

Computer Vision

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

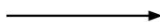
January 07, 2024

- ▶ Previously, we discussed Image Classification
- ▶ A core task in Computer Vision



This image by Nikita is
licensed under [CC-BY 2.0](#)

(assume given a set of possible labels)
{dog, cat, truck, plane, ...}



cat

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

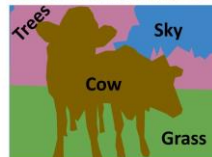
Instance Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

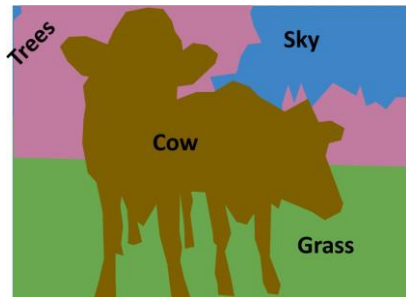
- ▶ **Things:** Object categories that can be separated into object instances (e.g. cats, cars, person)
- ▶ **Stuff:** Object categories that cannot be separated into instances (e.g. sky, grass, water, trees)



- ▶ **Object Detection:** Detects individual object instances, but only gives box(Only things!)



- ▶ **Semantic Segmentation:**
Gives per-pixel labels, but merges instances (Both things and stuff)





GRASS, CAT,
TREE, SKY, ...

Paired training data: for each training image, each pixel is labeled with a semantic category.



At test time, classify each pixel of a new image.

Full image

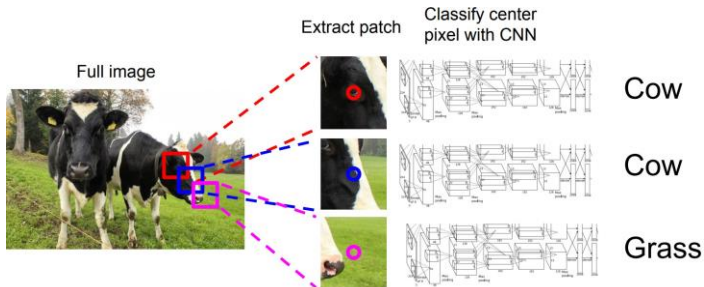


Full image



- ▶ Impossible to classify without context
- ▶ How do we include context?

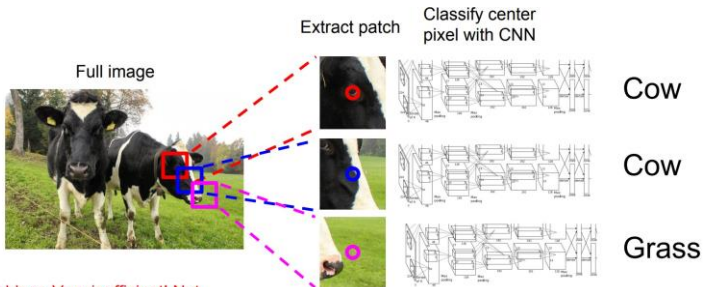
Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window (cont.)

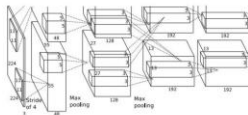


Problem: Very inefficient! Not reusing shared features between overlapping patches

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

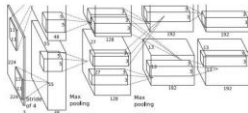
Semantic Segmentation Idea: Convolution

Full image



An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Full image

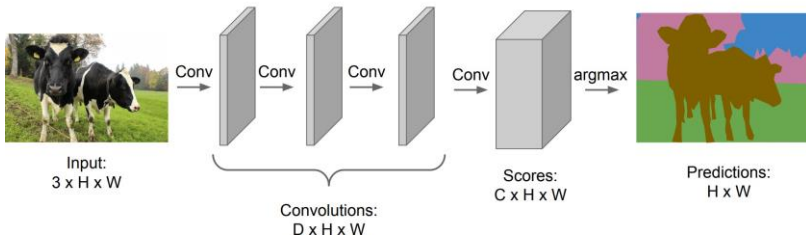


An intuitive idea: encode the entire image with conv net, and do semantic segmentation on top.

Problem: classification architectures often reduce feature spatial sizes to go deeper, but semantic segmentation requires the output size to be the same as input size.

