

Today: Outline

- ***Advanced Topics I (Slides from Stanford CS231; Kate Saenko DL)***
- **Announcements/Reminders:**
 - ***Advanced Topics II (Mon Jun 28)***
 - “Efficient Deep Learning in Production,
From Virtual Reality to Full Self-Driving”
Alvin Wan, Tesla, UC Berkeley
 - “Training GANs to be representative of humanity”
Tabitha Oanda, BU
 - *Five bonus points will be offered.*
 - *Problem Set 2 is graded
(regrade requests accepted until Mon Jun 28)*
 - *Problem Set 3, due Fri Jun 25*
 - *Last Pre-lec Material Posted, Due Mon Jun 28*

Today: Applications of CNNs to Computer Vision

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](#)

So far: CNN for classification

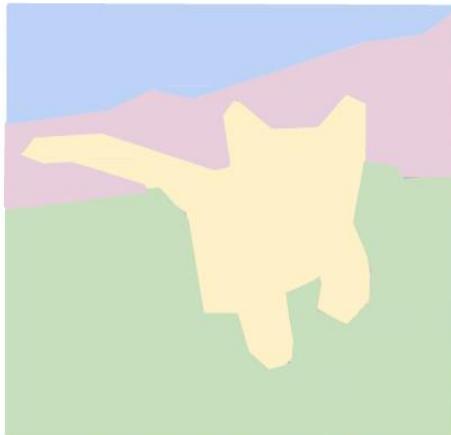
Classification



CAT

No spatial extent

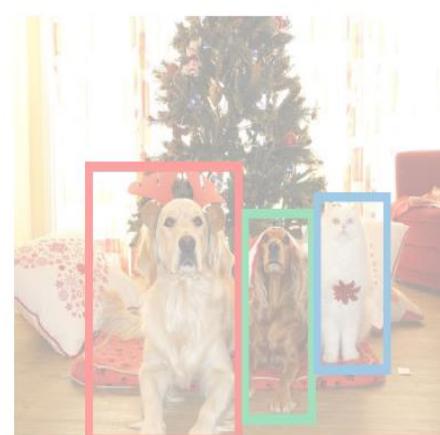
Semantic
Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Multiple Object

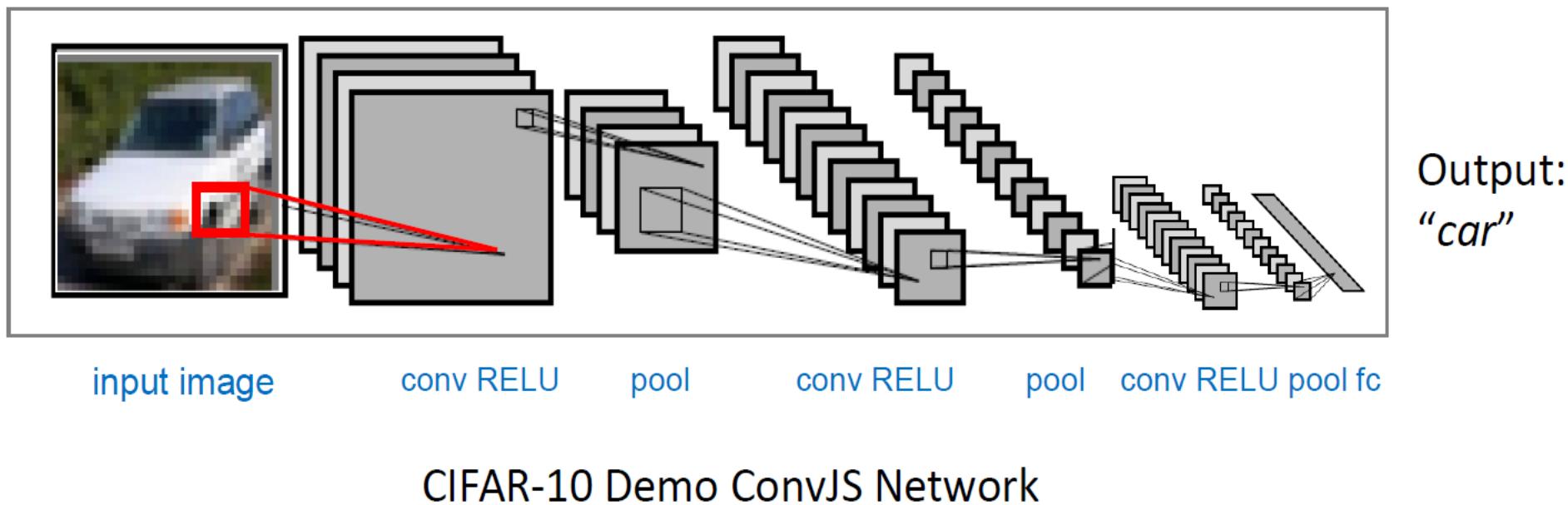
Instance
Segmentation



DOG, DOG, CAT

[This image is CC0 public domain](#)

Example: CNN architecture for classification



Today: Segmentation, Detection

Semantic Segmentation

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT,
TREE, SKY

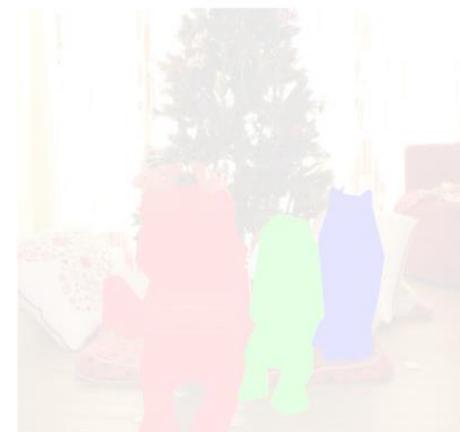
No objects, just pixels

Object
Detection



DOG, DOG, CAT

Instance
Segmentation



DOG, DOG, CAT

Multiple Object

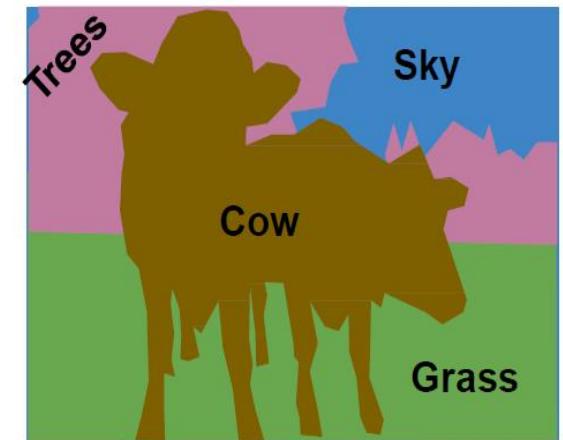
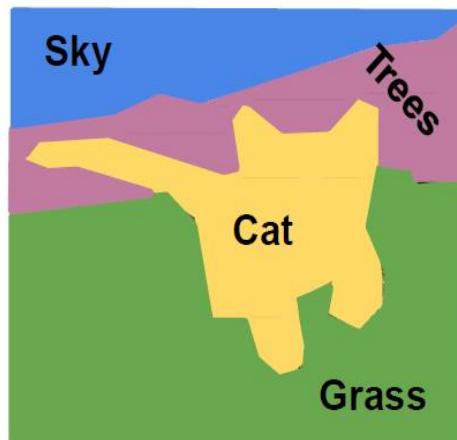
Semantic Segmentation

Label each pixel in the image with a category label

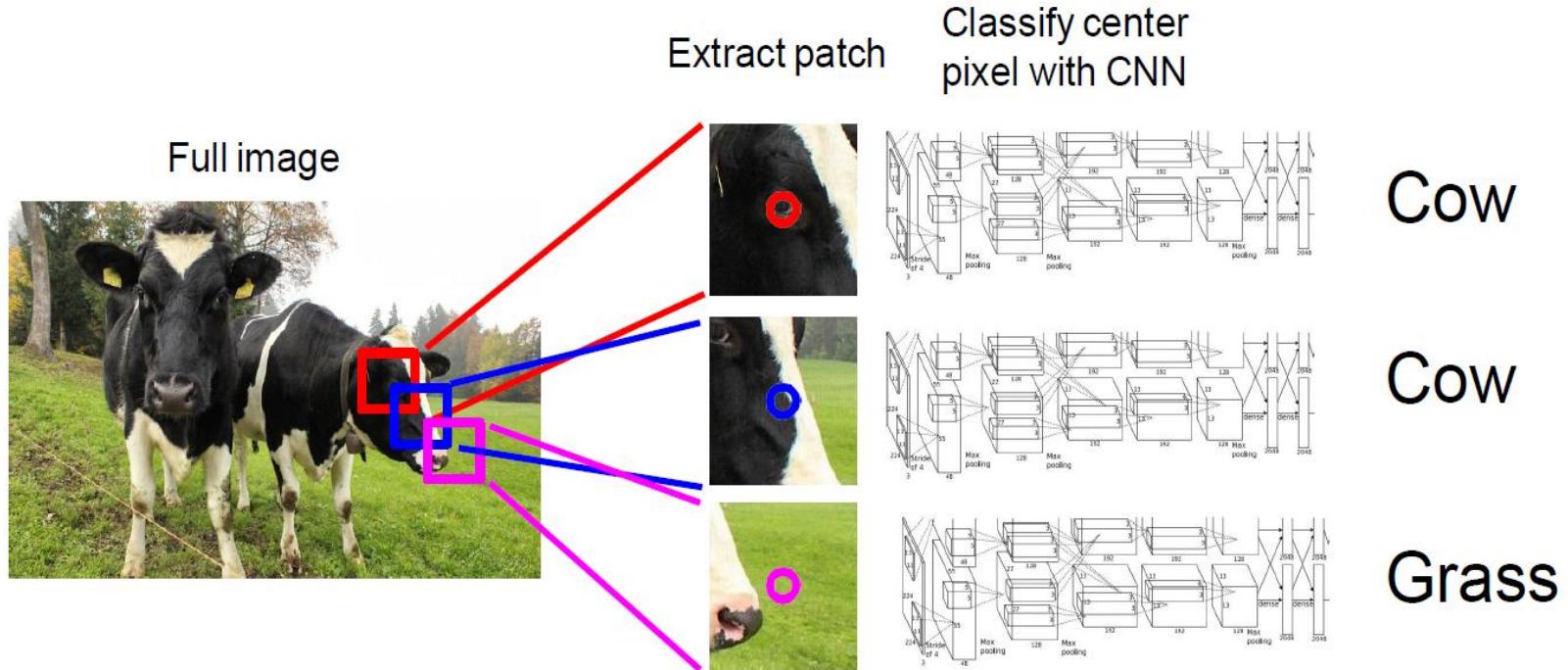
Don't differentiate instances, only care about pixels



[This image is CC0 public domain](#)



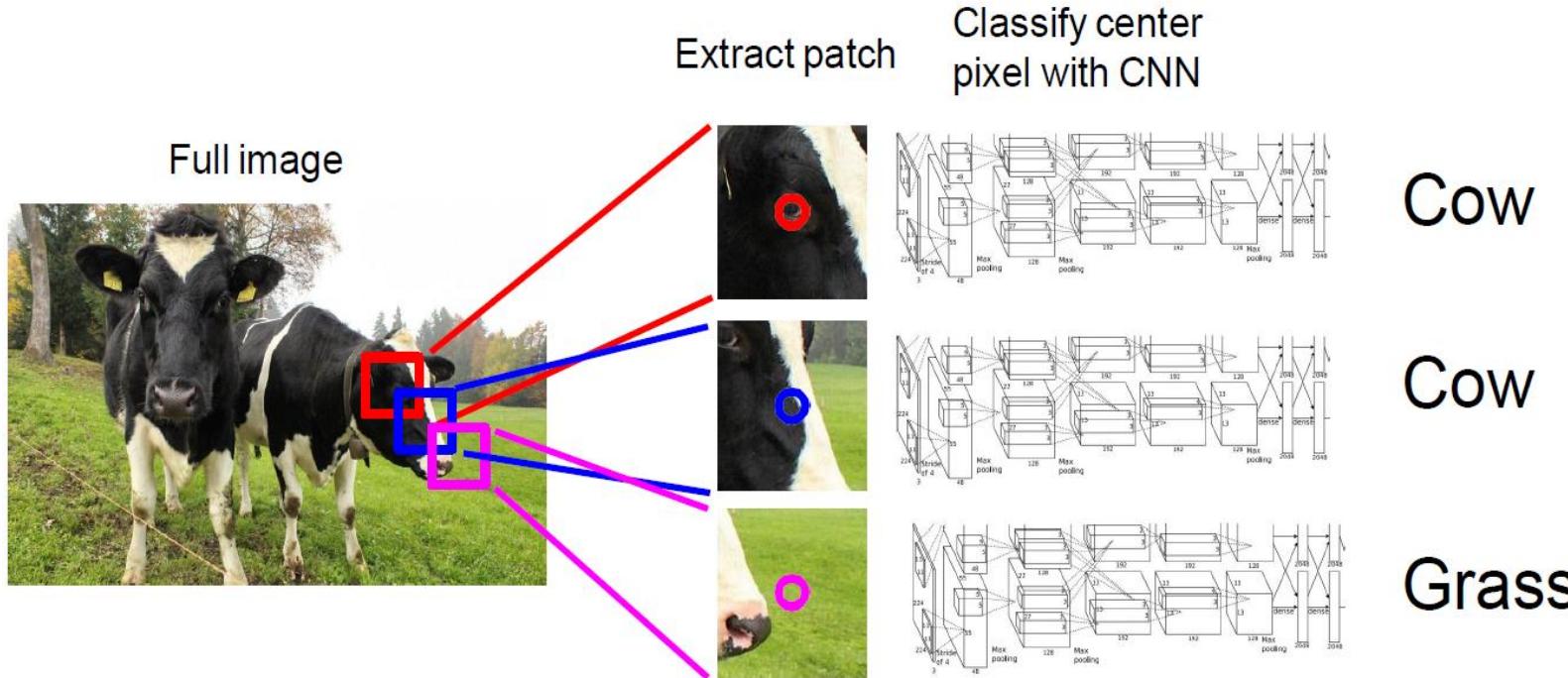
Semantic Segmentation Idea: Sliding Window



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

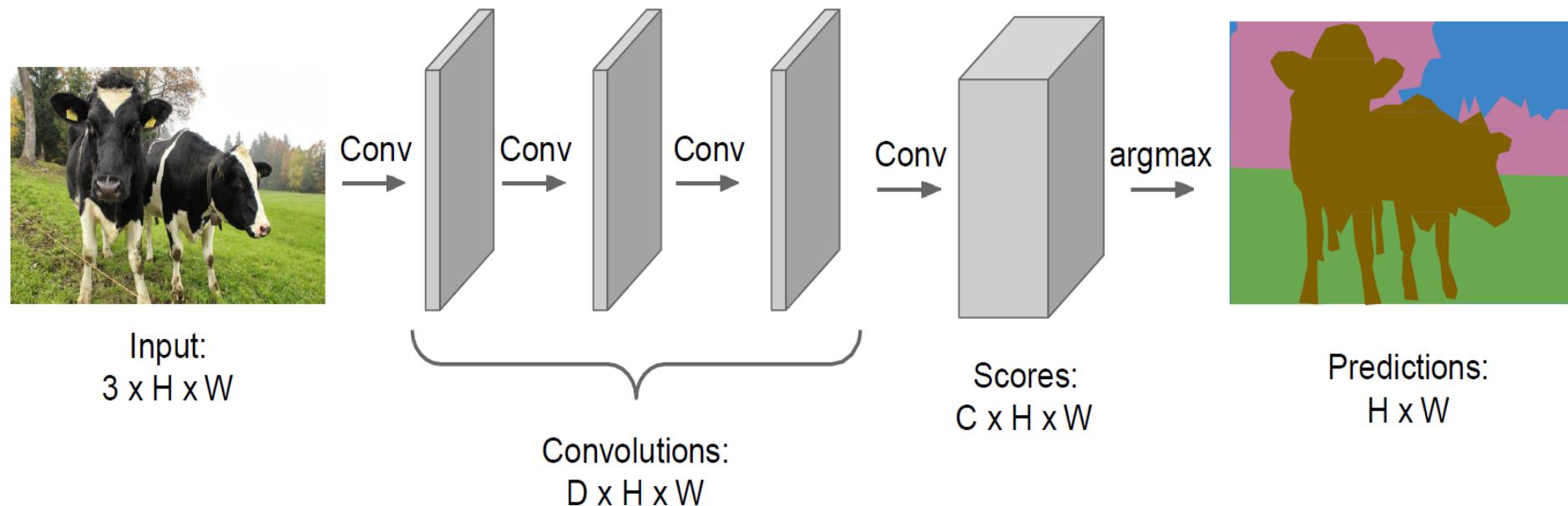


Problem: Very inefficient! Not reusing shared features between overlapping patches

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

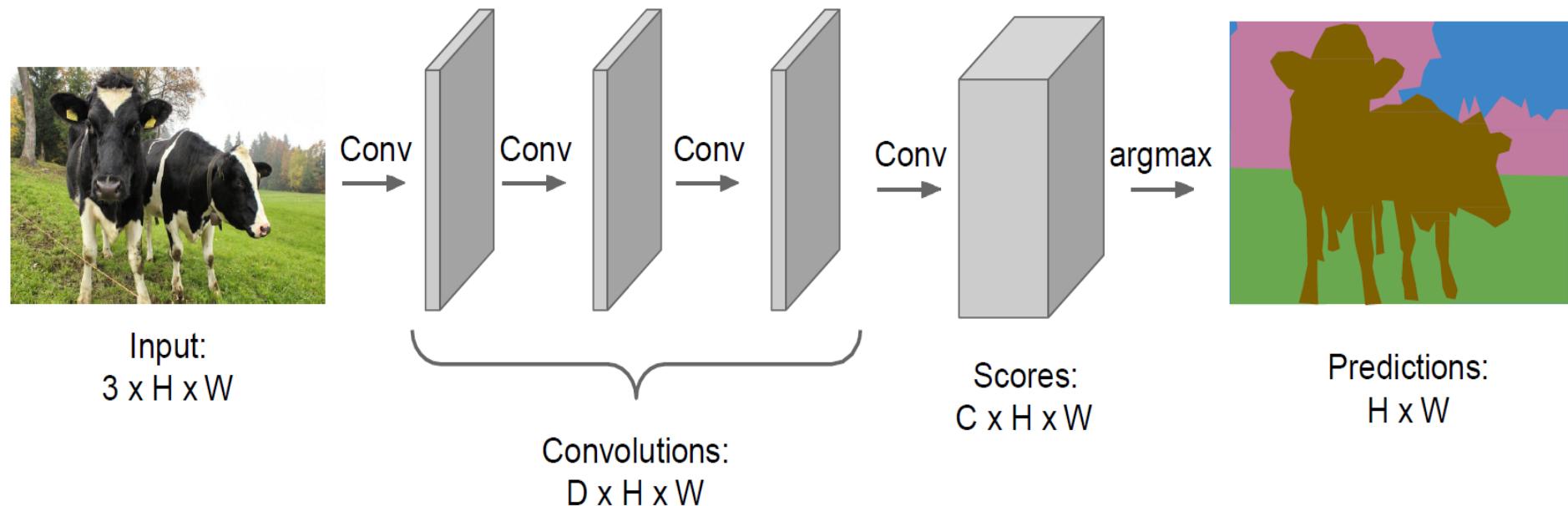
Semantic Segmentation Idea: Fully Convolutional

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!



