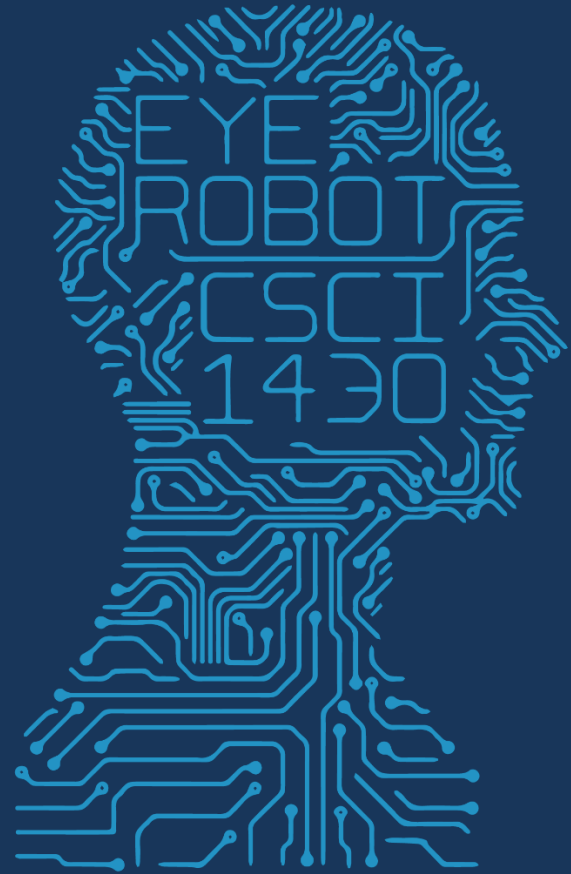




1950

FUTURE VISION



2017 MWF 1PM 368

COMPUTER VISION



[Boston March for Science 2017 – photo Hendrik Strobel]





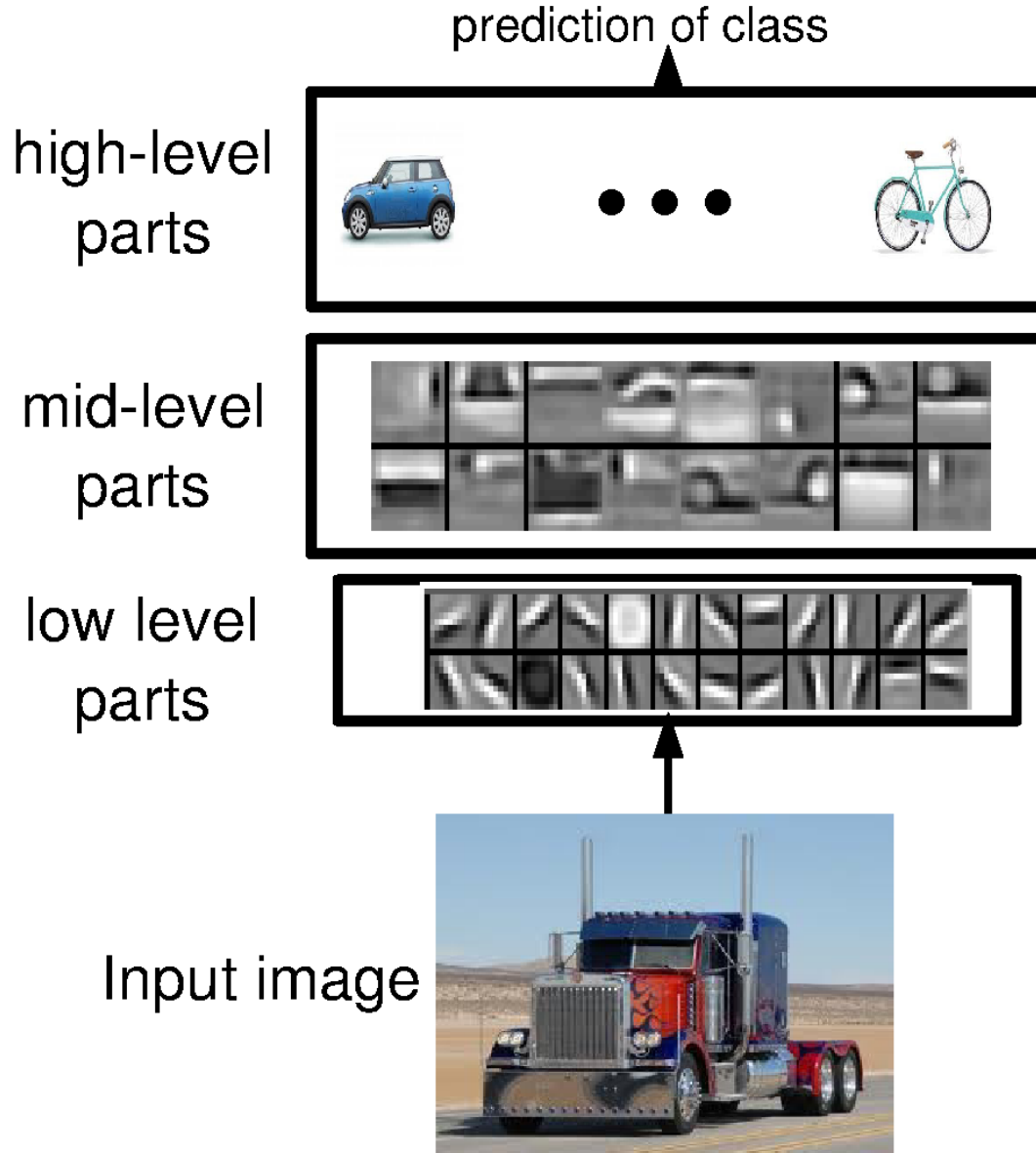






[Boston March for Science 2017]

# Interpretation



- distributed representations
- feature sharing
- compositionality

# Object Detectors Emerge in Deep Scene CNNs

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba

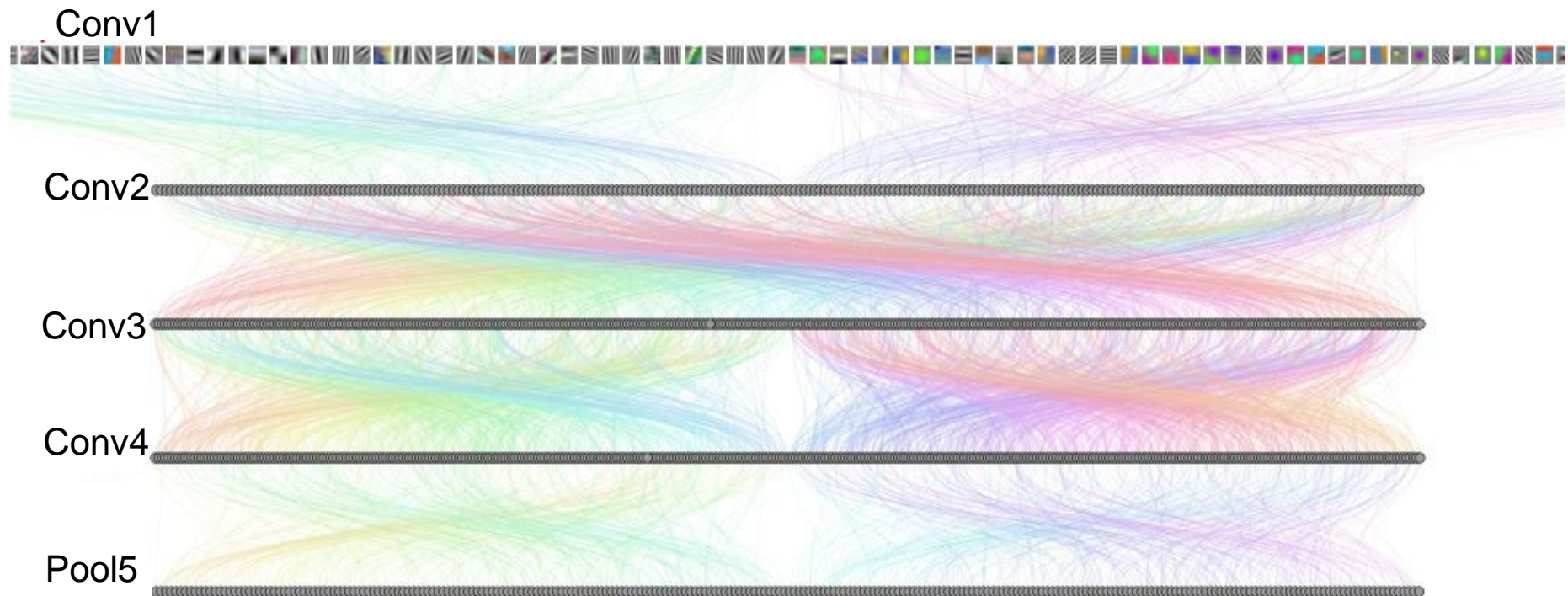


Massachusetts Institute of Technology



# How Objects are Represented in CNN?

CNN uses **distributed code** to represent objects.





# Estimating the Receptive Fields

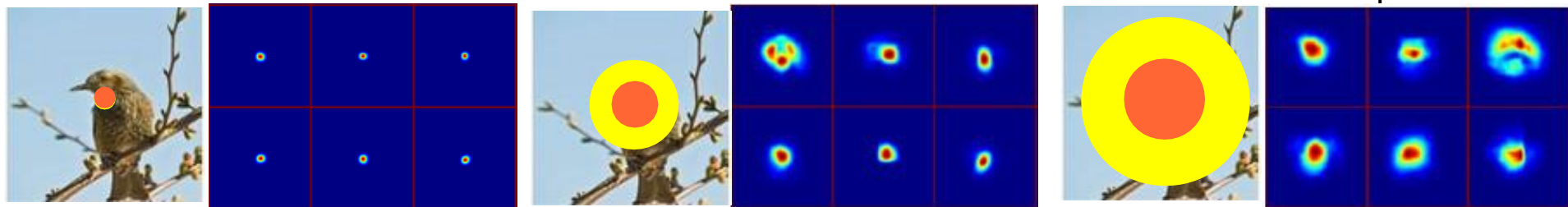
Estimated receptive fields

Actual size of RF is much smaller than the theoretic size

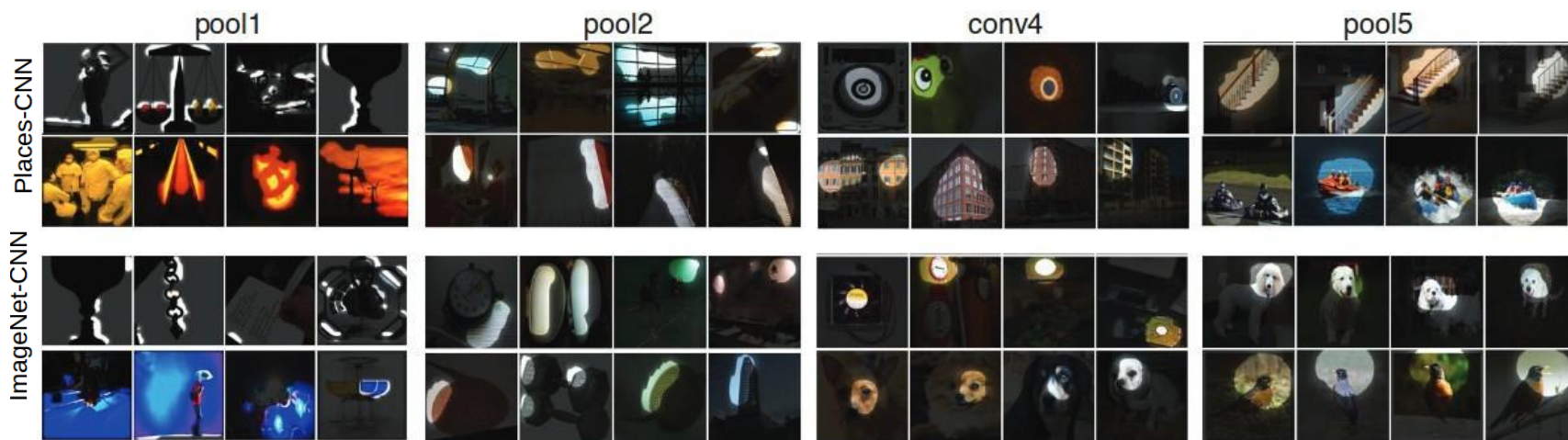
pool1

conv3

pool5



Segmentation using the RF of Units



More semantically meaningful →

# Annotating the Semantics of Units

Top ranked segmented images are cropped and sent to Amazon Turk for annotation.

## Task 1

Word/Short description:

tower

## Task 2

Mark (by clicking on them) the images which don't correspond to the short description you just wrote



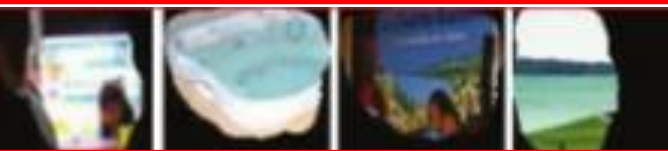
## Task 3

Which category does your short description mostly belong to?

- ☐ Scene (kitchen, corridor, street, beach, ...)
- ☐ Region or surface (road, grass, wall, floor, sky, ...)
- ☒ Object (bed, car, building, tree, ...)
- ☐ Object part (leg, head, wheel, roof, ...)
- ☐ Texture or material (striped, rugged, wooden, plastic, ...)
- ☐ Simple elements or colors (vertical line, curved line, color blue, ...)

# Annotating the Semantics of Units

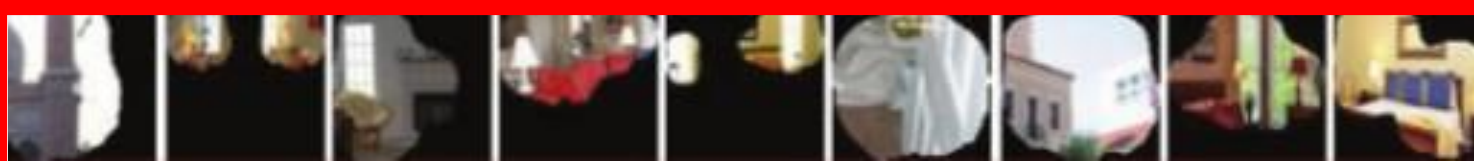
Pool5, unit 76; Label: ocean; Type: scene; Precision: 93%





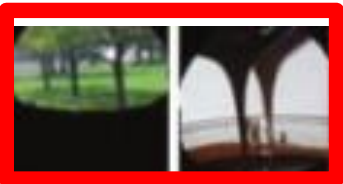
# Annotating the Semantics of Units

Pool5, unit 13; Label: Lamps; Type: object; Precision: 84%



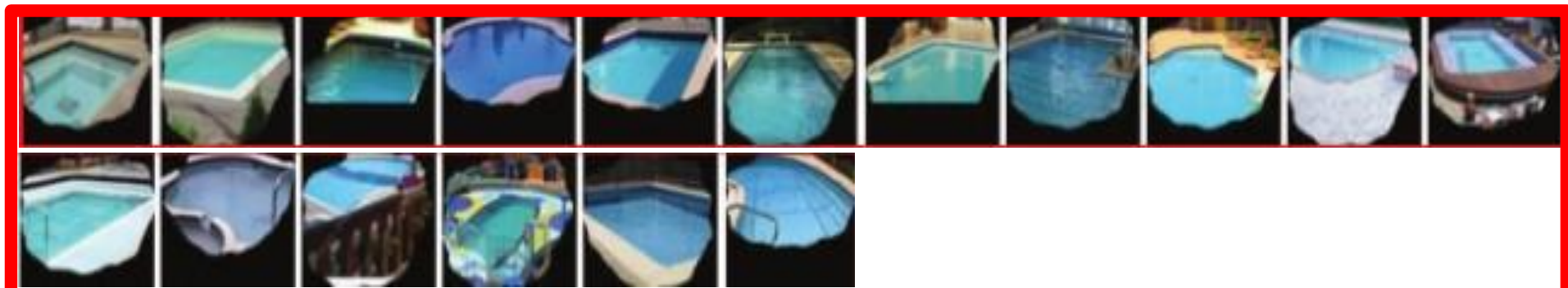
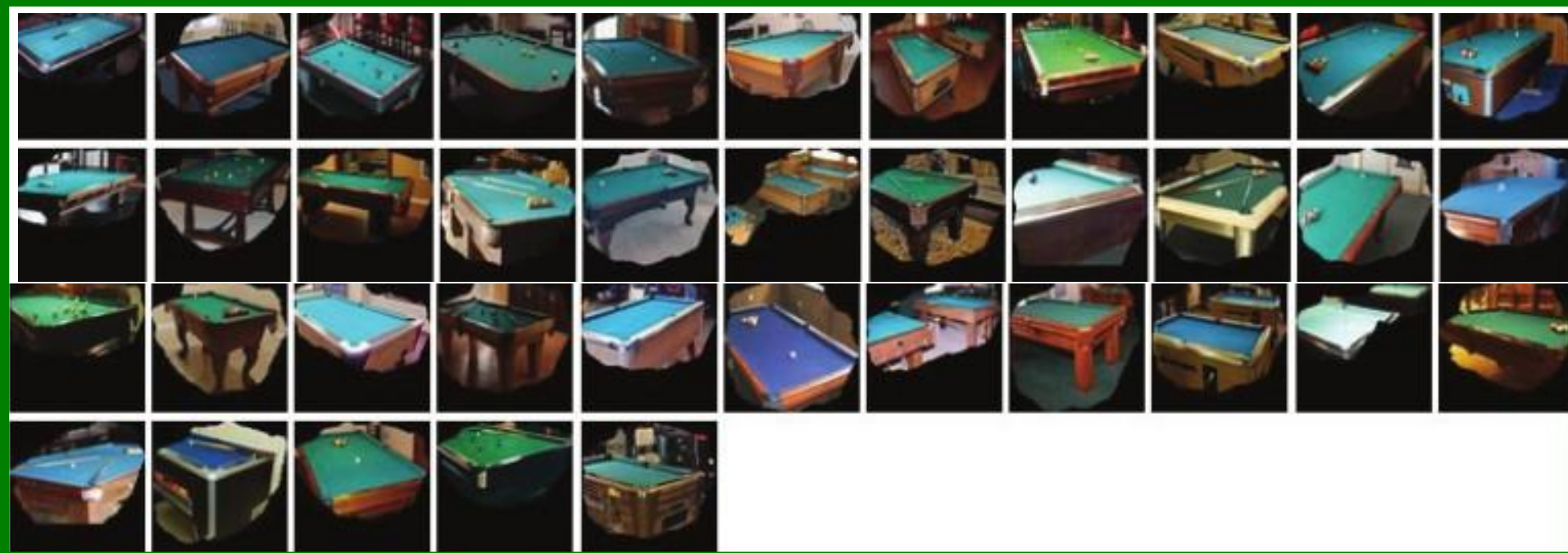
# Annotating the Semantics of Units