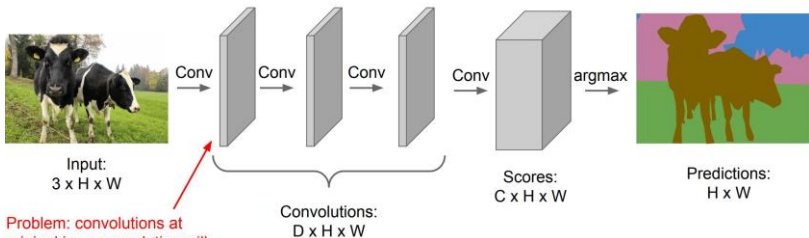
Semantic Segmentation Idea: Fully Convolutional

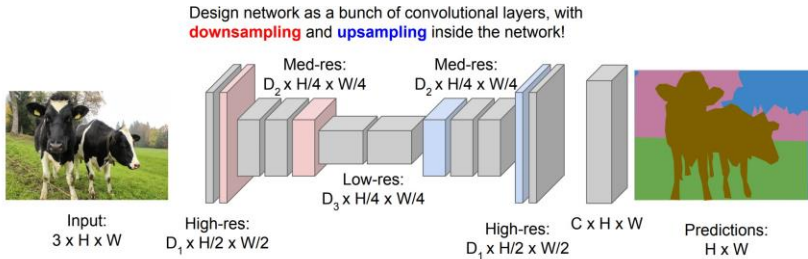
Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!



Semantic Segmentation Idea: Fully Convolutional (cont.)

Design a network with only convolutional layers without downsampling operators to make predictions for pixels all at once!





Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

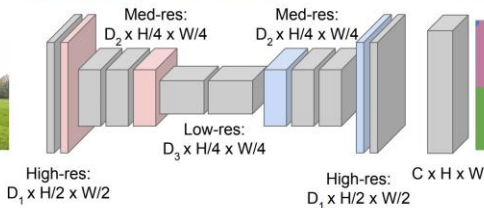
Semantic Segmentation Idea: Fully Convolutional (cont.)

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:
???

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network Upsampling: Unpooling

Nearest Neighbor

1	2
3	4

Input: 2 x 2



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output: 4 x 4

"Bed of Nails"

1	2
3	4

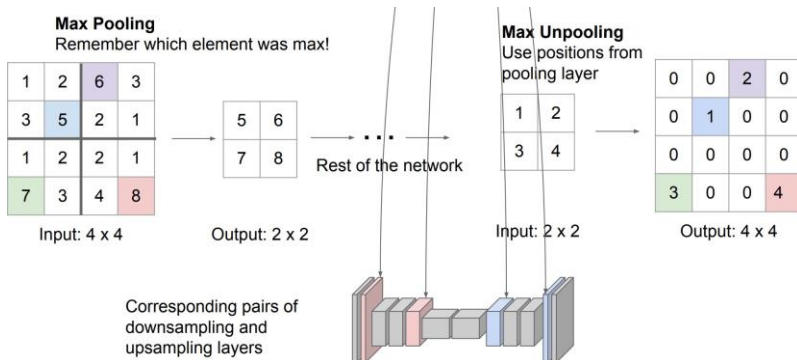
Input: 2 x 2



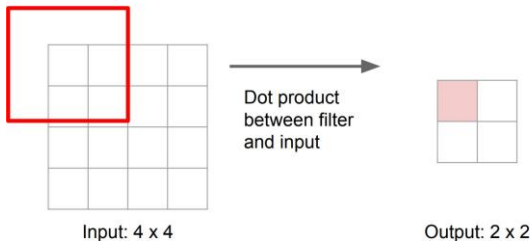
1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output: 4 x 4

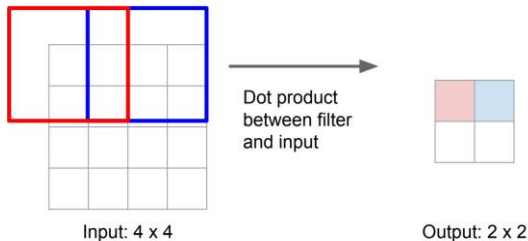
In-Network Upsampling: Max Unpooling



Recall: Normal 3 x 3 convolution, stride 2 pad 1



Recall: Normal 3 x 3 convolution, stride 2 pad 1



Filter moves 2 pixels in the input for every one pixel in the output

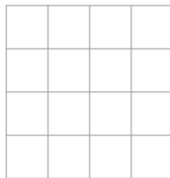
Stride gives ratio between movement in input and output

We can interpret strided convolution as "learnable downsampling".