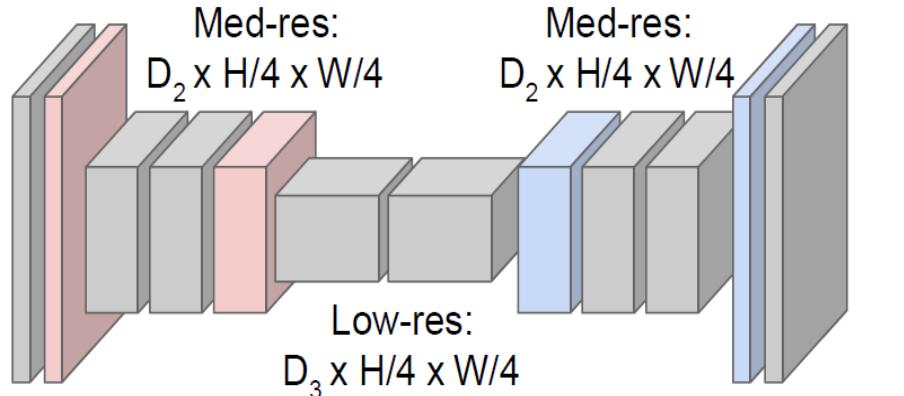
Problem: convolutions at
original image resolution will
be very expensive ...

Semantic Segmentation Idea: Fully Convolutional

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



Input:
 $3 \times H \times W$



High-res:
 $D_1 \times H/2 \times W/2$

High-res:
 $D_1 \times H/2 \times W/2$



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

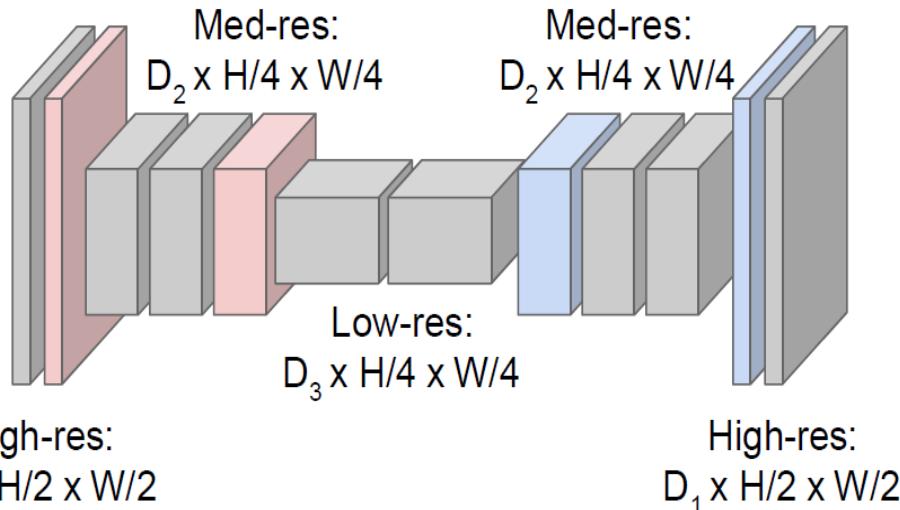
Semantic Segmentation Idea: Fully Convolutional

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:
 $H \times W$

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



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

Max Unpooling

Use positions from pooling layer

1	2
3	4

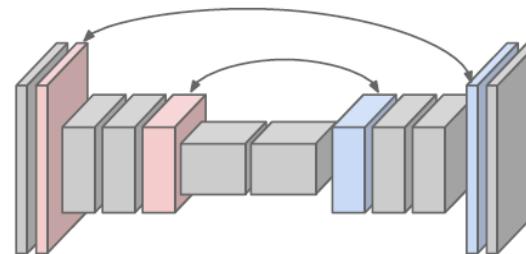
...

Rest of the network

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

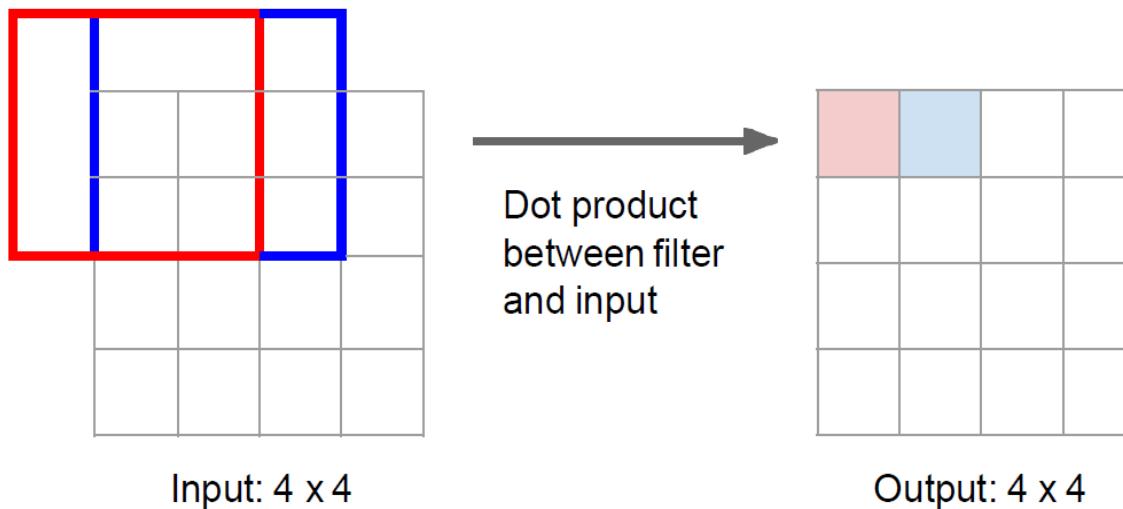
Output: 4 x 4

Corresponding pairs of
downsampling and
upsampling layers



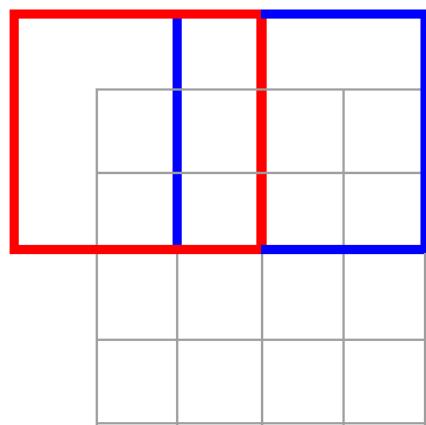
Learnable Upsampling: Transpose Convolution

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



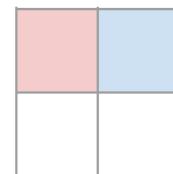
Learnable Upsampling: Transpose Convolution

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



Input: 4×4

Dot product
between filter
and input



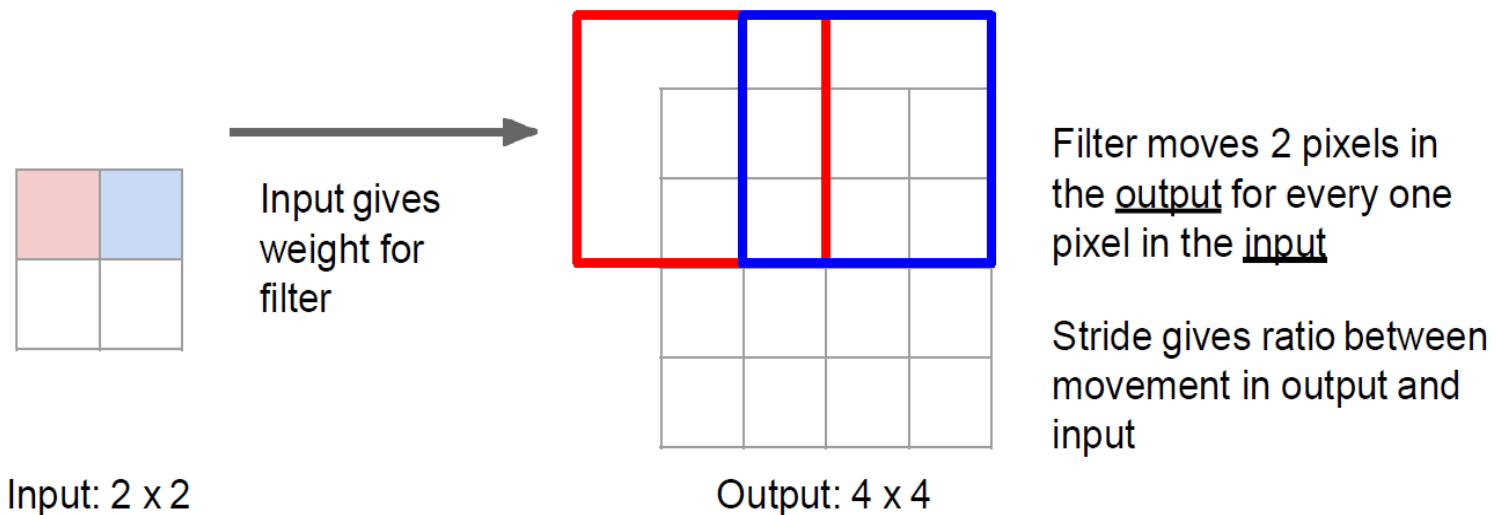
Output: 2×2

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

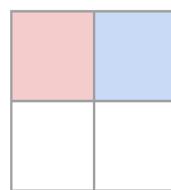
3 x 3 **transpose** convolution, stride 2 pad 1



Learnable Upsampling: Transpose Convolution

Other names:

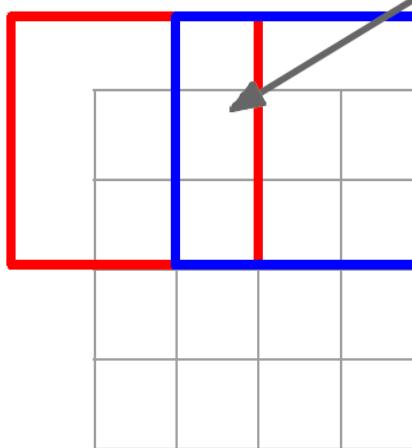
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution



Input: 2 x 2

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

Input gives weight for filter



Output: 4 x 4

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

Stride gives ratio between movement in output and input

Sum where output overlaps

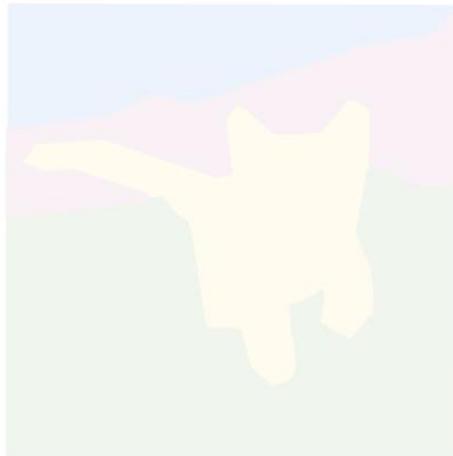
Object Detection

Classification



CAT

Semantic
Segmentation



GRASS, CAT,
TREE, SKY

Object
Detection



DOG, DOG, CAT

Instance
Segmentation



DOG, DOG, CAT

Multiple Object

Object Detection: Impact of Deep Learning

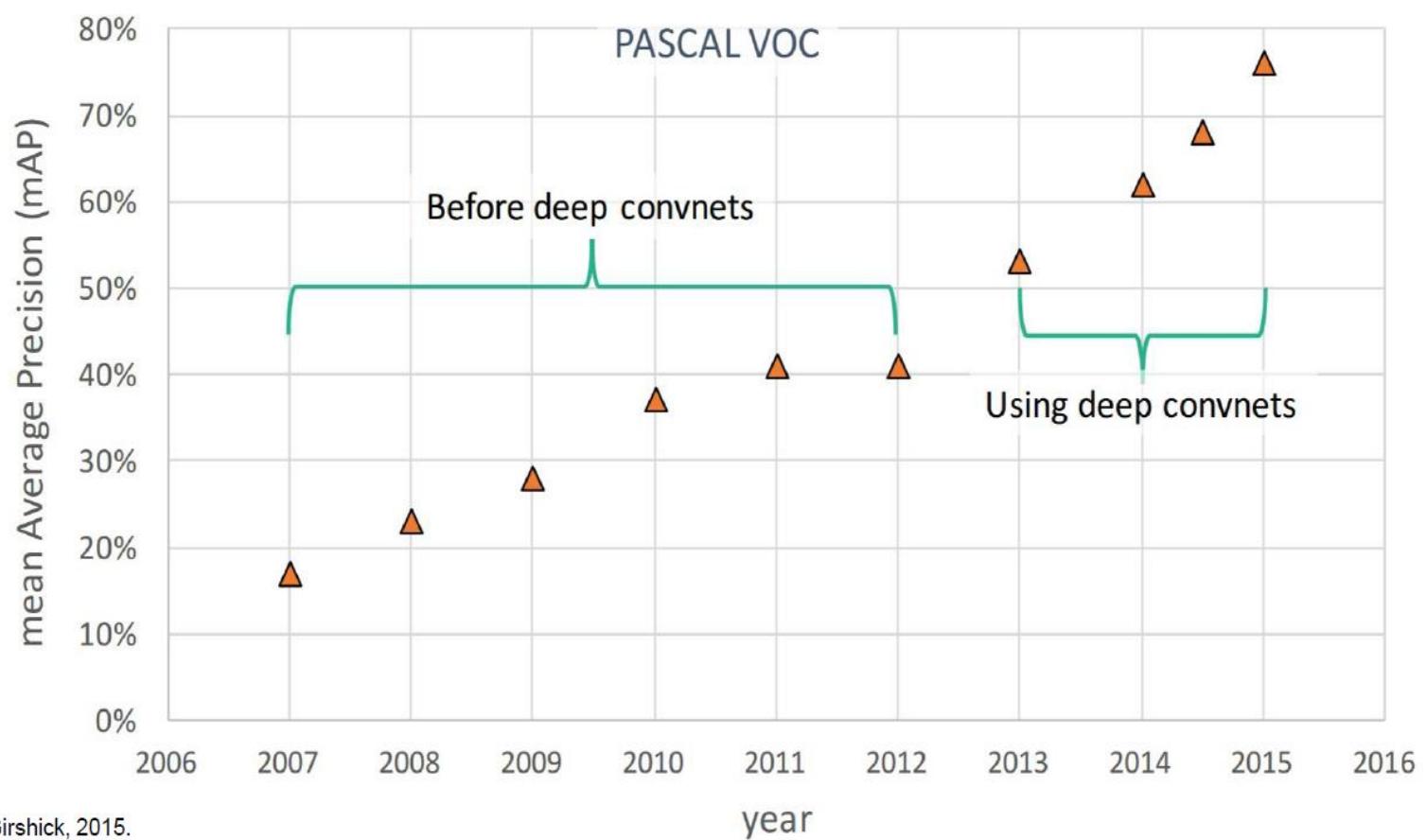


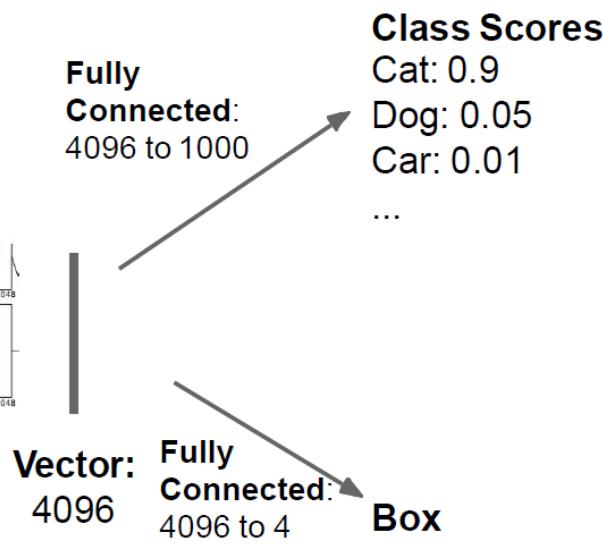
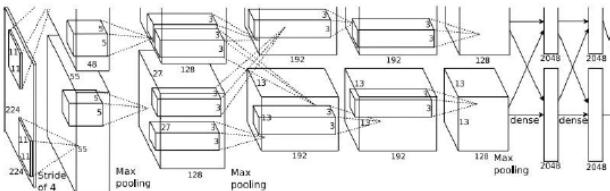
Figure copyright Ross Girshick, 2015.

