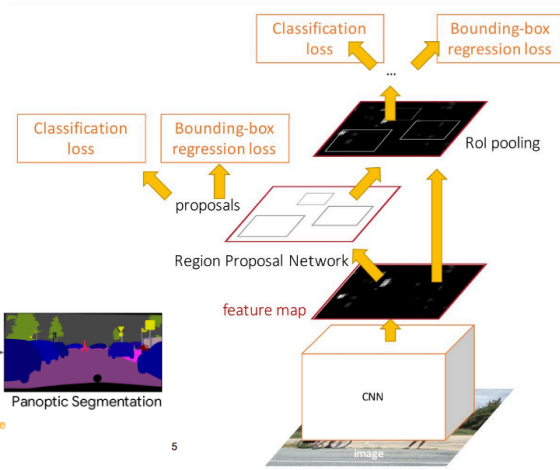
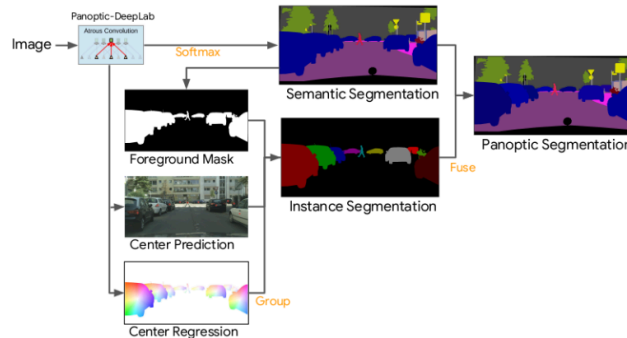
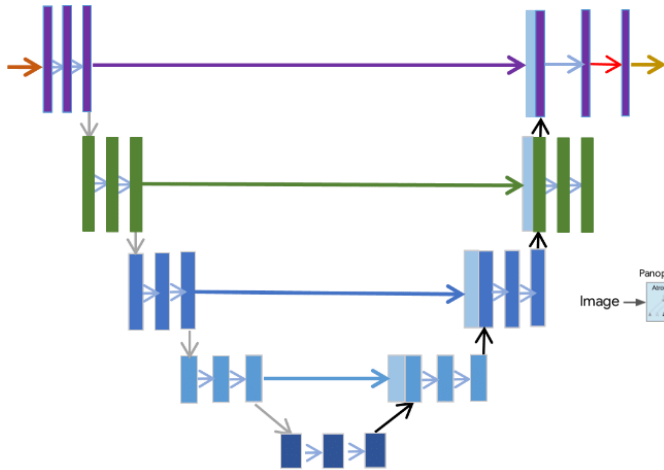


CSCI 497P/597P: Computer Vision

Convolutional Neural Networks: Other high-level problems



Announcements

- P3 grading underway, out this week
- Reminder: HW5
 - Lowest HW grade is dropped
 - Submit by tonight to guarantee grading before the final
 - Submit by Friday night (without penalty) to get credit
- P4 due tonight
- Final exam – takehome
 - Dates TBA in the next day or so.

out T 12/8
due W 12/9

Announcements

- Course evaluations
 - you have an email with a link
 - please, please, please fill it out!
 - closes “evening of” Sunday 12/6

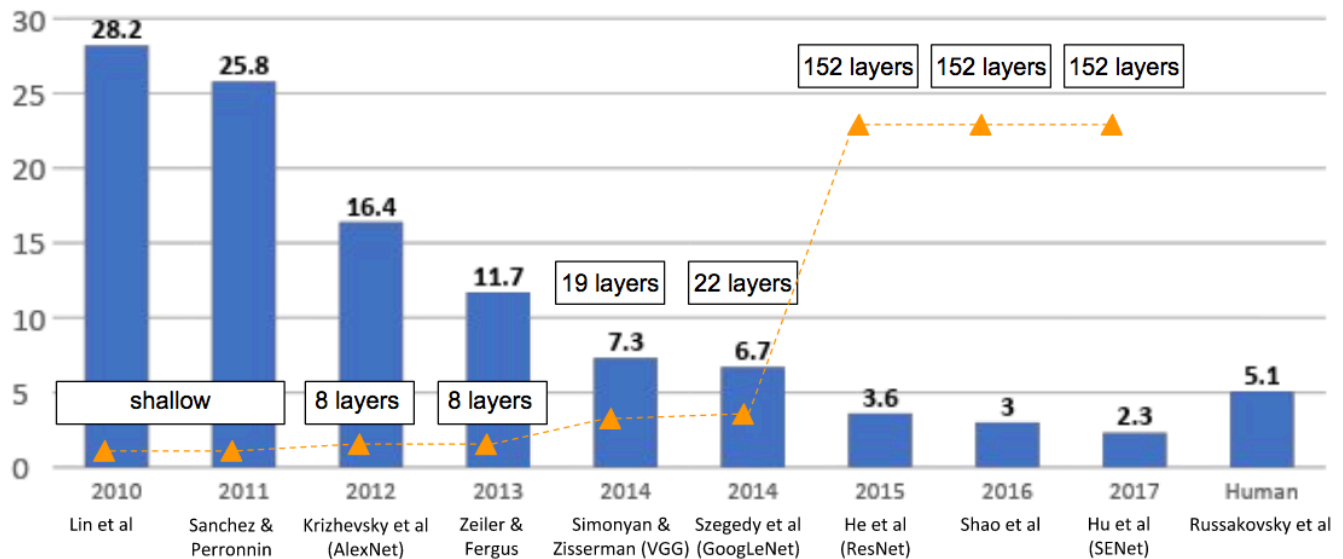
This week

- Review questions welcome
- Monday: CNNs for other high-level problems
 - semantic segmentation
 - object detection
 - panoptic segmentation
- Tuesday: generative models, deep dream, style transfer
- Wednesday: (fast) Bilateral Filter
- Friday: no new topics; AMA

And so on and so forth...

- So we've beat the crap out of ImageNet... what now?

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



And so on and so forth...

- So we've beat the crap out of ImageNet... what now?
 - **Can we do image classification on other datasets?**
 - Can we do things other than image classification?

Transfer Learning

“You need a lot of a data if you want to train/use CNNs”

Transfer Learning

“You need a lot of data if you want to train/use CNNs”

BUSTED
well... sort of

Transfer Learning with CNNs

1. Train on Imagenet



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

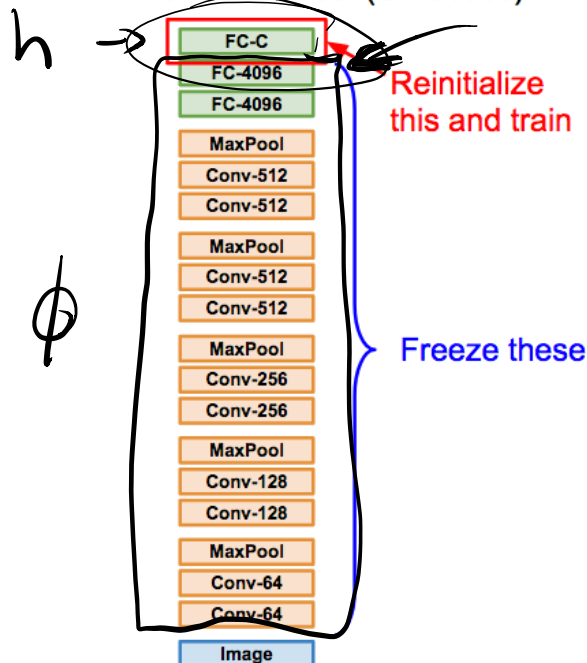
Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

