

# Deep Learning

## Lecture 8: Transfer learning & perceptual losses

Alexander Ecker  
Institut für Informatik, Uni Göttingen



<https://alexanderecker.wordpress.com>

– Credit: some of the slides based on Fei Fei Li, Justin Johnson and Serena Yeung's slides –

# Today

This is us



# Standing on the shoulders of giants

## Today's topics

1. Transfer learning
2. Neural style transfer

# Transfer learning

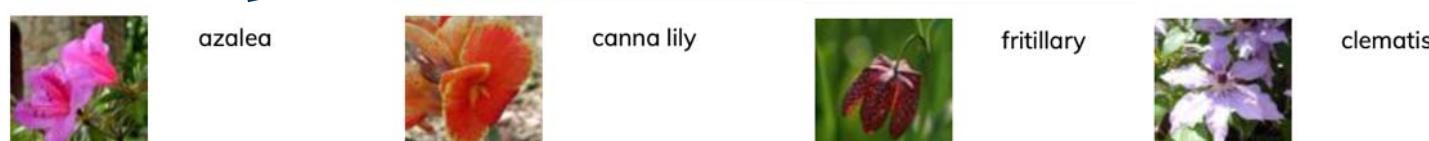
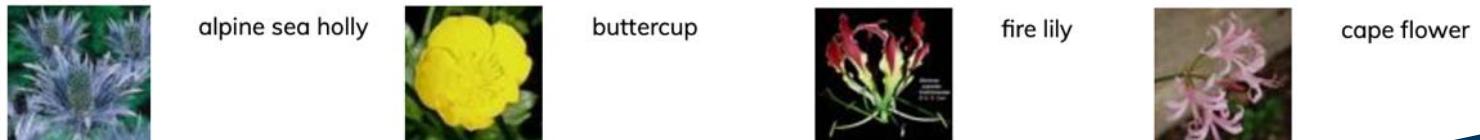
-

the accuracy will much lower than 99% or anything else  
even if 1000 images is small

# Transfer learning: motivation

Illustrative example

Suppose you want to build an app to recognize flowers...



Training a convolutional net from scratch won't work well

1000 images  
100 classes

# Transfer learning: the idea

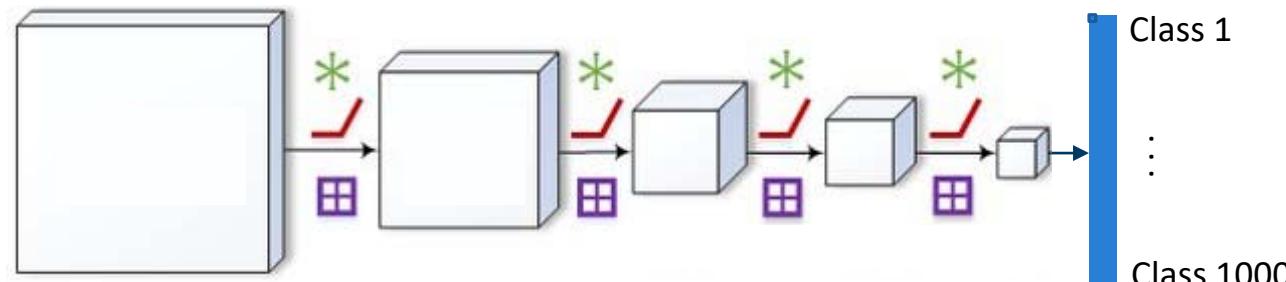
to solve the problem:

- train the classifier on much more images or reload it
  - afterwards fine tune it
- > should converge

1

**Pre-training**  
(massive data)

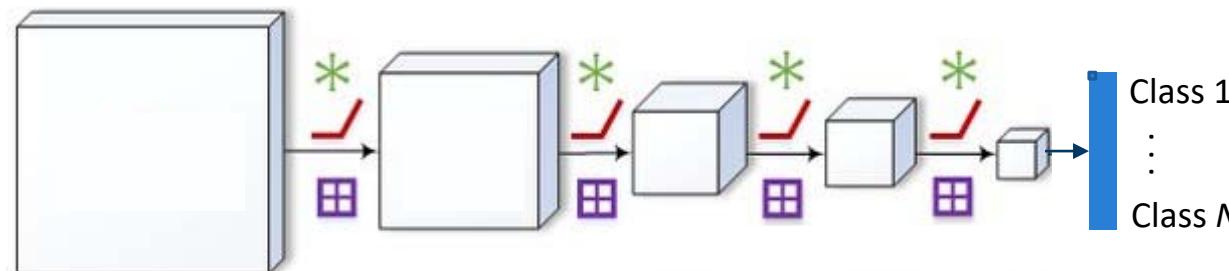
**IMAGENET**  
1.2 million images  
1000 object classes



IMAGENET

2

**Fine-tuning**  
(much less data)



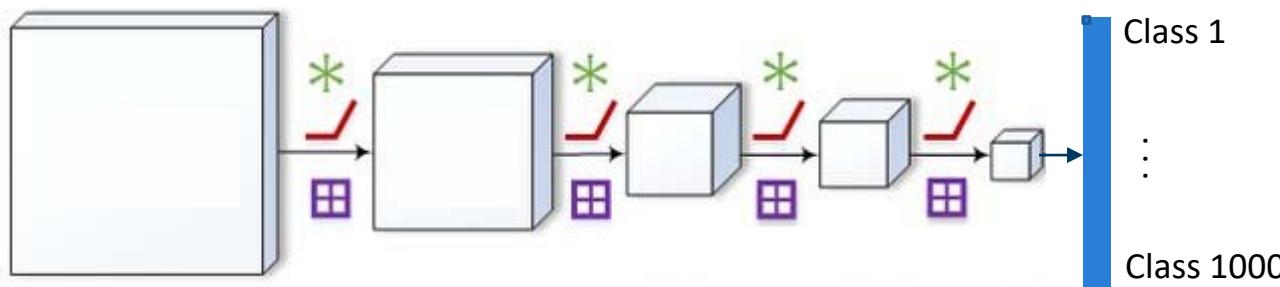
Recognize flowers

# Transfer learning: the idea

1

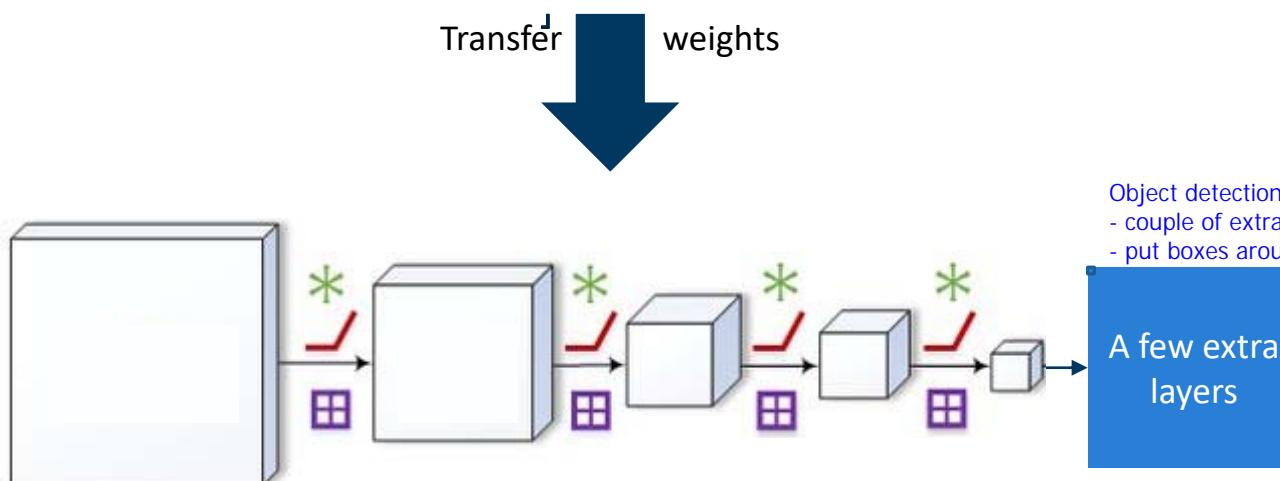
**Pre-training**  
(massive data)

**IMAGENET**  
1.2 million images  
1000 object classes



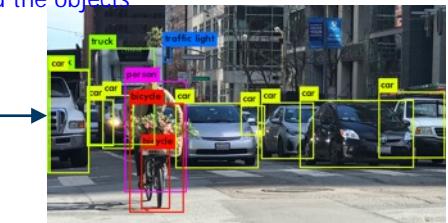
2

**Fine-tuning**  
(much less data)



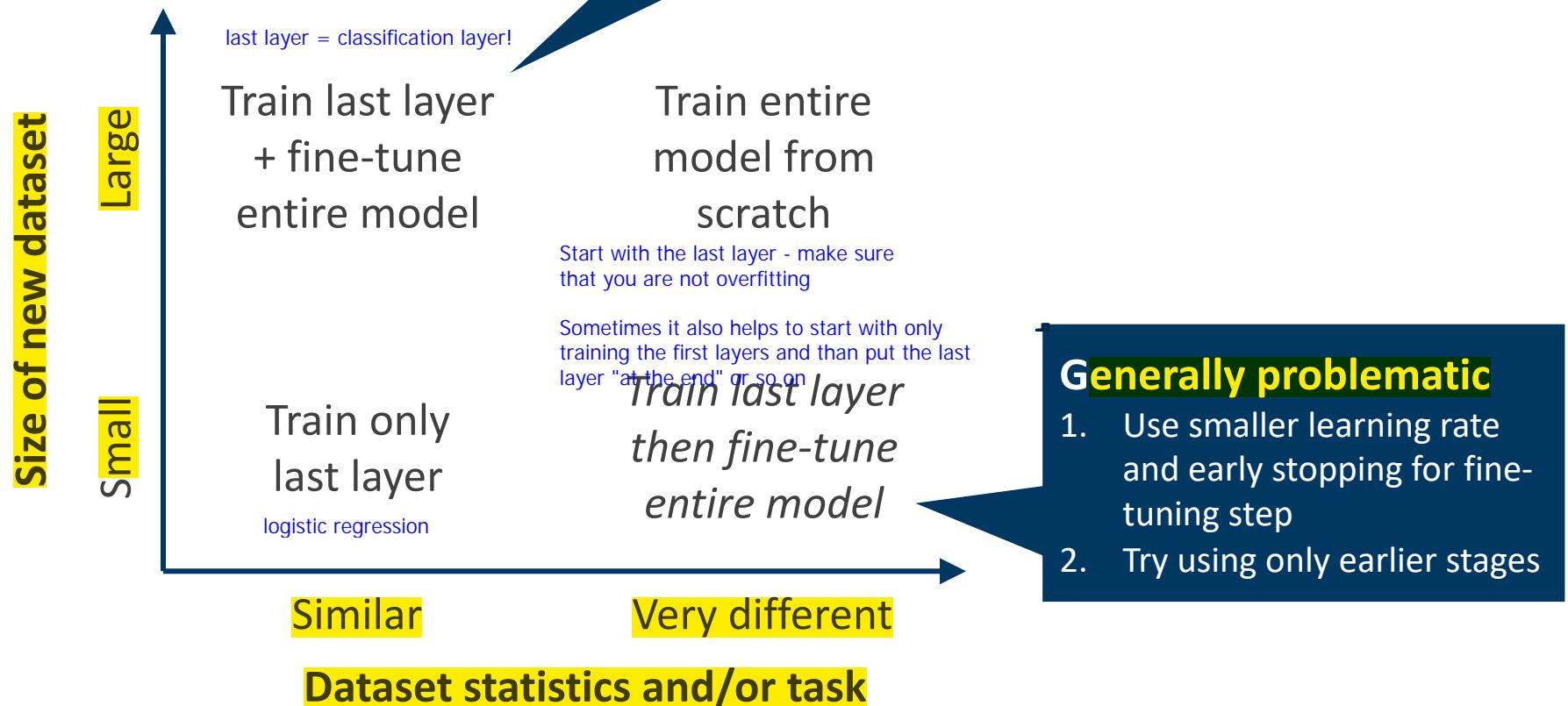
IMAGENET

Object detection:  
- couple of extra layers  
- put boxes around the objects



Object detection

# Strategies for training



# Better ImageNet lead to improvements on downstream tasks

the used the differnet architectures from the last lectures

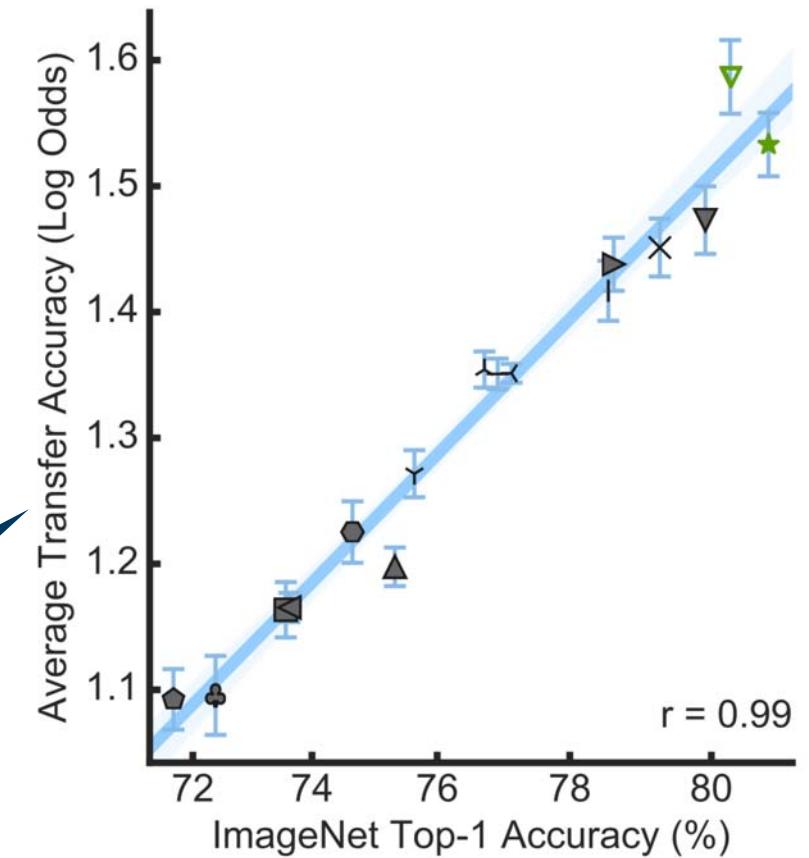
improvement by updating our classifier, no effort on the complexity

## Do Better ImageNet Models Transfer Better?

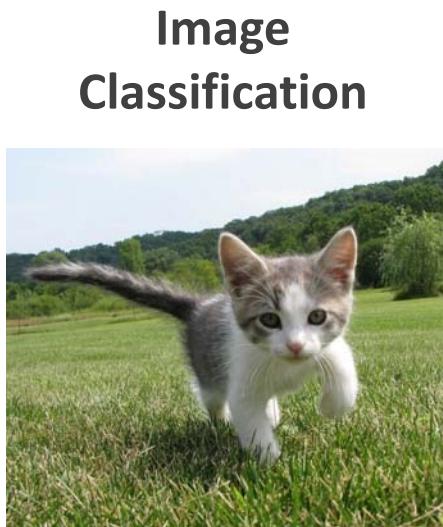
Simon Kornblith, Michael Shleifer, and Quoc V. Le  
Google Brain  
{skornblith, shlens, qvl}@google.com

Yes!

