

Deep Learning for Visual Computing (COMP0169)

# Convolutional Neural Networks

Niloy Mitra

**Tobias Ritschel**

# Introduction

- Idea
- Evolution

# Where do filters come from?

- So far we said we will code those filters ourselves
- Which is the right filter for a task?

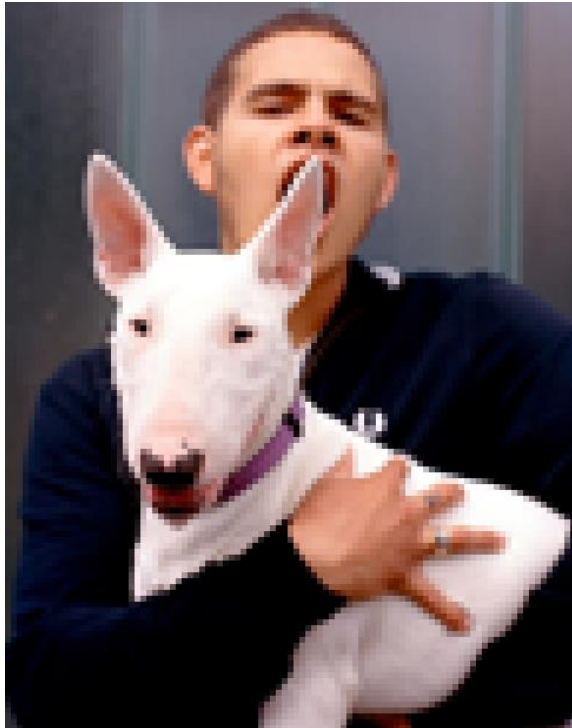
# Quiz

I say “



Input

“?



Output

You say  $(1, 1, 1)/3$

# More quiz

I say “



Input



Output

“?

You say  $(-1,0,1)/2$

# More quiz, computer helping

I say “



Input



“?

Computer says ... a lot of numbers

# Finding a filter by optimization

- Optimize the filter kernel to produce a desired result

$$\arg \min_{\theta} ||f \odot g_{\theta} - h||_2$$

Diagram illustrating the optimization process:

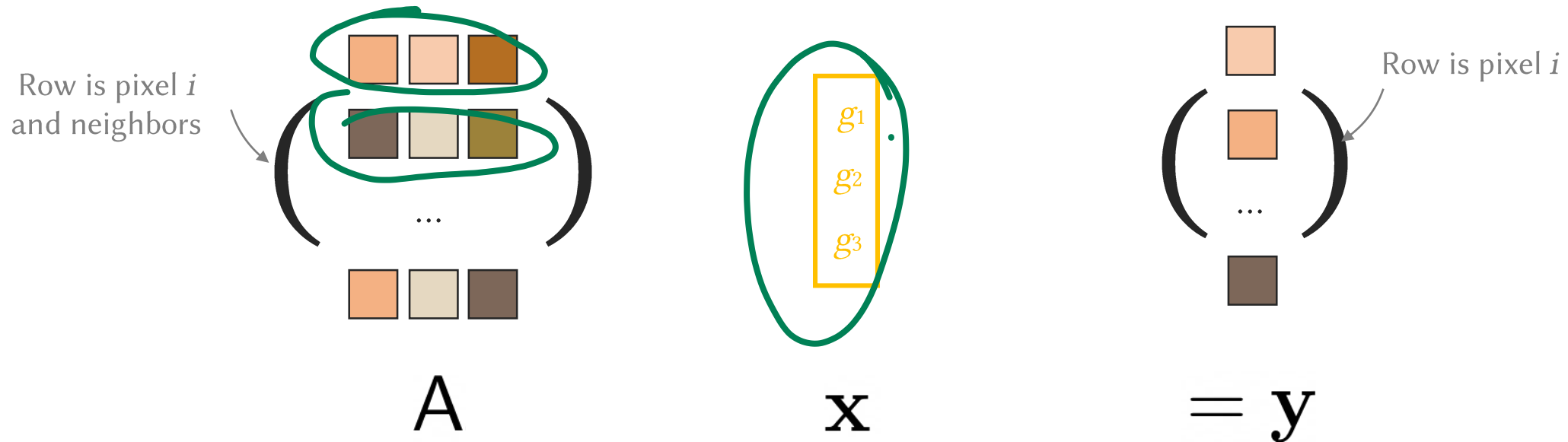
- The input image  $f$  (left) is a photo of a man holding a white dog.
- The filter kernel  $g_{\theta}$  is a  $3 \times 3$  matrix of parameters  $g_{ij}$ .
- The output image  $h$  (right) is the result of applying the kernel to the input image, showing a blurred version of the input.

$$\theta \in \mathbb{R}^{3 \times 3} = \{g_{11}, \dots, g_{33}\}$$

# Linear filtering

- Convolution is a linear operation, so this can be solved quite easily

$$\arg \min_{\theta} \|f \odot g_{\theta} - h\|_2$$





# The cheetah and the zebra

- Could we use this to classify images?