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

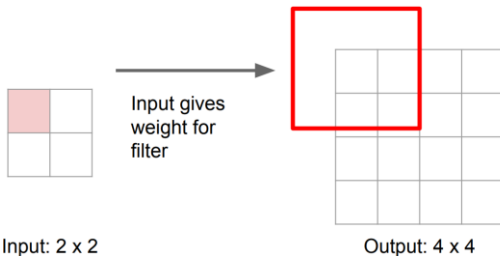


Input: 2 x 2

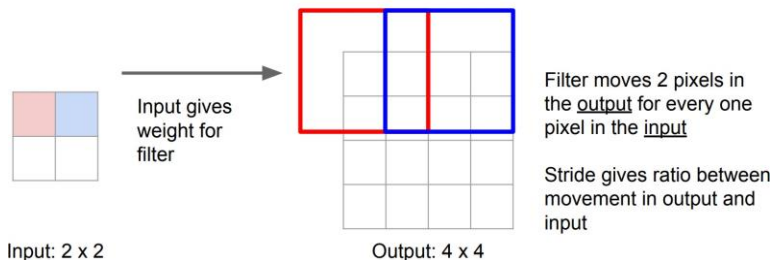


Output: 4 x 4

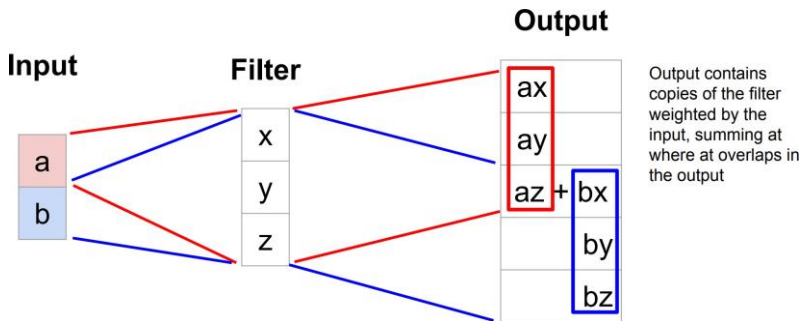
3 x 3 **transposed** convolution, stride 2 pad 1

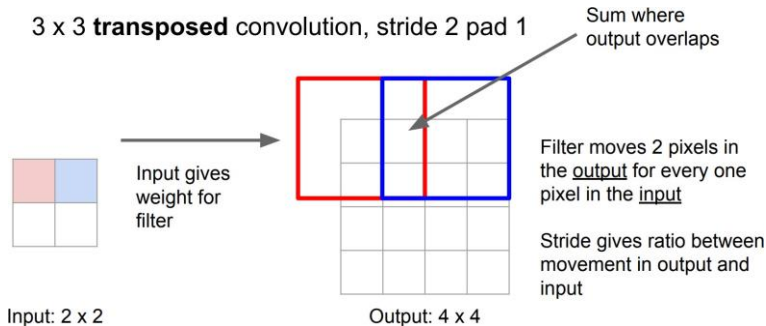


3 x 3 **transposed** convolution, stride 2 pad 1



Transposed Convolution: 1D Example





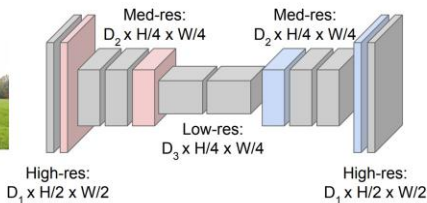
Semantic Segmentation Idea: Fully Convolutional

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!



