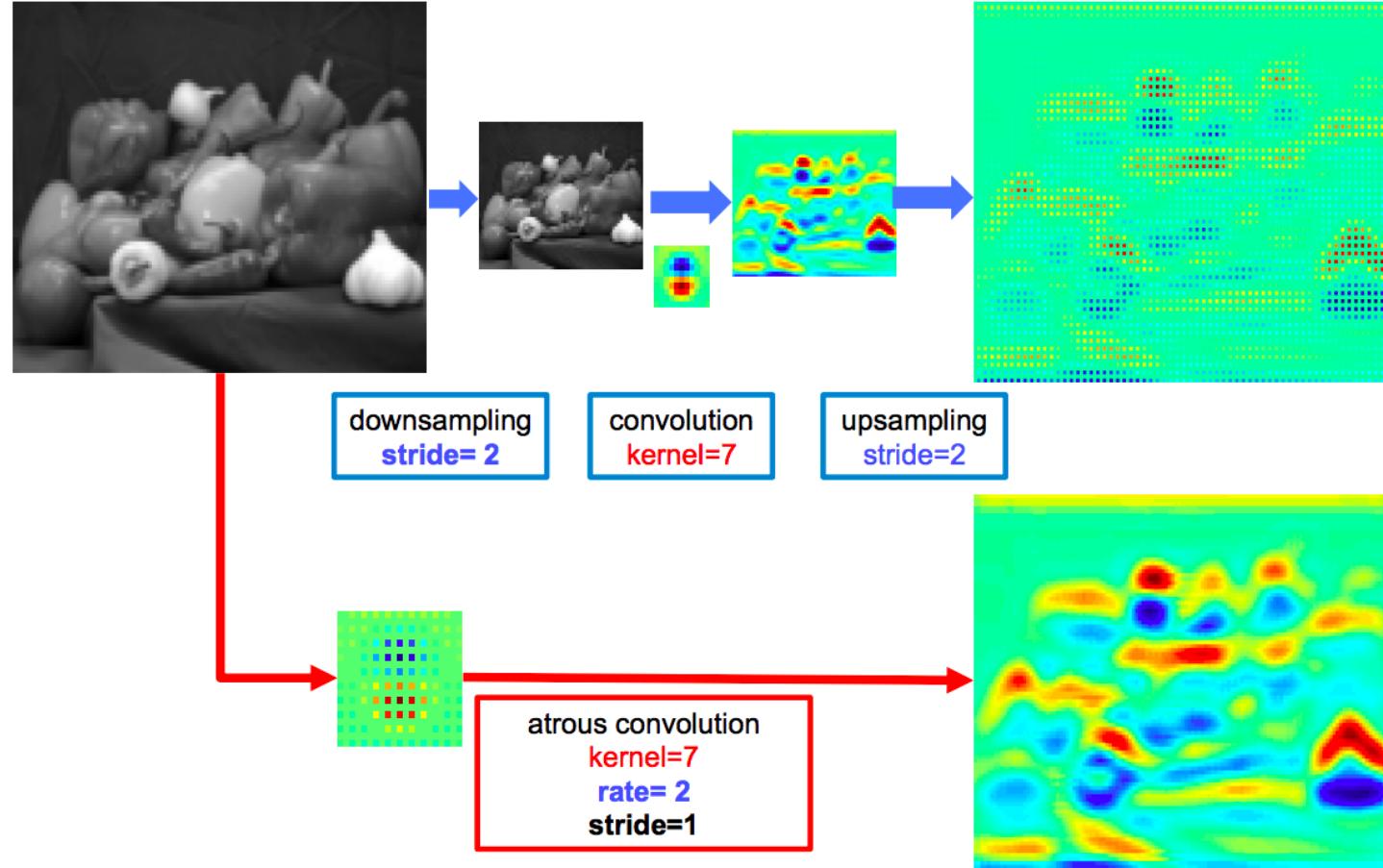
Dilated convolutions

- Idea: instead of reducing spatial resolution of feature maps, use a large sparse filter



Like 2x downsampling
followed by 3x3 convolution
followed by 2x upsampling

Dilated convolutions



L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. Yuille, [DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs](#), PAMI 2017

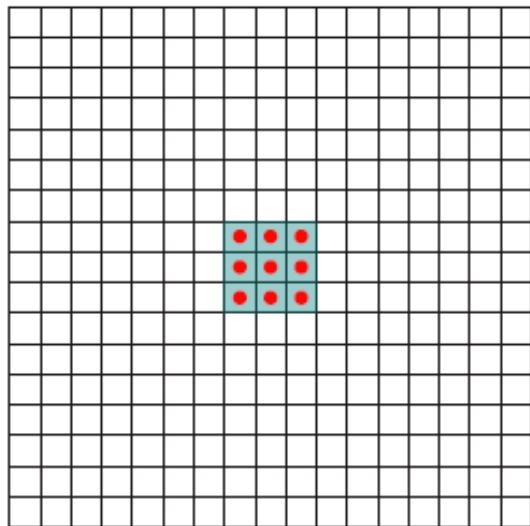
Dilated convolutions

- Can be used in FCN to remove downsampling:
change stride of max pooling layer from 2 to 1,
dilate subsequent convolutions by factor of 2.