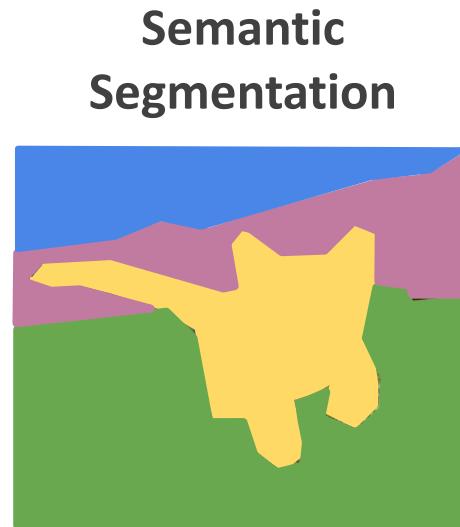
Tested on 12 different image classification sets



# Computer vision tasks



No spatial extent



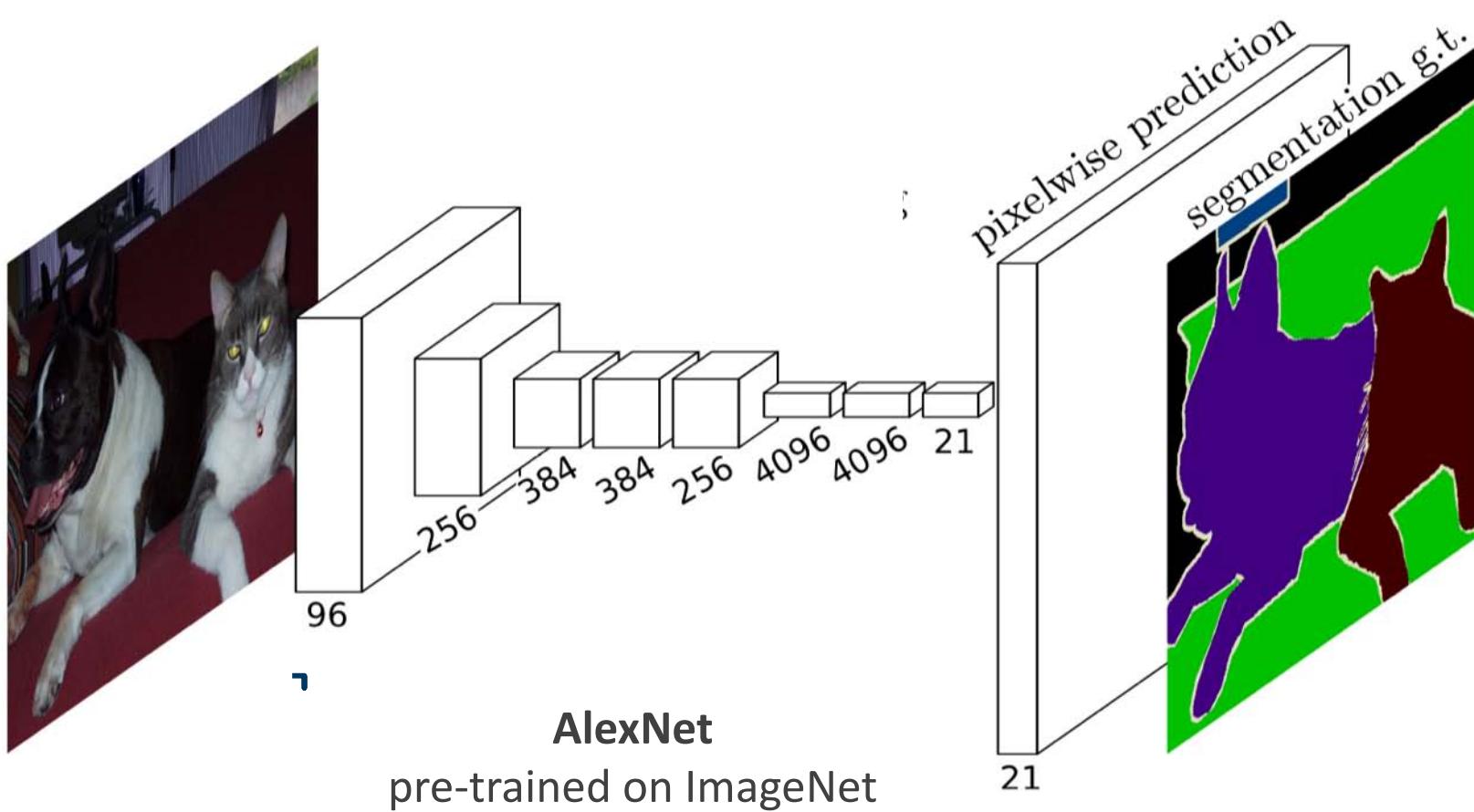
No objects, just pixels



Multiple objects

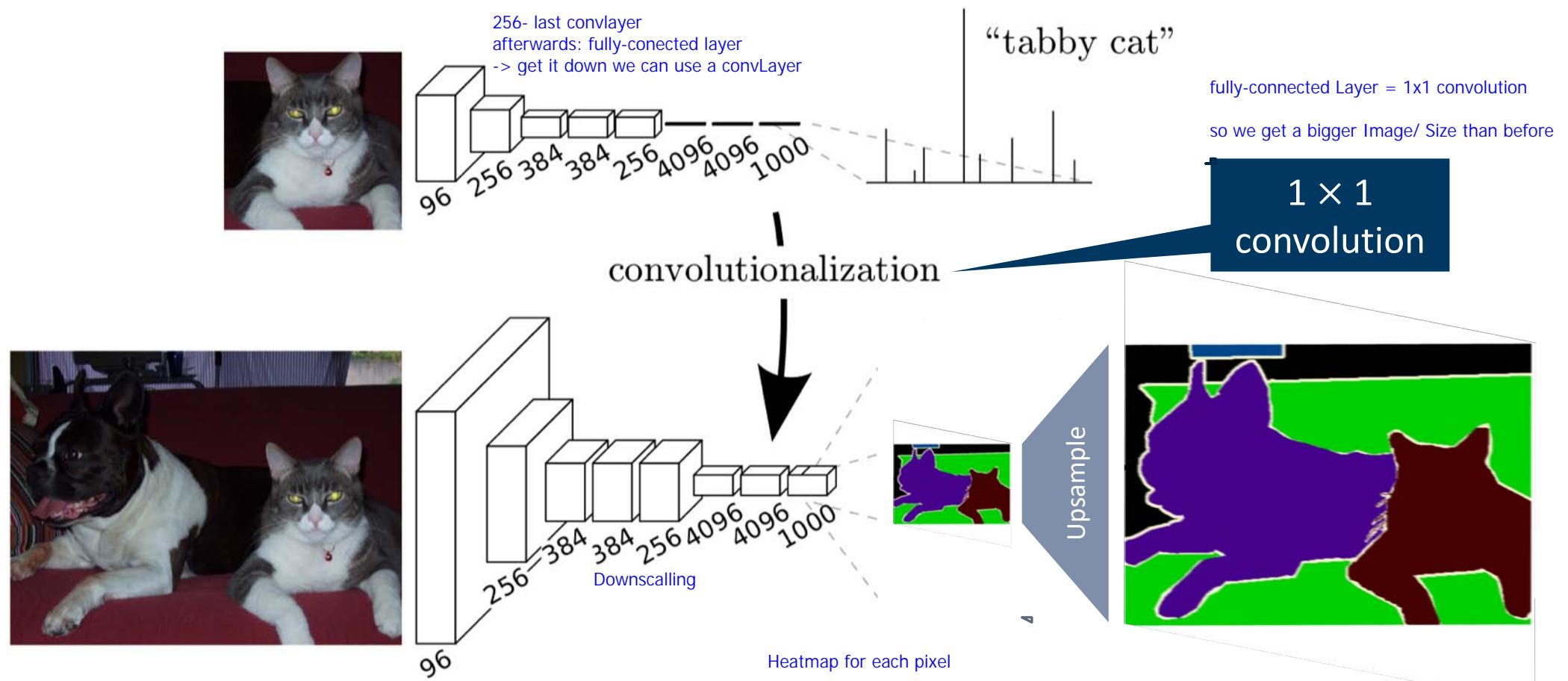


# Semantic segmentation

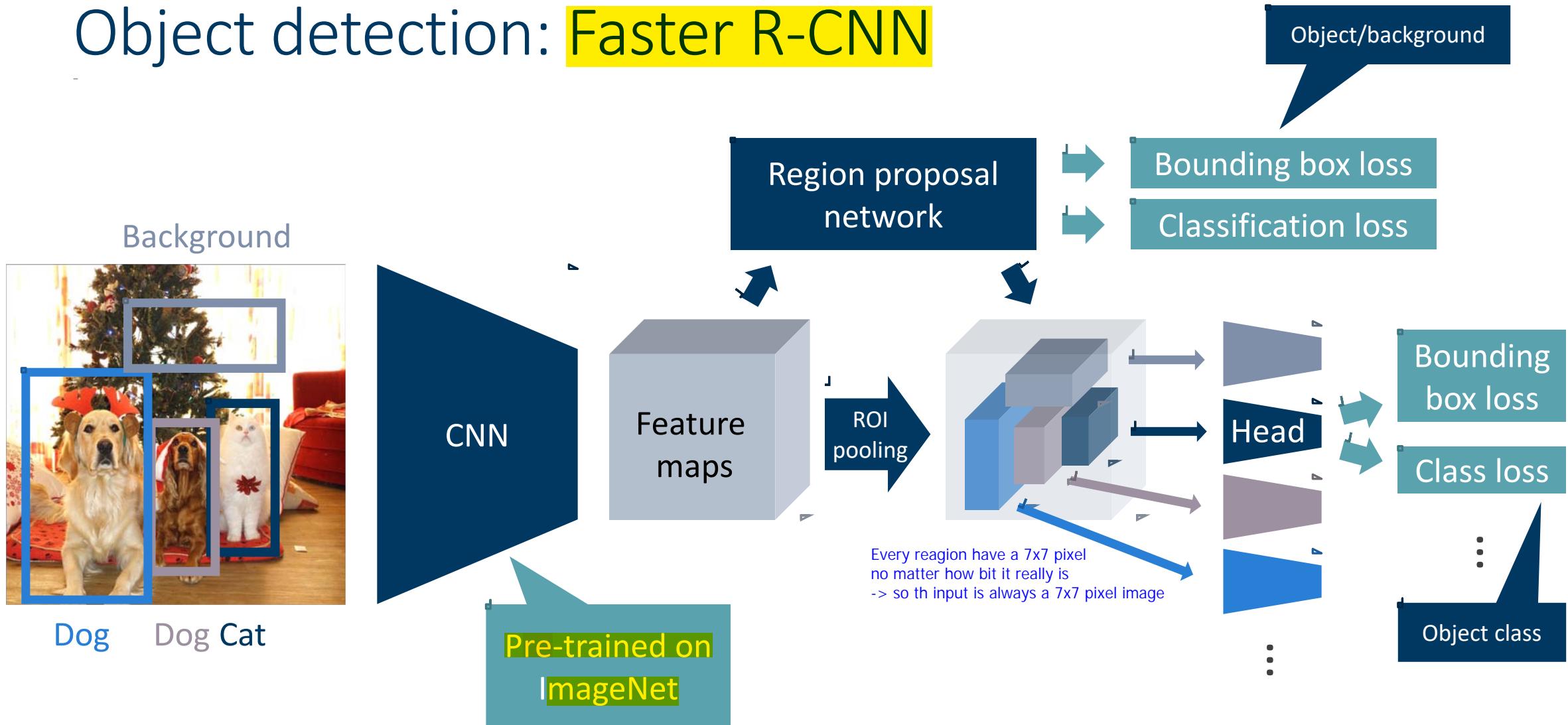


# Turning fully-connected layers into conv layers

Image -> feature map -> ... -> image -> feature map -> fully-connected layer

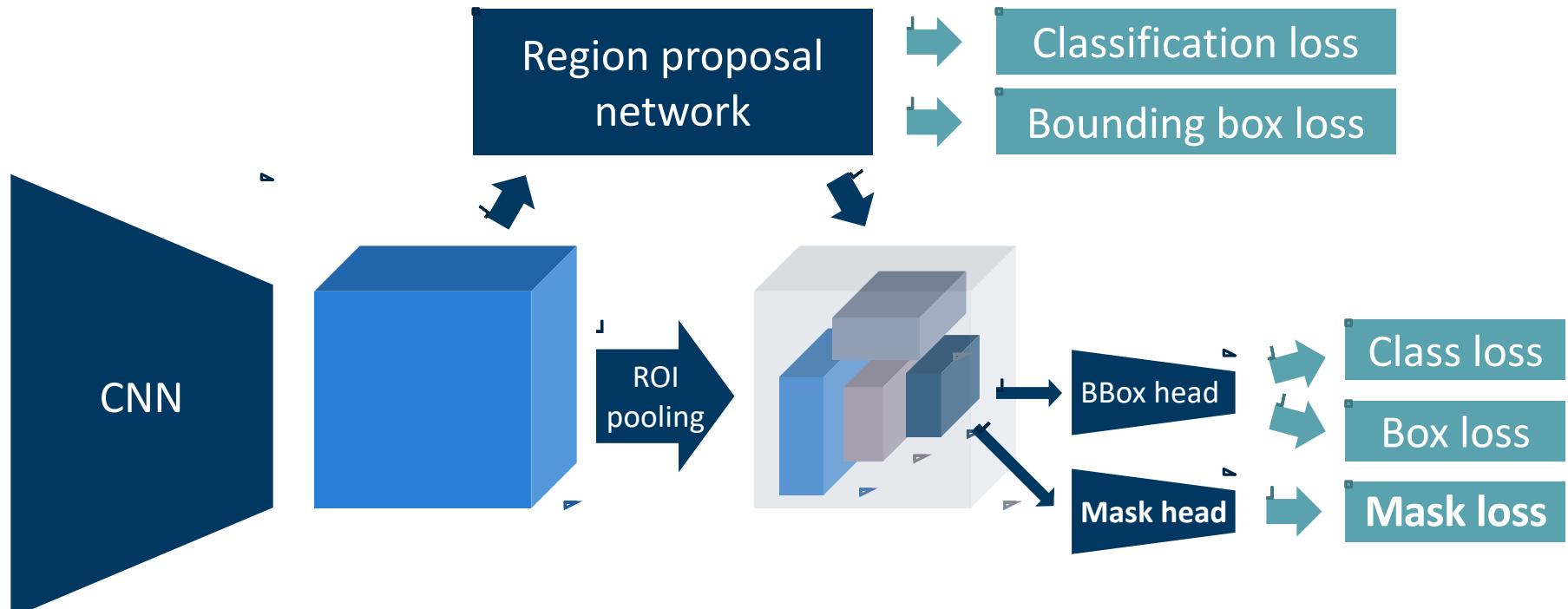
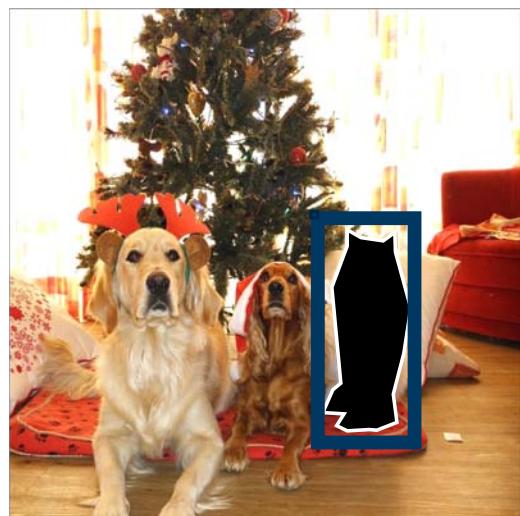


# Object detection: Faster R-CNN



which pixel is one object and which pixel is another object?!

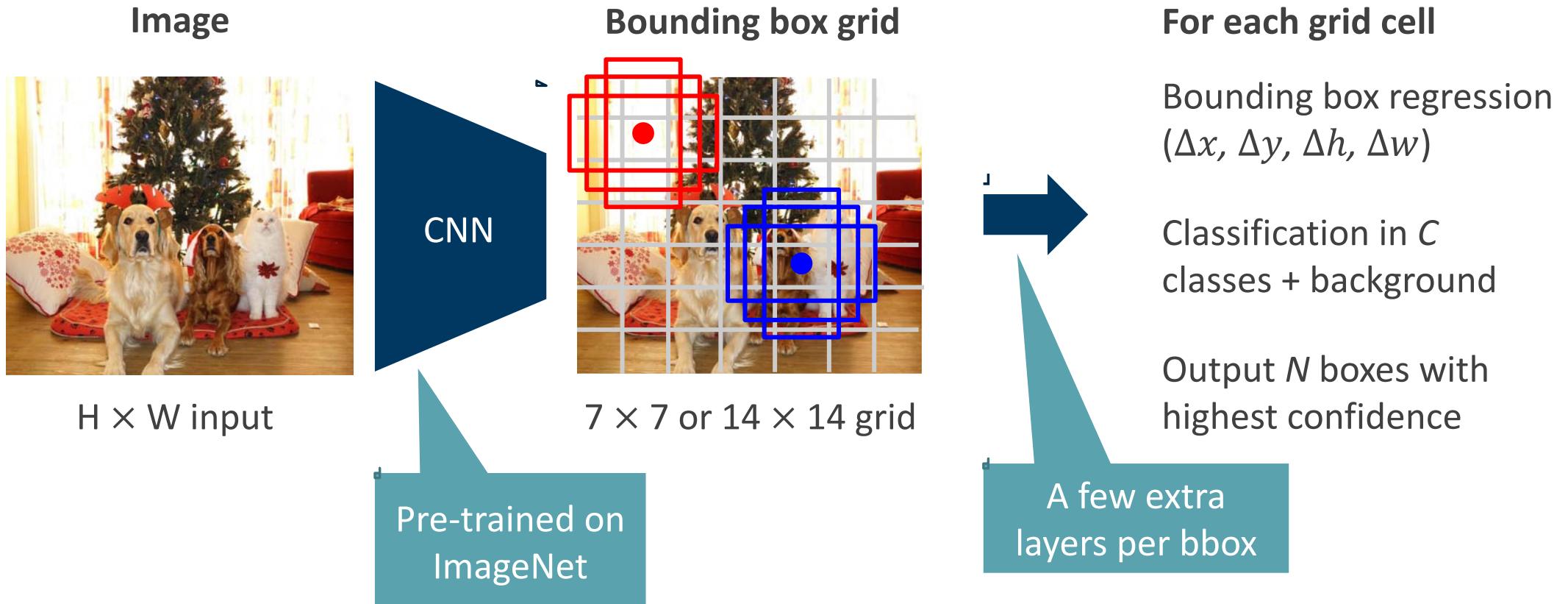
# Instance segmentation: Mask R-CNN



same idea/ aufbau, just add another head (Mask head) at the end  
-> some more convLayers etc.

# Single-stage detectors: SSD, YOLO, RetinaNet

Advantage: this systems are faster  
Disadvantage: but not so accurate



