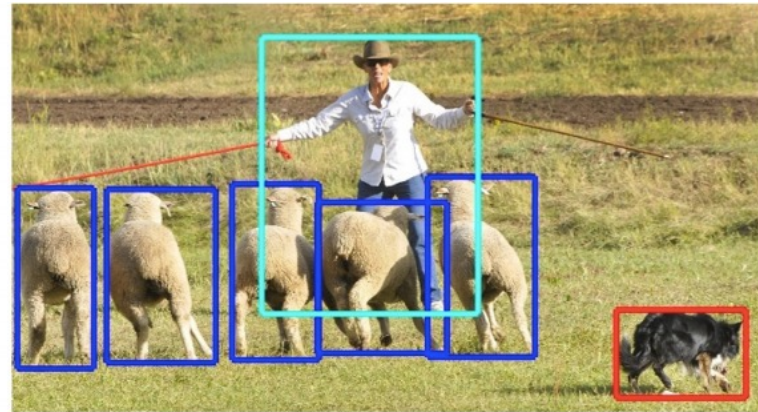


CNNs for dense image labeling



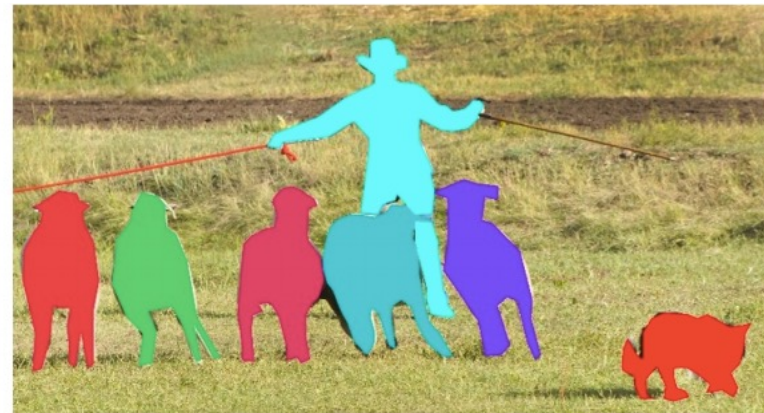
image classification



object detection



semantic segmentation



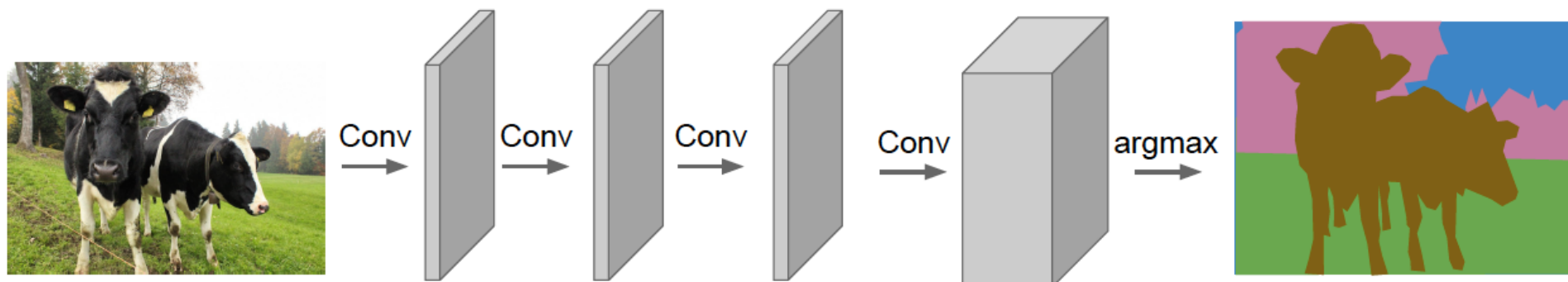
instance segmentation

Outline

- Fully convolutional networks
- Operations for dense prediction
 - Transposed convolutions, unpooling
- Architectures for dense prediction
 - DeconvNet, SegNet, U-Net
- Instance segmentation
 - Mask R-CNN
- Other dense prediction problems

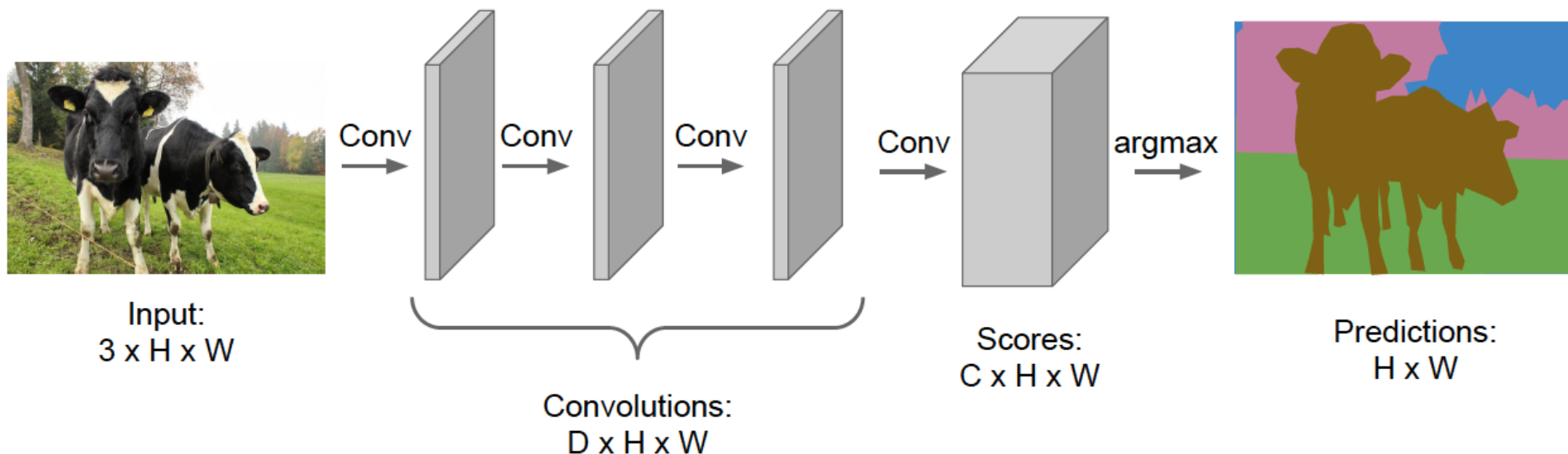
Fully convolutional networks

- Design a network with only convolutional layers, make predictions for all pixels at once



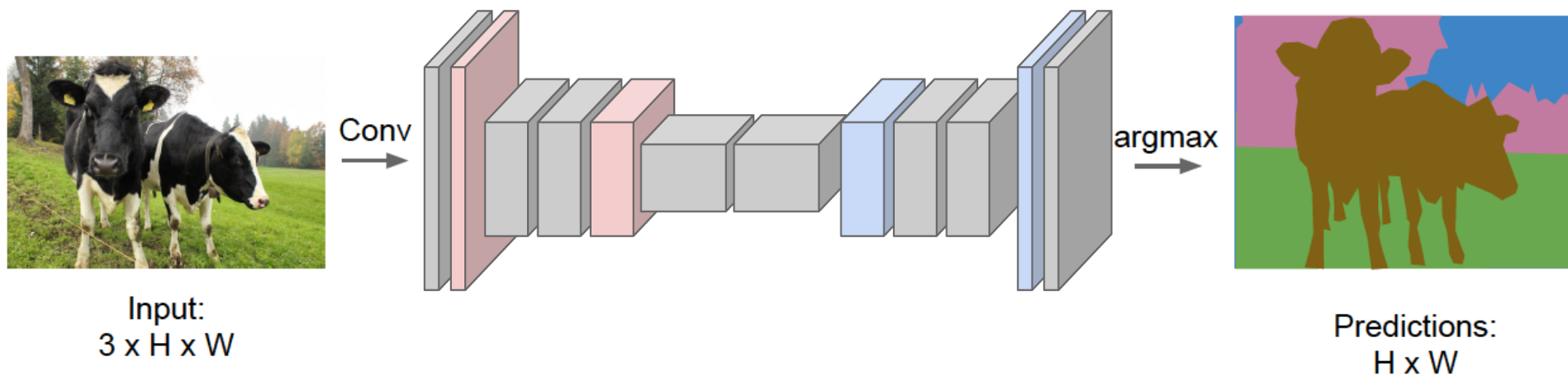
Fully convolutional networks

- Design a network with only convolutional layers, make predictions for all pixels at once
- Can the network operate at full image resolution?

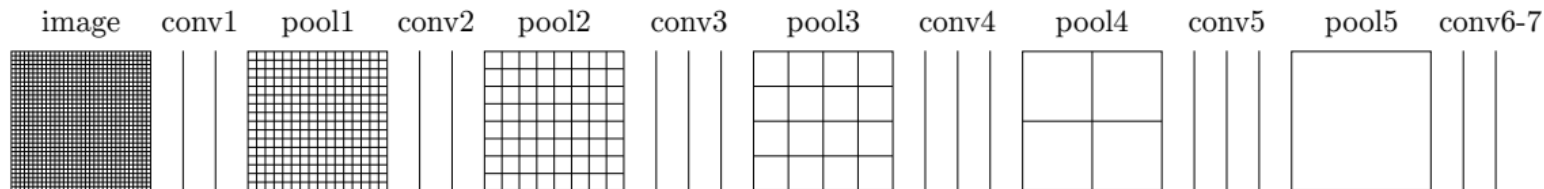


Fully convolutional networks

- Design a network with only convolutional layers, make predictions for all pixels at once
- Can the network operate at full image resolution?
- Practical solution: first downsample, then upsample

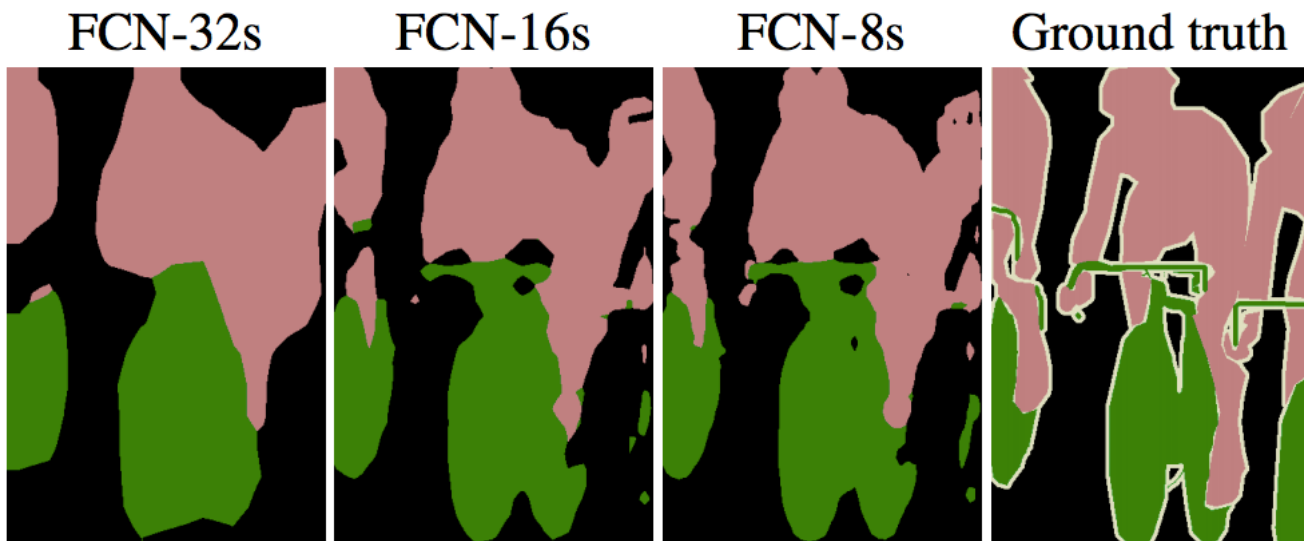


Fully convolutional networks (FCN)



- Predictions by 1x1 conv layers, bilinear upsampling to original image resolution
- Predictions by 1x1 conv layers, learned 2x upsampling using *transposed convolutions*, fusion by summing

Fully convolutional networks (FCN)



Comparison on a subset of PASCAL 2011 validation data:

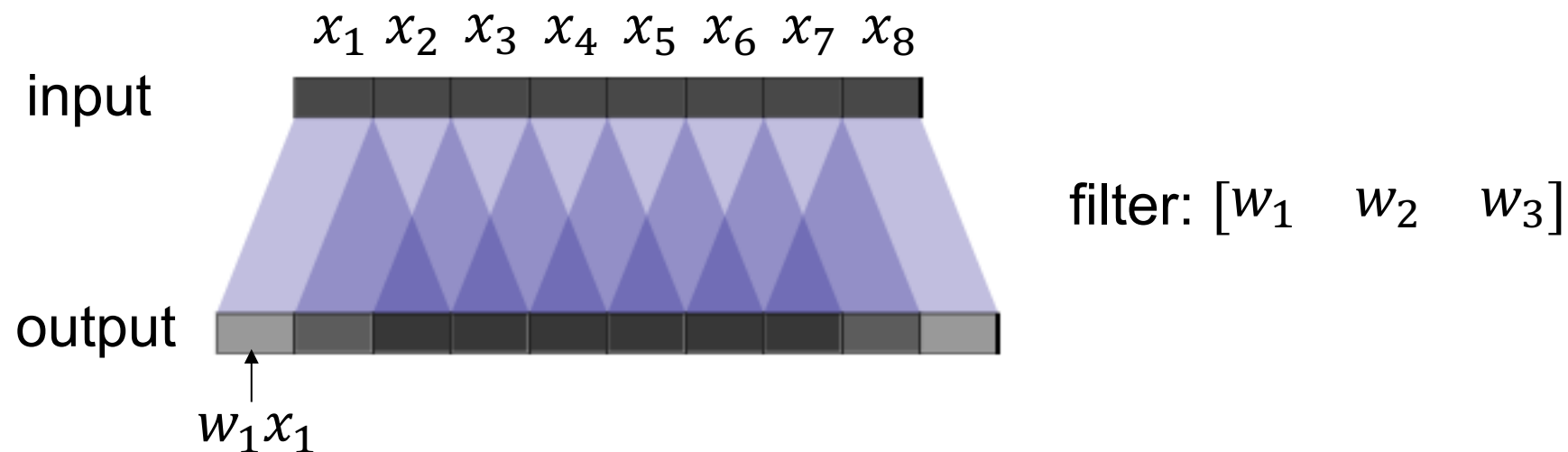
	pixel acc.	mean acc.	mean IU
FCN-32s-fixed	83.0	59.7	45.4
FCN-32s	89.1	73.3	59.4
FCN-16s	90.0	75.7	62.4
FCN-8s	90.3	75.9	62.7

Outline

- Fully convolutional networks
- Operations for dense prediction
 - Transposed convolutions, unpooling

Transposed convolution

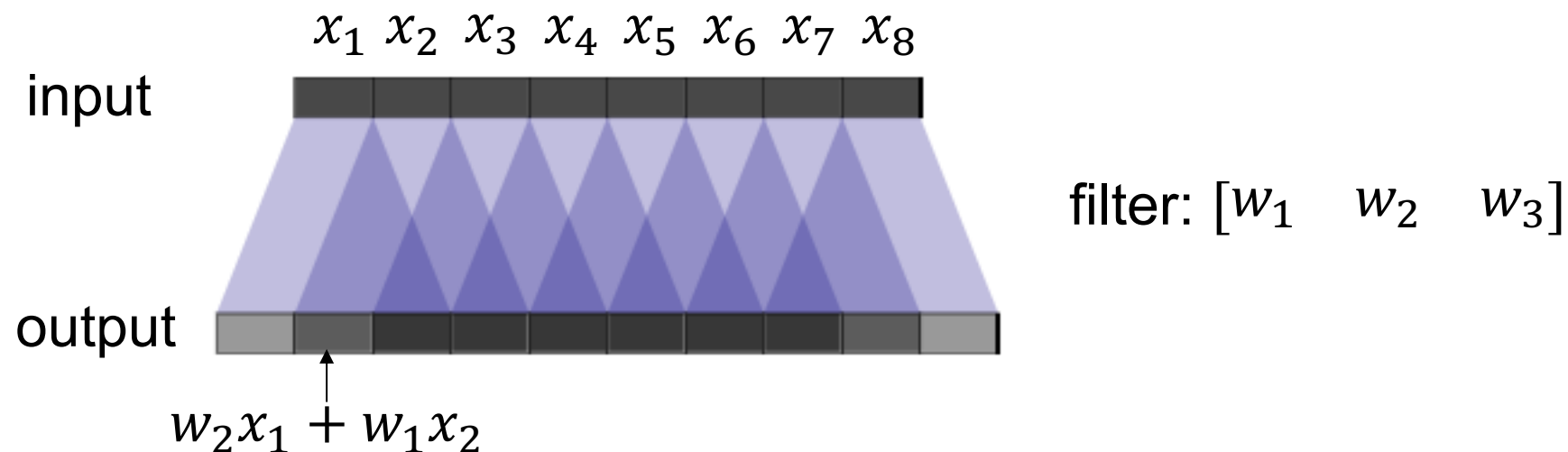
- Use the filter to “paint” in the output: place copies of the filter on the output, multiply by corresponding value in the input, sum where copies of the filter overlap
- 1D example:



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

Transposed convolution

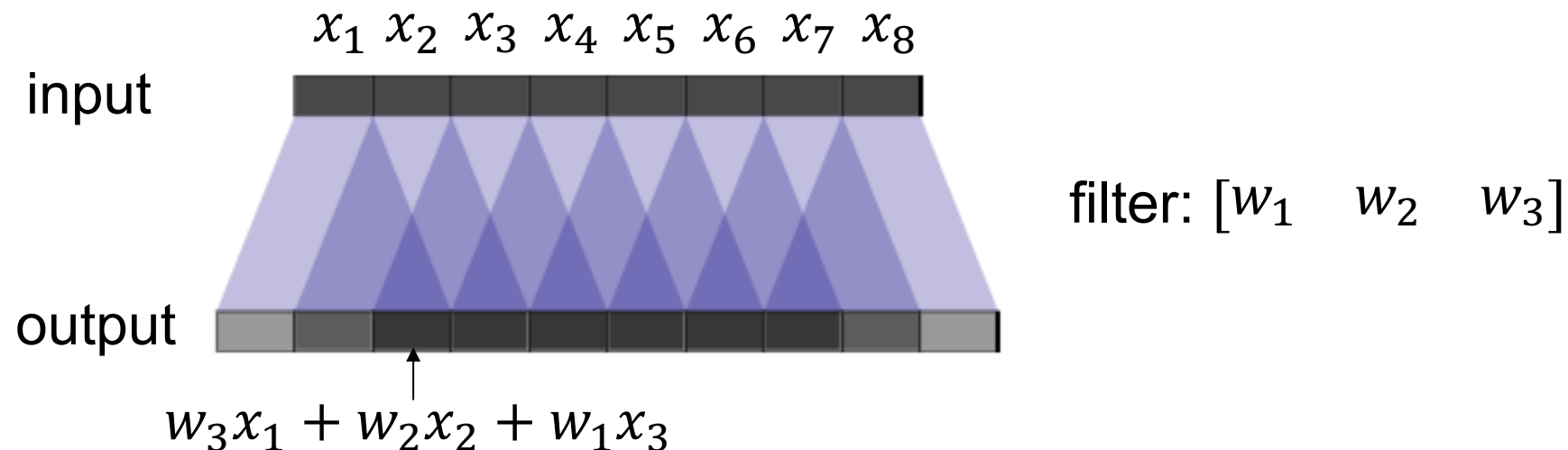
- Use the filter to “paint” in the output: place copies of the filter on the output, multiply by corresponding value in the input, sum where copies of the filter overlap
- 1D example:



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

Transposed convolution

- Use the filter to “paint” in the output: place copies of the filter on the output, multiply by corresponding value in the input, sum where copies of the filter overlap
- 1D example:

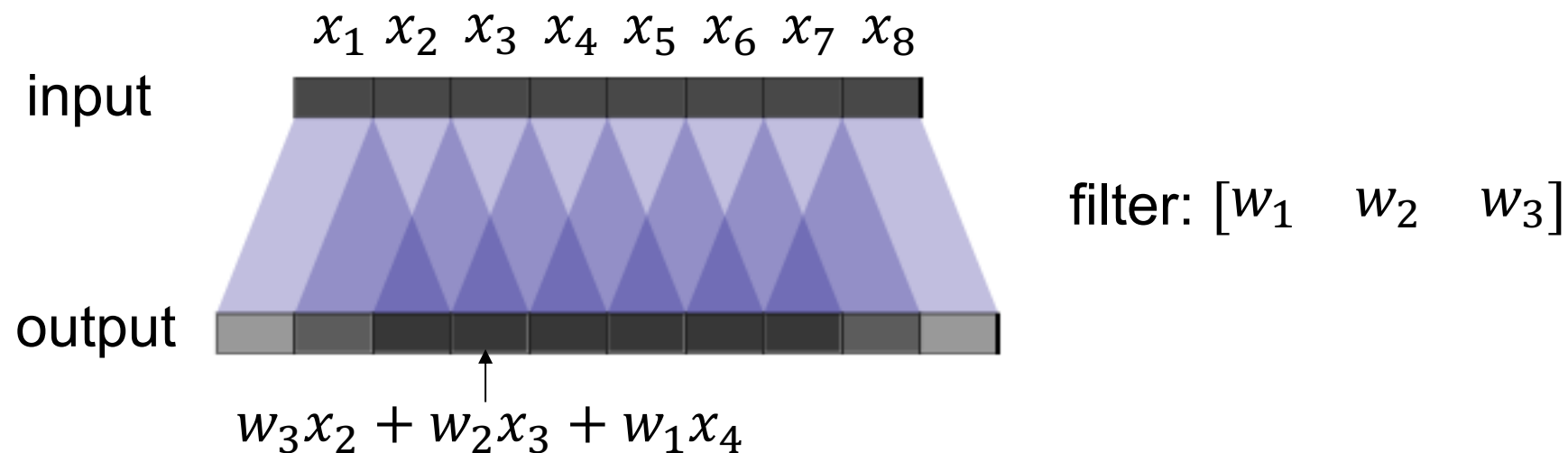


Same as convolution with a flipped filter!

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

Transposed convolution

- Use the filter to “paint” in the output: place copies of the filter on the output, multiply by corresponding value in the input, sum where copies of the filter overlap
- 1D example:

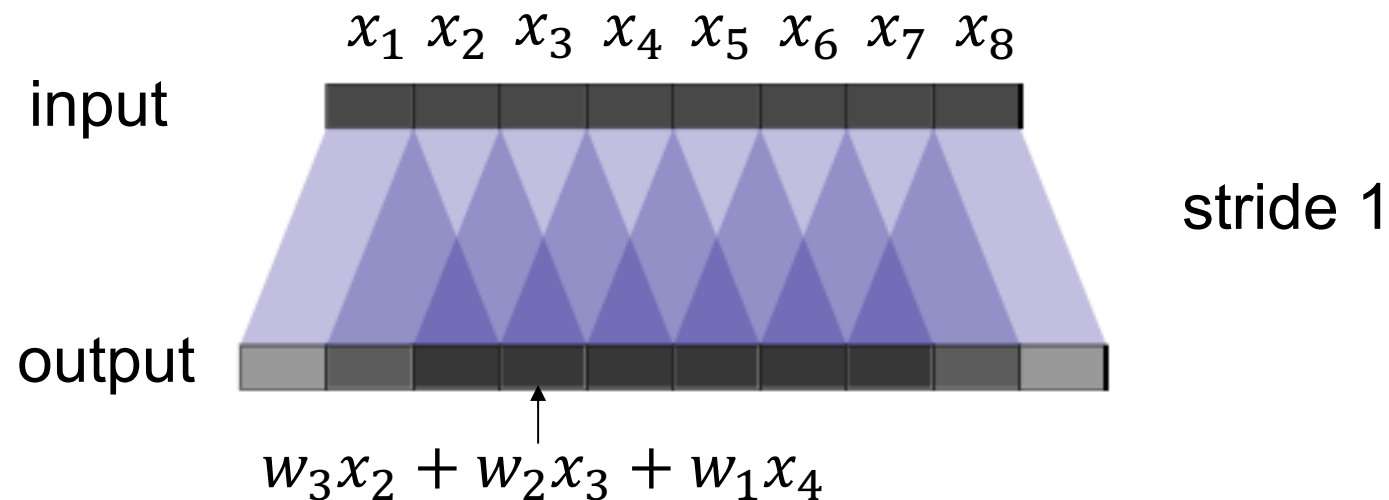


Same as convolution with a flipped filter!