Dilated convolutions

- Can increase receptive field size exponentially with a linear growth in the number of parameters

Feature map 1 (F_1)
produced from F_0 by 1-
dilated convolution



Receptive field: 3x3

Receptive field: 7x7

Receptive field: 15x15

Dilated convolutions

- Context module with dilation
 - Returns same number of feature maps at the same resolution as the input, so can be plugged in to replace components of existing dense prediction architectures
 - Requires identity initialization

| Layer | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-----------------|--------------|--------------|--------------|----------------|----------------|----------------|----------------|----------------|
| Convolution | 3×3 | 3×3 | 3×3 | 3×3 | 3×3 | 3×3 | 3×3 | 1×1 |
| Dilation | 1 | 1 | 2 | 4 | 8 | 16 | 1 | 1 |
| Truncation | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No |
| Receptive field | 3×3 | 5×5 | 9×9 | 17×17 | 33×33 | 65×65 | 67×67 | 67×67 |
| Output channels | | | | | | | | |
| Basic | C | C | C | C | C | C | C | C |
| Large | $2C$ | $2C$ | $4C$ | $8C$ | $16C$ | $32C$ | $32C$ | C |

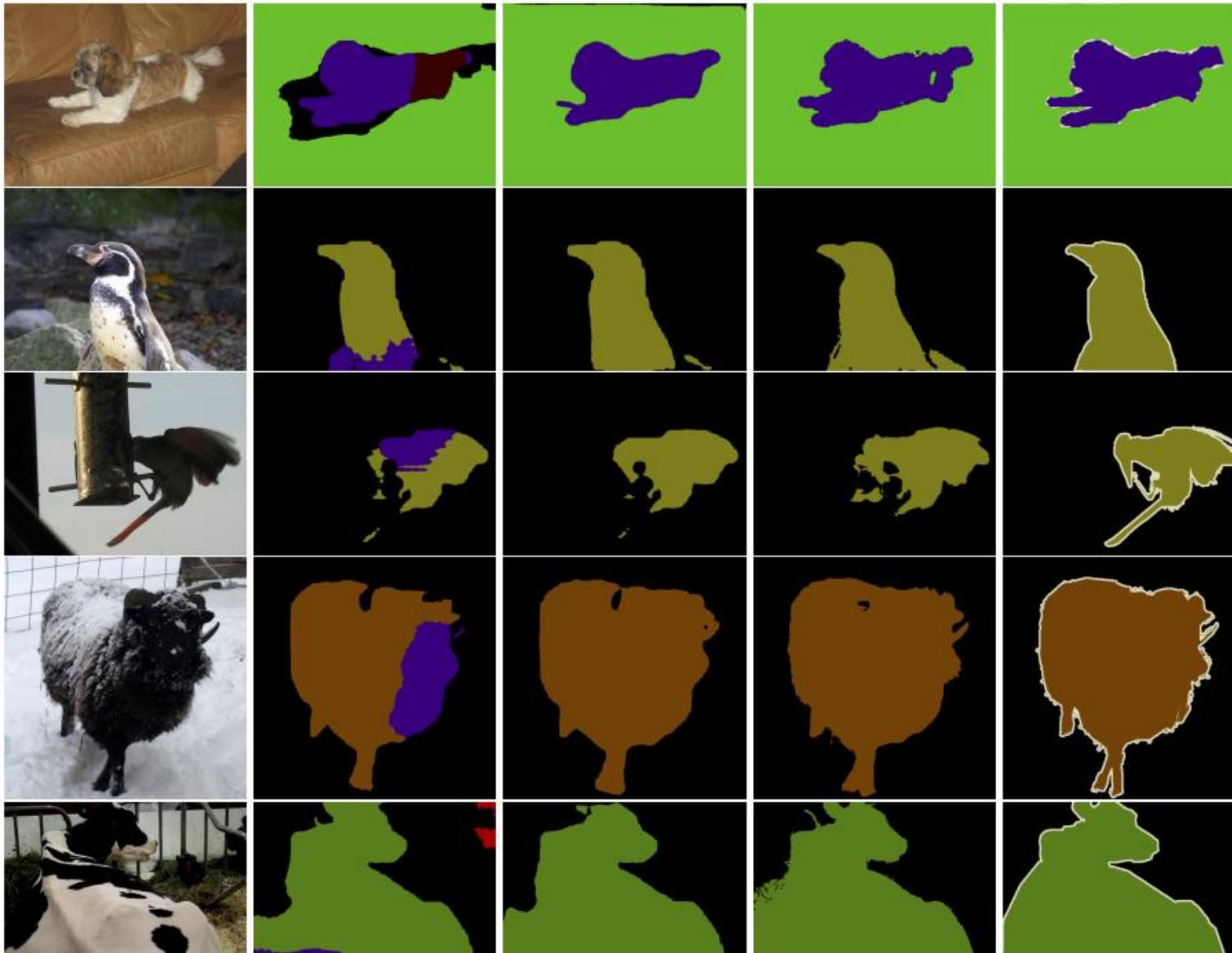
Dilated convolutions: Evaluation

Results on VOC 2012

| | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv | mean IoU |
|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| Front end | 86.3 | 38.2 | 76.8 | 66.8 | 63.2 | 87.3 | 78.7 | 82 | 33.7 | 76.7 | 53.5 | 73.7 | 76 | 76.6 | 83 | 51.9 | 77.8 | 44 | 79.9 | 66.3 | 69.8 |
| Front + Basic | 86.4 | 37.6 | 78.5 | 66.3 | 64.1 | 89.9 | 79.9 | 84.9 | 36.1 | 79.4 | 55.8 | 77.6 | 81.6 | 79 | 83.1 | 51.2 | 81.3 | 43.7 | 82.3 | 65.7 | 71.3 |
| Front + Large | 87.3 | 39.2 | 80.3 | 65.6 | 66.4 | 90.2 | 82.6 | 85.8 | 34.8 | 81.9 | 51.7 | 79 | 84.1 | 80.9 | 83.2 | 51.2 | 83.2 | 44.7 | 83.4 | 65.6 | 72.1 |

*Front end: re-implementation of FCN-8 with last two pooling layers dropped (5% better than original FCN-8)