Reproduced with permission.

Object Detection: Single Object (Classification + Localization)



[This image is CC0 public domain](#)

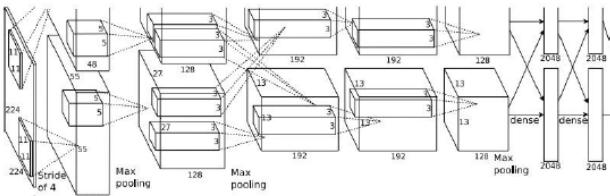


Treat localization as a
regression problem!

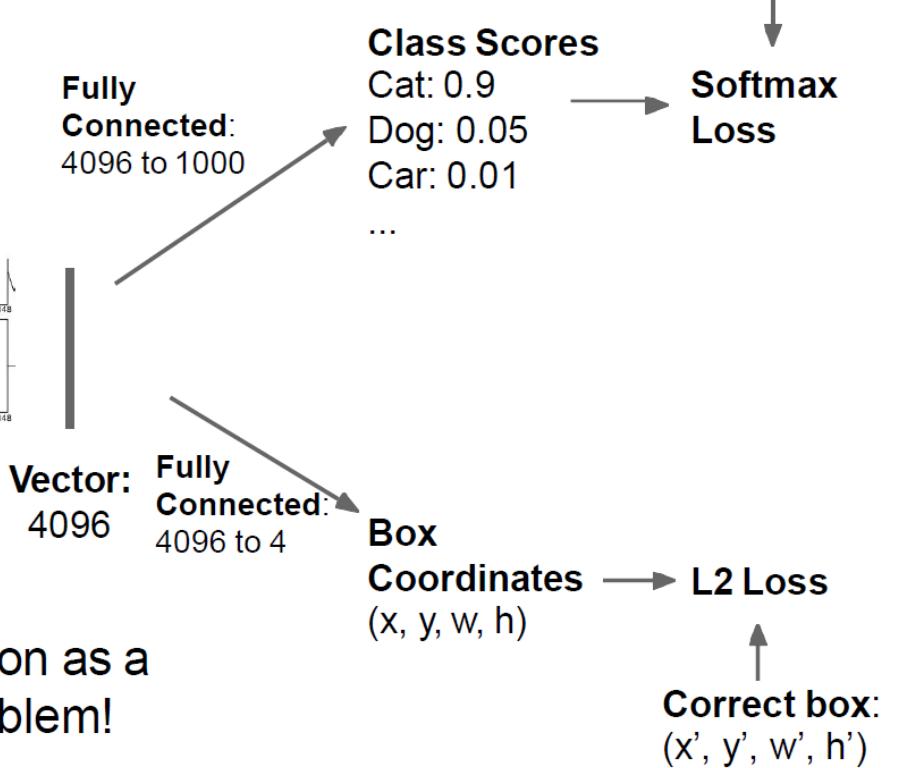
Object Detection: Single Object (Classification + Localization)



This image is CC0 public domain.



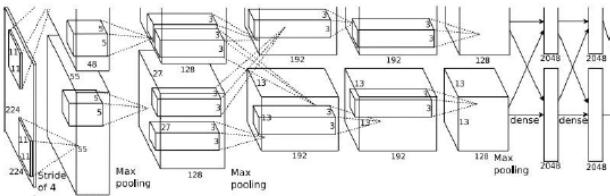
Treat localization as a
regression problem!



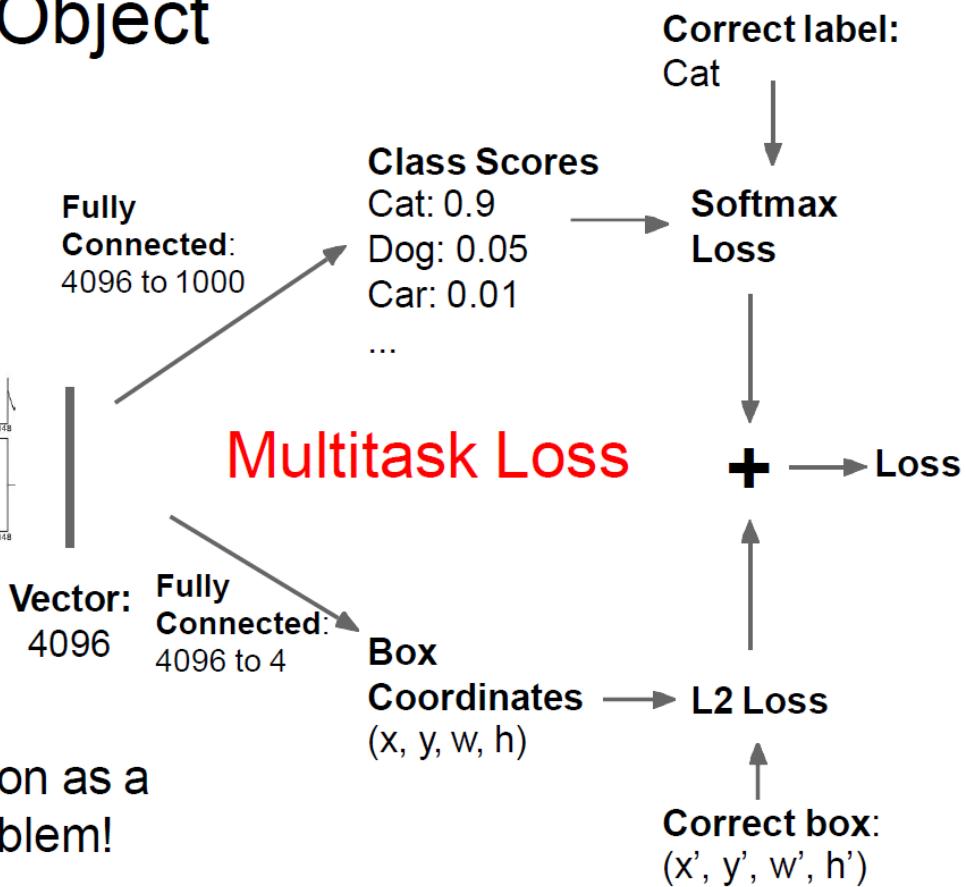
Object Detection: Single Object (Classification + Localization)



This image is CC0 public domain



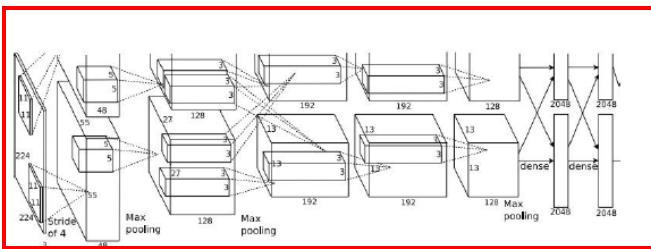
Treat localization as a
regression problem!



Object Detection: Single Object (Classification + Localization)

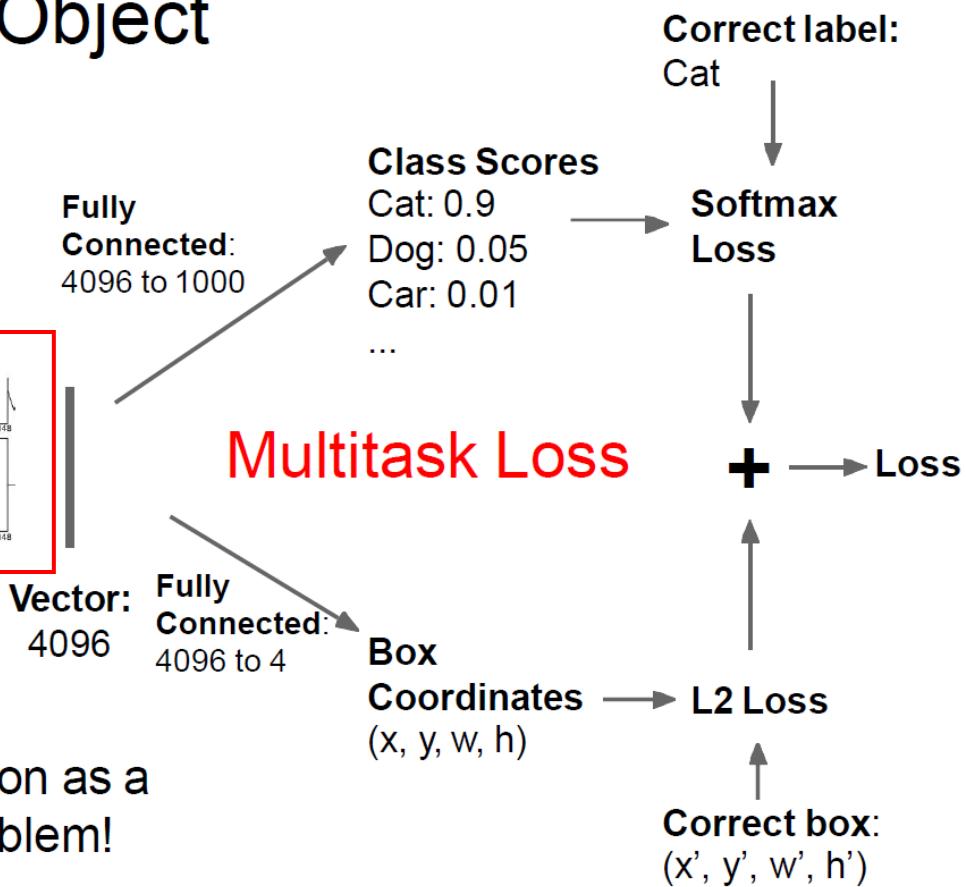


This image is CC0 public domain

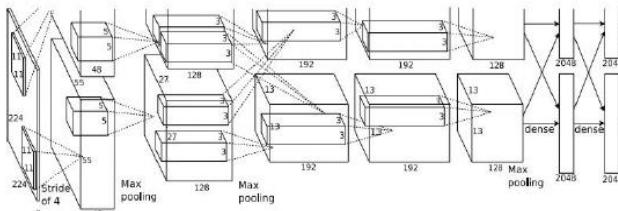


Often pre-trained on ImageNet
(Transfer Learning)

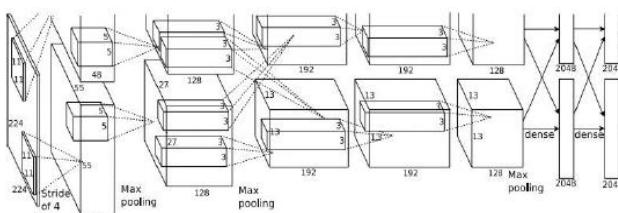
Treat localization as a
regression problem!



Object Detection: Multiple Objects



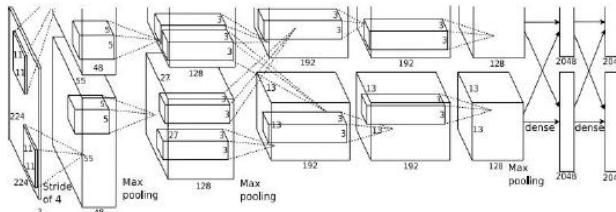
CAT: (x, y, w, h)



DOG: (x, y, w, h)

DOG: (x, y, w, h)

CAT: (x, y, w, h)



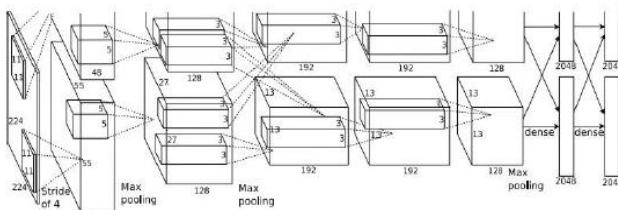
DUCK: (x, y, w, h)

DUCK: (x, y, w, h)

...

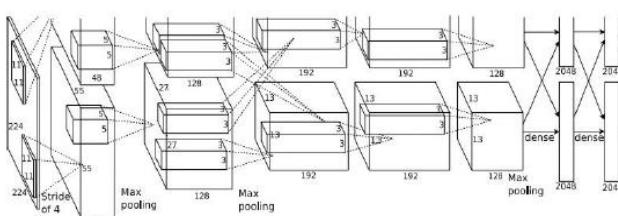
Object Detection: Multiple Objects

Each image needs a different number of outputs!



CAT: (x, y, w, h)

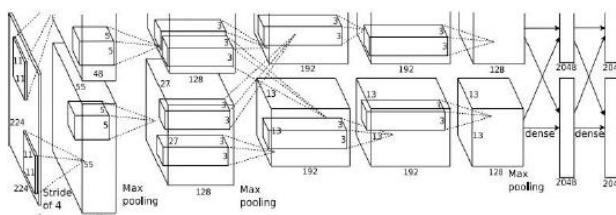
4 numbers



DOG: (x, y, w, h)

16 numbers

CAT: (x, y, w, h)



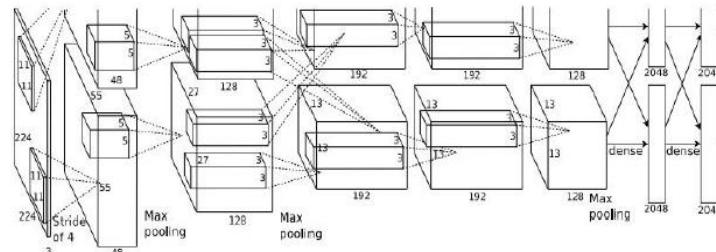
DUCK: (x, y, w, h)

Many numbers!

...

Object Detection: Multiple Objects

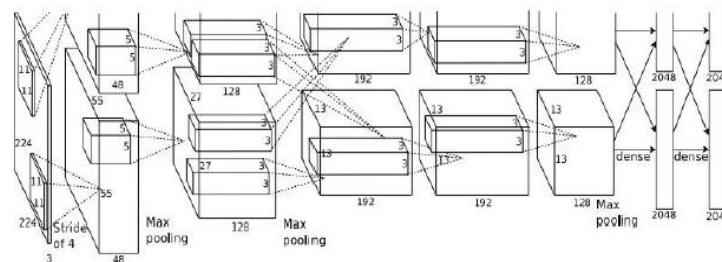
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

Object Detection: Multiple Objects

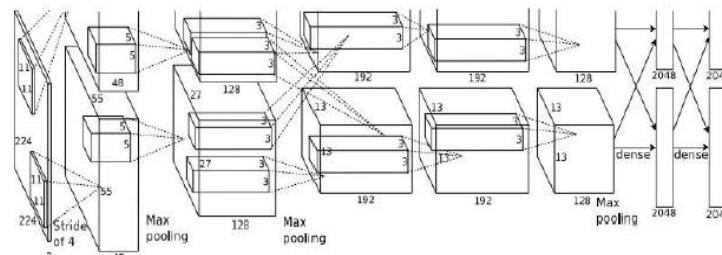
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection: Multiple Objects

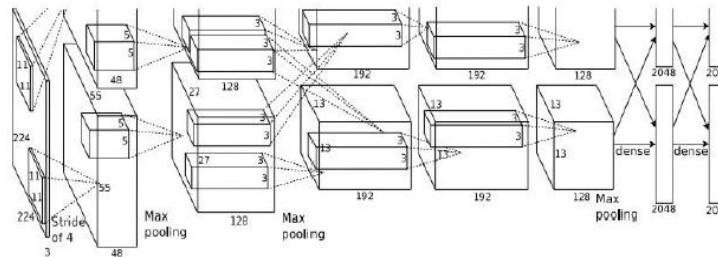
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection: Multiple Objects

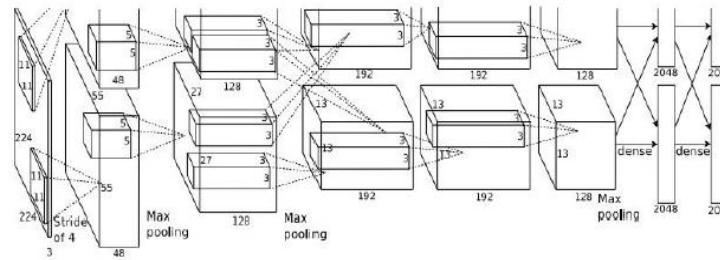
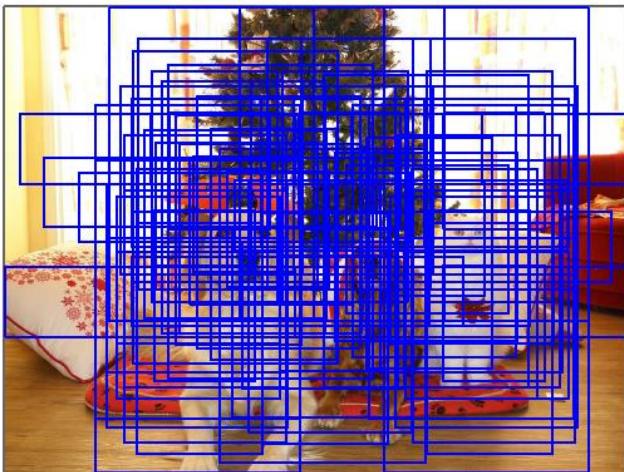
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Object Detection: Multiple Objects