cheetah



zebra

# Probably relevant



cheetah



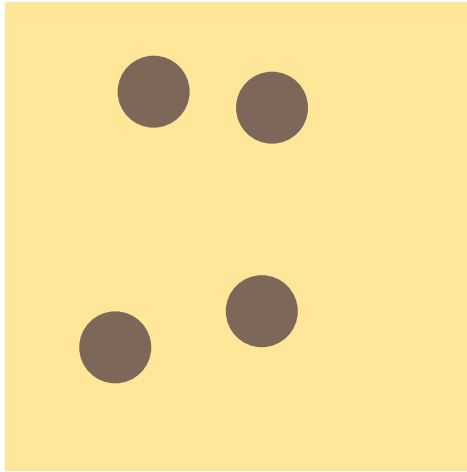
zebra



# The cheetah and the zebra



# Points, edges, non-points, non-edges

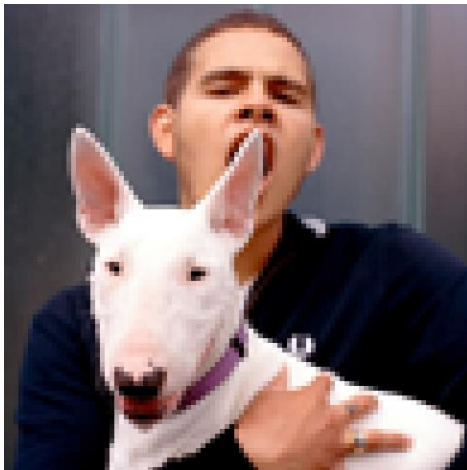


$$\odot f_{\theta_1} = 1$$

是 chaitin

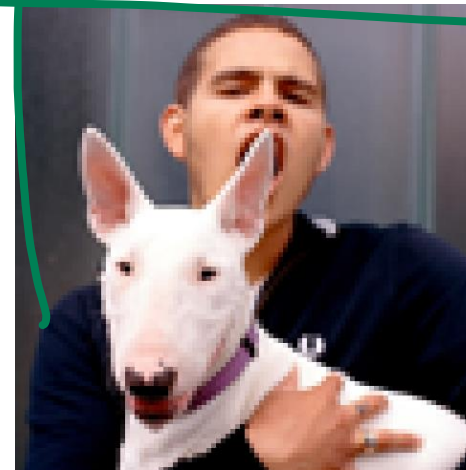


$$\odot f_{\theta_2} = 1$$



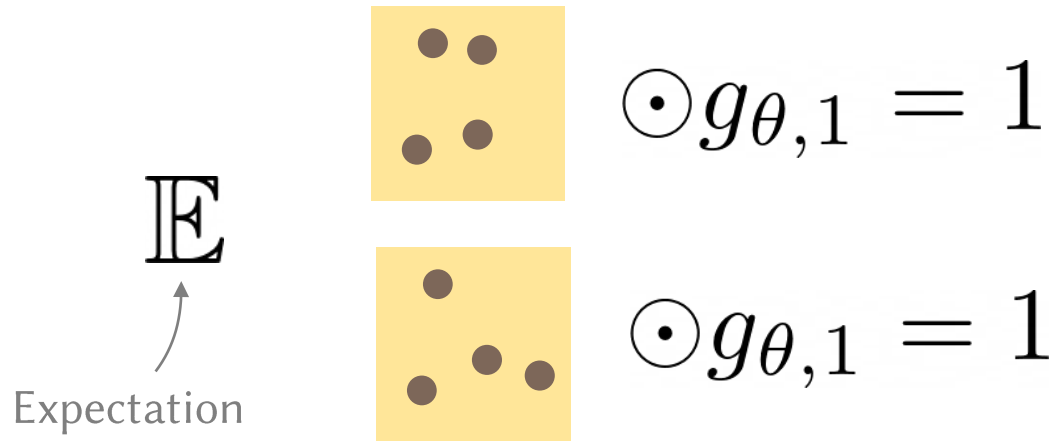
$$\odot f_{\theta_1} = -1$$

不是 chaitin



$$\odot f_{\theta_2} = -1$$

# Example solution



-1	-1	-1
-1	8	-1
-1	-1	-1



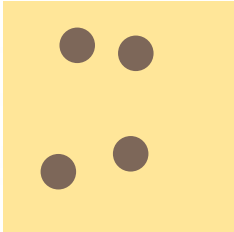
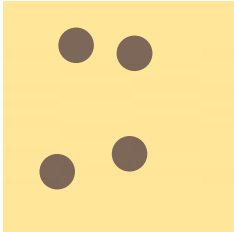
-1	2	-1
-1	2	-1
-1	2	-1

# Non-linearity

- As we saw in the classification and NN lecture, non-linearity is important

$$\mathbb{E}[\phi(\text{image} \odot g_{\theta})] \neq \phi(\mathbb{E}[\text{image} \odot g_{\theta}])$$

Non-linearity



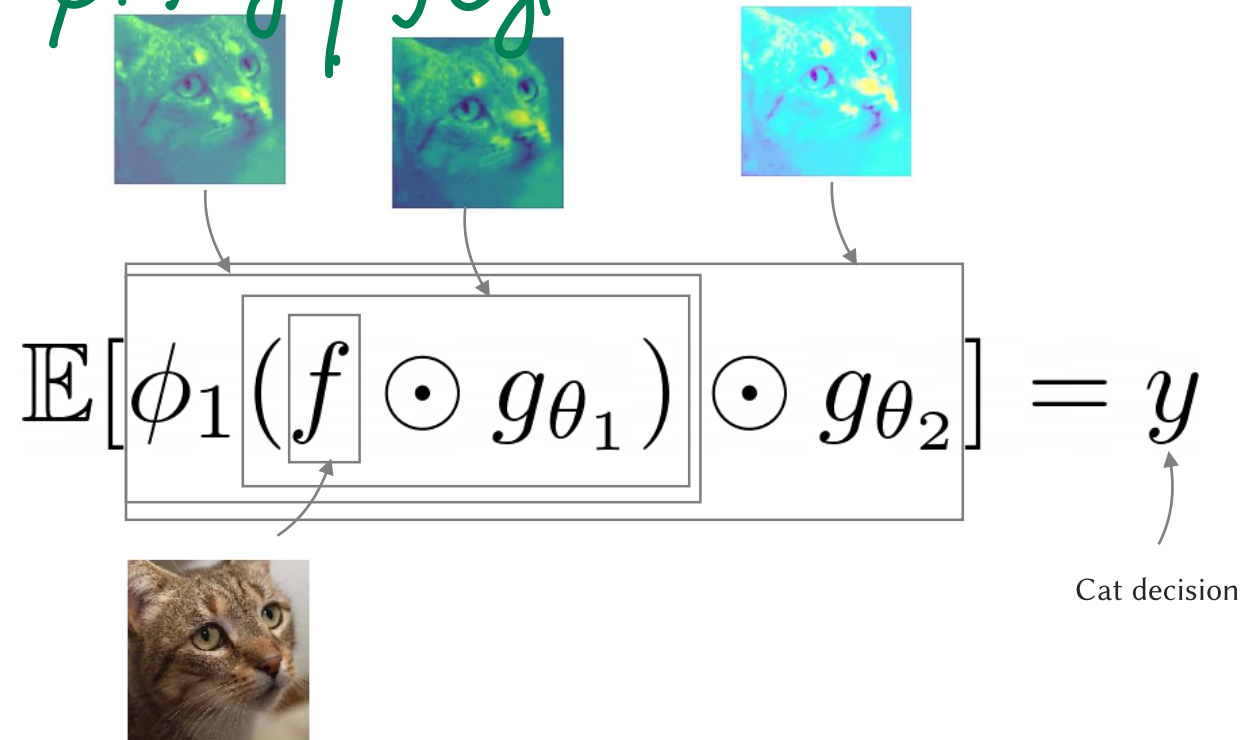
# Non-linearity

- And after that no-linearity, we could again do a convolution
- We call this **layers**

$$\mathbb{E}[\phi_1(\text{cat} \odot g_{\theta_1}) \odot g_{\theta_2}] \neq \mathbb{E}[\phi_1(\text{cat} \odot (g_{\theta_1} \odot g_{\theta_2}))]$$