Dilated convolutions: Evaluation



(a) Image

(b) Front end

(c) + Context

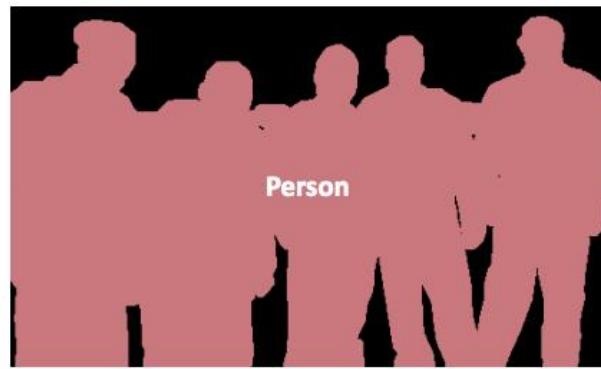
(d) + CRF-RNN

(e) Ground truth

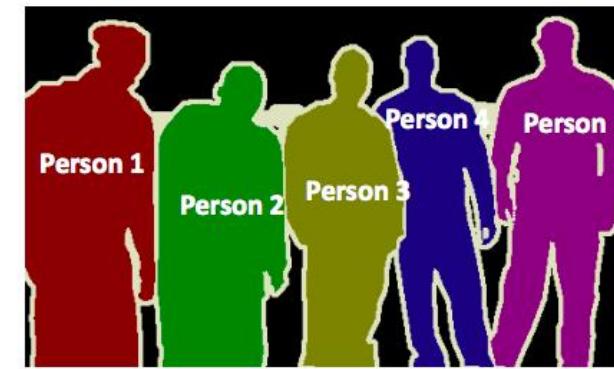
Instance segmentation



Object Detection



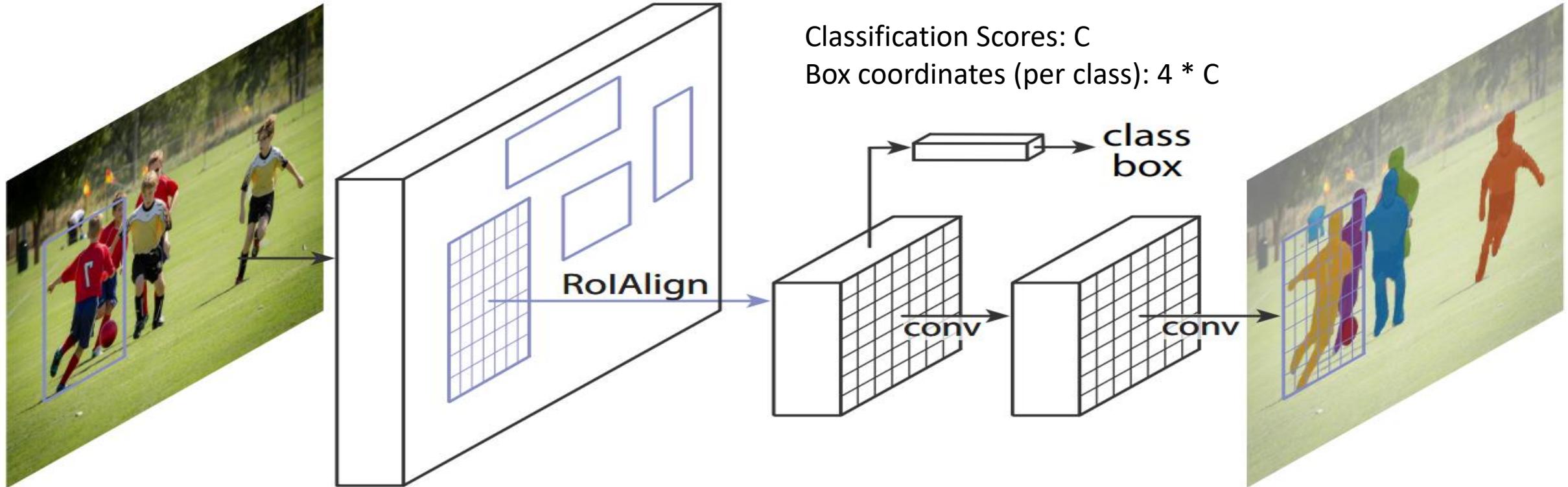
Semantic Segmentation



Instance Segmentation



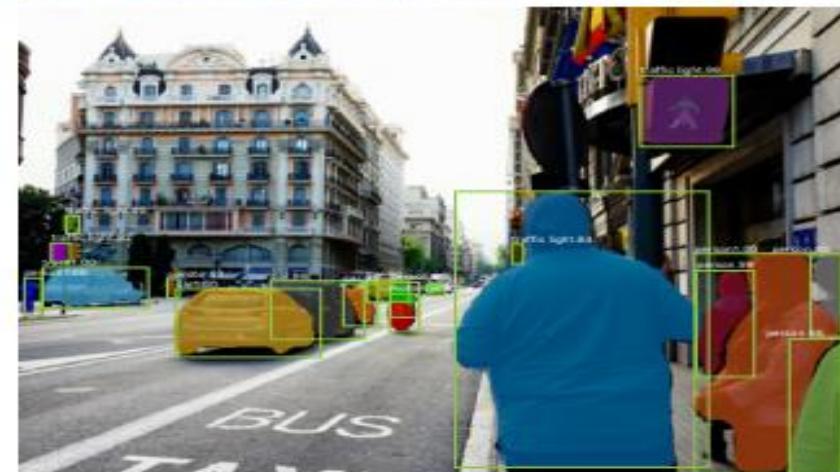
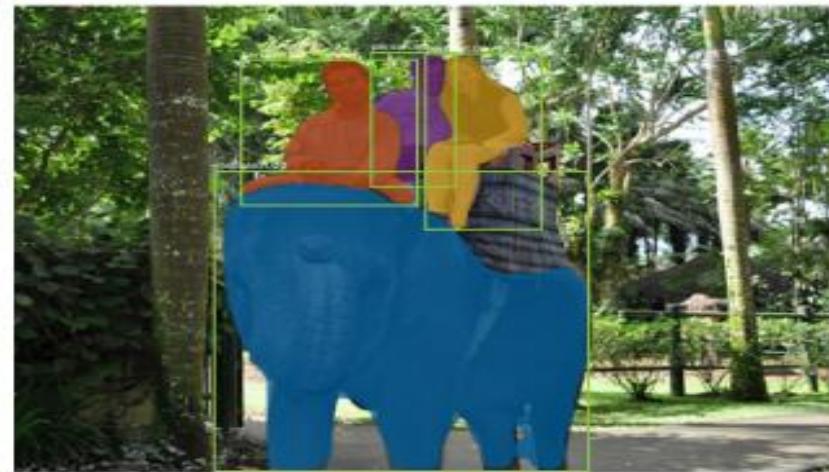
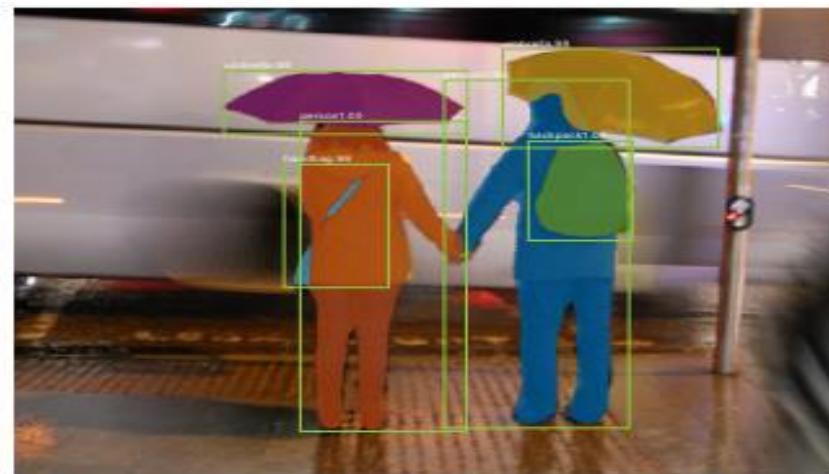
Instance Segmentation: Mask R-CNN (He et al. ICCV 2017)



Mask R-CNN - extends Faster R-CNN by adding a branch for predicting an object mask in parallel with the existing branch for bounding box recognition.

Mask R-CNN is simple to train and adds only a small overhead to Faster R-CNN, running at 5 fps.

Instance Segmentation: Mask R-CNN (He et al. ICCV 2017)



Mask R-CNN results on the COCO test set. These results are based on ResNet-101, achieving a mask AP of 35.7 and running at 5 fps. Masks are shown in color, and bounding box, category, and confidences are also shown.

