

Advanced Image Processing

Exploring Architectures

Big space of designs!

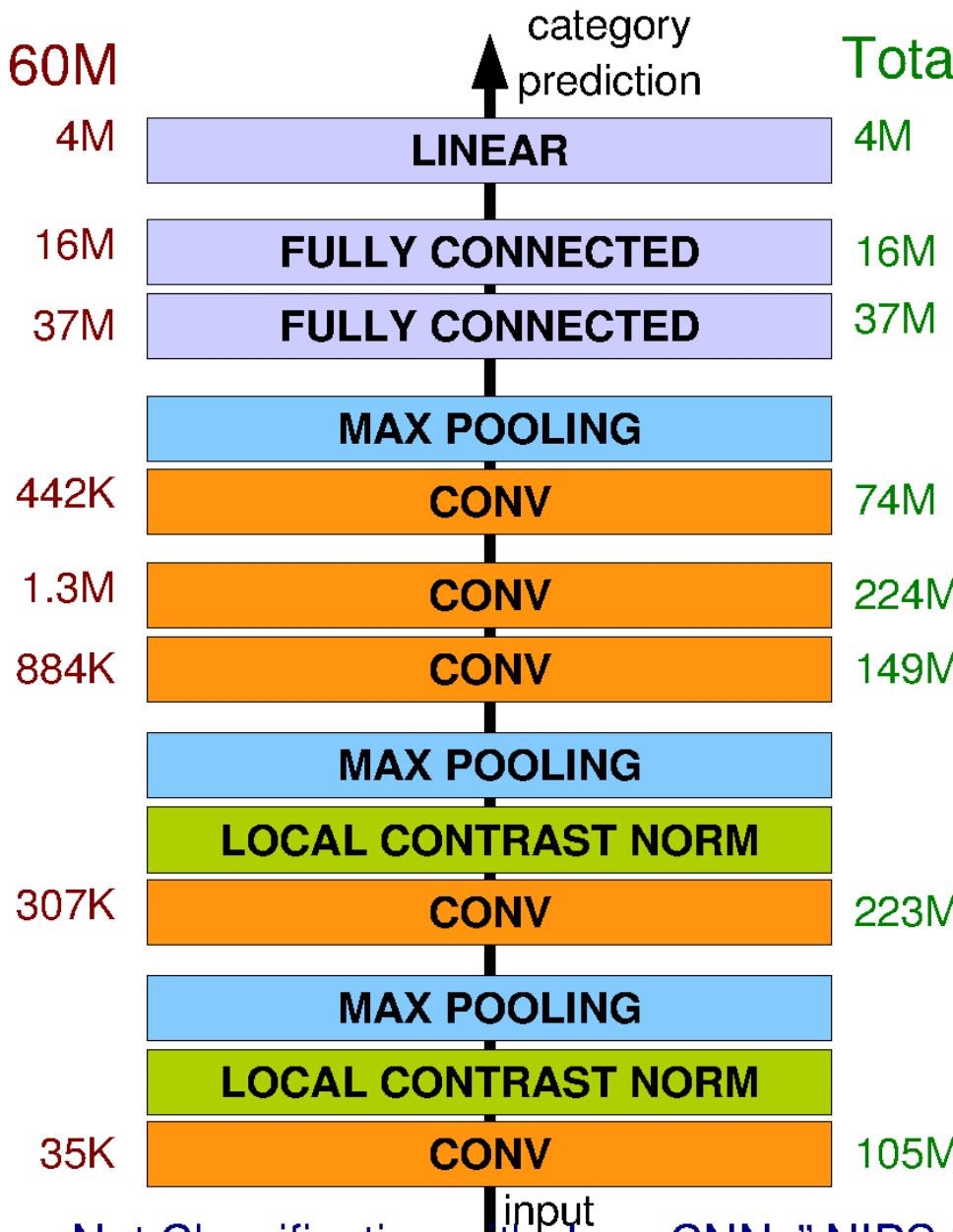
We still don't even know how many layers we need.

Deep-CNN

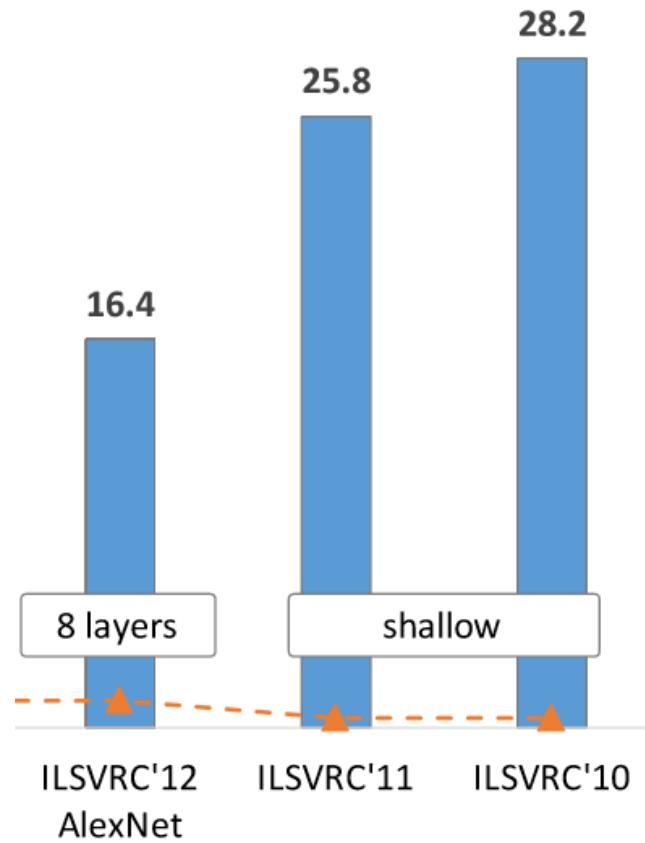
AlexNet

Total nr. params: 60M

Total nr. flops: 832M



Revolution of Depth



ImageNet Classification top-5 error (%)



Beyond AlexNet

VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

Karen Simonyan & Andrew Zisserman 2015

These are the pre-trained “VGG” networks
that you use in assignment 4

VGG = Visual Geometry Group at Oxford, UK

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224×224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

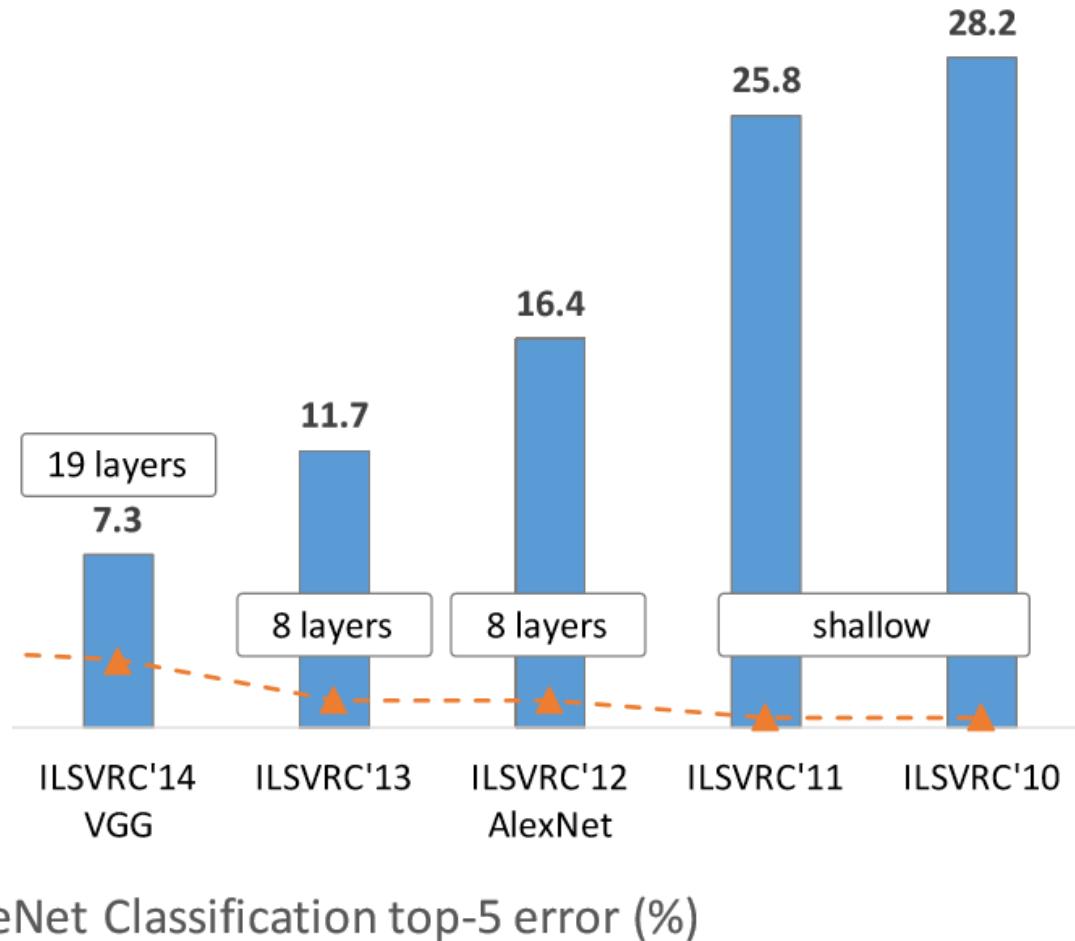
Table 2: **Number of parameters** (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

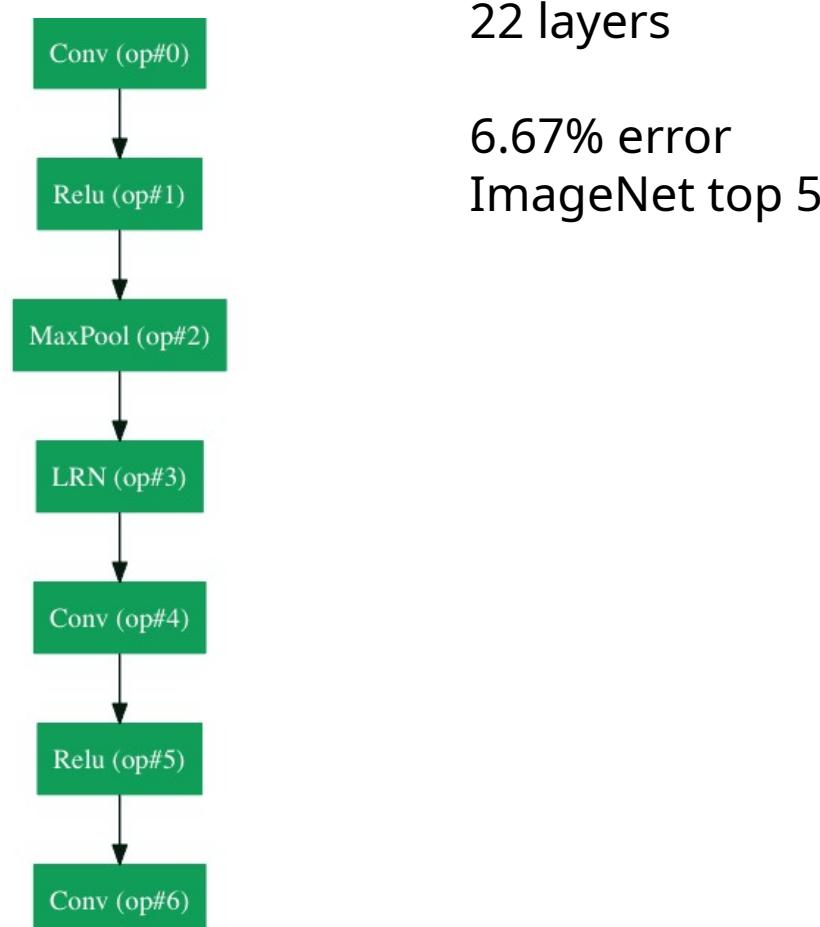
Table 4: ConvNet performance at multiple test scales.

ConvNet config. (Table 1)	smallest image side		top-1 val. error (%)	top-5 val. error (%)
	train (S)	test (Q)		
B	256	224,256,288	28.2	9.6
C	256	224,256,288	27.7	9.2
	384	352,384,416	27.8	9.2
	[256; 512]	256,384,512	26.3	8.2
D	256	224,256,288	26.6	8.6
	384	352,384,416	26.5	8.6
	[256; 512]	256,384,512	24.8	7.5
E	256	224,256,288	26.9	8.7
	384	352,384,416	26.7	8.6
	[256; 512]	256,384,512	24.8	7.5

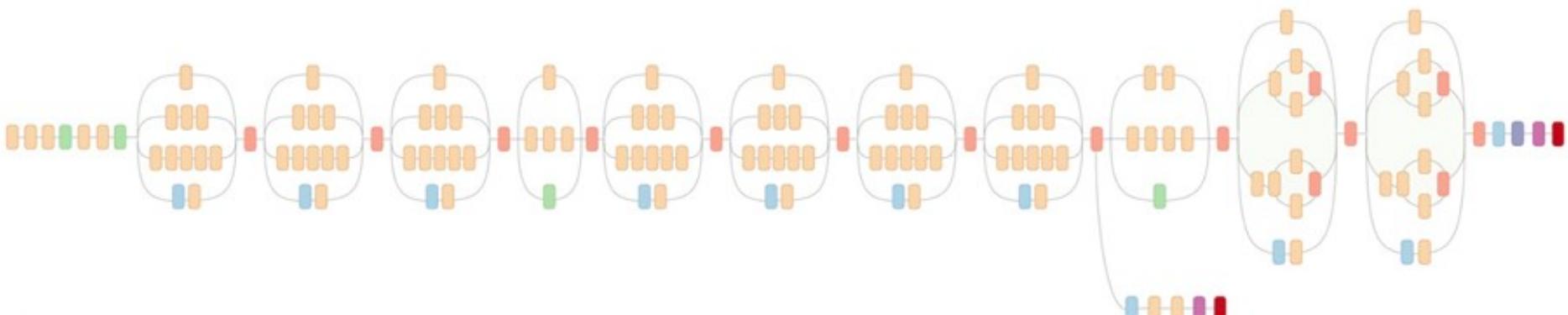
Revolution of Depth



Google LeNet (inception 2014)



Google LeNet (inception 2014) – another view



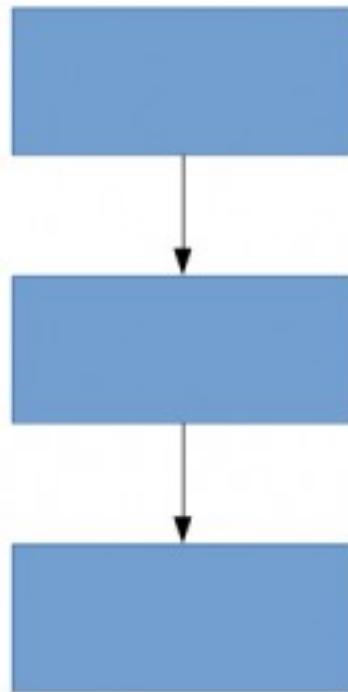
[Another view of GoogLeNet's architecture.](#)

Inception model idea:

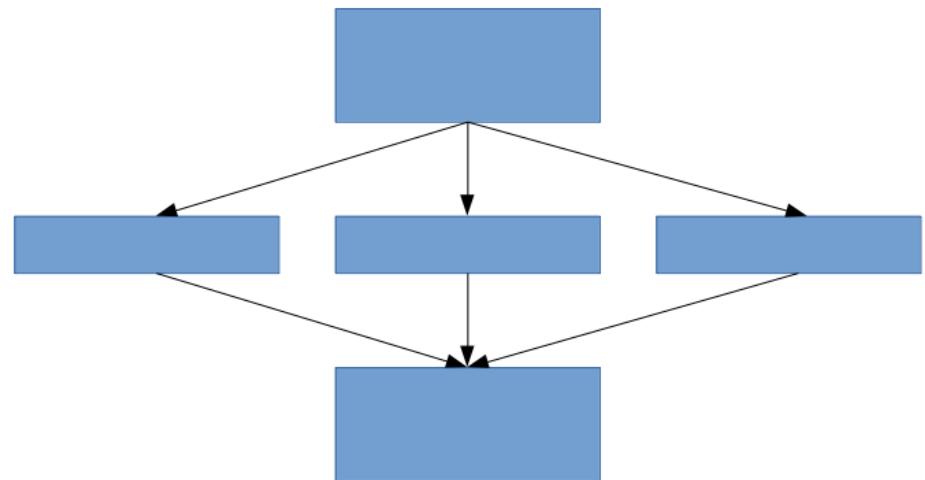
Go deeper while maintaining computational cost:

Achieved by replacing fully connected architectures with sparsely connected ones

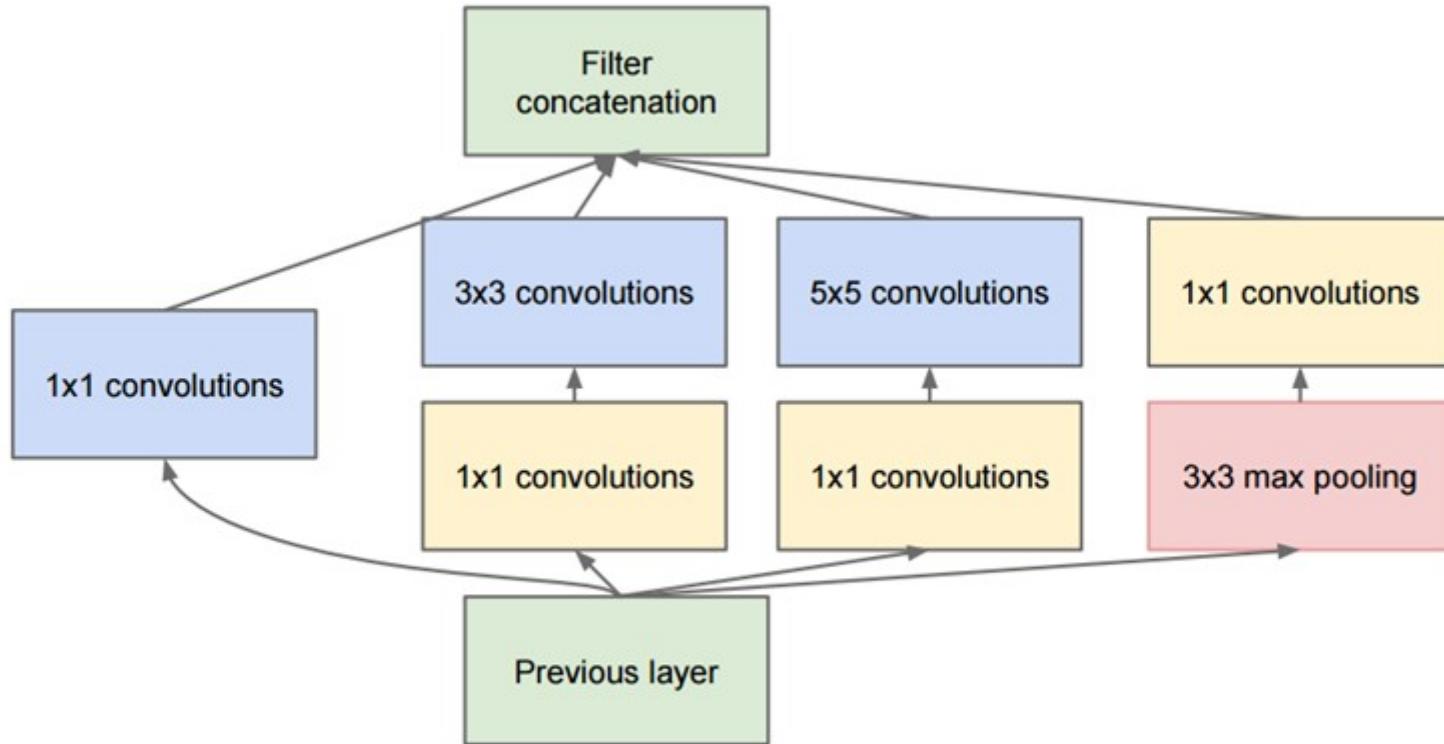
Densely
connected



Sparsely
connected

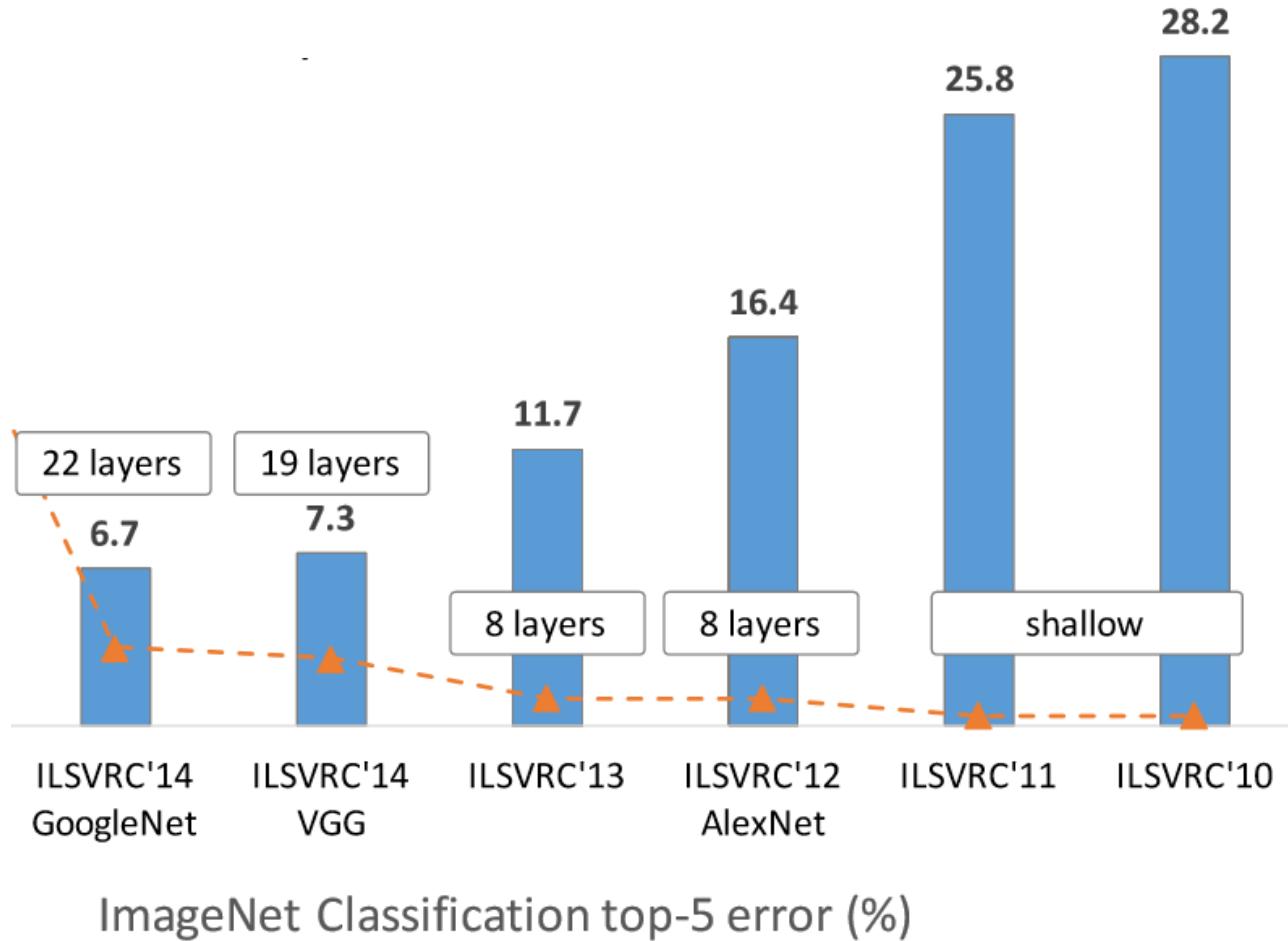


Sparsely connected layers: Parallel layers and 1x1 convolutions



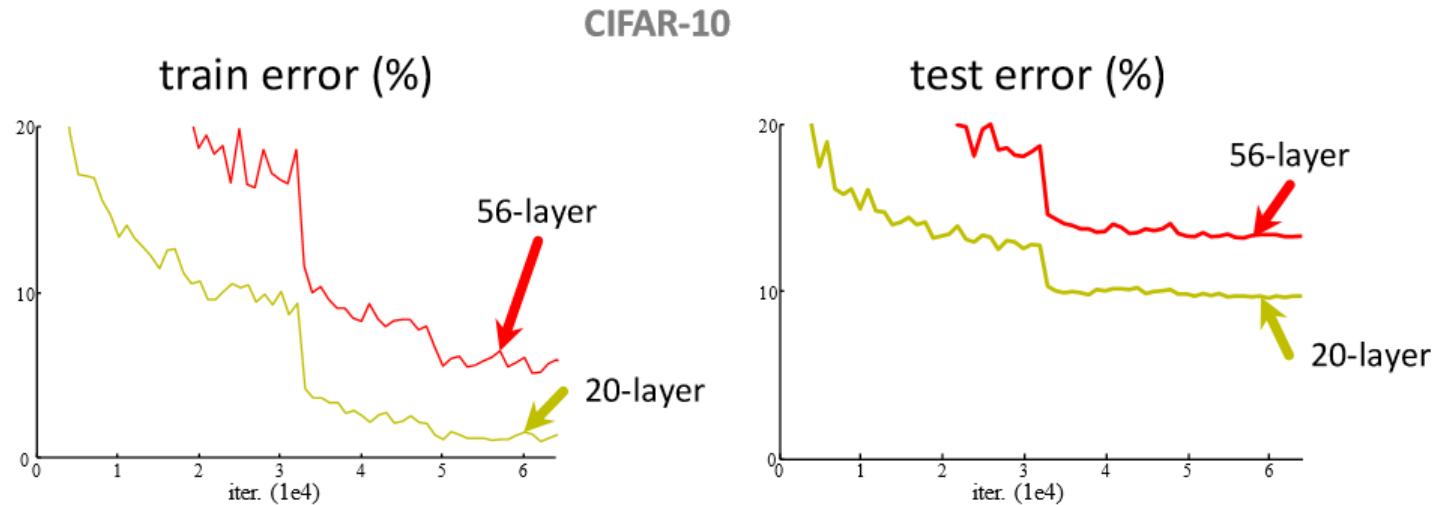
Full Inception module

Revolution of Depth



Problems with going deeper

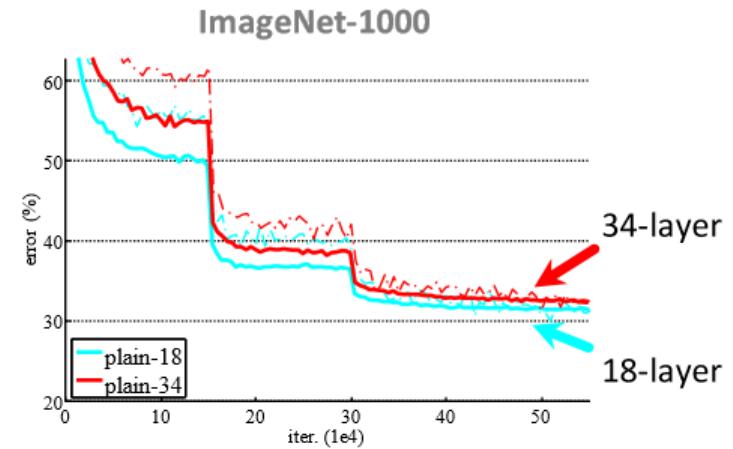
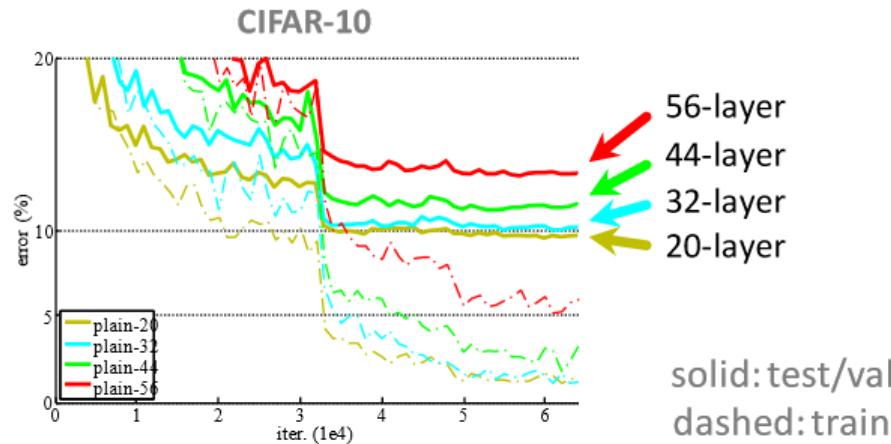
Simply stacking layers?



- Plain nets: stacking 3x3 conv layers...
- 56-layer net has **higher training error** and test error than 20-layer net

Problems with going deeper

Simply stacking layers?



- “Overly deep” plain nets have **higher training error**
- A general phenomenon, observed in many datasets

Why simply stacking layers does not work:

Vanishing/exploding gradient problem

Backpropagation:

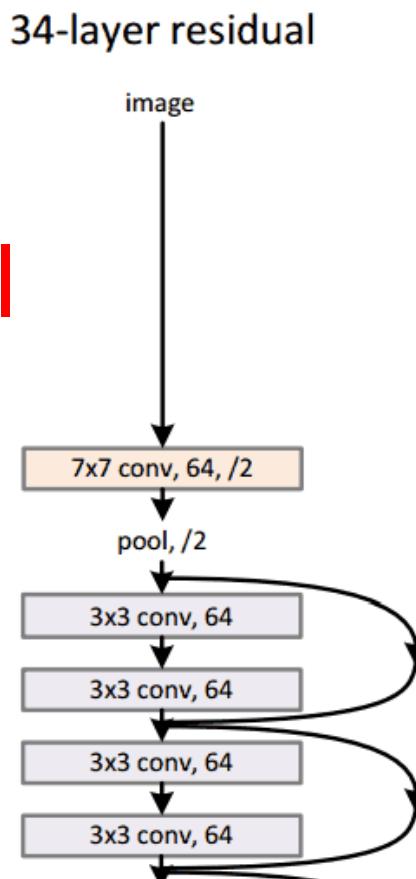
- Compute gradient update for every neuron which was involved in the output across layers

Involves chaining partial derivates over many layers!

- If derivative < 1 , gradient gets smaller and smaller as we go deeper and deeper -> *vanishing gradients!*
- If derivative > 1 , gradient gets larger and larger as we go deeper and deeper -> *exploding gradients!*

ResNet (He et al., 2015)

152-layer model

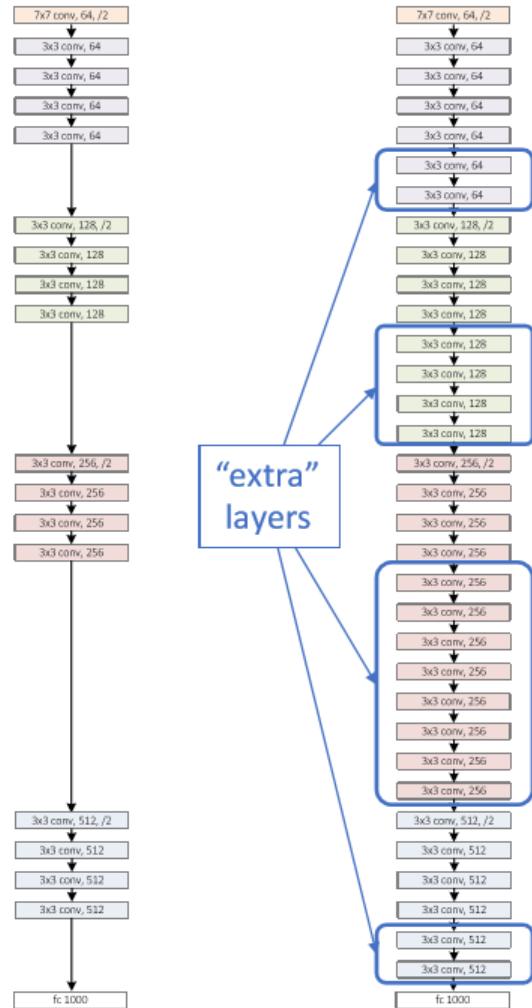


ResNet won ILSVRC 2015 with a top-5 error rate of 3.6%

Depending on their skill and expertise, humans generally hover around a 5-10% error.

Superhuman performance!
But the task is arguably not well defined.

a shallower
model
(18 layers)

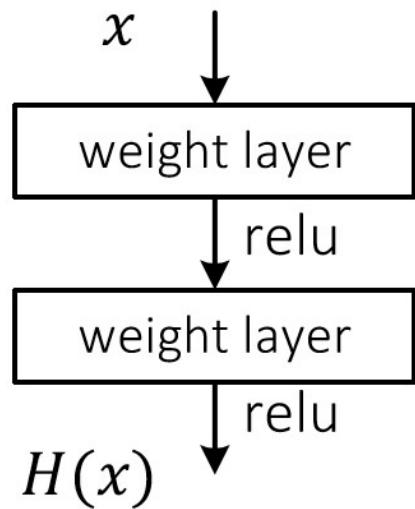


a deeper
counterpart
(34 layers)

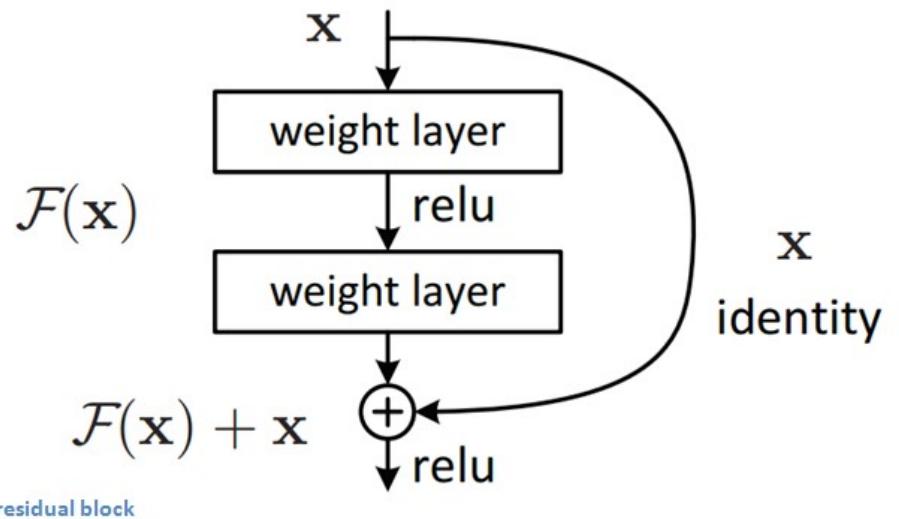
- Richer solution space
- A deeper model should not have **higher training error**
- A solution *by construction*:
 - original layers: copied from a learned shallower model
 - extra layers: set as **identity**
 - at least the same training error
- **Optimization difficulties**: solvers cannot find the solution when going deeper...

Vanishing/exploding gradient solution

Regular Net

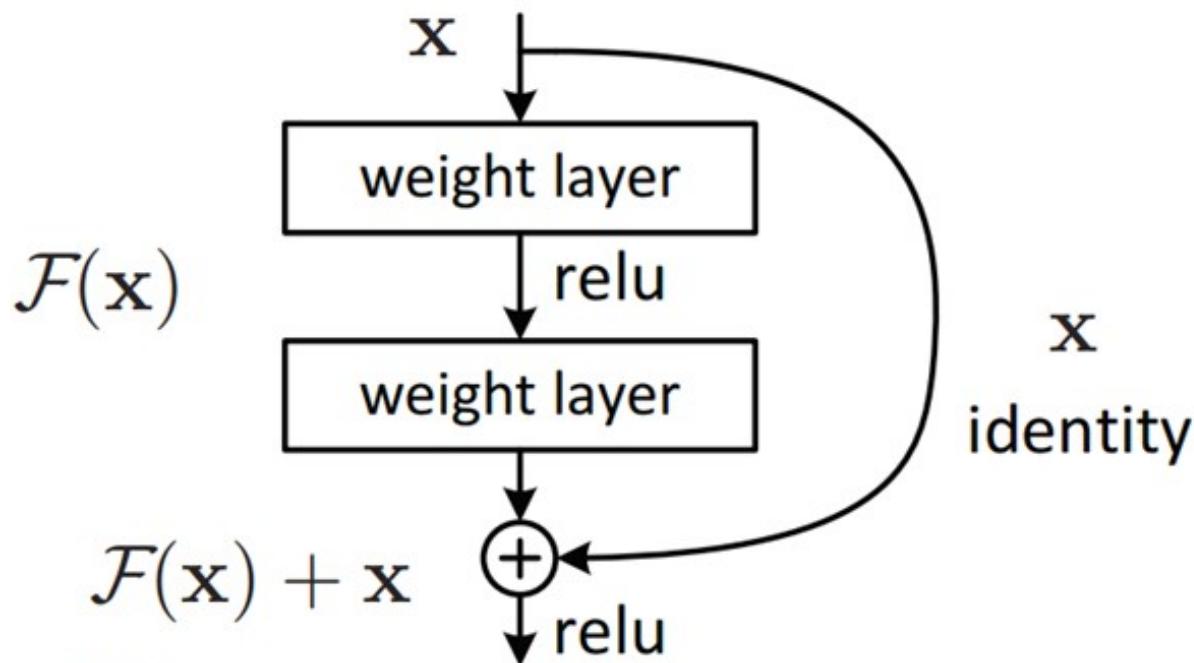


Residual Unit



Residual Unit

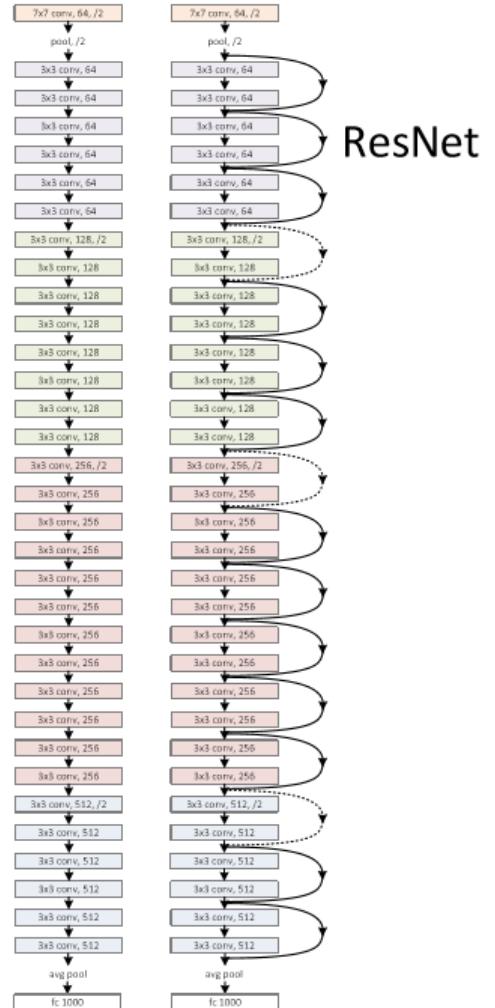
The inputs of a lower layer is made available to a node in a higher layer.