# Pre-training on ImageNet is everywhere ...

Human pose estimation

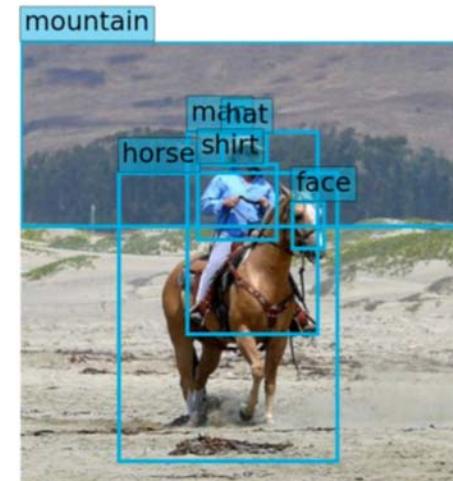


Image captioning



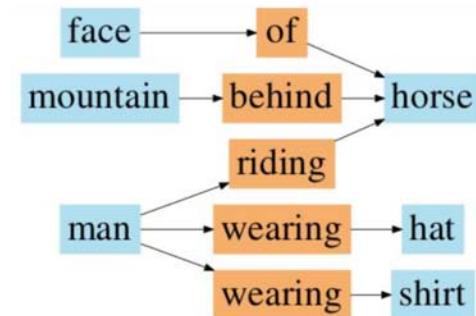
"two young girls are playing with  
lego toy."

Scene graph prediction



...

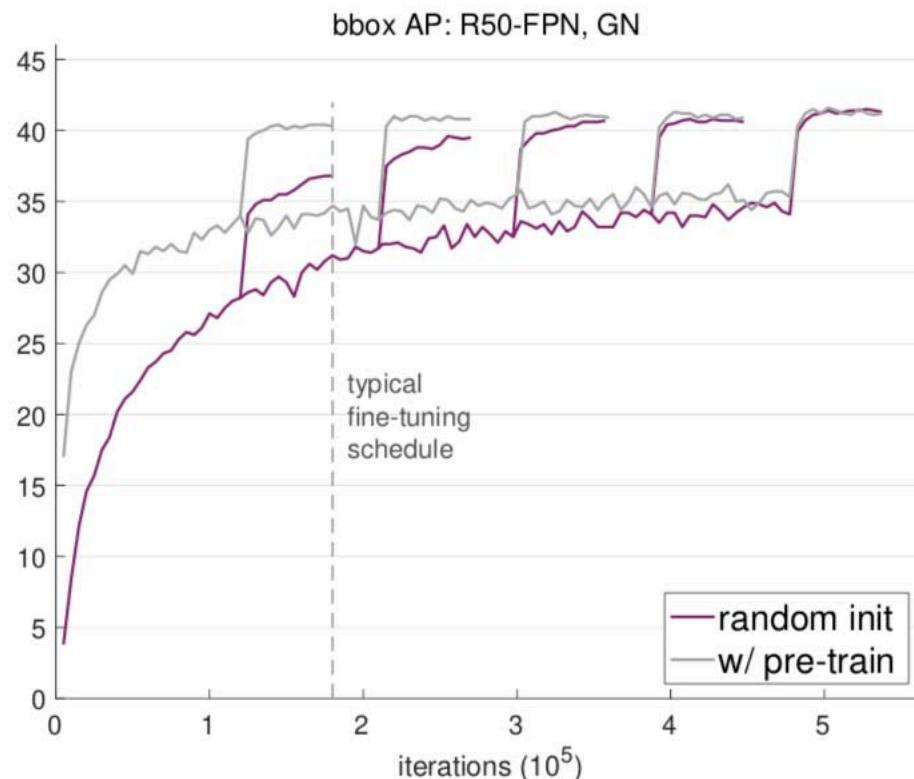
Many,  
many  
more



# Transfer learning from ImageNet is the standard...

... but may not always be necessary

Not always necessary to pre-train on image net,  
but is faster



Object detection  
on MS COCO dataset  
123,287 images  
886,284 instances

# Visual saliency

PREDICTING WHERE PEOPLE LOOK

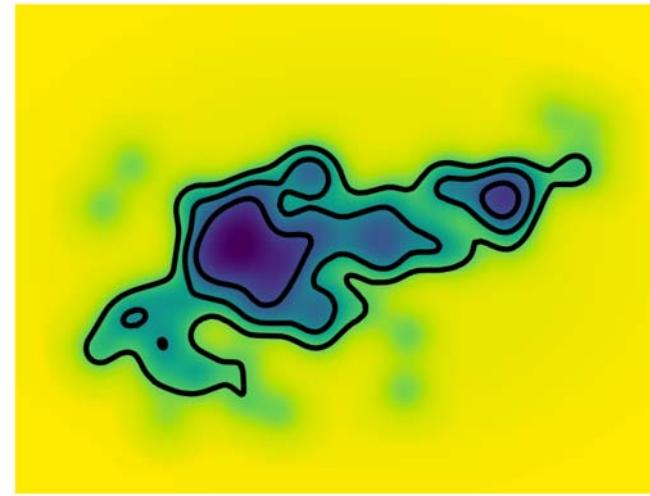
# Image saliency: predicting where people look

Image shown



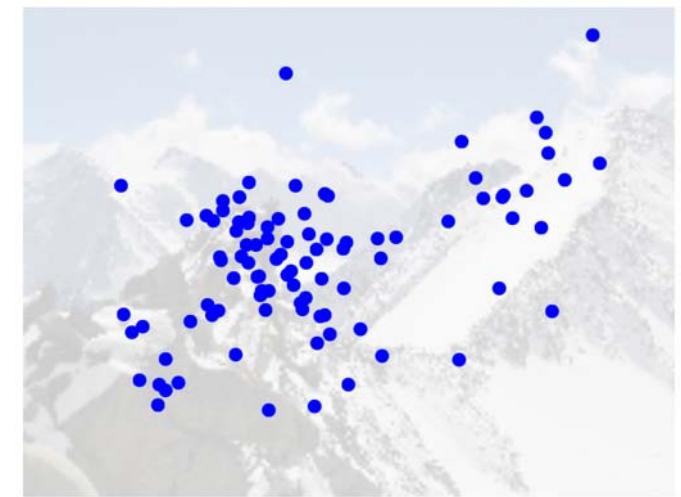
Predict

Goal:  
Saliency map

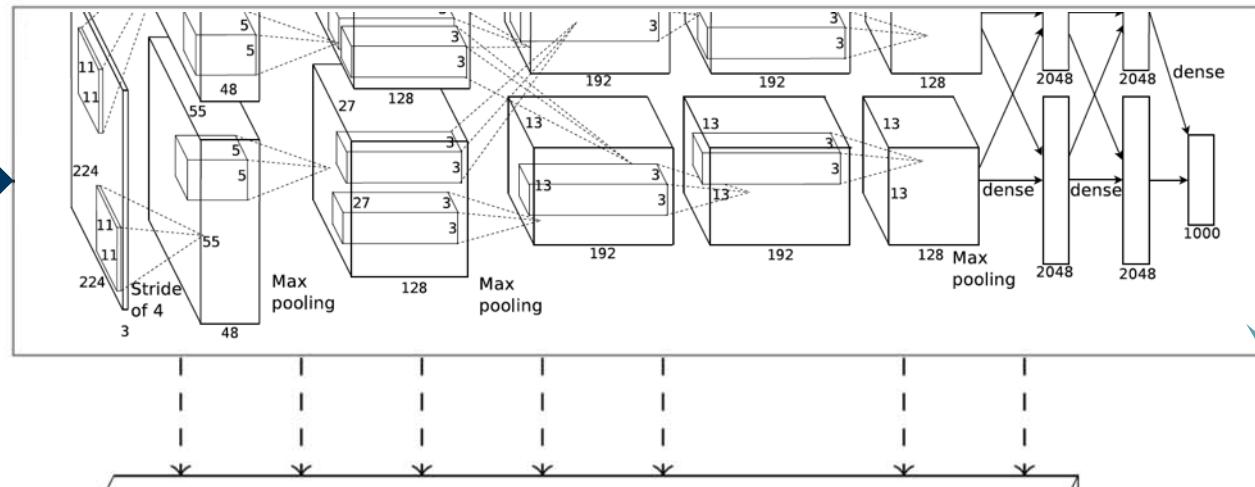


$$P(x, y \mid \text{image})$$

Measured eye fixations



Loss



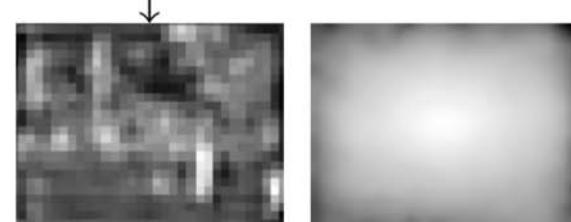
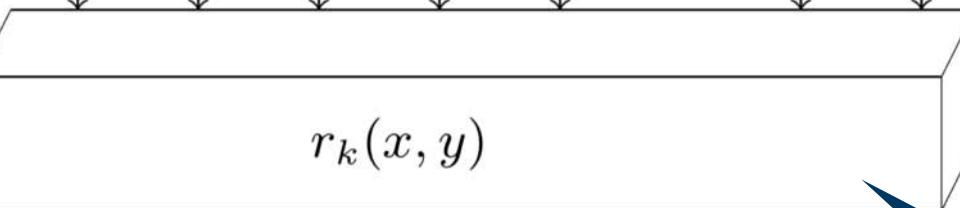
One slide is missing!  
Deep Gaze: loss fct!!