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

Upsampling by transposed convolution

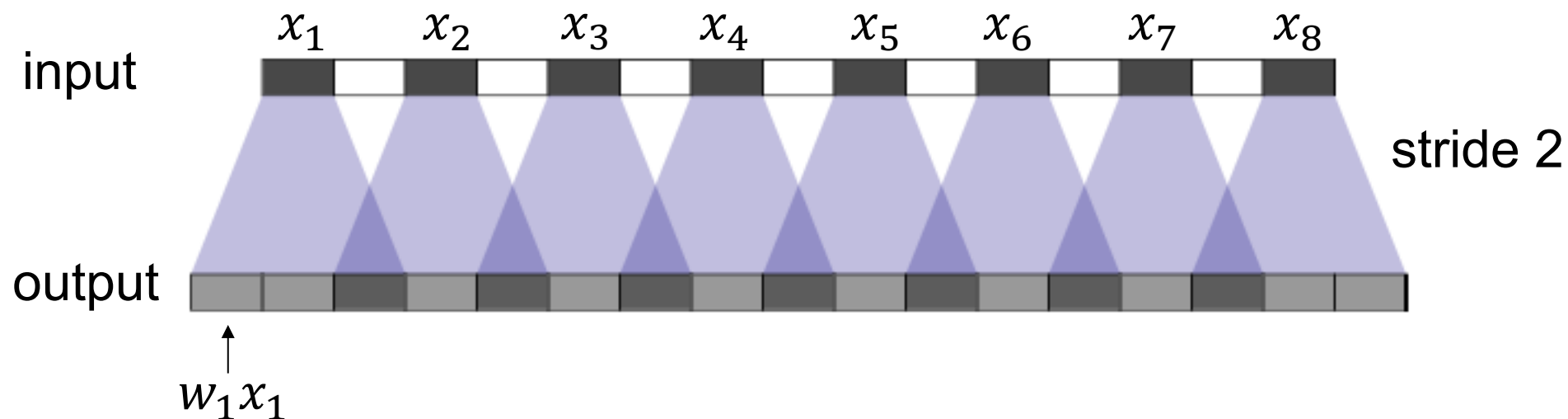
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

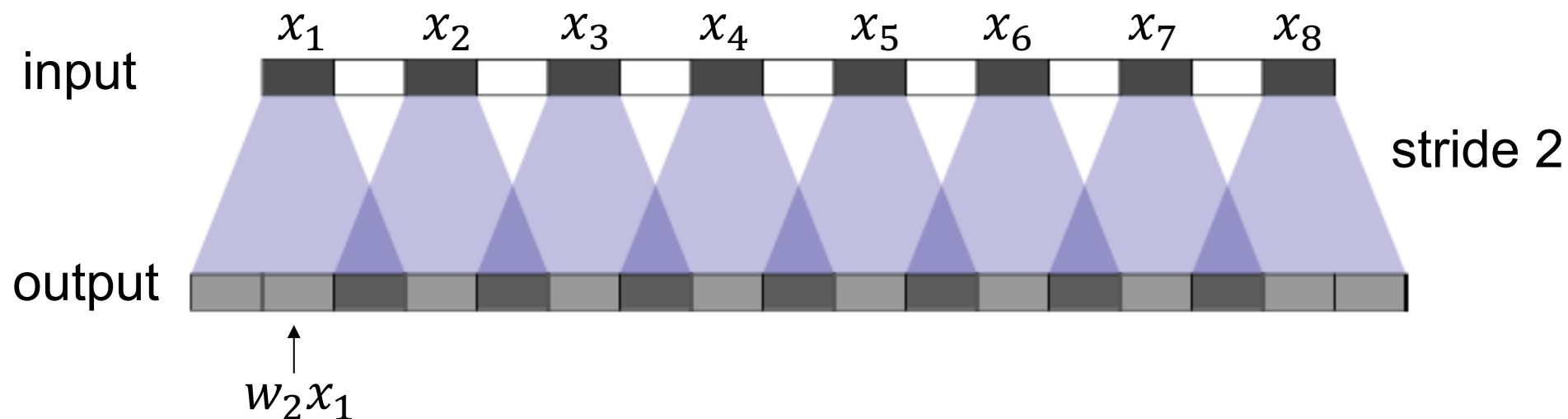
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

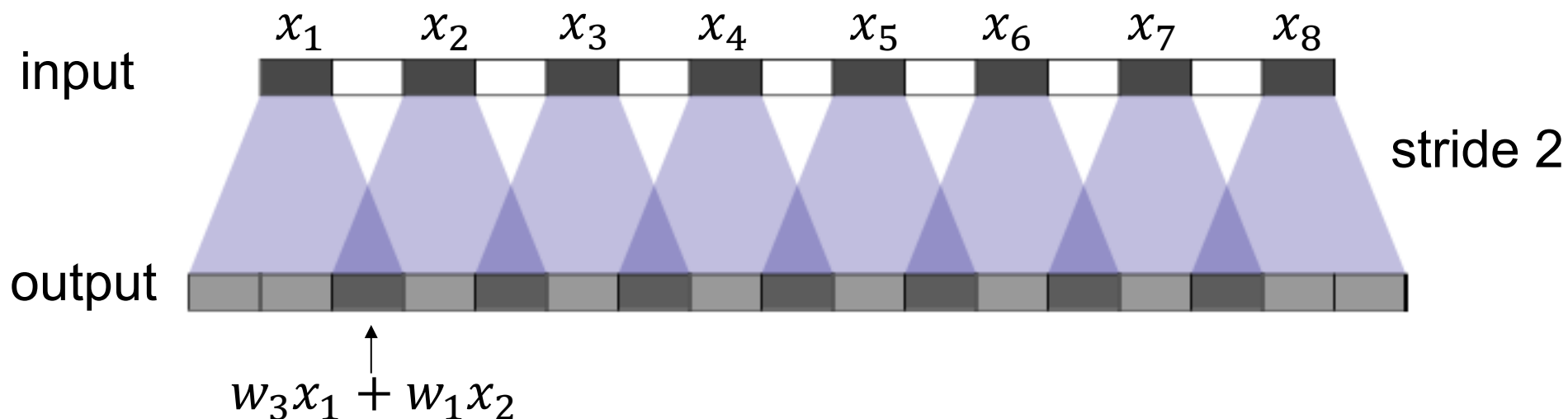
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

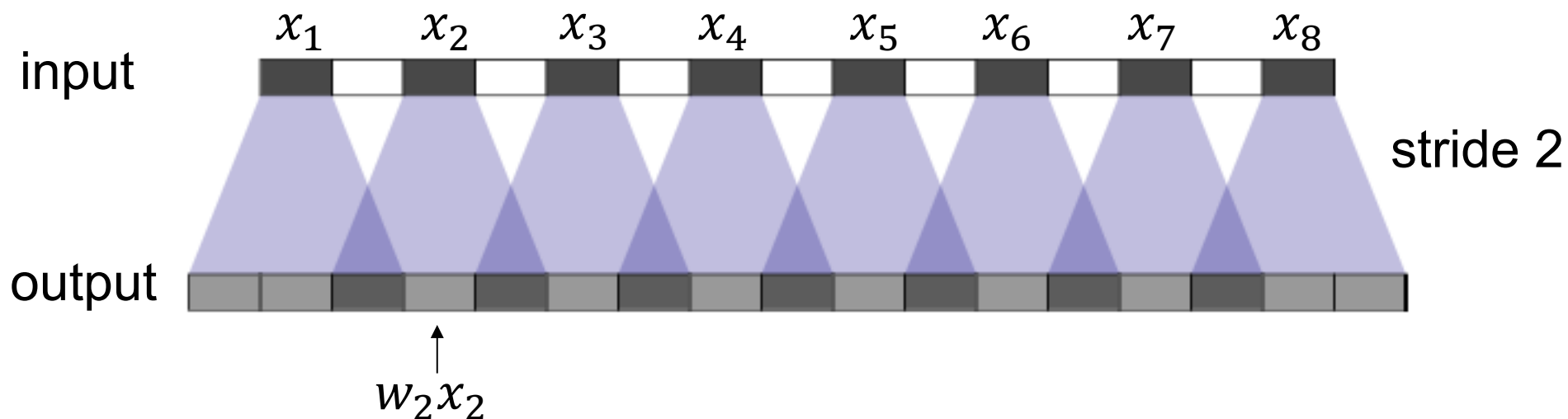
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

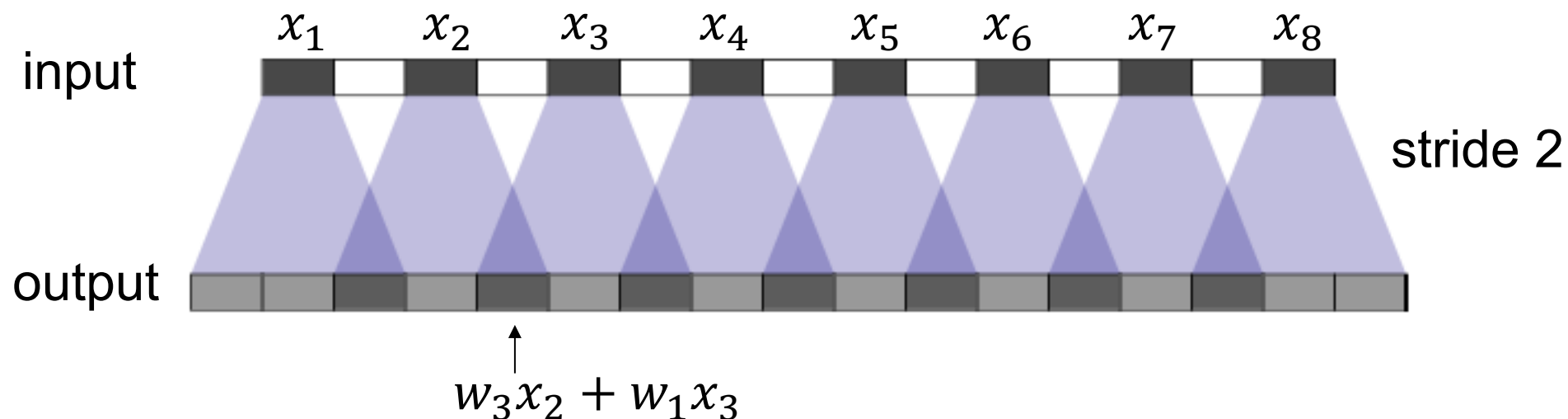
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

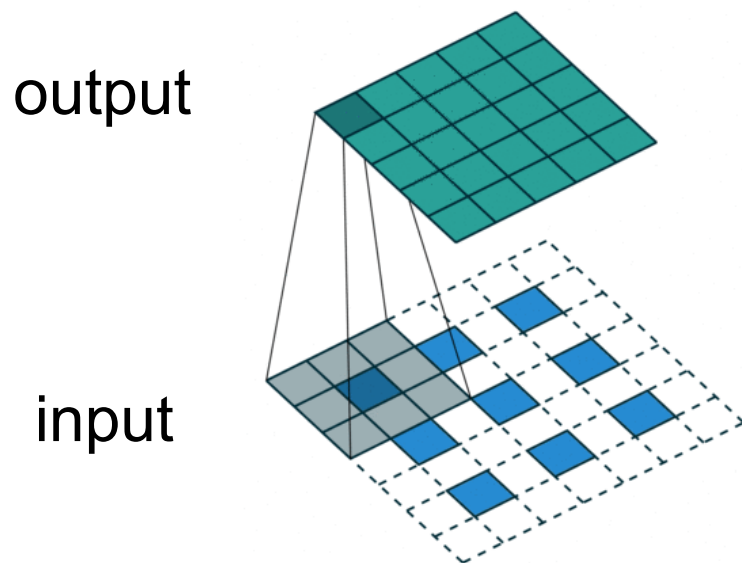
- *Backwards-strided convolution*: to increase resolution, use *output stride* > 1



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

Upsampling by transposed convolution

- *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, convolve with flipped filter
 - Sometimes called convolution with *fractional input stride* $1/2$