A residual block

Network “Design”

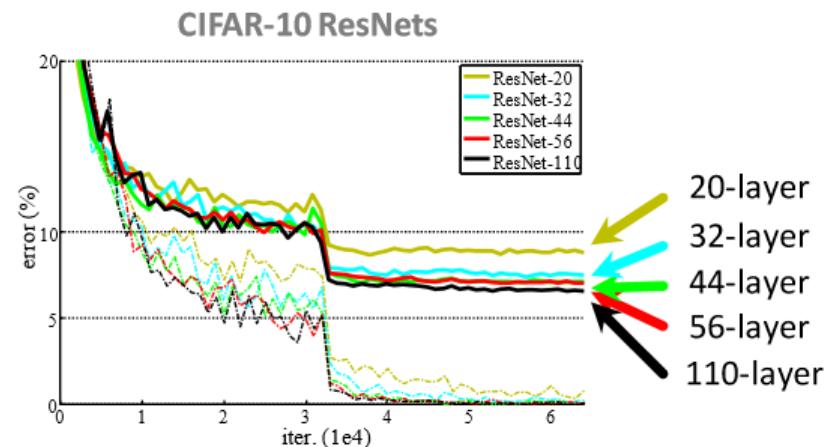
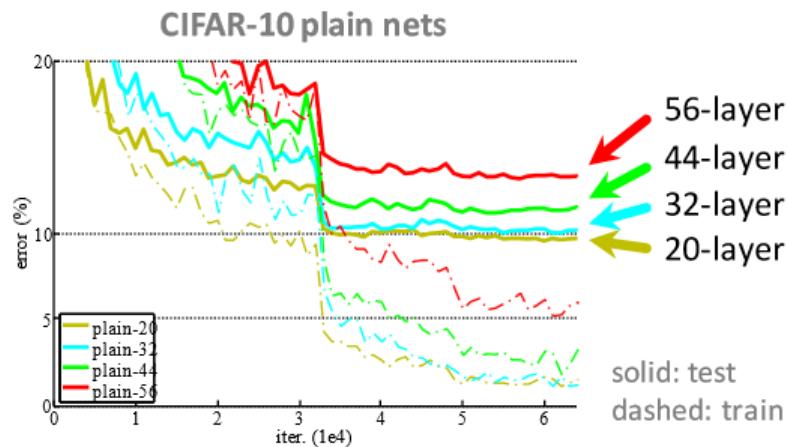
plain net



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

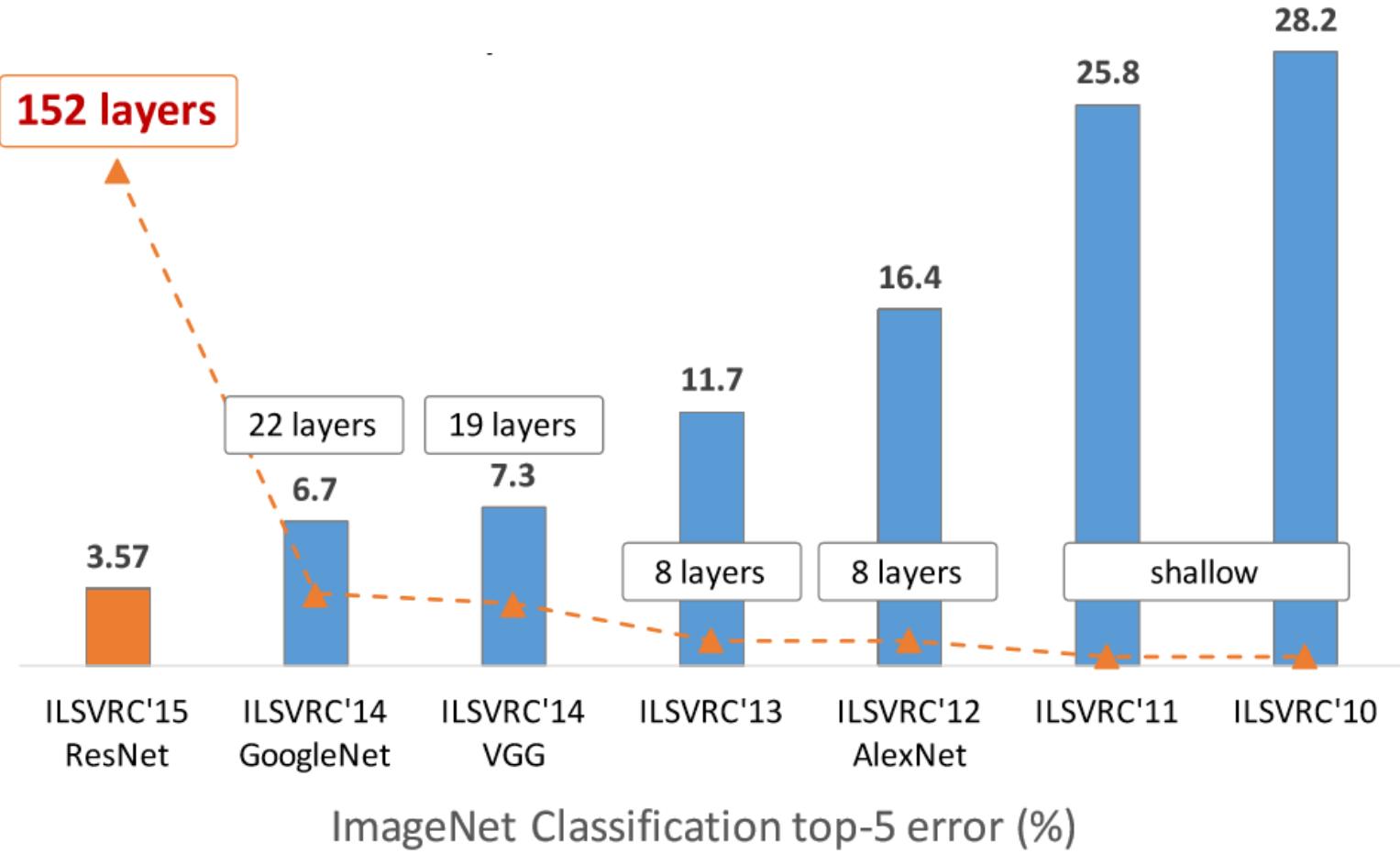
ResNet results

CIFAR-10 experiments



- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

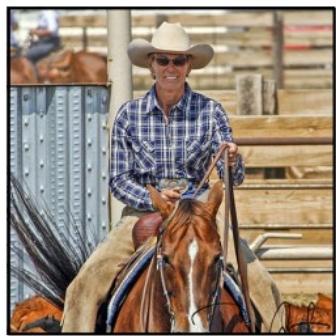
Revolution of Depth



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

R-CNN: Region-based CNN

R-CNN: Region-based CNN



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

4. Classify regions

Figure: Girshick et al.

Reading material:

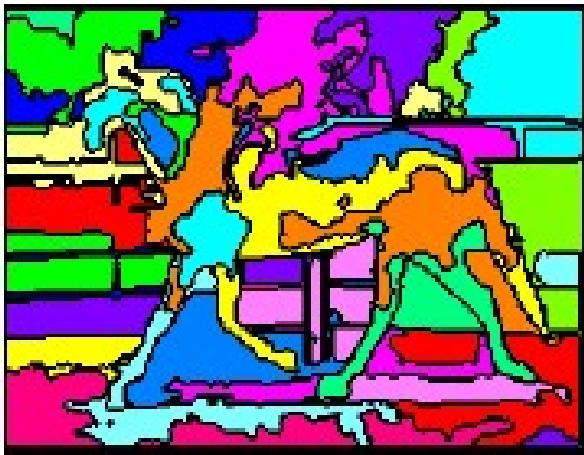
[https://d2l.ai/chapter_computer-vision/
rcnn.html](https://d2l.ai/chapter_computer-vision/rcnn.html)

Efficient region proposals?

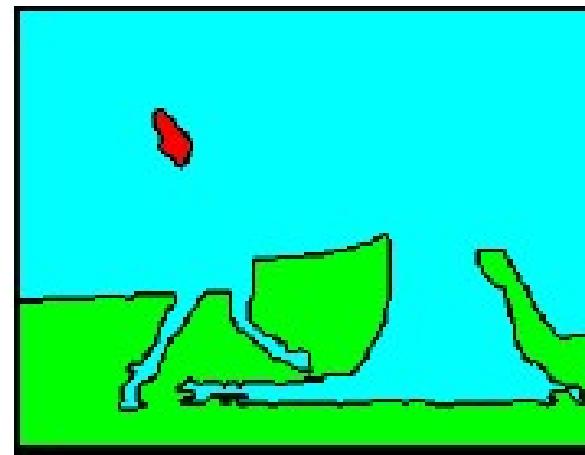
- Brute force on $1000 \times 1000 = 250$ billion rectangles
 - Testing the CNN over each one is too expensive
- Let's use B.C. vision! Before CNNs
 - Hierarchical clustering for segmentation

Use segmentation results for region proposals (instead of exhaustive search).

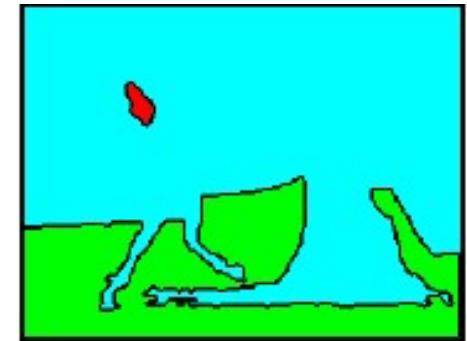
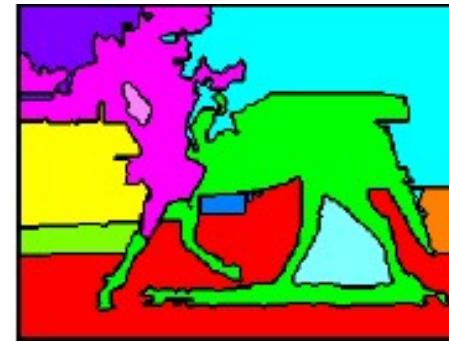
Remember clustering for segmentation!



Oversegmentation



Undersegmentation



Hierarchical Segmentations

Going through the details of selective search:

Original paper:

<http://www.huppelen.nl/publications/selectiveSearchDraft.pdf>

presentations about the paper:

[http://vision.stanford.edu/teaching/cs231b_spring1415/slides/
search_schuyler.pdf](http://vision.stanford.edu/teaching/cs231b_spring1415/slides/search_schuyler.pdf)

<http://www.cs.cornell.edu/courses/cs7670/2014sp/slides/VisionSeminar14.pdf>

Cluster low-level features

- Define similarity on color, texture, size, ‘fill’
- Greedily group regions together by selecting the pair with highest similarity
 - Until the whole image become a single region
- Draw a bounding box around each one
 - Into a hierarchy

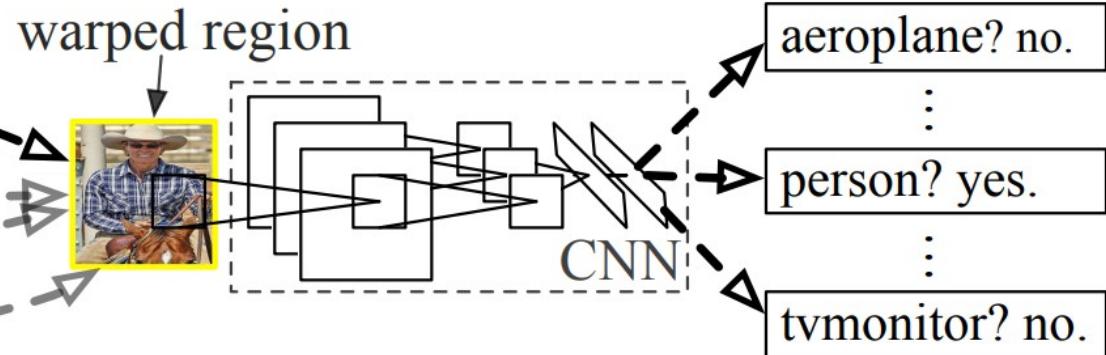
R-CNN: Region-based CNN



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

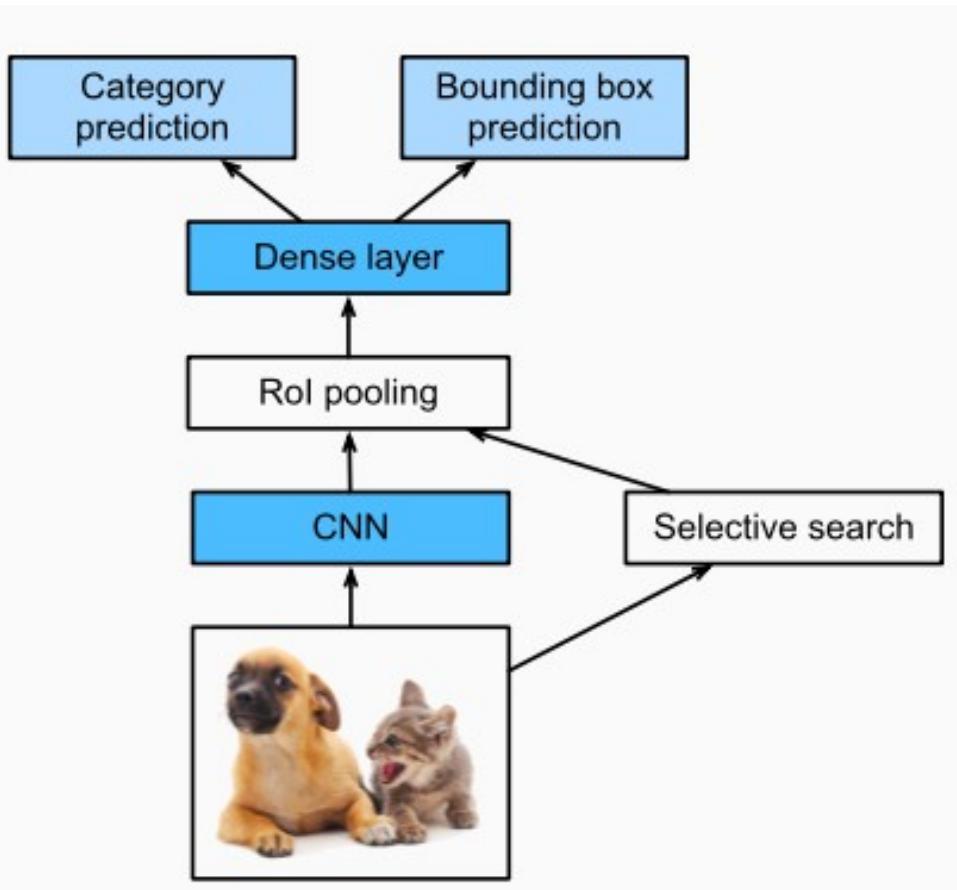
4. Classify regions

Figure: Girshick et al.

10,000 proposals with recall 0.991 is better than ~billions but still takes 17 seconds per image to generate them.

Then we have to perform classification on each one!

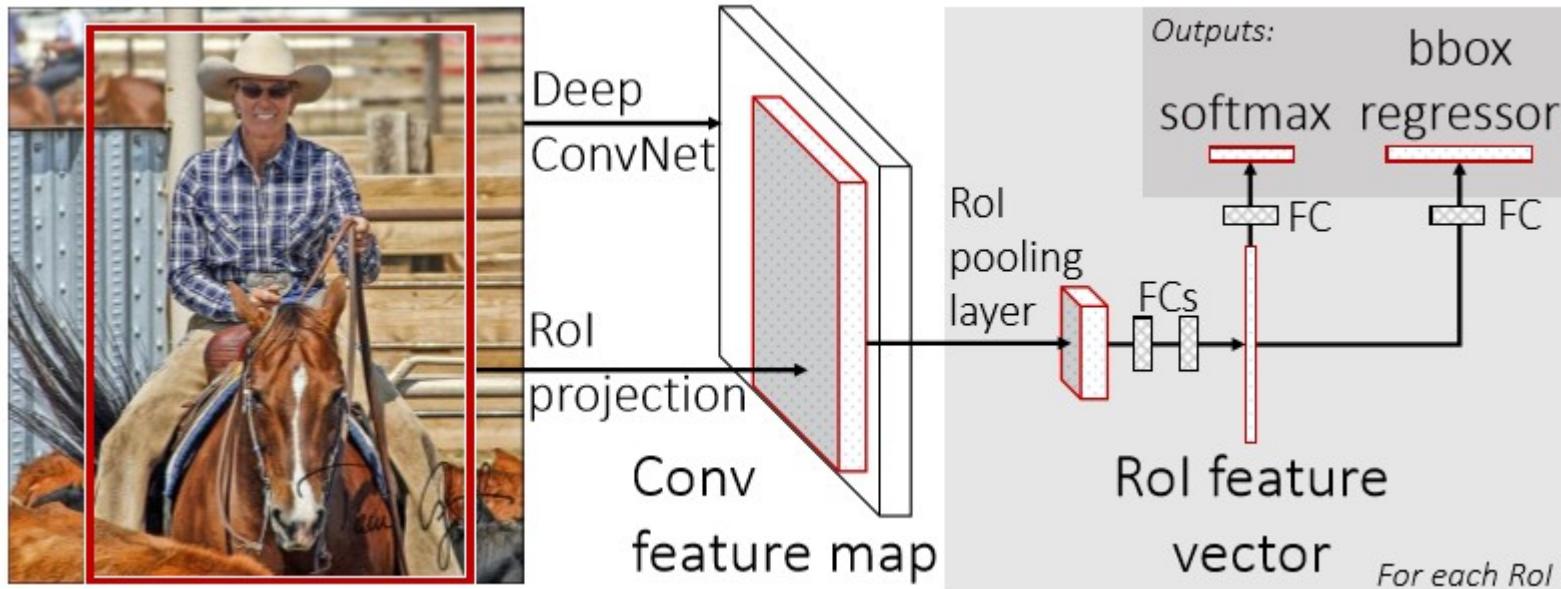
Fast R-CNN



Note: Calling your thing ‘fast’ in computing is not descriptive. In three years, it will no longer be fast, and you will look silly.

Fast R-CNN

Figure: Girshick et al.