# Deep Gaze I

Kümmerer, Theis, Bethge,  
*ICLR Workshops 2015*

just linear comb., no normalization

Linear prediction



Blurring  
(as regularization)

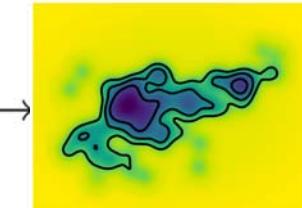


like a logLikelihood

Large filterbank

A-priori saliency map  
(center bias)

softmax



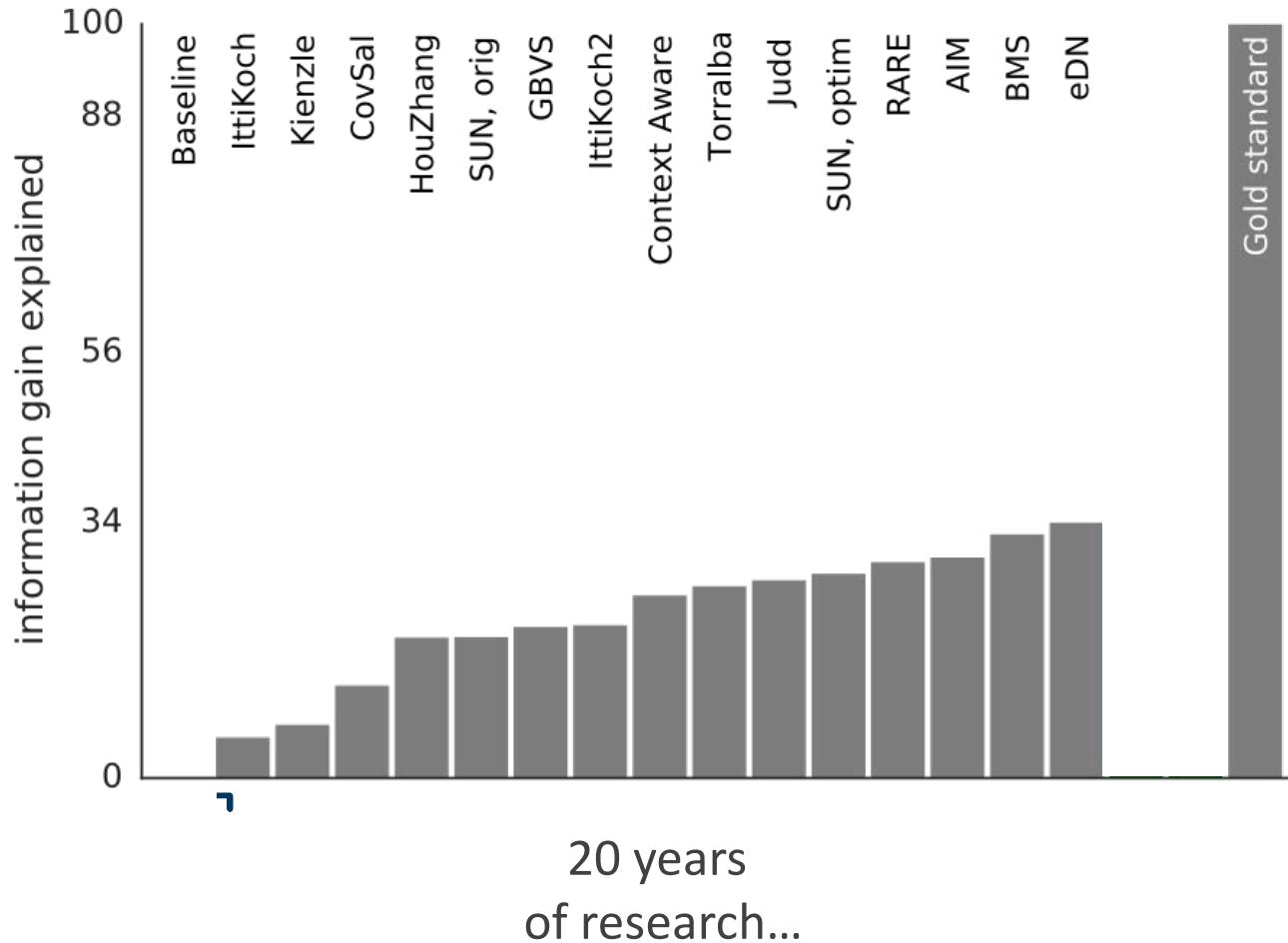
Saliency map  
(probability distribution)

for the homework we use instead of AlexNet: VGG

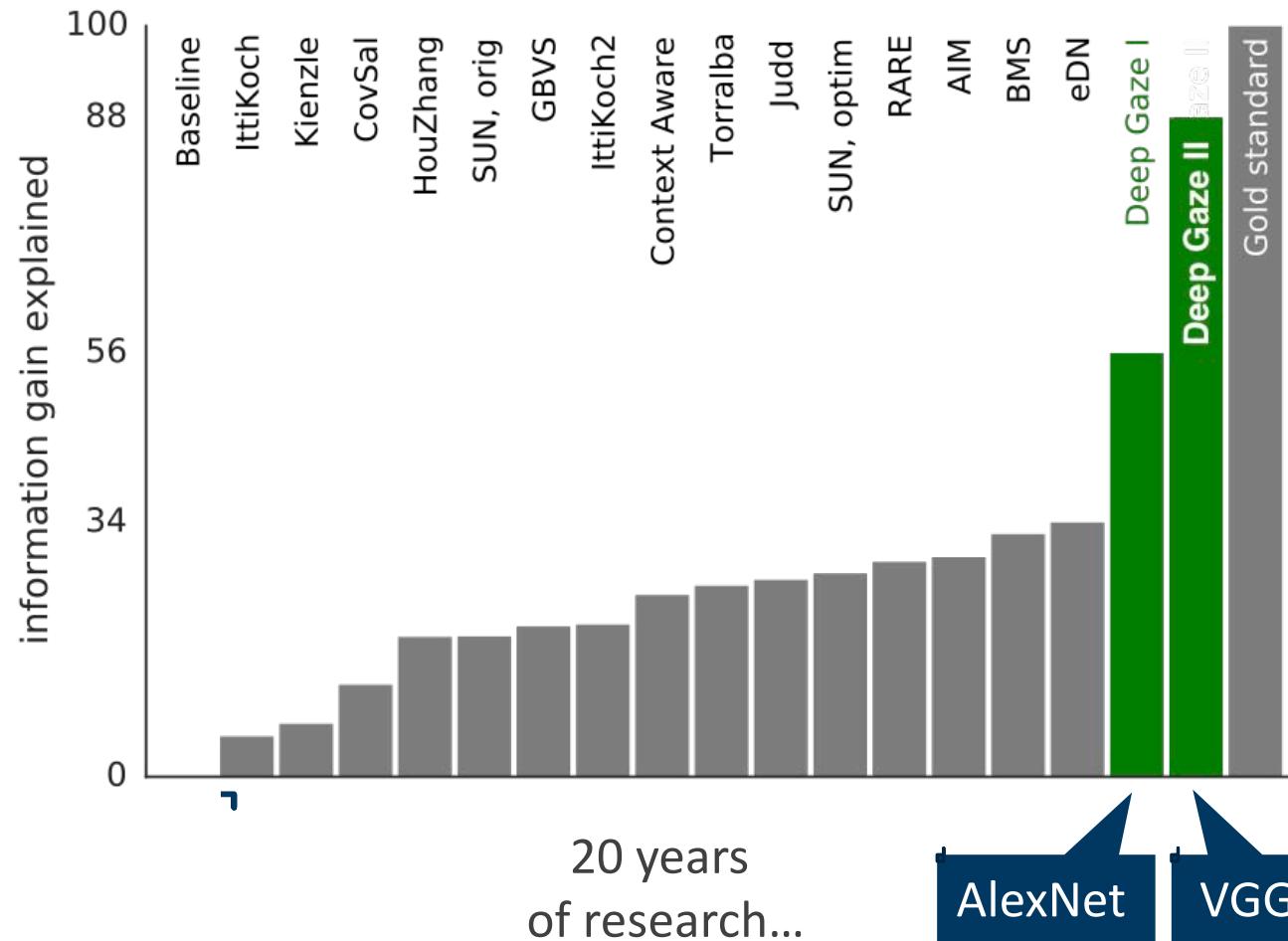
**AlexNet**  
pre-trained on  
ImageNet

it is an additional step  
-> when people look at image the  
look firstly at the center, because of  
the pos of the head

# Deep Gaze I + II



# Deep Gaze I + II



Kümmerer, Theis, Bethge,  
ICLR Workshops 2015

Kümmerer, Wallis, Bethge, arXiv 2016

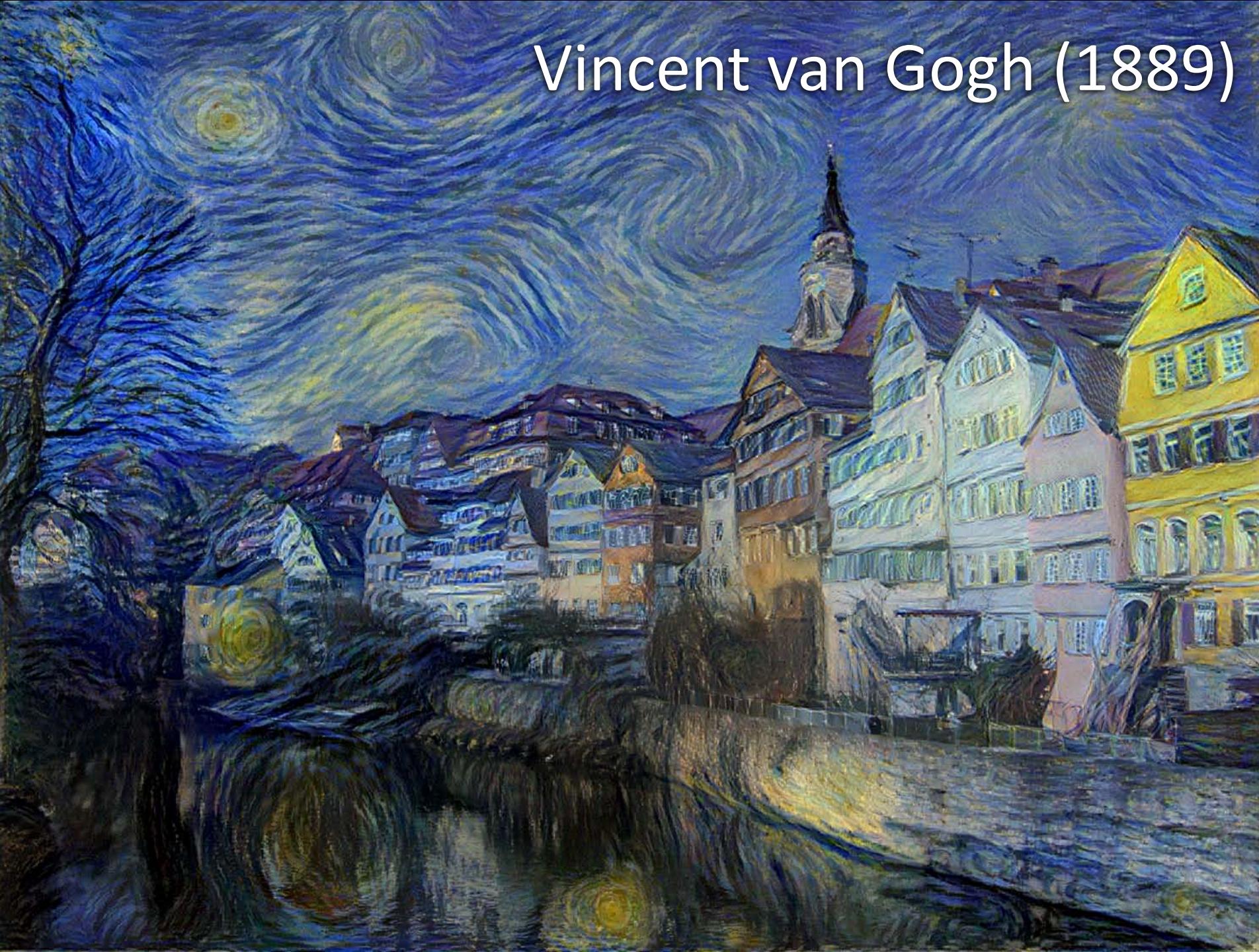
# Neural style transfer

TEXTURE SYNTHESIS AND PERCEPTUAL LOSS FUNCTIONS

# History of art – Tübingen



# Vincent van Gogh (1889)



William Turner (1805)