Pool5, unit 77; Label:legs; Type: object part; Precision: 96%



# Annotating the Semantics of Units

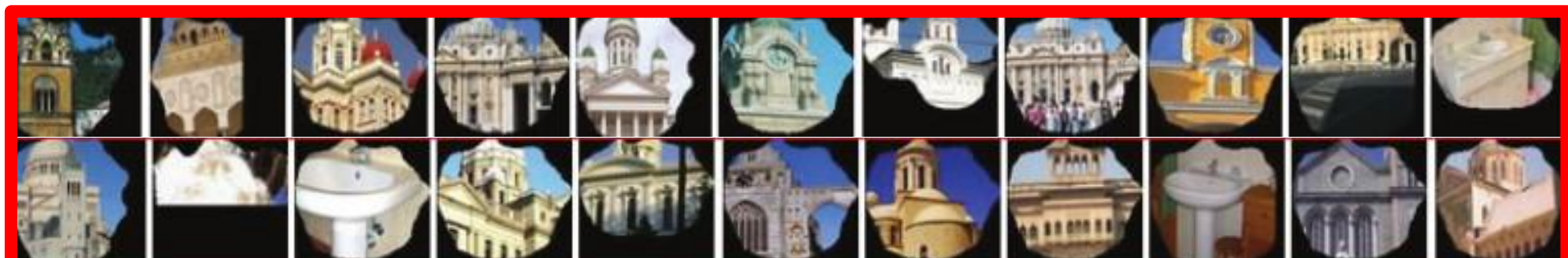
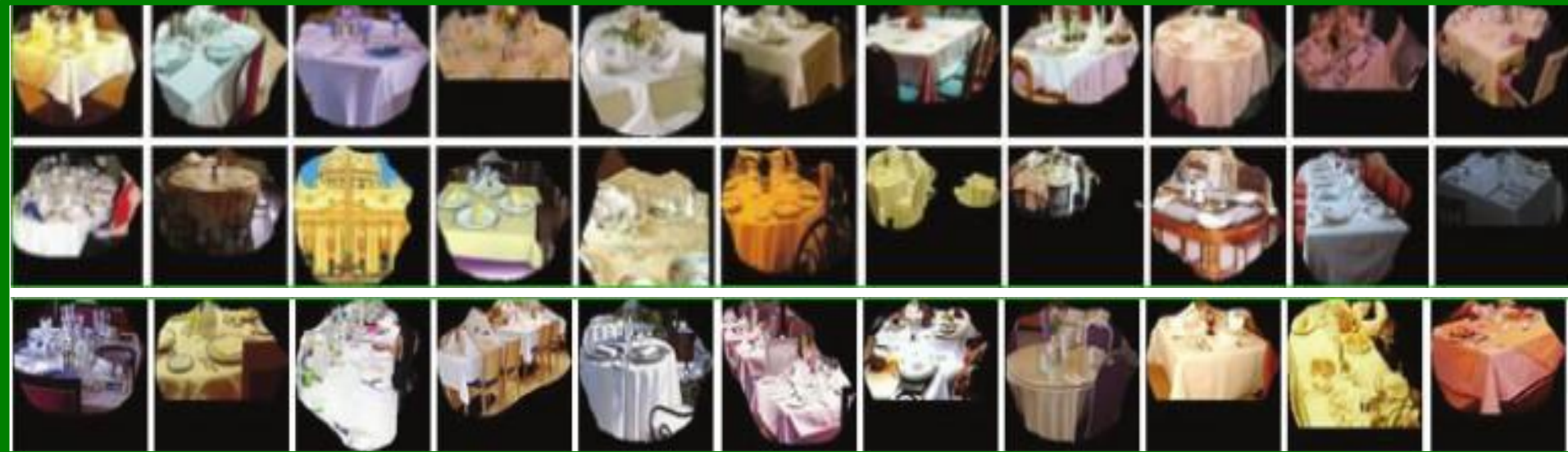
Pool5, unit 112; Label: pool table; Type: object; Precision: 70%



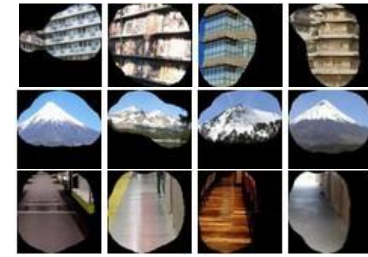
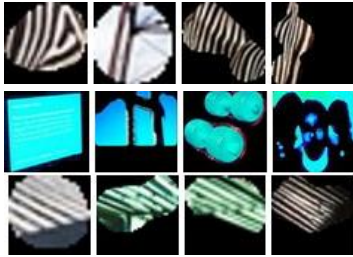


# Annotating the Semantics of Units

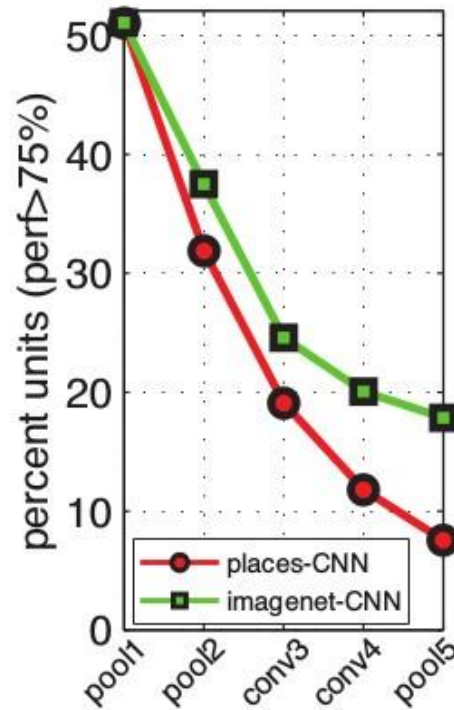
Pool5, unit 22; Label: dinner table; Type: scene; Precision: 60%



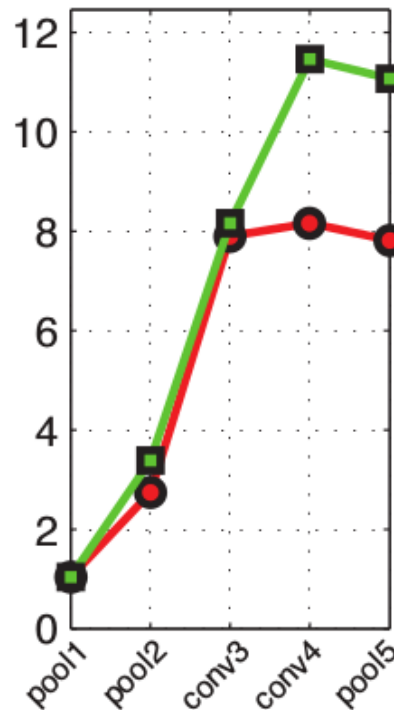
# Distribution of Semantic Types at Each Layer



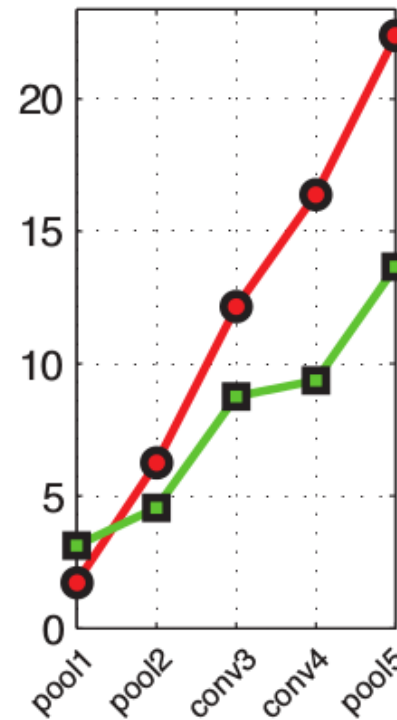
Simple elements & colors



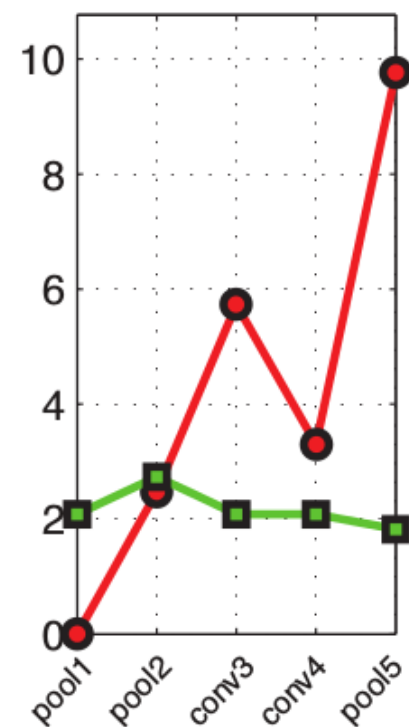
Object part



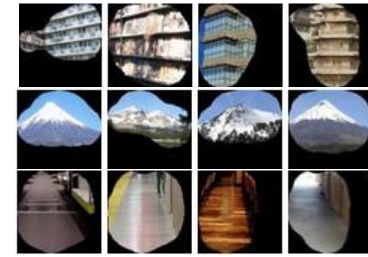
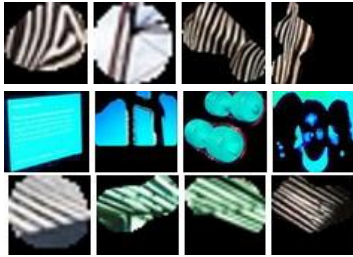
Object



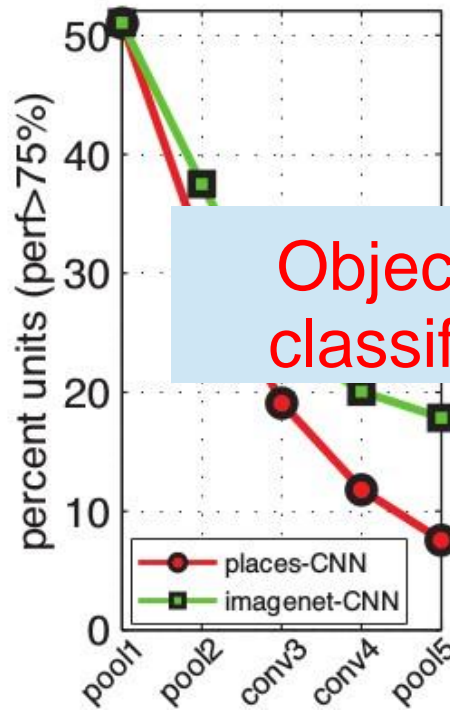
Scene



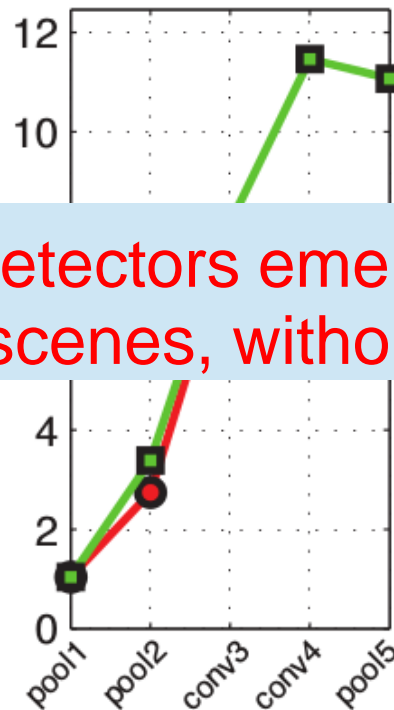
# Distribution of Semantic Types at Each Layer



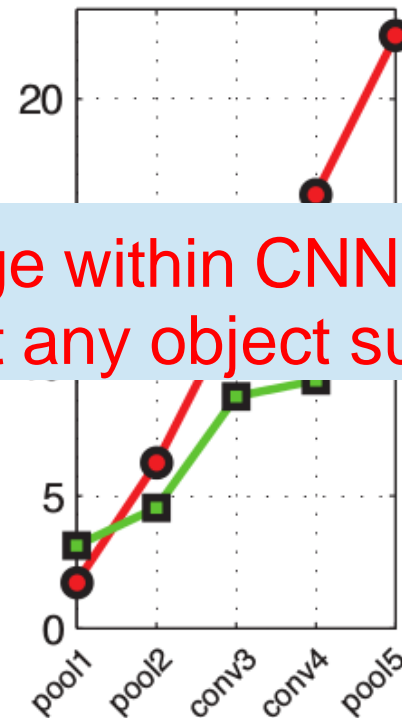
Simple elements & colors



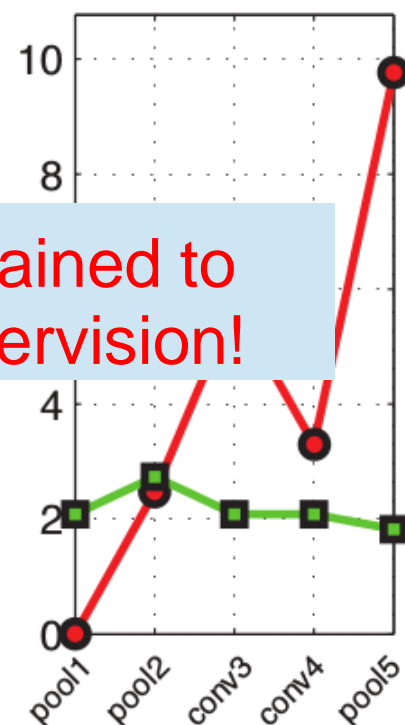
Object part



Object



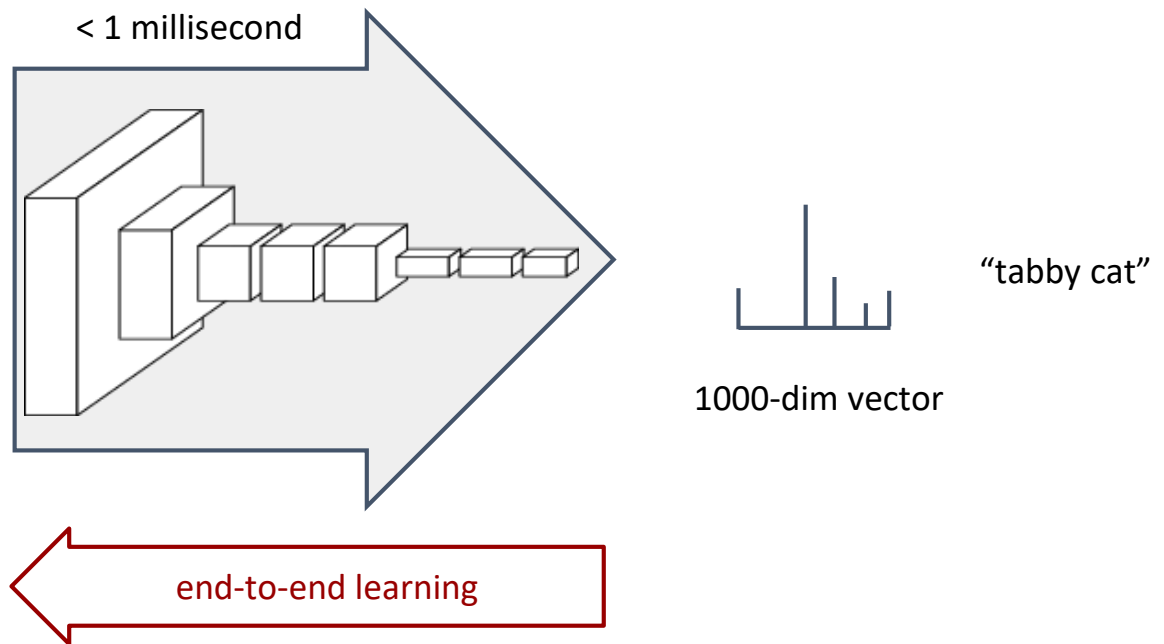
Scene



Object detectors emerge within CNN trained to classify scenes, without any object supervision!