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

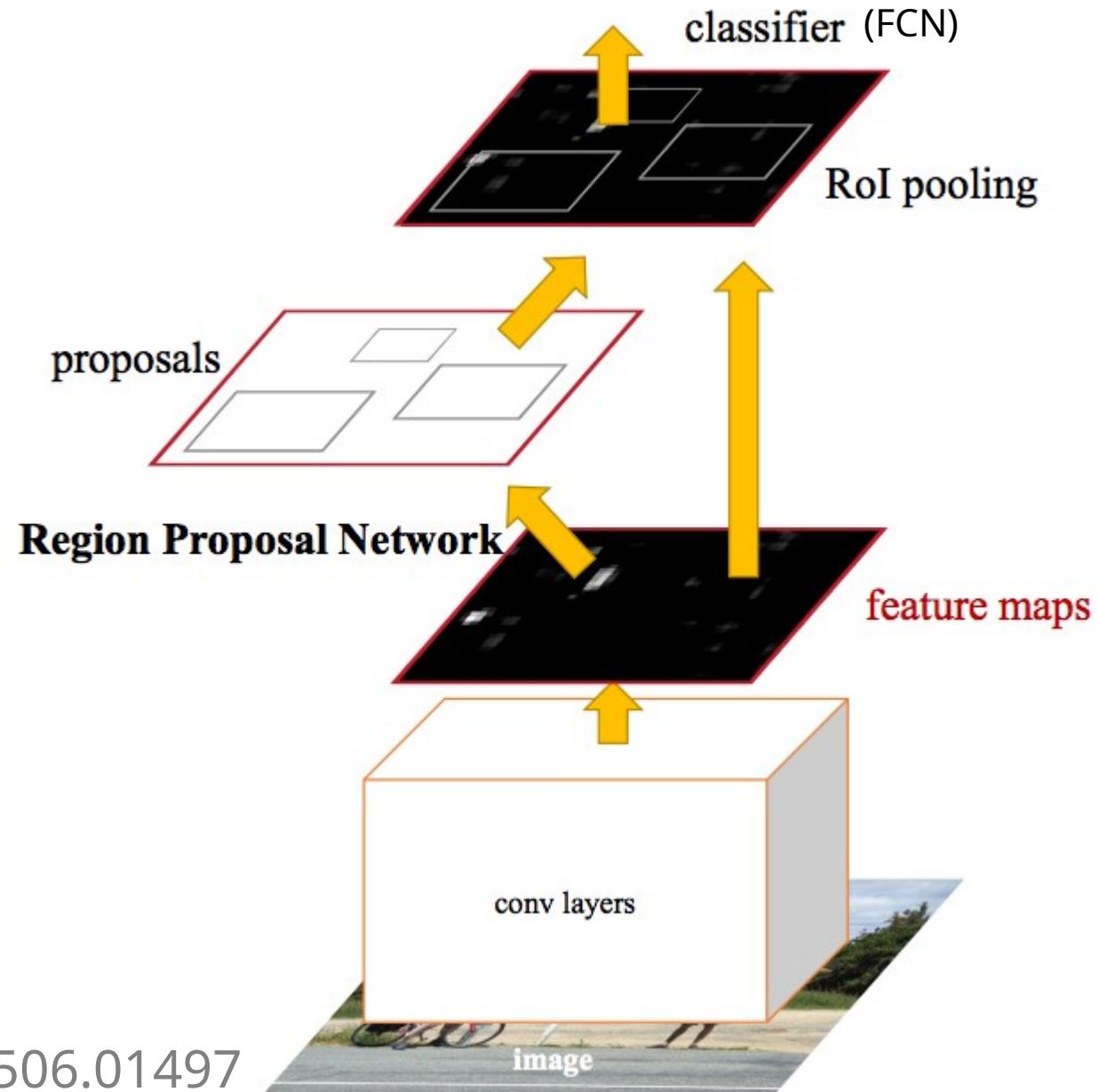
The network has two output vectors per RoI: softmax probabilities and per-class bounding-box regression offsets.

Even more: Faster R-CNN

'Region Proposal Network' uses CNN feature maps.

Then, FCN on top to classify.

End to end object detection.



Ren et al. 2016

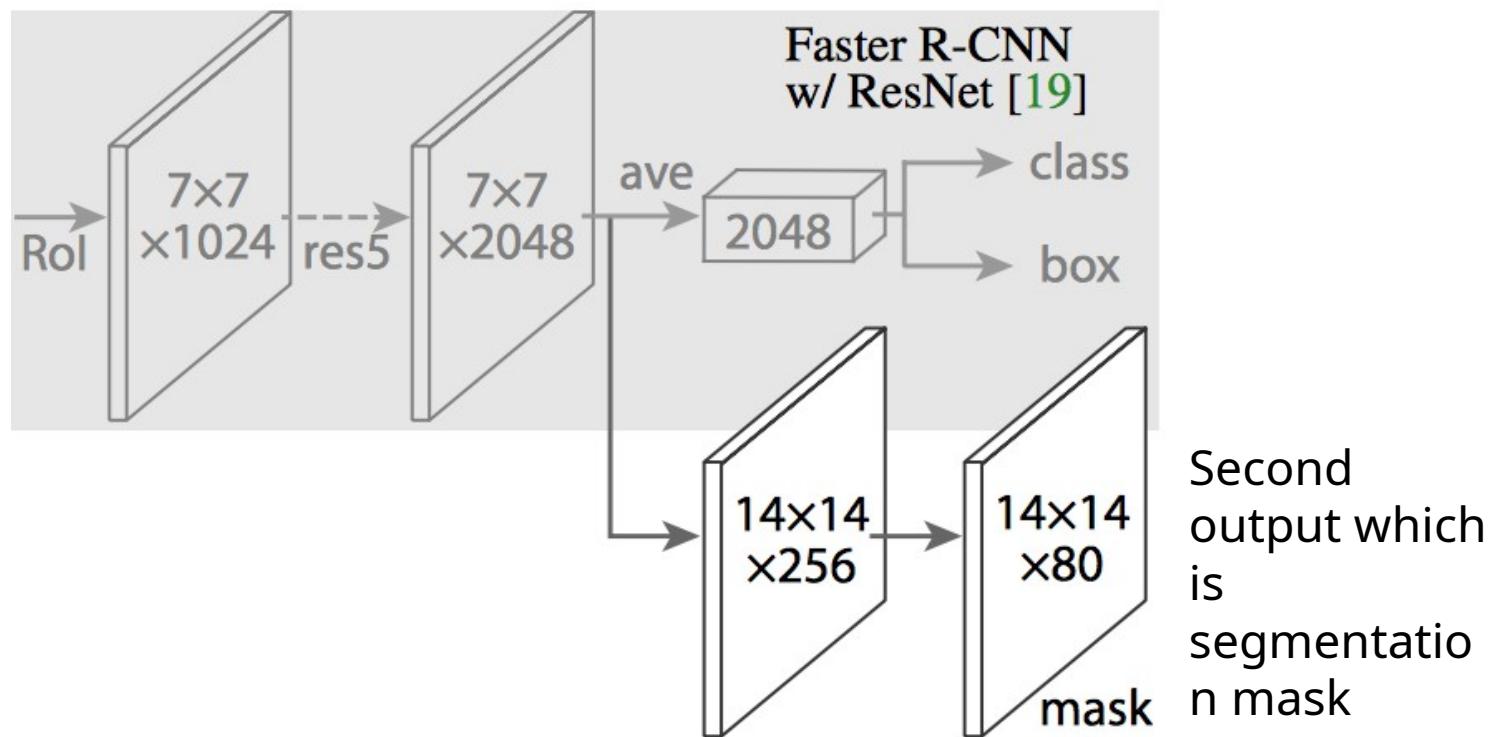
<https://arxiv.org/abs/1506.01497>

Even more! Mask R-CNN

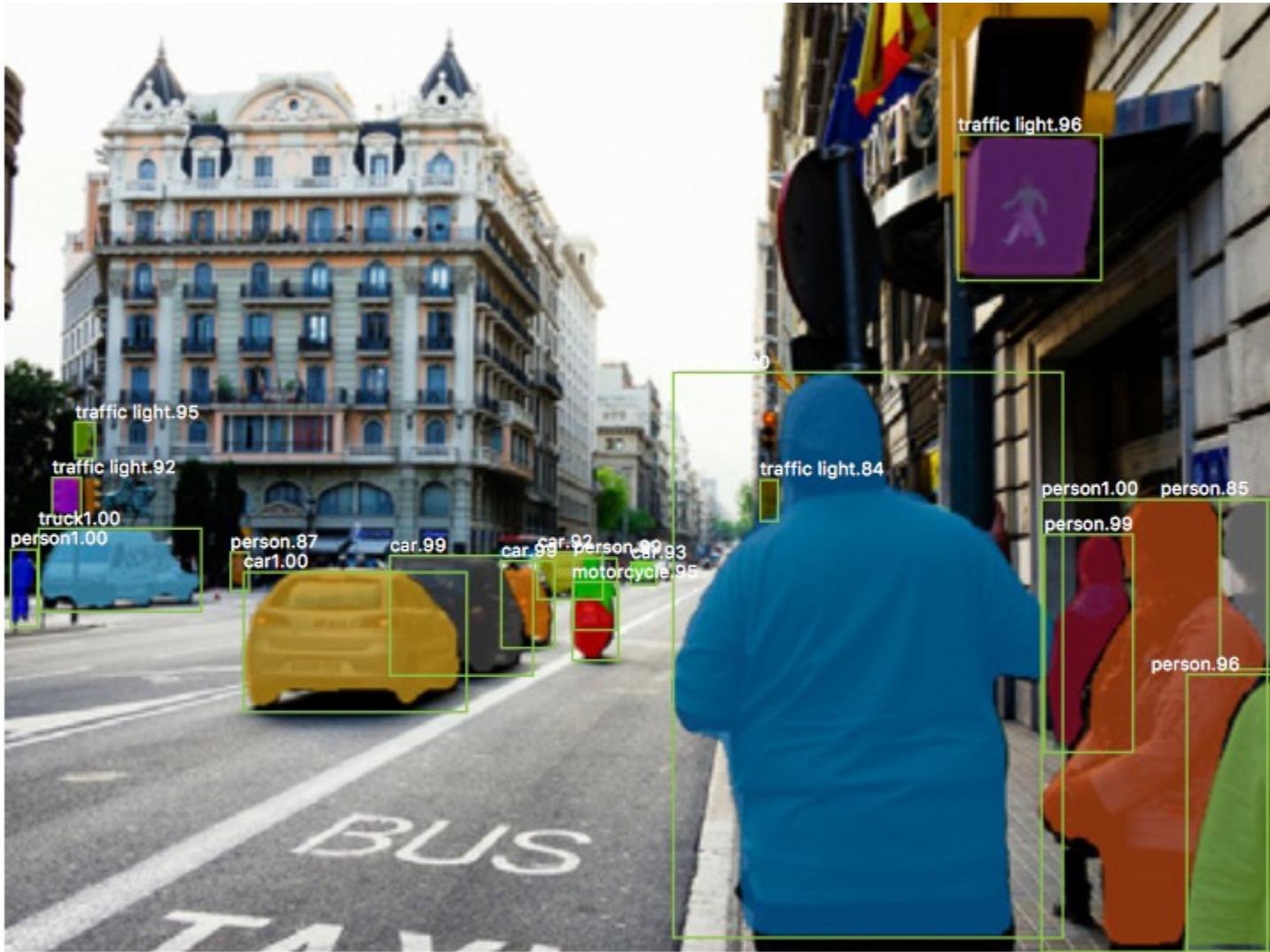
Extending Faster R-CNN for Pixel Level Segmentation

He et al. - <https://arxiv.org/abs/1703.06870>

Add new
training
data:
segmentatio
n masks

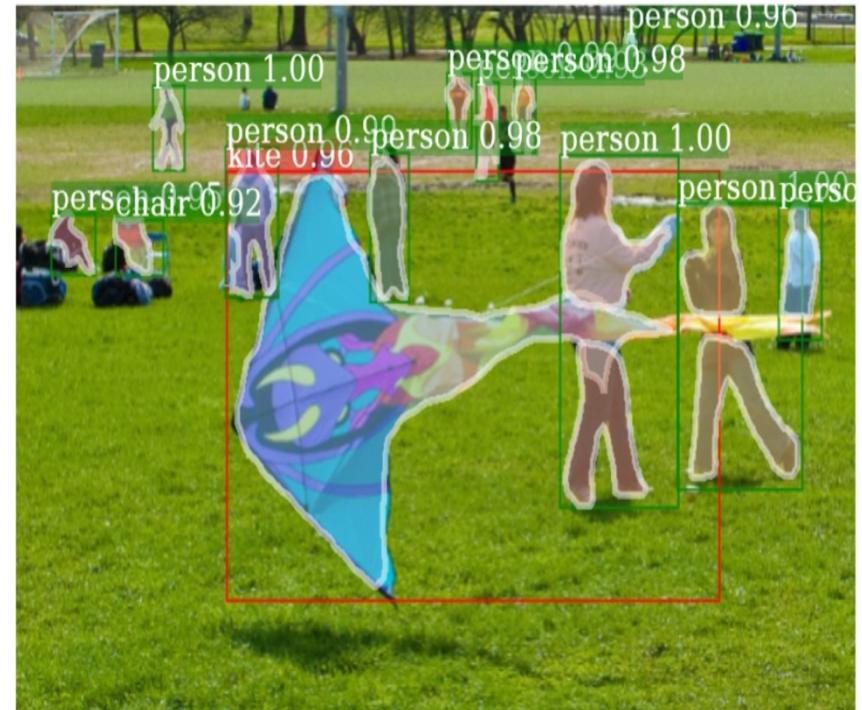
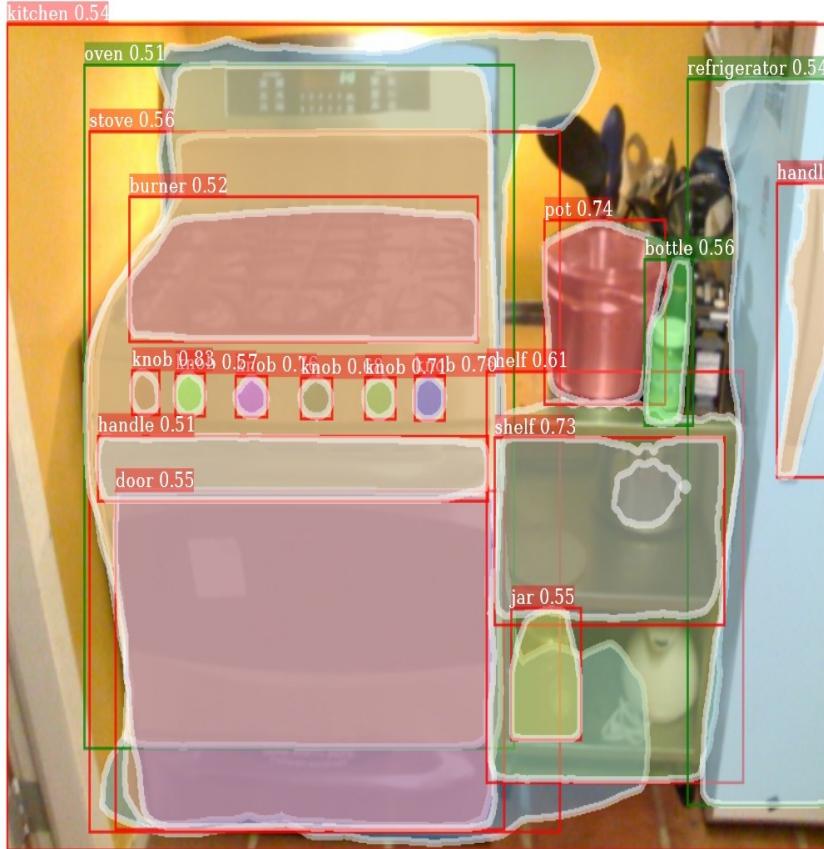


Second
output which
is
segmen
tatio
n mask



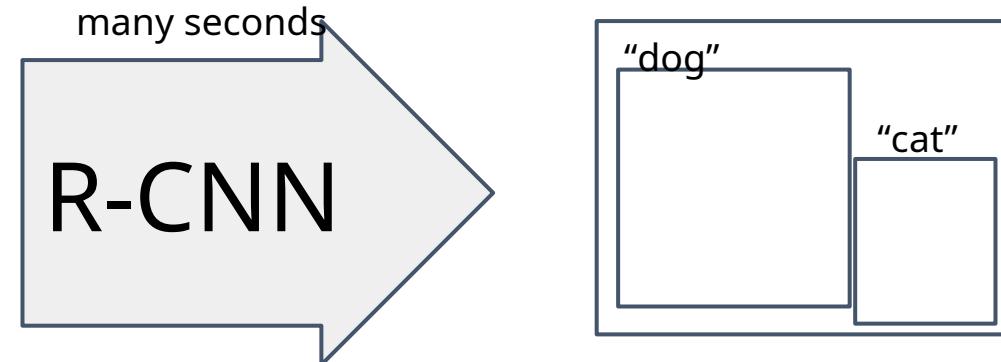


Learning to Segment Everything Mask X R-CNN



FCN: FULLY CONVOLUTIONAL NETWORKS

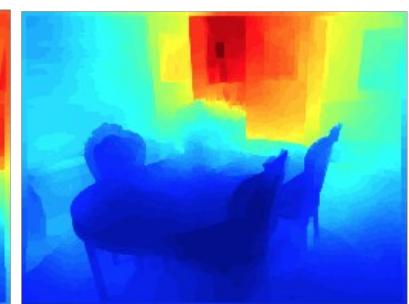
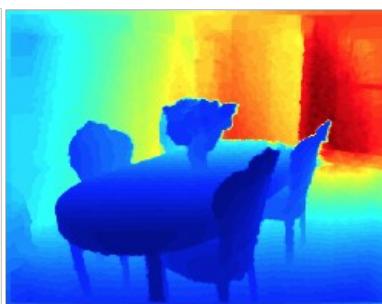
R-CNN does detection



What if we want pixels out?

monocular depth estimation Eigen & Fergus 2015

semantic segmentation

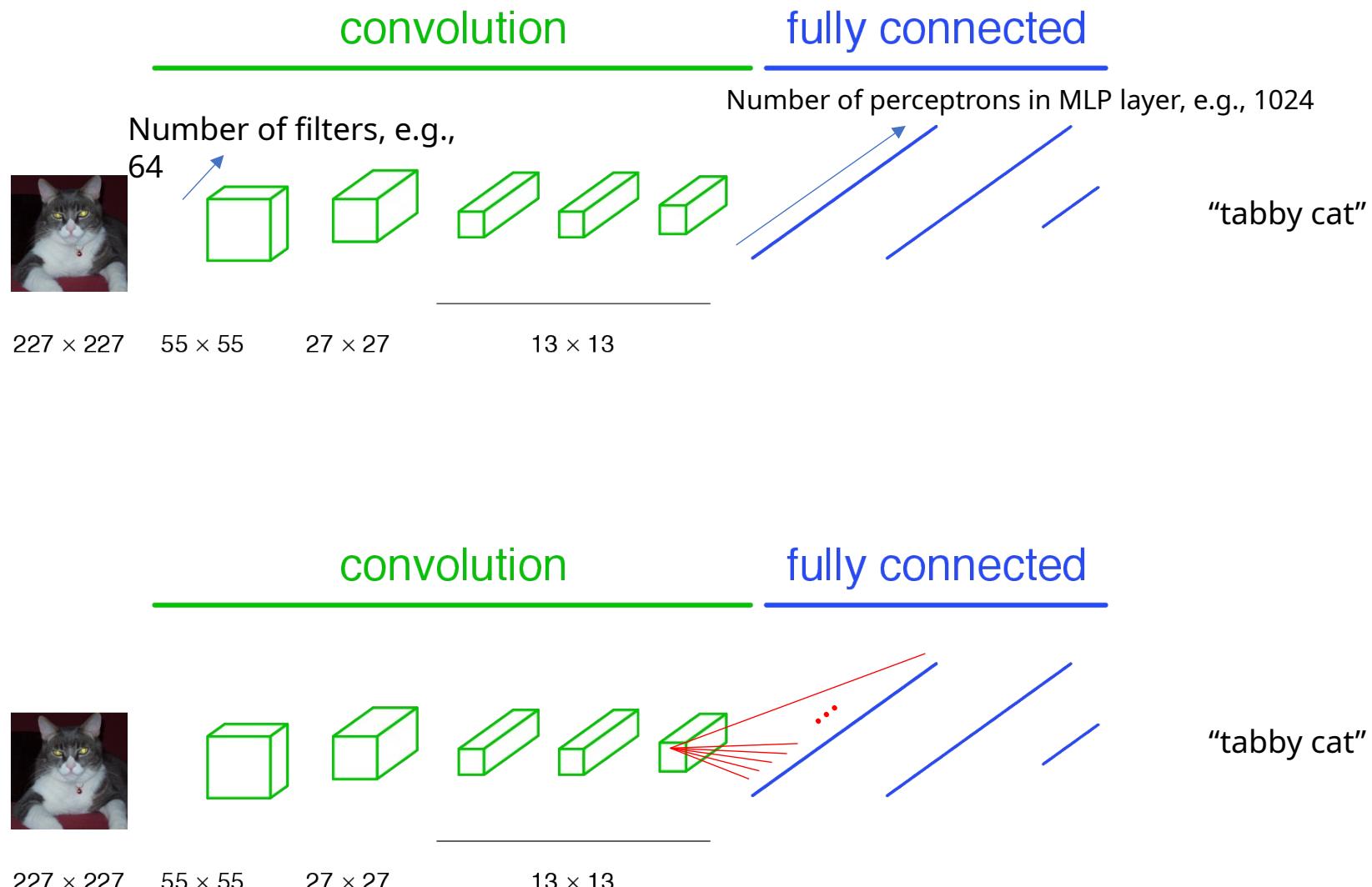


convolutional network



boundary prediction Xie & Tu 2015
[Long et al.]