# Depth + non-linearity $\Phi$

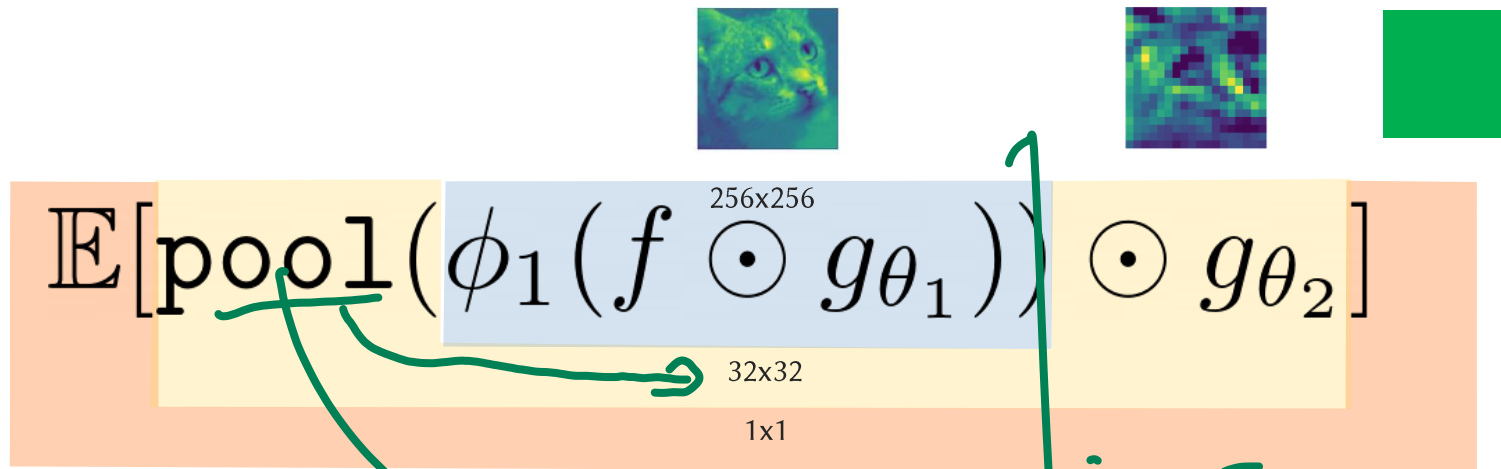
- What this looks like  $\phi, f \odot g_{\theta_1}, f \odot g_{\theta_2}$





# Multi-resolution

- We recall that finding edges and any feature is best done on a pyramid
- CNN so far one resolution
- Solution: insert **pooling** reducing resolution after each convolution



提取这E域的主要特征

# Fully-connected layer

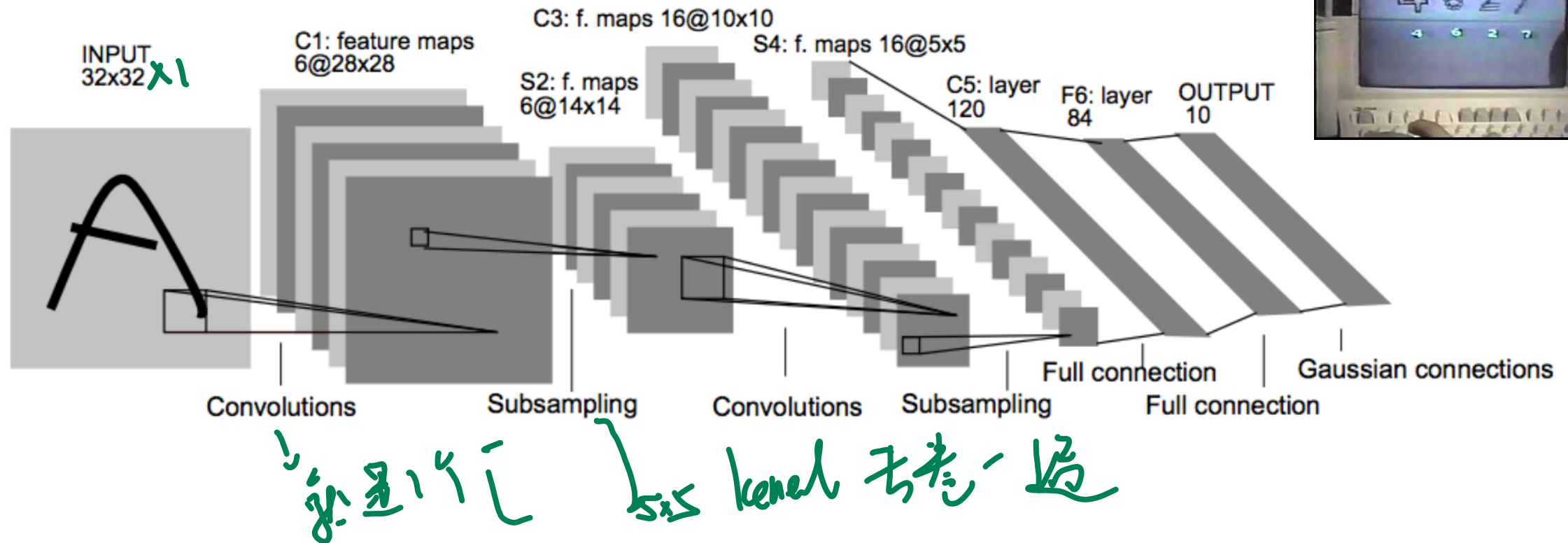
- Instead of expectation, take all pixels and feed them to a classic NN

Map c-times-n vector to number

$$\psi(\text{stack}(\text{pool}(\phi_1(f \odot g_{\theta_1})) \odot g_{\theta_2})) = y$$

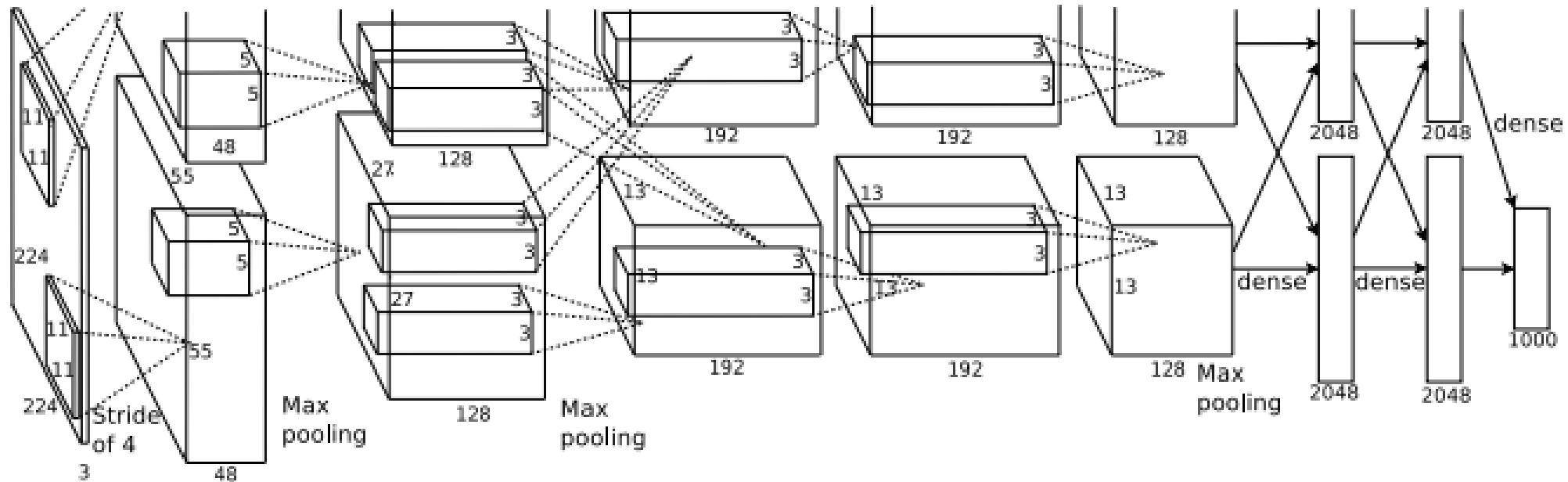
Make c-times-n-vector from all n pixels and all c chans