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

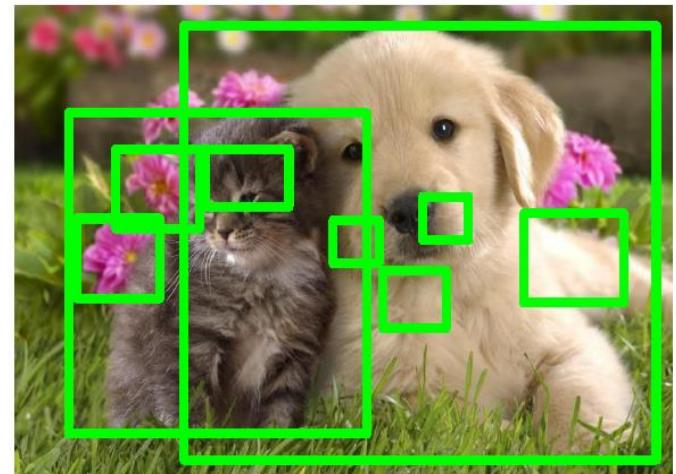


Dog? NO
Cat? YES
Background? NO

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

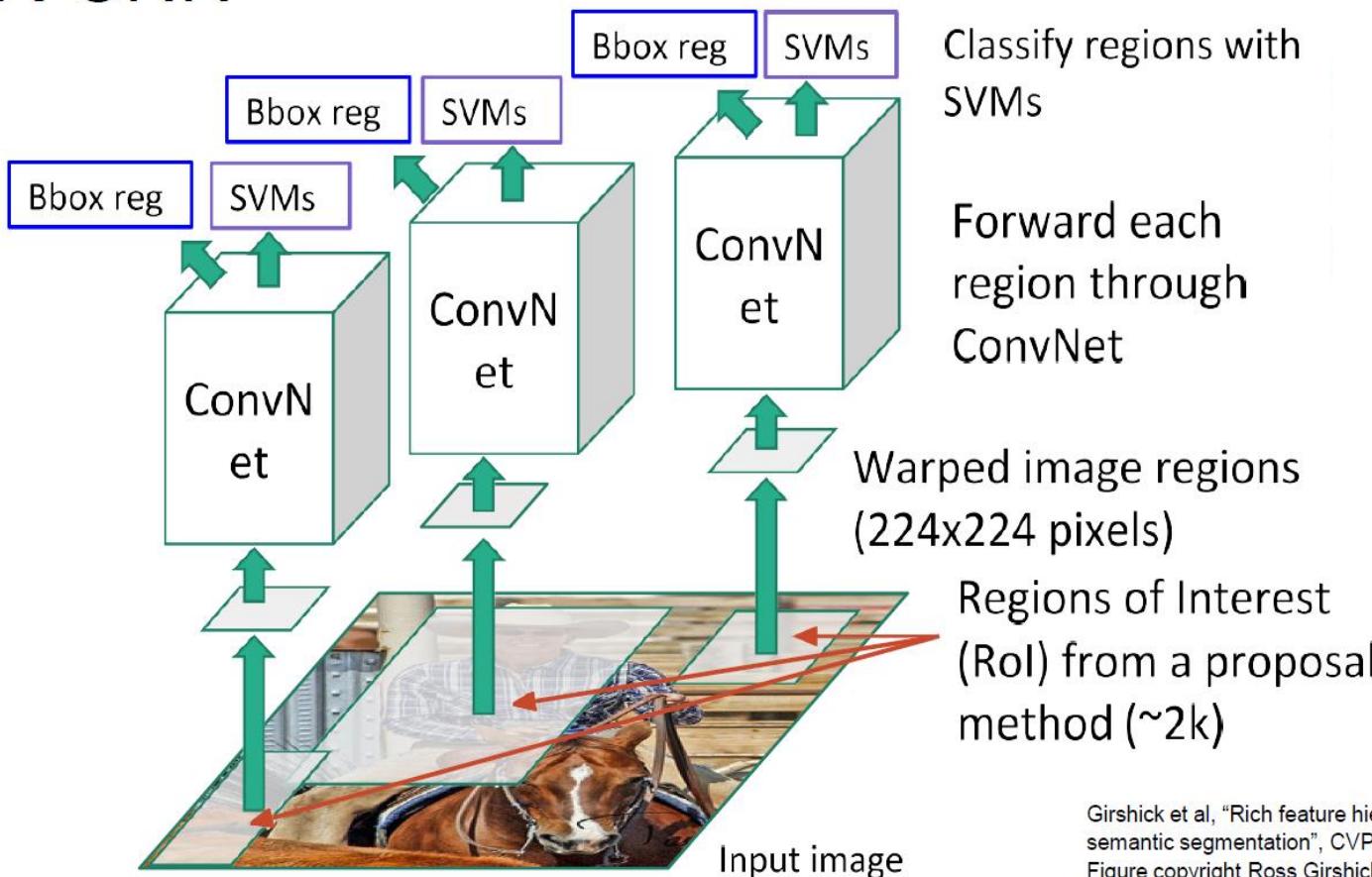
Uijlings et al, "Selective Search for Object Recognition", IJCV 2013

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

Predict “corrections” to the RoI: 4 numbers: (dx , dy , dw , dh)

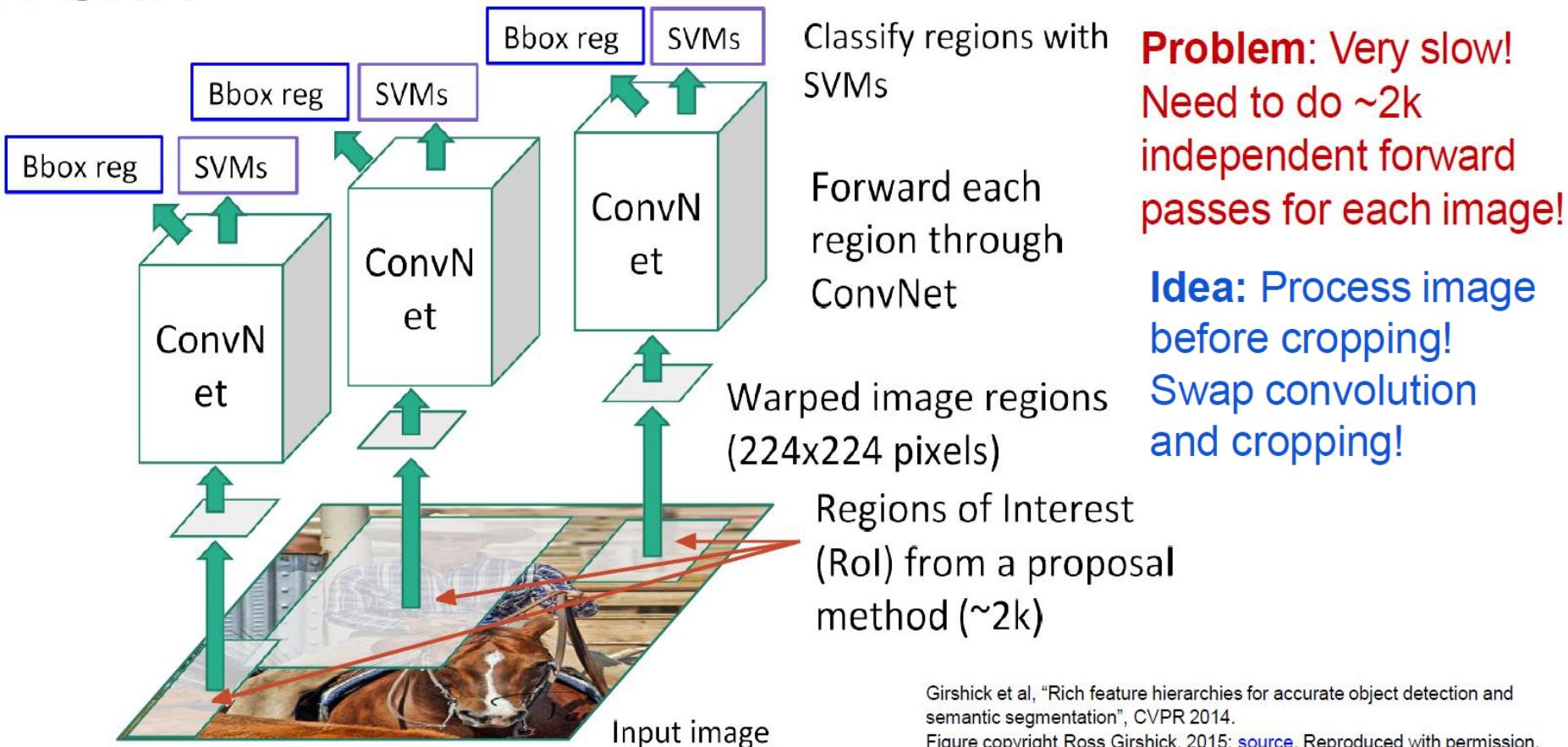


Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

R-CNN

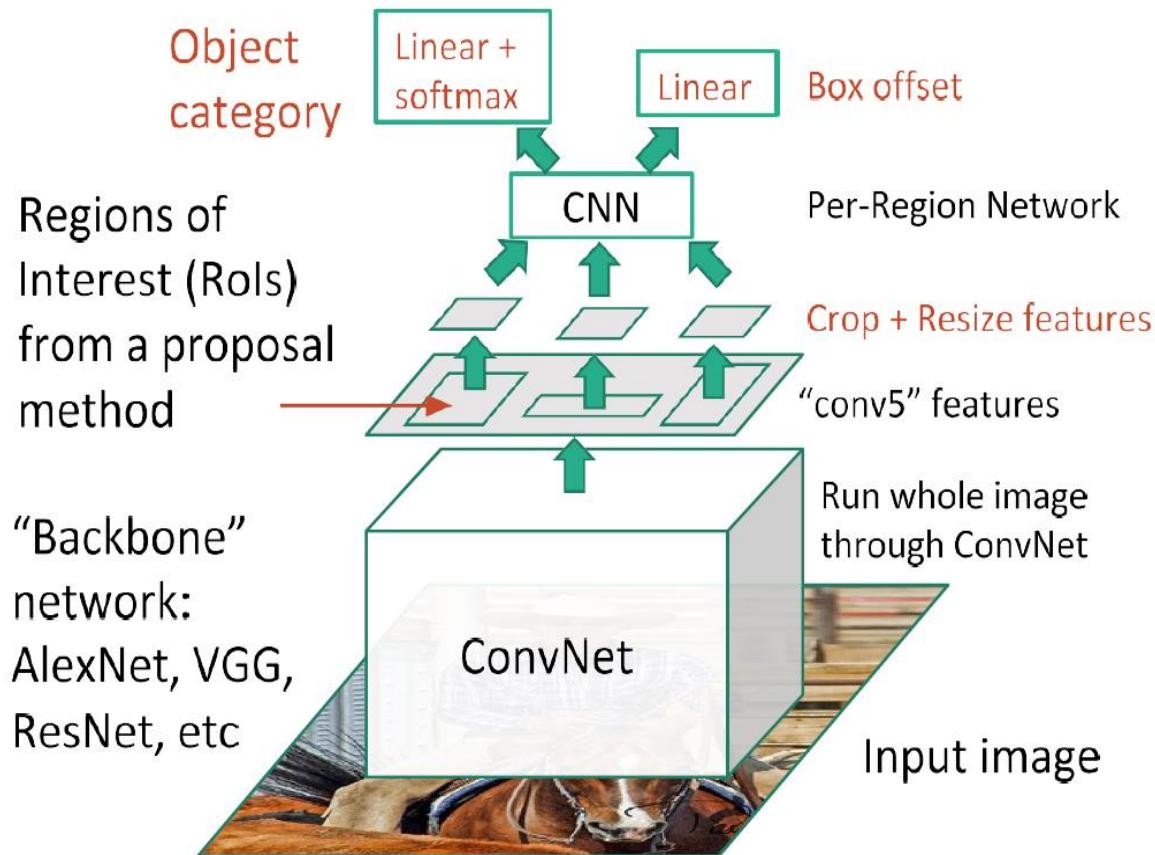
Predict “corrections” to the RoI: 4 numbers: (dx, dy, dw, dh)



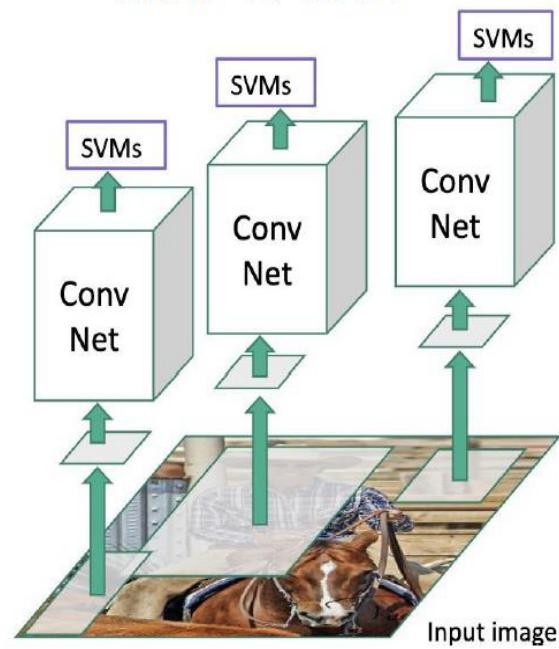
Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fast R-CNN

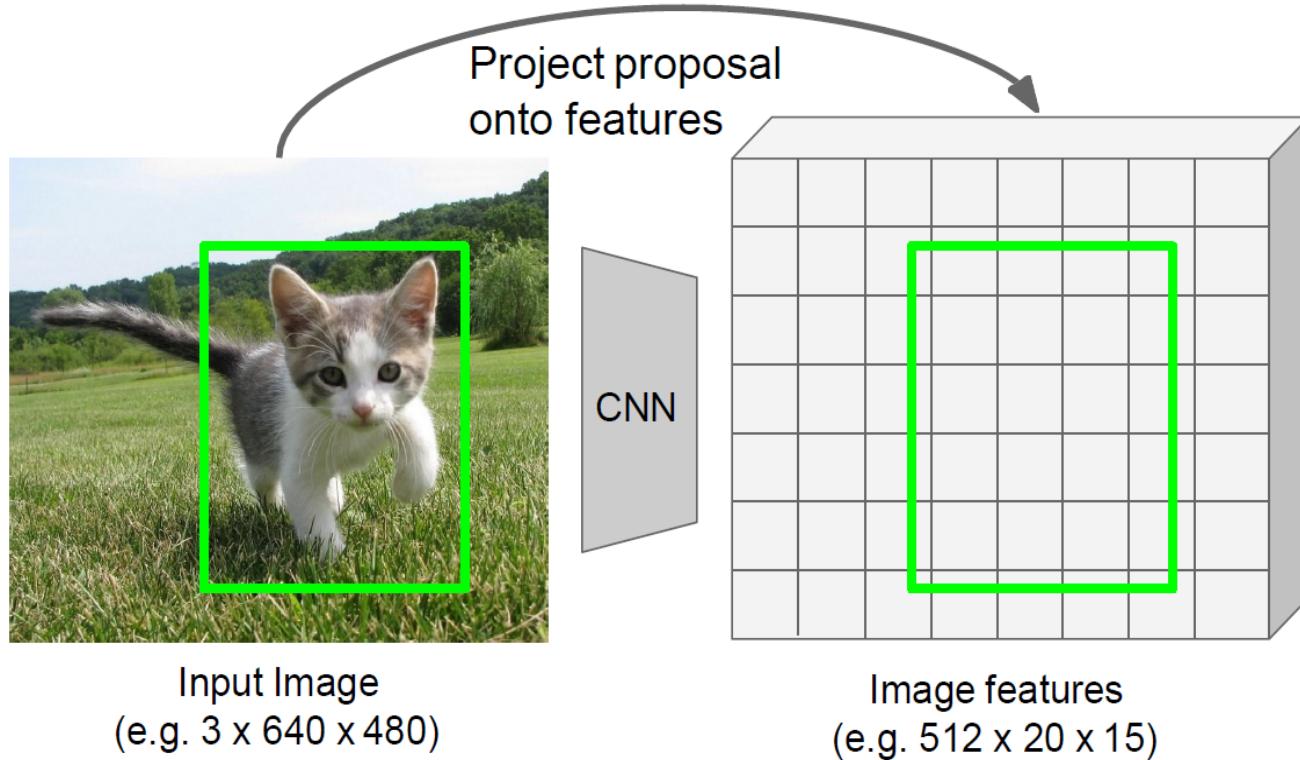


"Slow" R-CNN



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

Cropping Features: RoI Pool



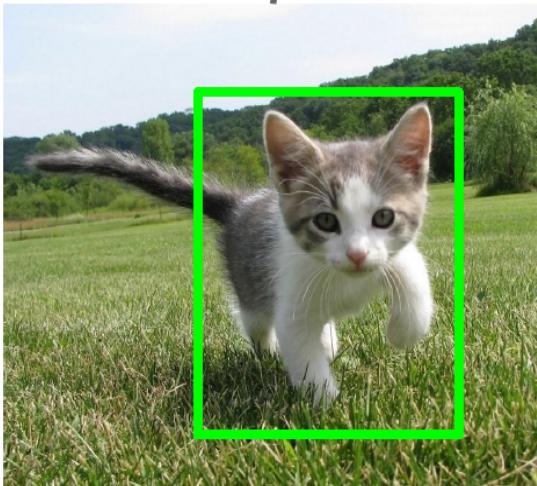
Girshick, "Fast R-CNN", ICCV 2015.

Girshick, "Fast R-CNN", ICCV 2015.