A classification network...



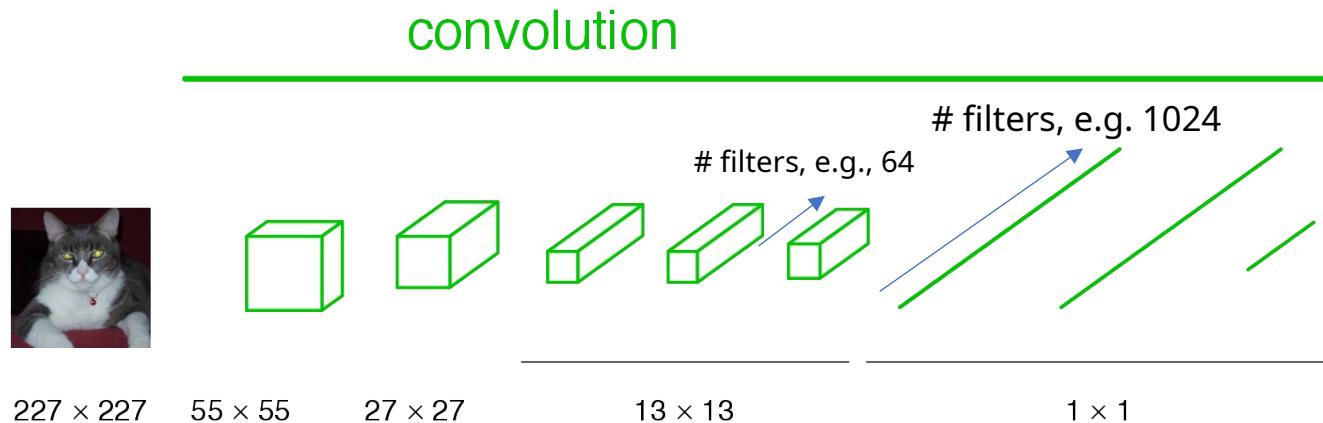
The response of every kernel across all positions are attached densely to the array of perceptrons in the fully-connected layer.

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

AlexNet parameters

Layer	Units	Weights
L_1 (Conv)	290,400	34,848
L_2 (Conv)	186,624	307,200
L_3 (Conv)	64,896	884,736
L_4 (Conv)	64,869	663,552
L_5 (Conv)	43,264	442,368
L_6 (Dense)	4096	37,748,736
L_7 (Dense)	4096	16,777,216
L_8 (Dense)	1000	4,096,000
Conv Subtotal	650,080	2,332,704
Dense Subtotal	9192	58,621,952
Total	659,272	60,954,656

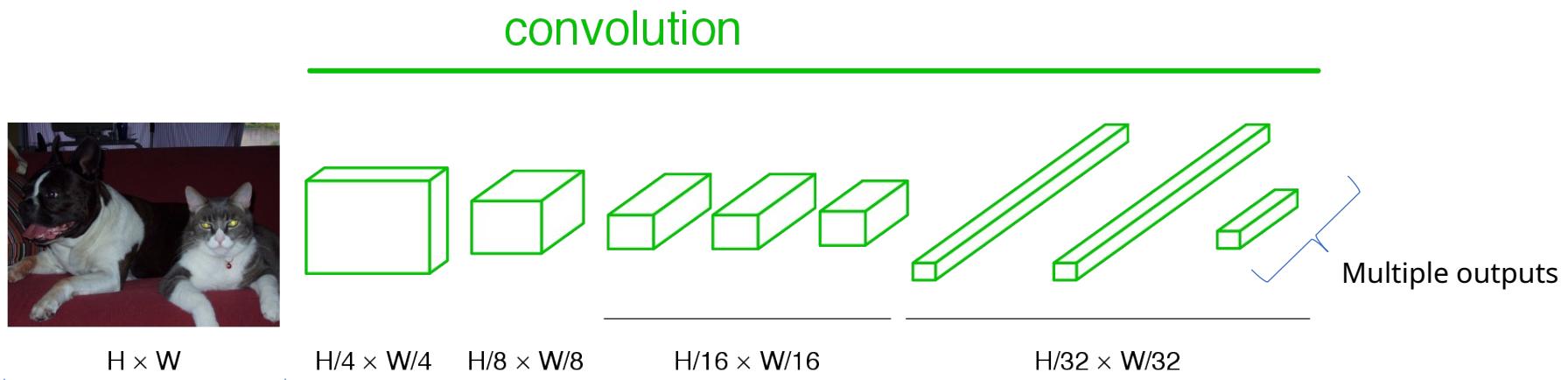
Convolutionalization



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

e.g., $64 \times 1 \times 1$ kernel, with shared weights over 13×13 output, $\times 1024$ filters = 11mil params.

Becoming fully convolutional



When we turn these operations into a convolution, the 13×13 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.

Fully Convolutional Networks



Yann LeCun

6 April 2015 ·

Follow

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

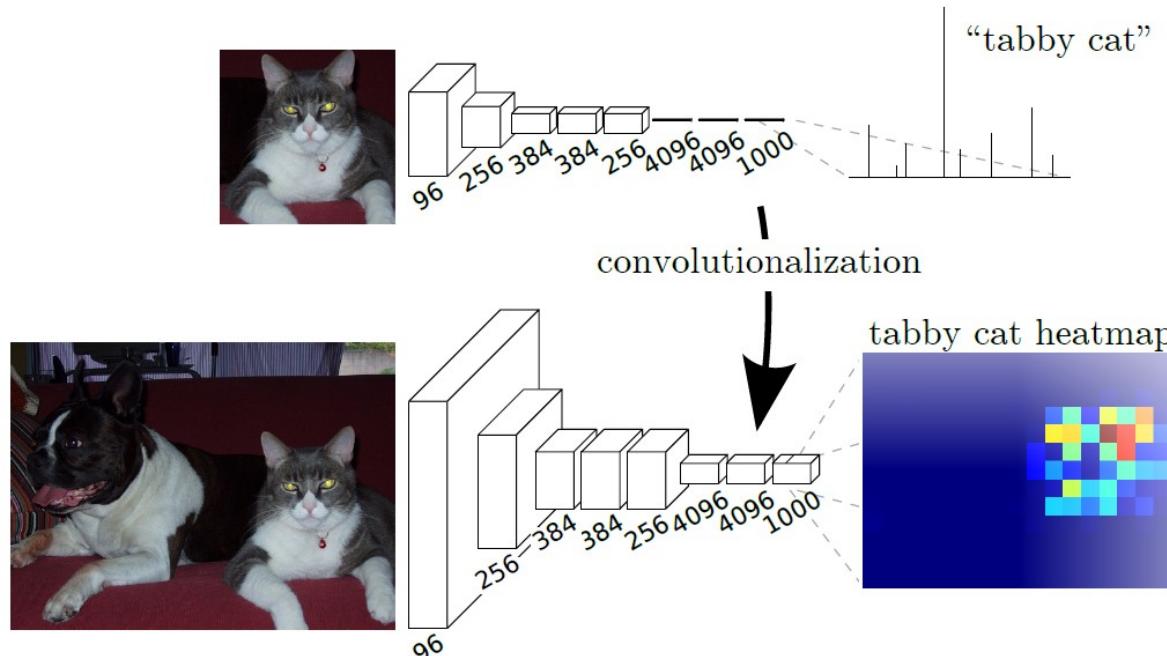
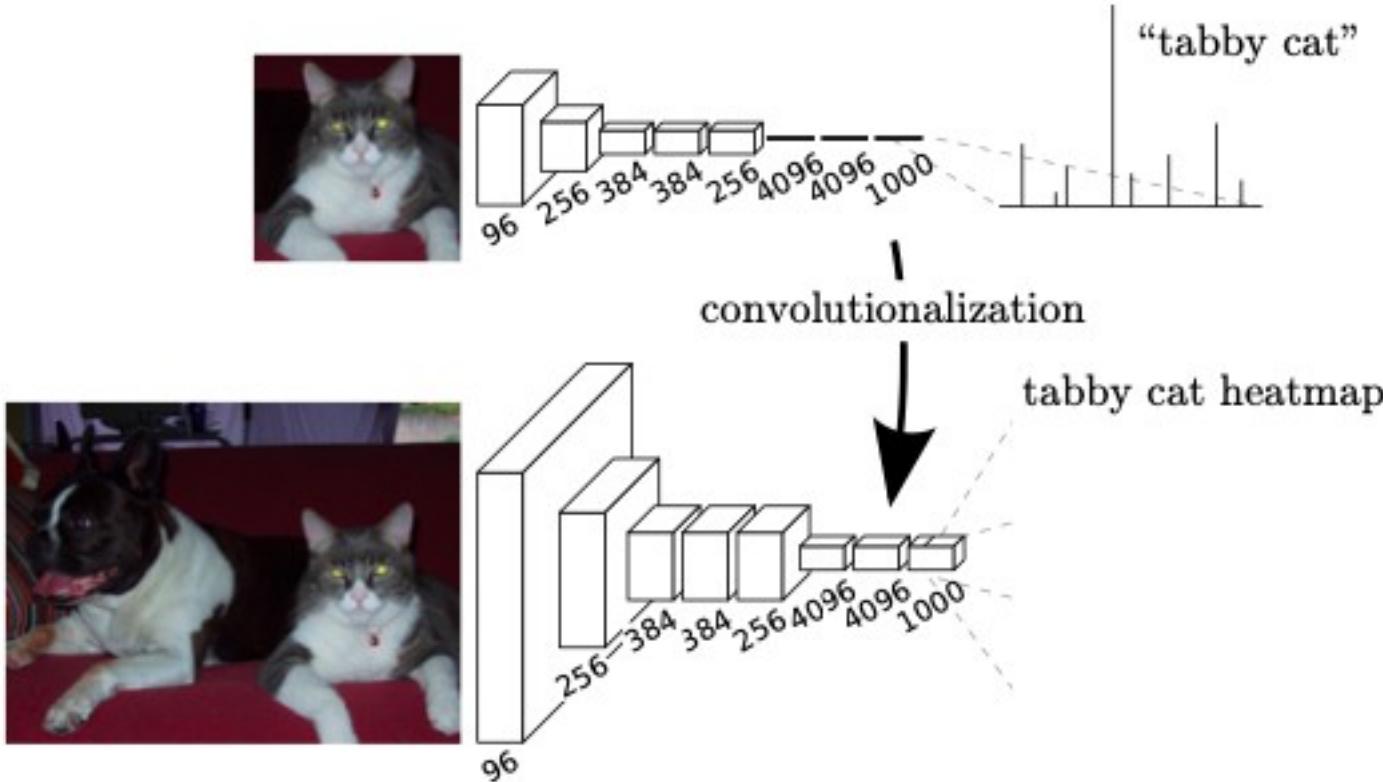


Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

Long, Shelhamer, and Darrell 2014

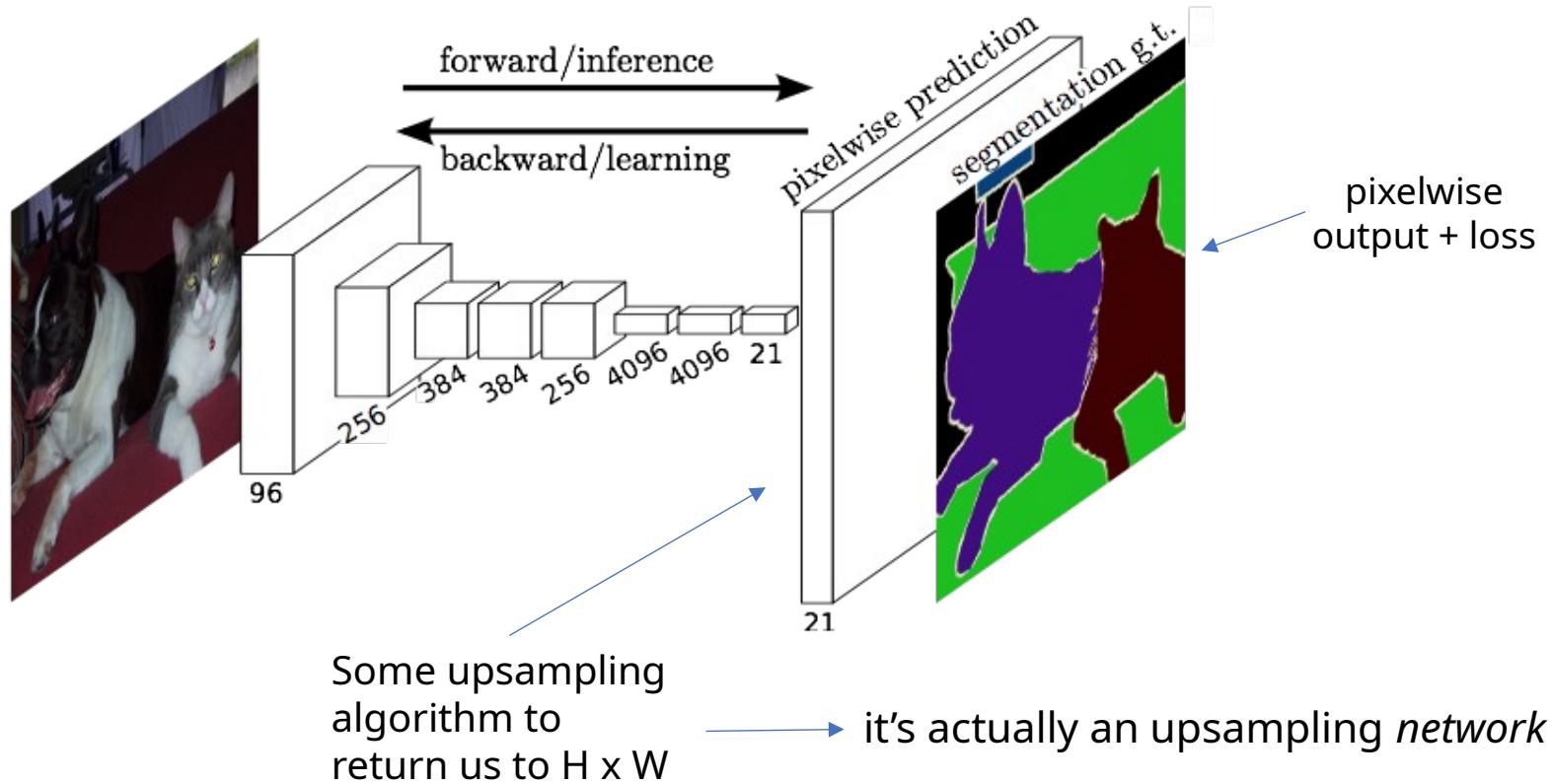
Fully Convolutional Networks



Differences from regular CNN

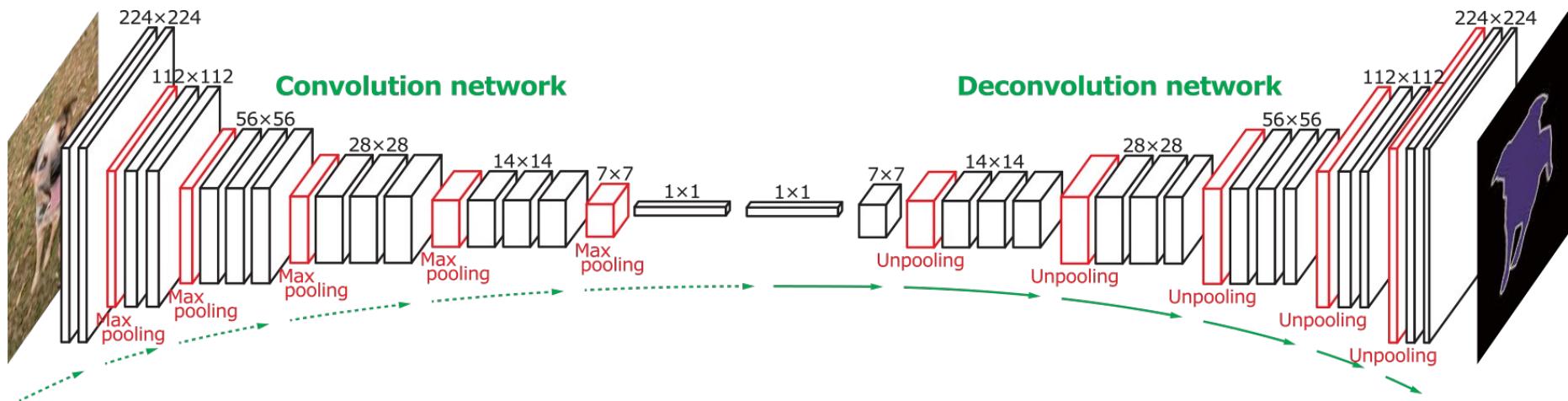
- Can naturally handle varying input image sizes
- Does not forego spatial information
- Output heat maps (the classification parameters is along the kernel's depth)

Fully Convolutional Networks for Semantic Segmentation



Jonathan Long* Evan Shelhamer* Trevor Darrell
UC Berkeley

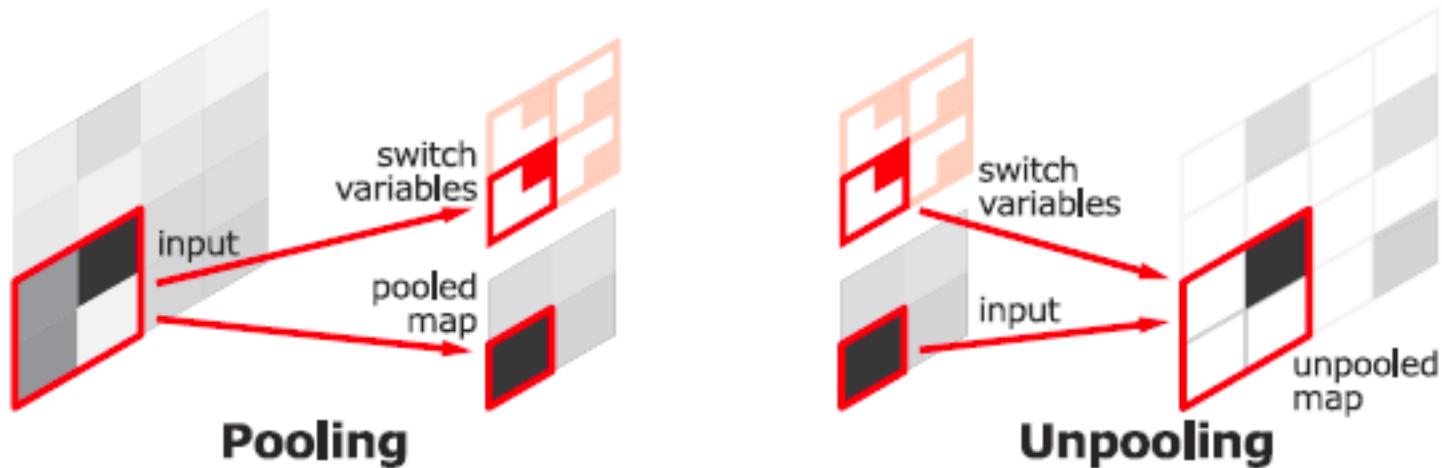
'Deconvolution' networks *learn* to upsample



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

‘Transposed convolution’ is better.