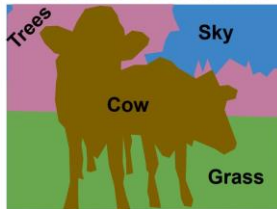
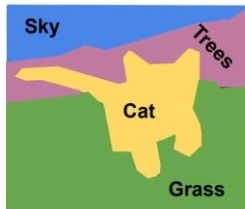
Upsampling:
Unpooling or strided
transposed convolution



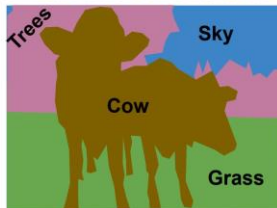
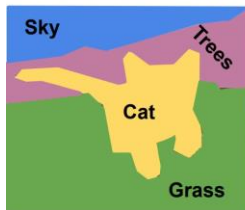
Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

- ▶ Label each pixel in the image with a category label

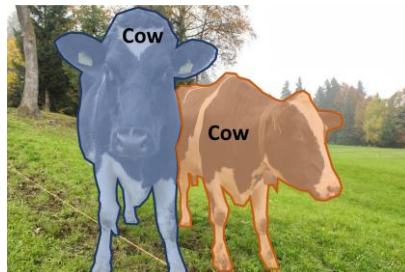


- ▶ Label each pixel in the image with a category label

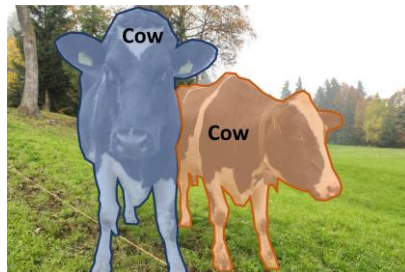


- ▶ Does not differentiate instances, only care about pixels

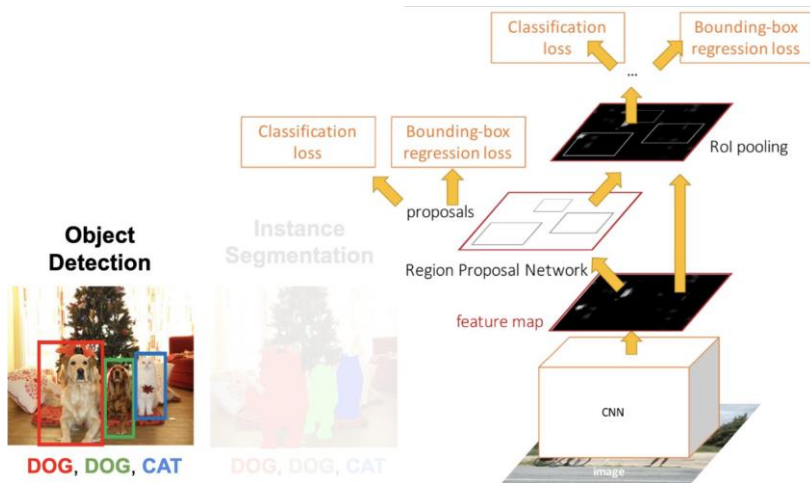
- Detect all objects in the image, and identify the pixels that belong to each object (Only things!)



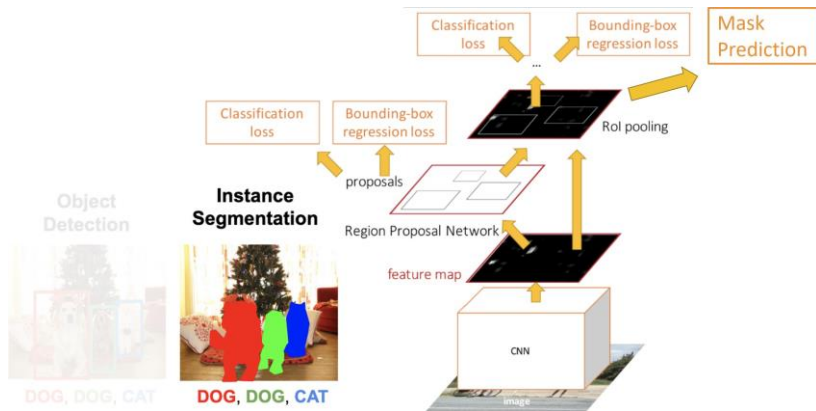
- ▶ Detect all objects in the image, and identify the pixels that belong to each object (Only things!)
- ▶ **Approach:** Perform object detection, then predict a segmentation mask for each object!