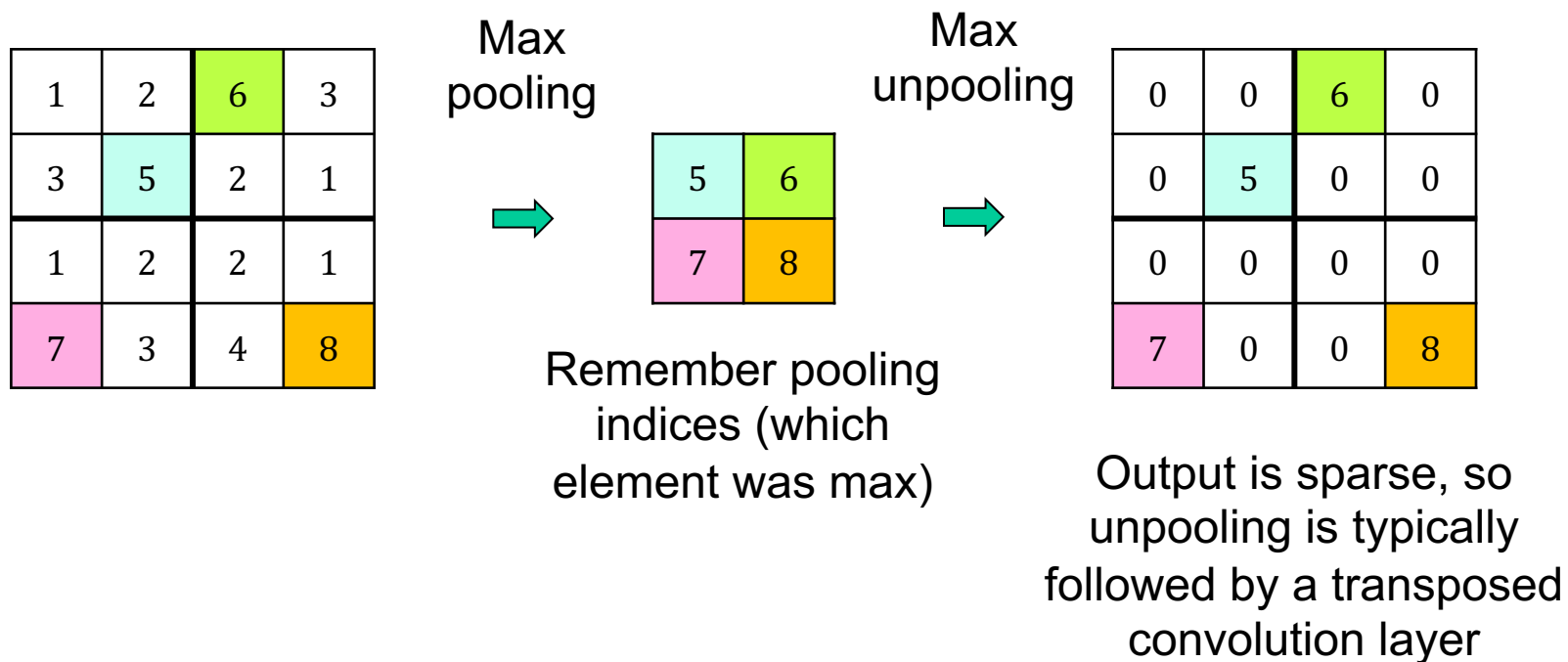


Q: What 3x3 filter would correspond to bilinear upsampling?

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

Upsampling by unpooling

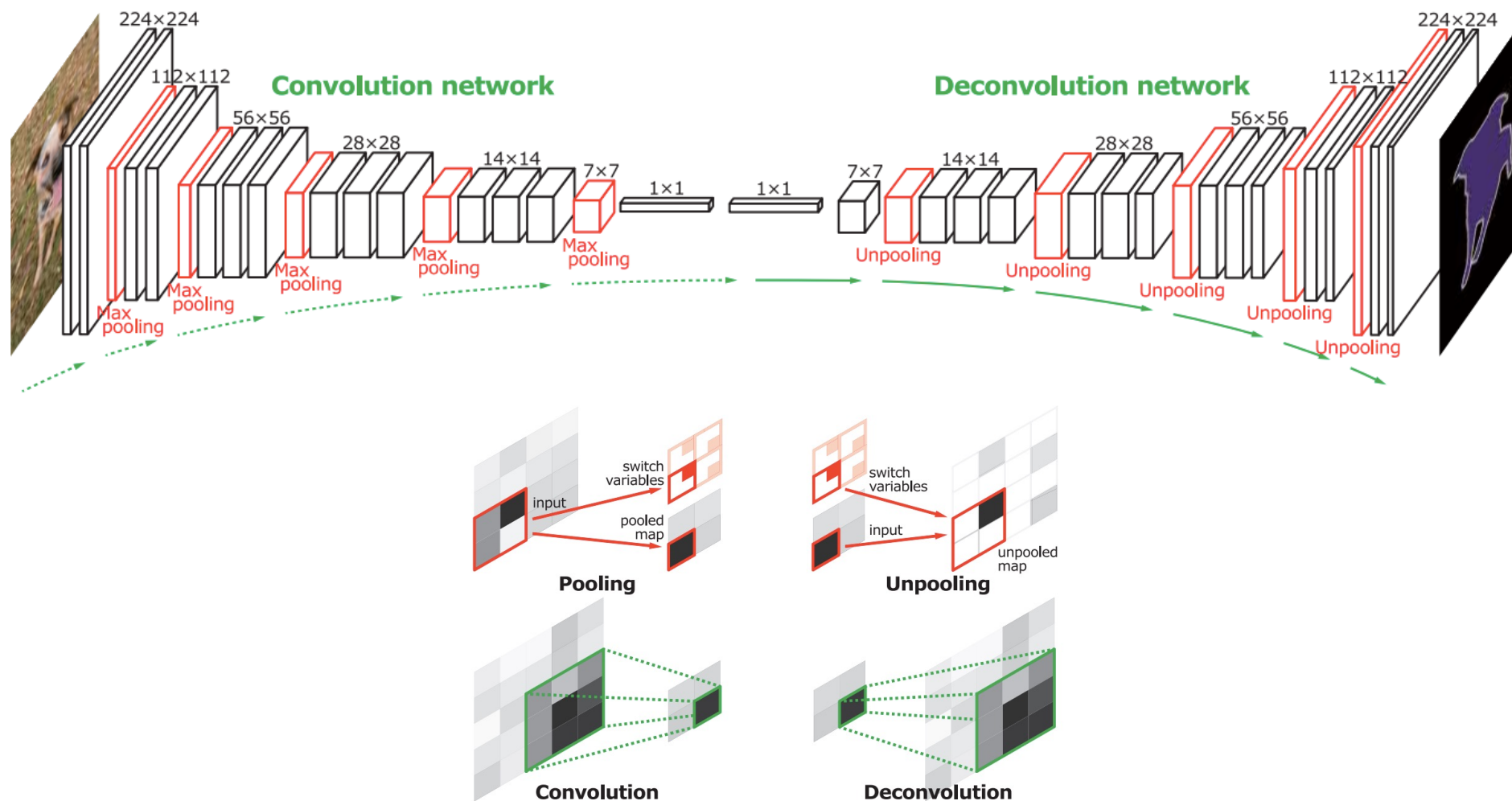
- Alternative to transposed convolution: max unpooling



Dense prediction: Outline

- Fully convolutional networks
- Operations for dense prediction
 - Transposed convolutions, unpooling
- Architectures for dense prediction
 - **DeconvNet, SegNet, U-Net**
- Instance segmentation
 - Mask R-CNN
- Other dense prediction problems