# Handwriting recognition: LeNet



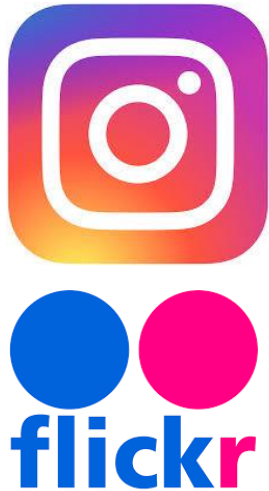
LeCun et al, 1990: Handwritten digit recognition with a back-propagation network

# Image classification: AlexNet



Krizhevsky et al, 2012: *ImageNet Classification with Deep Convolutional Neural Networks*

# What happened between 1992 and 2012?



**Image data**

60k vs  
60M



**GPUs**

Pentium II 0.2 GFLOPS  
Nvidia GTX 680 2600 GFLOPS

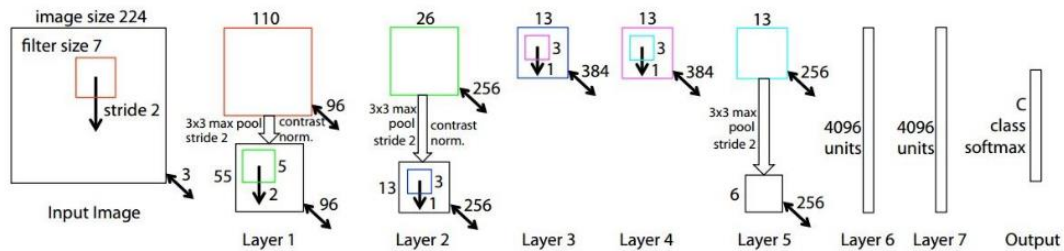


**Scripting**

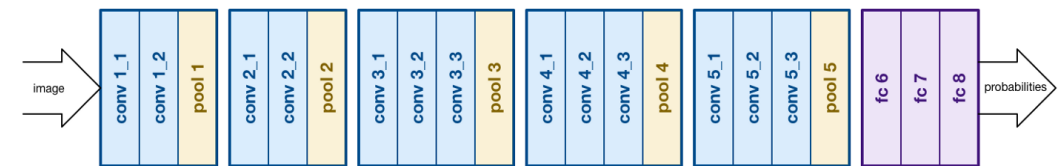
C++ vs  
CUDA/Python

# More architectures

- Not always relevant or clear what happened



GoogLeNet/Inception (2014)

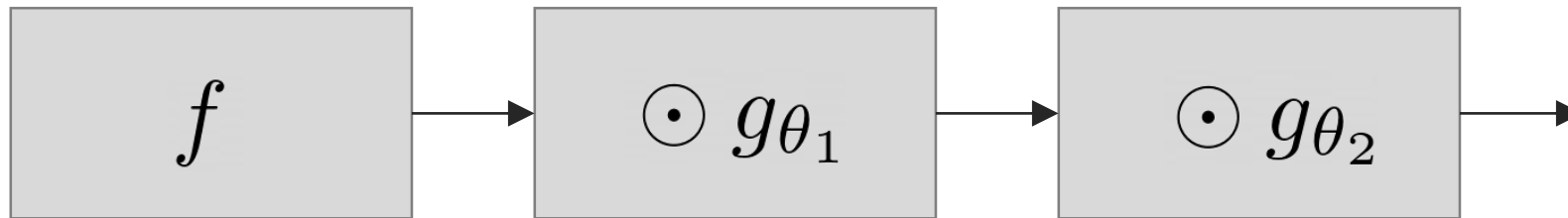


VGG Net (2014)

# ResNet, 2015

- Usual sequence makes propagating gradients hard

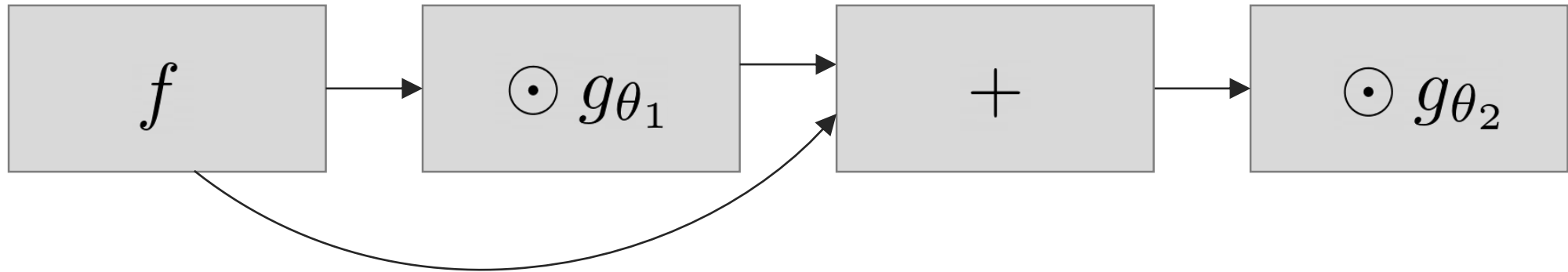
$$\psi(\phi_1(f \odot g_{\theta_1}) \odot g_{\theta_2})$$