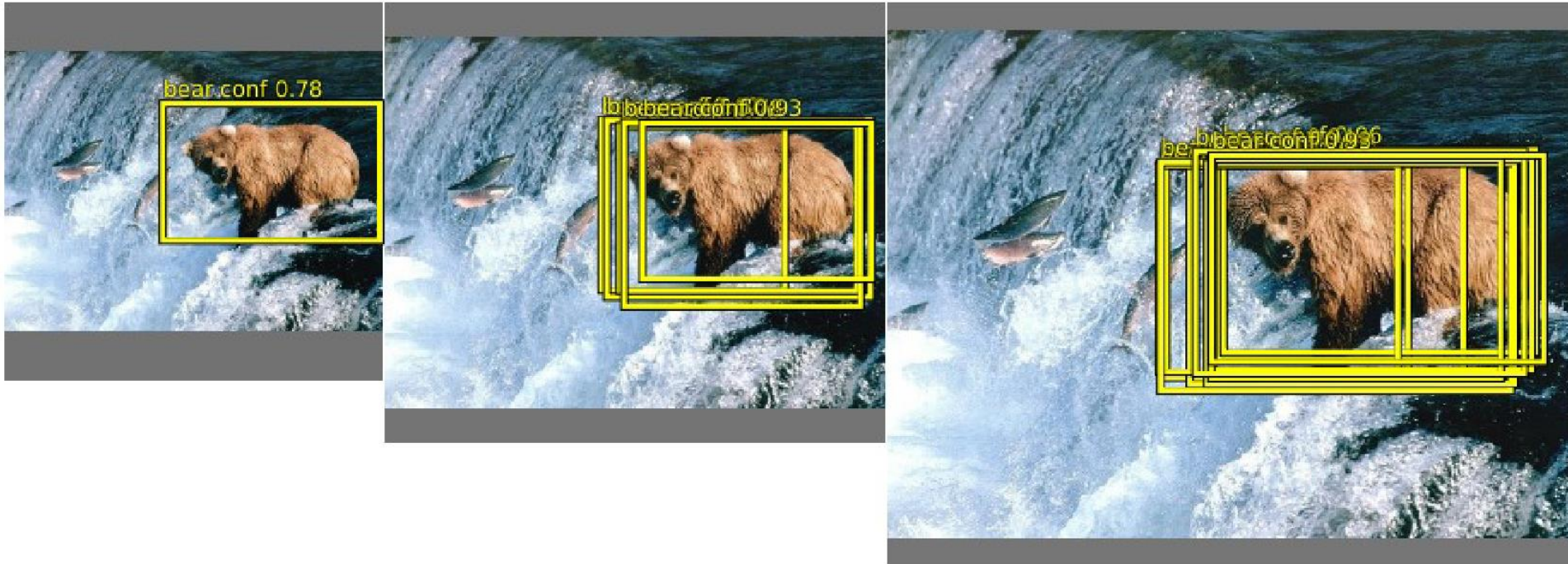


# ConvNets perform classification



# CONV NETS: EXAMPLES

## - Object detection



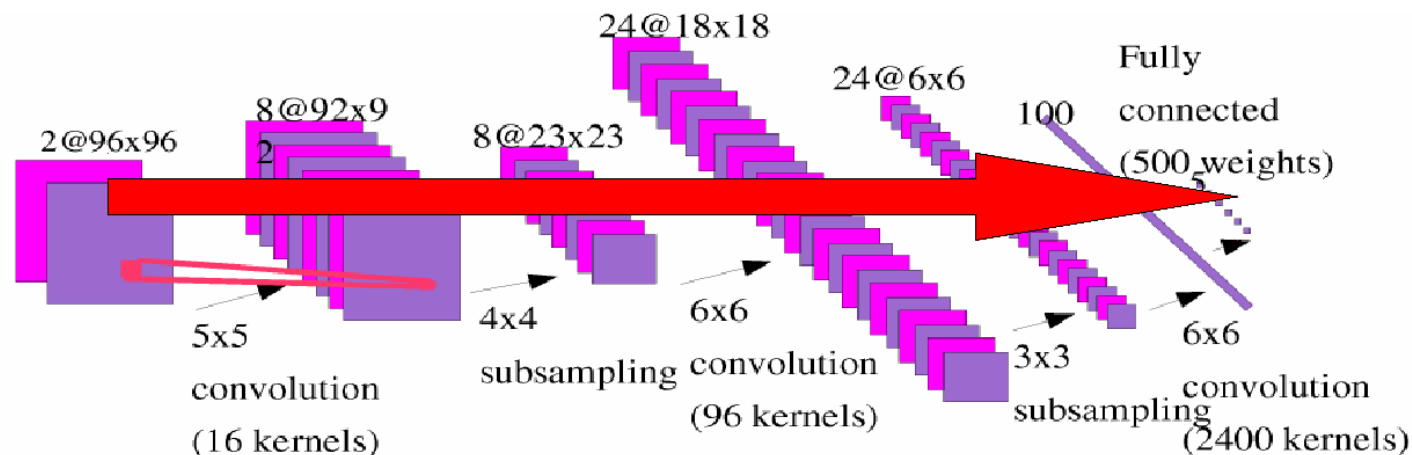
Sermanet et al. "OverFeat: Integrated recognition, localization, ..." arxiv 2013

Girshick et al. "Rich feature hierarchies for accurate object detection..." arxiv 2013 <sup>91</sup>

Szegedy et al. "DNN for object detection" NIPS 2013

# ConvNets: Test

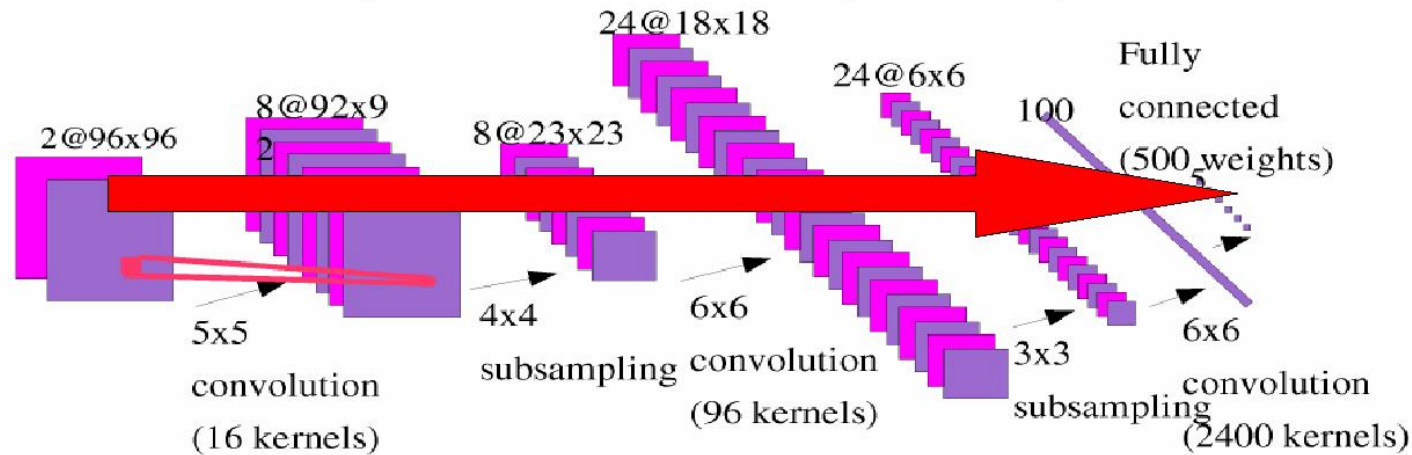
At test time, run only is forward mode (FPROP).





# ConvNets: Test

At test time, run only is forward mode (FPROP).



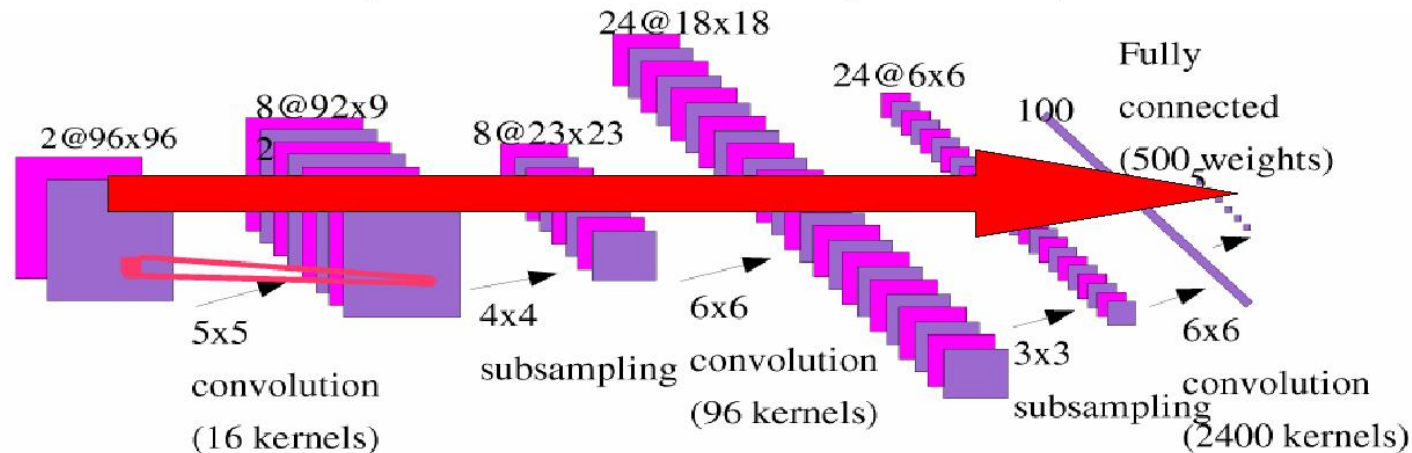
Naturally, convnet can process larger images



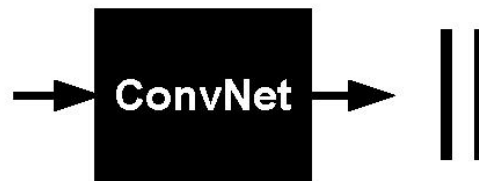
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



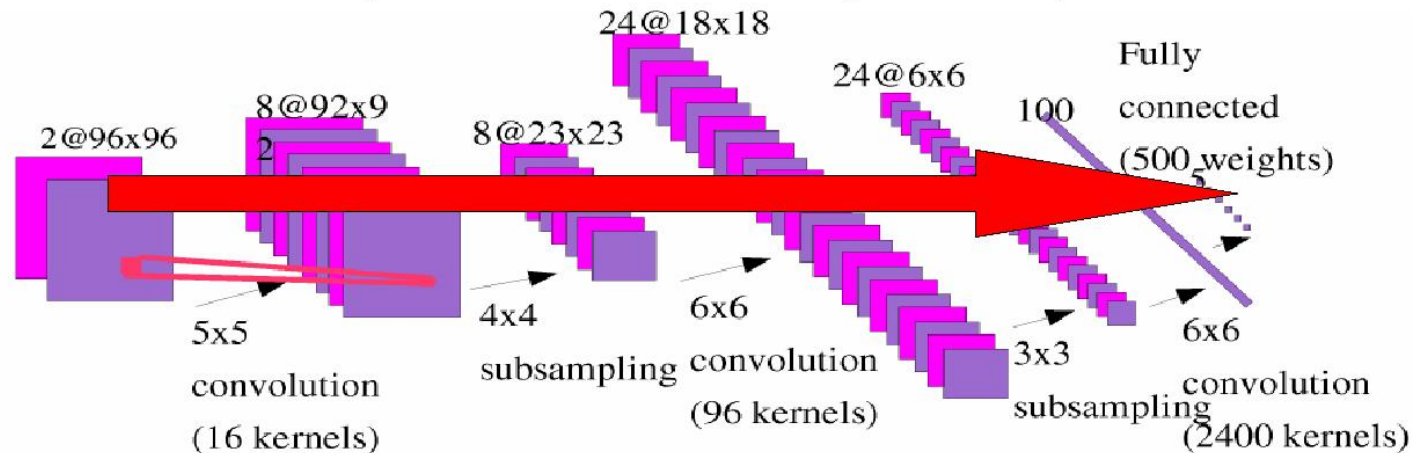
Naturally, convnet can process larger images



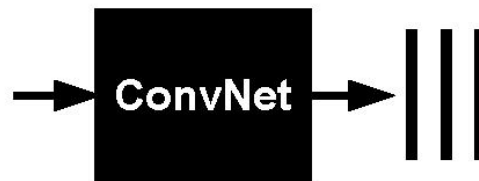
Traditional methods use inefficient sliding windows.

# ConvNets: Test

At test time, run only is forward mode (FPROP).



Naturally, convnet can process larger images

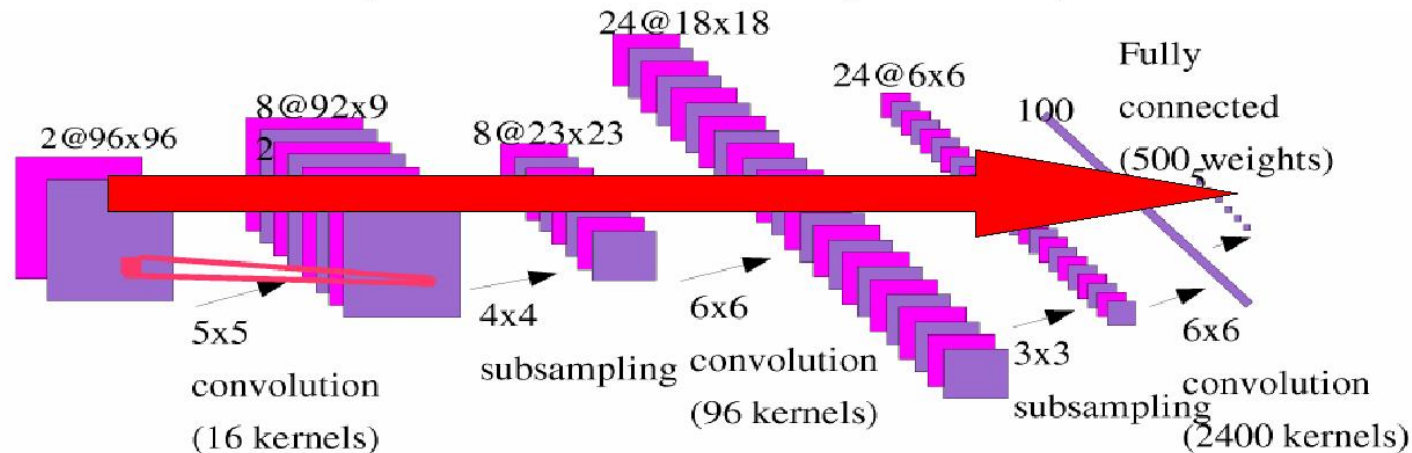


Traditional methods use inefficient sliding windows.

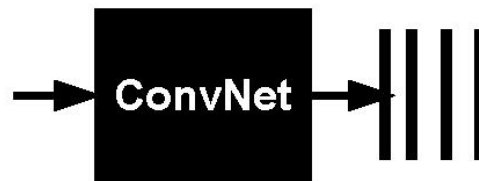


# ConvNets: Test

At test time, run only is forward mode (FPROP).



Naturally, convnet can process larger images

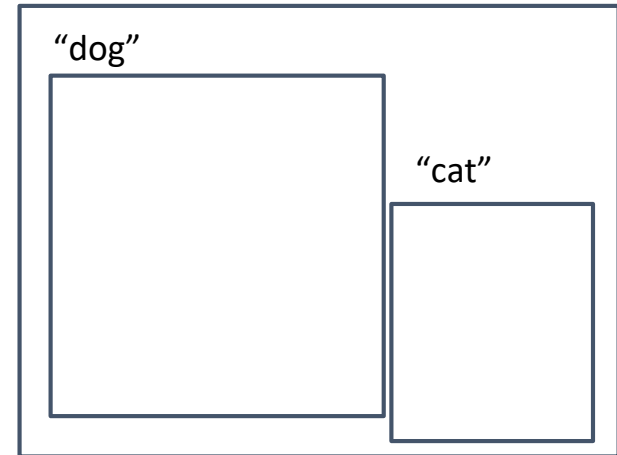
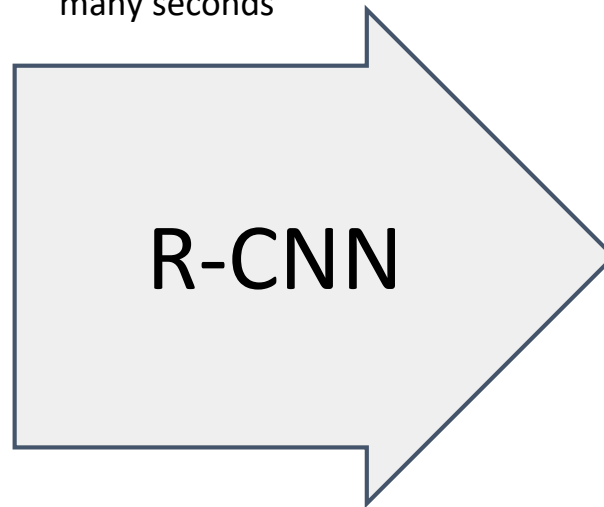


Traditional methods use inefficient sliding windows.