Cropping Features: RoI Pool

“Snap” to grid cells

Project proposal
onto features



Input Image
(e.g. $3 \times 640 \times 480$)

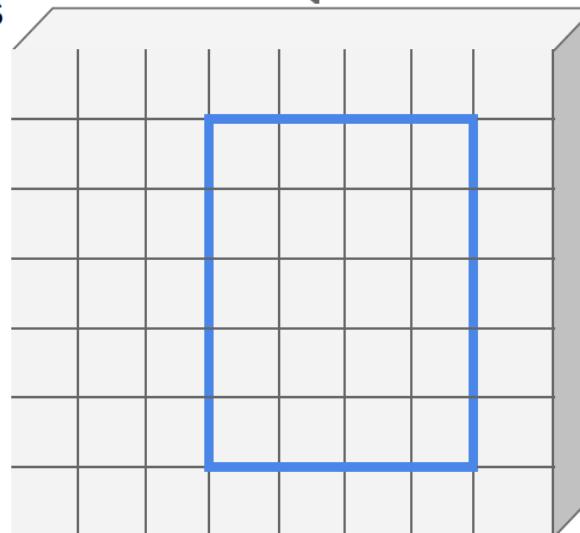
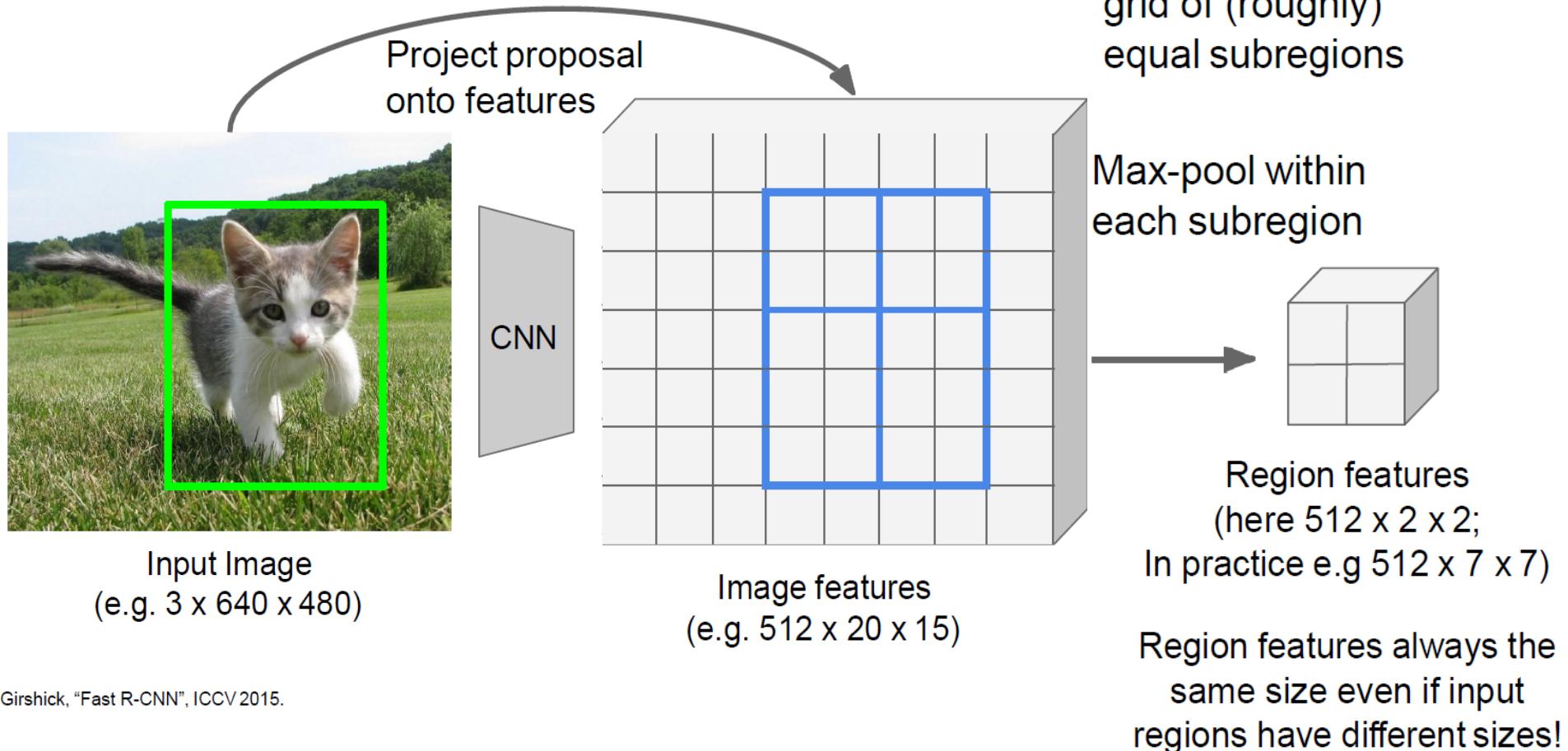


Image features
(e.g. $512 \times 20 \times 15$)

Girshick, “Fast R-CNN”, ICCV 2015.

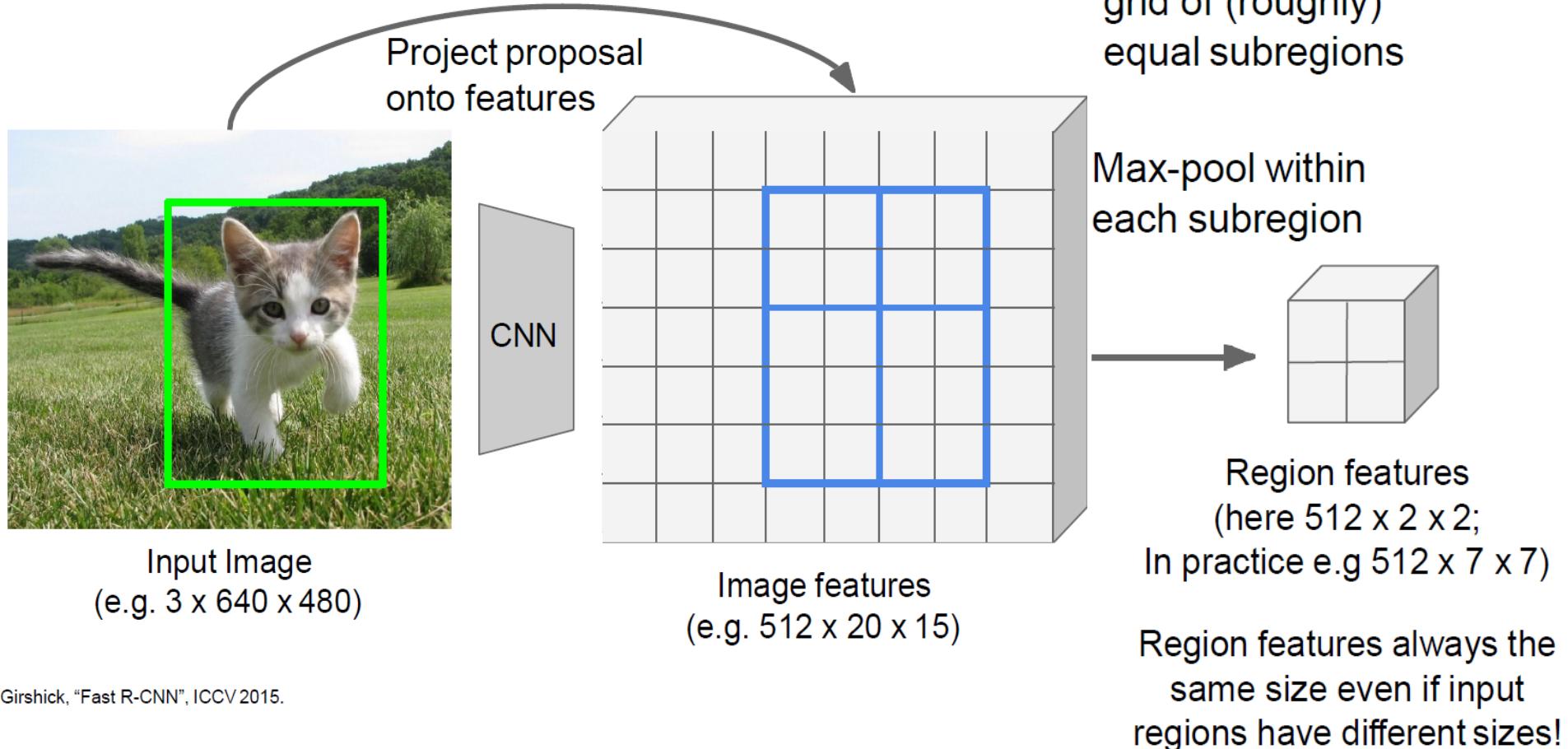
Girshick, “Fast R-CNN”, ICCV 2015.

Cropping Features: RoI Pool



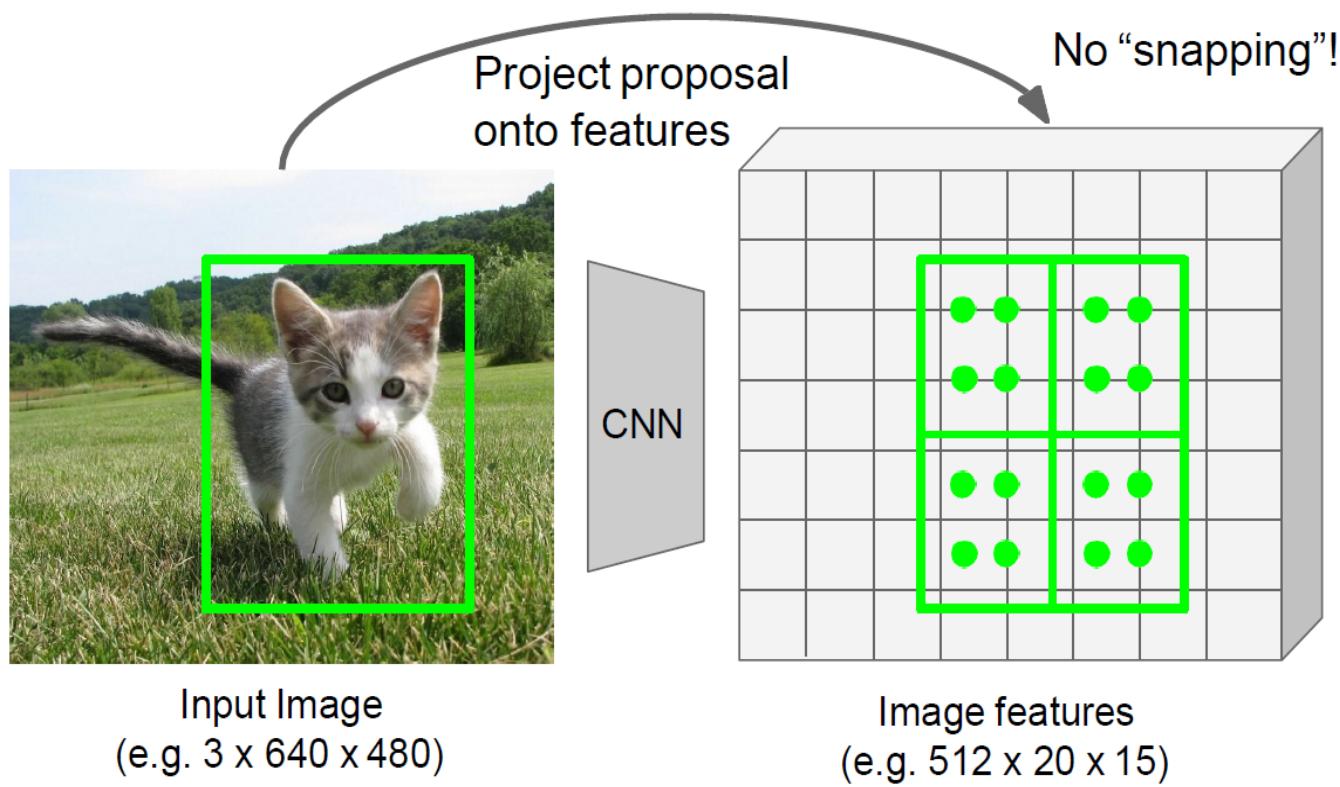
Girshick, "Fast R-CNN", ICCV 2015.

Cropping Features: RoI Pool

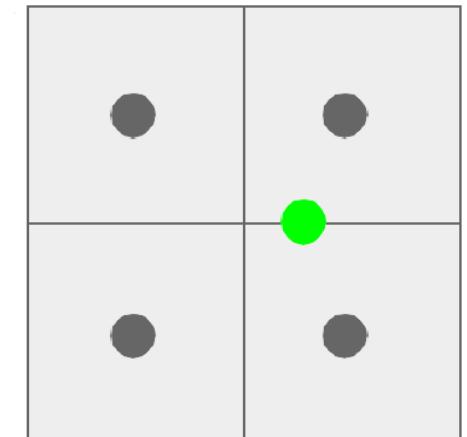


Problem: Region features slightly misaligned

Cropping Features: RoI Align

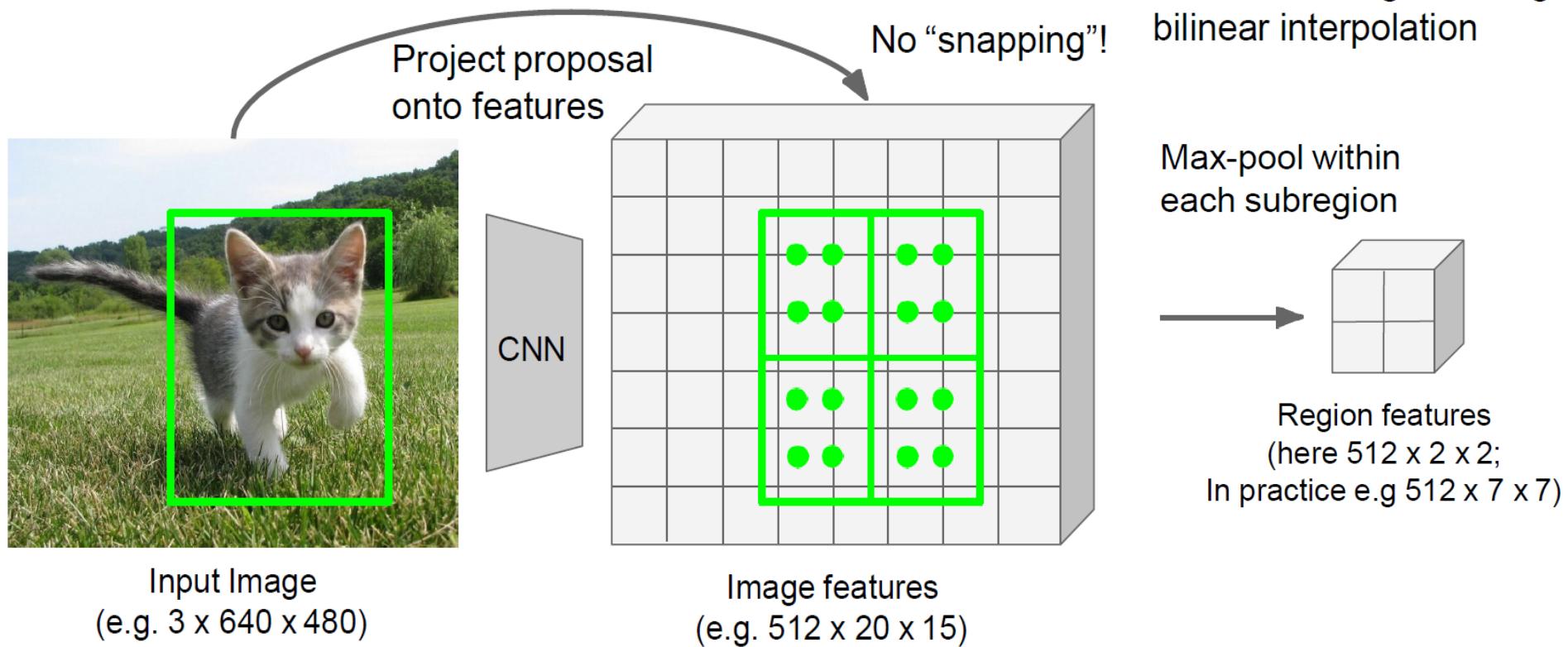


Sample at regular points
in each subregion using
bilinear interpolation



Feature f_{xy} for point (x, y)
is a linear combination of
features at its four
neighboring grid cells:

Cropping Features: RoI Align

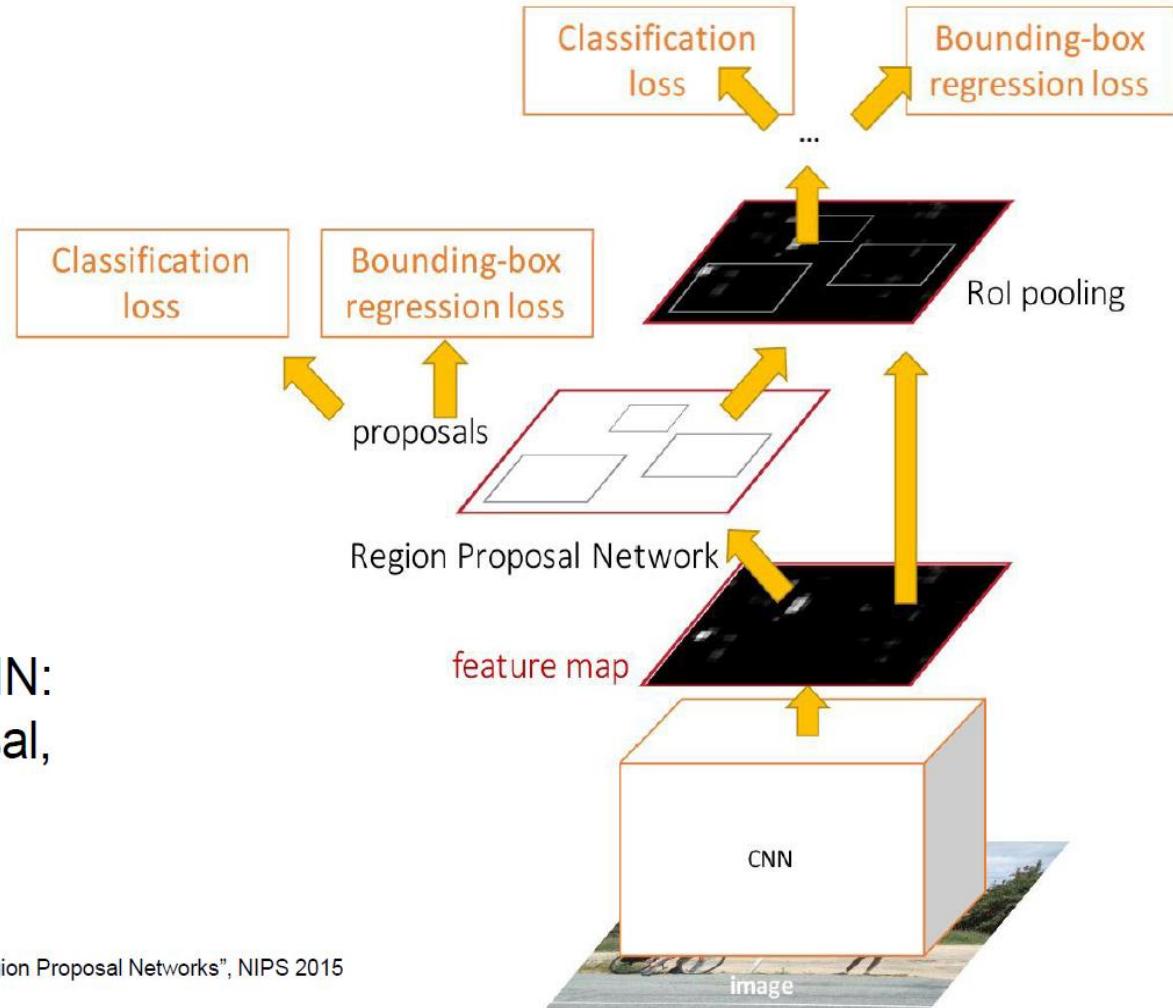


Faster R-CNN:

Make CNN do proposals!