DeconvNet

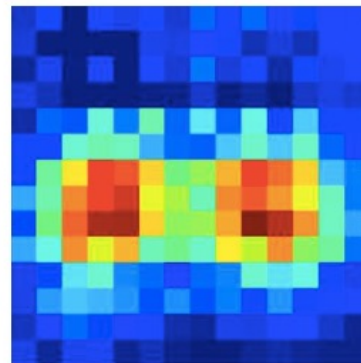


H. Noh, S. Hong, and B. Han, [Learning Deconvolution Network for Semantic Segmentation](#), ICCV 2015

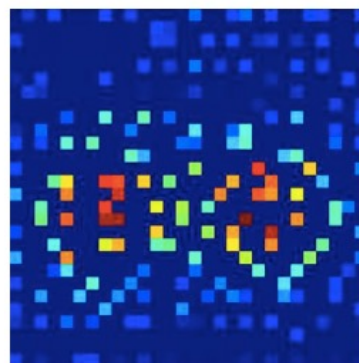
DeconvNet



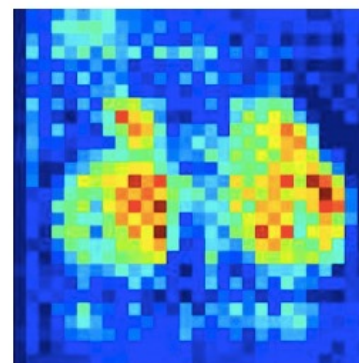
Original image



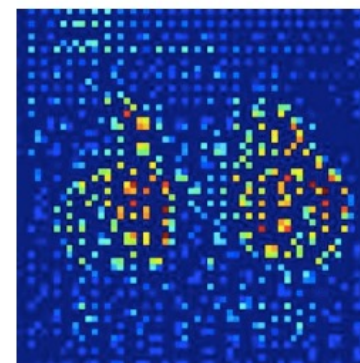
14x14 deconv



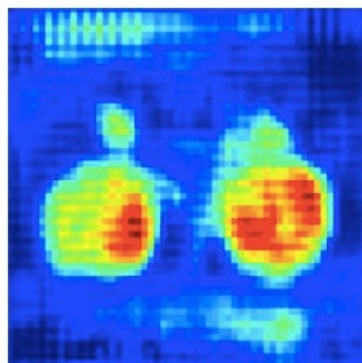
28x28 unpooling



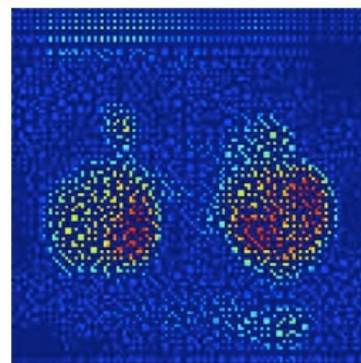
28x28 deconv



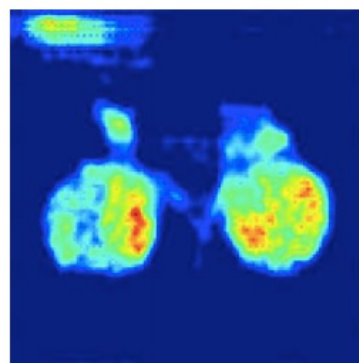
54x54 unpooling



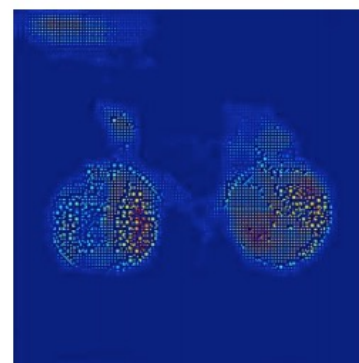
54x54 deconv



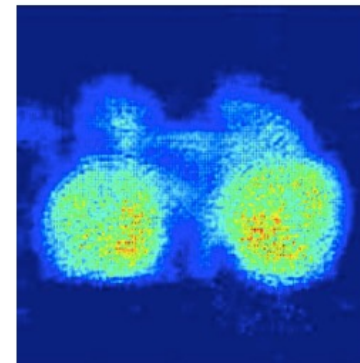
112x112 unpooling



112x112 deconv



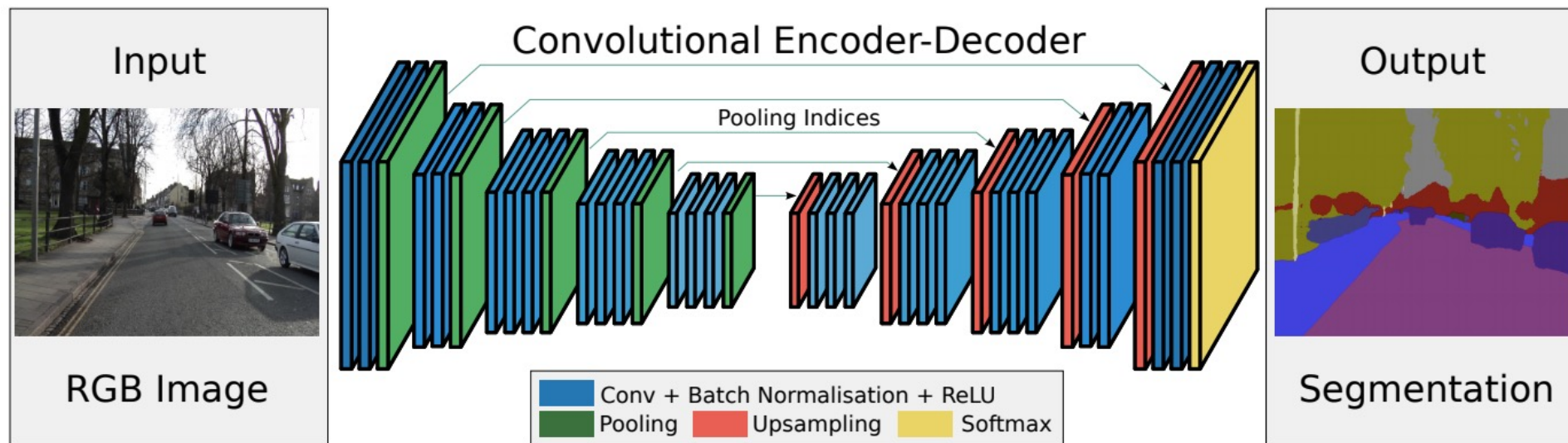
224x224 unpooling



224x224 deconv

H. Noh, S. Hong, and B. Han, [Learning Deconvolution Network for Semantic Segmentation](#), ICCV 2015

Similar architecture: SegNet

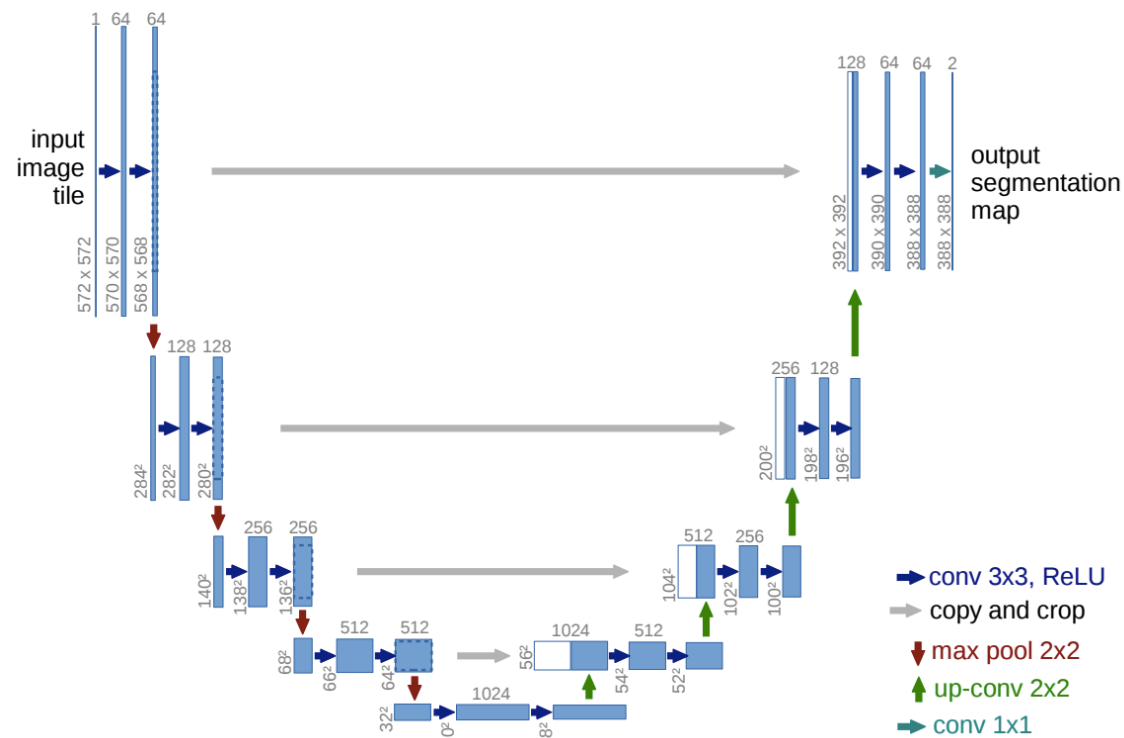


Drop the FC layers,
get better results

V. Badrinarayanan, A. Kendall and R. Cipolla, [SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation](#), PAMI 2017

U-Net

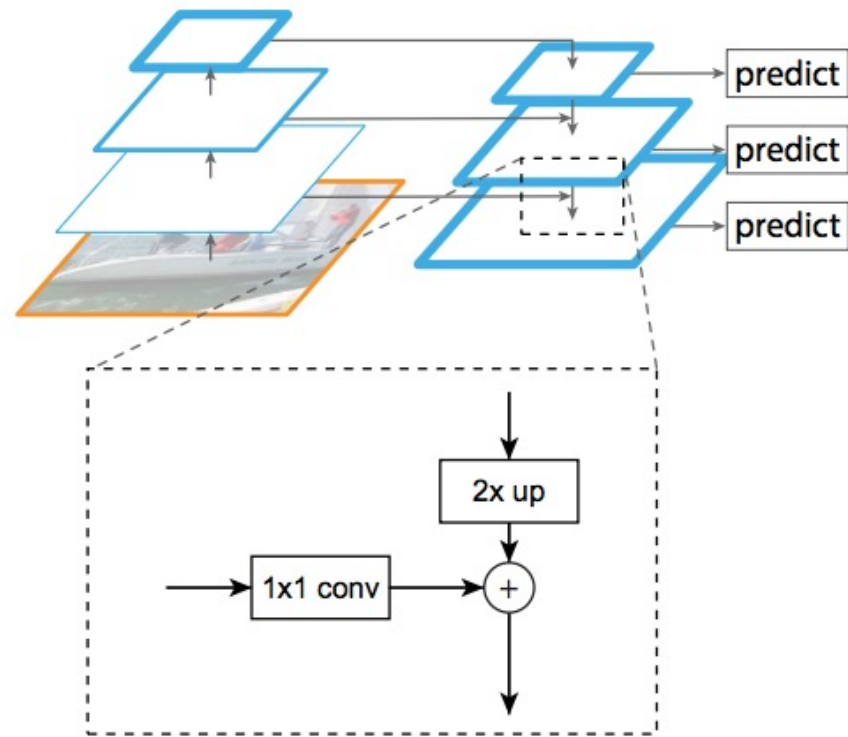
- Like FCN, fuse upsampled higher-level feature maps with higher-res, lower-level feature maps
- Unlike FCN, fuse by concatenation, predict at the end



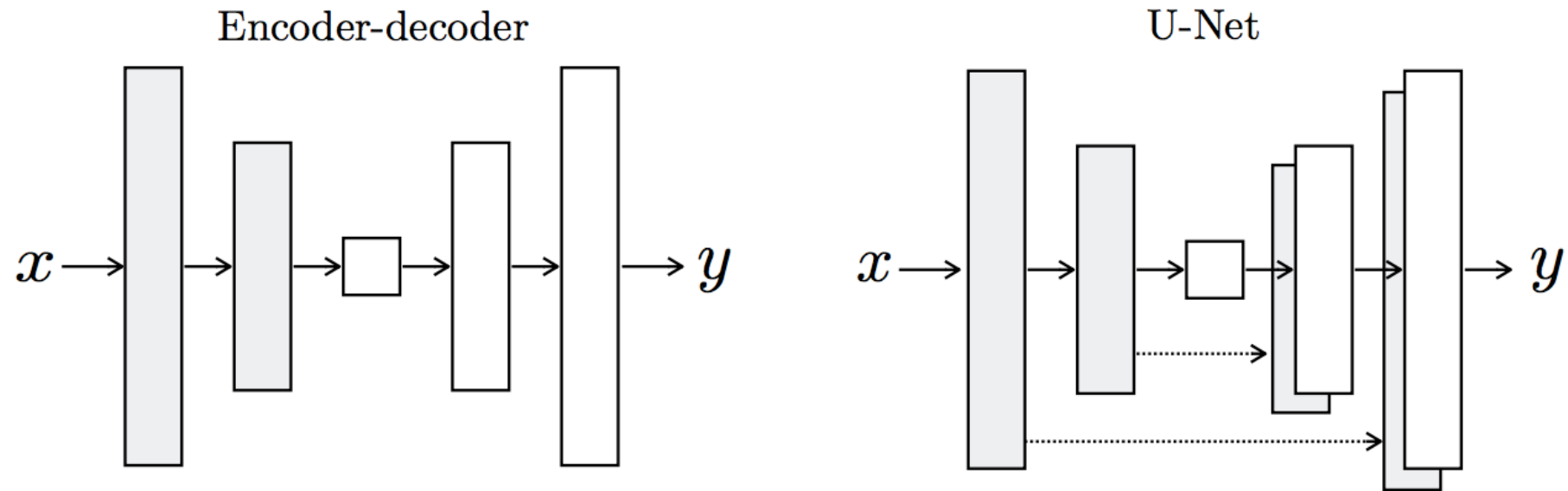
O. Ronneberger, P. Fischer, T. Brox, [U-Net: Convolutional Networks for Biomedical Image Segmentation](#), MICCAI 2015

Recall: Feature pyramid networks

- Improve predictive power of lower-level feature maps by adding contextual information from higher-level feature maps
- Predict different sizes of bounding boxes from different levels of the pyramid (but share parameters of predictors)



Summary of dense prediction architectures



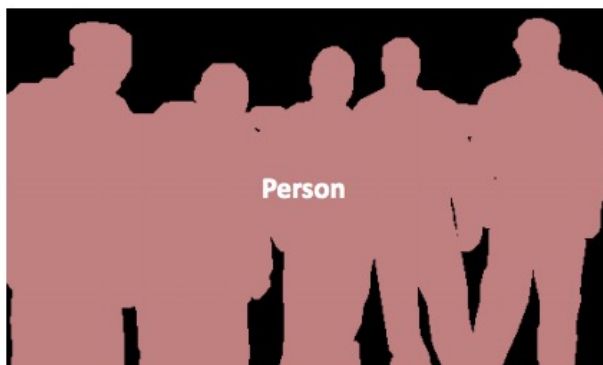
Outline

- Fully convolutional networks
- Operations for dense prediction
 - Transposed convolutions, unpooling
- Architectures for dense prediction
 - DeconvNet, SegNet, U-Net
- **Instance segmentation**
 - Mask R-CNN

