# Segmentation

Computer Vision – Lecture 11

## **Further Reading**

- Slides from <u>S Lazebnik</u>
- Slides from <u>Johnson</u>
- Slides from <u>A Geiger</u>



### Image classification



Image-scale

### Object detection



Localization (Regions) Spatial Labelling

Segmentation



Pixel-level

## Semantic Segmentation

- Label each pixel in the image with a category label
- Do not differentiate instances, only care about pixels



4 Slide: J Johnson

## Evaluation

- Per-class IoU.
- Since there are no instances, evaluation is much easier than detection.
- Compare the binary masks of prediction and GT for each class.
- Average over classes.
- Note: IoU is favourable to large objects. (harder to intersect small objects)

# Sliding Window

Similar to detection:

- Slide a window over the whole image.
- Classify the centre pixel of the window.



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

# Sliding Window Segmentation

- Very inefficient!
- Very noisy: independent decision for each patch.
- Observation: we are recomputing the same features for overlapping patches.
- How can we share computation?
- Convolutions!

## Fully Convolutional Networks

If we do not down-sample (and always pad appropriately) we can design a network with convolutions that has the same input and output size.





## Fully Convolutional Networks

If we do not down-sample (and always pad appropriately) we can design a network with convolutions that has the same input and output size.



<sup>9</sup> Source: <u>Stanford CS231n</u>

## FCN: Fully Convolutional Networks

If we do not down-sample (and always pad appropriately) we can design a network with convolutions that has the same input and output size.

Very costly at full resolution. Down-sample, then up-sample!





J. Long, E. Shelhamer, and T. Darrell, <u>Fully Convolutional Networks for Semantic Segmentation</u>, CVPR 2015

# Upsampling: Unpooling

- "Inverse" max/avg-pooling operation.
- Pooling is not an invertible function: information is lost and cannot be recovered.
- Several options to approximate it.

## Unpooling: Nearest Neighbour



Input C x H x W

Output C x 2H x 2W

## Unpooling: Bed of Nails



Input C x H x W

Output C x 2H x 2W

## **Unpooling: Bilinear Interpolation**



Input C x H x W

Output C x 2H x 2W

## Max-Unpooling

While down-sampling: remember locations









While up-sampling: use remembered locations



15 Slide: J Johnson

## Convolution with Stride

3x3 Convolution with stride 2





## **Transposed Convolution**

3x3 convolution transpose, stride 2

- 9 weights
- Scale by input value
- Sum in overlapping regions



## **Transposed Convolution**



<i>a</i> <sub>1</sub> <i>w</i> <sub>1</sub>	$a_1w_2$	$\begin{array}{c}a_1w_3+\\a_2w_1\end{array}$	
$a_1w_4$	<i>a</i> <sub>1</sub> <i>w</i> <sub>5</sub>	$\begin{array}{c}a_1w_6+\\a_2w_4\end{array}$	
$a_1w_7 + a_3w_1$	$a_1w_8 + a_3w_2$	$a_1w_9 + a_2w_7 + a_3w_3 + a_4w_1$	

## Transposed Convolution

- Learnable upsampling.
- Outputs are the sum of 1-4 values: often grid pattern in output. Add normal convolution to learn smoothing.

Other names (can be confusing)

- Deconvolution
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

## U-Net

- Skip-connections.
- Concatenate upsampled higher-level feature maps with higher-res, lower-level feature maps.
- Low-level feature maps: details.
- High-level feature maps: semantics.



### **Dense Prediction Architectures**





## **Transformer Architectures**



## Interactive Segmentation

- User specifies what to segment.
- Input:
  - Seed points
  - Scribbles
  - Bounding box
  - Text
  - ...
- Model segments corresponding object(s).



## SAM: Segment Anything Model

Kirilliov et al., 2023

- Trained a model on various input modalities.
- Large scale supervision
  - 1B masks
  - 11M images



## Training with Humans in the Loop

- 1. Annotate data.
- 2. Train model.
- 3. Label more data with the model.
- 4. Humans fix, improve labels.
- 5. Goto 2.



## Large-Scale Annotations



## Interactive Segmentation

- User centric.
- SAM: no class labels, just binary segmentation.

Other types:

- Foreground/background segmentation
- Referring expressions segmentation ("The man with the blue hat")
- Saliency segmentation (what stands out in the image)

# Things and Stuff

**Thing**: An object with a specific size and shape.

**Stuff**: Material defined by a homogeneous or repetitive pattern of fine-scale properties, but has no specific or distinctive spatial extent or shape

- Object detection: things (instances)
- Semantic segmentation: things and stuff (but no instances)



#### Things:

- Dog
- Tree
- Lantern

• •••

#### Stuff:

- Sky
- Snow

• ...

## Instance Segmentation

- Semantic segmentation does not separate objects of the same class.
- Object detection finds individual objects (instances).
- Instance segmentation: segmentation at instance-level.



## Instance Segmentation

- Detect all objects in the image, and identify the pixels that belong to each object (only things, not stuff)
- Intuitive approach:
  - Detect objects
  - predict a segmentation mask for each object
- Practice: add another branch to your detector that predicts a mask for each box.

## Mask R-CNN

#### Faster R-CNN + FCN on Rols



Mask branch: separately predict segmentation for each possible class

## RolPool

Nearest neighbor quantization results in small errors



## RolAlign vs. RolPool

RoIPool: nearest neighbor quantization RoIAlign: bilinear interpolation



## Mask R-CNN

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



Classification head

Regression head

Separately predict binary mask (28x28) for each class with perpixel sigmoids, use average binary cross-entropy loss (80 classes)

## **Mask R-CNN Prediction**



#### 28x28 soft prediction



#### Resized soft prediction



Final mask



## Mask R-CNN Prediction





#### Resized Soft prediction



Final mask







## Panoptic Segmentation

- Combine semantic segmentation for stuff,
- with instance segmentation for things.
- Can be solved separately.
- More efficient: share some computation between tasks.



(a) Image



(b) Semantic Segmentation



(c) Instance Segmentation



(d) Panoptic Segmentation

## Keypoints

Instead of predicting masks, we can predict other things such as keypoints.

Keypoints here: object-specific landmarks, e.g. joints.



## Human Pose



## Keypoint Classification Loss

- Turn the GT location into a class.
- As many classes as pixels in the heatmap.
- Train with softmax cross-entropy loss.
- Output resolution limited to heatmap resolution (28x28 for Mask R-CNN)

## Differentiable Keypoint Regression

- We can use softmax to formulate the keypoint regression task as a heatmap prediction problem.
- Compute softmax over heatmap:  $H = \operatorname{softmax}(h)$

$$(p_x, p_y)^T = \sum_{u, v} {\binom{u}{v}} H(u, v)$$

- Output location is the weighted sum of pixel locations.
- Use temperature to make it sharper.

## Combining Tasks

Heads can be combined to solve multiple tasks simultaneously with minimal overhead.





## **Dense Captioning**

#### Add a text-output head: predict captions.



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016

## 3D Shape

Add a mesh prediction head.

