

Deep Learning for Visual Computing (COMP0169)

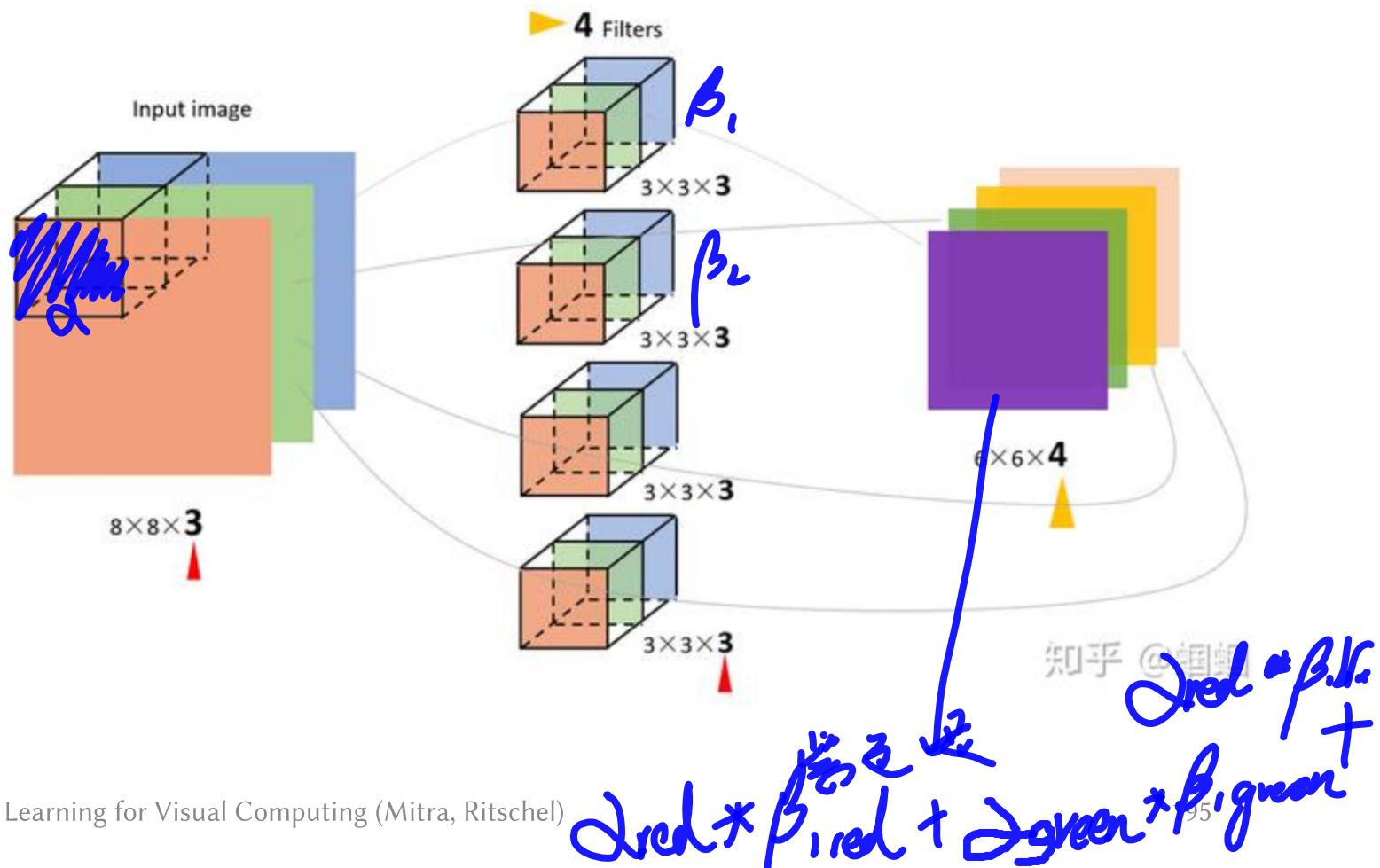
Convolutional Neural Networks

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

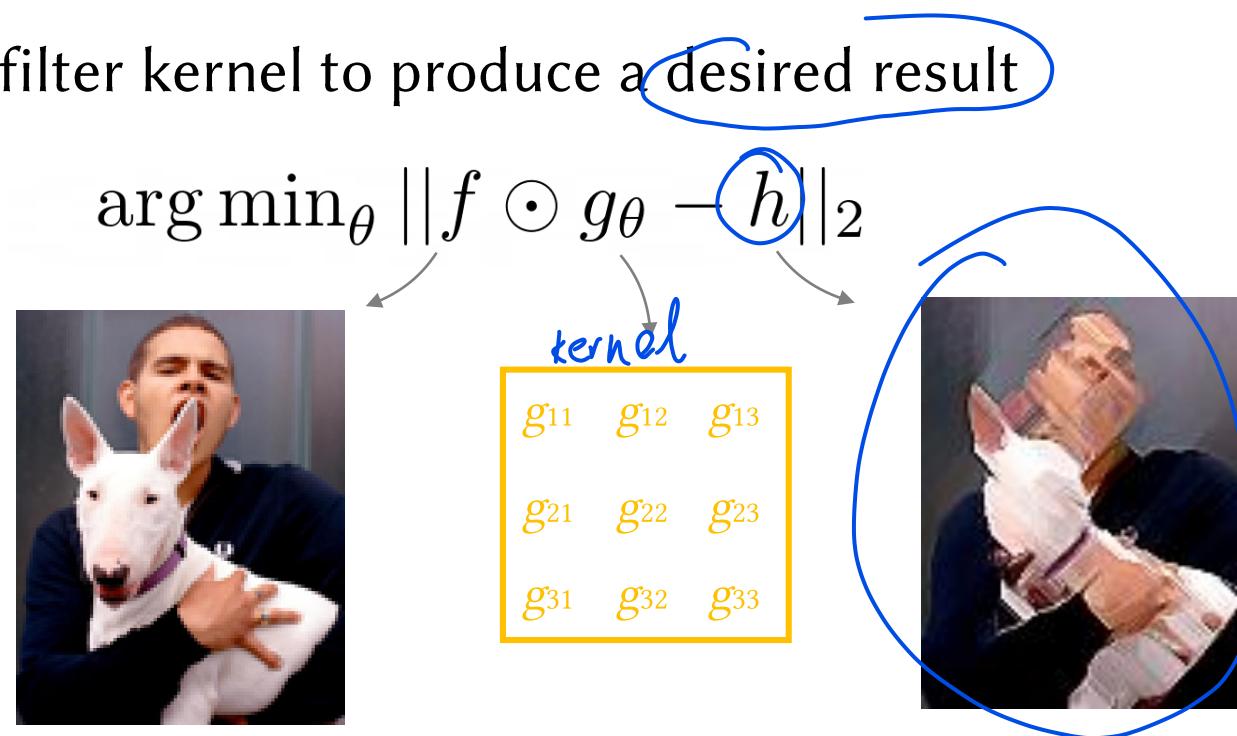


“?

Output