# ResNet, 2015

- Usual sequence makes propagating gradients hard

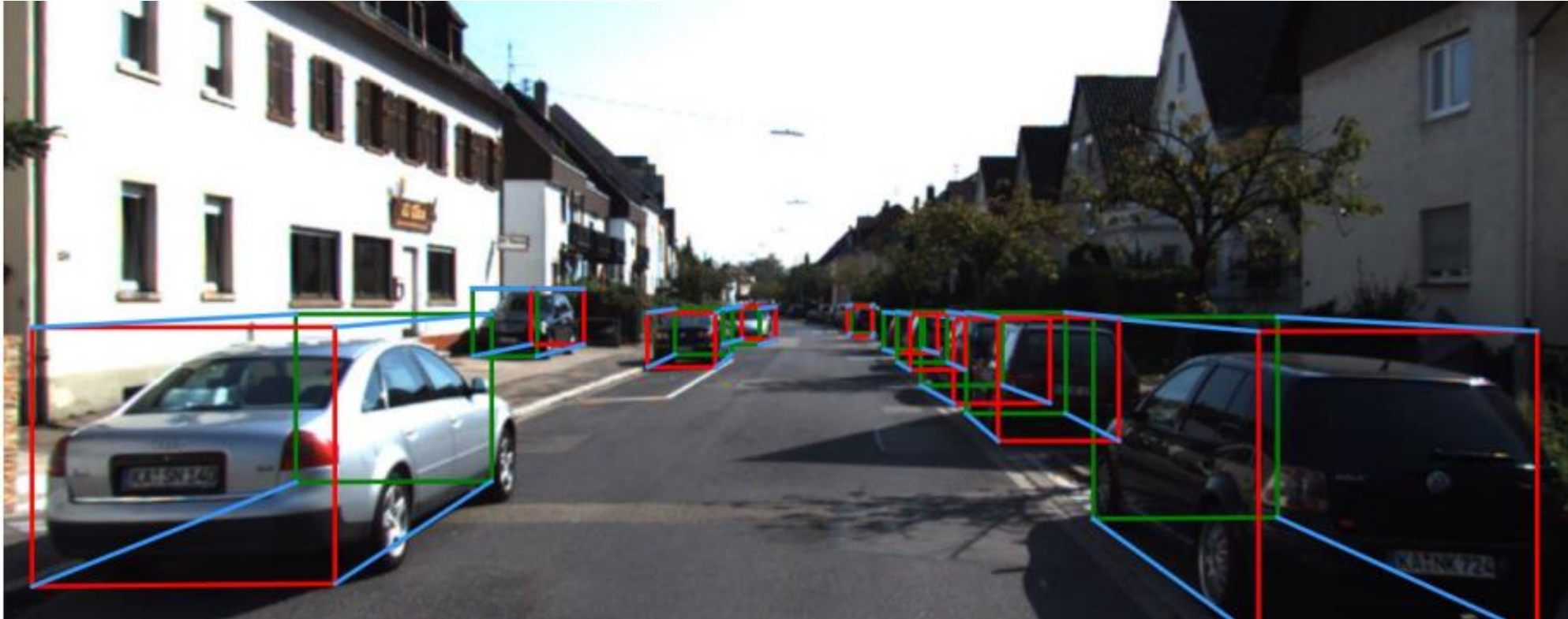
$$\psi((\phi_1(f \odot g_{\theta_1}) + f) \odot g_{\theta_2})$$