# R-CNN does detection



many seconds



# R-CNN: Region-based CNN

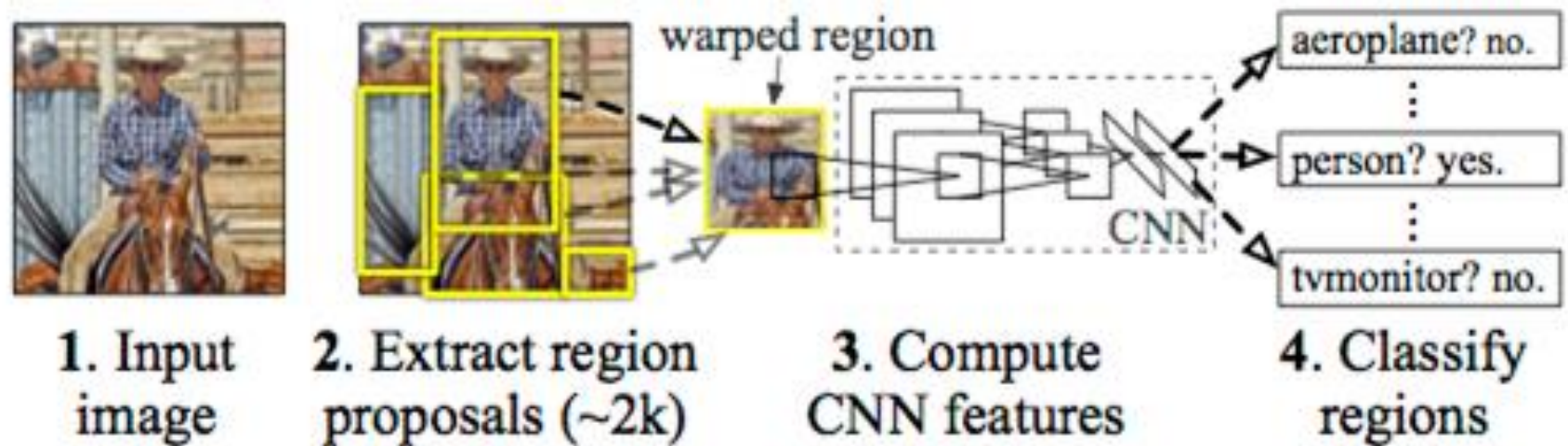
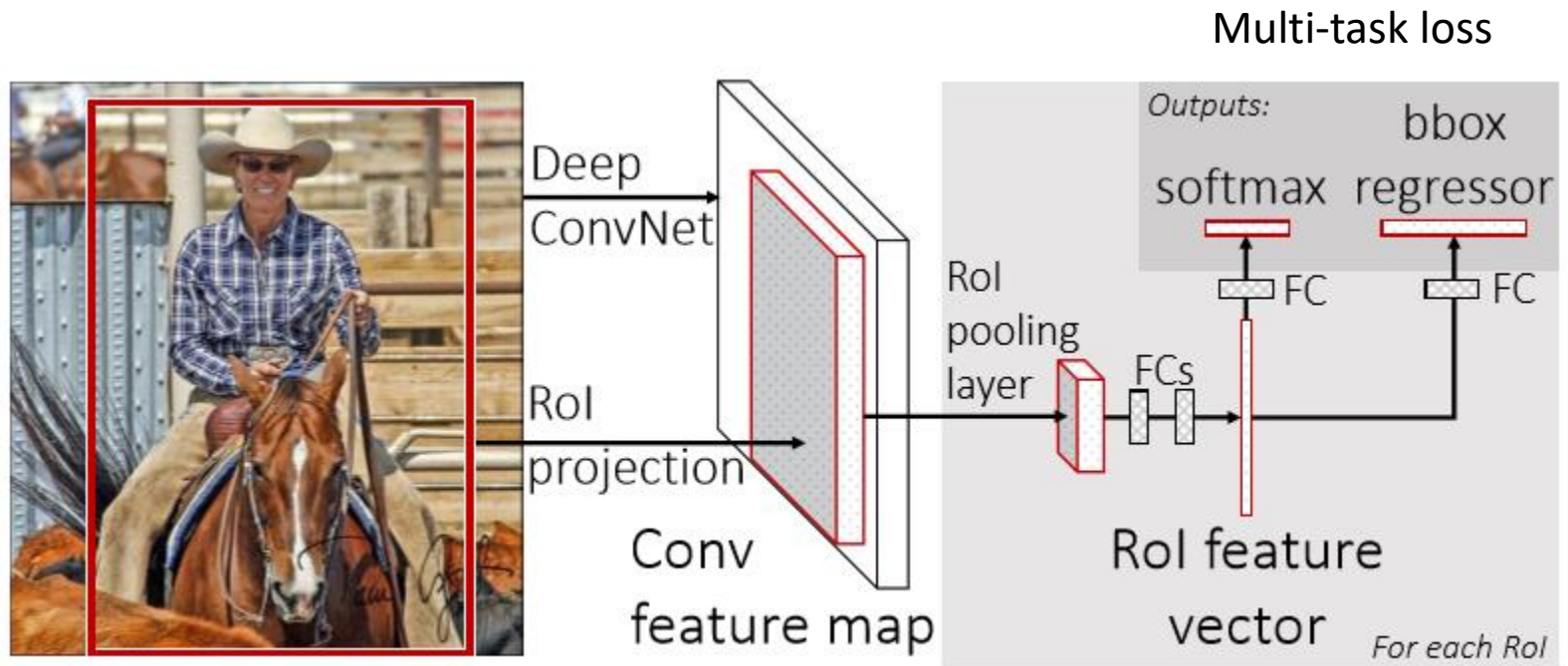


Figure: Girshick et al.



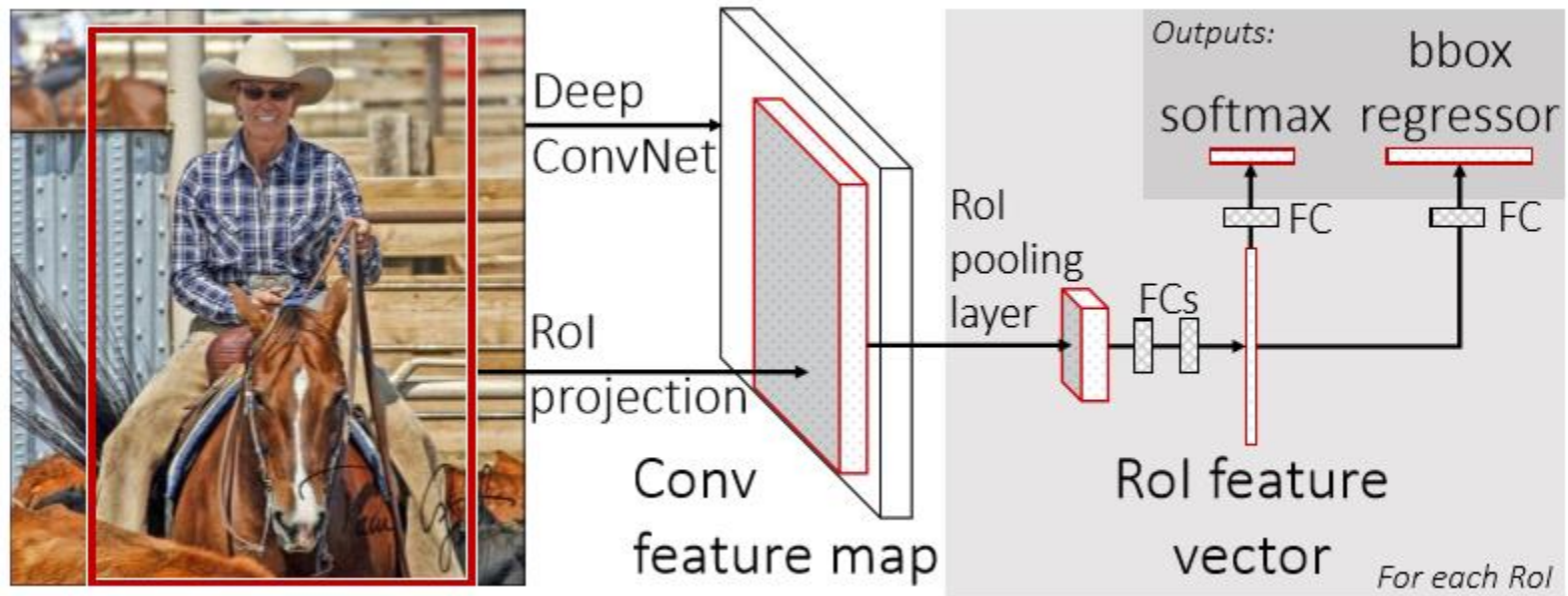
# Fast R-CNN



RoI = Region of Interest

Figure: Girshick et al.

# Fast R-CNN



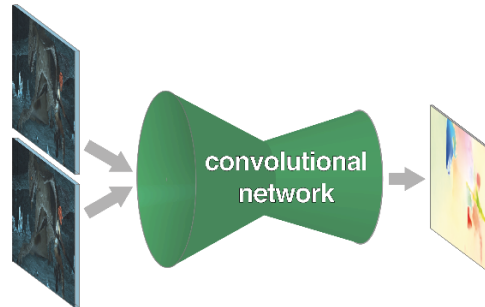
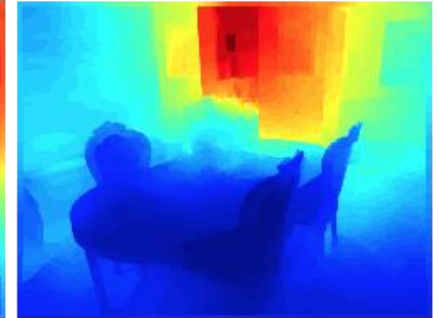
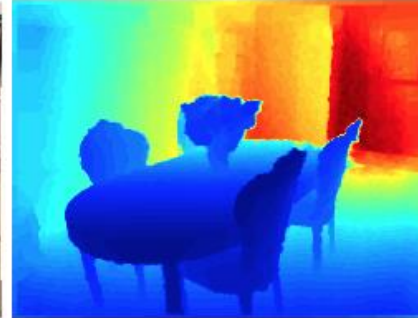
- Convolve whole image into feature map (many layers; abstracted)
- For each candidate RoI:
  - Squash feature map weights into fixed-size 'RoI pool' – adaptive subsampling!
    - Divide RoI into  $H \times W$  subwindows, e.g.,  $7 \times 7$ , and max pool
  - Learn classification on RoI pool with own fully connected layers (FCs)
  - Output classification (softmax) + bounds (regressor)

# What if we want pixels out?

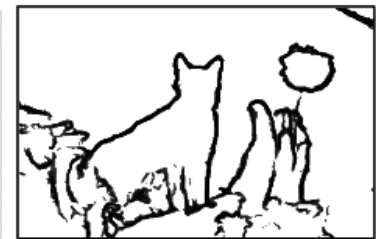
semantic  
segmentation



monocular depth estimation Eigen & Fergus 2015



optical flow Fischer et al. 2015

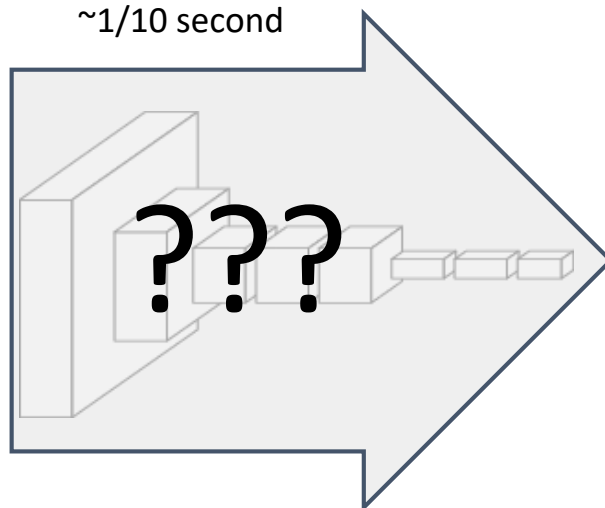


boundary prediction Xie & Tu 2015



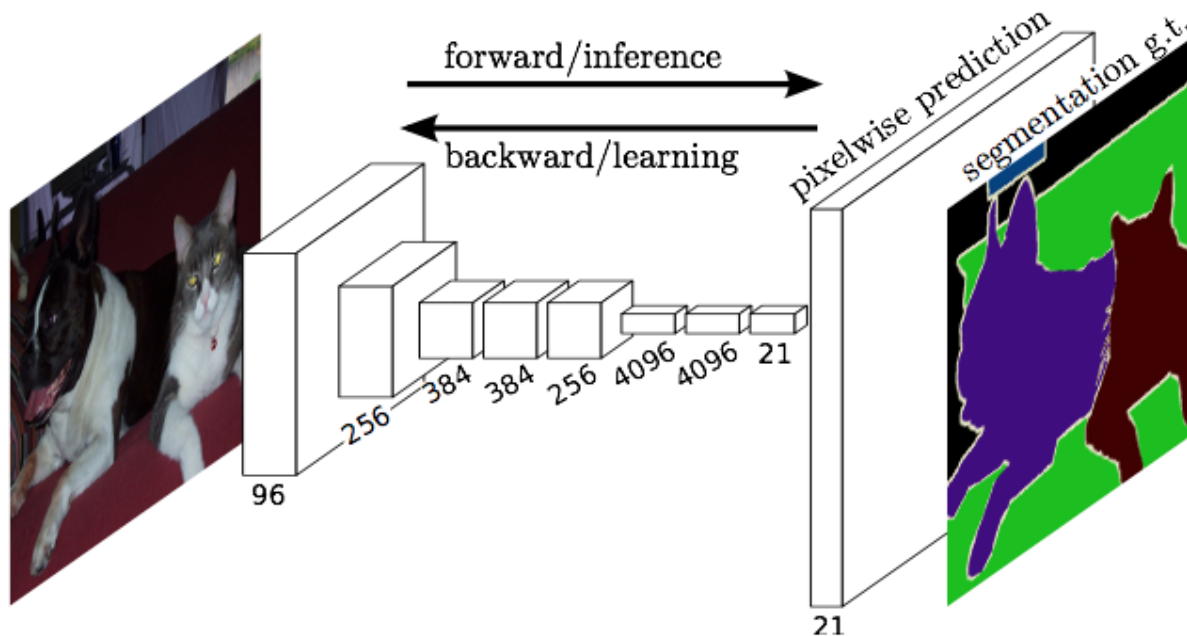


~1/10 second



end-to-end learning

# Fully Convolutional Networks for Semantic Segmentation



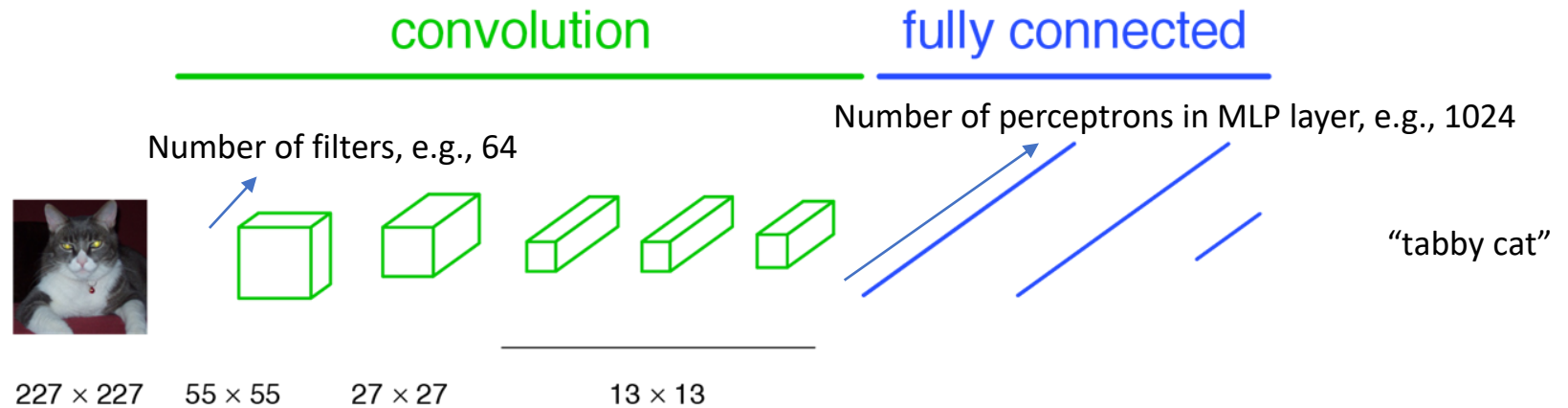
Jonathan Long\*

Evan Shelhamer\*

Trevor Darrell

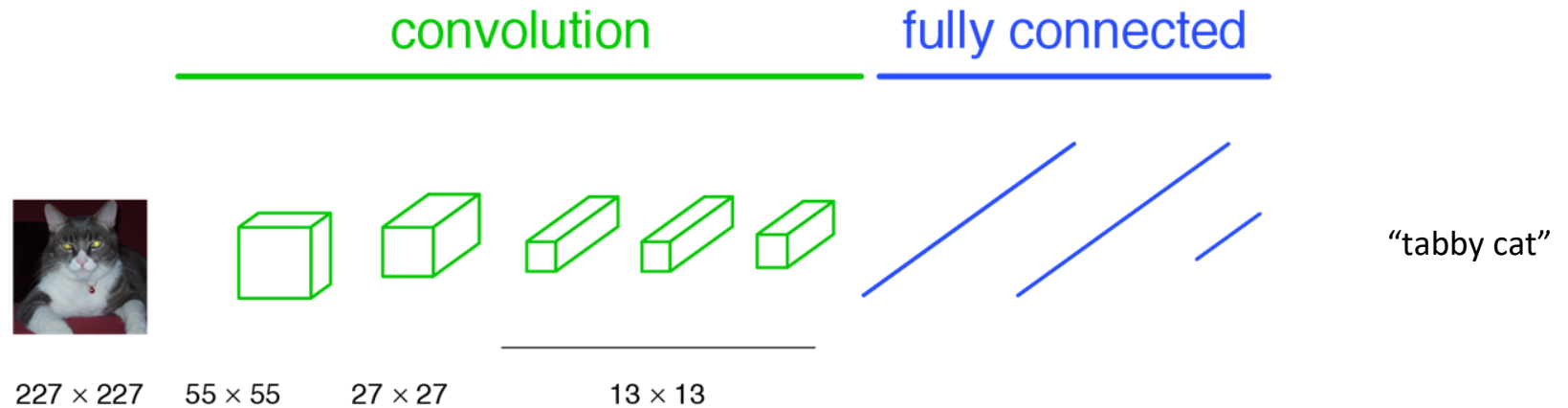
UC Berkeley

# A classification network...

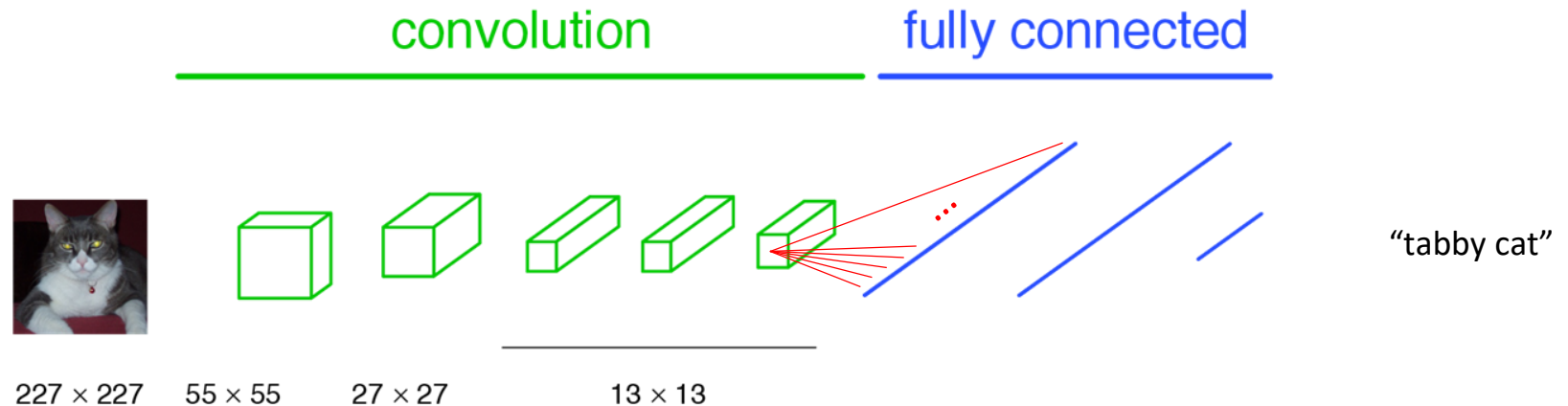




# A classification network...

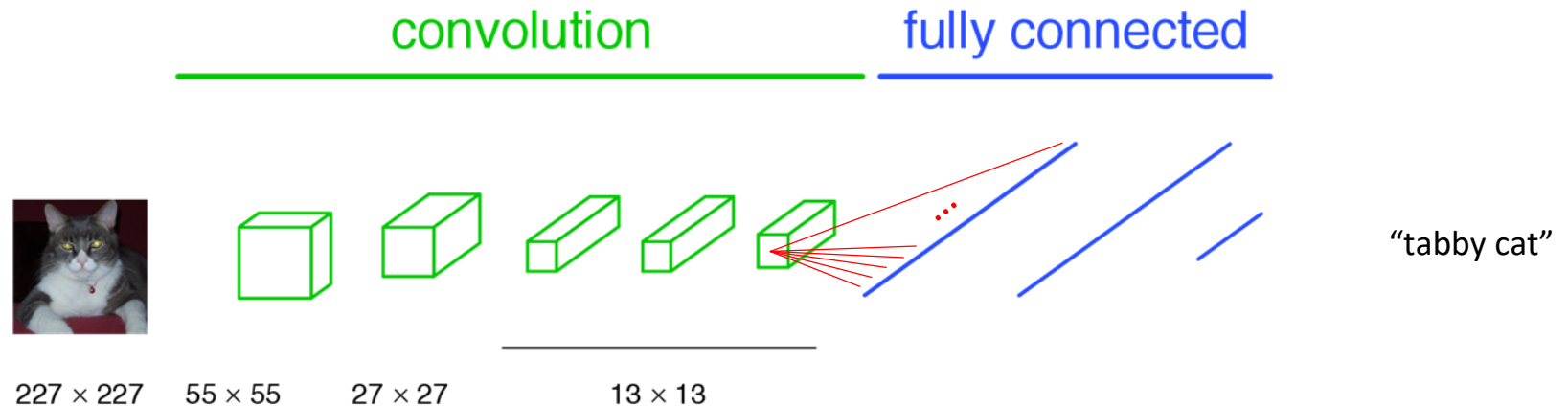


# A classification network...



The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

# A classification network...



The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

AlexNet: 256 filters over  $6 \times 6$  response map

Each 2,359,296 response is attached to one of 4096 perceptrons, leading to 37 mil params.