Finding a filter by optimization

- Optimize the filter kernel to produce a desired result

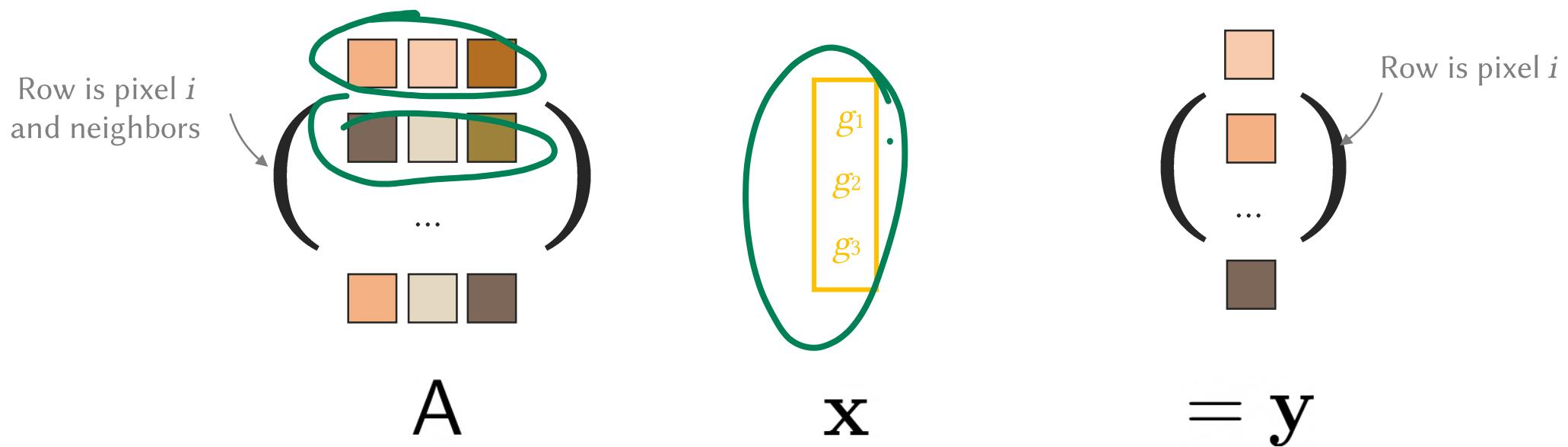


$$\theta \in \mathbb{R}^{3 \times 3} = \{g_{11}, \dots, g_{33}\}$$

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?