Object Detection: Faster R-CNN

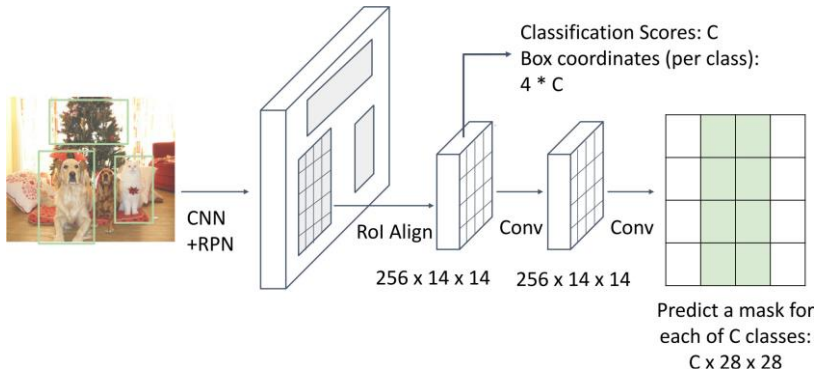


Instance Segmentation: Mask R-CNN

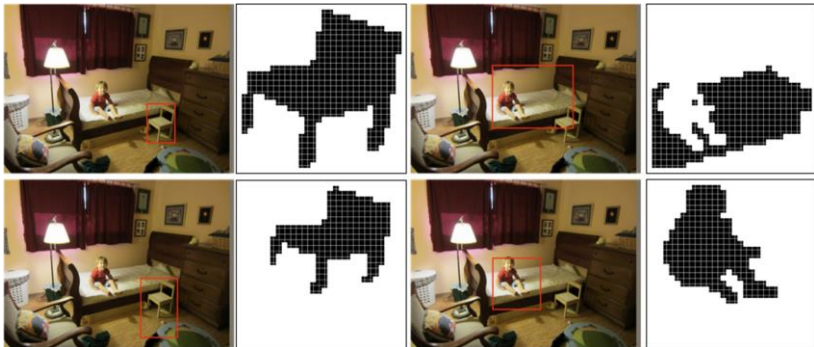


⁰He et al, "Mask R-CNN", ICCV 2017

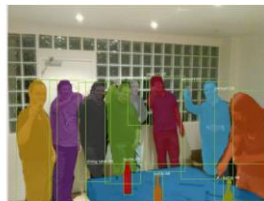
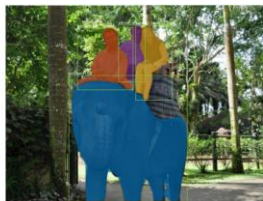
Instance Segmentation: Mask R-CNN



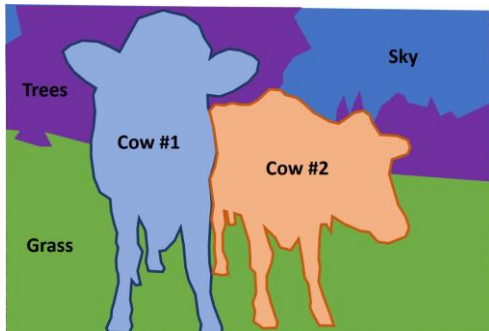
Mask R-CNN: Example Training Targets



Mask R-CNN: Very Good Results!



- ▶ Label all pixels in the image (both things and stuff)
- ▶ For “thing” categories also separate into instances



Beyond Instance Segmentation: Panoptic Segmentation

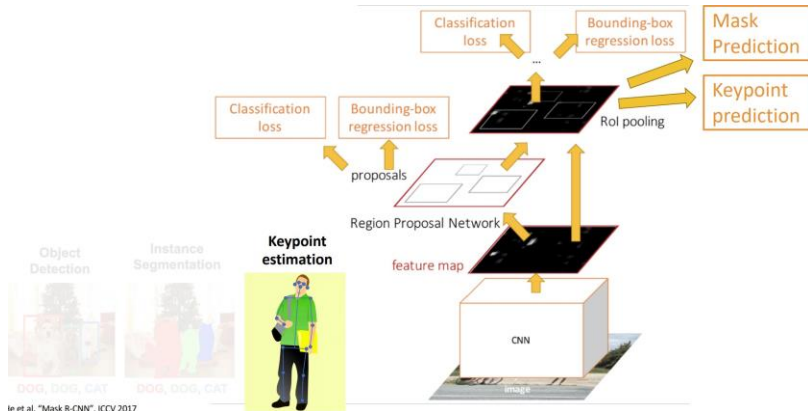


⁰Kirillov et al, "Panoptic Feature Pyramid Networks", [CVPR 2019](#)