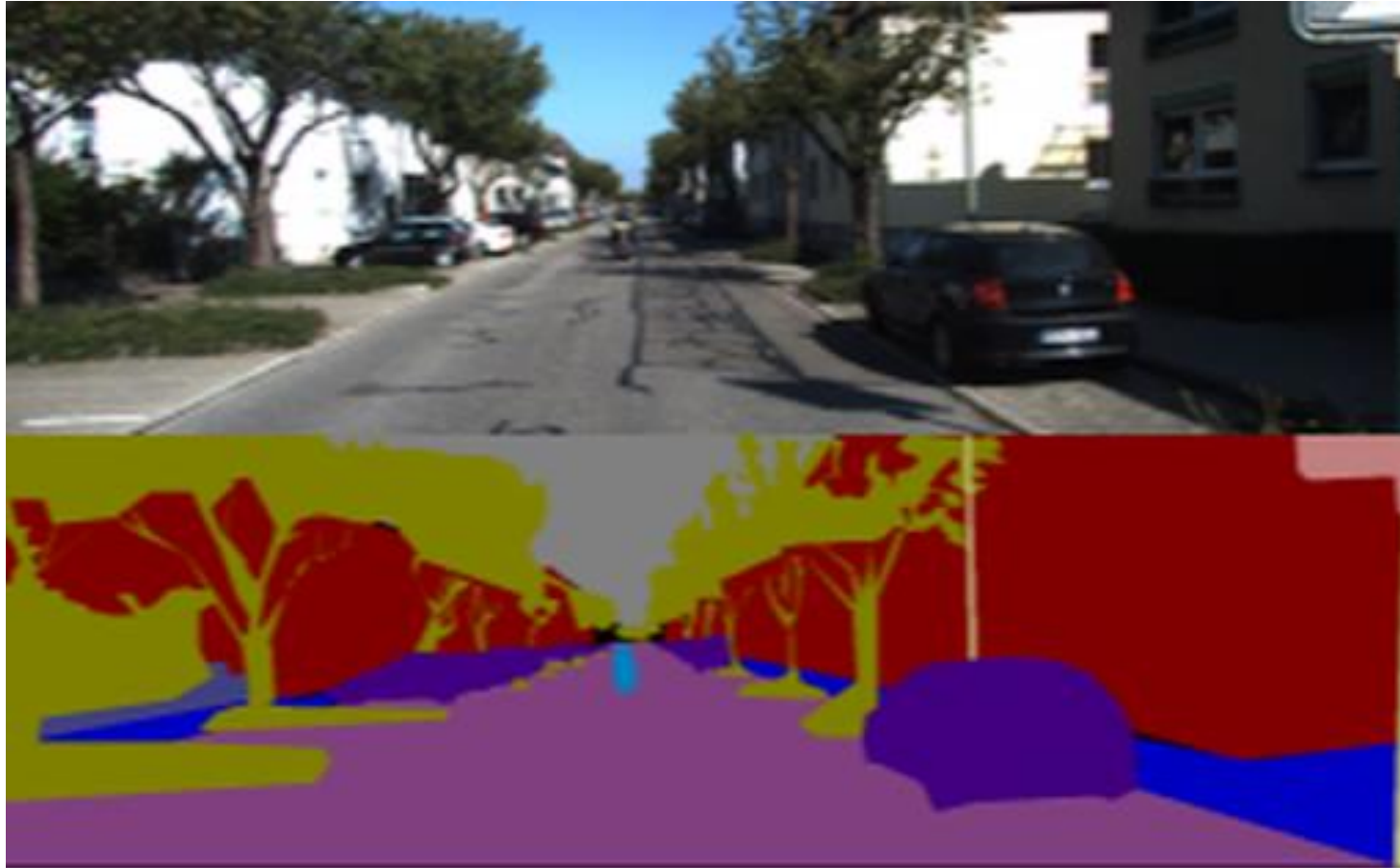


# 3D Object detection



# Image-to-image: Segmentation

RGB



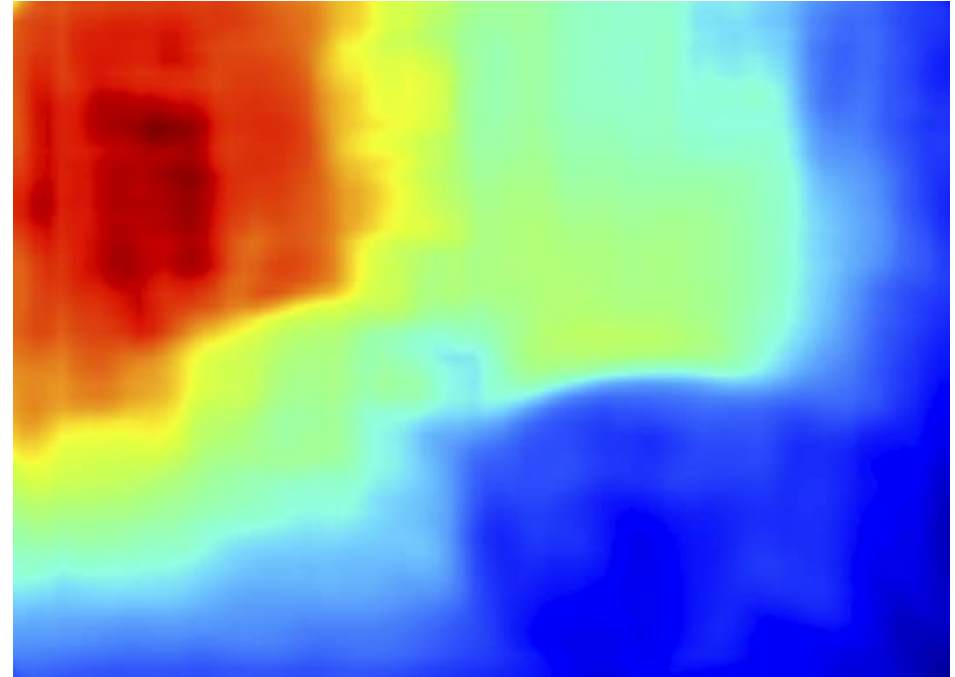
Semantic  
segments



# Image-to-image



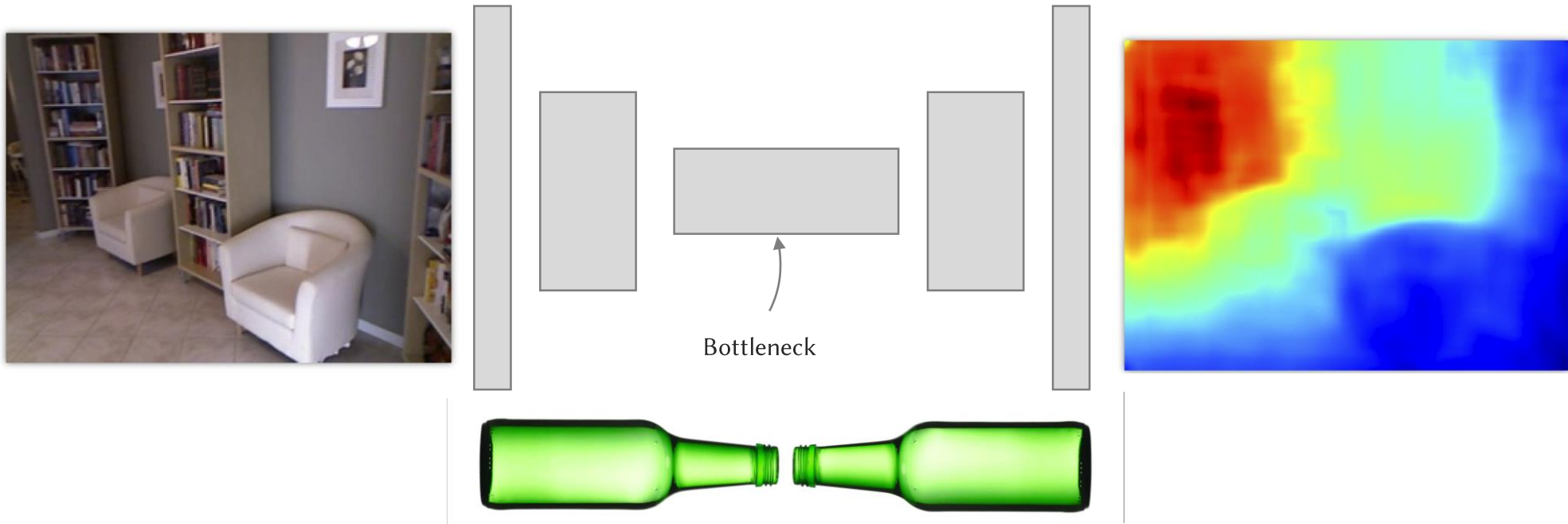
RGB



Depth

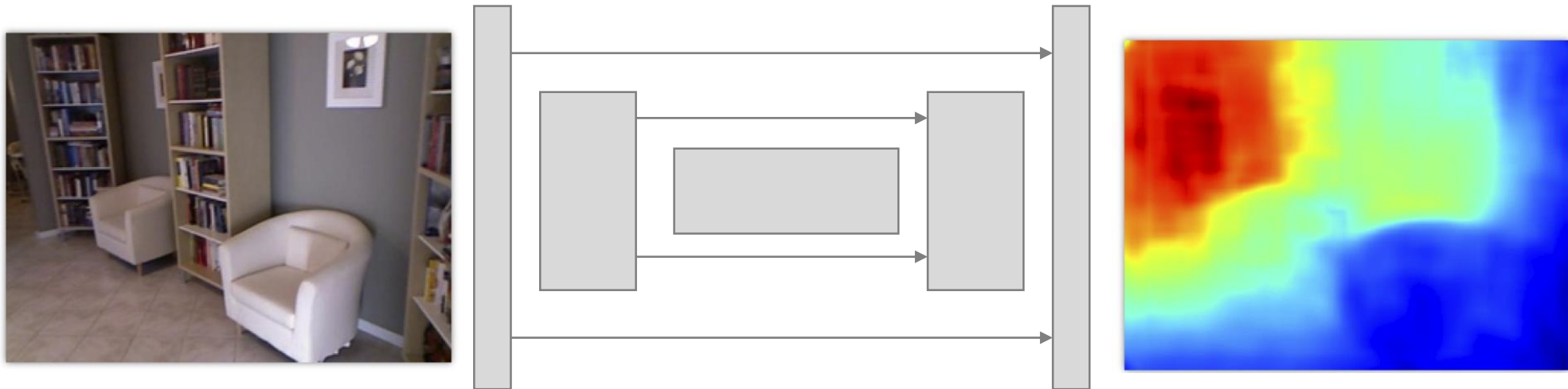
# Image-to-image

- Image in, number out: Encoder-decoder
- How about image out?