# Problem

- We want a label at every pixel
- Current network gives us a label for the whole image.
- We want a matrix of labels
- Approach:
  - Make CNN for sub-image size
  - 'Convolutionalize' all layers of network, so that we can treat it as one (complex) filter and slide around our full image.

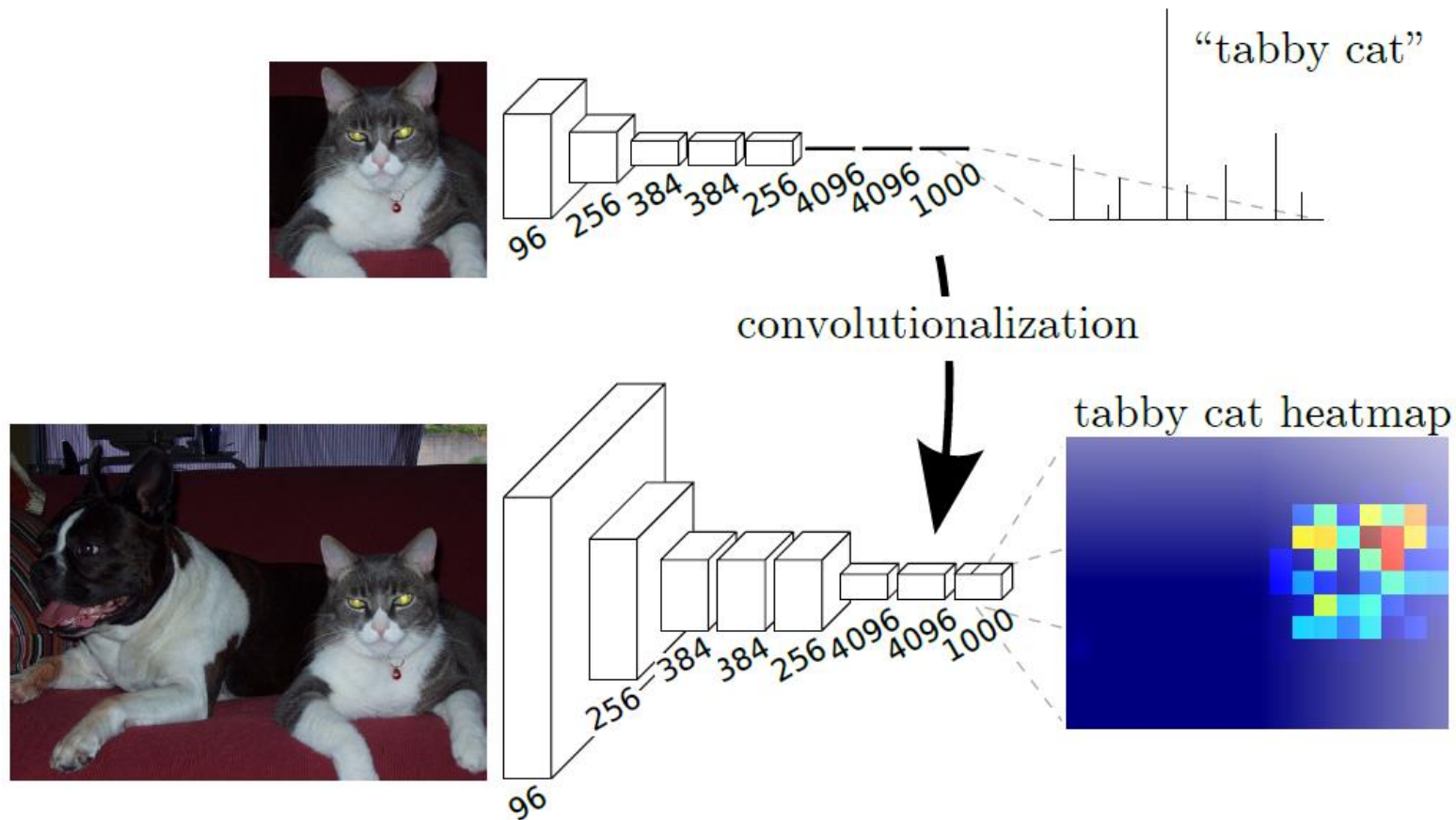
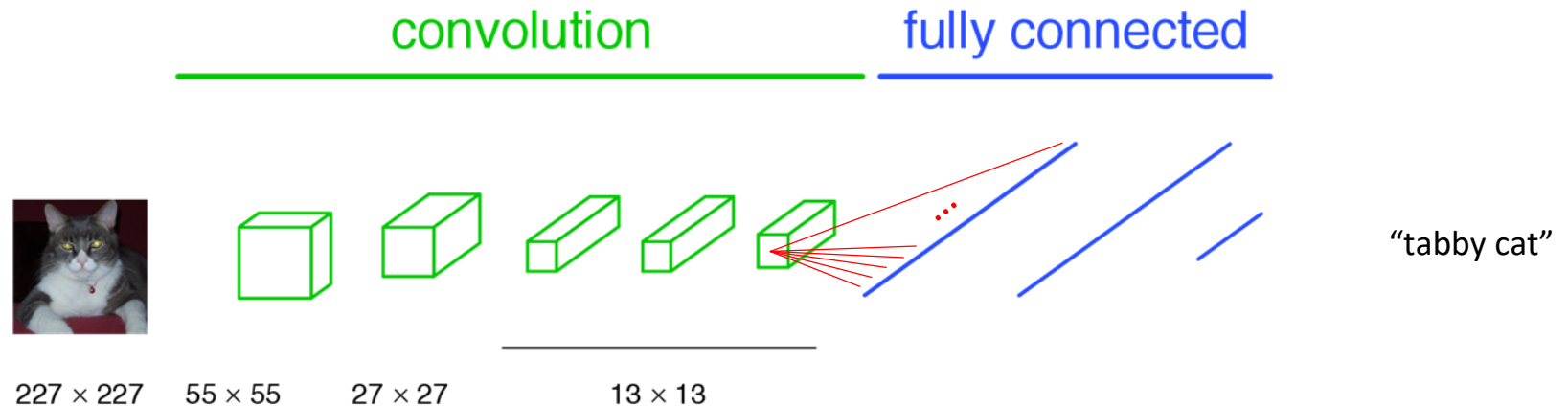


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

# A classification network...



The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

AlexNet: 256 filters over  $6 \times 6$  response map

Each 2,359,296 response is attached to one of 4096 perceptrons, leading to 37 mil params.





**Yann LeCun**

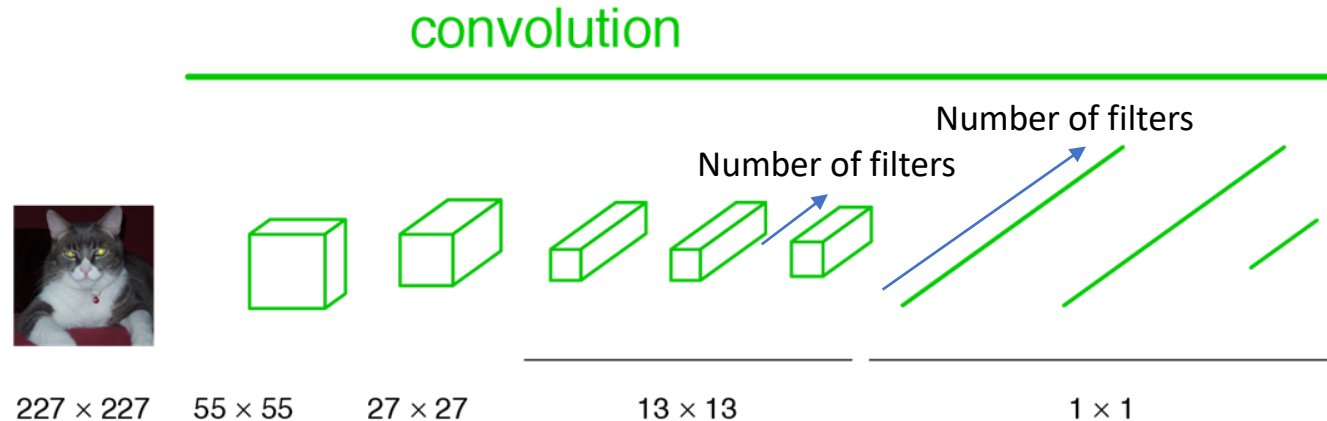
6 April 2015 · 

 Follow



In Convolutional Nets, there is no such thing as "fully-connected layers". There are only convolution layers with 1x1 convolution kernels and a full connection table.

# Convolutionalization

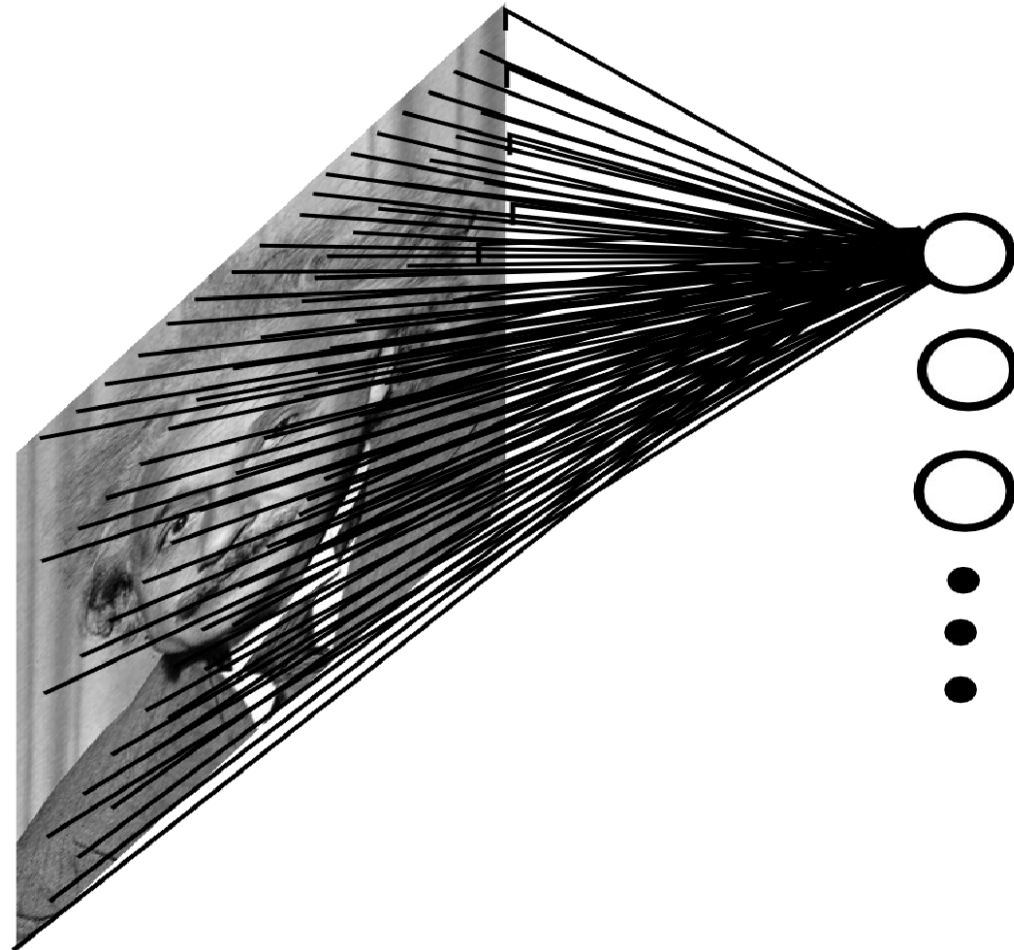


1x1 convolution operates across all filters in the previous layer, and is slid across all positions.

# Back to the fully-connected perceptron...

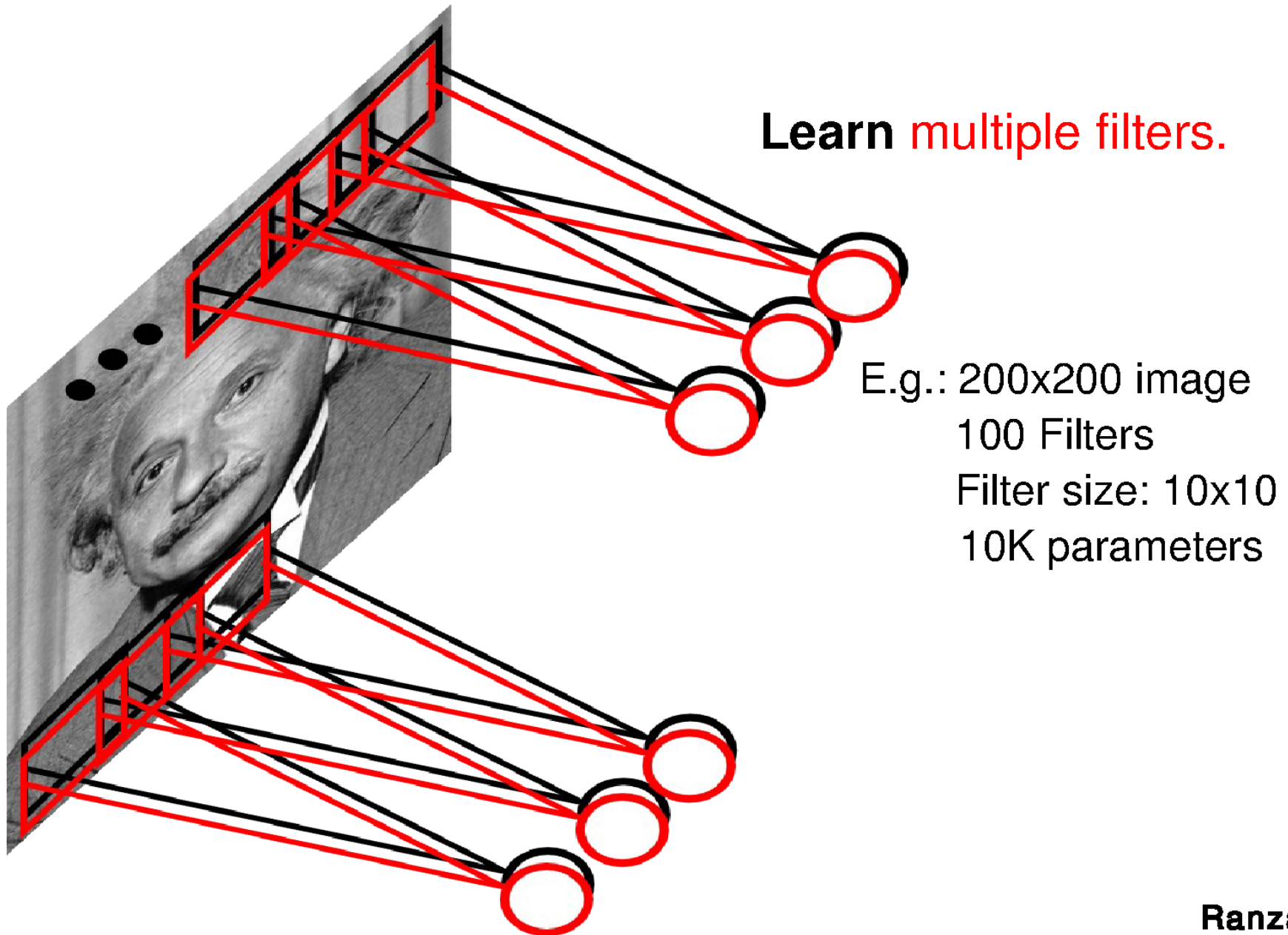
$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x \leq 0 \\ 1 & \text{if } w \cdot x > 0 \end{cases}$$