Pablo Picasso (1910)



Edvard Munch (1893)



Wassily Kandinsky (1913)



# A Neural Algorithm of Artistic Style

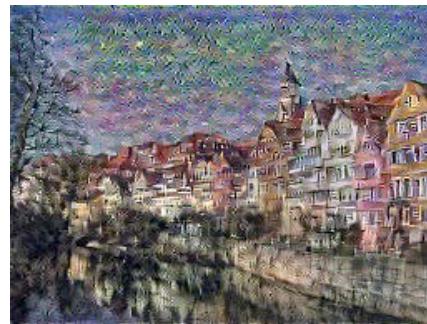


Gatys, Ecker, Bethge (CVPR 2016)



# Neural Style Transfer

Content



Style (→ texture)



# A brief history of parametric texture modeling

Julesz (1962)

- Conjecture:  $N^{\text{th}}$  order summary statistics determine texture perception

Portilla & Simoncelli (2000)

- Summary statistics of features from early visual system (steerable pyramid)

Gatys, Ecker Bethge (NIPS 2015)

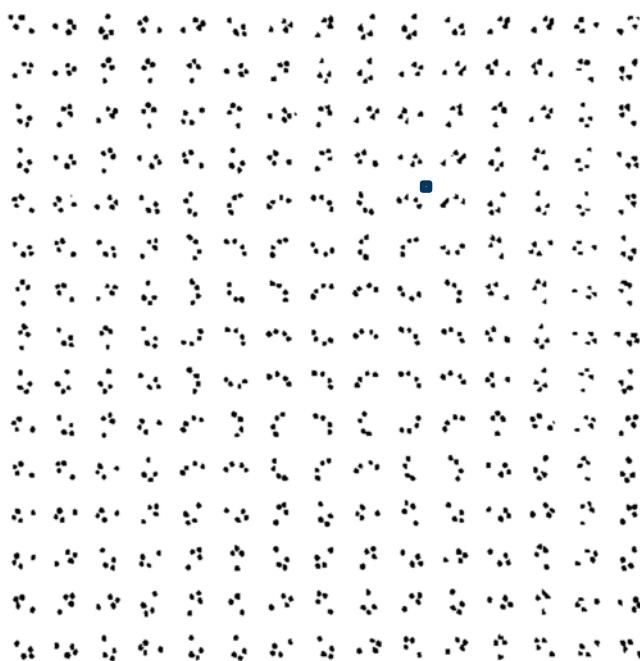
- Summary statistics of features from CNNs trained on ImageNet

# Julesz' conjecture (1962)

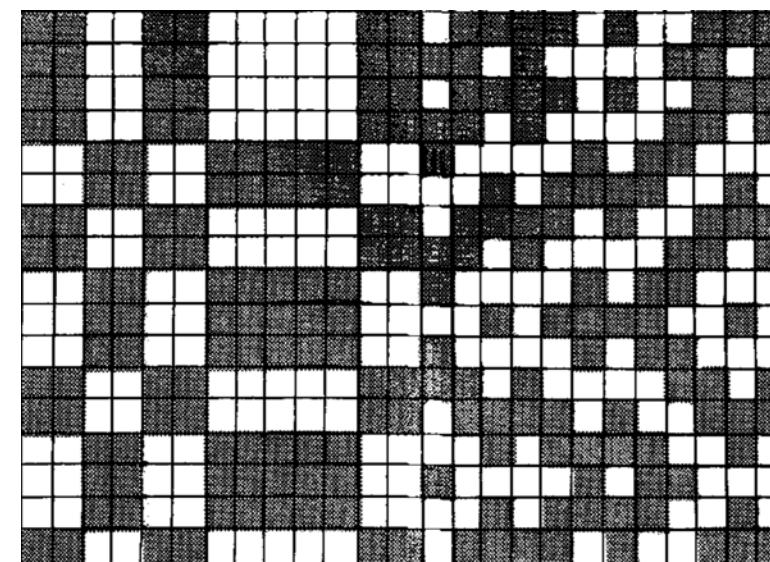
2nd and 3rd order statistic is not enough, as we can see  
-> so it was a nice idea, but not really works

**$N^{\text{th}}$  order summary statistics determine texture perception.**

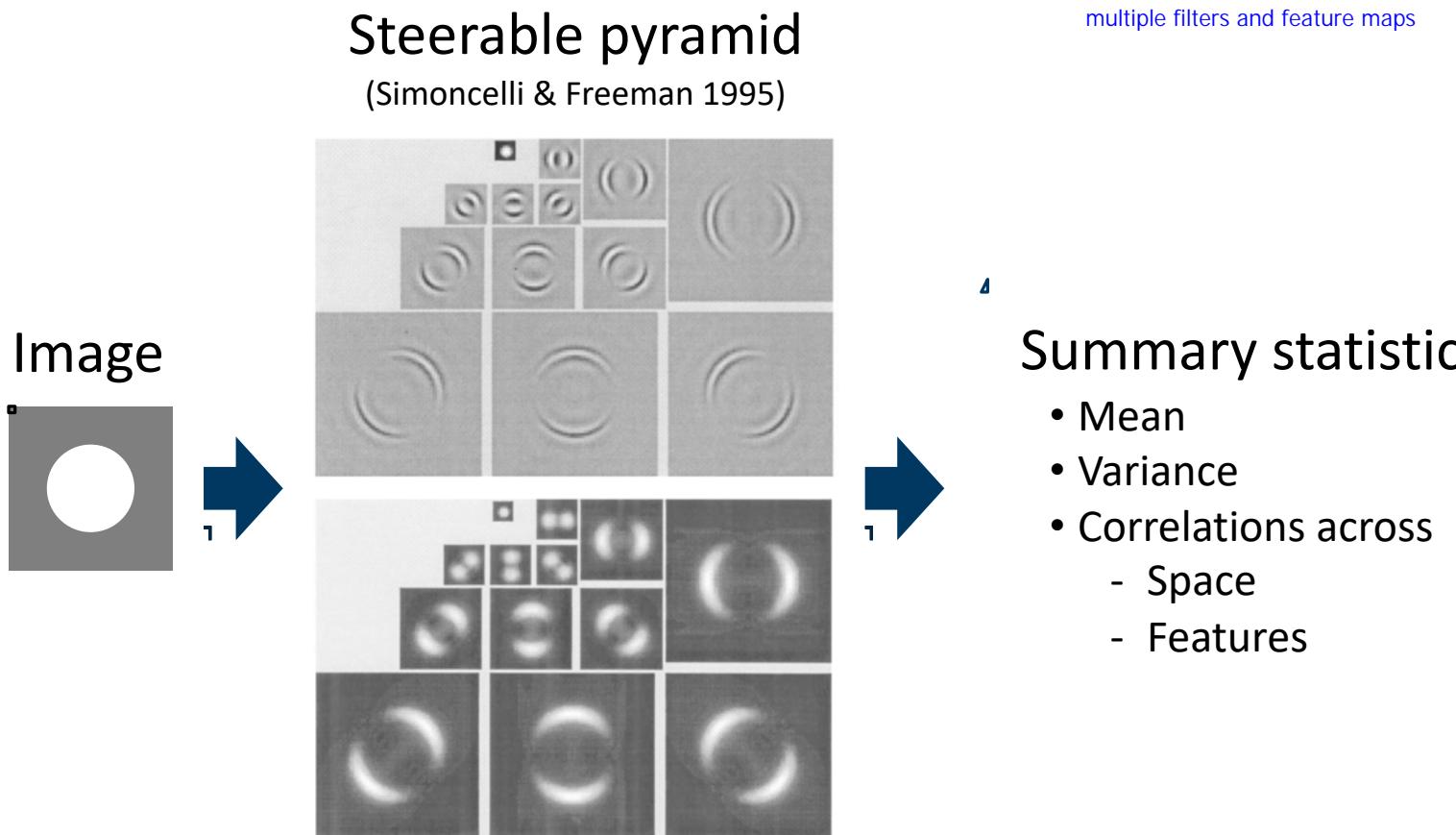
Uniform 2<sup>nd</sup> order statistics  
(power spectrum)



Uniform 3<sup>rd</sup> order statistics



# Portilla & Simoncelli texture model

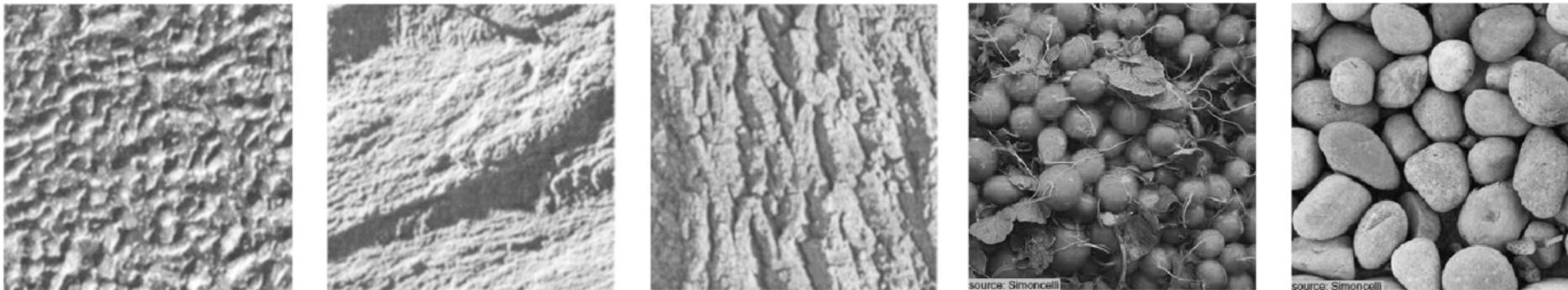


- ## Summary statistics
- Mean
  - Variance
  - Correlations across
    - Space
    - Features

# Portilla & Simoncelli: results

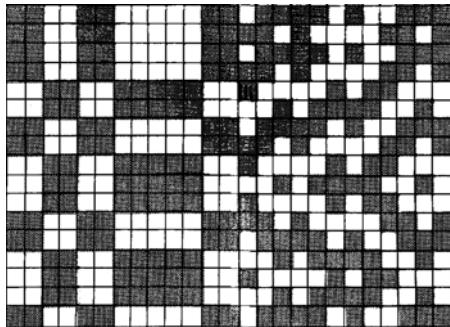
not the same in terms of pixels, but the same in terms of texture  
-> works for some well, but not for all