1. Train on Imagenet



2. Small Dataset (C classes)



Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

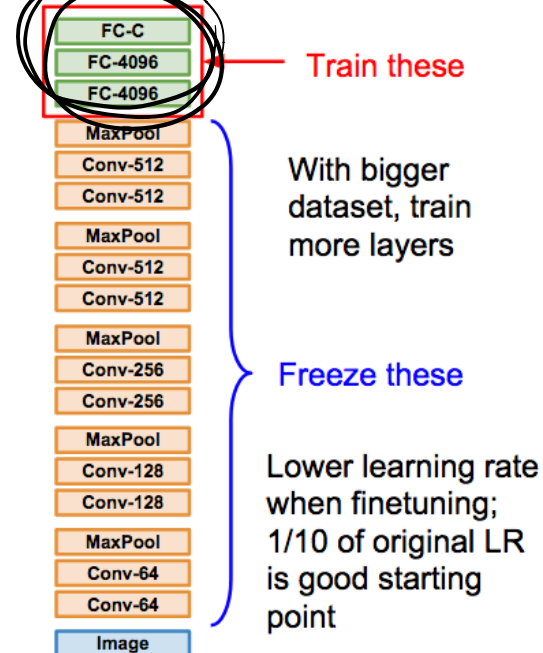
1. Train on Imagenet



2. Small Dataset (C classes)



3. Bigger dataset



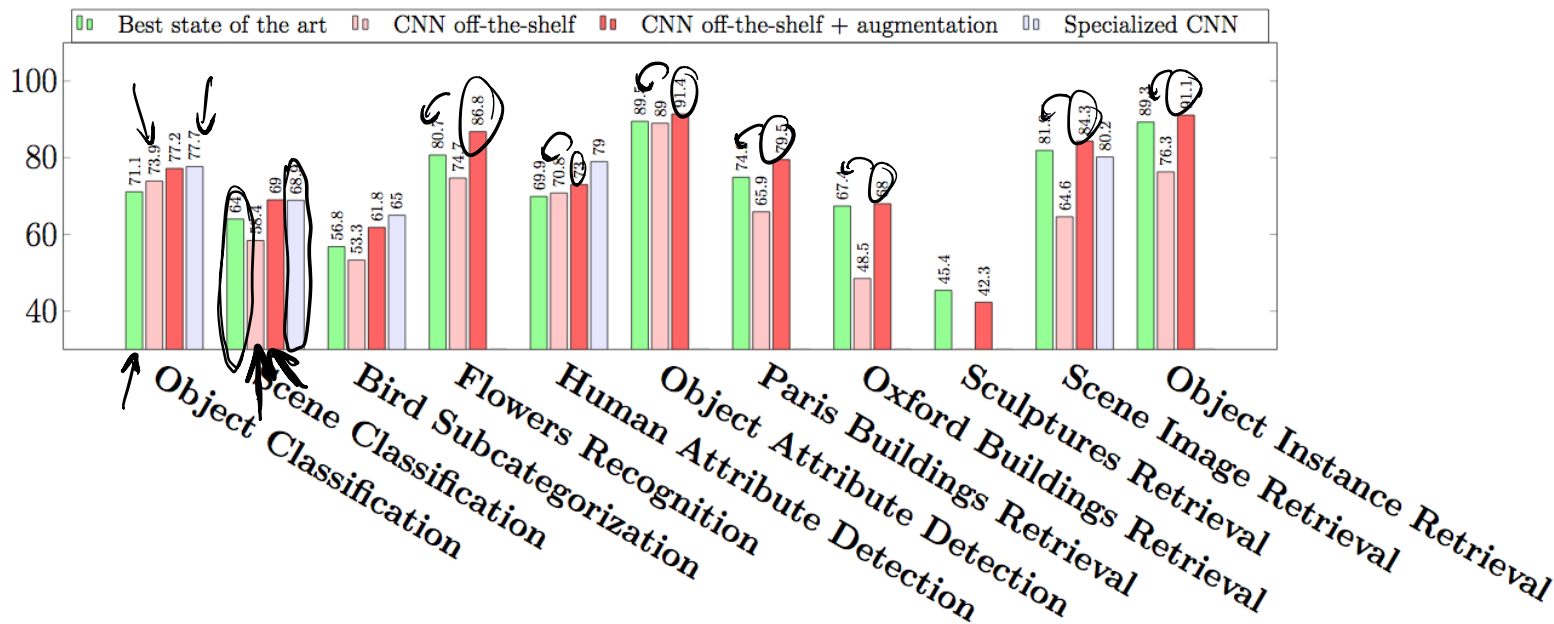
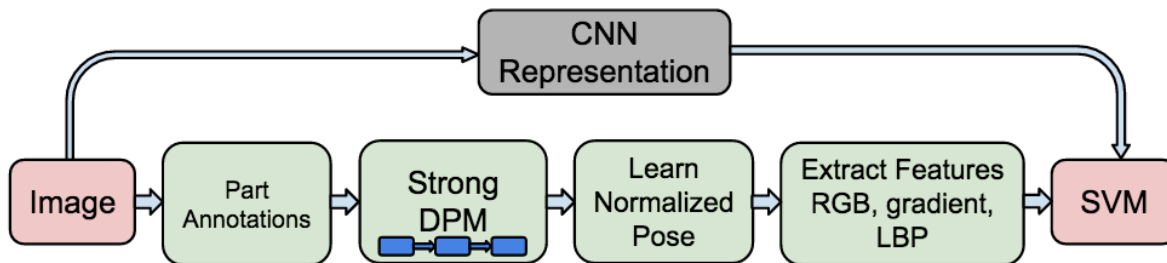


Figure: Razavian et al.: CNN Features off-the-shelf: an Astounding Baseline for Recognition

<https://arxiv.org/pdf/1403.6382.pdf>

Transfer learning with CNNs is pervasive... (it's the norm, not an exception)

Object Detection (Fast R-CNN)

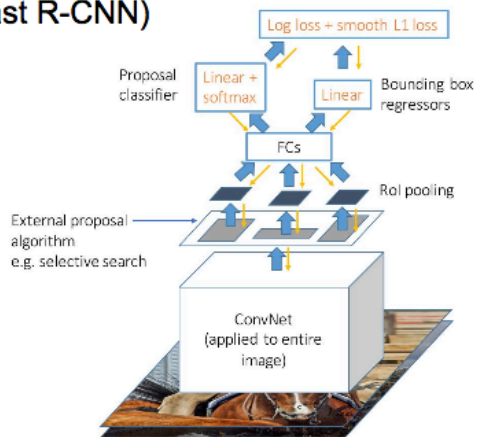
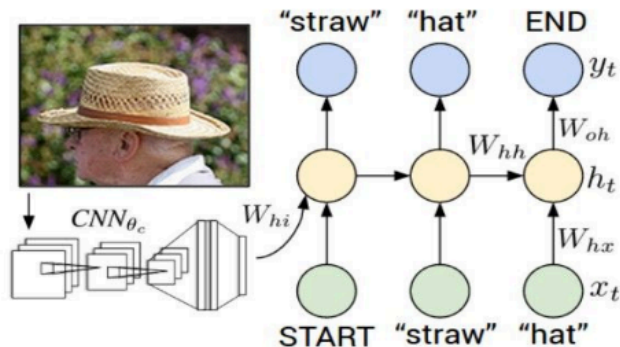


Image Captioning: CNN + RNN

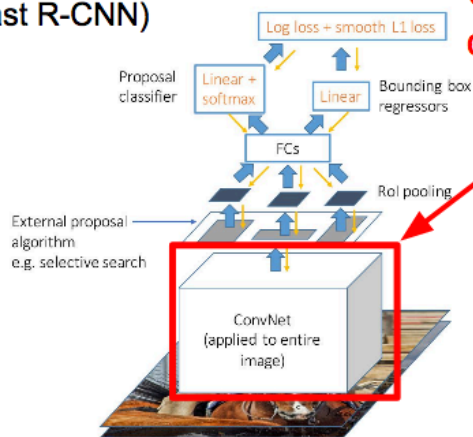


Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for
Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

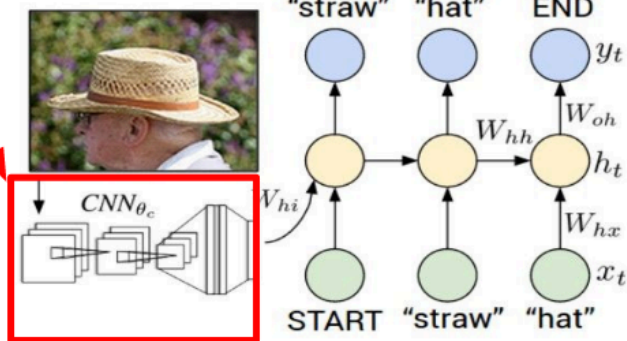
Transfer learning with CNNs is pervasive... (it's the norm, not an exception)

Object Detection (Fast R-CNN)



**CNN pretrained
on ImageNet**

Image Captioning: CNN + RNN






And so on and so forth...

- So we've beat the crap out of ImageNet... what now?
 - Can we do image classification on other datasets?
 - **Can we do things other than image classification?**

But first...

- A brief note about datasets.
- ImageNet is a collection of images with labels
 - The 1000 classes used for evaluation are a tiny subset of the tags available.
 - The labels were produced by humans.



newsreader, news reader



microeconomist, microeconomic expert



Lil Uzi Hurt at Home

@lostblackboy



No matter what kind of image I upload, ImageNet Roulette, which categorizes people based on an AI that knows 2500 tags, only sees me as Black, Black African, Negroid or Negro.

Some of the other possible tags, for example, are “Doctor,” “Parent” or “Handsome.”



♡ 511 5:08 PM - Sep 17, 2019 · Brooklyn, NY



🗨️ 184 people are talking about this

- **The viral selfie app ImageNet Roulette seemed fun – until it called me a racist slur**

The Guardian, September 2019

<https://www.theguardian.com/technology/2019/sep/17/imagenet-roulette-asian-racist-slur-selfie>

600,000 Images Removed from AI Database After Art Project Exposes Racist Bias

The image tagging system that went viral on social media was part of artist Trevor Paglen and AI researcher Kate Crawford's attempts to publicize how prejudiced technology can be.

<https://hyperallergic.com/518822/600000-images-removed-from-ai-database-after-art-project-exposes-racist-bias/>

Dataset bias

LFW

[Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Huang et al.]

77.5% male
83.5% white

IJB-A

[Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus benchmark. Klare et al.]

79.6% lighter-skinned

Adience




[Age and gender classification using convolutional neural networks. Levi and Hassner.]