$$w \cdot x \equiv \sum_j w_j x_j$$



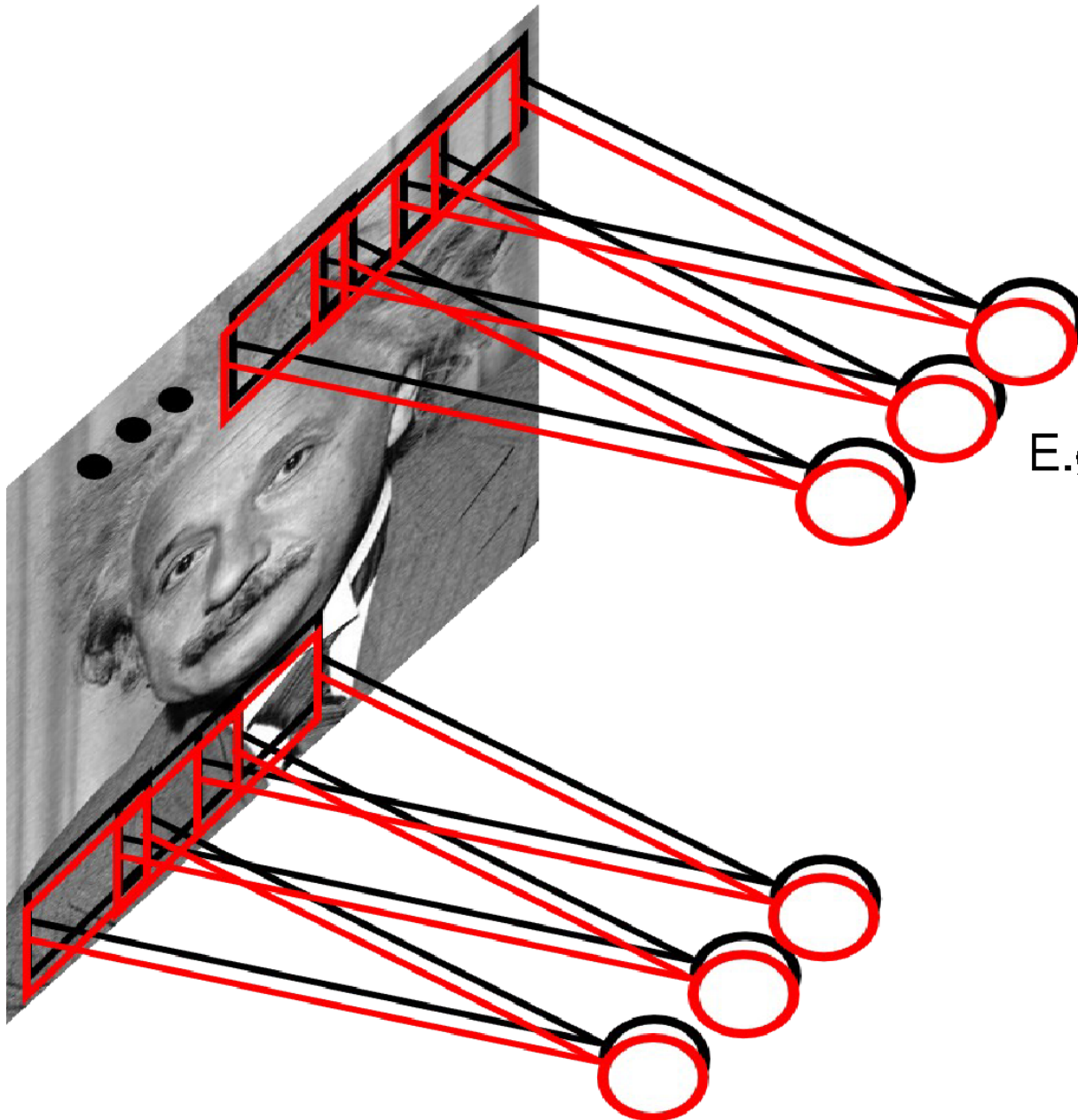
Perceptron is connected to every value in the previous layer (across all channels; 1 visible).

# Convolutional Layer



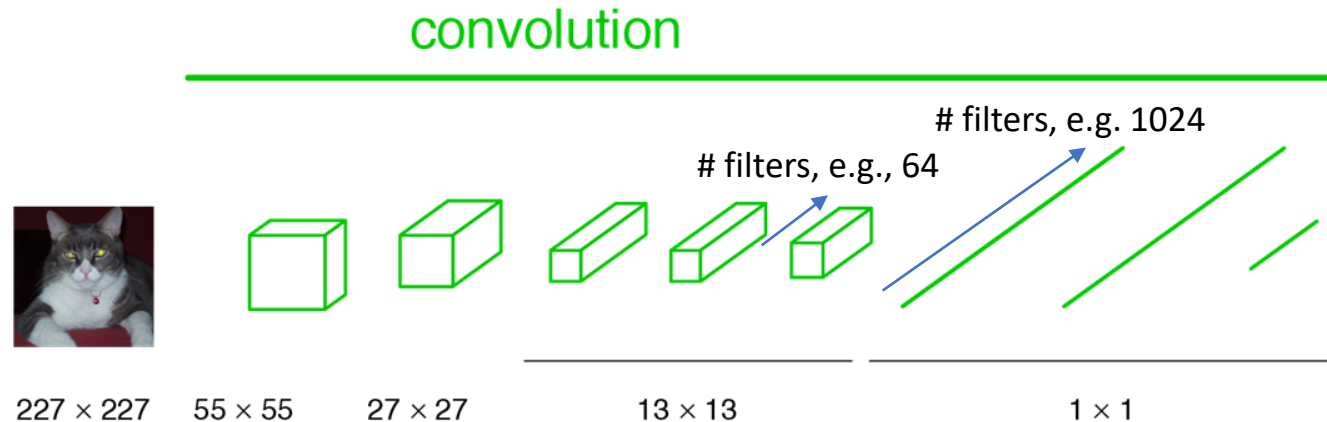


# Convolutional Layer



E.g.: 200x200 image  
100 Filters  
Filter size: 1x1  
100 parameters

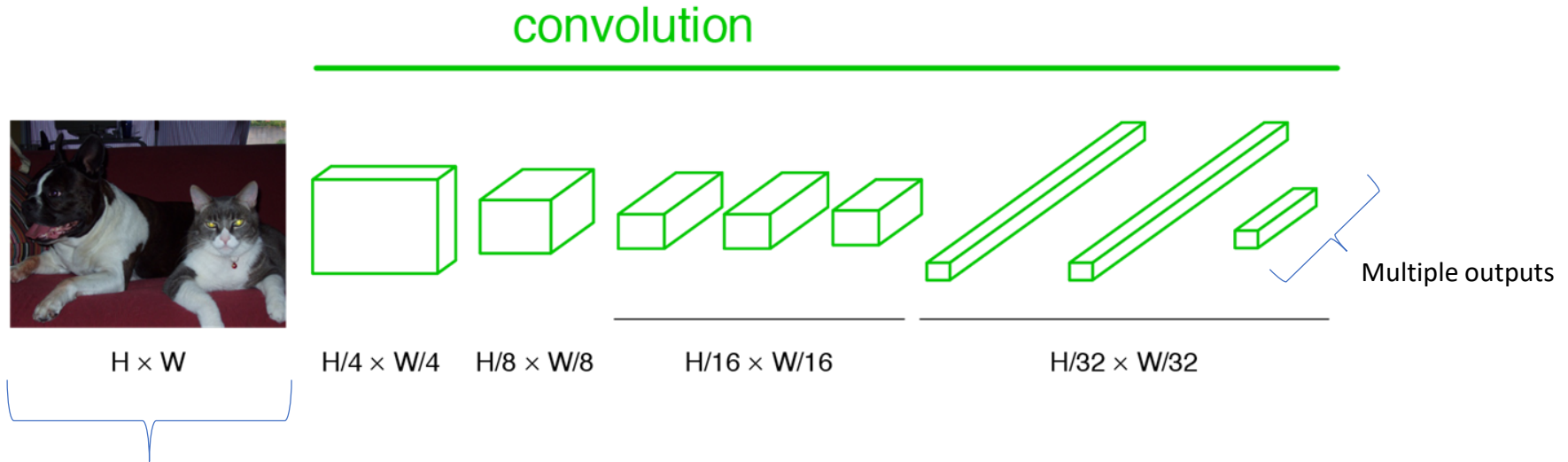
# Convolutionalization



1x1 convolution operates across all filters in the previous layer, and is slid across all positions.

e.g., 64x1x1 kernel, with shared weights over 13x13 output, x1024 filters = 11mil params.

# Becoming fully convolutional



Arbitrary-  
sized image

When we turn these operations into a convolution, the  $13 \times 13$  just becomes another parameter and our output size adjust dynamically.

Now we have a *vector/matrix* output, and our network acts itself like a complex filter.

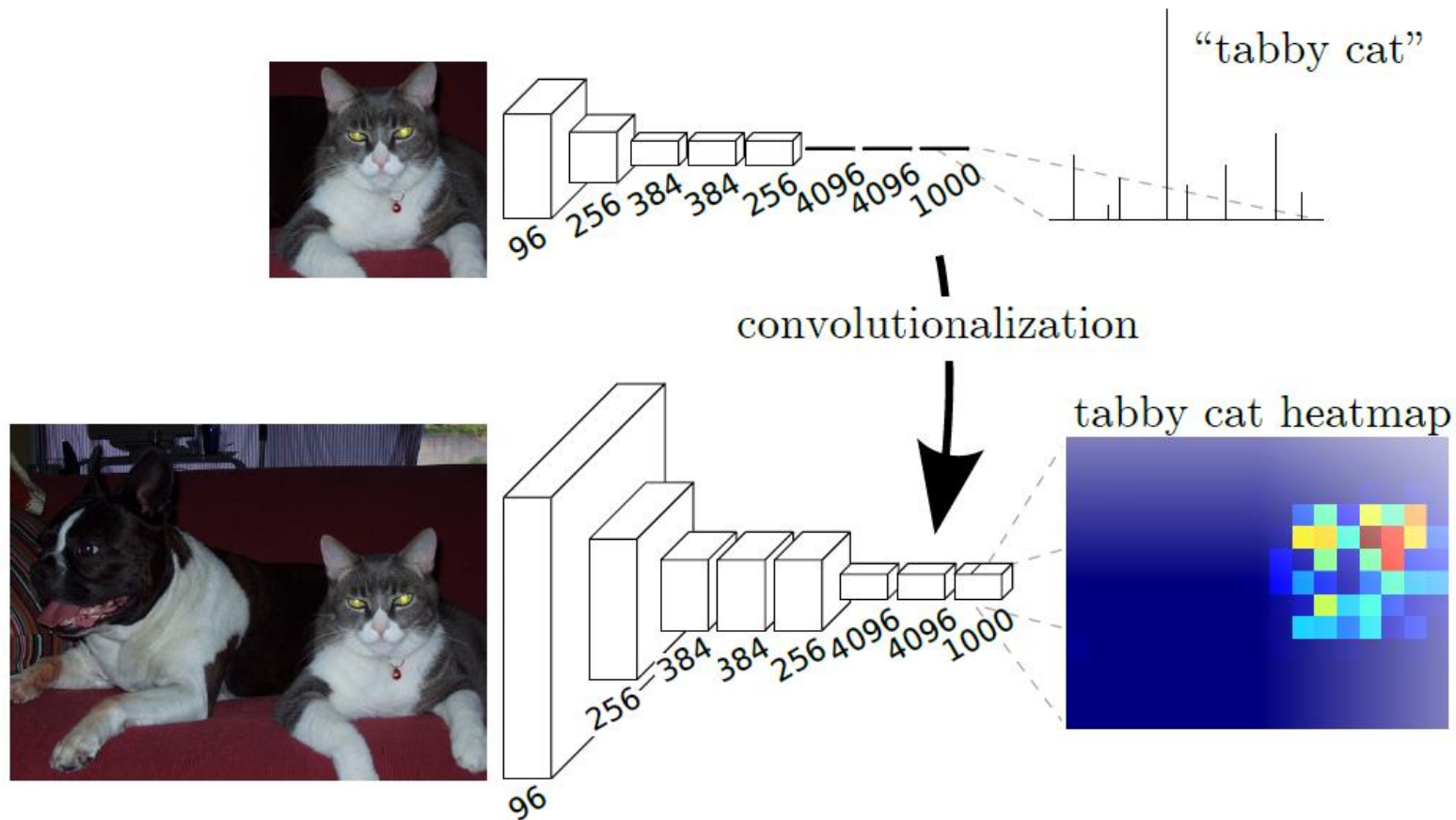
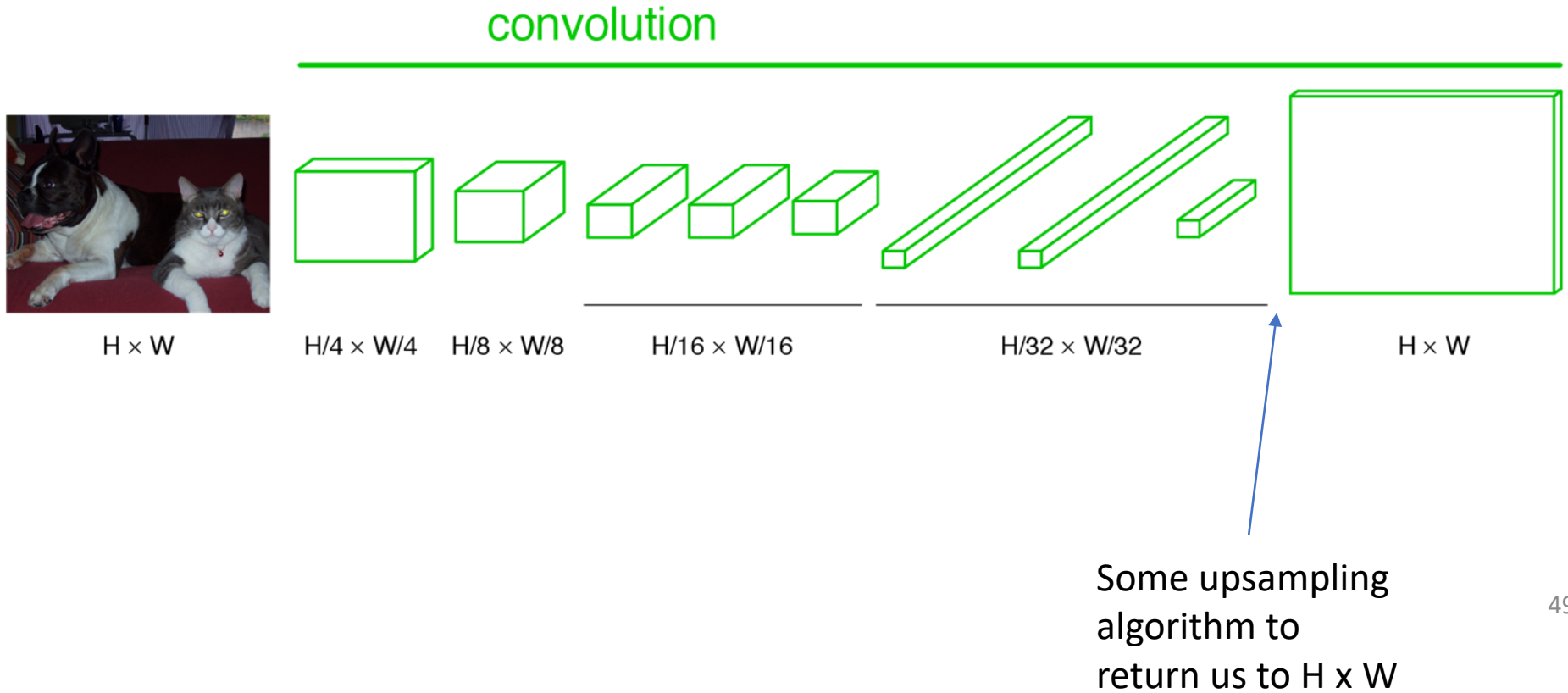


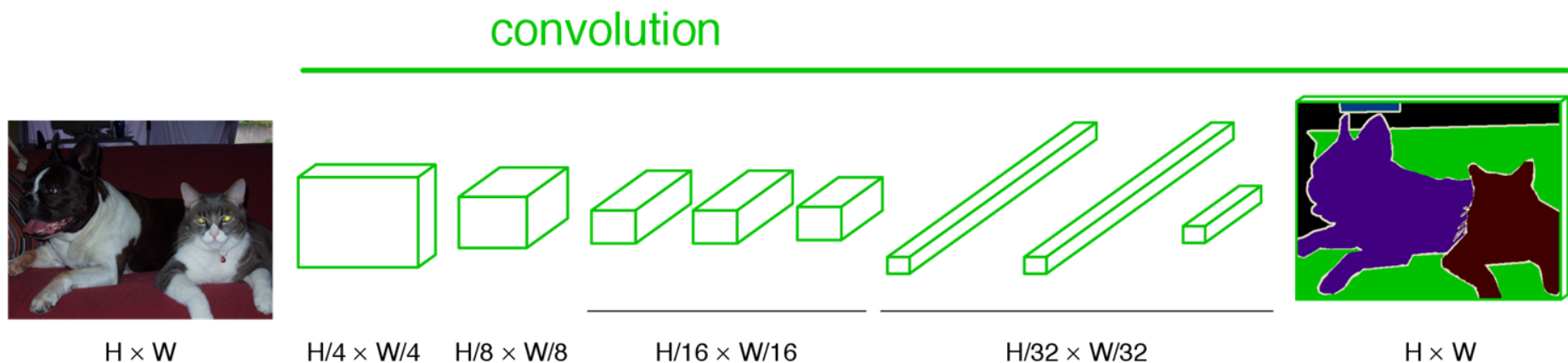
Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.