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

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one

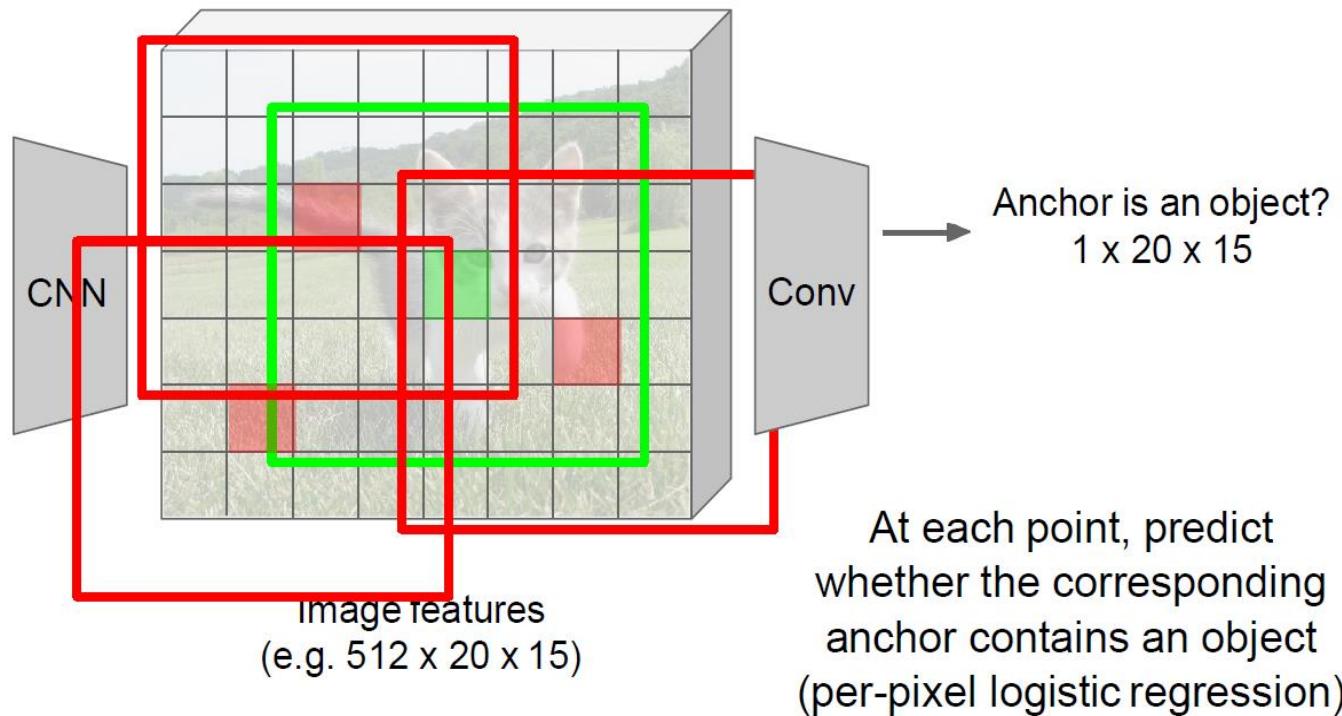


Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

Region Proposal Network



Input Image
(e.g. $3 \times 640 \times 480$)



Imagine an **anchor box** of fixed size at each point in the feature map

Region Proposal Network



Input Image
(e.g. $3 \times 640 \times 480$)

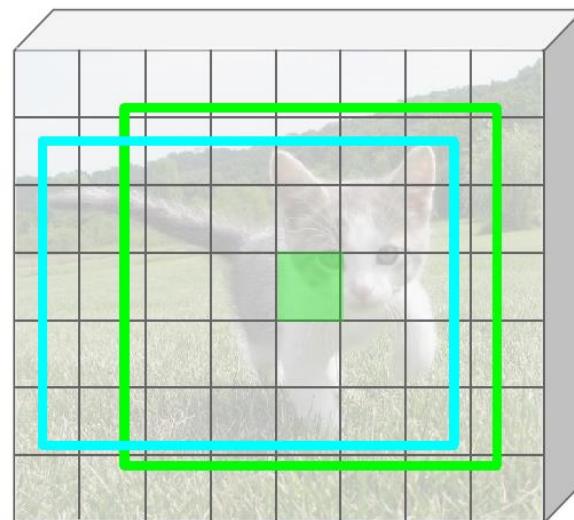


Image features
(e.g. $512 \times 20 \times 15$)



Imagine an **anchor box** of fixed size at each point in the feature map

Anchor is an object?
 $1 \times 20 \times 15$

Box transforms
 $4 \times 20 \times 15$

For positive boxes, also predict a transformation from the anchor to the ground-truth box (regress 4 numbers per pixel)

Faster R-CNN:

Make CNN do proposals!

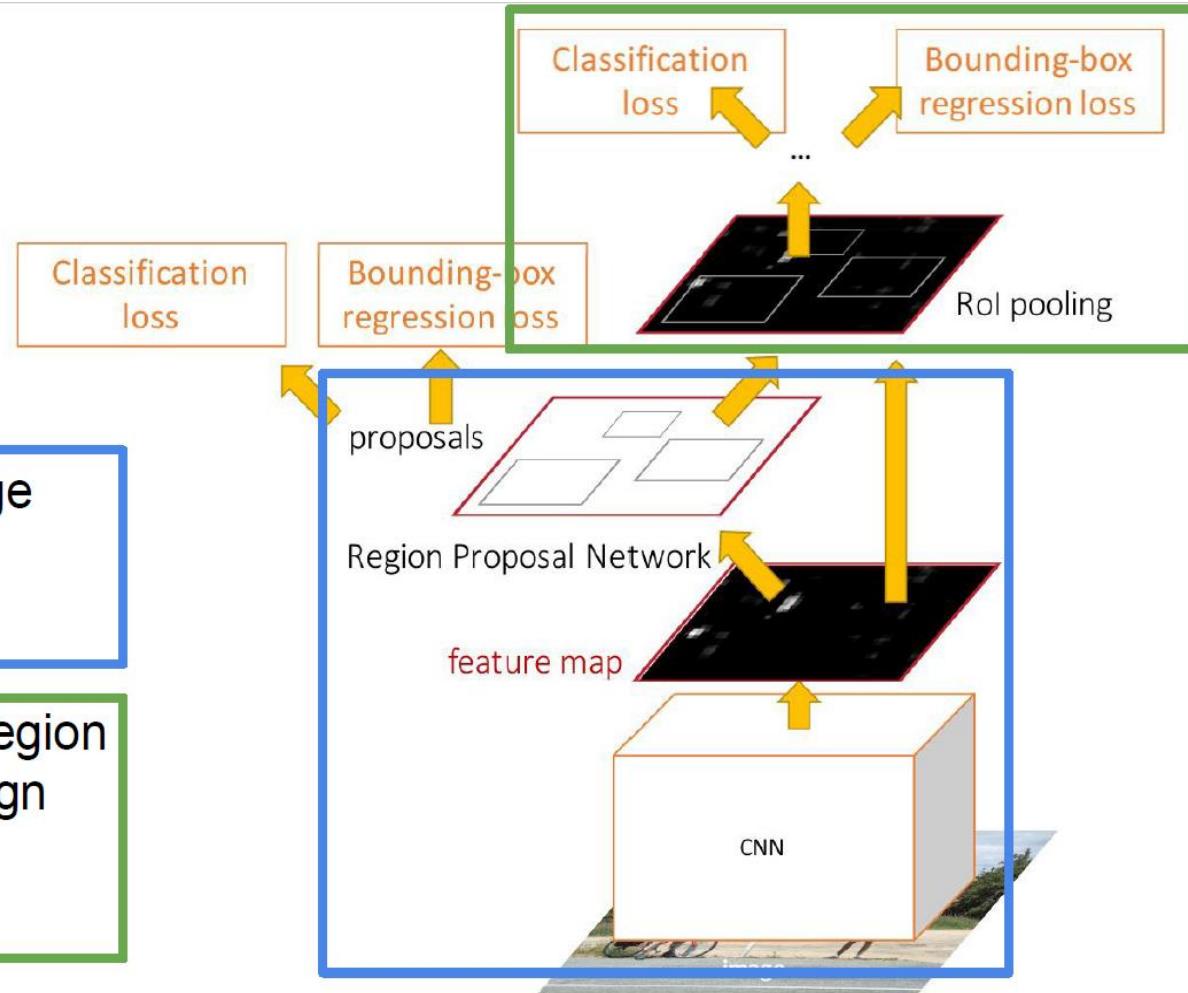
Faster R-CNN is a
Two-stage object detector

First stage: Run once per image

- Backbone network
- Region proposal network

Second stage: Run once per region

- Crop features: RoI pool / align
- Predict object class
- Prediction bbox offset