86.2% lighter-skinned

[Buolamwini and Gebru. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification]

Error Rate_(1-PPV) By Female x Skin Type



	TYPE I	TYPE II	TYPE III	TYPE IV	TYPE V	TYPE VI
	1.7%	1.1%	3.3%	0%	23.2%	25.0%
	11.9%	9.7%	8.2%	13.9%	32.4%	46.5%
	5.1%	7.4%	8.2%	8.3%	33.3%	46.8%

Buolamwini & Gebru FAT* 2018, Slides from Joy Buolamwini

See also

- Writeup on ImageNet Bias:
<https://www.excavating.ai/>
- ACM Conference on Fairness, Accountability, and Transparency:
<https://facctconference.org/index.html>
- CVPR 2020 Tutorial on Fairness, Accountability, Transparency, and Ethics in Vision:
<https://sites.google.com/view/fatecv-tutorial>

Other Computer Vision Tasks

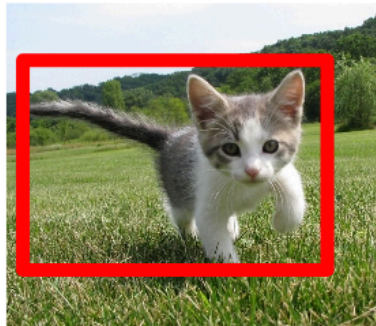
Semantic Segmentation



**GRASS, CAT,
TREE, SKY**

No objects, just pixels

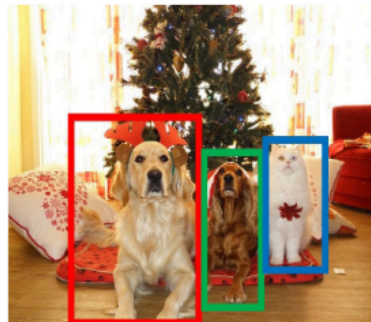
**Classification
+ Localization**



CAT

Single Object

**Object
Detection**



DOG, DOG, CAT

Multiple Object

**Instance
Segmentation**



DOG, DOG, CAT

This image is CC0 public domain

Other Computer Vision Tasks

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



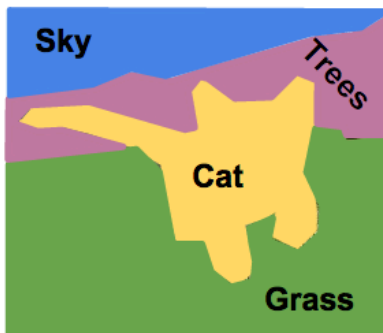
DOG, DOG, CAT

This image is CC0 public domain

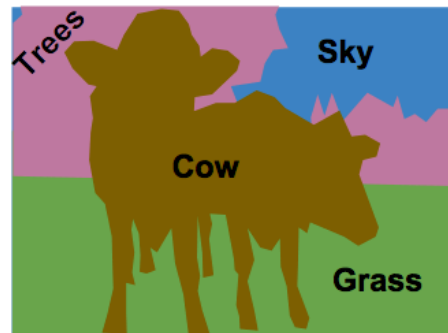
Semantic Segmentation

Label each pixel in the image with a category label

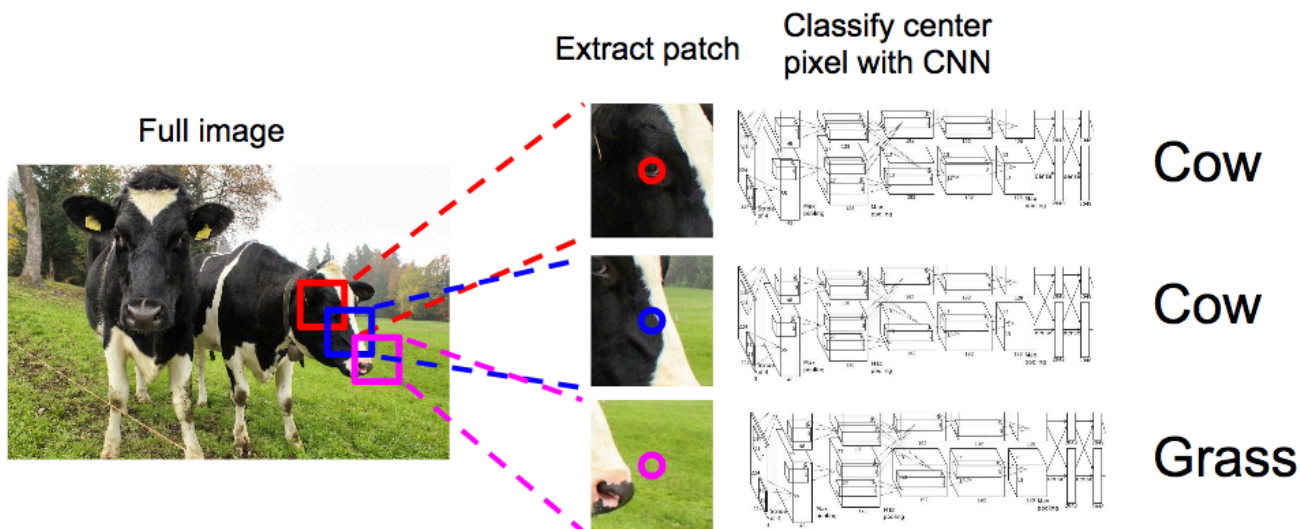
Don't differentiate instances, only care about pixels



This image is CC0 public domain



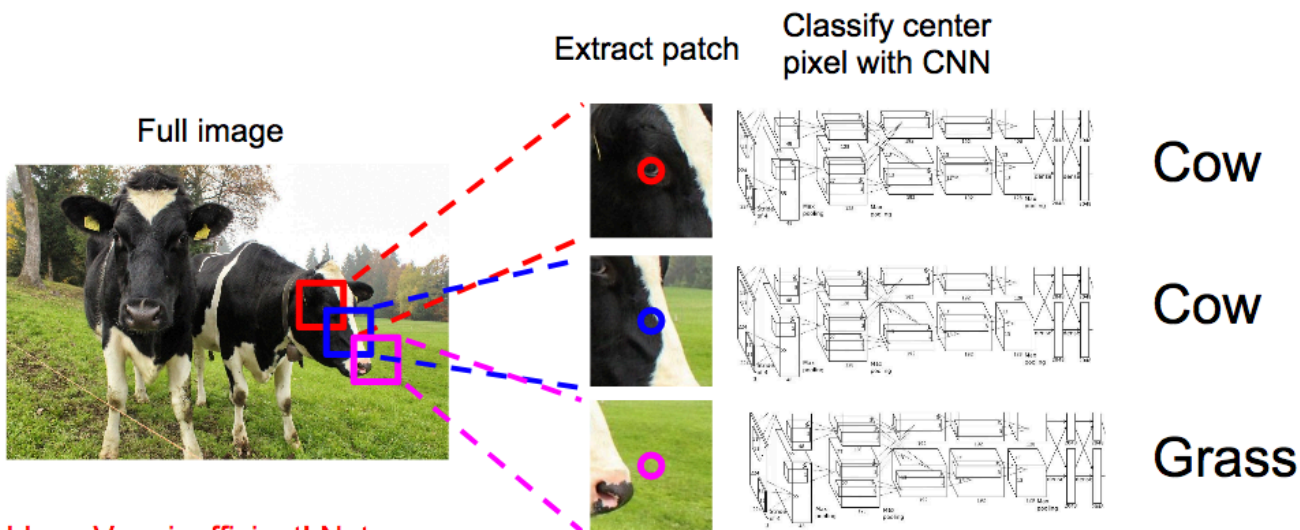
Semantic Segmentation Idea: Sliding Window



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

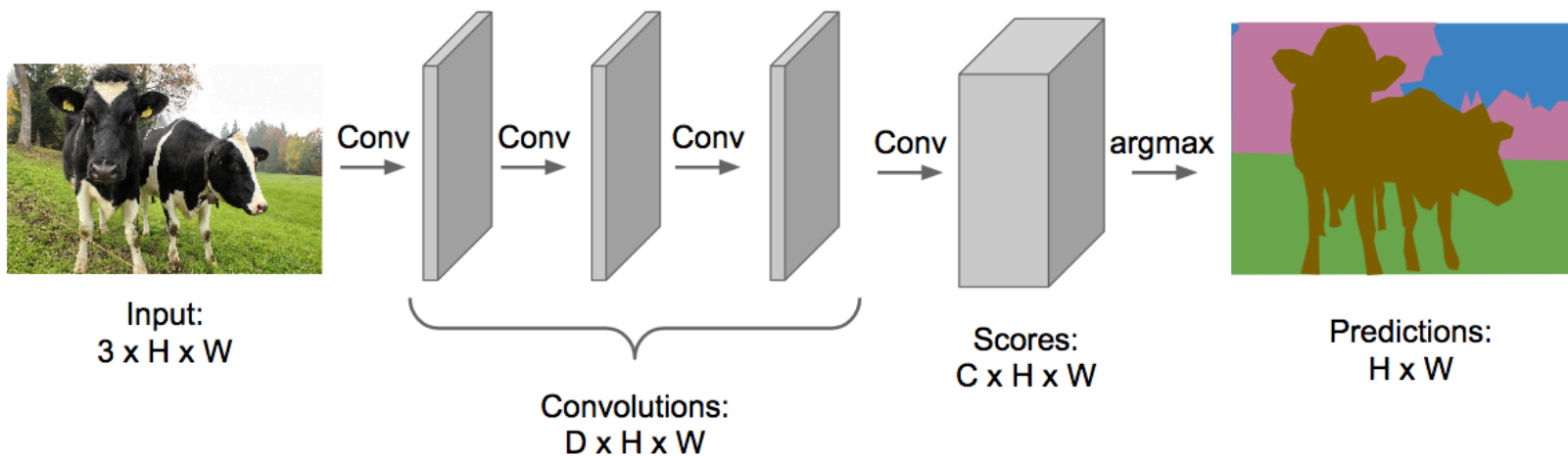


Problem: Very inefficient! Not reusing shared features between overlapping patches