Original images

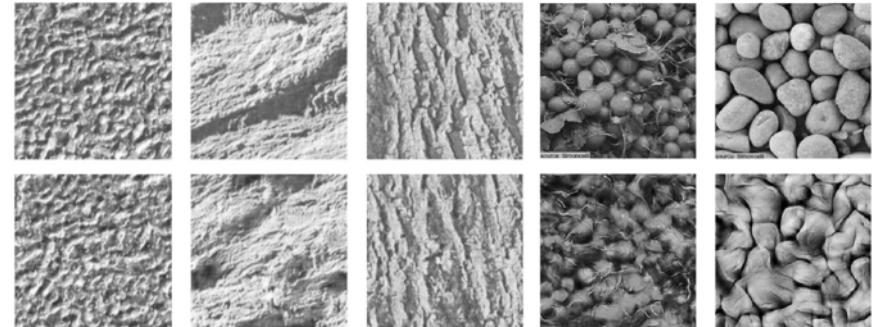


Synthesized by Portilla & Simoncelli

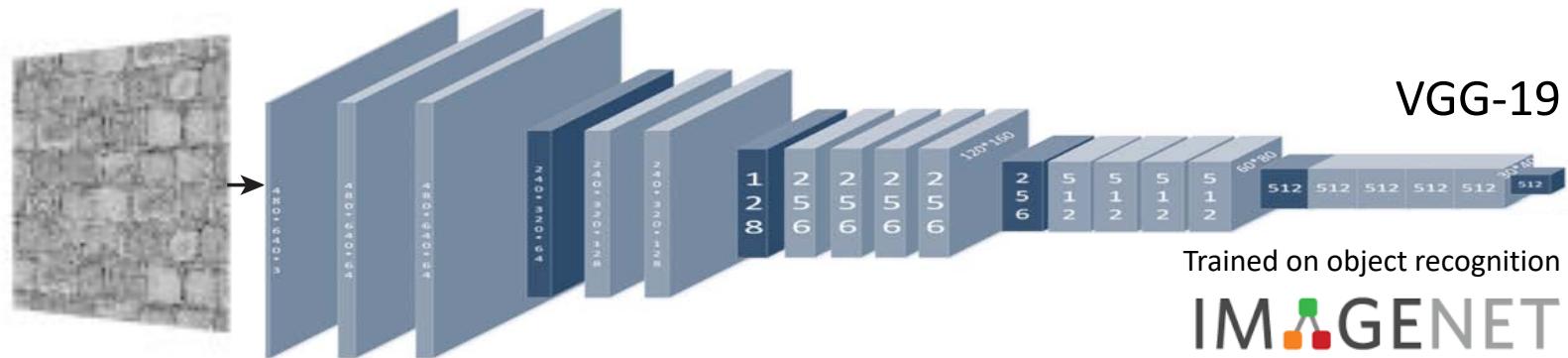
# Texture models: summary statistics



Julesz (1962):  
 $N^{\text{th}}$  order pixel stats

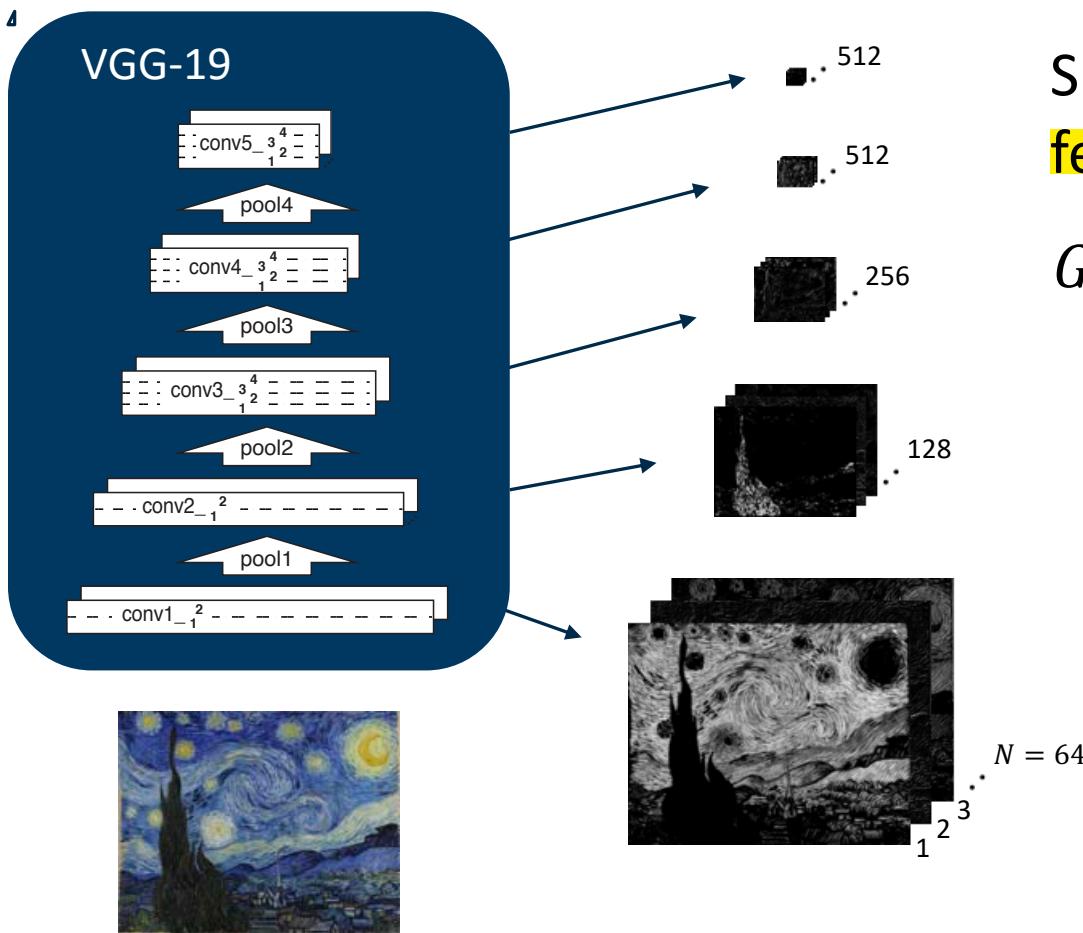


Portilla & Simoncelli (2000):  
V1 summary statistics → Steerable Pyramid



Gatys, Ecker, Bethge (2015): Entire visual pathway → CNN features

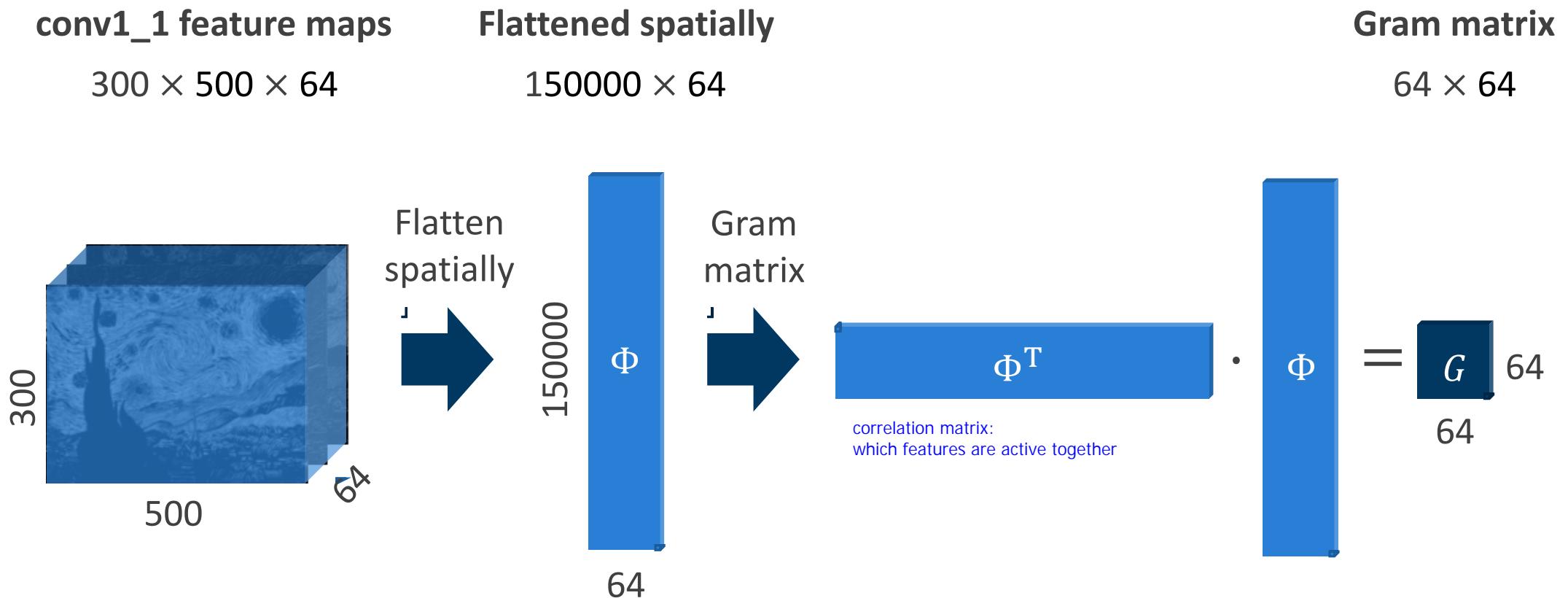
# CNN: Multiscale nonlinear filter bank



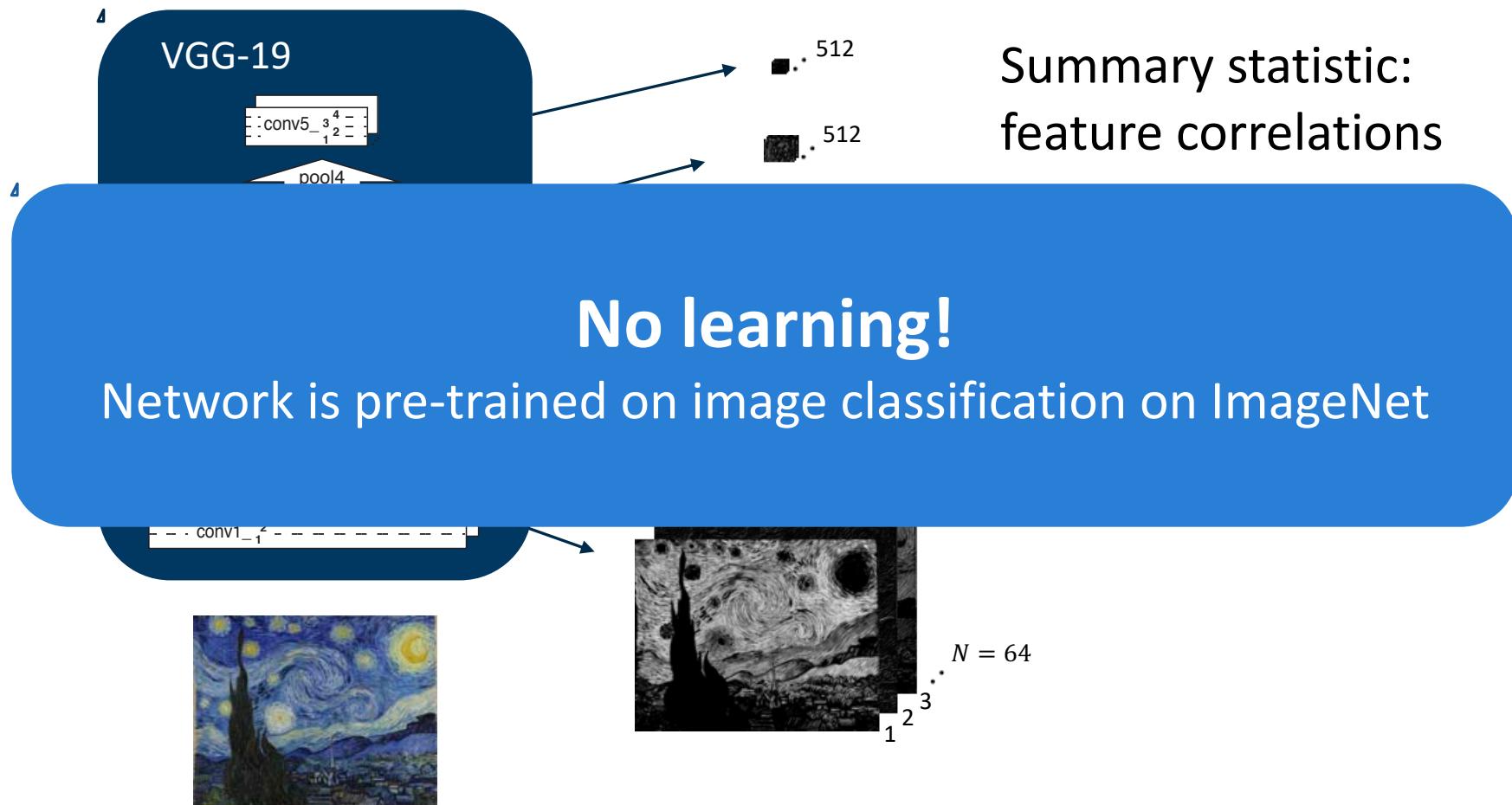
Summary statistic:  
feature correlations

$$G = \Phi^T \Phi$$

# Gram matrix computation



# CNN: Multiscale nonlinear filter bank

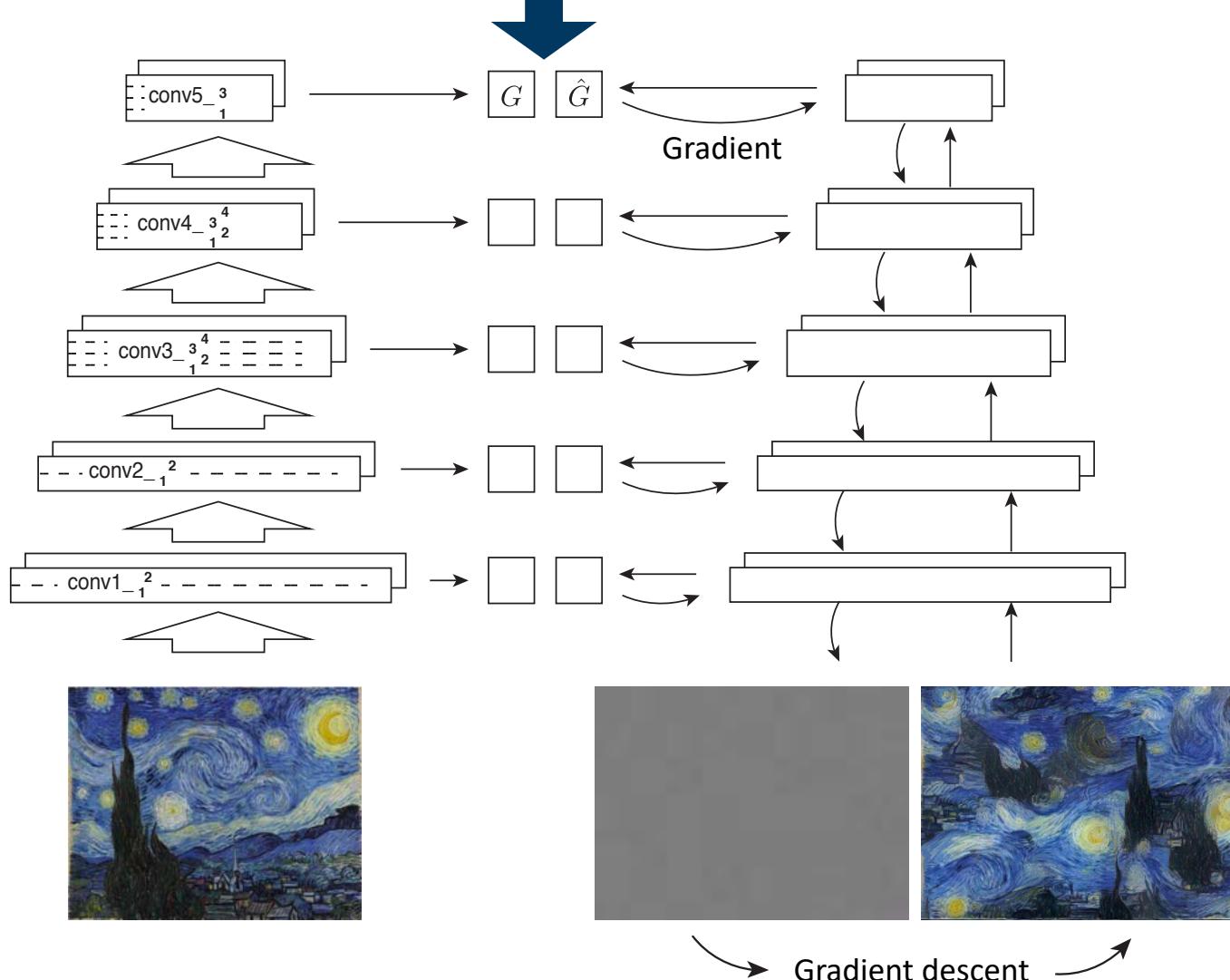


## Texture loss

$$\mathcal{L}_{\text{style}} = \sum_{ij} (G_{ij} - \hat{G}_{ij})^2$$

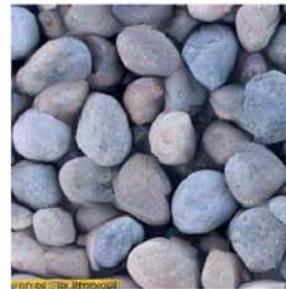
MSE between the target image and the real image