Instance segmentation results on COCO

| | backbone | AP | AP ₅₀ | AP ₇₅ | AP _S | AP _M | AP _L |
|--------------------|-----------------------|-------------|------------------|------------------|-----------------|-----------------|-----------------|
| MNC [10] | ResNet-101-C4 | 24.6 | 44.3 | 24.8 | 4.7 | 25.9 | 43.6 |
| FCIS [26] +OHEM | ResNet-101-C5-dilated | 29.2 | 49.5 | - | 7.1 | 31.3 | 50.0 |
| FCIS+++ [26] +OHEM | ResNet-101-C5-dilated | 33.6 | 54.5 | - | - | - | - |
| Mask R-CNN | ResNet-101-C4 | 33.1 | 54.9 | 34.8 | 12.1 | 35.6 | 51.1 |
| Mask R-CNN | ResNet-101-FPN | 35.7 | 58.0 | 37.8 | 15.5 | 38.1 | 52.4 |
| Mask R-CNN | ResNeXt-101-FPN | 37.1 | 60.0 | 39.4 | 16.9 | 39.9 | 53.5 |

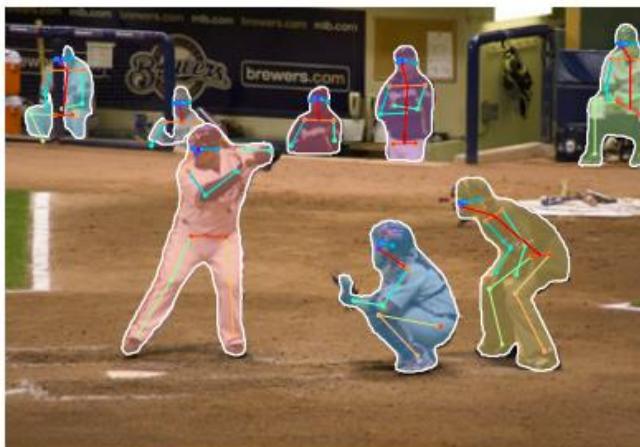
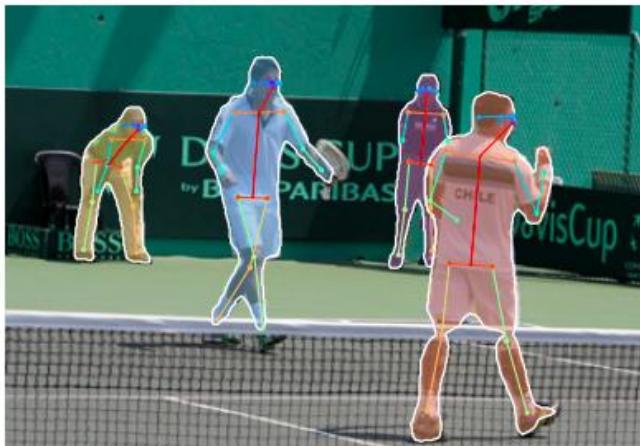
AP at different IoU
thresholds