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

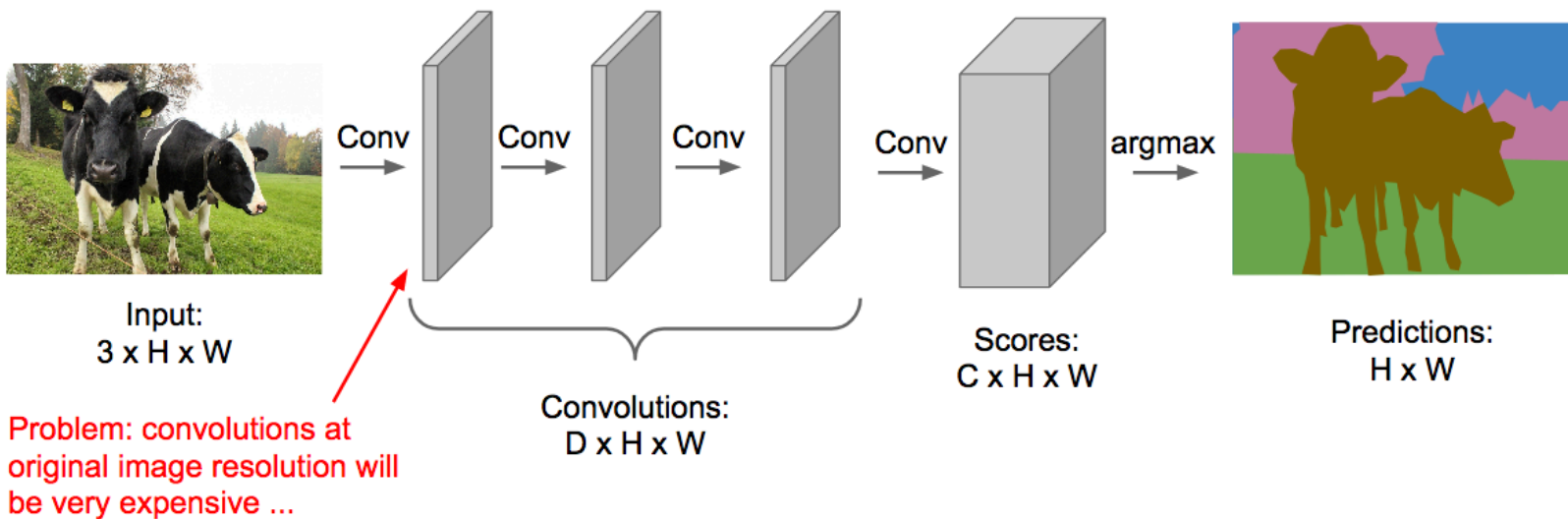
Semantic Segmentation Idea: Fully Convolutional

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



Semantic Segmentation Idea: Fully Convolutional

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

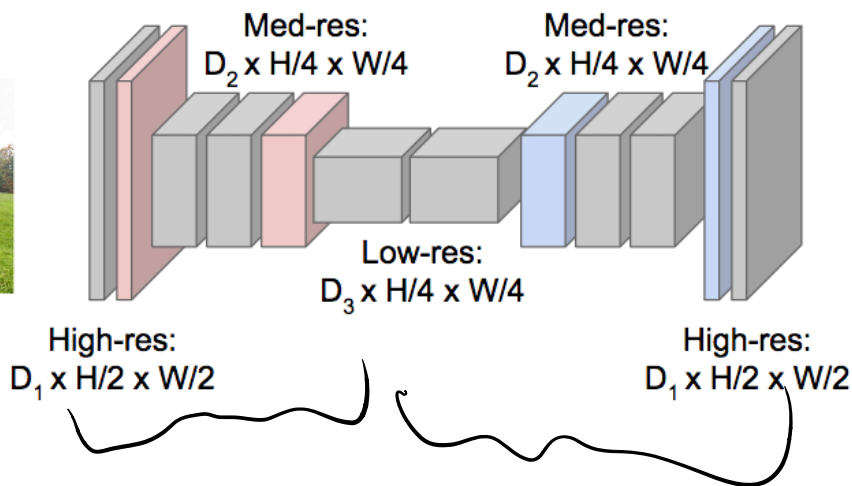


Semantic Segmentation Idea: Fully Convolutional

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



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

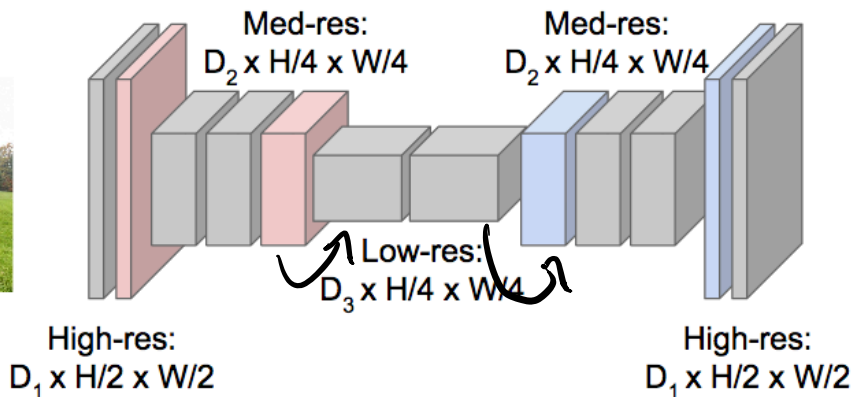
Semantic Segmentation Idea: Fully Convolutional

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:
???



Predictions:
 $H \times W$

In-Network 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 x 2

Output: 4 x 4

“Bed of Nails”

1	2
3	4



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

Input: 2 x 2

Output: 4 x 4

In-Network 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 x 4



5	6
7	8

Output: 2 x 2



Rest of the network

Max Unpooling

Use positions from pooling layer

1	2
3	4

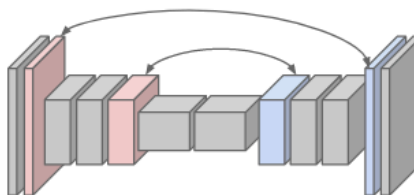
Input: 2 x 2



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

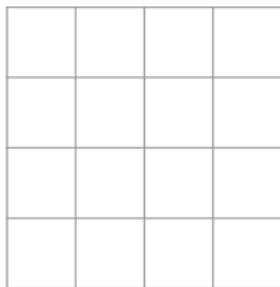
Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers

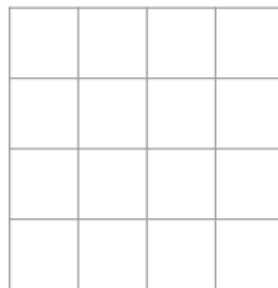


Learnable Upsampling: Transpose Convolution

Recall: Typical 3 x 3 convolution, stride 1 pad 1



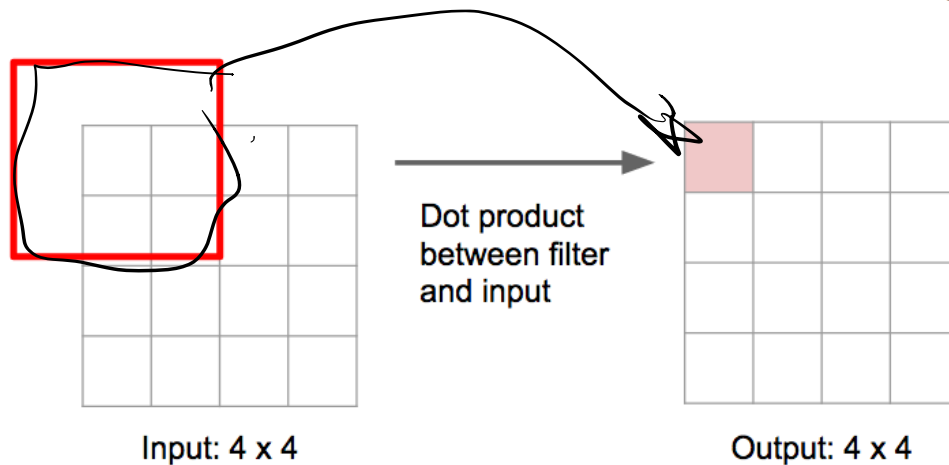
Input: 4 x 4



Output: 4 x 4

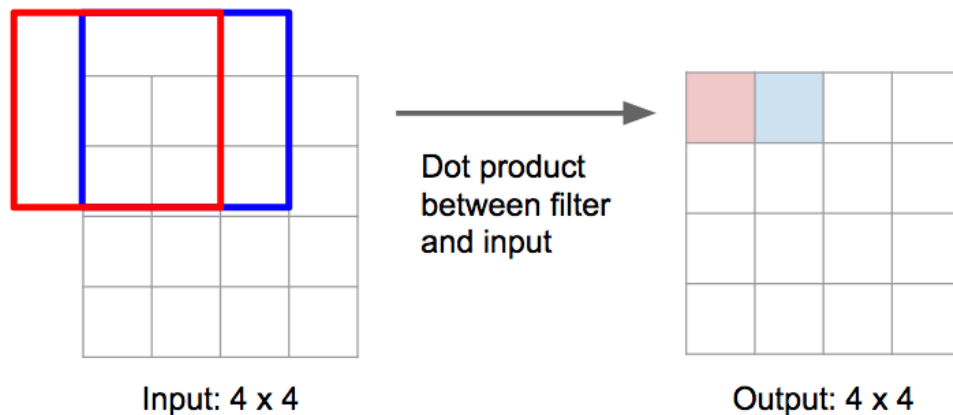
Learnable Upsampling: Transpose Convolution

