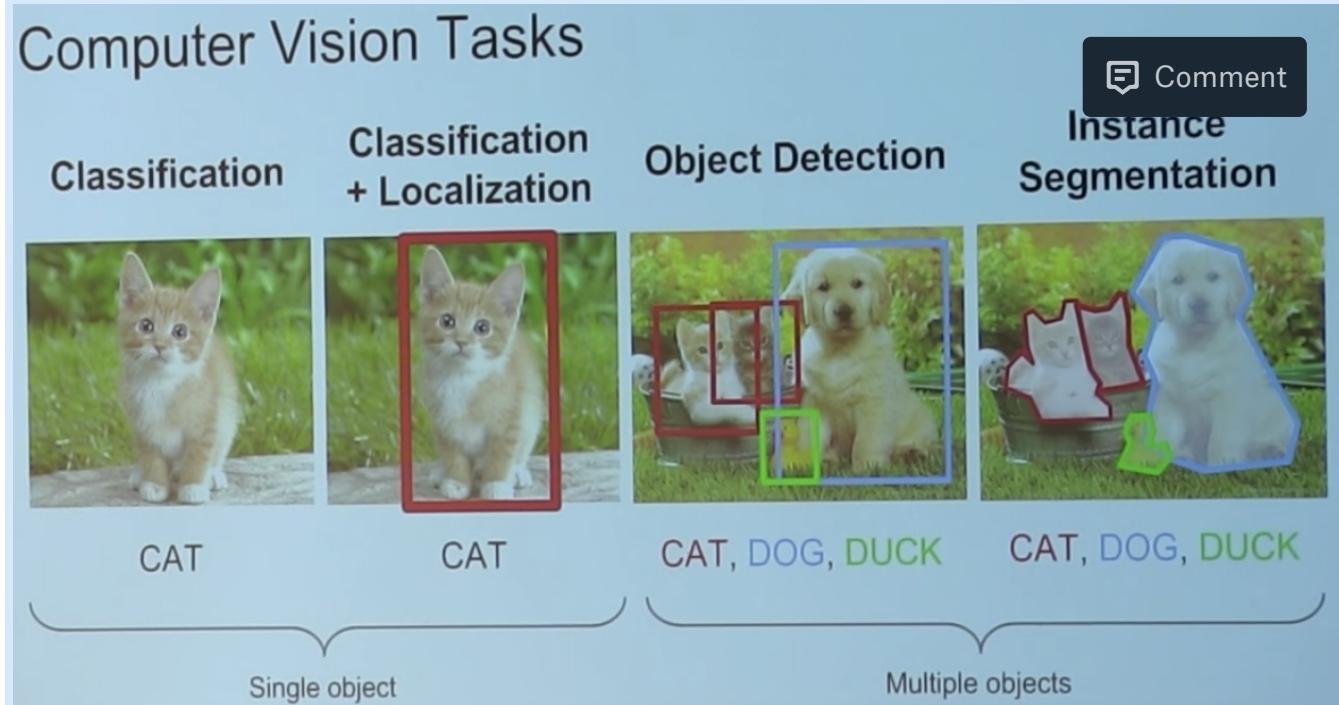


# segmentation & attention



## Classification + Localization: ImageNet

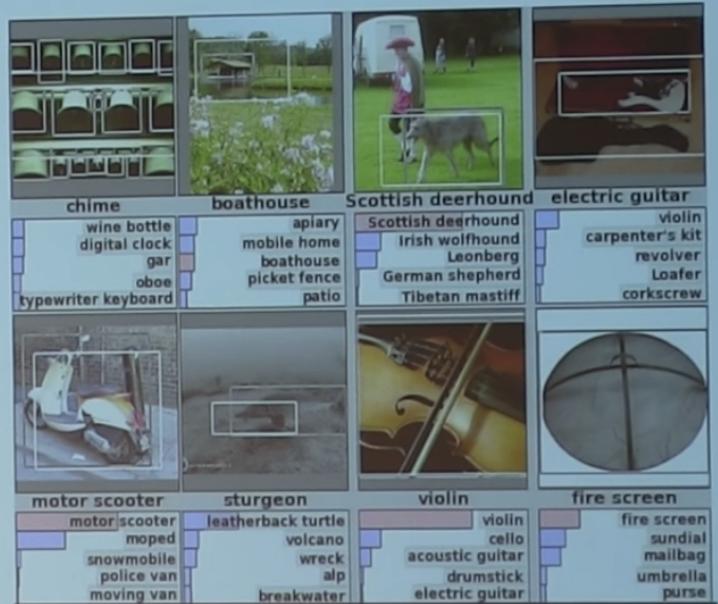
1000 classes (same as classification)