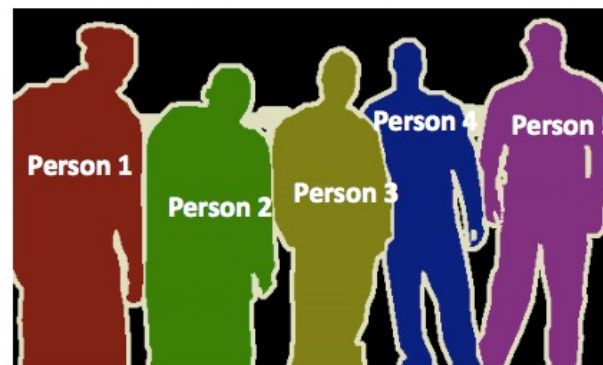
Instance segmentation



Object Detection



Semantic Segmentation

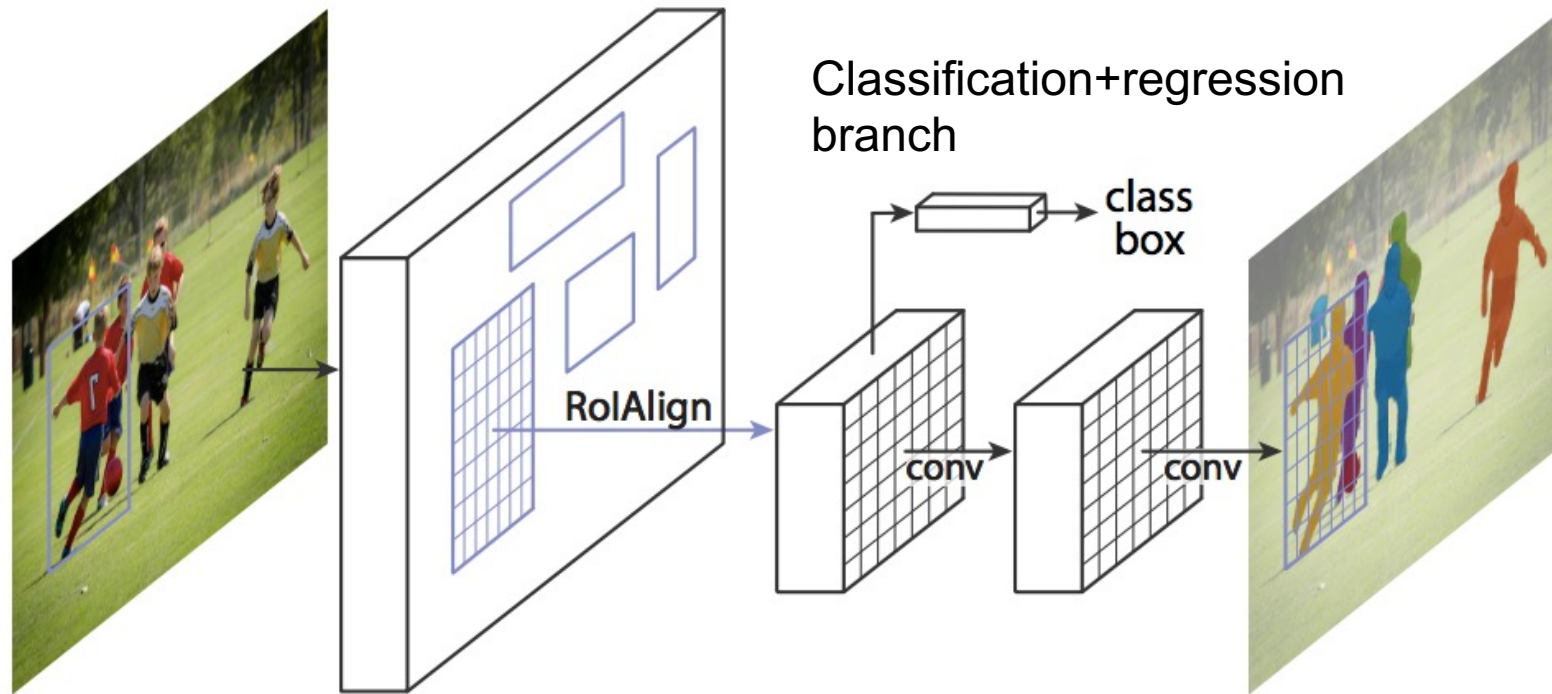


Instance Segmentation



Mask R-CNN

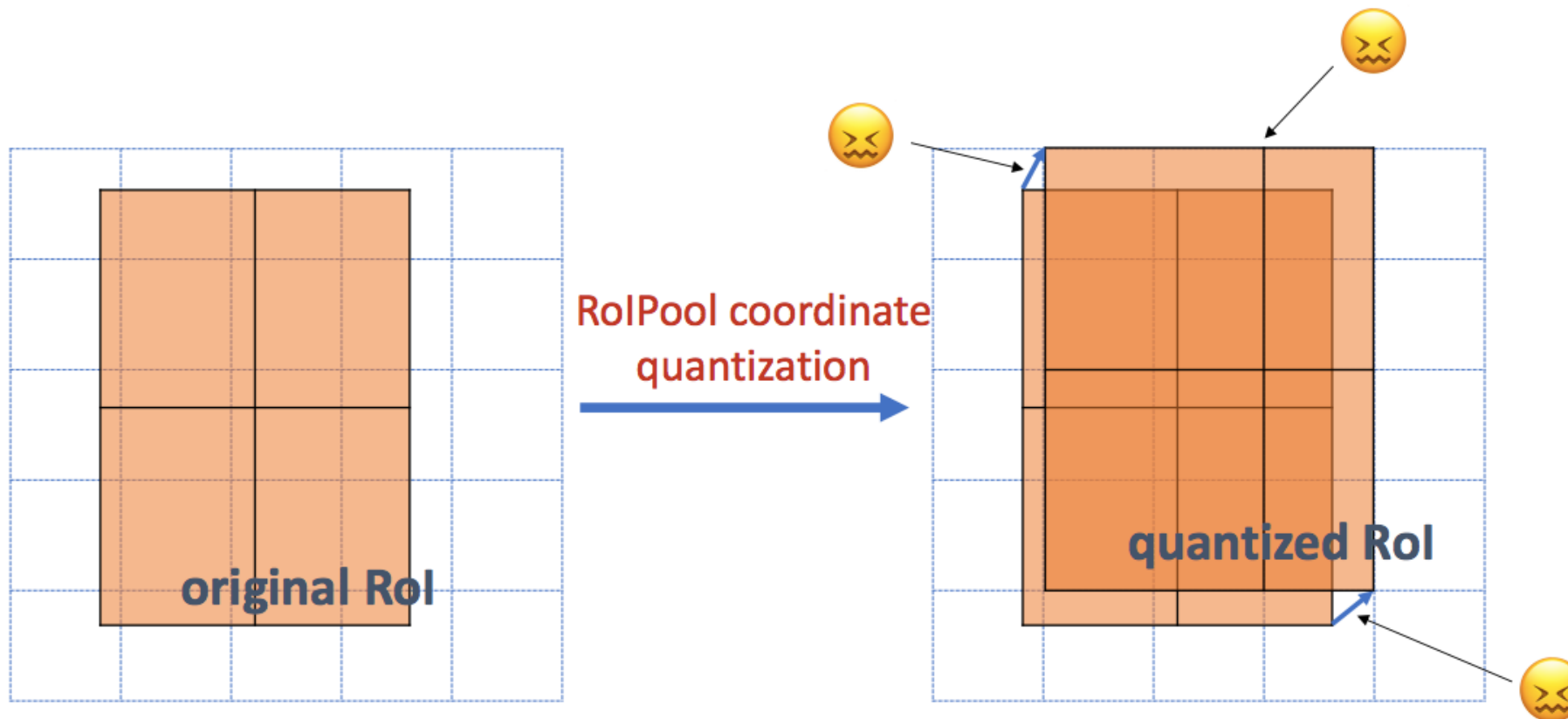
- Mask R-CNN = Faster R-CNN + FCN on Rols



Mask branch: separately predict segmentation for each possible class

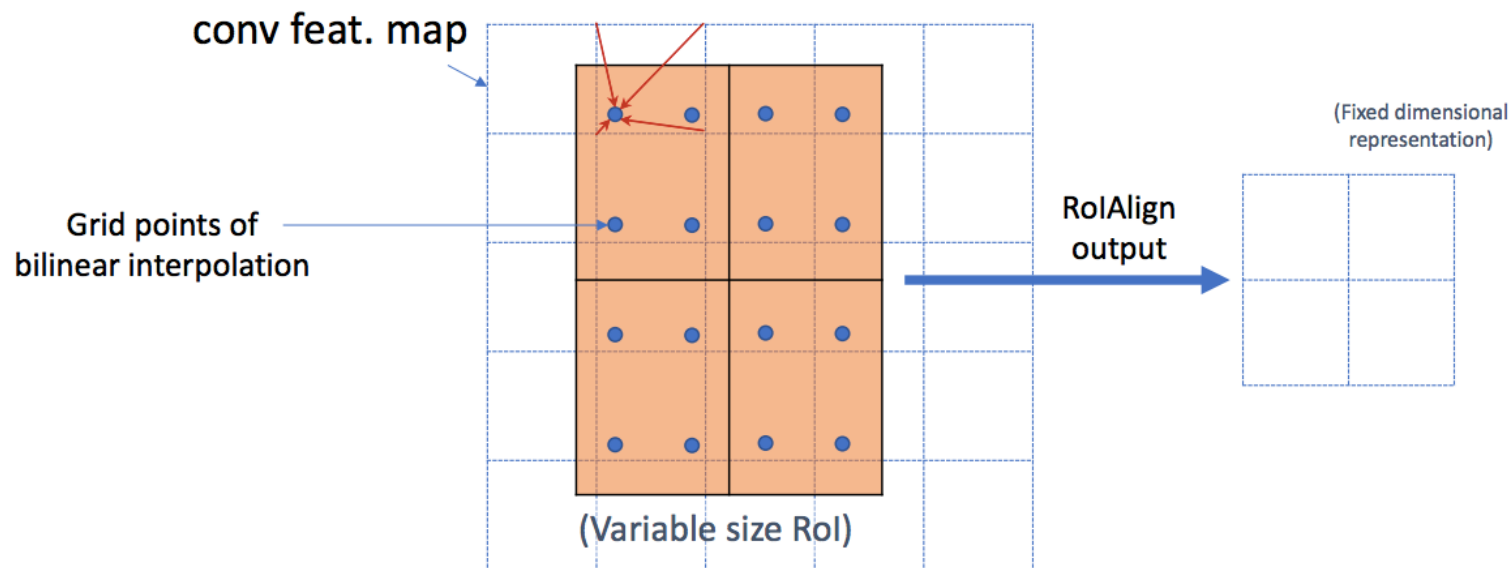
RoIAlign vs. RoIPool

- RoIPool: nearest neighbor quantization



RoIAlign vs. RoIPool

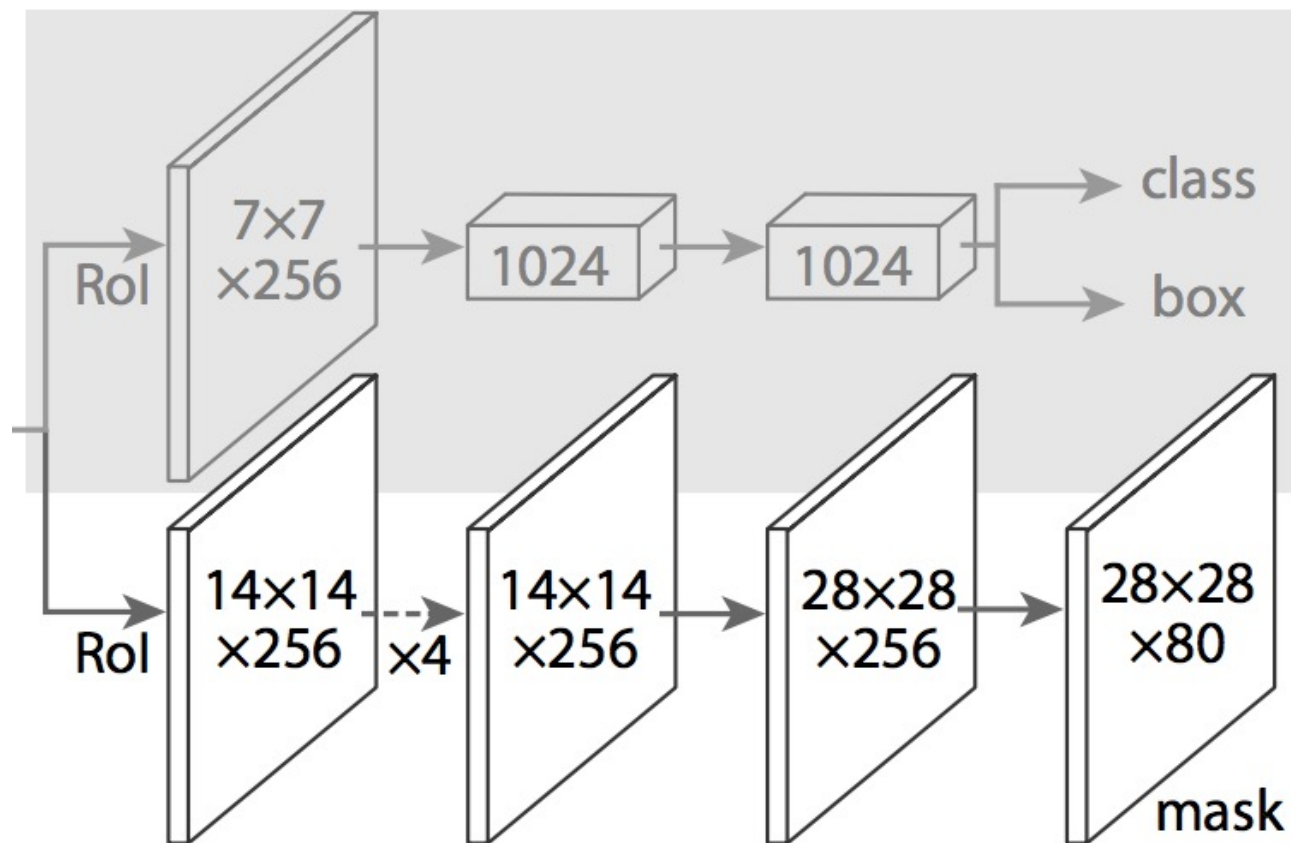
- RoIPool: nearest neighbor quantization
- RoIAlign: bilinear interpolation



K. He, G. Gkioxari, P. Dollar, and R. Girshick, [Mask R-CNN](#),
ICCV 2017 (Best Paper Award)

Mask R-CNN

- From RoIAlign features, predict class label, bounding box, and segmentation mask



Classification/regression head from an established object detector (e.g., FPN)