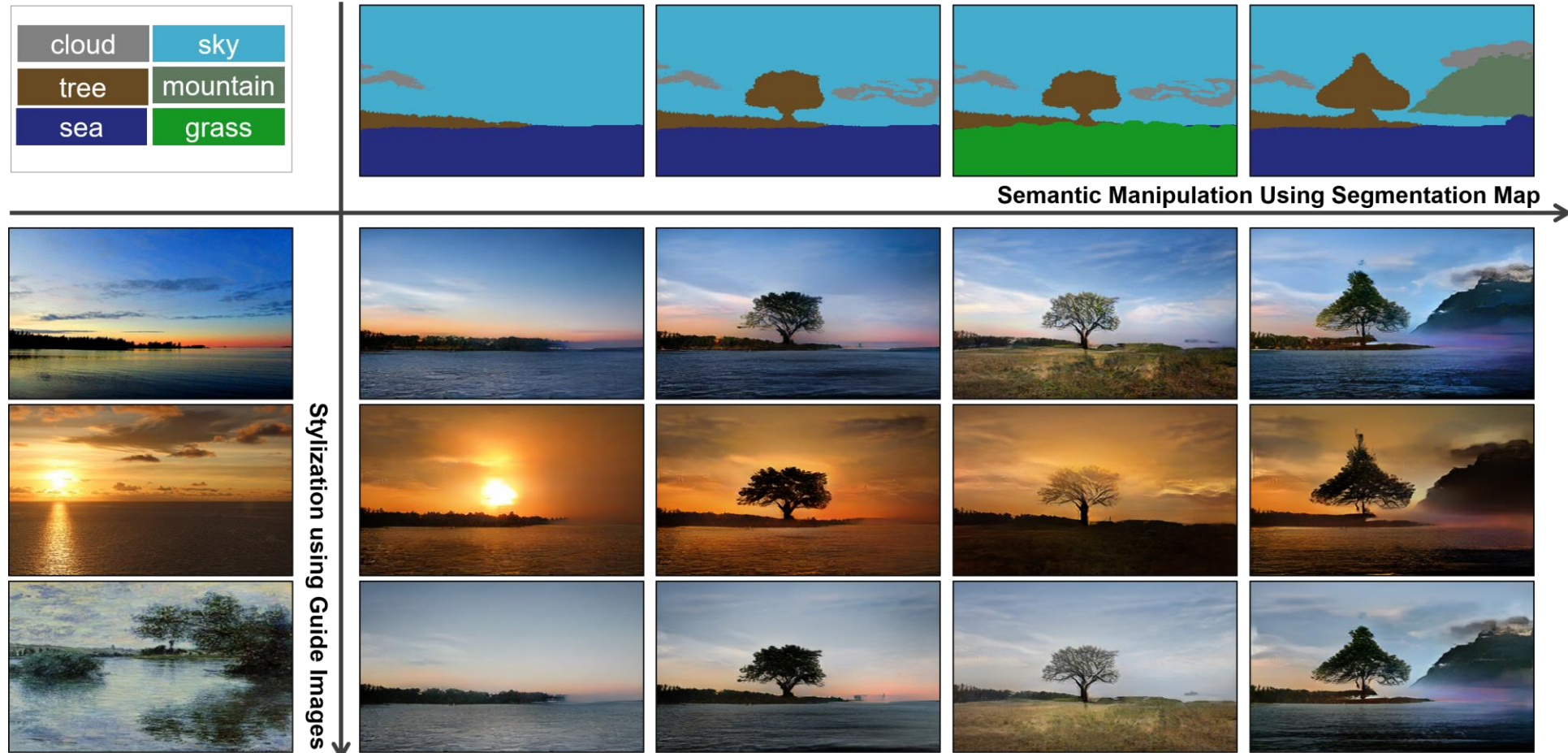


# Image-to-image

- No chance to get back details
- Idea: **skip connections**, allow decoder to access encoder state (Unet)