Instance Segmentation

Classification



CAT

Semantic
Segmentation



GRASS, CAT,
TREE, SKY

Object
Detection



DOG, DOG, DOG

Instance
Segmentation

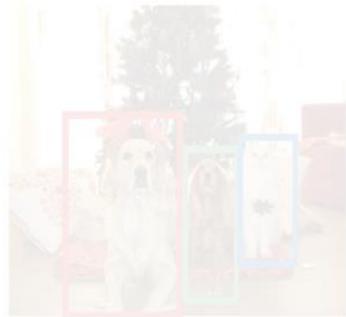


DOG, DOG, DOG, CAT

Multiple Object

Instance Segmentation: Mask R-CNN

Object
Detection

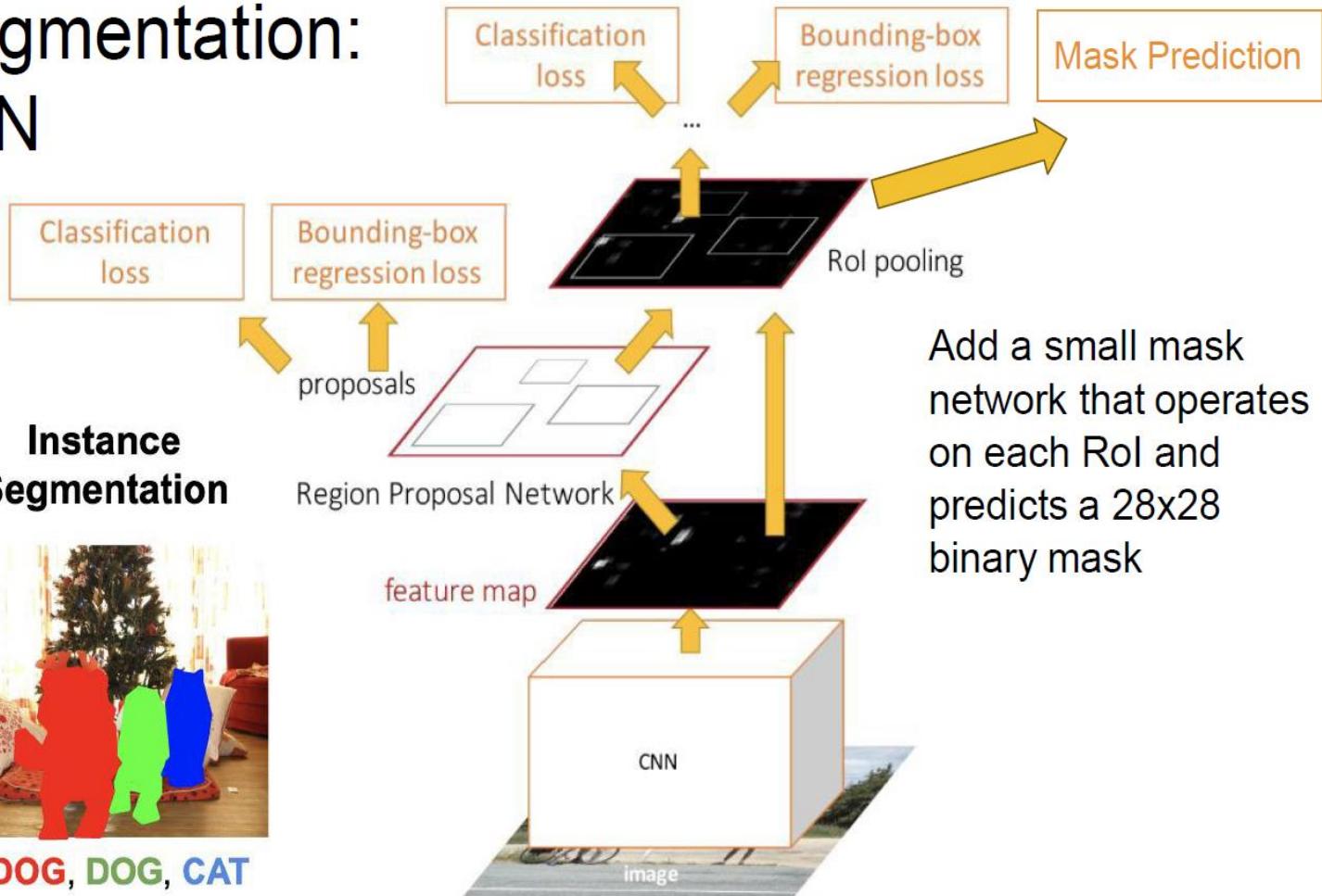


DOG, DOG, CAT

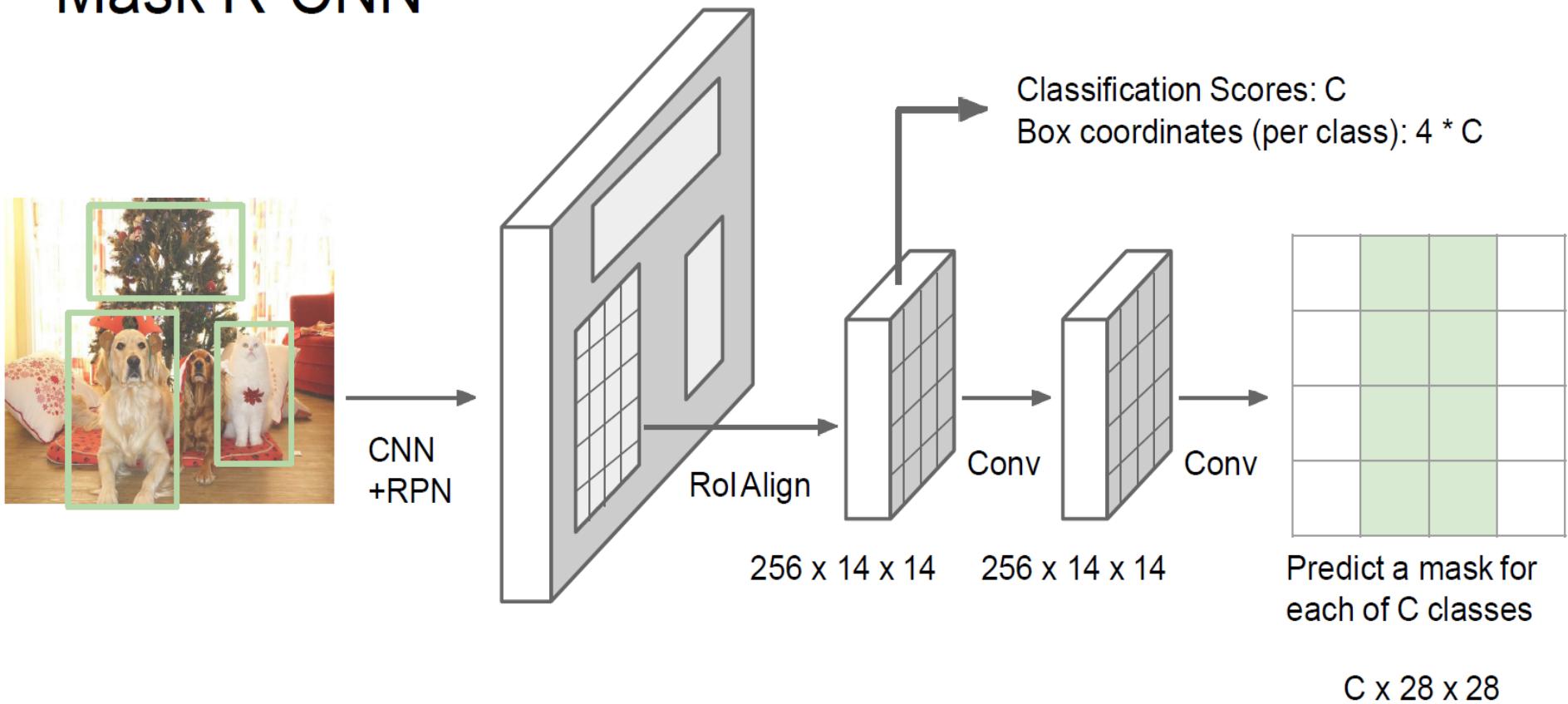


DOG, DOG, CAT

Instance Segmentation

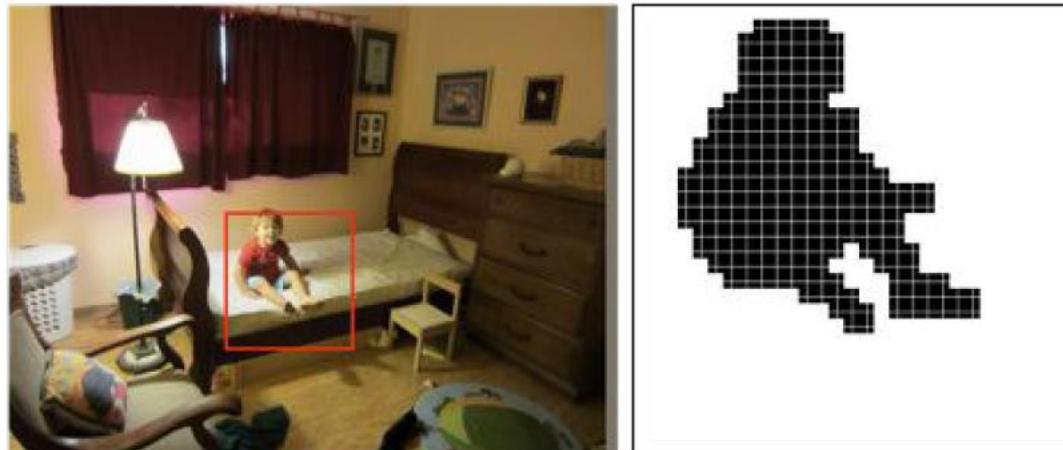


Mask R-CNN

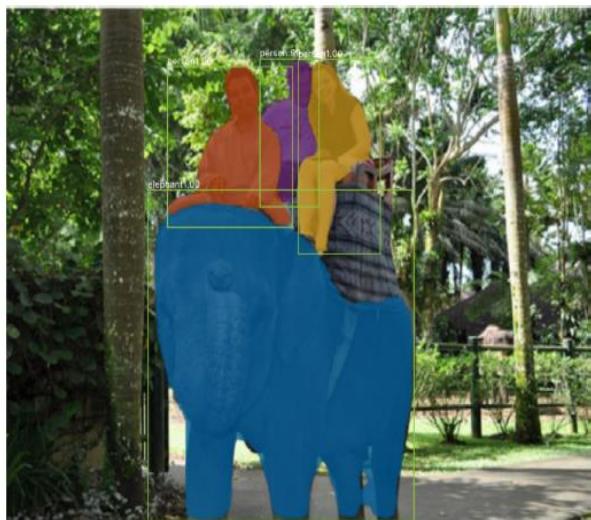


He et al, "Mask R-CNN", arXiv 2017

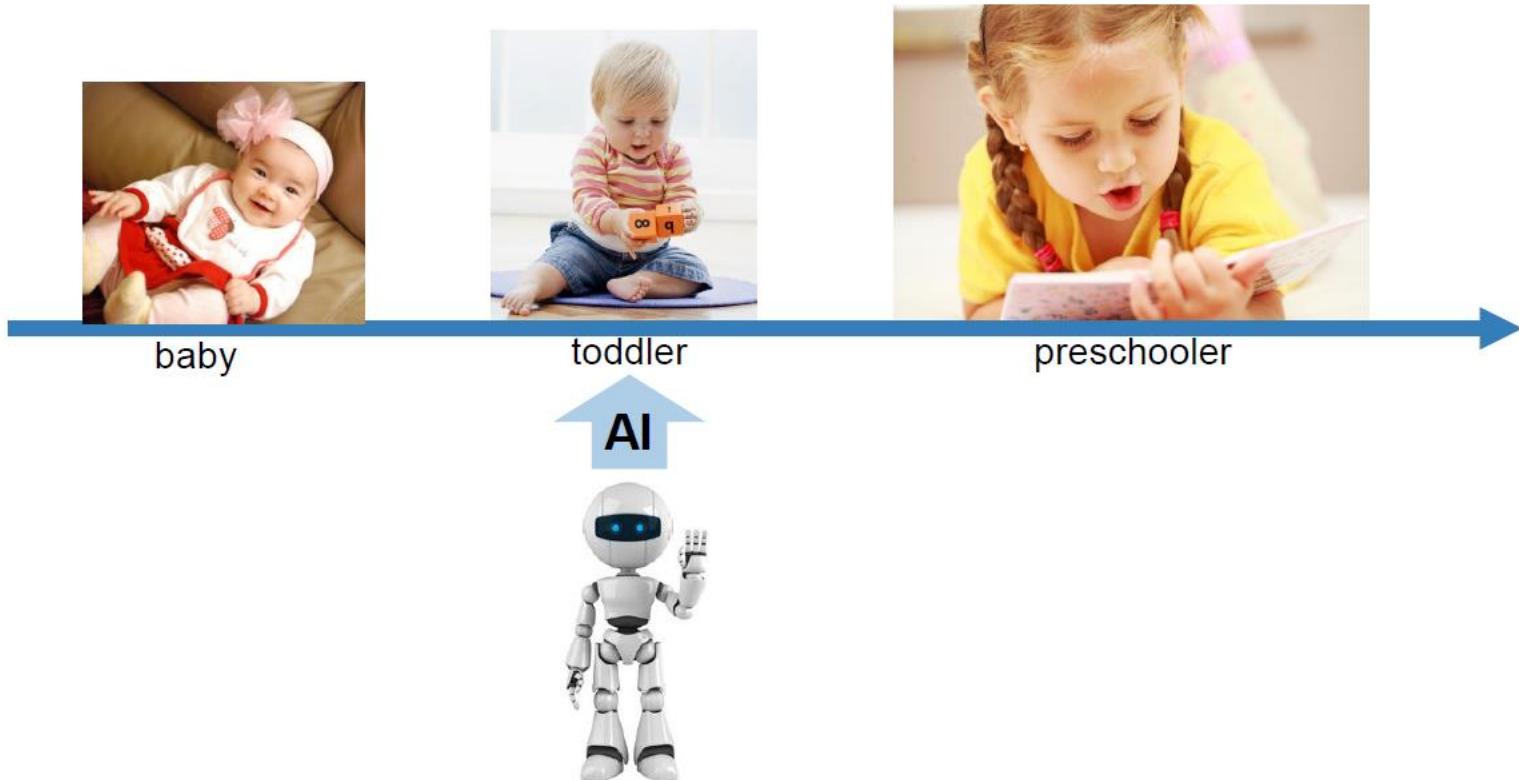
Example



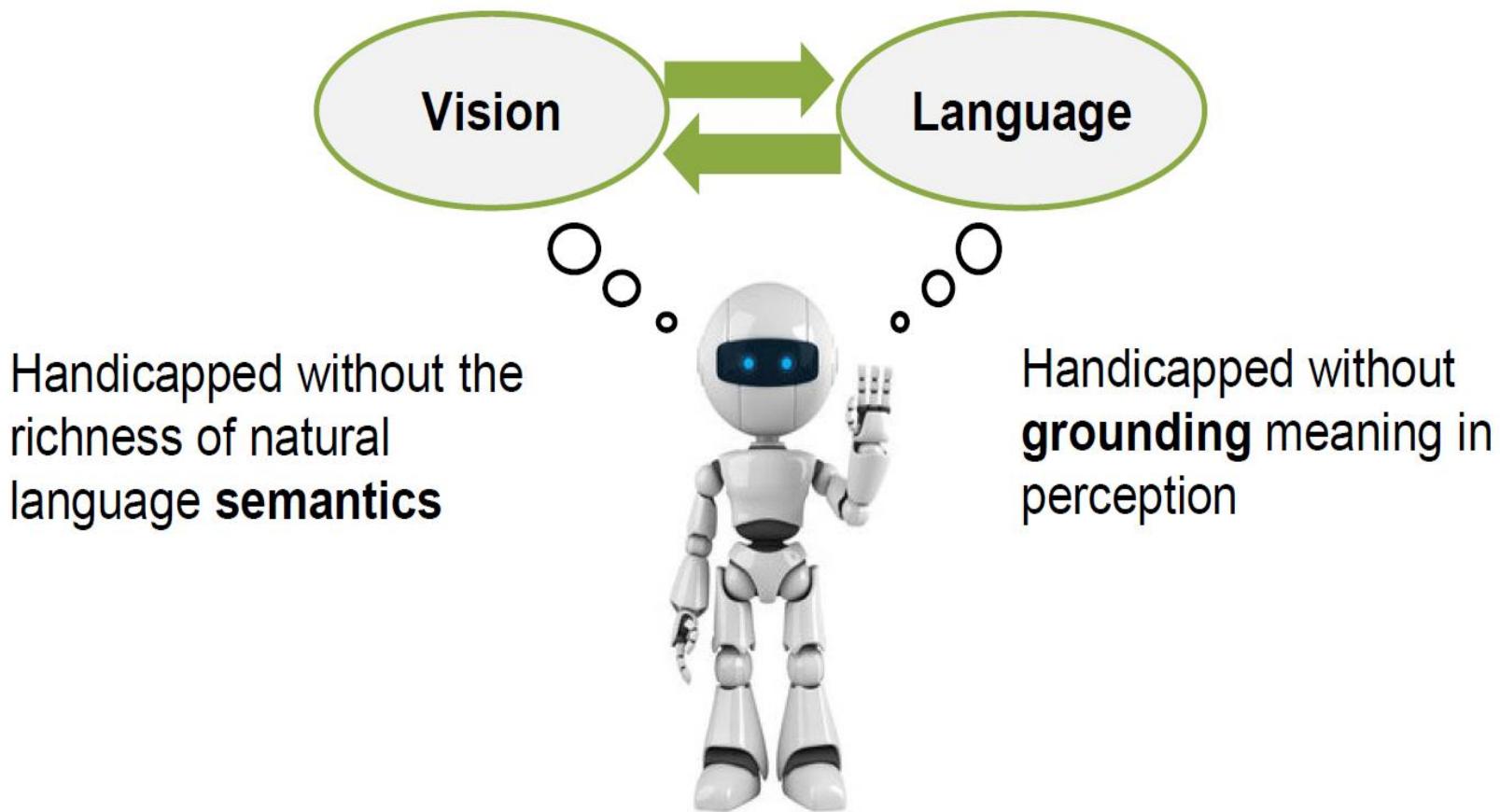
Mask R-CNN: Very Good Results!



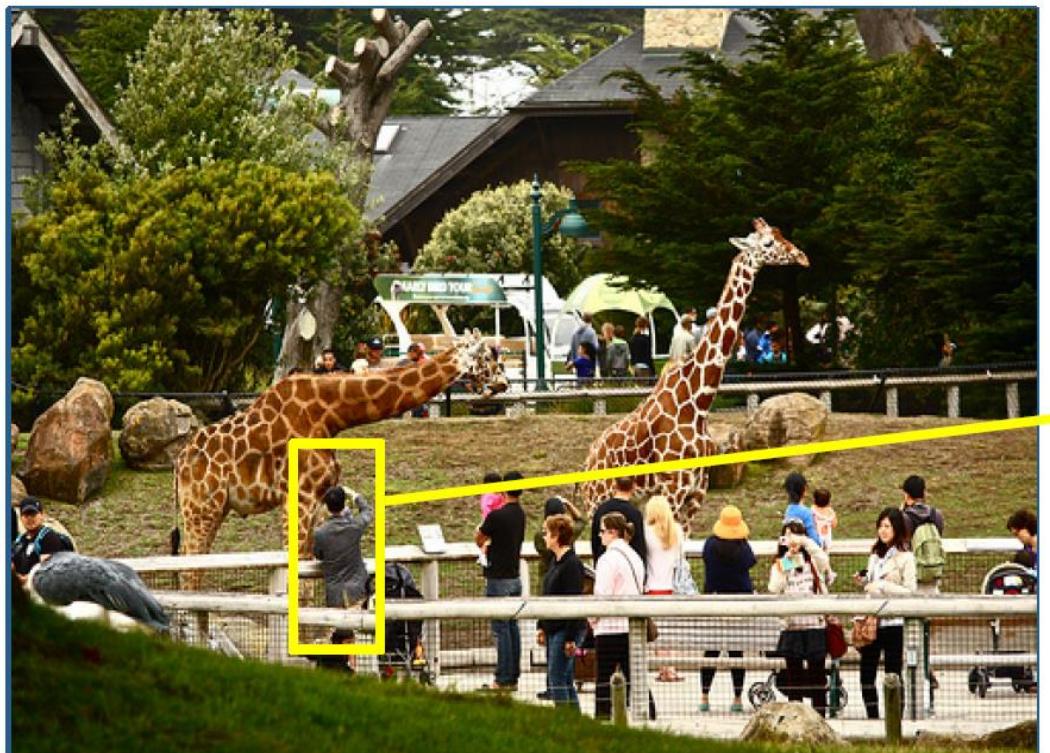
Problem: the “Toddler” AI



Why combine vision and language?



How can we connect vision and language?



Captioning Hendricks et al, CVPR16
Ramanishka et al, ACMM16

A crowd of people looking at giraffes in a zoo.

Referring Expressions Hu et al., CVPR16
Person taking a photo?

Question Answering Xu and Saenko ECCV16
What time of year is it?
Answer: summer

Today: connecting language to vision

Image Description

Input image



Output: A close up of a hot dog on a bun.

Video Description

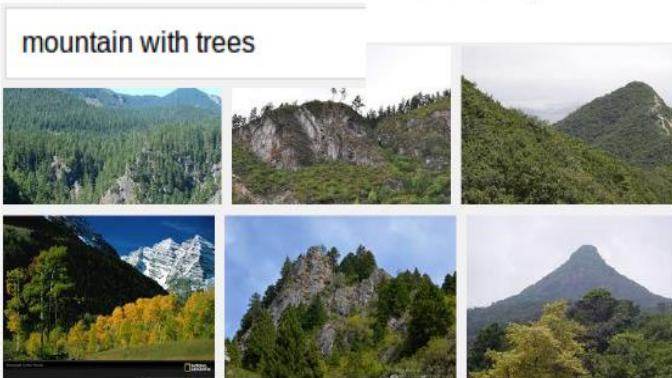
Input video



Output: A woman shredding chicken in a kitchen

Applications

Image and video retrieval by content.



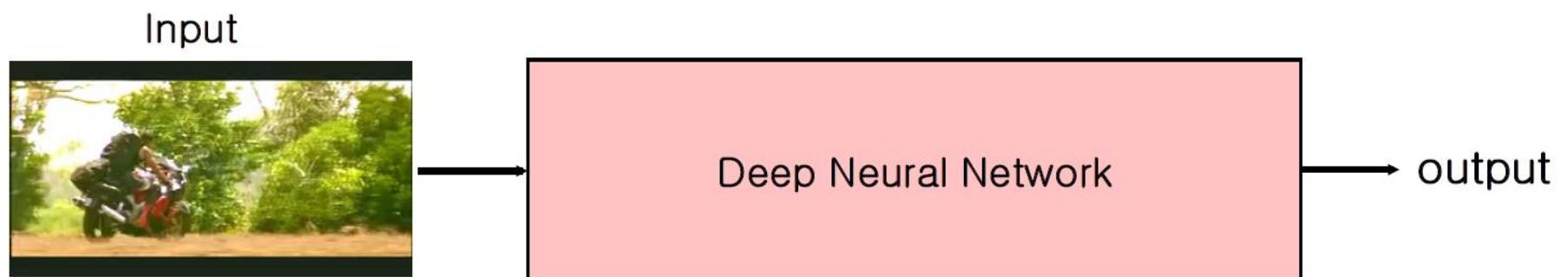
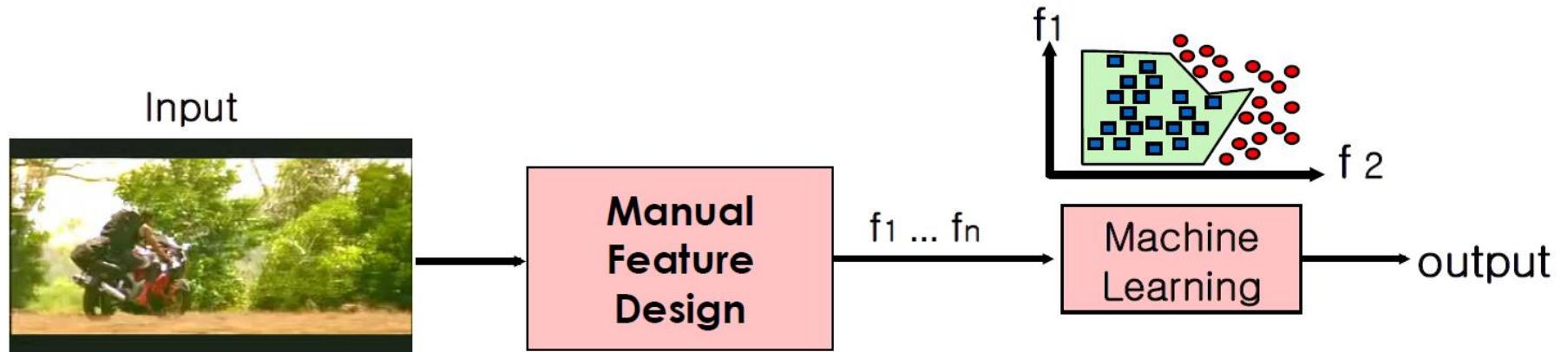
Human Robot Interaction

Video description service.