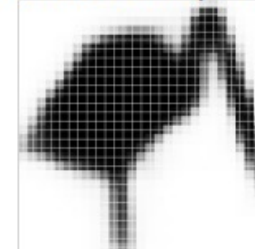
Separately predict binary mask for each class with per-pixel sigmoids, use average binary cross-entropy loss

Mask R-CNN



Validation image with box detection shown in red

28x28 soft prediction



Resized Soft prediction



Final mask

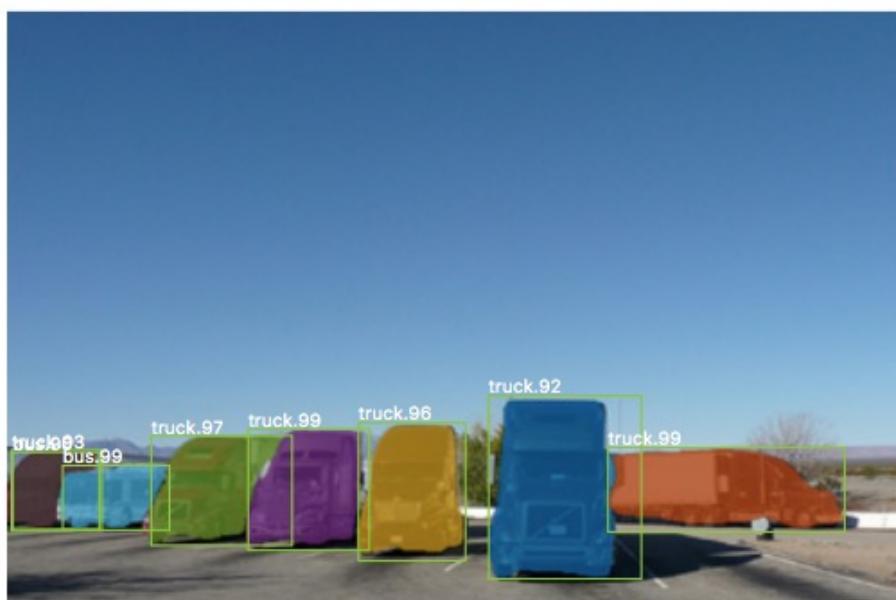
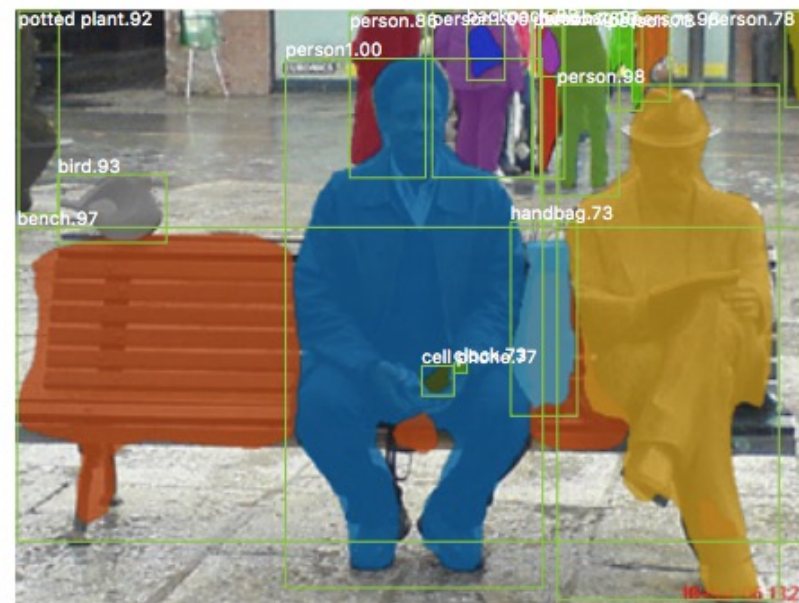
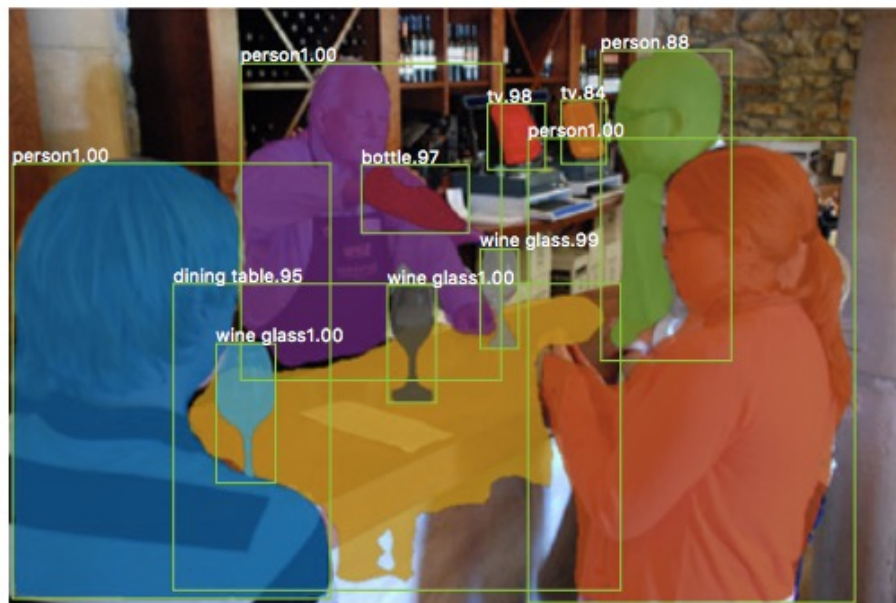


K. He, G. Gkioxari, P. Dollar, and R. Girshick, [Mask R-CNN](#),
ICCV 2017 (Best Paper Award)

Example results



Example results



Instance segmentation results on COCO

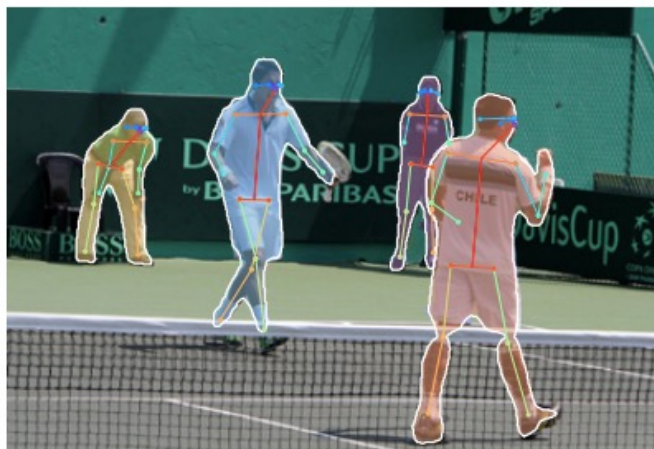
	backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
MNC [10]	ResNet-101-C4	24.6	44.3	24.8	4.7	25.9	43.6
FCIS [26] +OHEM	ResNet-101-C5-dilated	29.2	49.5	-	7.1	31.3	50.0
FCIS+++ [26] +OHEM	ResNet-101-C5-dilated	33.6	54.5	-	-	-	-
Mask R-CNN	ResNet-101-C4	33.1	54.9	34.8	12.1	35.6	51.1
Mask R-CNN	ResNet-101-FPN	35.7	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

Keypoint prediction

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



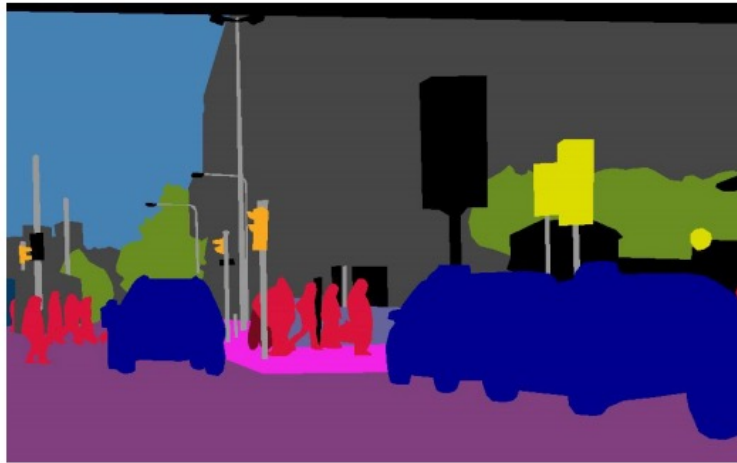
Outline

- Fully convolutional networks
- Operations for dense prediction
 - Transposed convolutions, unpooling
- Architectures for dense prediction
 - DeconvNet, SegNet, U-Net
- Instance segmentation
 - Mask R-CNN
- **Other dense prediction problems**