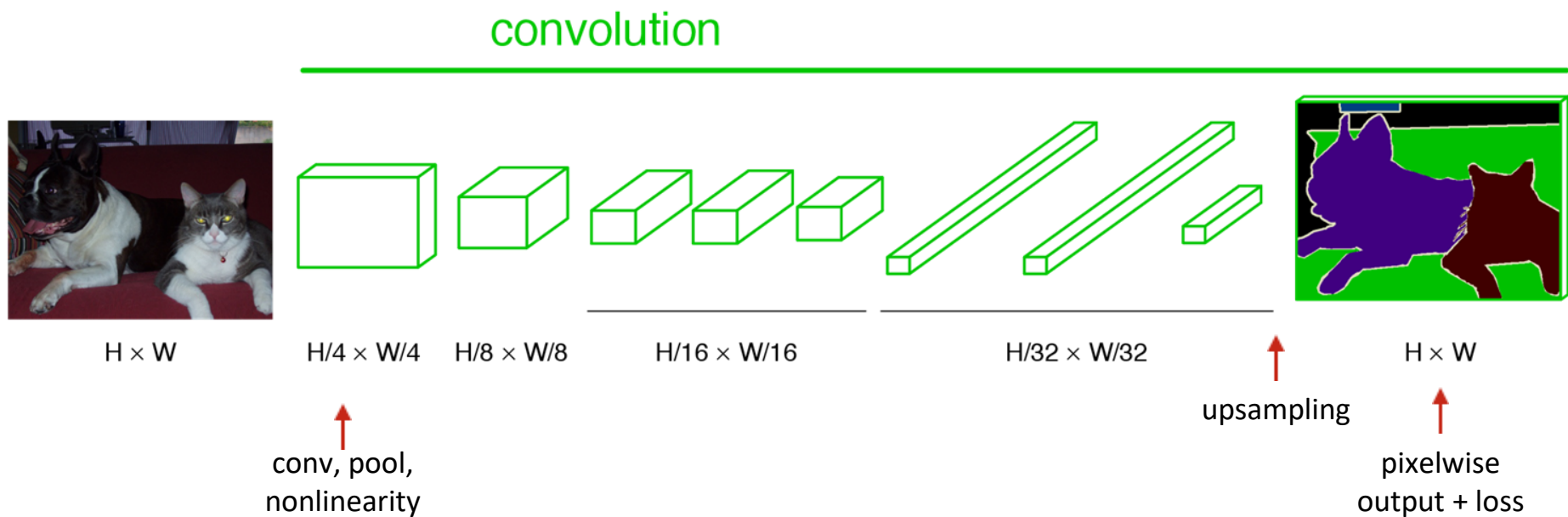
# Upsampling the output



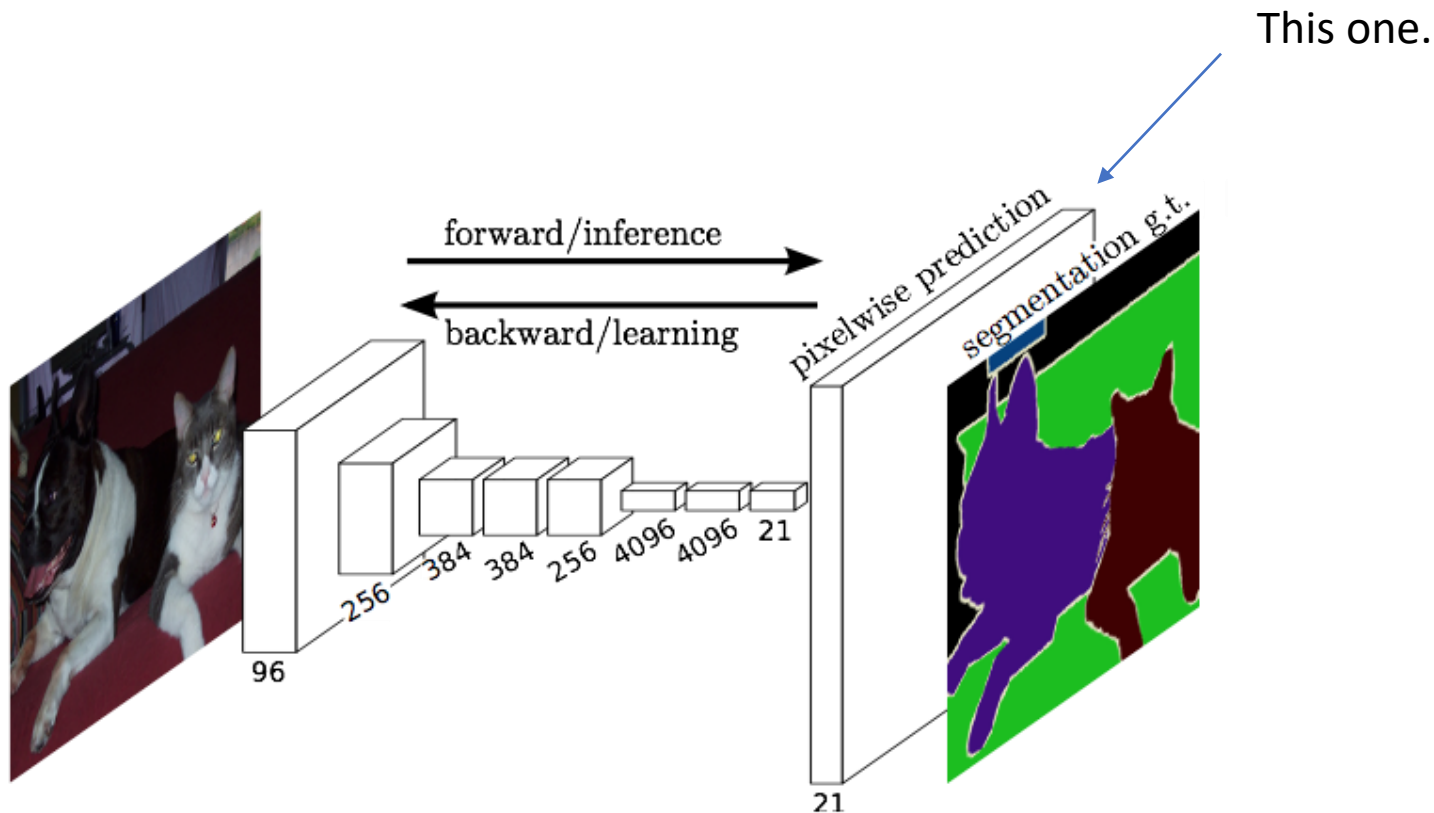
# End-to-end, pixels-to-pixels network



# End-to-end, pixels-to-pixels network



# What is the upsampling layer?

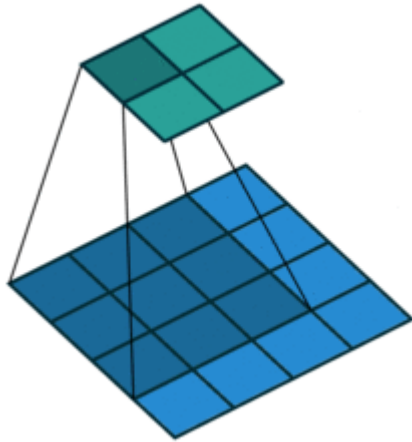


Hint: it's actually an upsampling\_network\_

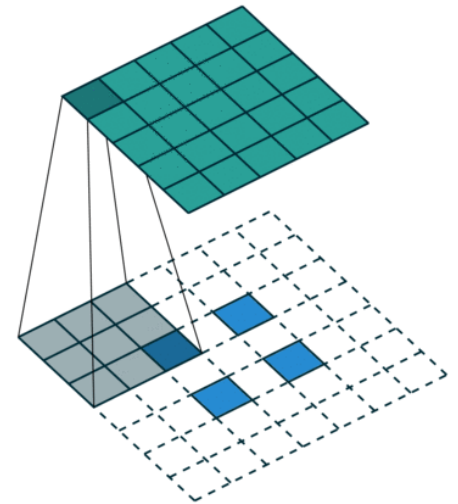
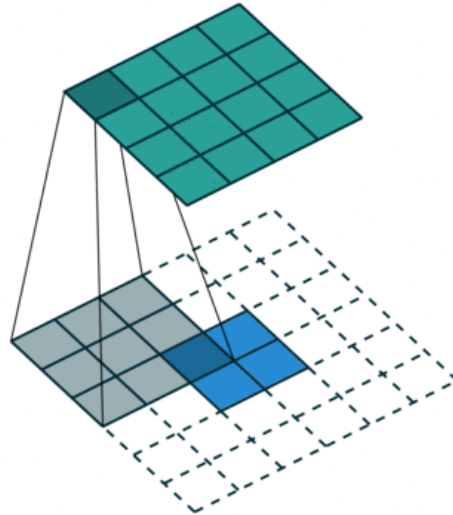


# Upsampling with convolution

Convolution



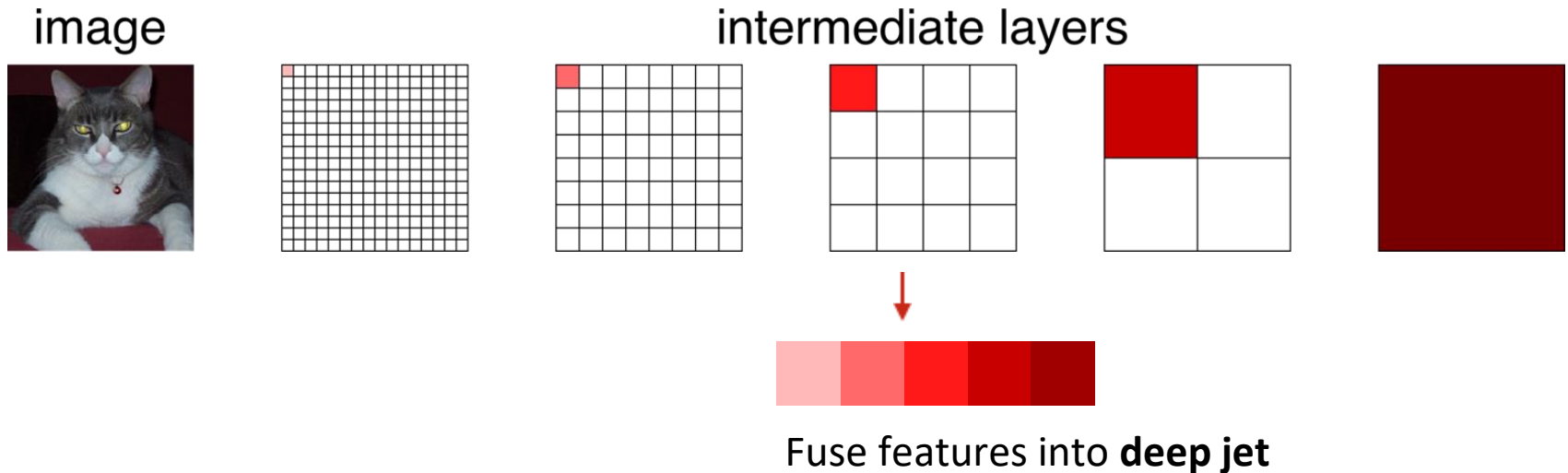
Transposed convolution =  
weighted kernel 'stamp'



Often called “deconvolution”,  
but not actually the deconvolution  
that we previously saw in deblurring ->  
that is division in the Fourier domain.

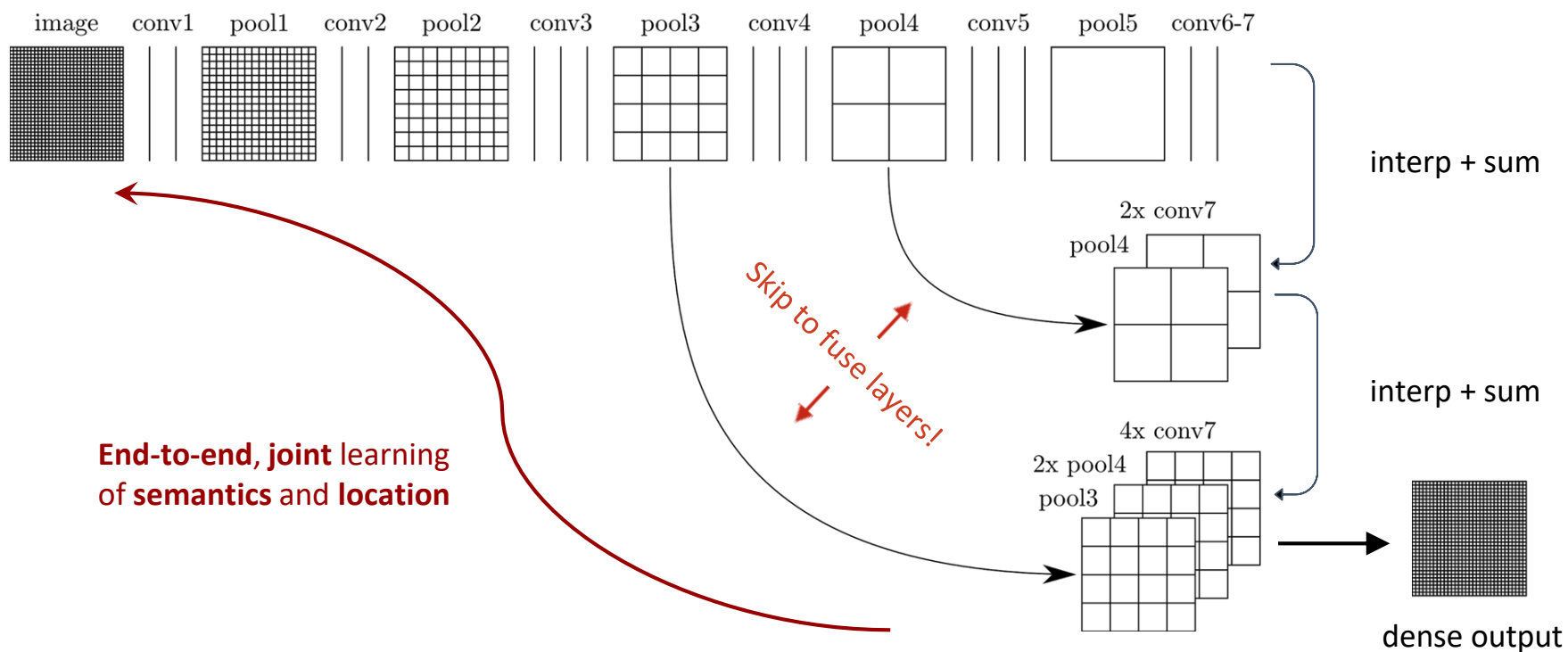
# Spectrum of deep features

Combine *where* (local, shallow) with *what* (global, deep)



(cf. Hariharan et al. CVPR15 “hypercolumn”)