AP for different size
instances

Other dense prediction problems

Keypoint prediction

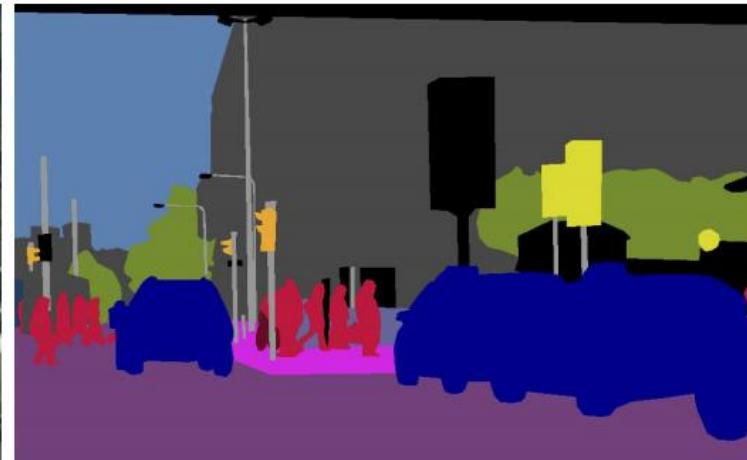
- Given K keypoints, train model to predict $K m \times m$ one-hot maps with cross-entropy losses over m^2 outputs