Recently proposed task: Panoptic segmentation



(a) image



(b) semantic segmentation

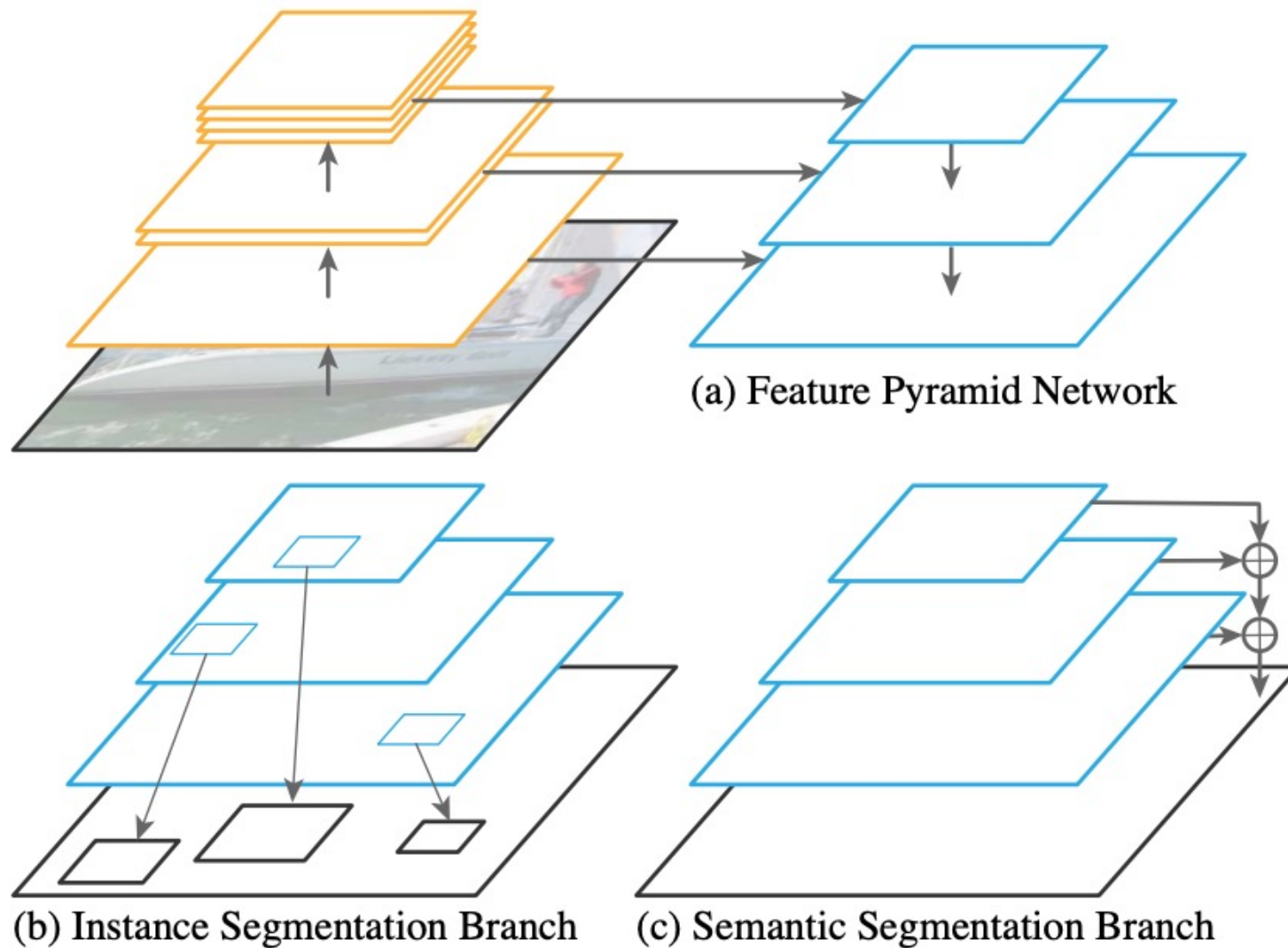


(c) instance segmentation



(d) panoptic segmentation

Panoptic feature pyramid networks



Panoptic feature pyramid networks

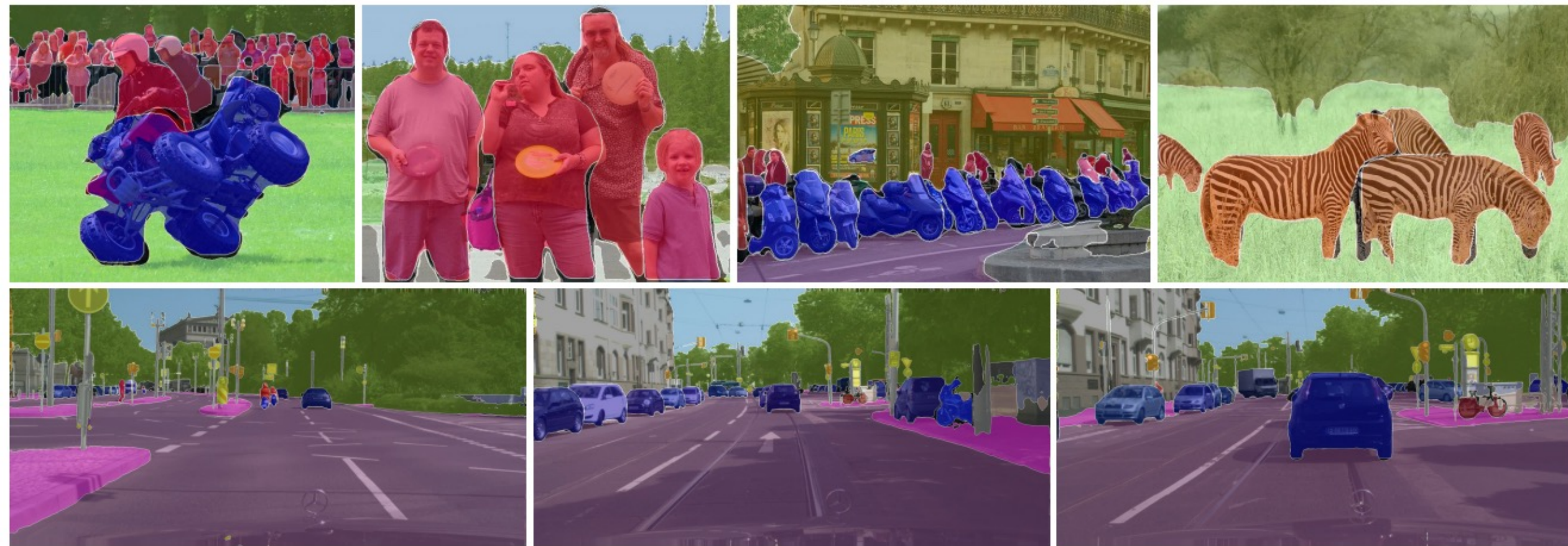
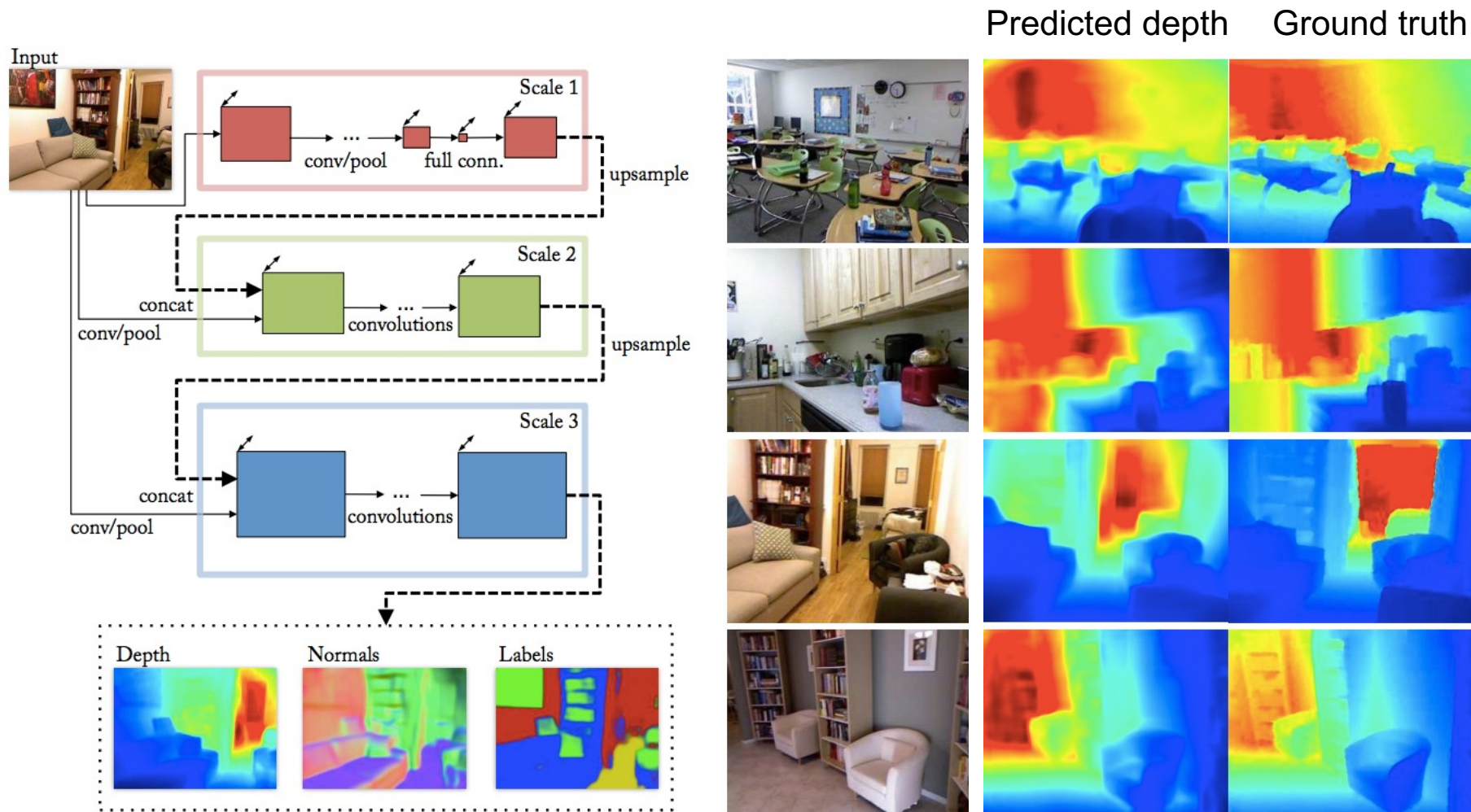


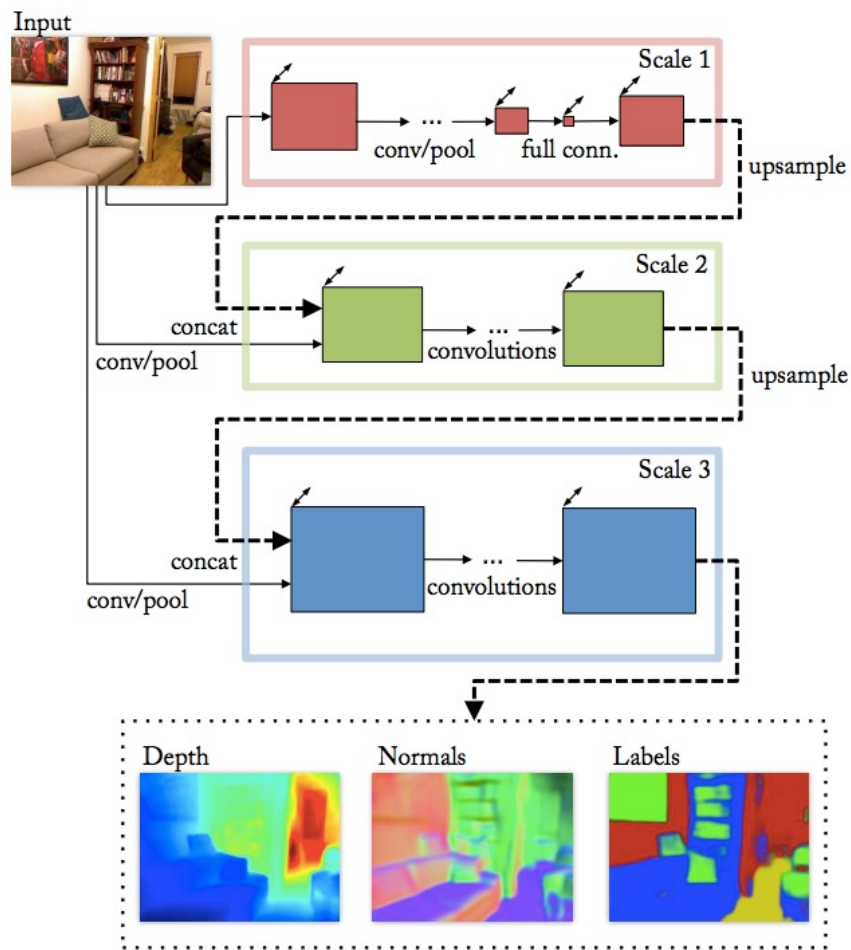
Figure 2: Panoptic FPN results on COCO (top) and Cityscapes (bottom) using a single ResNet-101-FPN network.

Depth and normal estimation

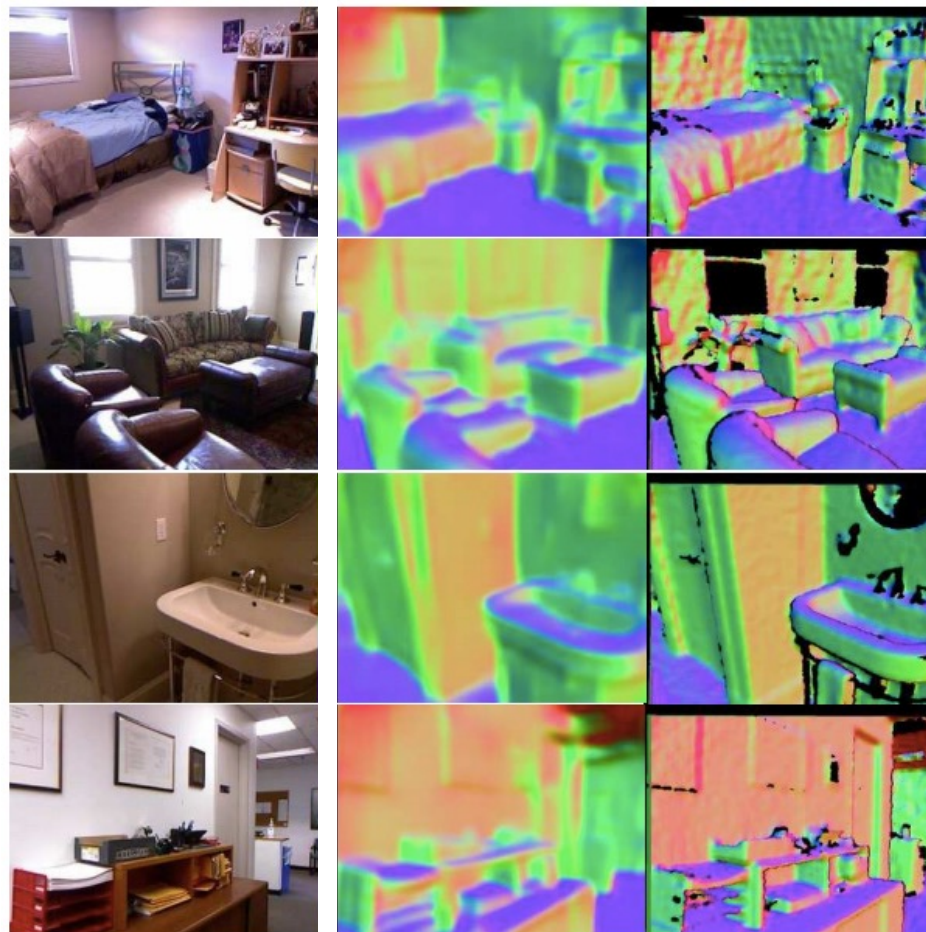


D. Eigen and R. Fergus, [Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture](#), ICCV 2015

Depth and normal estimation



Predicted normals Ground truth



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

Colorization