Unpooling layer

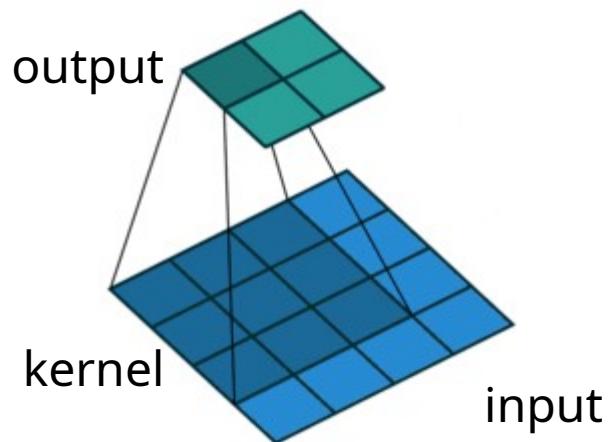


Remember positions when Pooling (Left), Reuse the position information during Unpooling (right)

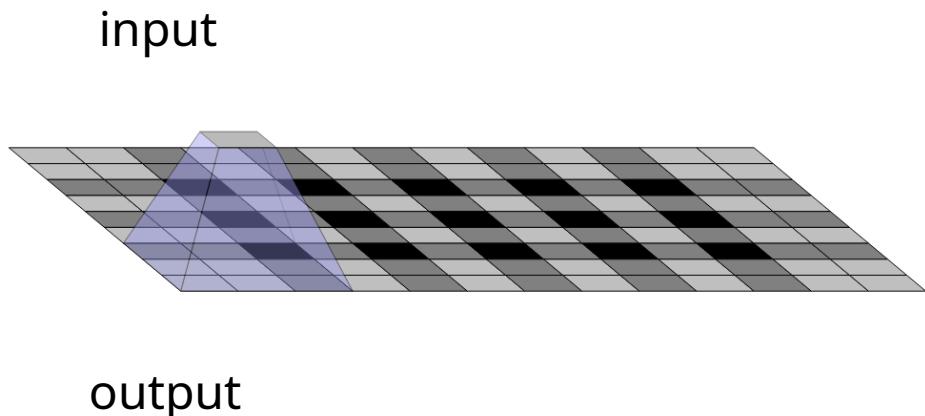
<https://towardsdatascience.com/review-deconvnet-unpooling-layer-semantic-segmentation-55cf8a6e380e>

Upsampling with transposed convolution

Convolution



Transposed convolution smaller image
weighted sum of input x filter: 'stamping'
kernel

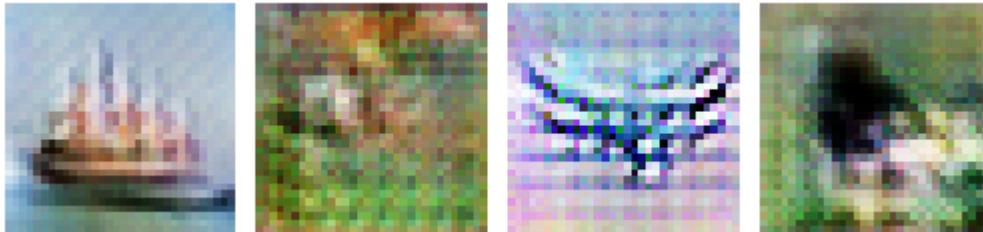


<https://distill.pub/2016/deconv-checkerboard/>

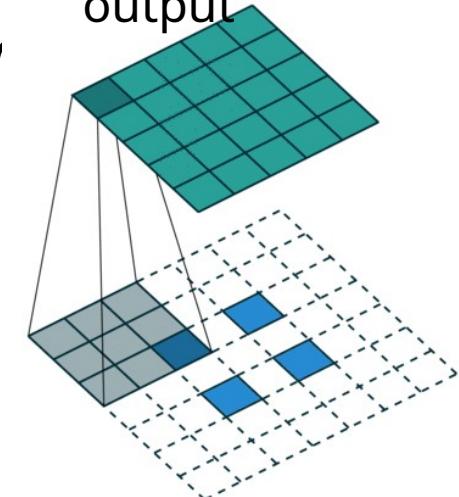
http://deeplearning.net/software/theano/tutorial/conv_arithmetic.html

Is uneven overlap a problem

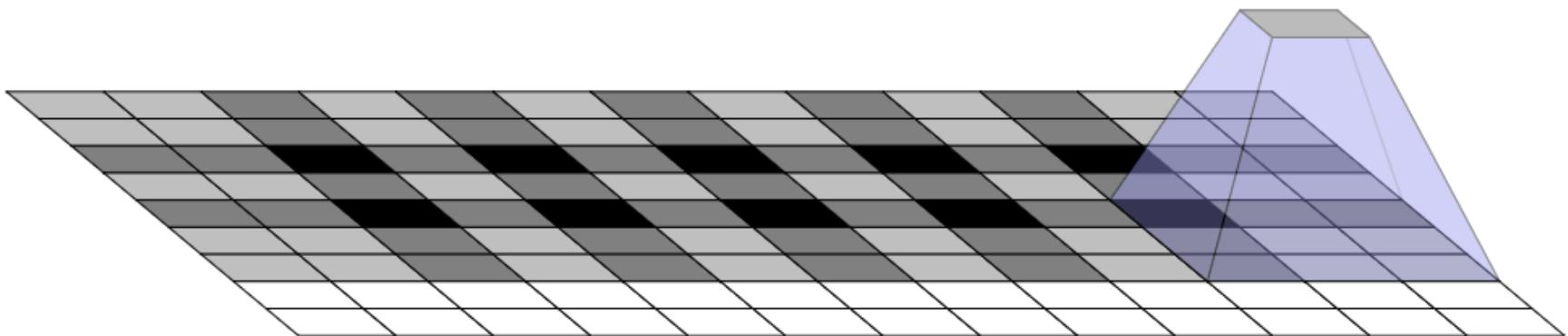
Yes = causes grid artifacts



Uneven
overlap across
output

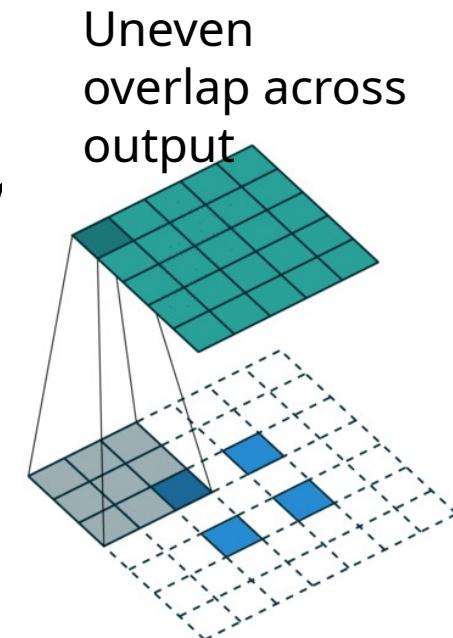
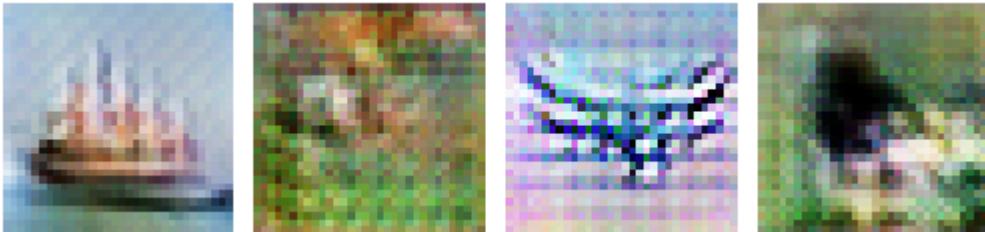


Could fix it by picking stride/kernel numbers which have no overlap...



Is uneven overlap a problem

Yes = causes grid artifacts



Could fix it by picking stride/kernel numbers which have no overlap...

Or...

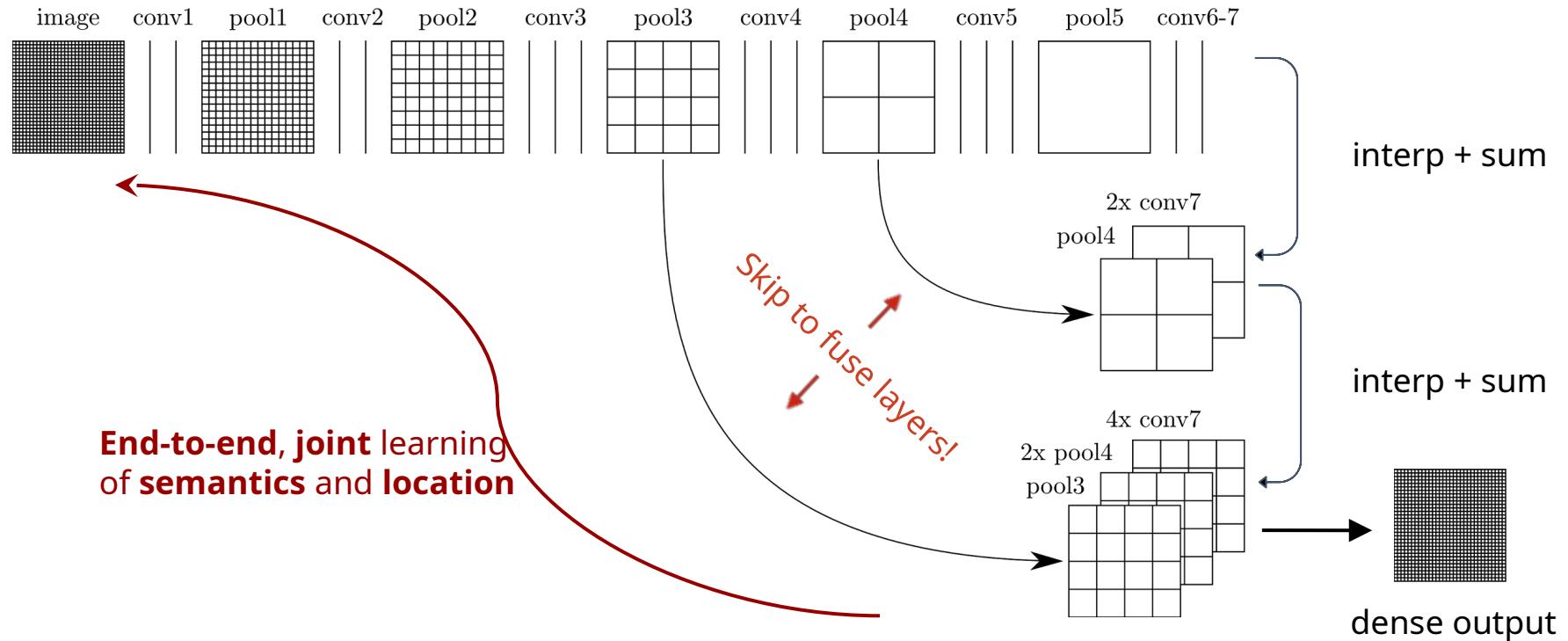
perform explicit bilinear upsampling before transpose convolution; let kernels of transpose convolution learn to fill in only high-frequency detail.

<https://distill.pub/2016/deconv-checkerboard/>

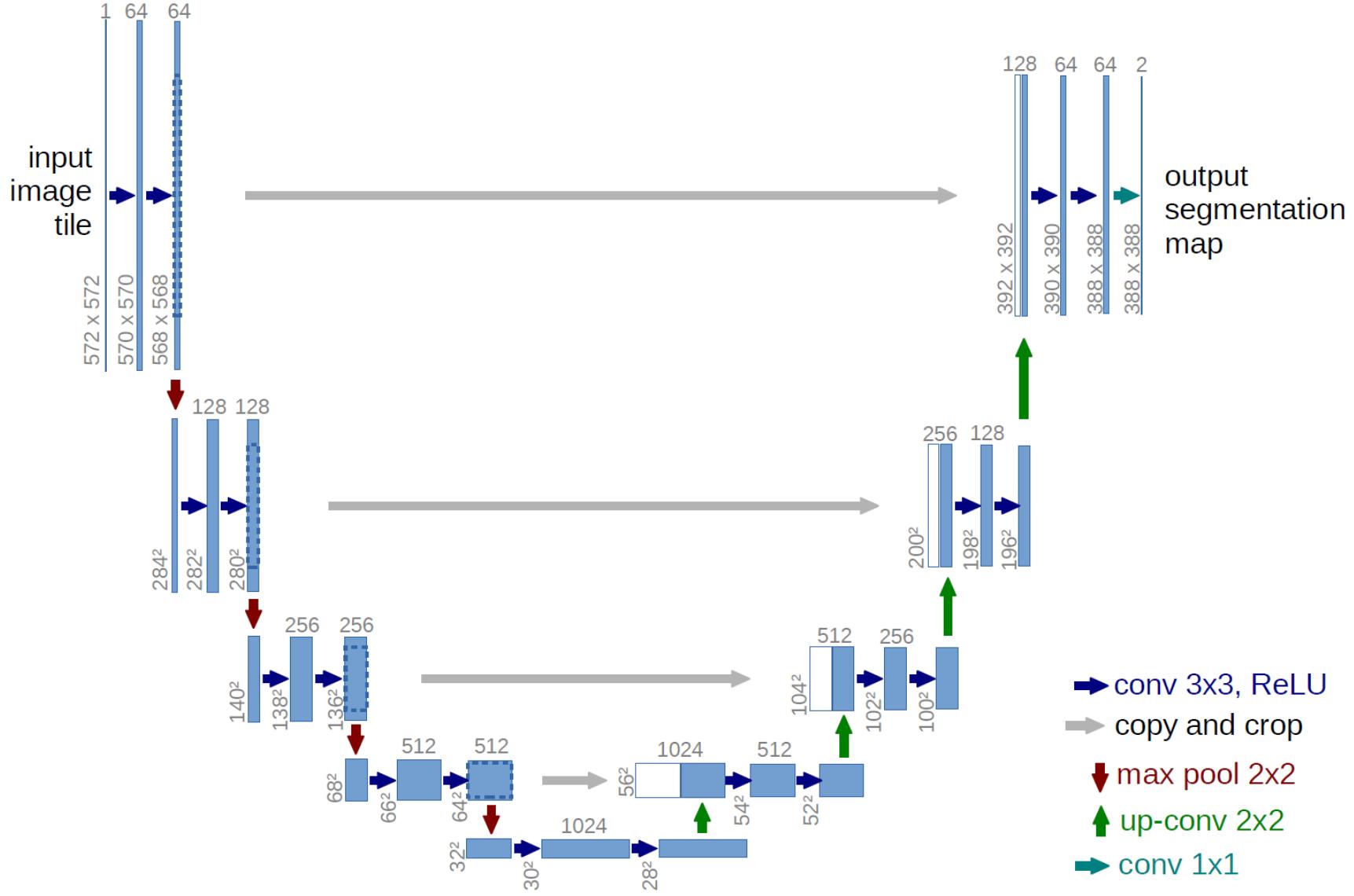
BUT...

How do we 'learn to create' or 'learn to restore'
new high frequency detail?

Learning upsampling kernels with skip layer refinement



UNet [Ronneberger et al., 2015]

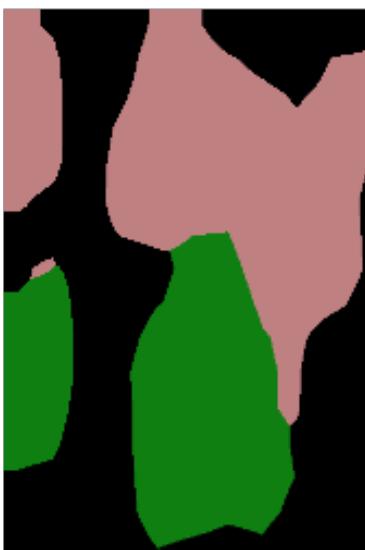


Skip layer refinement

input



stride 32



stride 16



stride 8



ground truth



no skips

1 skip

2 skips

Other usage example of an FCN revisiting im2gps

How much can an image tell about its geographic location?