- > no updating the weights
- > updating the images

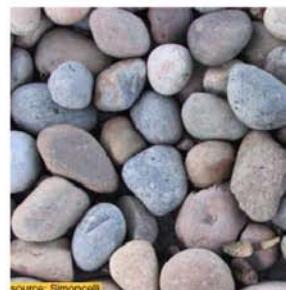


# Textures based on deep features

CNN-based model



Original image

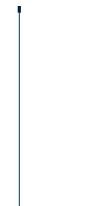
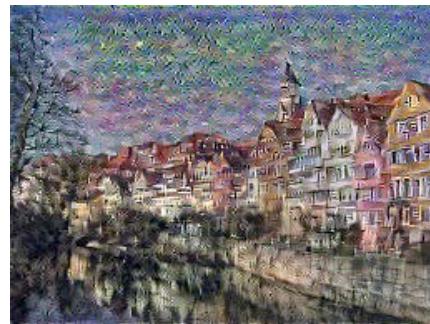


Portilla & Simoncelli



# Neural Style Transfer

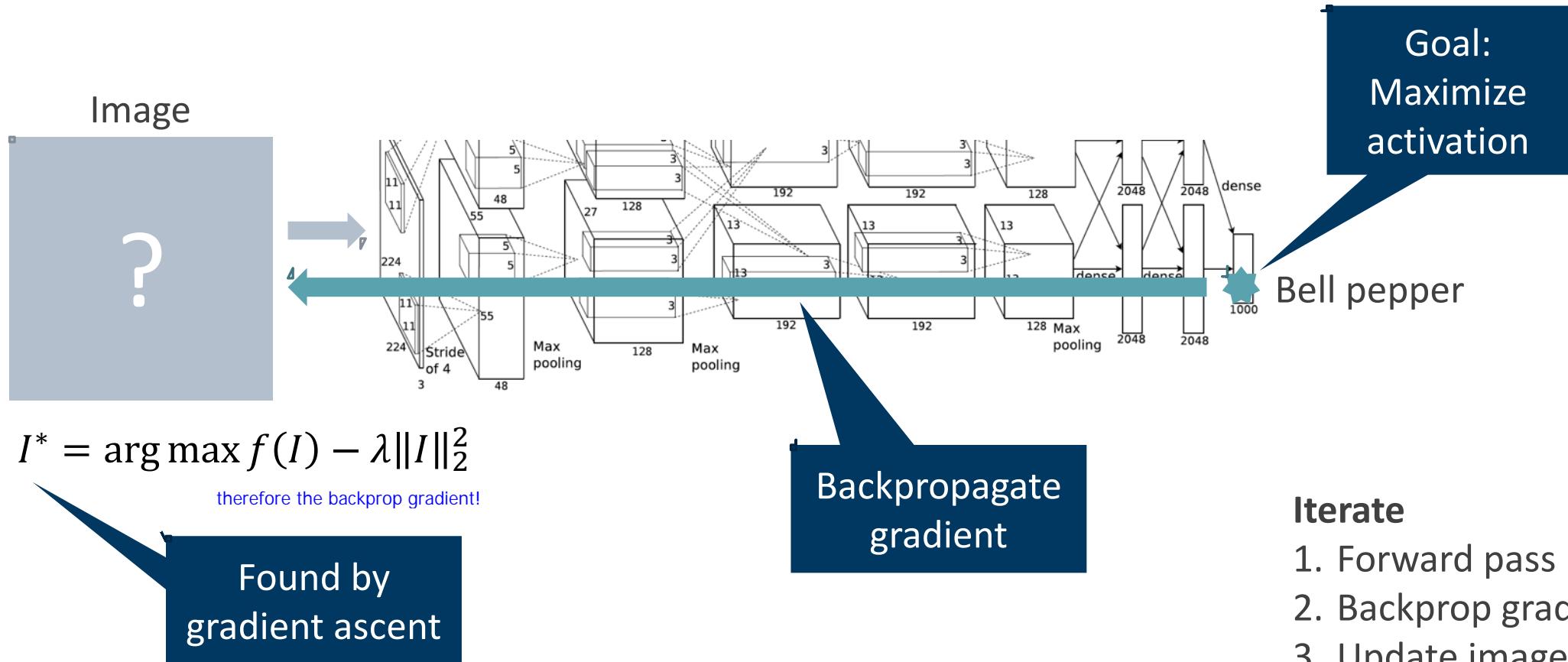
Content



Style (→ texture)



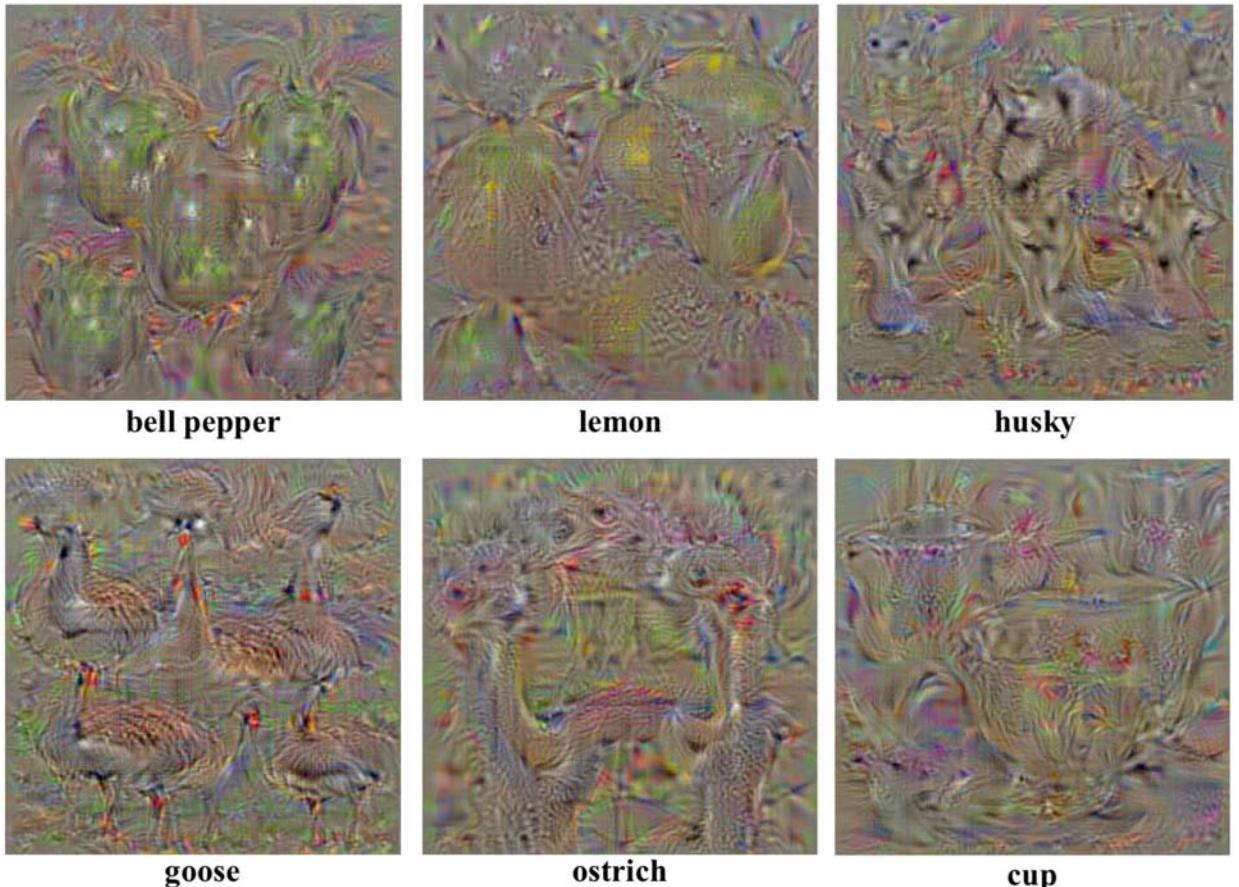
# Visualizing CNN features by activity maximization



# Visualizing classification layer

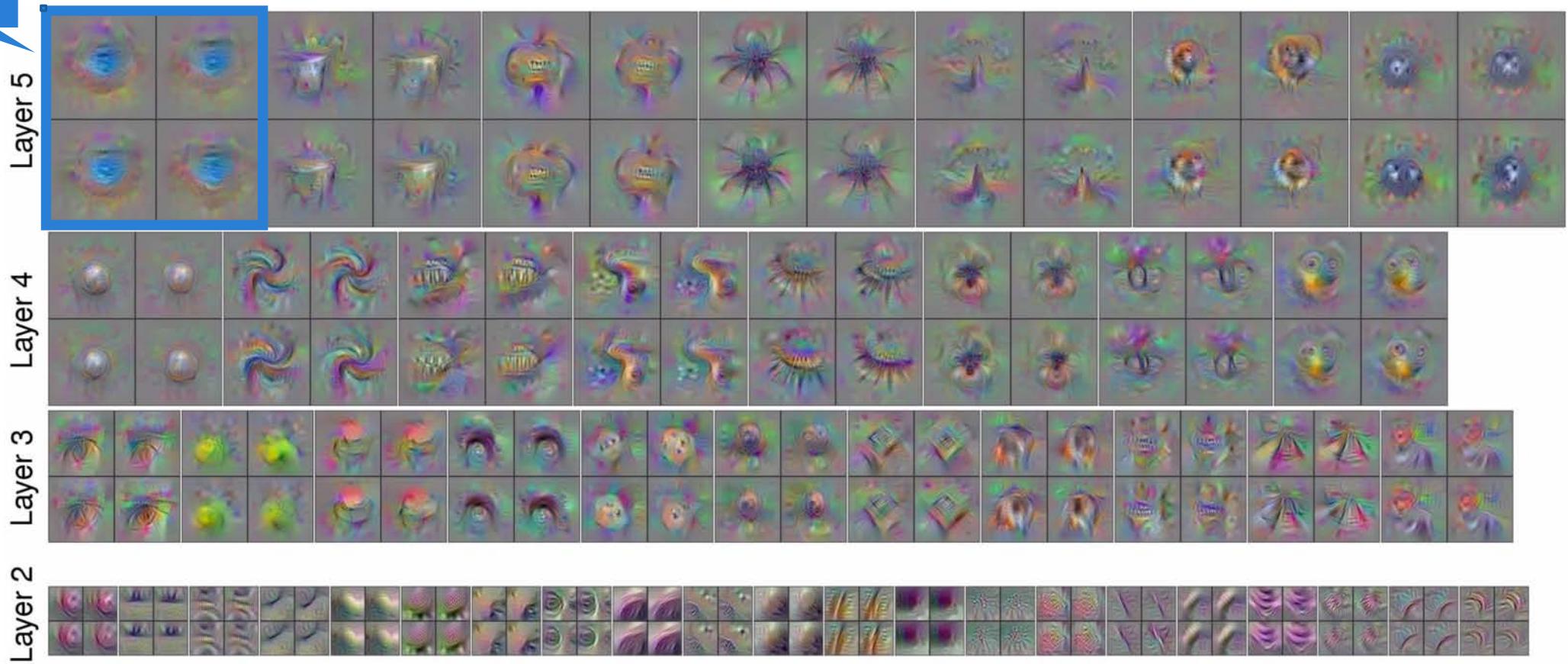
$$I^* = \arg \max f(I) - \lambda \|I\|_2^2$$

Regularizer to get  
“natural” images



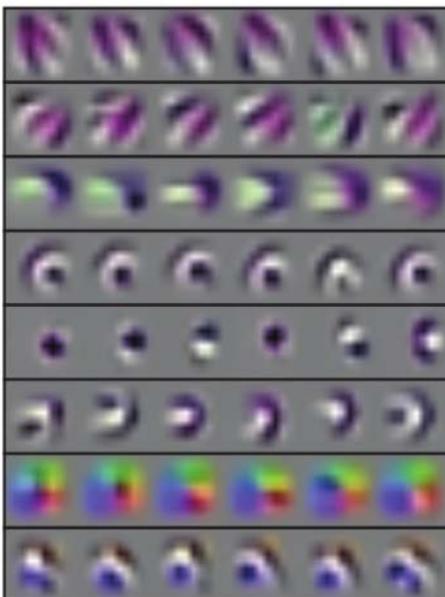
# Visualization of intermediate layers of AlexNet