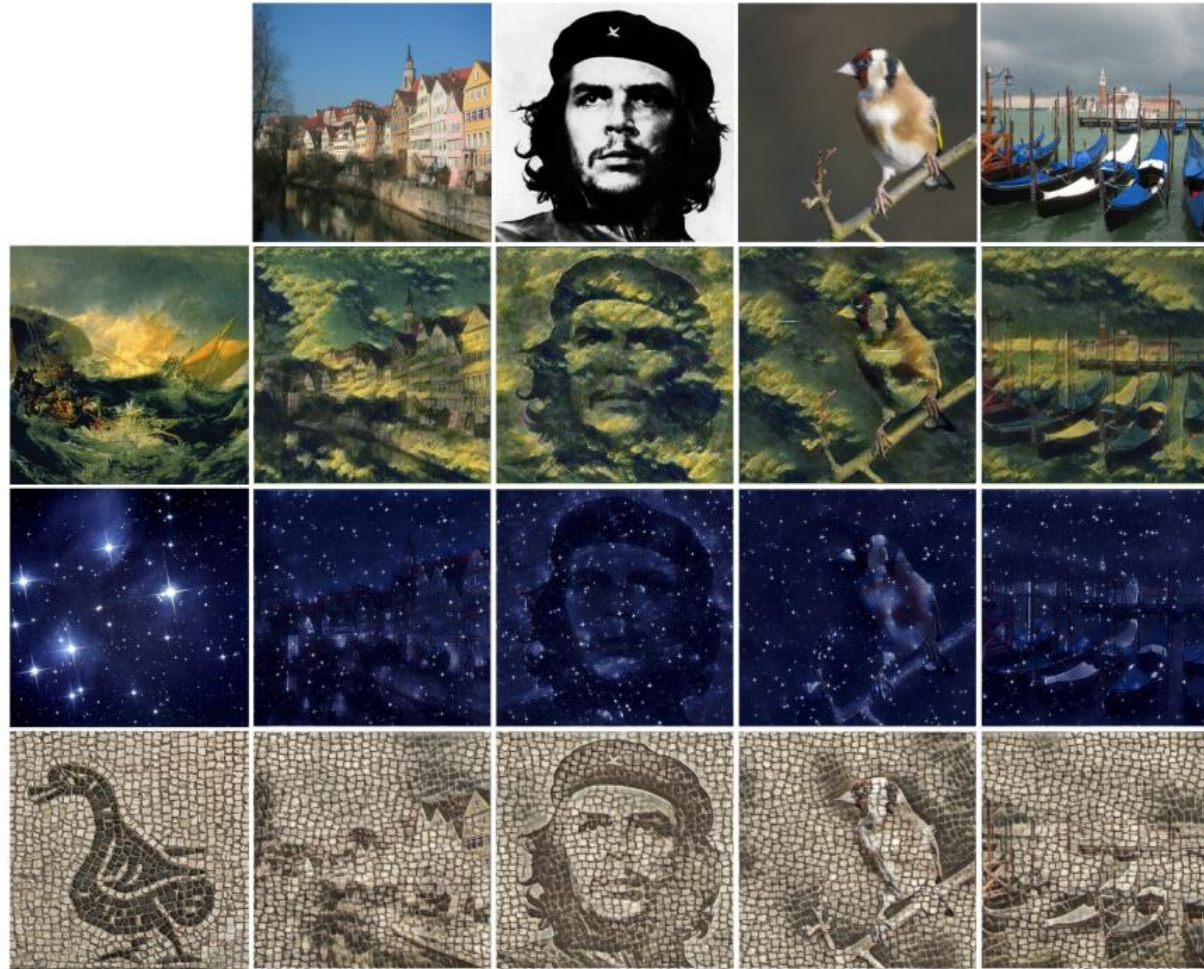


# Semantic image synthesis





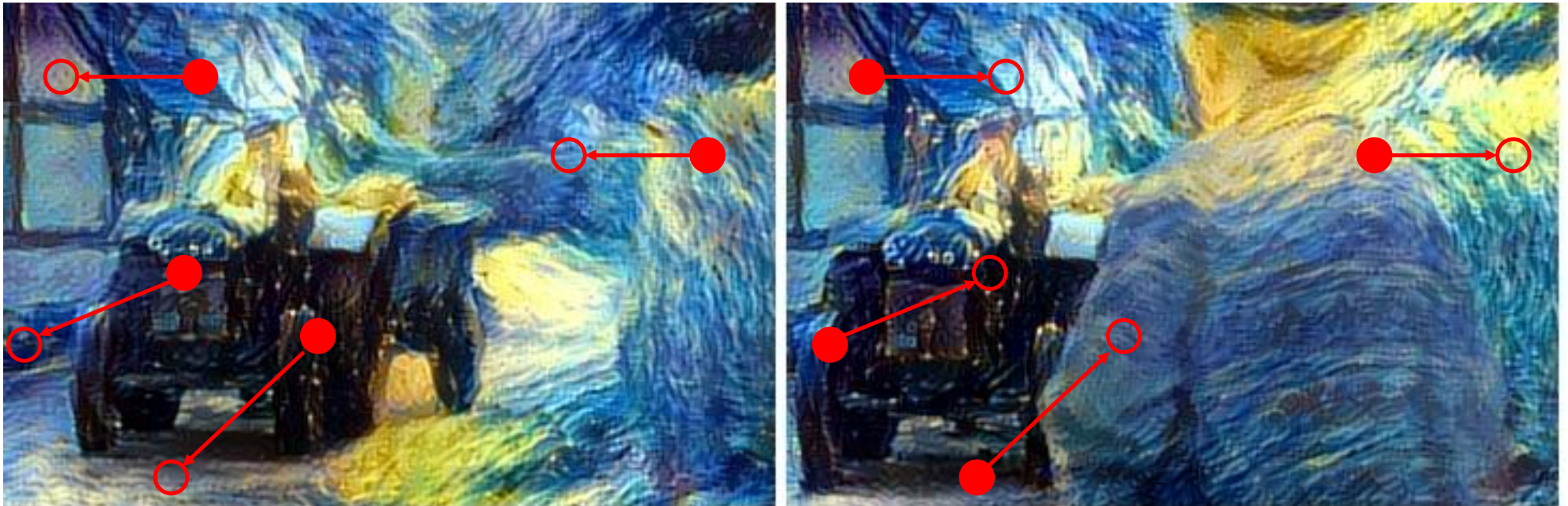
# Style transfer





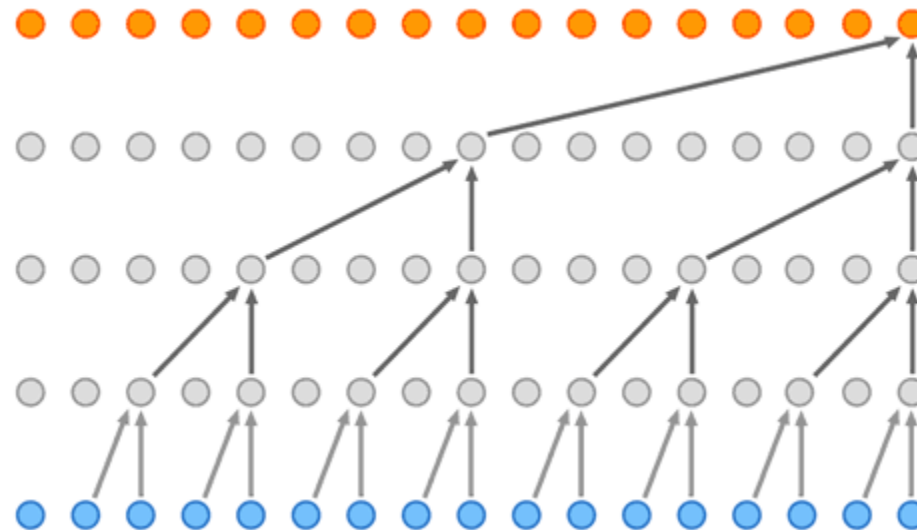
# Encoder-decoder for video

- Insert flow-compensated difference into loss



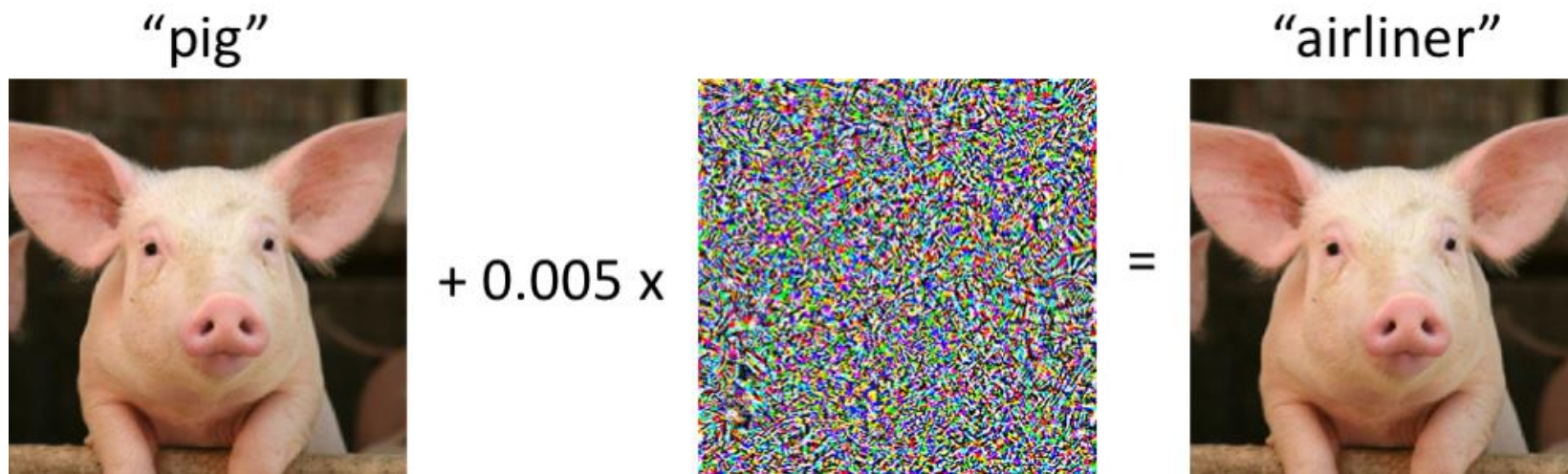
# Encoder-Decoder for music

- Wavenet
- Like U-Net, just 1D



# Failure

- Adversarial examples
- Small perturbation put detection off
- Found by optimizing for it





# Dreaming



# Conclusion

- Don't write filters yourself
- Let the computer find them
- Stack them with non-linearities
- Add a few tricks (skip, res)
- And it does many things