# Learning upsampling kernels with skip layer refinement



# Skip layer refinement

input image



stride 32



no skips

stride 16



1 skip

stride 8



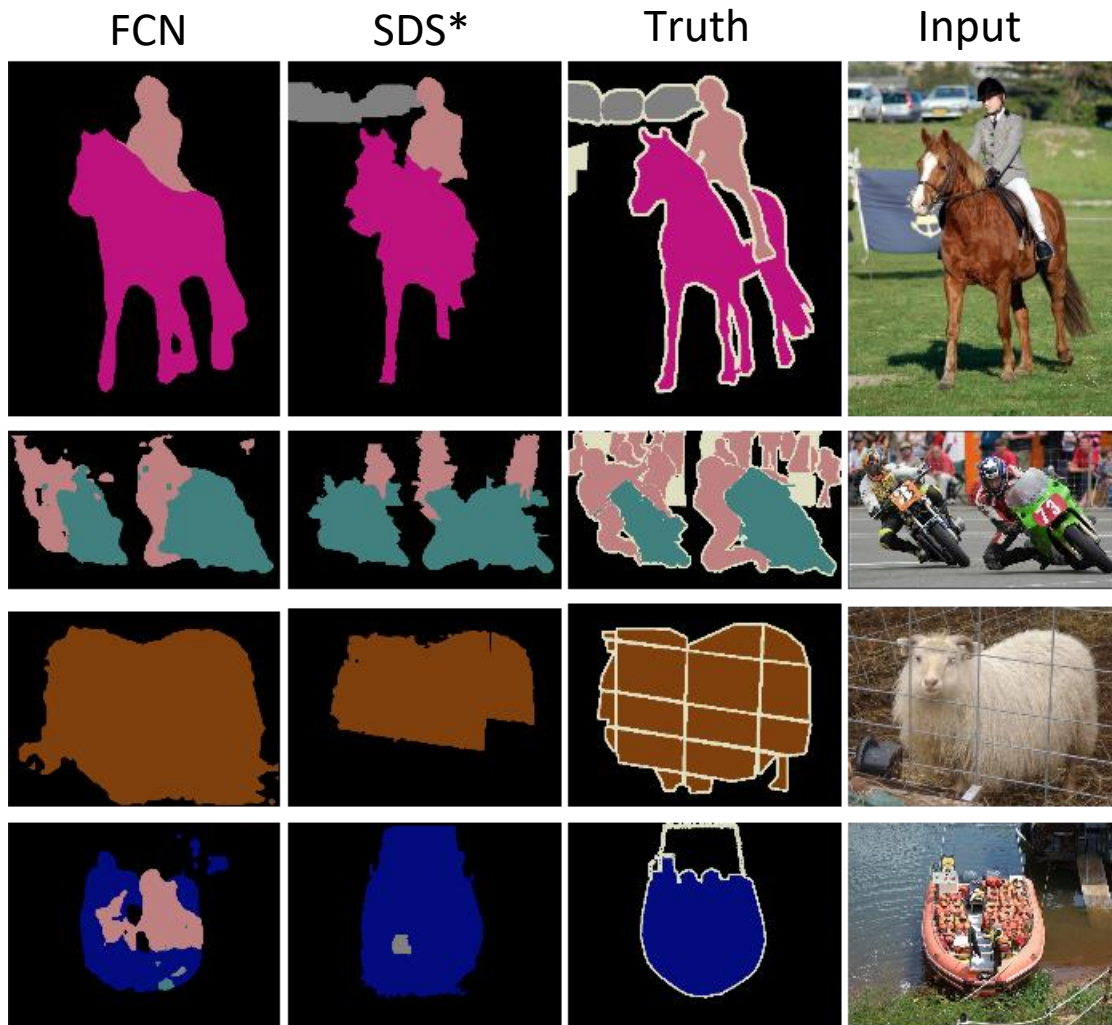
2 skips

ground truth





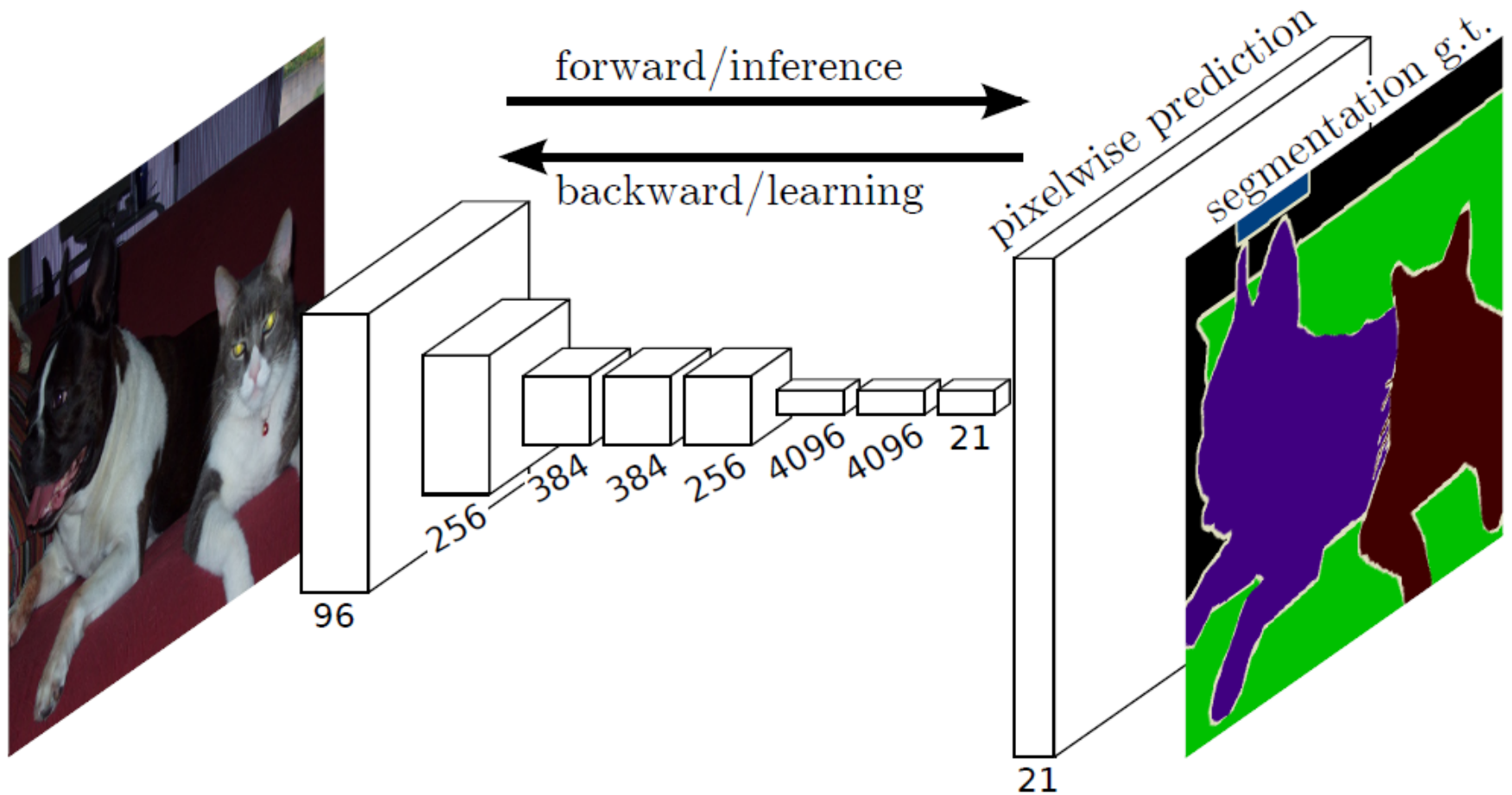
# Results



Relative to prior state-of-the-art SDS:

- 30% relative improvement for mean IoU
- 286× faster

\*Simultaneous Detection and Segmentation  
Hariharan et al. ECCV14



What can we do with an FCN?

# How much can an image tell about its geographic location?



6 million geo-tagged Flickr images

<http://graphics.cs.cmu.edu/projects/im2gps/>

[im2gps](#) (Hays & Efros, CVPR 2008)



# Nearest Neighbors according to gist + bag of SIFT + color histogram + a few others



Paris



Paris



Paris



Paris



Paris



Paris



Paris



Madrid



Rome



Paris



Cuba



Paris



Paris



Poland



Paris



Paris





# PlaNet - Photo Geolocation with Convolutional Neural Networks

Tobias Weyand, Ilya Kostrikov, James Philbin

ECCV 2016

# Discretization of Globe

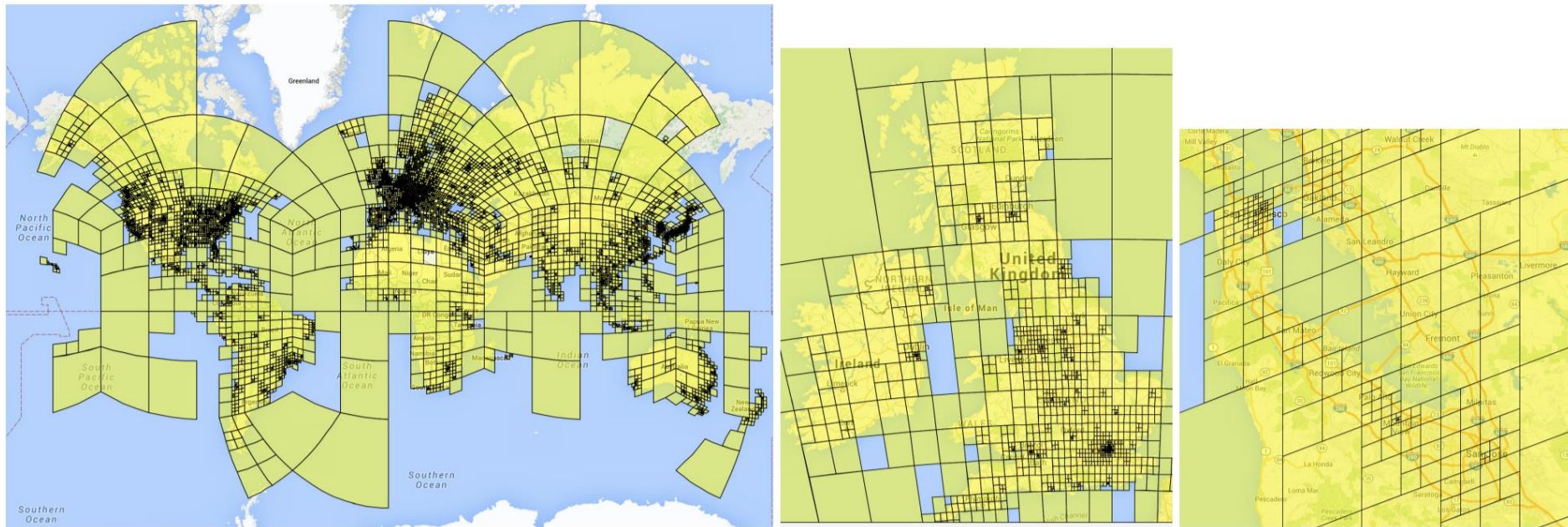


Figure 2. Left: Adaptive partitioning of the world into 26,263 S2 cells. Right: Detail views of Great Britain and Ireland and the San

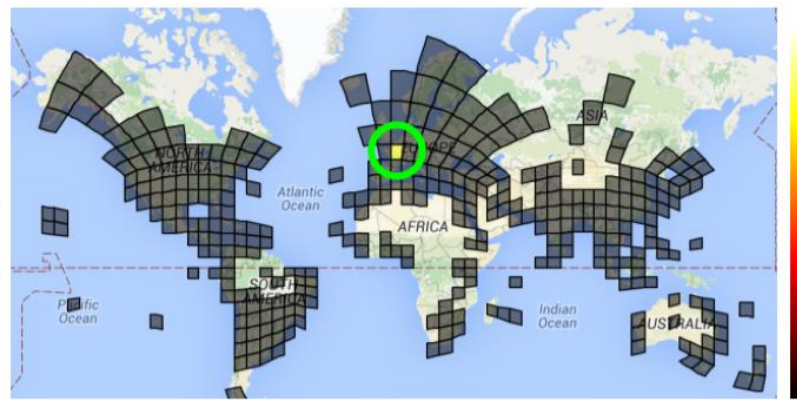
# Network and Training

- Network Architecture: Inception with 97M parameters
- 26,263 “categories” – places in the world
- 126 Million Web photos
- 2.5 months of training on 200 CPU cores





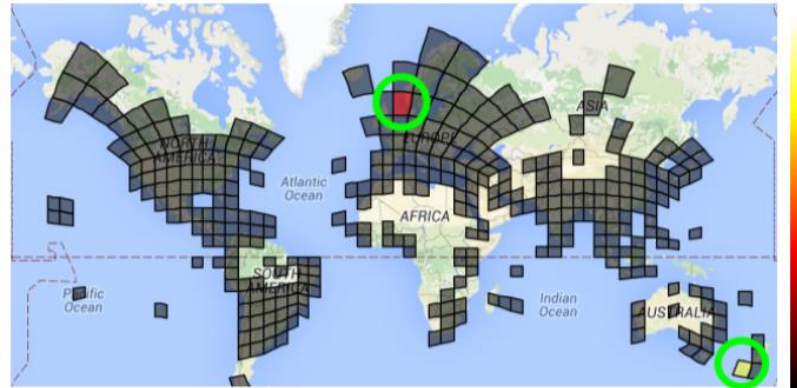
Photo CC-BY-NC by stevekc



(a)



Photo CC-BY-NC by edwin.11



(b)



Photo CC-BY-NC by jonathanfh







Namibia / Botswana



Photo by jamie.loveclark / CC BY NC Photo by MongoosePhotography / CC BY NC



Photo by Mister-E / CC BY NC



Photo by dalangalma / CC BY NC



Photo by slamjack / CC BY NC



Kauai, Hawaii



Photo by ryan + sarah / CC BY NC



Photo by stuartchambers / CC BY NC



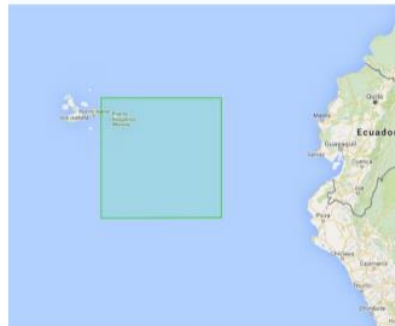
Photo by samgrover / CC BY NC



Photo by steuben / CC BY NC



Photo by steve-stevens / CC BY NC



Galapagos Islands



Photo by p.j.k. / CC BY NC



Photo by victor408 / CC BY NC



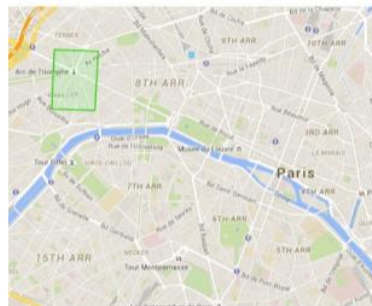
Photo by Domen jakus / CC BY NC



Photo by cvanholder / CC BY NC



Photo by rwoan / CC BY NC



Paris



Photo by feliven / CC BY NC



Photo by fred\_v / CC BY NC



Photo by Turansa Tours / CC BY NC



Photo by JA\_FS / CC BY NC



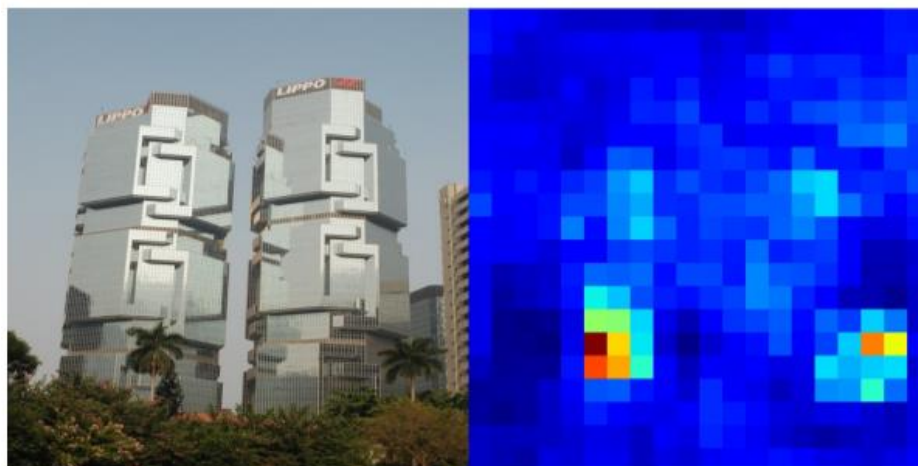
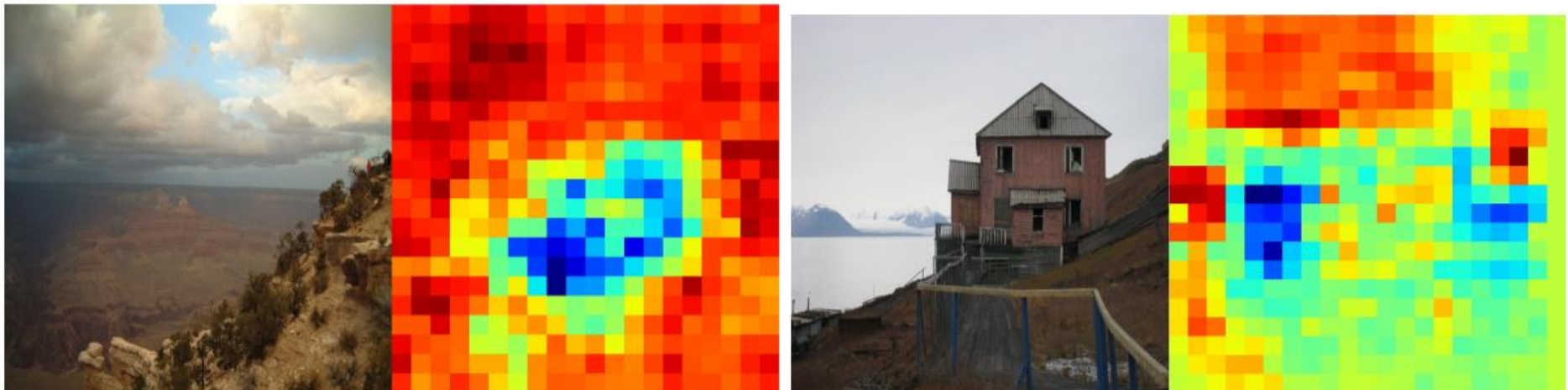
Photo by CedEm photographs / CC BY NC

# PlaNet vs im2gps (2008, 2009)

<b>Method</b>	<b>Street 1 km</b>	<b>City 25 km</b>	<b>Region 200 km</b>	<b>Country 750 km</b>	<b>Continent 2500 km</b>
Im2GPS (orig) [17]		12.0%	15.0%	23.0%	47.0%
Im2GPS (new) [18]	2.5%	21.9%	32.1%	35.4%	51.9%
PlaNet	<b>8.4%</b>	<b>24.5%</b>	<b>37.6%</b>	<b>53.6%</b>	<b>71.3%</b>

<b>Method</b>	<b>Manmade Landmark</b>	<b>Natural Landmark</b>	<b>City Scene</b>	<b>Natural Scene</b>	<b>Animal</b>
Im2GPS (new)	61.1	37.4	3375.3	5701.3	6528.0
PlaNet	74.5	61.0	212.6	1803.3	1400.0

# Spatial support for decision

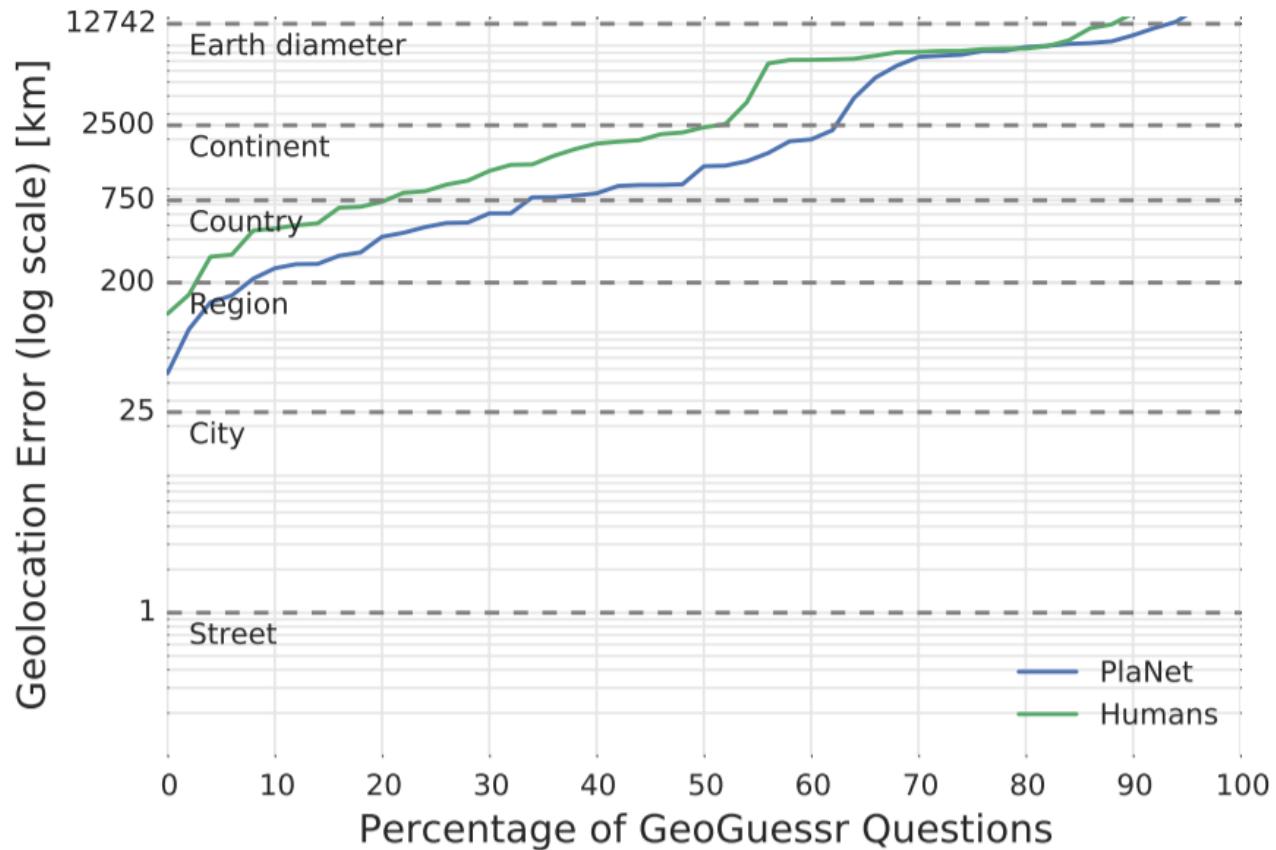




# PlaNet vs Humans



# PlaNet vs. Humans





# PlaNet summary

- Very fast geolocalization method by categorization.
- Uses far more training data than previous work (im2gps)
- Better than humans!