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



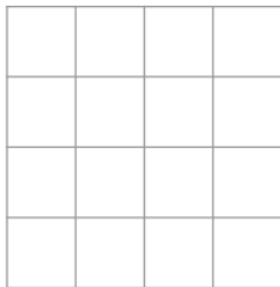
Learnable Upsampling: Transpose Convolution

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



Learnable Upsampling: Transpose Convolution

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



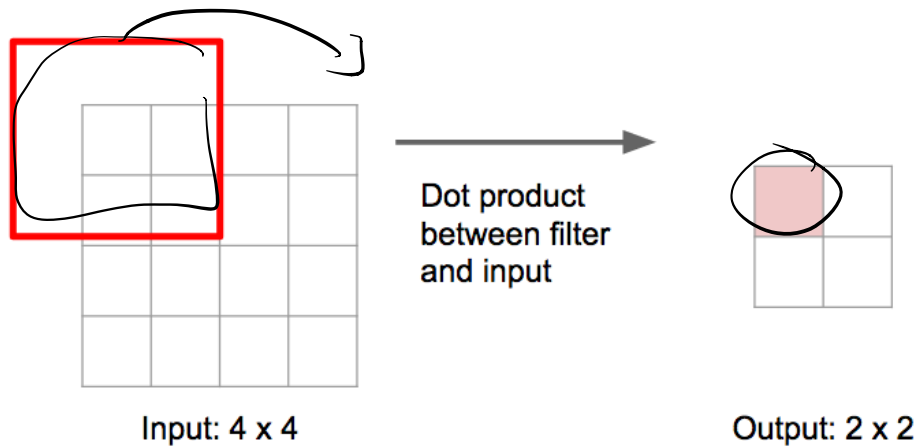
Input: 4 x 4



Output: 2 x 2

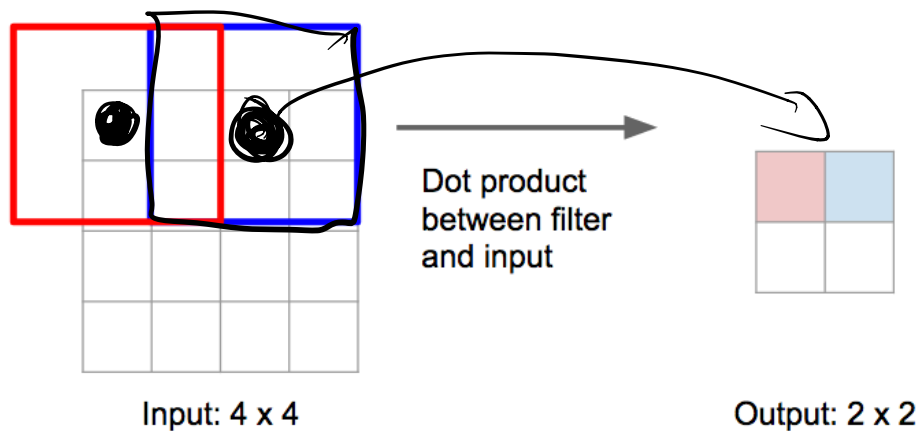
Learnable Upsampling: Transpose Convolution

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



Learnable Upsampling: Transpose Convolution

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

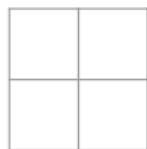


Filter moves 2 pixels in the input for every one pixel in the output

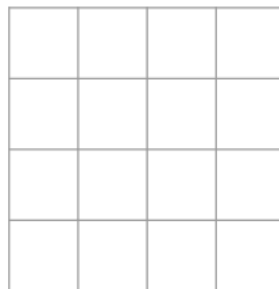
Stride gives ratio between movement in input and output

Learnable Upsampling: Transpose Convolution

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



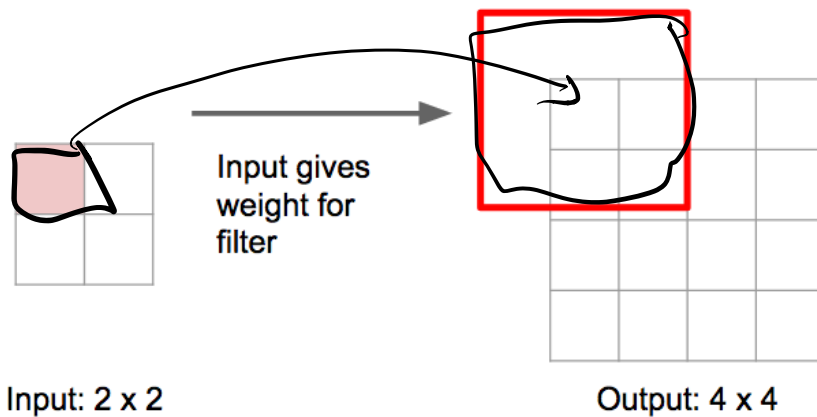
Input: 2 x 2



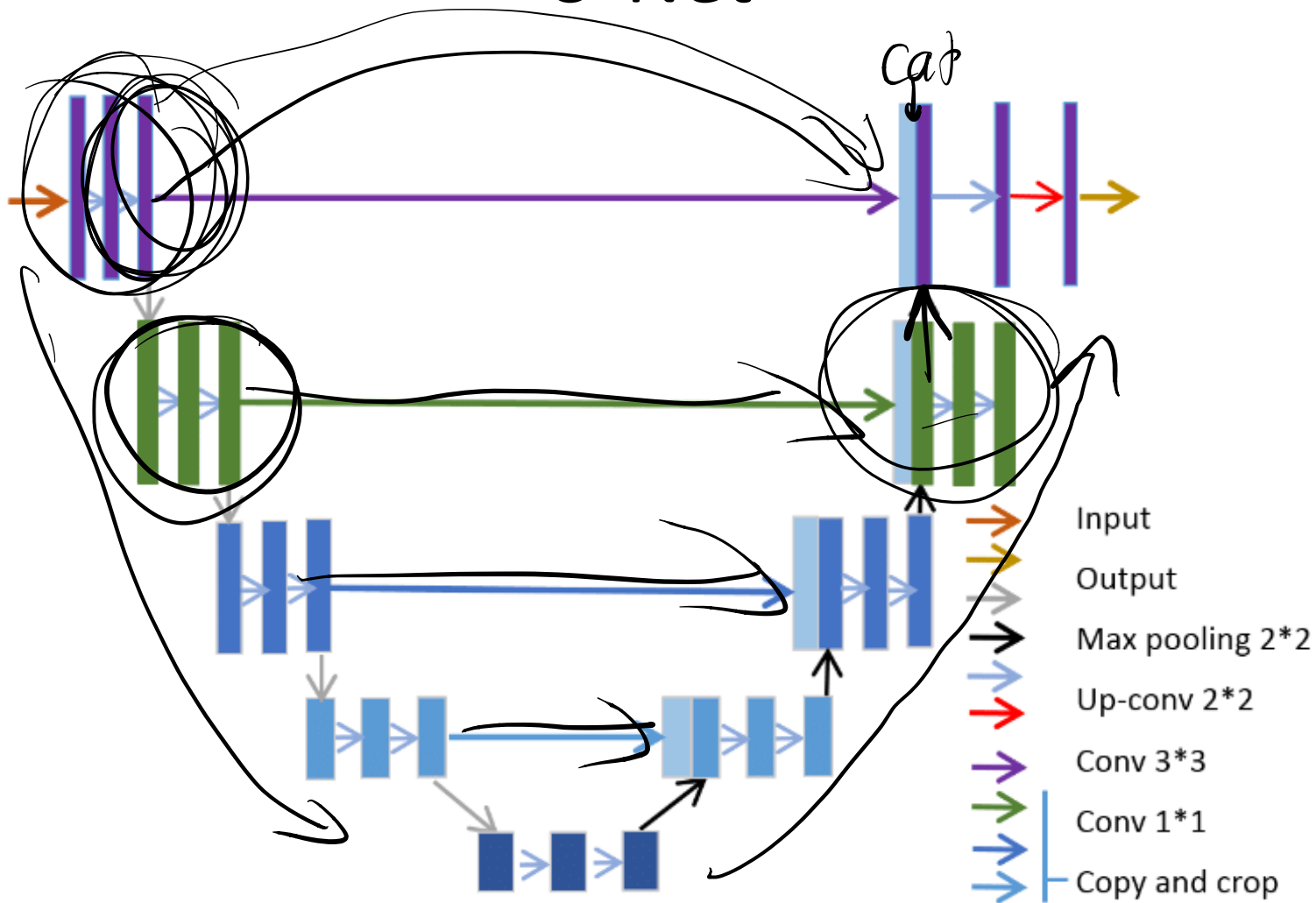
Output: 4 x 4

Learnable Upsampling: Transpose Convolution

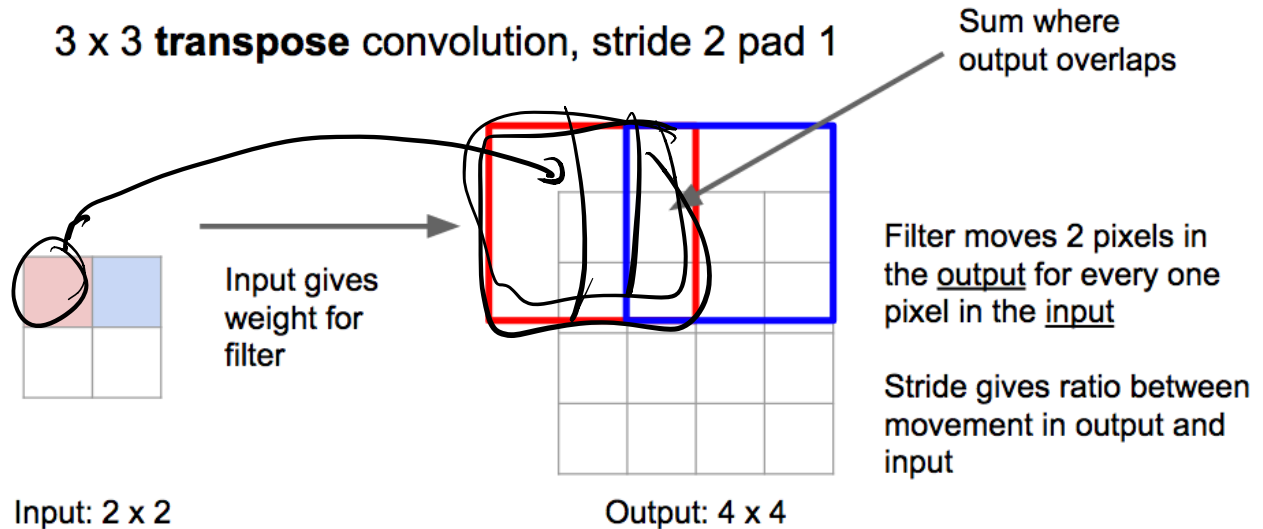
3 x 3 **transpose** convolution, stride 2 pad 1