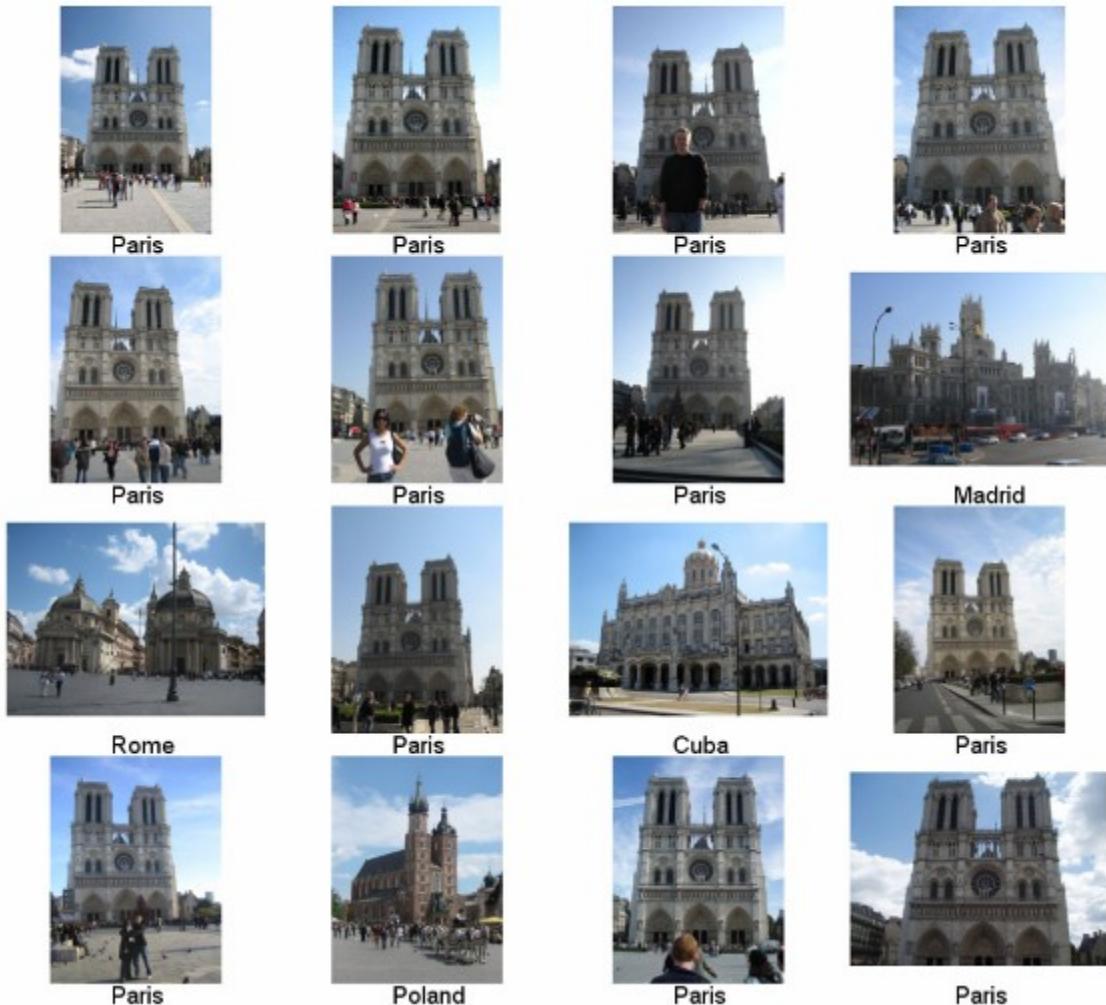


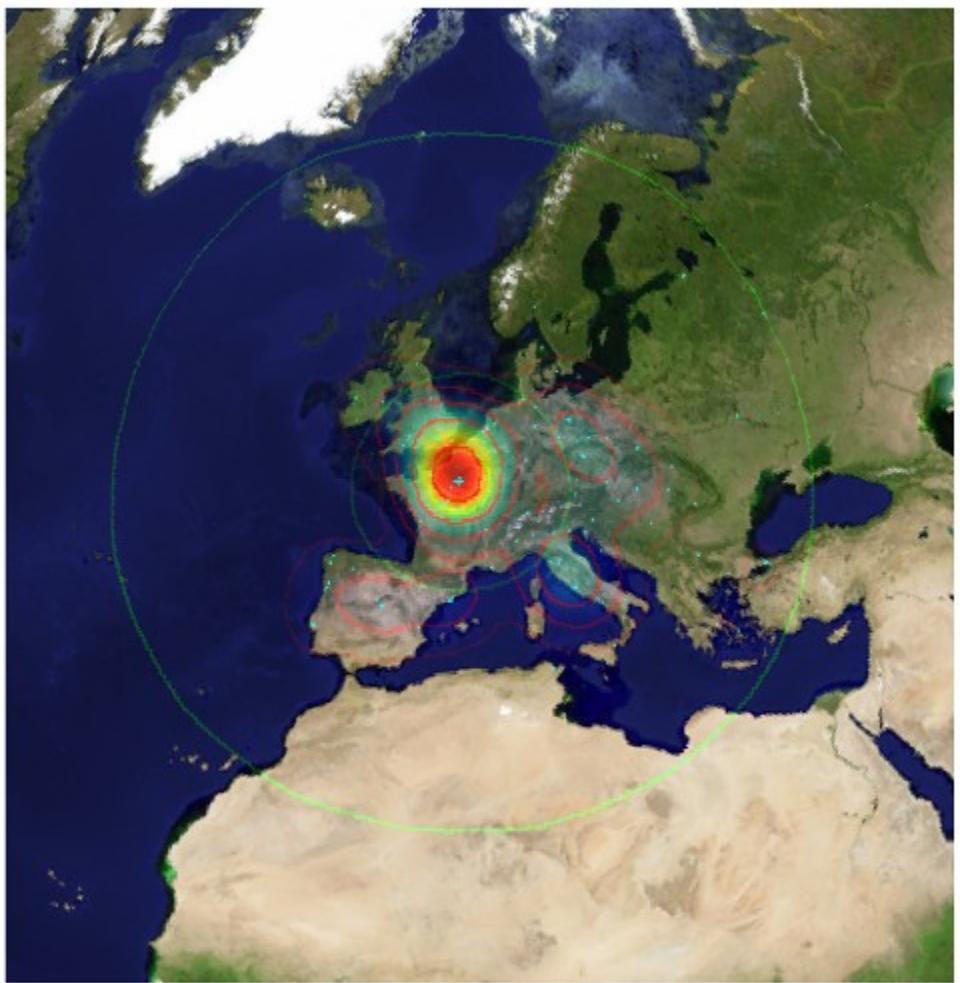
6 million geo-tagged Flickr images

<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





PlaNet - Photo Geolocation with Convolutional Neural Networks

Tobias Weyand, Ilya Kostrikov, James Philbin

ECCV 2016

Discretization of Globe

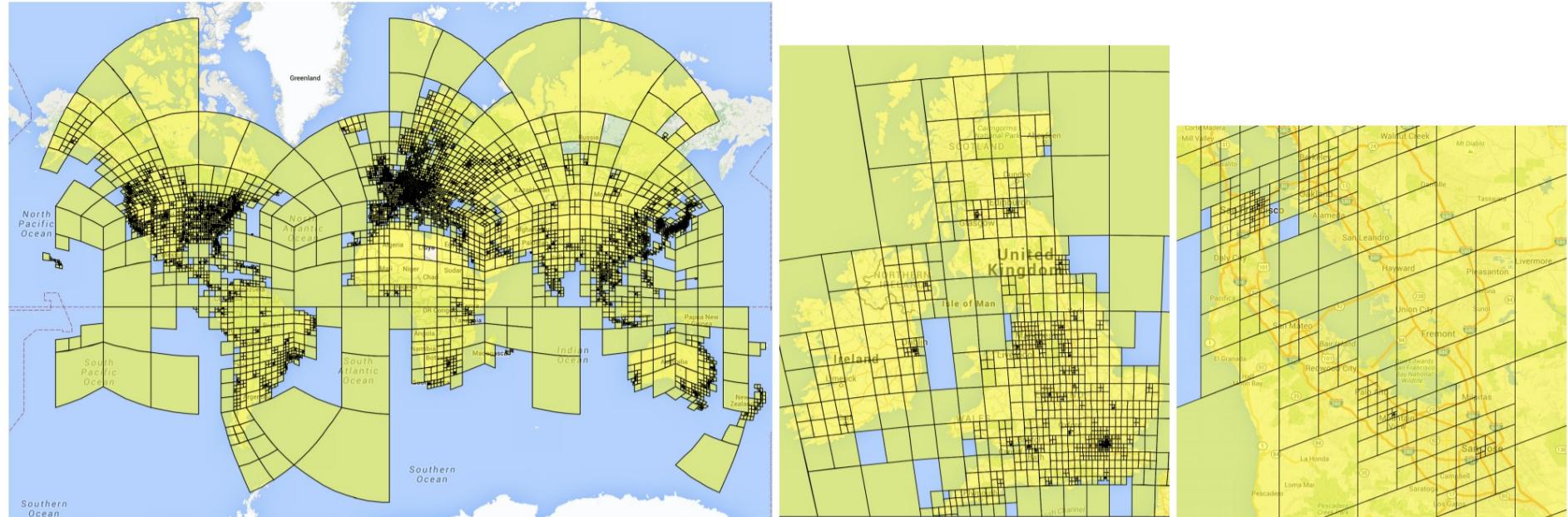


Figure 2. Left: Adaptive partitioning of the world into 26,263 S2 cells. Right: Detail views of Great Britain and Ireland and the San

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)

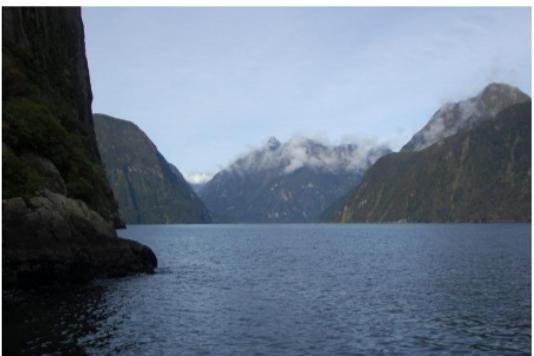


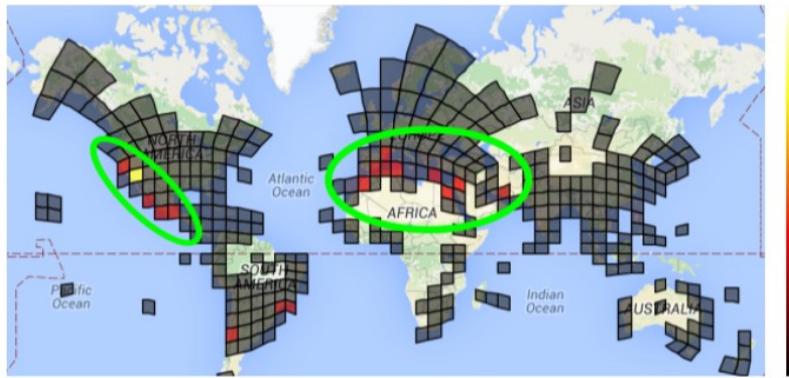
Photo CC-BY-NC by edwin.11



(b)



Photo CC-BY-NC by jonathanfh

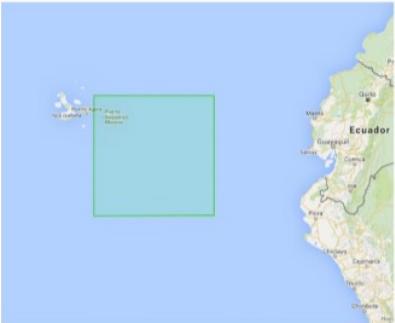




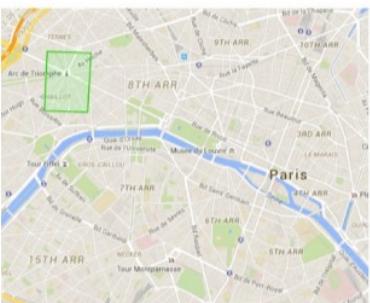
Namibia / Botswana



Kauai, Hawaii



Galapagos Islands



Paris



PlaNet (2016) vs im2gps (2008,2009)

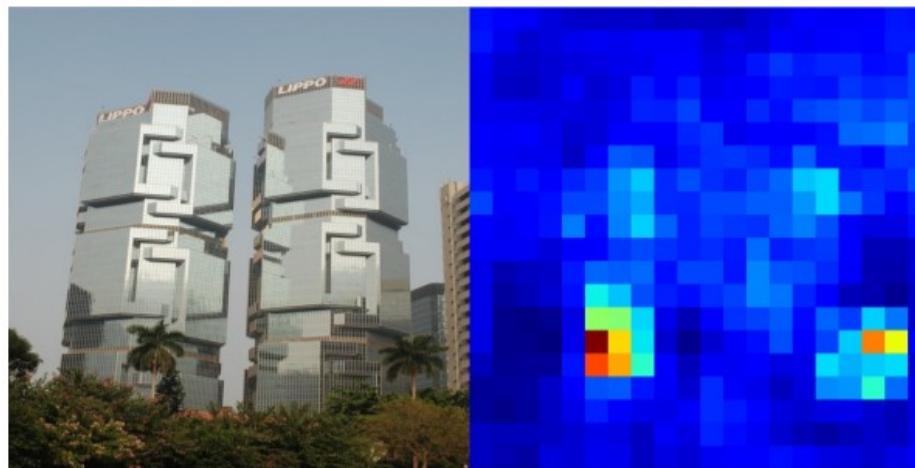
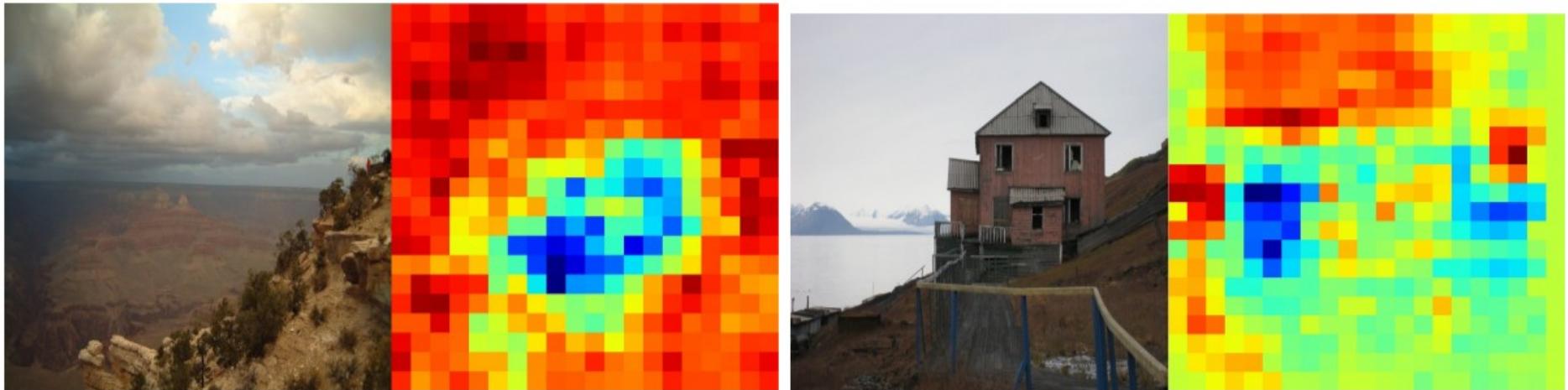
Percentages are the fraction of images from the Im2GPS test set that were localized within the given radius.

Method	Street	City	Region	Country	Continent
	1 km	25 km	200 km	750 km	2500 km
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	8.4%	24.5%	37.6%	53.6%	71.3%

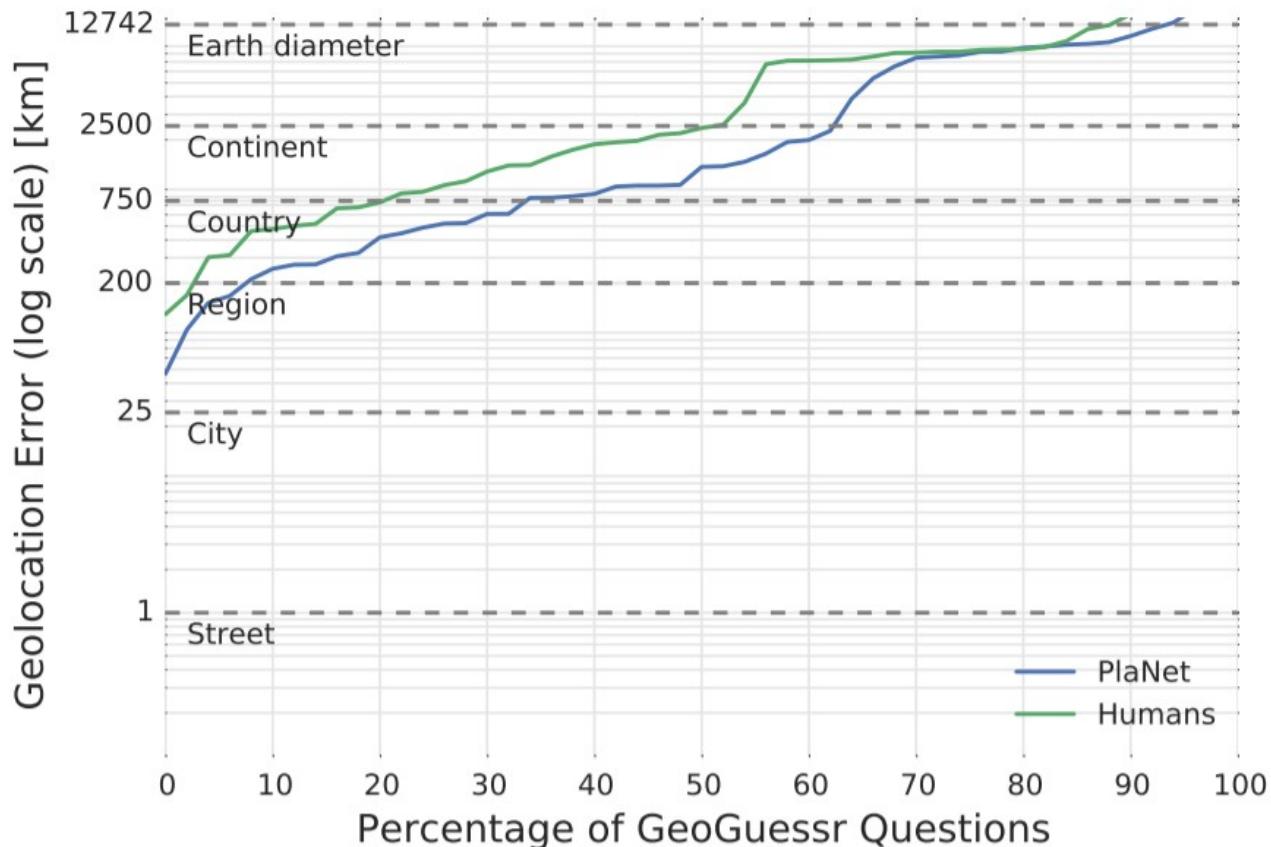
Median localization error (km) by image category.