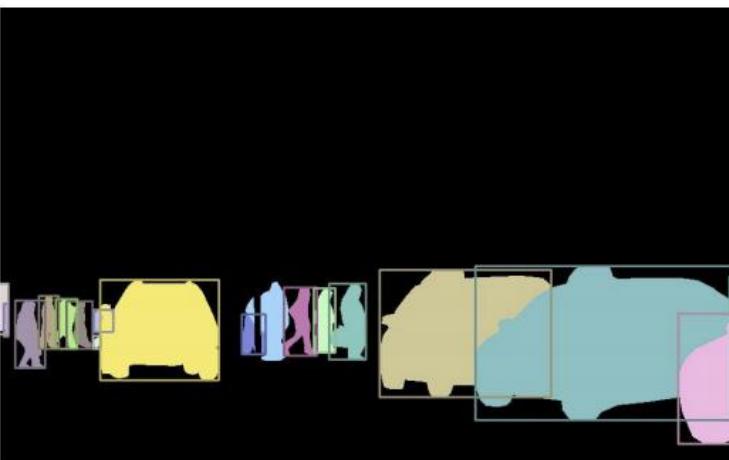
Panoptic Segmentation



(a) image



(b) semantic segmentation



(c) instance segmentation



(d) panoptic segmentation

Panoptic feature pyramid networks

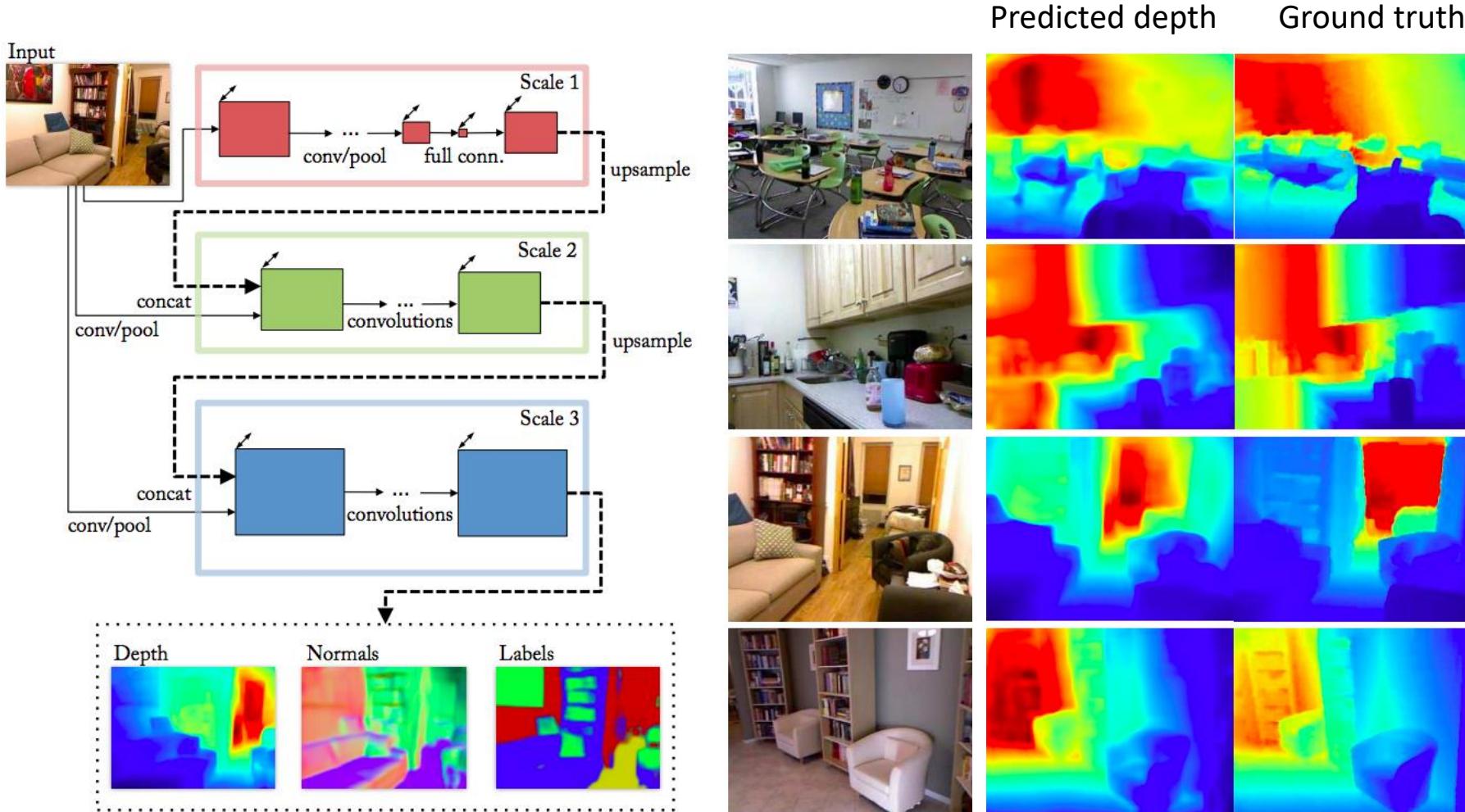


Figure 2: Panoptic FPN results on COCO (top) and Cityscapes (bottom) using a single ResNet-101-FPN network.

Even more dense prediction problems

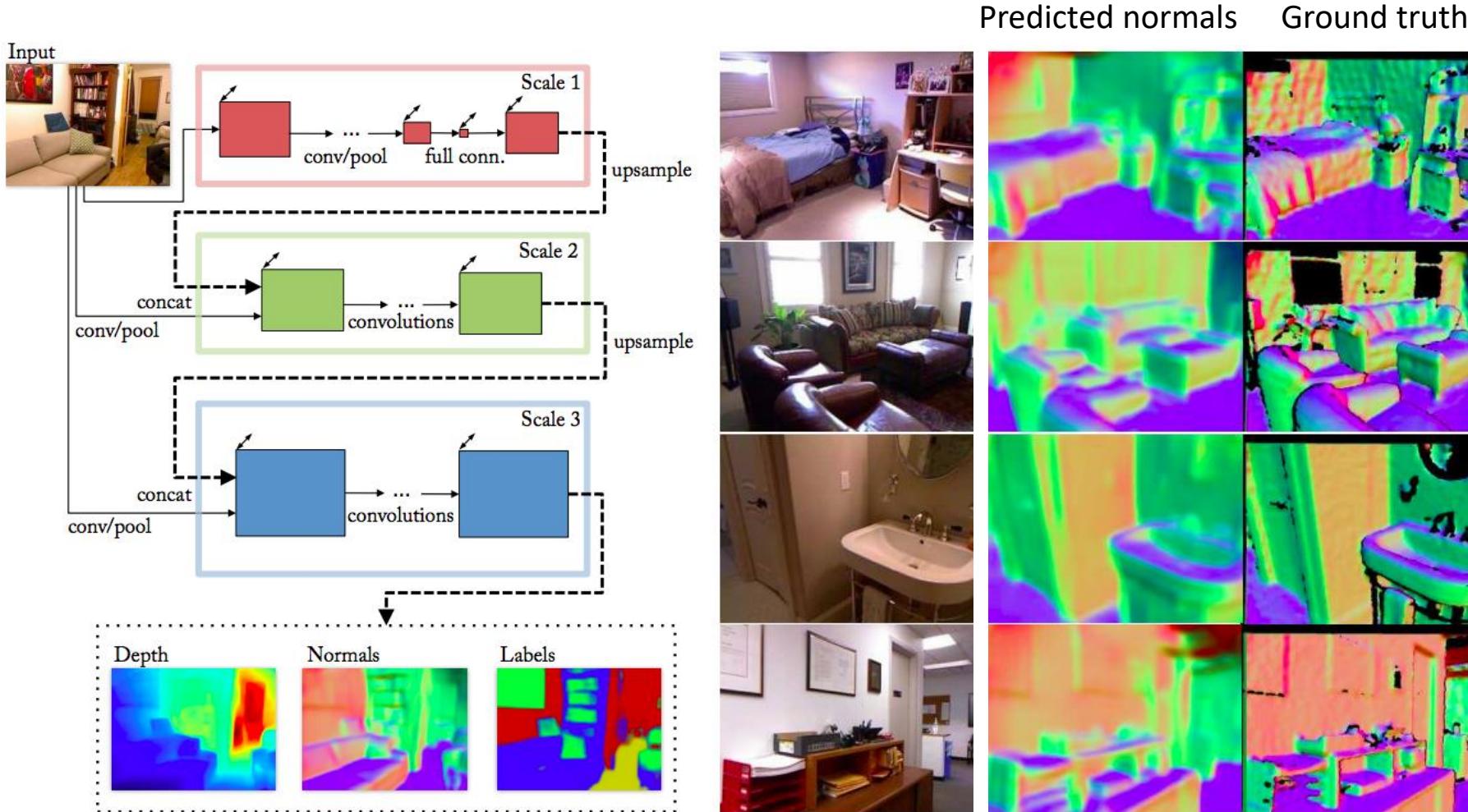
- Depth estimation
- Surface normal estimation
- Colorization
-