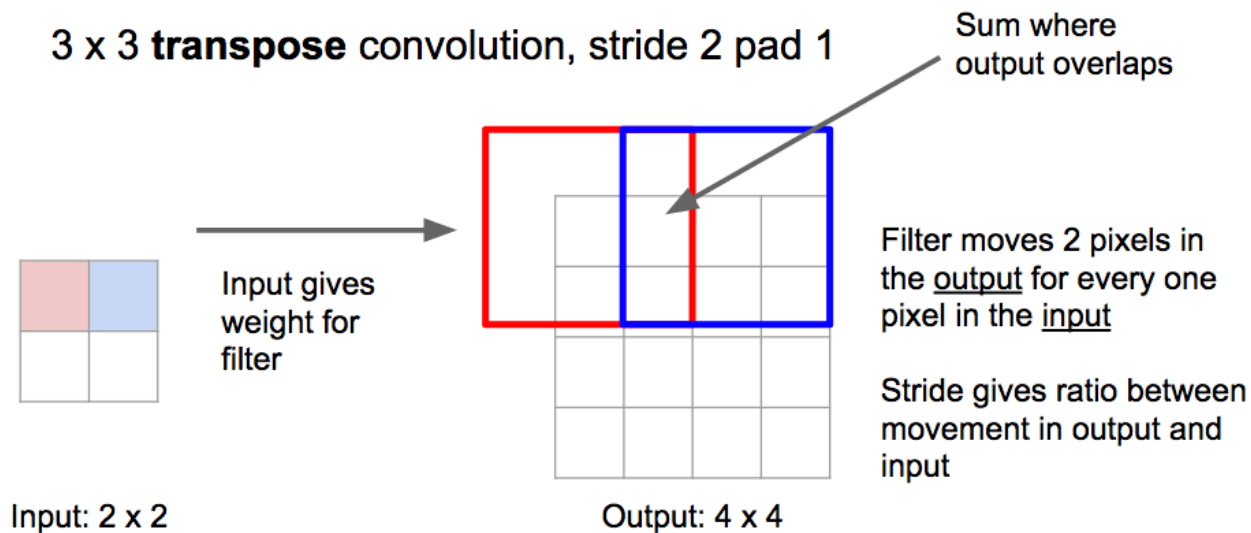
U-Net



Learnable Upsampling: Transpose Convolution



Learnable Upsampling: Transpose Convolution

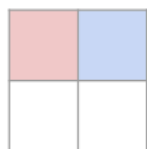


Learnable Upsampling: Transpose Convolution

Other names:

- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

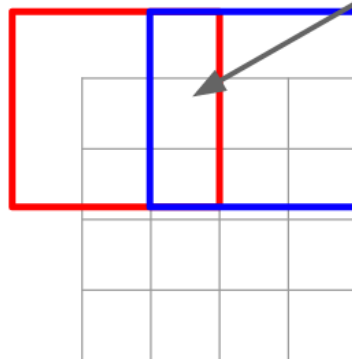
3 x 3 **transpose** convolution, stride 2 pad 1



Input: 2 x 2



Input gives weight for filter



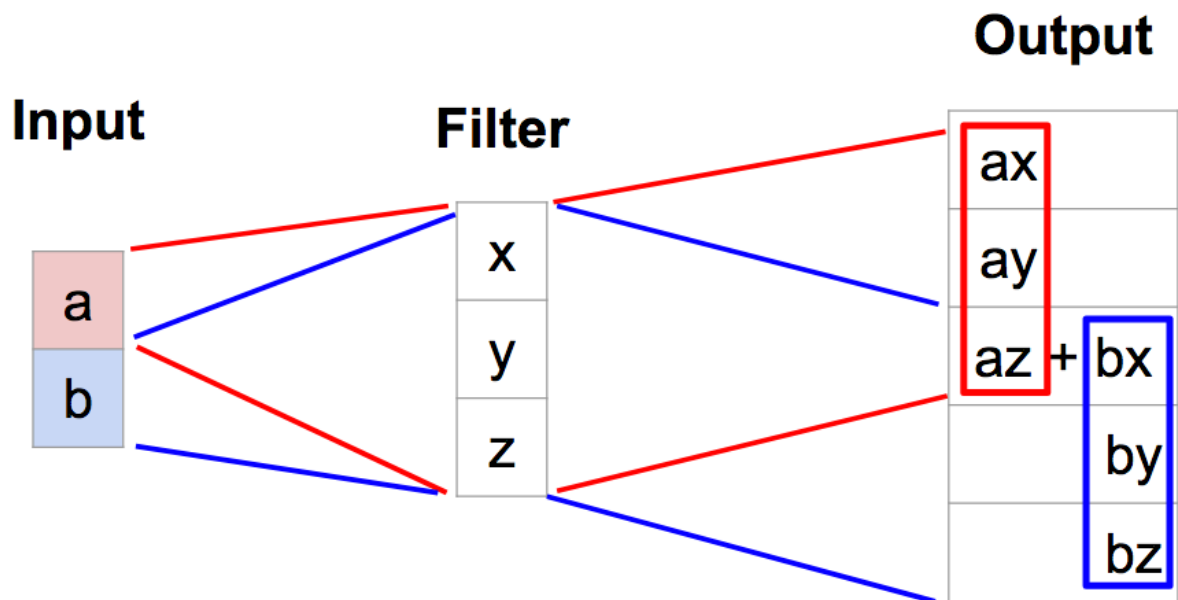
Output: 4 x 4

Sum where output overlaps

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

Stride gives ratio between movement in output and input

Learnable Upsampling: 1D Example



Output contains copies of the filter weighted by the input, summing at where it overlaps in the output

Need to crop one pixel from output to make output exactly 2x input



2D Object Detection

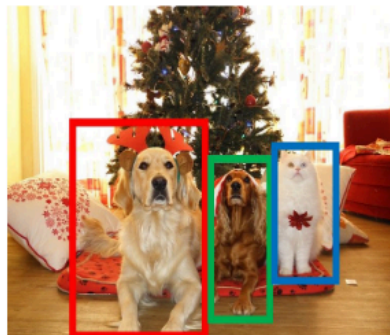
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

2D Object Detection



DOG, DOG, CAT

Object categories +
2D bounding boxes

3D Object Detection



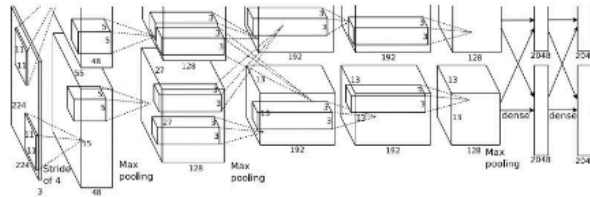
Car

Object categories +
3D bounding boxes

This image is CC0 public domain

Object Detection as Classification: Sliding Window

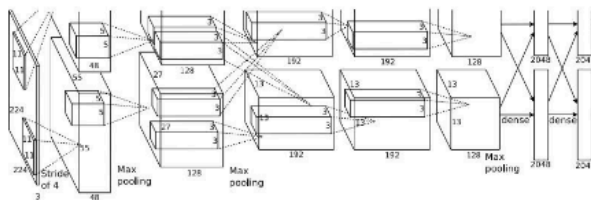
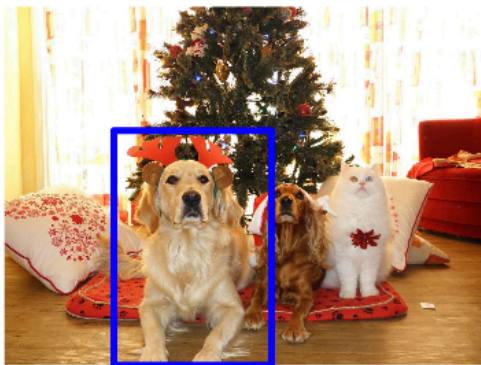
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

Object Detection as Classification: Sliding Window

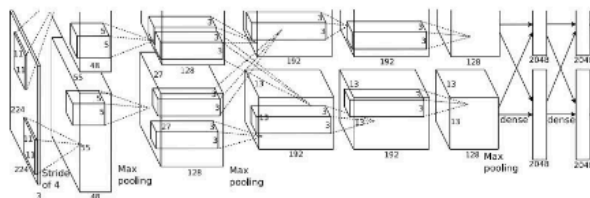
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection as Classification: Sliding Window