Each image has 1 class, at least one bounding box

~800 training images per class

Algorithm produces 5 (class, box) guesses

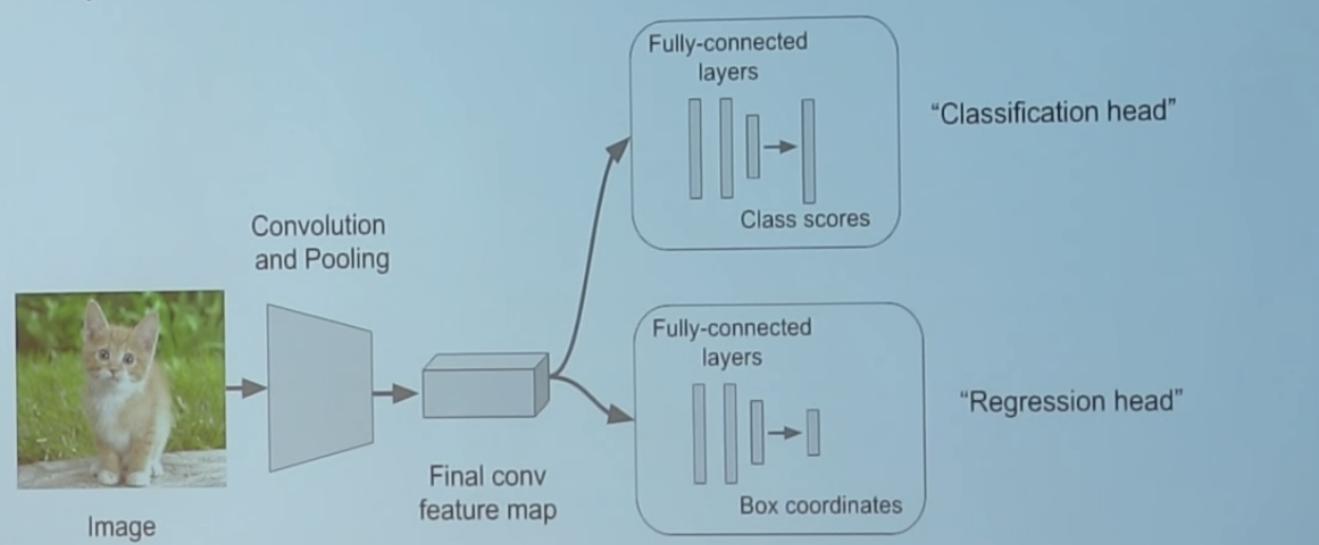
Example is correct if at least one guess has correct class AND bounding box at least 0.5 intersection over union (IoU)



Krahenbuhl et al. 2012

# Simple Recipe for Classification + Localization

Step 2: Attach new fully-connected “regression head” to the network

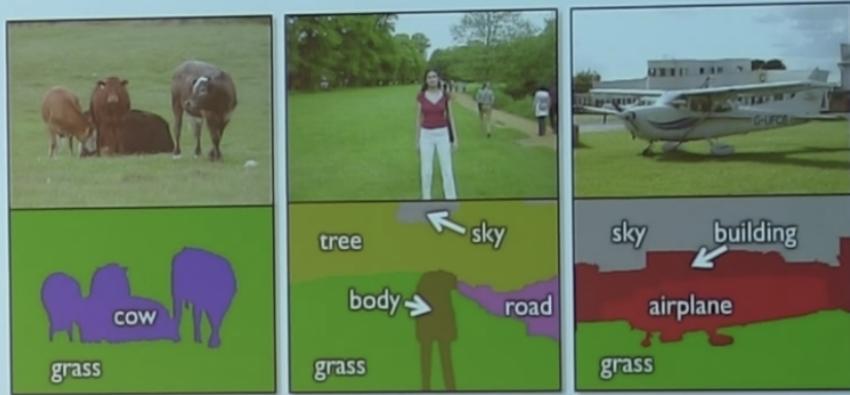


## Semantic Segmentation

Label every pixel!

Don't differentiate instances (cows)

Classic computer vision problem



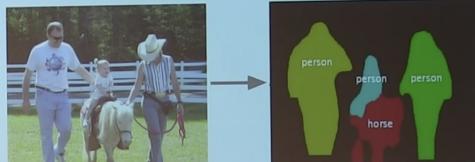
object classes	building	grass	tree	cow	sheep	sky	airplane	water	face	car
bicycle	flower	sign	bird	book	chair	road	cat	dog	body	boat

## Instance Segmentation

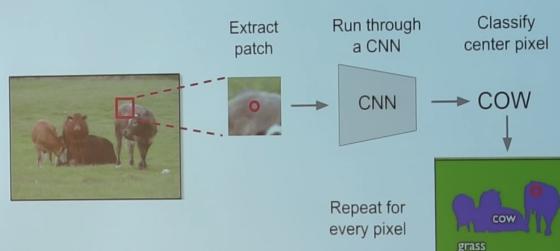
Detect instances,  
give category, label  
pixels