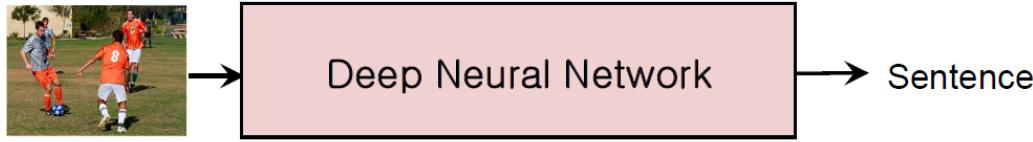


Video surveillance

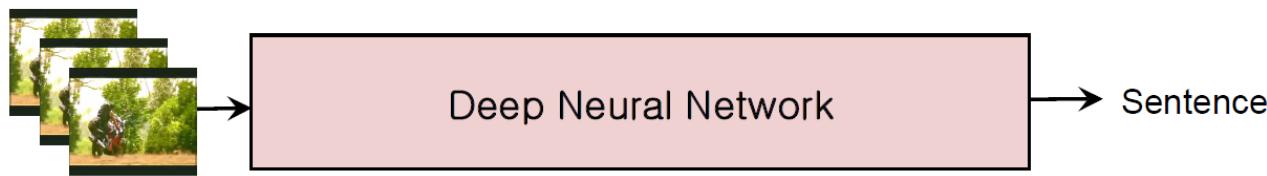
Deep Learning Revolution



Deep End-to-End Neural Models based on Recurrent Nets

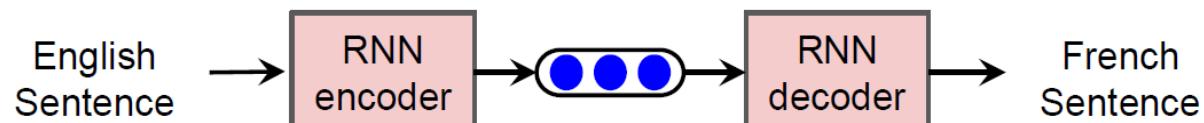


[Donahue et al. CVPR'15]
(our work)
[Vinyals et al. CVPR'15]

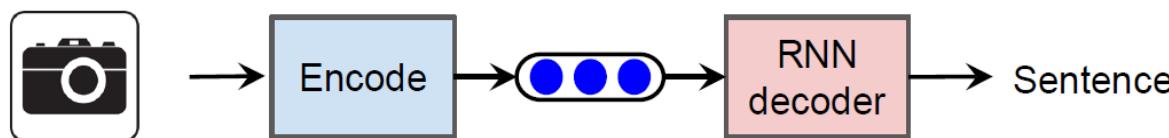


[Venugopalan et. al.
NAACL'15]
[Venugopalan et. al.
ICCV'15] (our work)

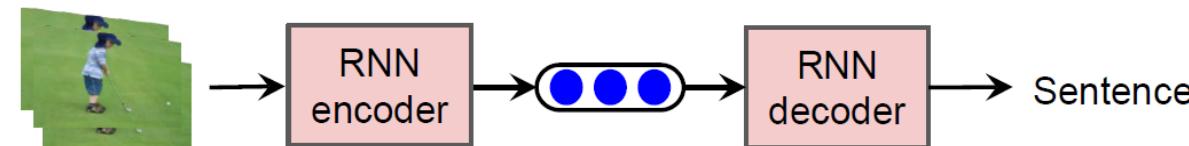
Deep End-to-End Neural Models based on Recurrent Nets



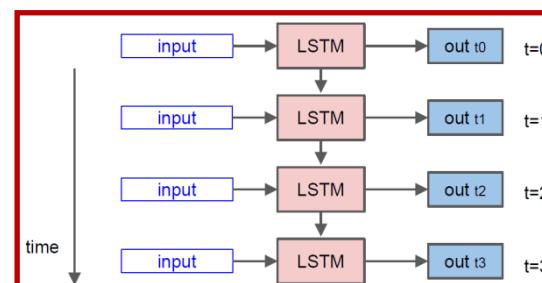
[Sutskever et al. NIPS'14]

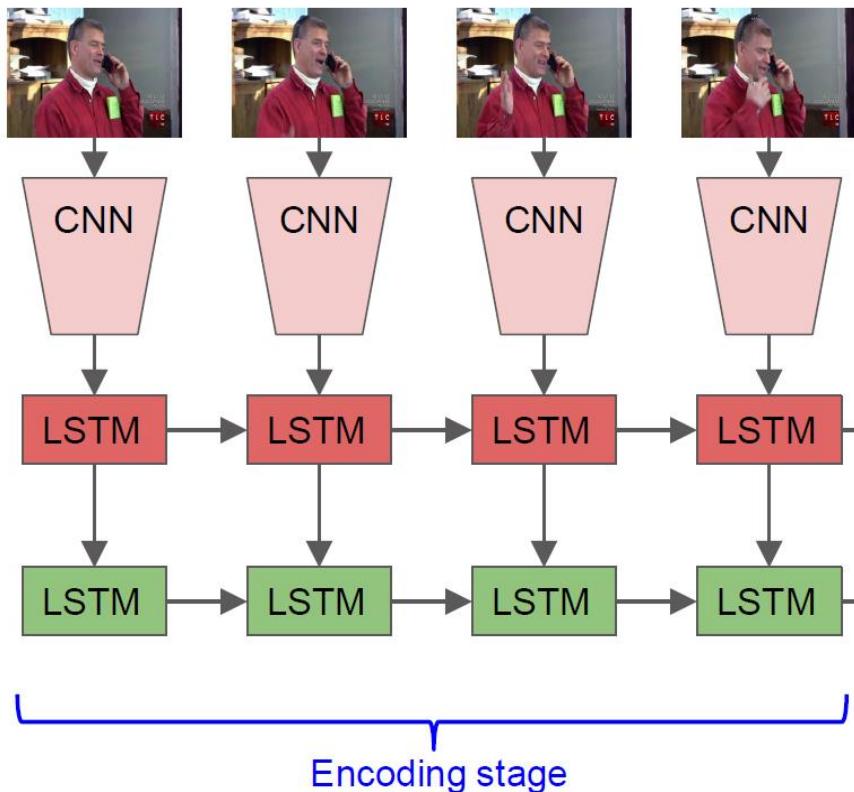


[Donahue et al. CVPR'15]
[Vinyals et al. CVPR'15]



[Venugopalan et. al. NAACL'15]
[Venugopalan et. al. ICCV'15]





S2VT: Sequence to Sequence Video to Text

Now decode it to a sentence!

...

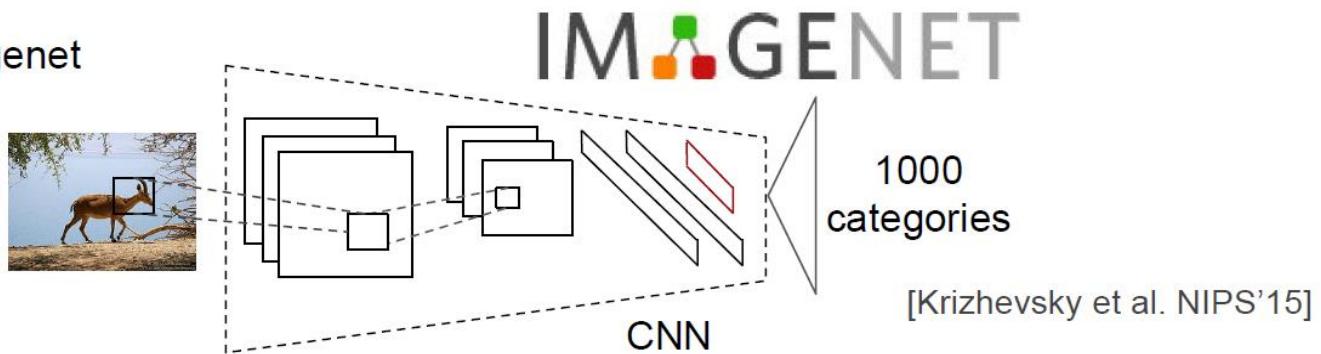
...

...

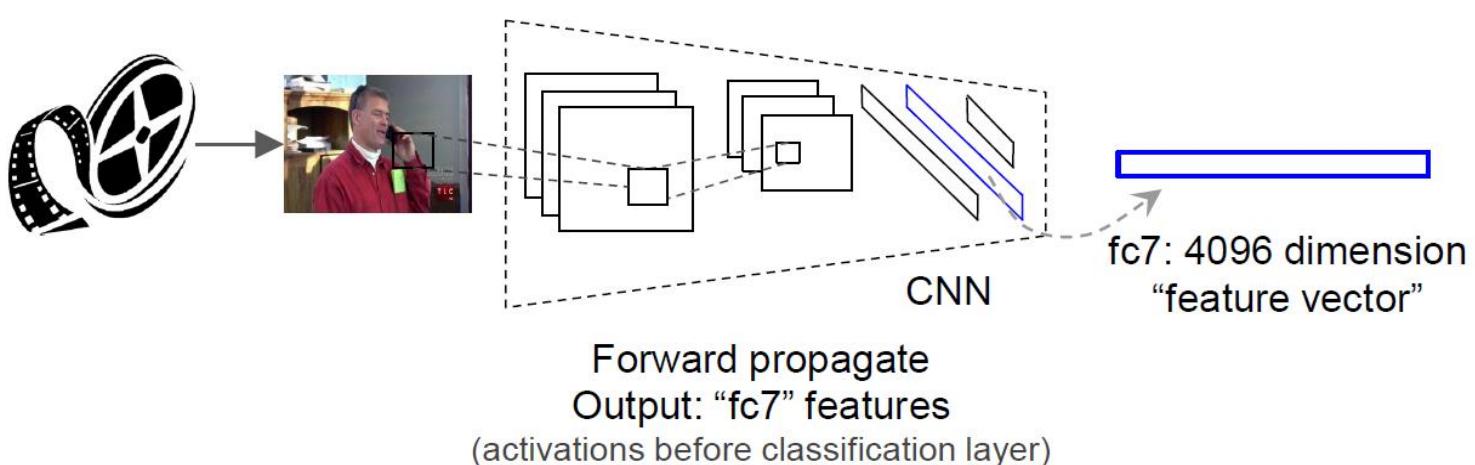
[Venugopalan et. al. ICCV'15]

Objects

1. Train on Imagenet

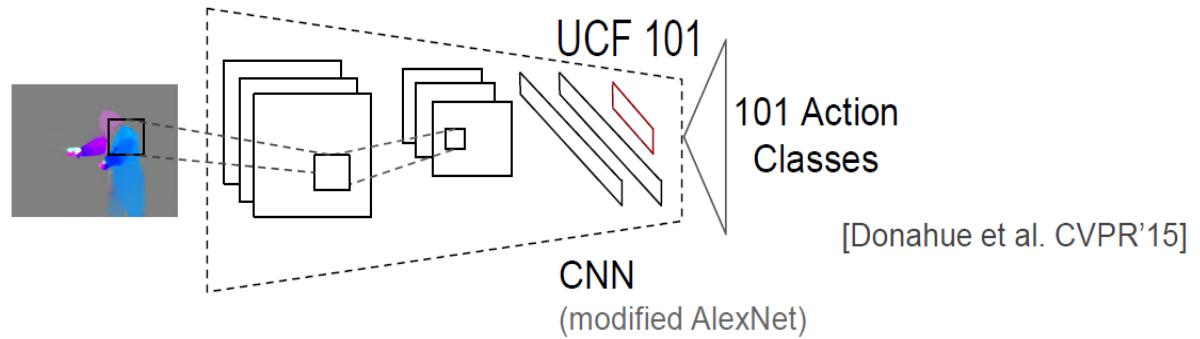


2. Take activations from layer before classification

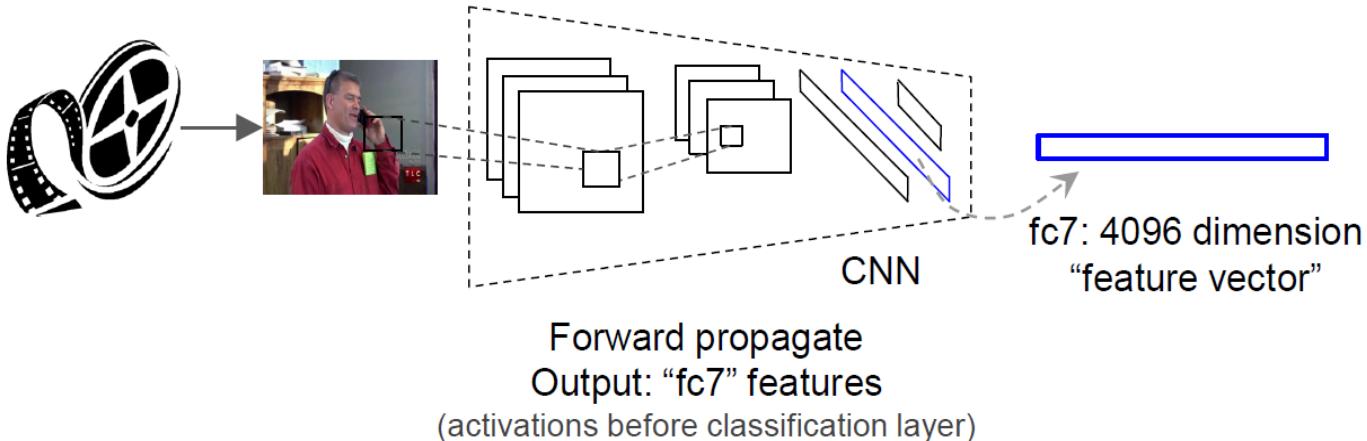


Actions

1. Train CNN on Activity classes



2. Take activations from layer before classification



Explaining the network's captions

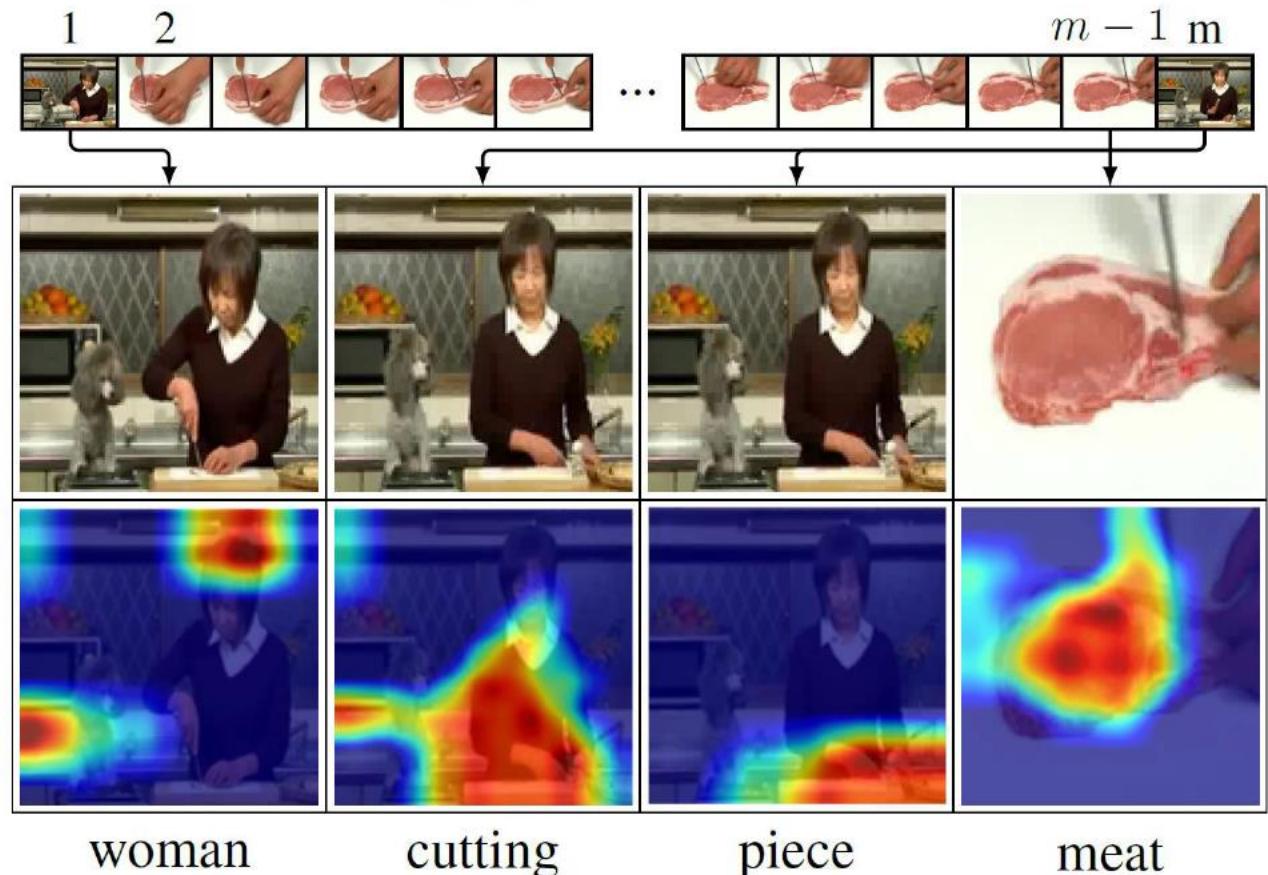
Predicted sentence: A woman is cutting a piece of meat



can the network
localize objects?

Spatiotemporal saliency

Predicted sentence: A woman is cutting a piece of meat



Pointing game in Flickr30kEntities

An elderly man sleeps sitting up on the end of a red couch .



An elderly man



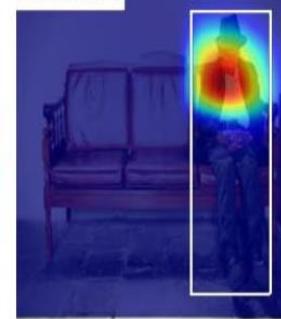
the end of a red couch



An old man is sitting alone on a couch and sleeping .



An old man



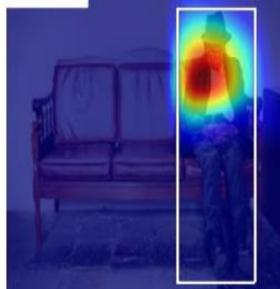
a couch



Old man wearing a hat and coat sleeping sitting up on a sofa .



Old man



a hat



coat



a sofa