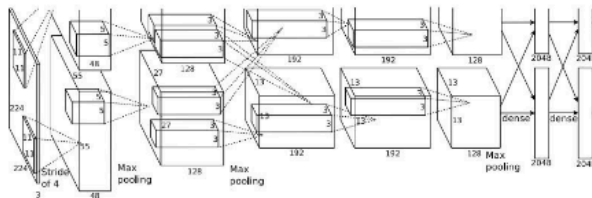
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection as Classification: Sliding Window

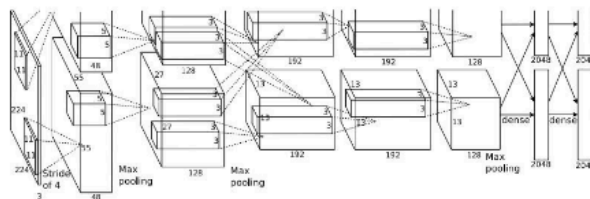
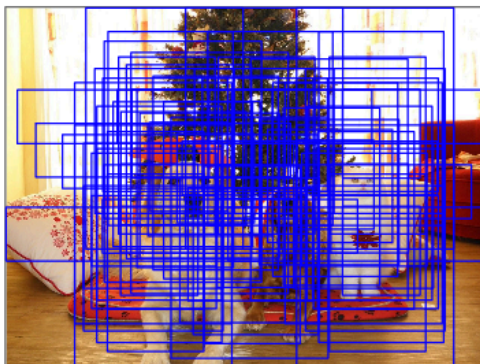
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Object Detection as Classification: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

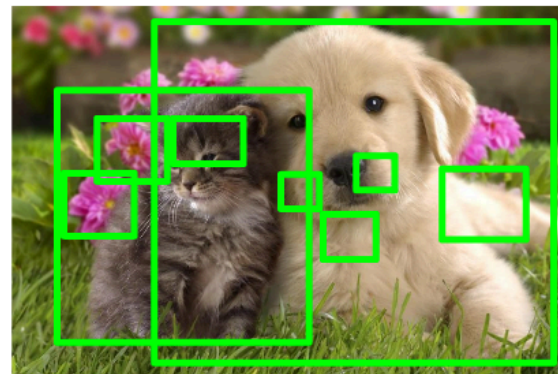


Dog? NO
Cat? YES
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

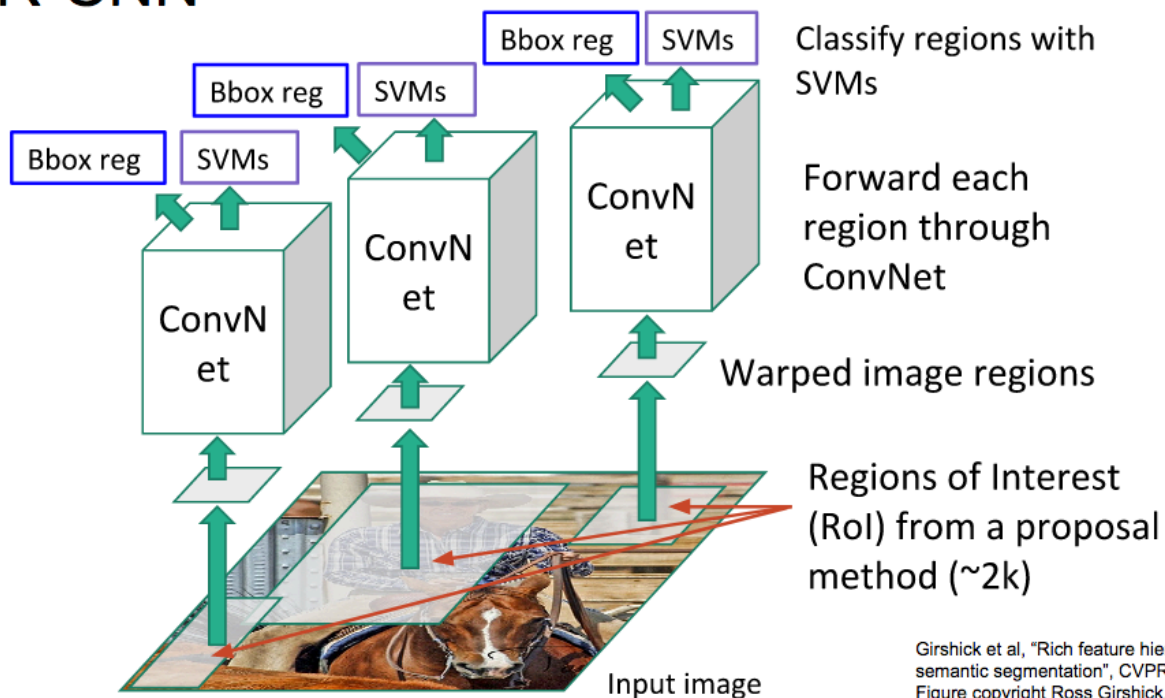
Region Proposals / Selective Search

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



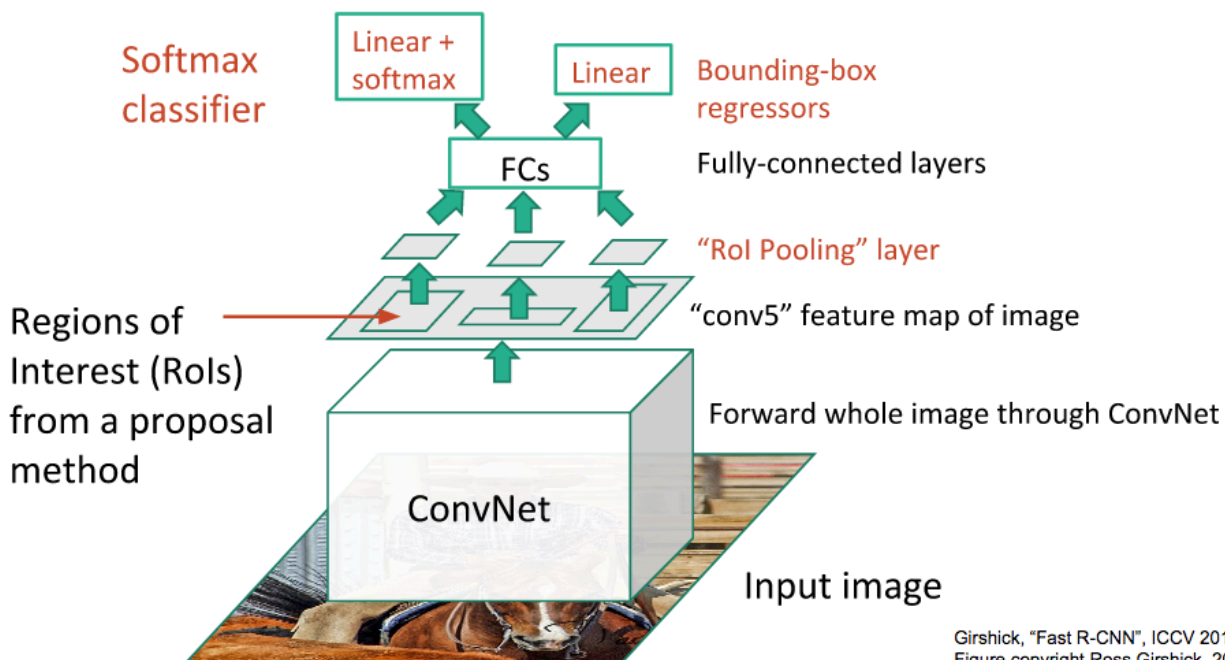
Alexe et al, "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013
Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

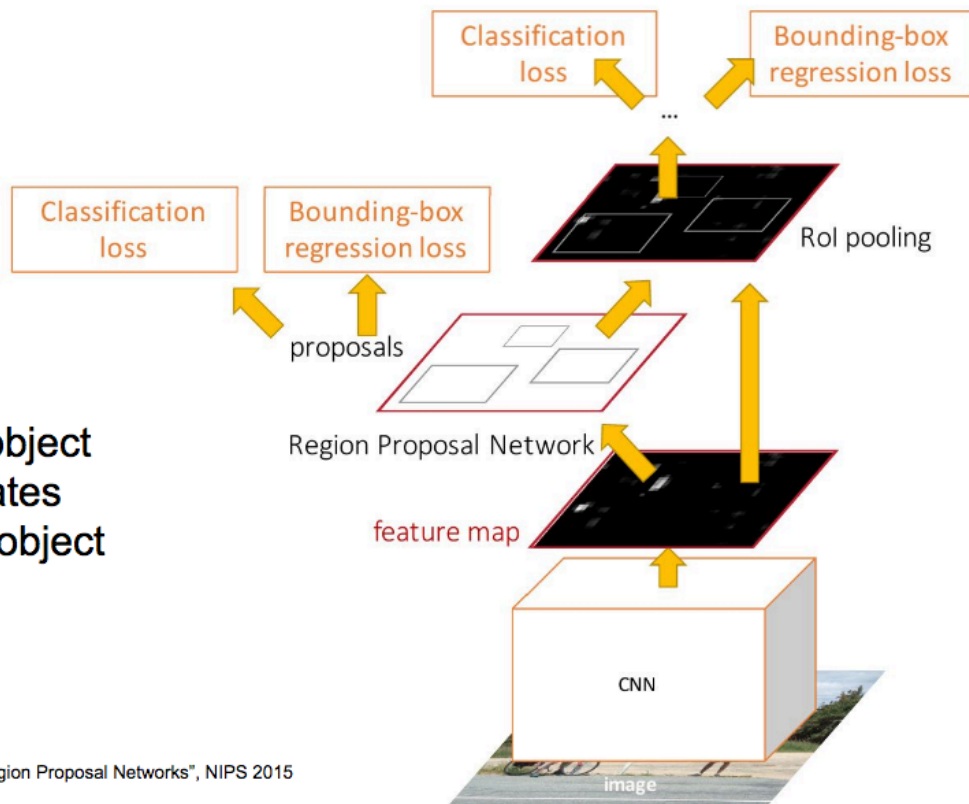
Faster R-CNN:

Make CNN do proposals!