“simultaneous  
detection and  
segmentation” (SDS)

Lots of recent work  
(MS-COCO)

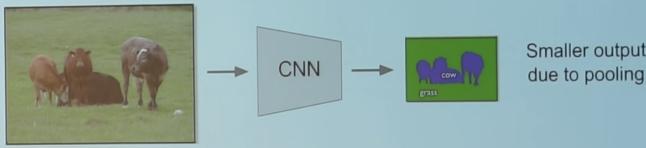


## Semantic Segmentation



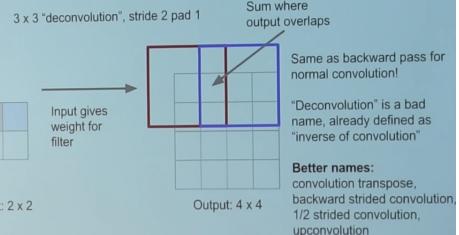
## Semantic Segmentation

Run "fully convolutional" network  
to get all pixels at once

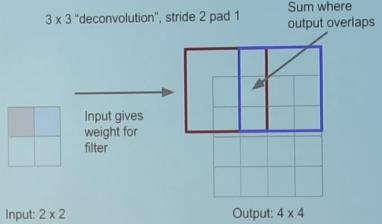


Smaller output  
due to pooling

## Learnable Upsampling: "Deconvolution"



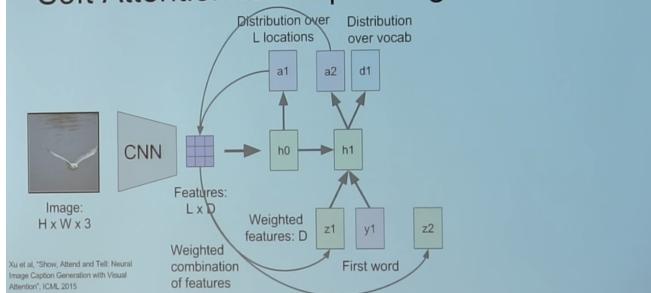
## Learnable Upsampling: "Deconvolution"



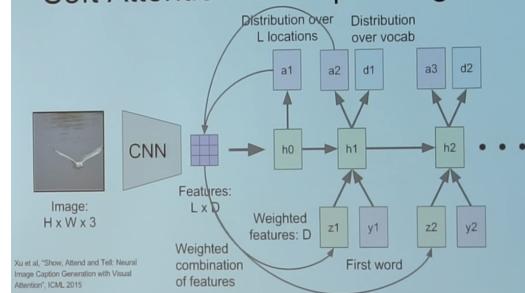
## Segmentation Overview

- Semantic segmentation
  - Classify all pixels
  - Fully convolutional models, downsample then upsample
  - Learnable upsampling: fractionally strided convolution
  - Skip connections can help
- Instance Segmentation
  - Detect instance, generate mask
  - Similar pipelines to object detection

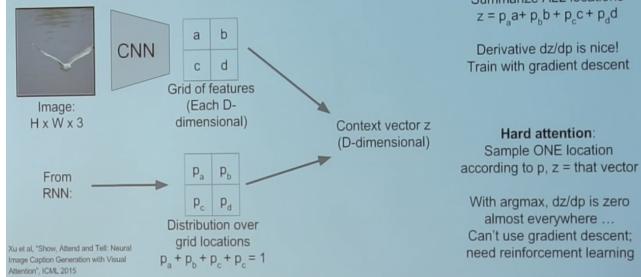
## Soft Attention for Captioning



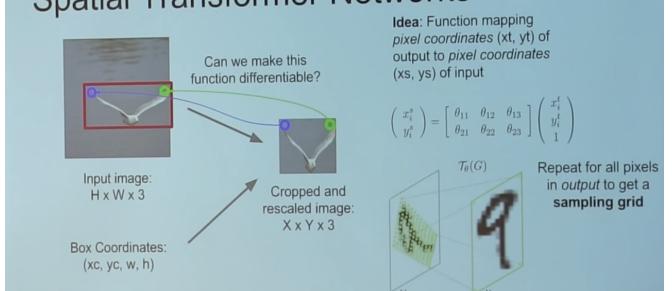
## Soft Attention for Captioning



## Soft vs Hard Attention



## Spatial Transformer Networks



# Attention Recap

- Soft attention:
  - Easy to implement: produce distribution over input locations, reweight features and feed as input
  - Attend to arbitrary input locations using spatial transformer networks
- Hard attention:
  - Attend to a single input location
  - Can't use gradient descent!
  - Need **reinforcement learning!**