- ▶ Represent the pose of a human by locating a set of keypoint se.g. 17 keypoints:
- ▶ Nose
- ▶ Left / Right eye
- ▶ Left / Right earLeft/ Right shoulder
- ▶ Left / Right elbow
- ▶ Left / Right wrist



Mask R-CNN: Keypoint Estimation



Mask R-CNN: Keypoint Estimation

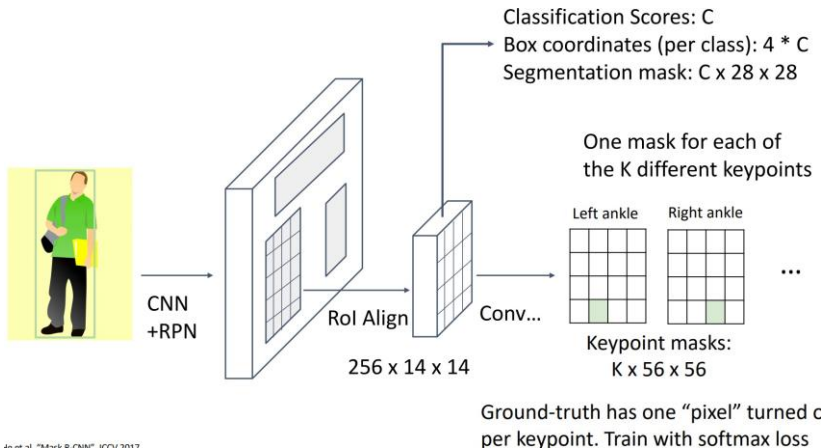


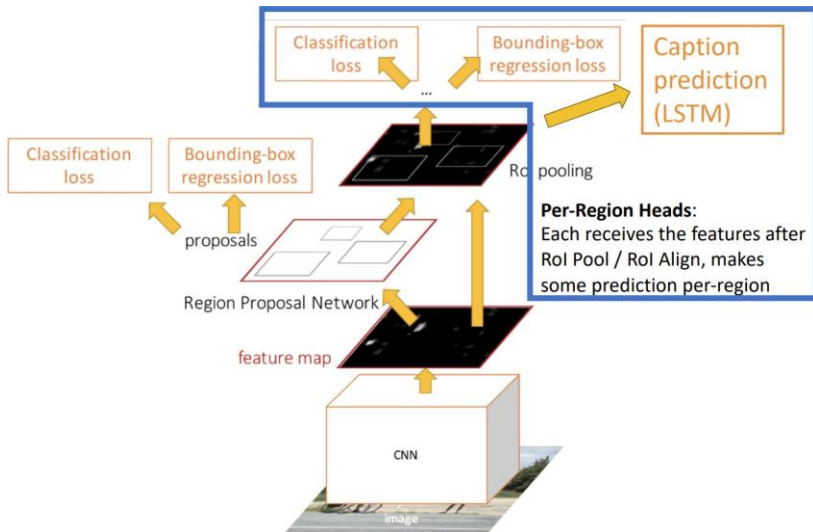
Image of "Person B-CNN" 10/10/2017

Joint Instance Segmentation and Pose Estimation

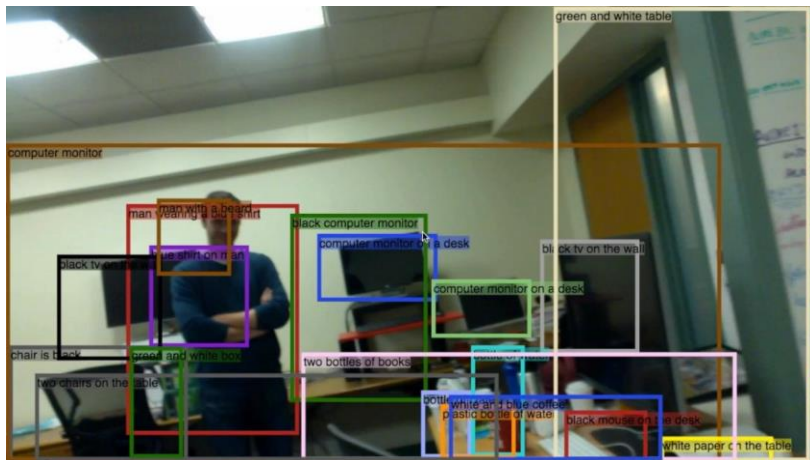


⁰He et al, "Mask R-CNN", ICCV 2017

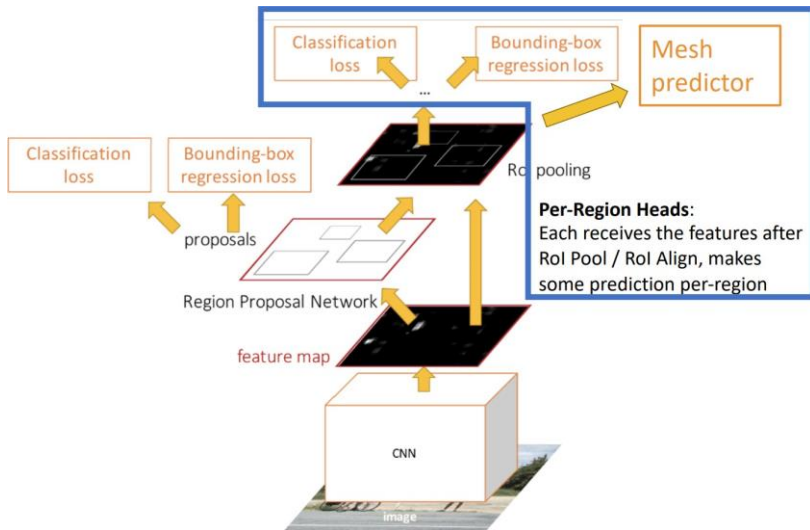
Captioning: Predict a caption per region!



Captioning: Predict a caption per region!



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



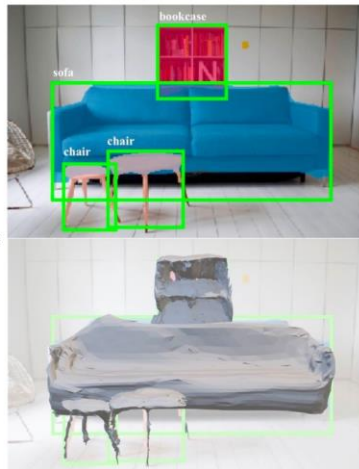
Mask R-CNN:

2D Image -> 2D shapes