Depth and normal estimation



D. Eigen and R. Fergus, [Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture](#), ICCV 2015

Depth and normal estimation



D. Eigen and R. Fergus, [Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture](#), ICCV 2015

Estimation of everything at the same time

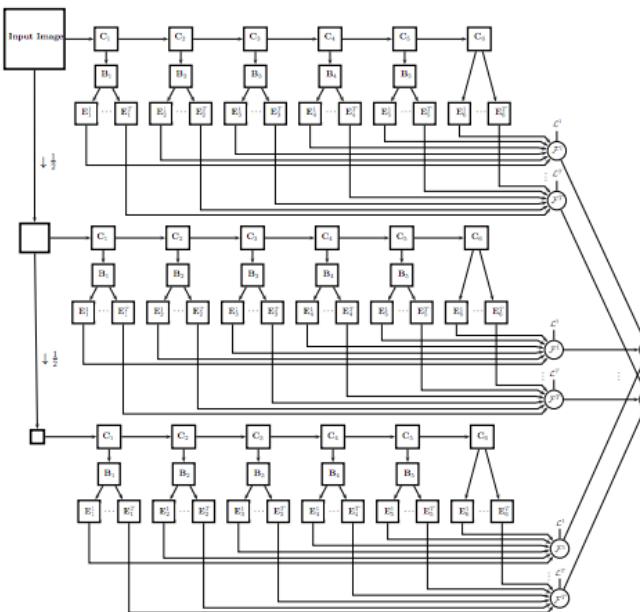
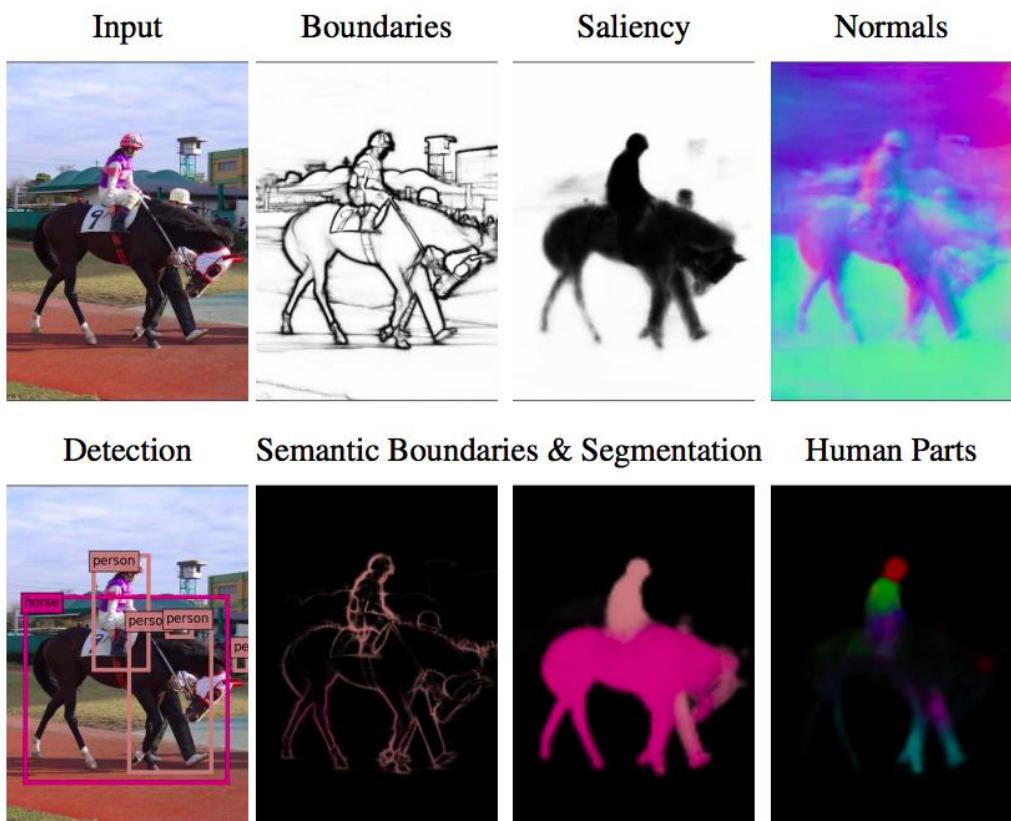
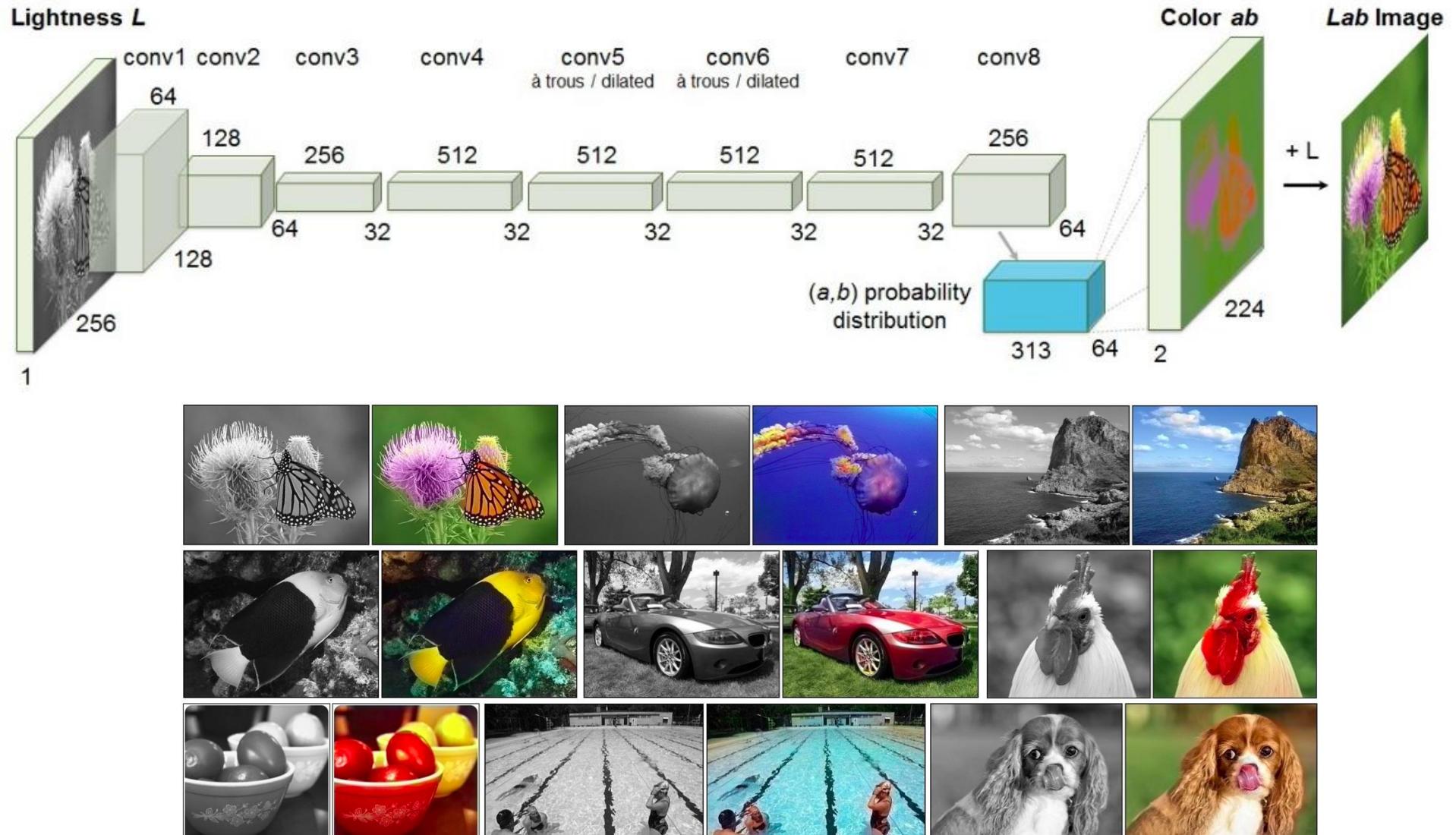