Method	Manmade	Natural	City	Natural	
	Landmark	Landmark	Scene	Scene	Animal
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



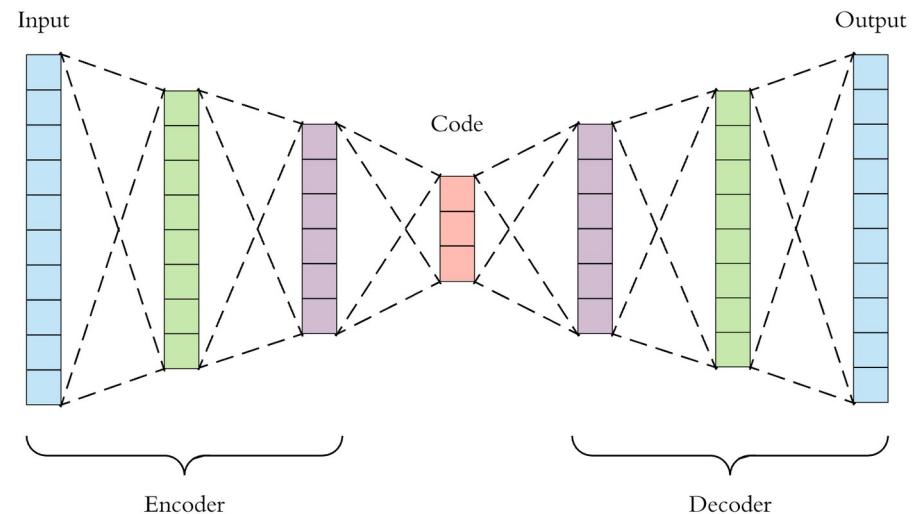
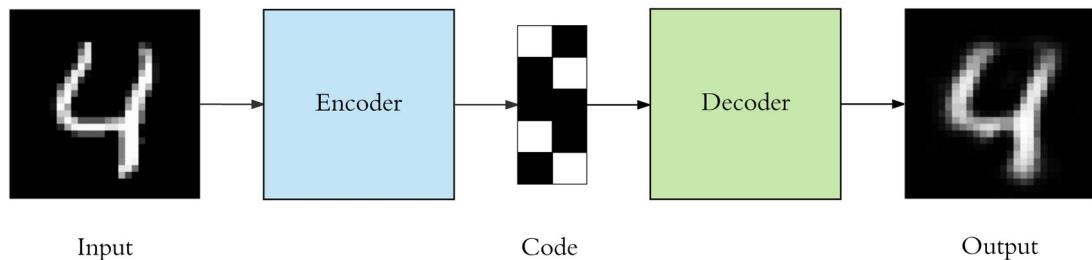
PlaNet vs. Humans



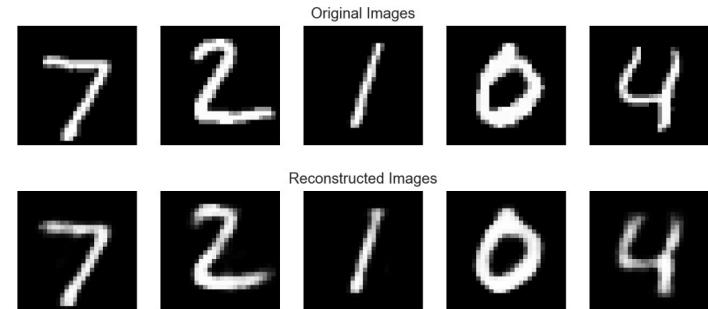
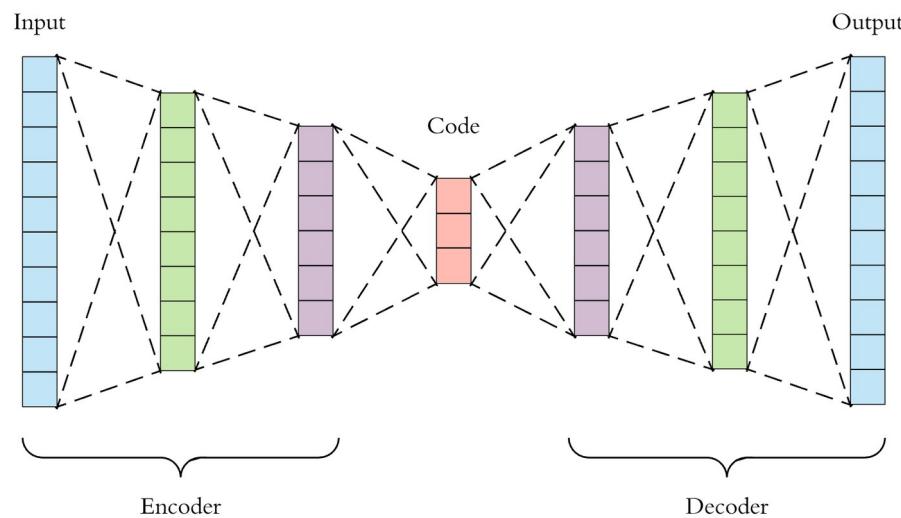
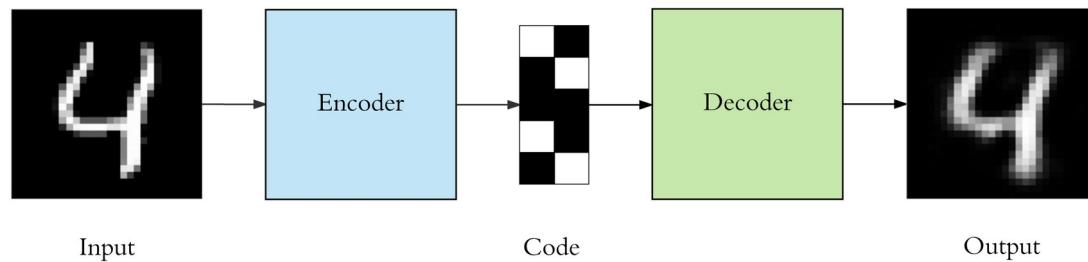
PlaNet vs Humans



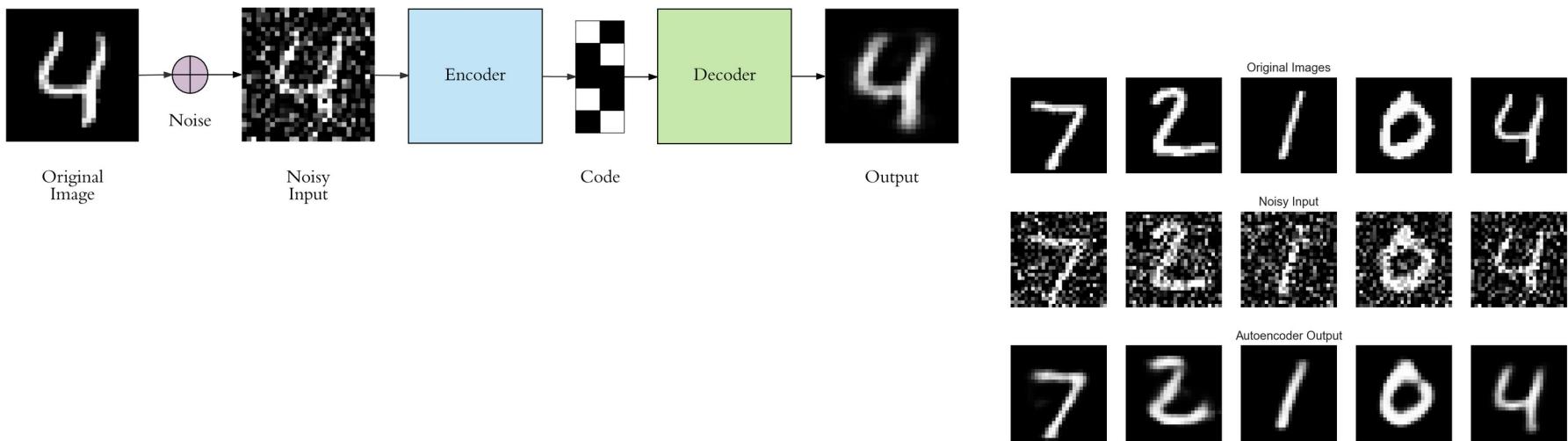
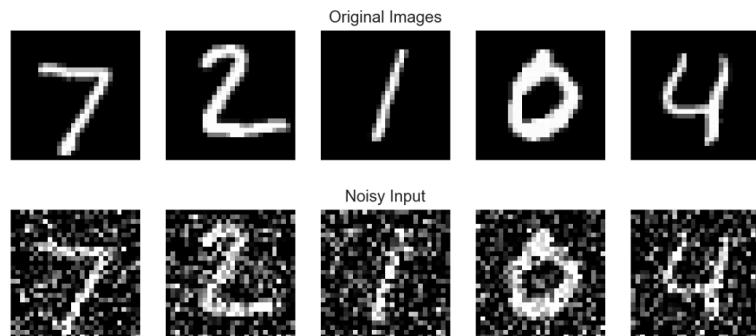
Autoencoders



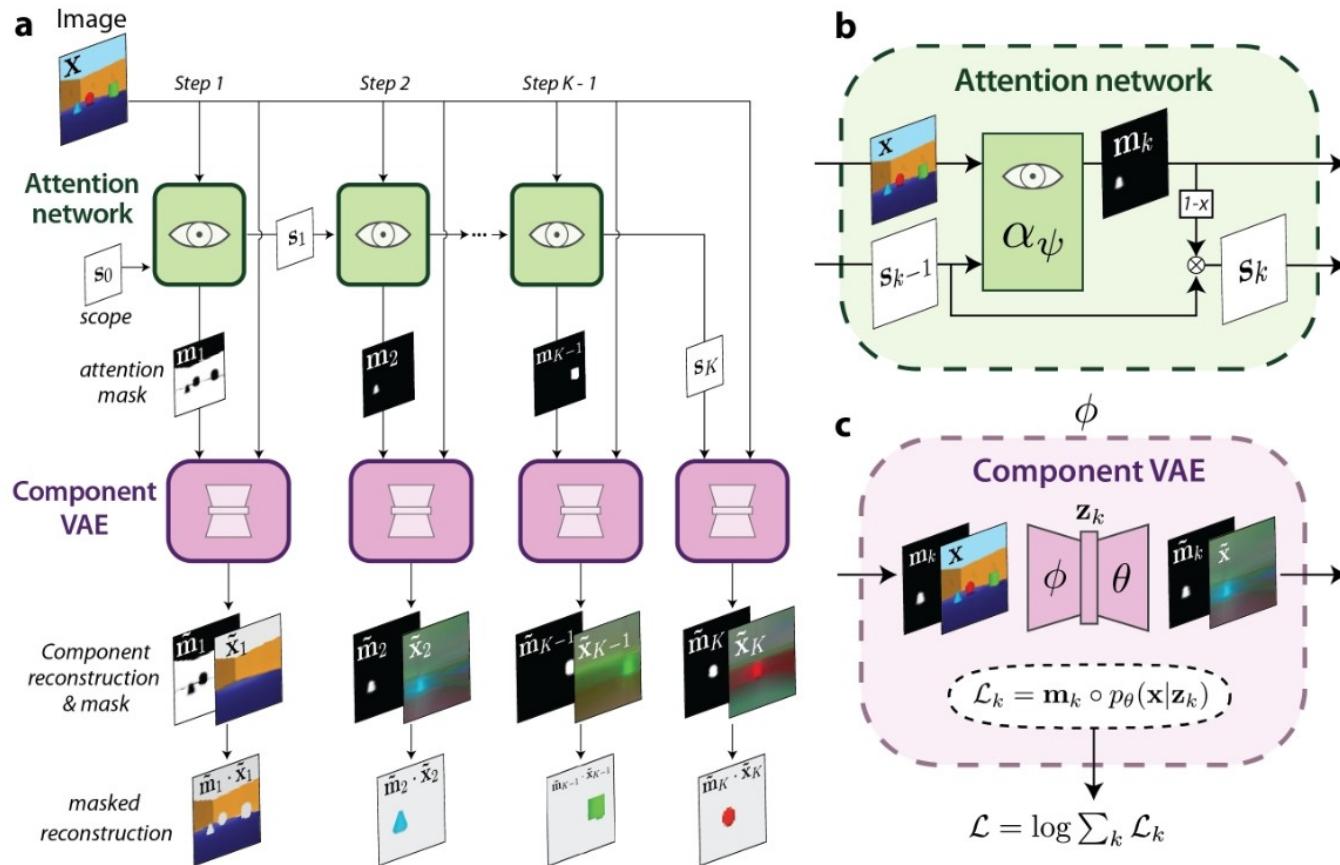
Autoencoders



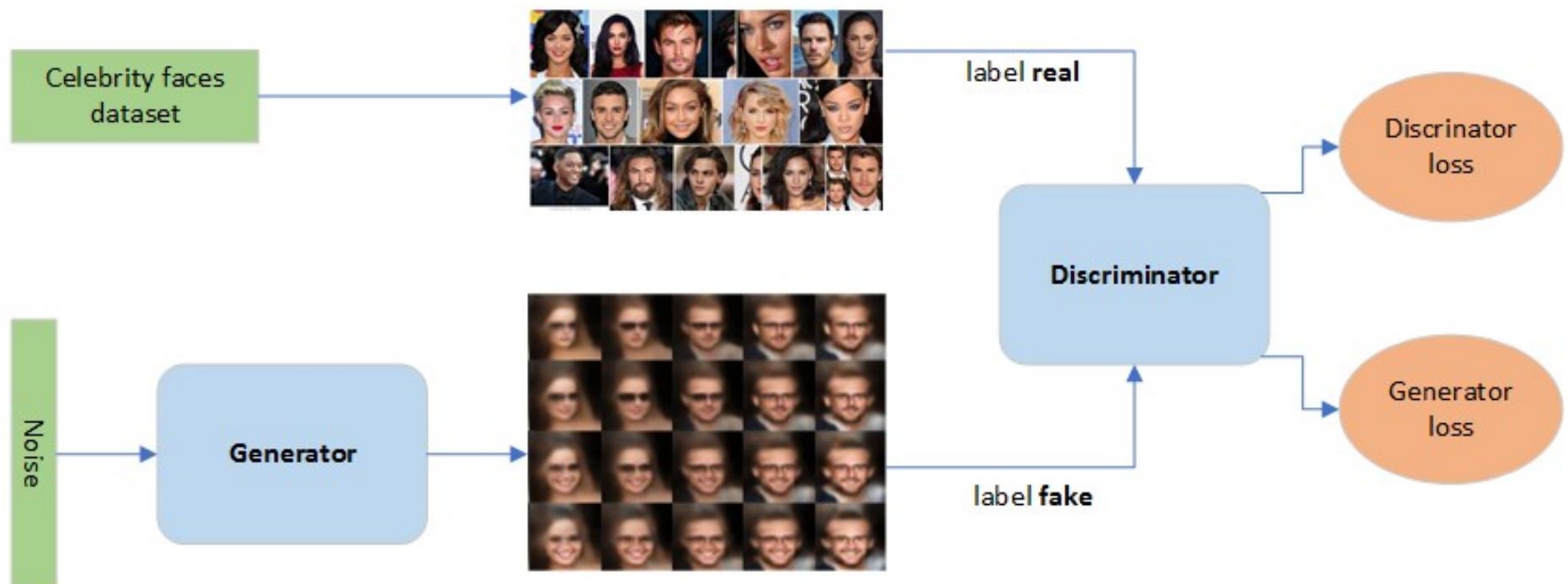
Autoencoders



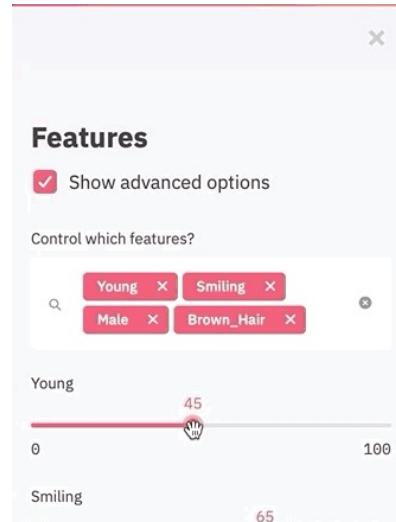
Autoencoders - MOnet



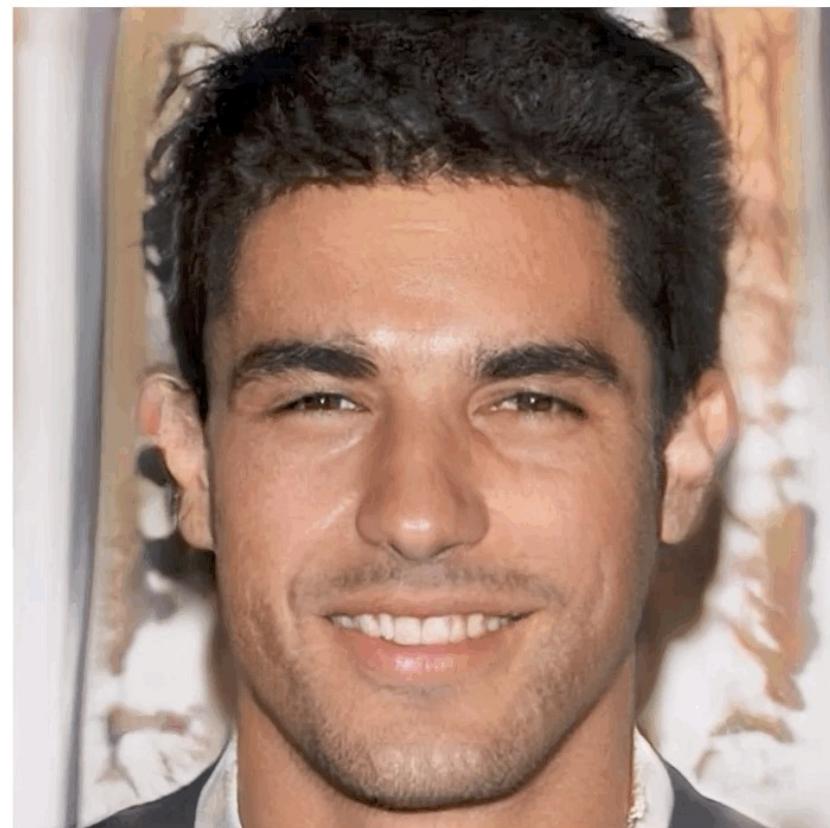
Generative Adversarial Networks



Generative Adversarial Networks

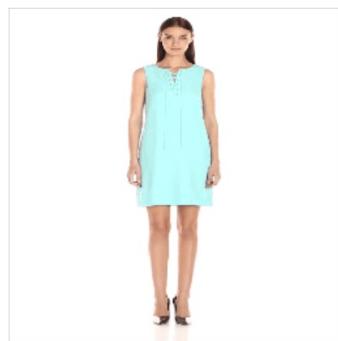


Streamlit Face-GAN Demo



<https://towardsdatascience.com/build-an-app-to-synthesize-photorealistic-faces-using-tensorflow-and-streamlit-dd2545828021>

Generative Models - Example



(Siarohin et al., 2019)

Some refs

Intro: <https://cnl.salk.edu/~schraudo/teach/NNcourse/intro.html>

Deep learning rules of thumb (Jeff Macaluso) summary of Goodfellow et al book:
<https://jeffmacaluso.github.io/post/DeepLearningRulesOfThumb/>

<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-convolutional-neural-networks>

Recipe for training nn: <http://karpathy.github.io/2019/04/25/recipe/>

Useful basic guidelines: <https://yerevann.com/a-guide-to-deep-learning/>

Lots of interesting stuff out there for example on:

generalization of networks: <https://singularityhub.com/2019/10/03/deep-learning-networks-can-t-generalize-but-theyre-learning-from-the-brain/>

top 10 architectures: <https://medium.com/cracking-the-data-science-interview/a-gentle-introduction-to-neural-networks-for-machine-learning-d5f3f8987786>

backprop alternatives: <https://medium.com/@sallyrobotics.blog/backpropagation-and-its-alternatives-c09d306aae4c>

how deep should we go: <https://link.springer.com/article/10.1007/s12559-019-09667-7>