Mesh R-CNN:

2D Image -> **3D** shapes



Gkioxari, Malik, and Johnson, "Mesh R-CNN", ICCV 2019

- ▶ **Goal:** Track objects over a sequence of photos or a video
- ▶ Exceedingly challenging in multi-object tracking scenarios
- ▶ Need to take care of not mixing up or losing objects midway
- ▶ **One Solution:** Perform object detection and assign IDs to each object and store its feature vector. Then track the objects based on its ID and feature vector

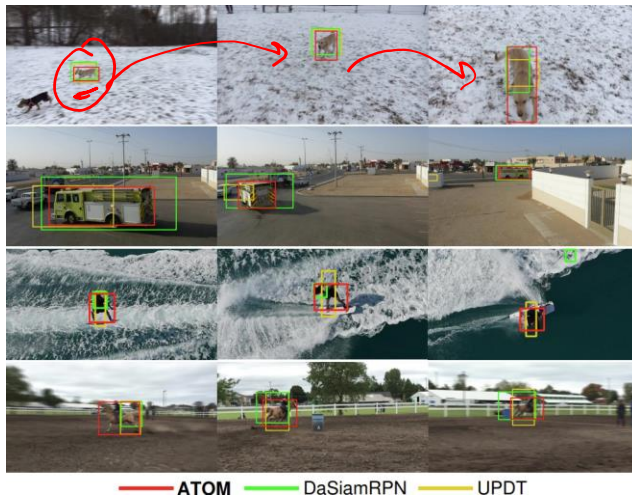


Figure 2: Comparison of 3 approaches for object tracking

These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAICDL4CV [Deep Learning for Computer Vision: Fundamentals and Applications](#)
- ▶ Justin Johnson, UMich EECS 498.008/598.008: [Deep Learning for Computer Vision](#)