Figure 2: UberNet architecture: an image pyramid is formed by successive down-sampling operations, and each image is processed by a CNN with tied weights; the responses of the network at consecutive layers (C_i) are processed with Batch Normalization (B_i) and then fed to task-specific skip layers (E_i^t); these are combined across network layers (\mathcal{F}^t) and resolutions (\mathcal{S}^t) and trained using task-specific loss functions (\mathcal{L}^t), while the whole architecture is jointly trained end-to-end. For simplicity we omit the interpolation and detection layers mentioned in the text.

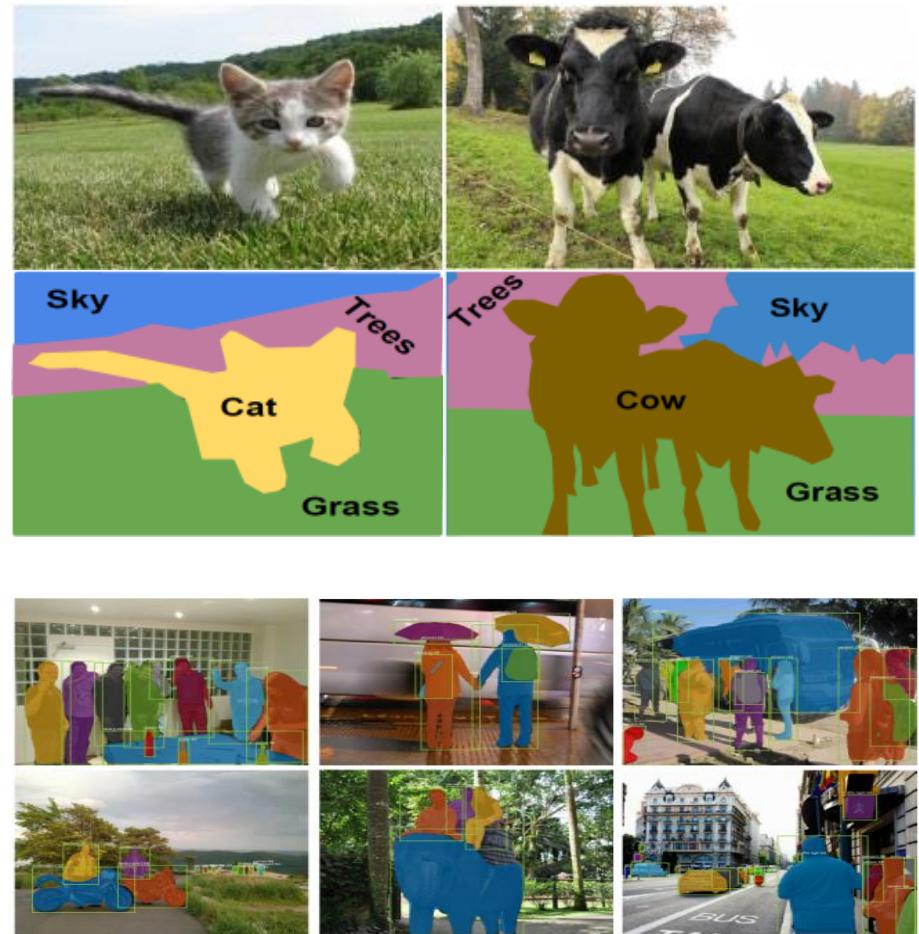
I. Kokkinos, [UberNet: Training a Universal Convolutional Neural Network for Low-, Mid-, and High-Level Vision using Diverse Datasets and Limited Memory](#),

Colorization



Things to remember

- Semantic Segmentation
 - Sliding window inefficient
 - Fully Convolutional: Convolutions at original image resolution is very expensive
 - Fully Convolutional: Down sampling and Up sampling
 - Unsampling usning Unpooling and Transpose Convolution
 - DeconvNet: deeper network
- Instance Segmentation
 - Mak-RCNN: Extension of faster RCNN



Acknowledgement

Thanks to the following courses and corresponding researchers for making their teaching/research material online

- Deep Learning, Stanford University
- Introduction to Deep Learning, University of Illinois at Urbana-Champaign
- Introduction to Deep Learning, Carnegie Mellon University
- Convolutional Neural Networks for Visual Recognition, Stanford University
- Natural Language Processing with Deep Learning, Stanford University
- And Many More