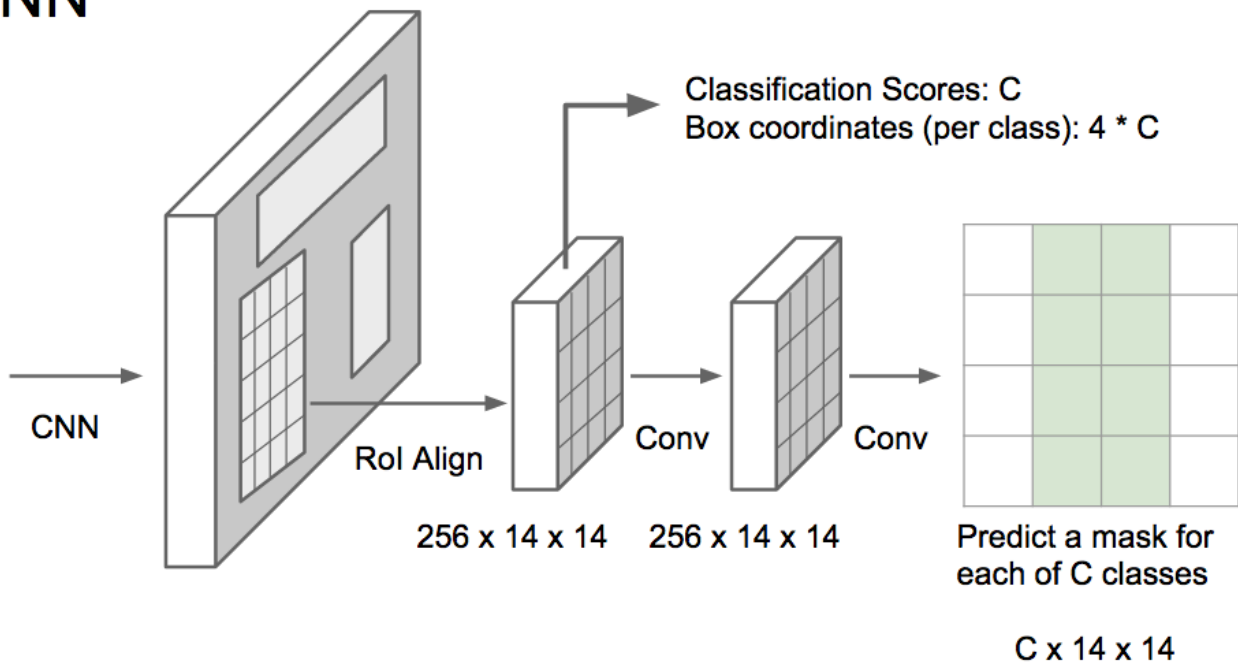
Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:

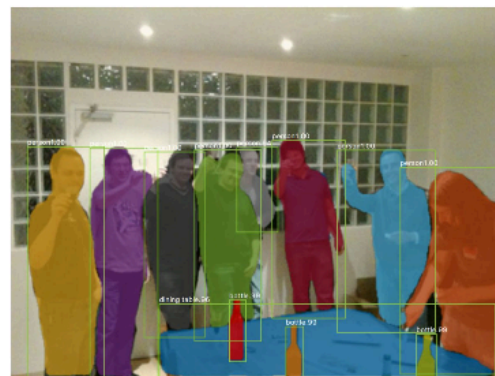
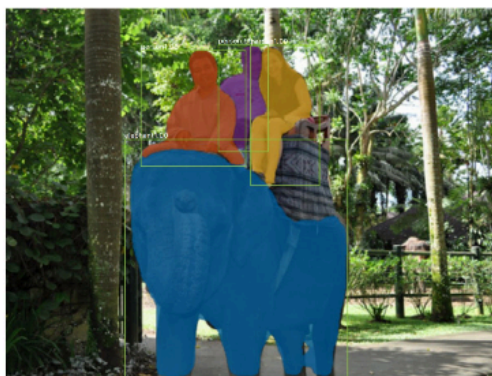
1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Mask R-CNN



Mask R-CNN: Very Good Results!



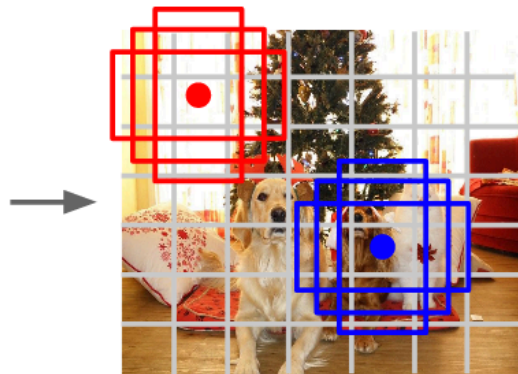
He et al, "Mask R-CNN", arXiv 2017
Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, 2017.
Reproduced with permission.

Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network! →



Input image
 $3 \times H \times W$



Divide image into grid
 7×7

Image a set of **base boxes**
centered at each grid cell
Here $B = 3$

Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers: $(dx, dy, dh, dw, \text{confidence})$
- Predict scores for each of C classes (including background as a class)

Output:
 $7 \times 7 \times (5 * B + C)$



Panoptic Segmentation

- Mash together object detection and semantic segmentation

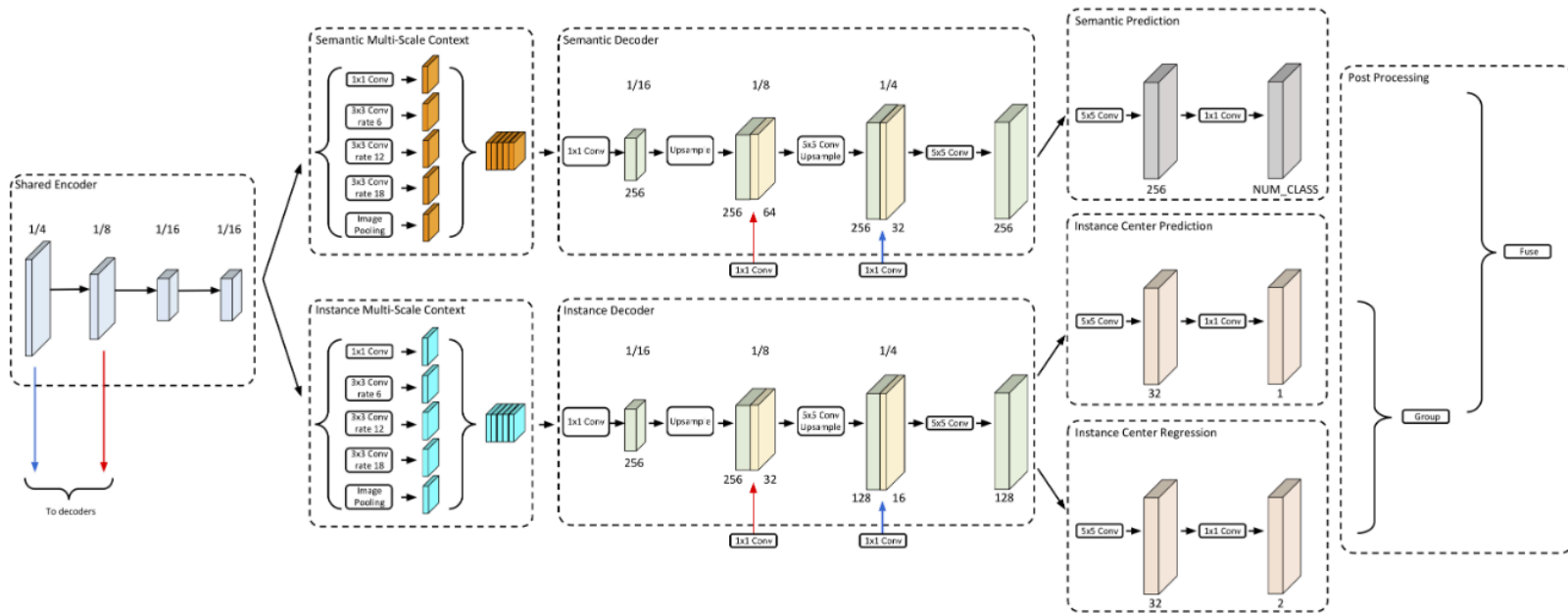
semantic (“stuff”)



instance (“things”)



Example architecture: Panoptic DeepLab



Autoencoders and Generative Models

Generative Adversarial Networks

Other Problems

- Fine-grained recognition (e.g., dog/bird species)
- Instance segmentation
- Face detection and recognition
- Motion estimation
- Feature detection and description
- Depth estimation
- Novel view synthesis
- ...and many others