1 unit



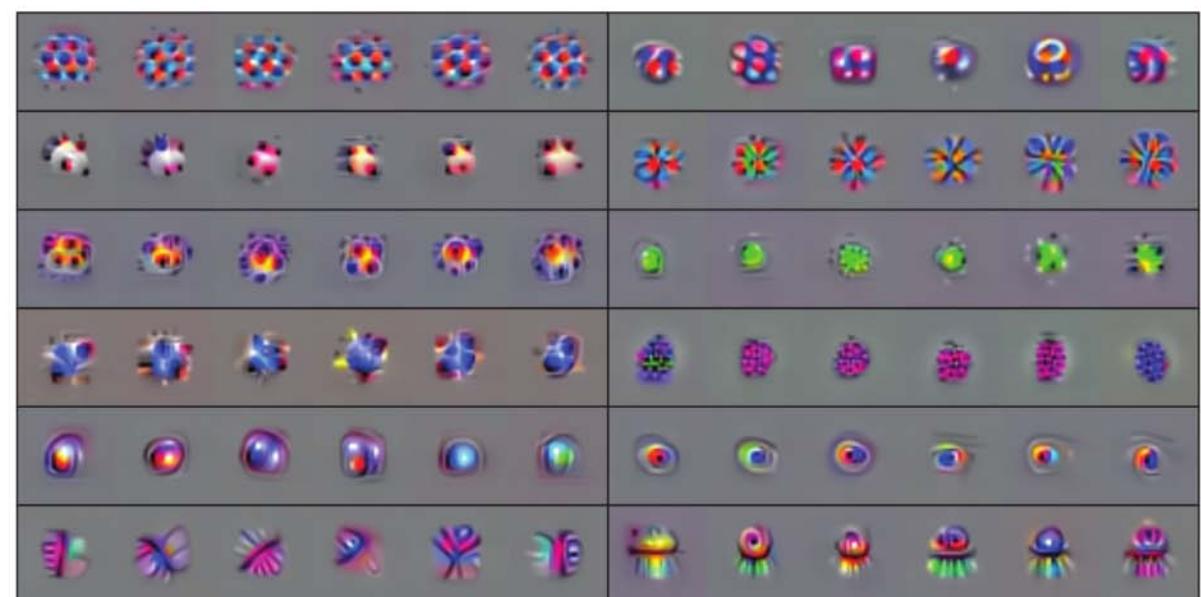
# Visualization of intermediate layers of VGG-19

conv2\_1



Early: selective

conv3\_3



Intermediate/late: Increased invariance

**Optimization target here:** obtain diverse set of images that maximize overall activity

## Style loss

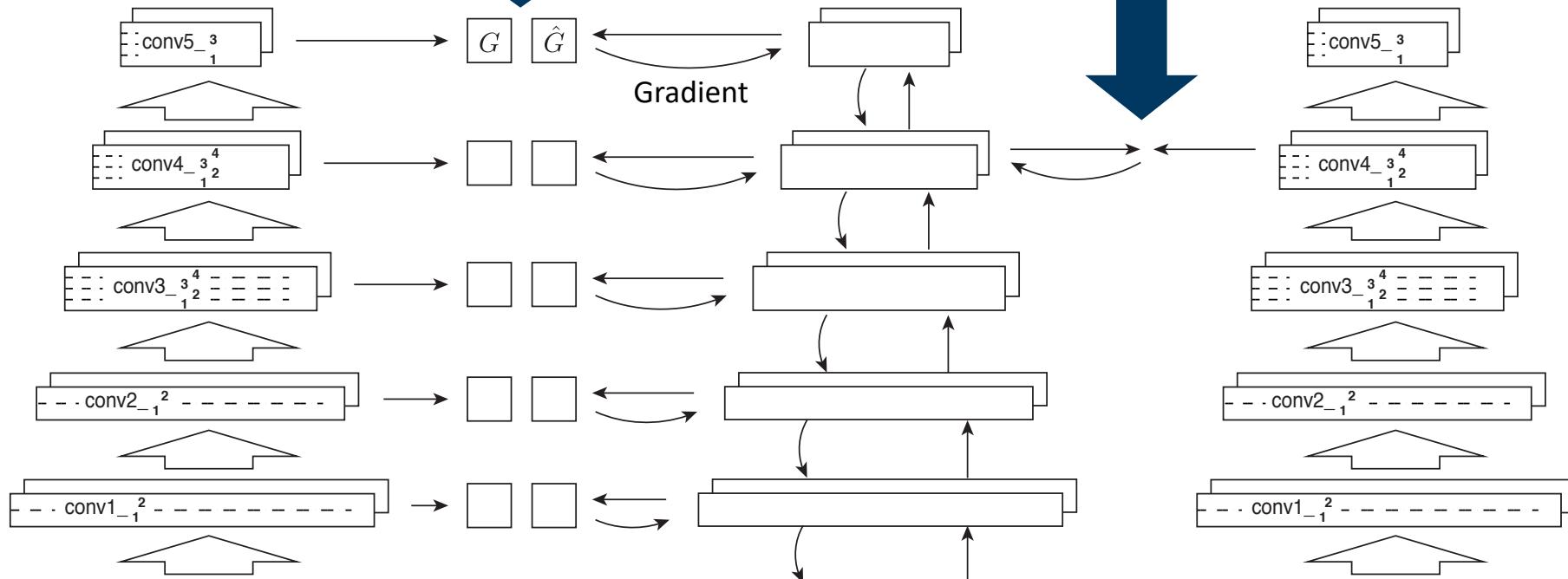
$$\mathcal{L}_{\text{style}} = \sum_{ij} (G_{ij} - \hat{G}_{ij})^2$$

## Content loss

$$\mathcal{L}_{\text{content}} = \sum_{ijk} (\Phi_{ijk} - \hat{\Phi}_{ijk})^2$$

Matching the individual activations

Content Loss:  
MSE between feature activations



Gradient descent

Gatys, Ecker Bethge, CVPR 2016

# Timelapse

---





Create your own artwork  
<https://dedepart.io>

## TURN ANY PHOTO INTO AN ARTWORK – FOR FREE!

We use an algorithm inspired by the human brain. It uses the stylistic elements of one image to draw the content of another. Get your own artwork in just three steps.

### 1 Upload photo

The first picture defines the scene you would like to have painted.



### 2 Choose style

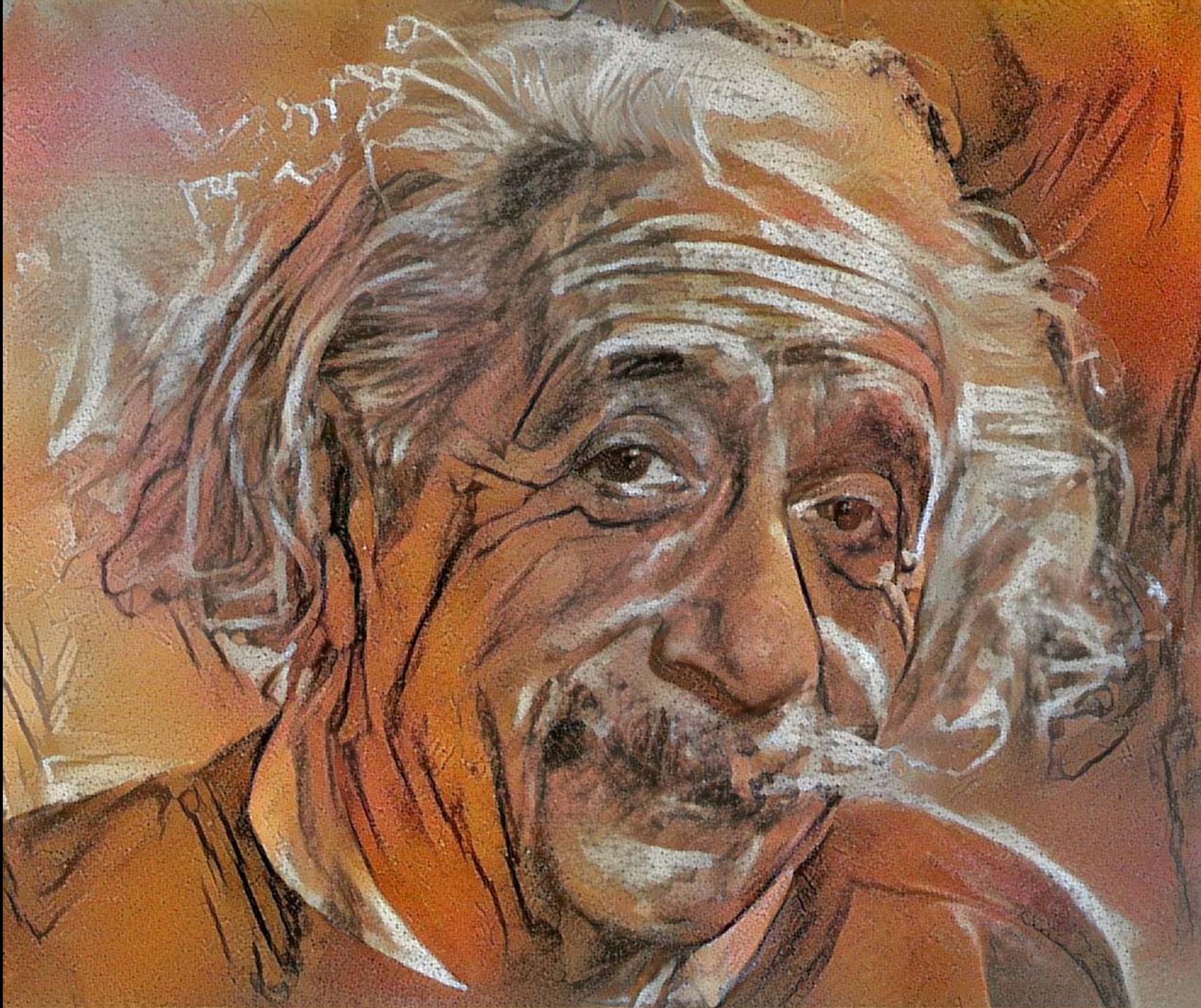
Choose among predefined styles or upload your own style image.



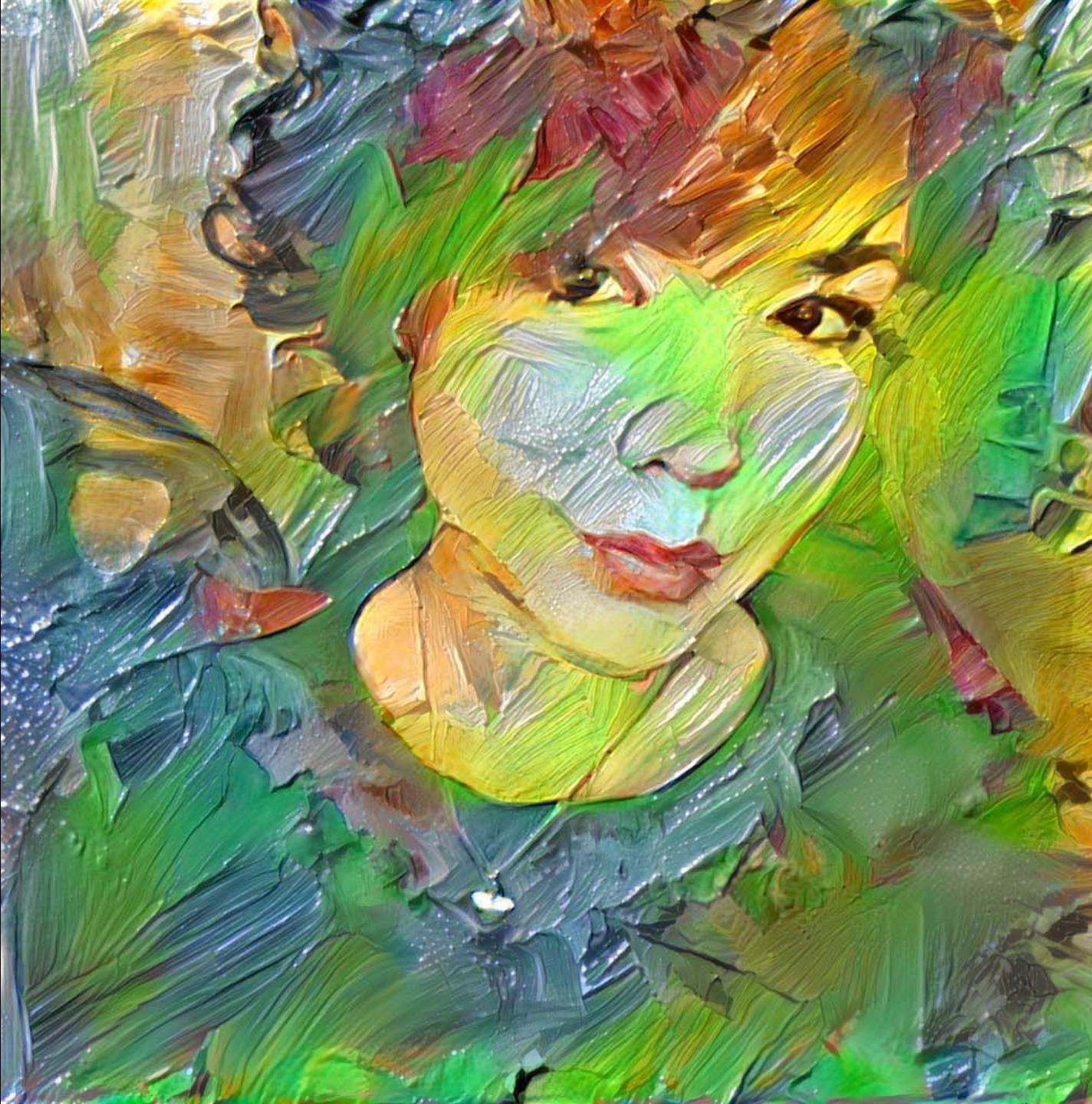
### 3 Submit

Our servers paint the image for you. You get an email when it's done.













# Significance of neural style transfer

Comp. Neuroscience

Texture models  
as models of  
primate vision

Transfer learning  
from deep nets  
trained on ImageNet

Machine Learning

Neural Style  
Transfer  
(Gatys, Ecker, Bethge  
CVPR 2016)

Perceptual loss functions  
for image processing  
(Gatys, Ecker, Bethge CVPR 2017,  
Current Opinion Neurobiol. 2017)

Disentangling  
content and style  
of images

Weaknesses of current  
deep learning systems

(Geirhos et al. NeurIPS 2018, Geirhos et  
al. ICLR 2019)



# Perceptual loss functions

Euclidean distance in pixel space does not cover semantics of images



# Perceptual loss functions

Euclidean distance in pixel space does not cover perceptual differences



Blurred



Original



Shifted by 20 pixels



Compute distances in high-level feature space  
of deep neural network trained on ImageNet

# Neural style transfer implementations

Leon Gatys' PyTorch implementation:

<https://github.com/leongatys/PytorchNeuralStyleTransfer>

Justin Johnson's excellent and very popular Lua/Torch implementation:

*(lots of features, but may be hard to run these days)*

<https://github.com/jcjohnson/neural-style>

Neural Style in official PyTorch tutorials:

[https://pytorch.org/tutorials/advanced/neural\\_style\\_tutorial.html](https://pytorch.org/tutorials/advanced/neural_style_tutorial.html)

# Questions!

-