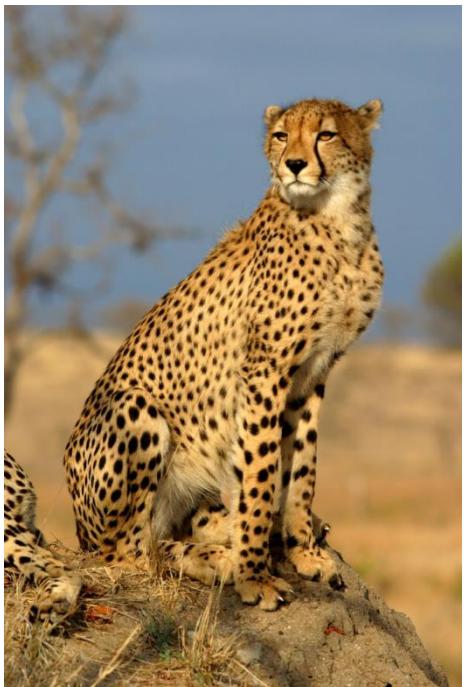


cheetah



zebra

Probably relevant

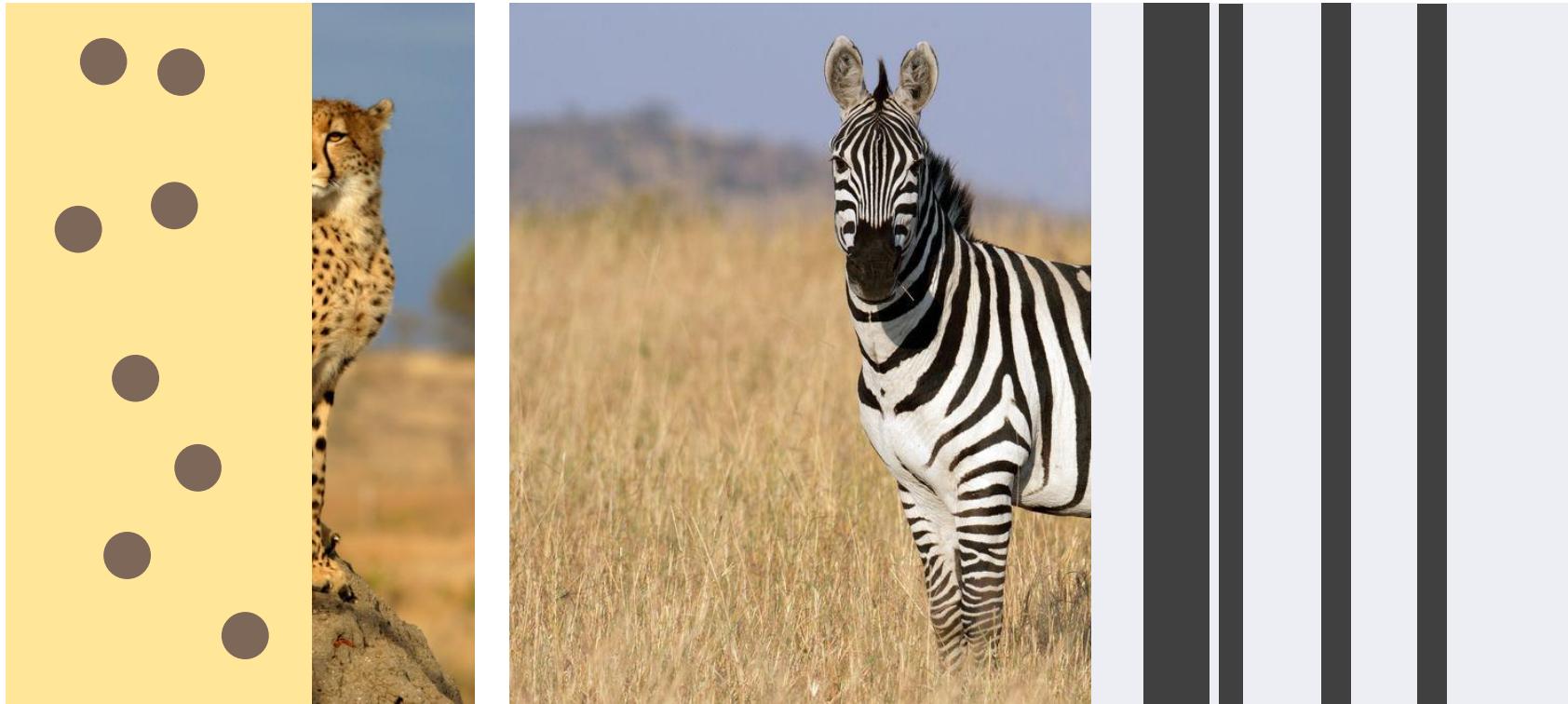


cheetah

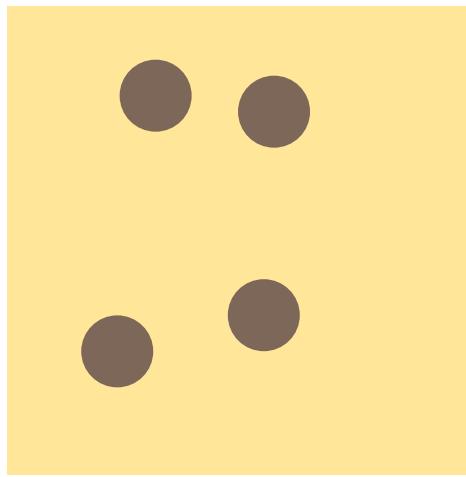


zebra

The cheetah and the zebra



Points, edges, non-points, non-edges



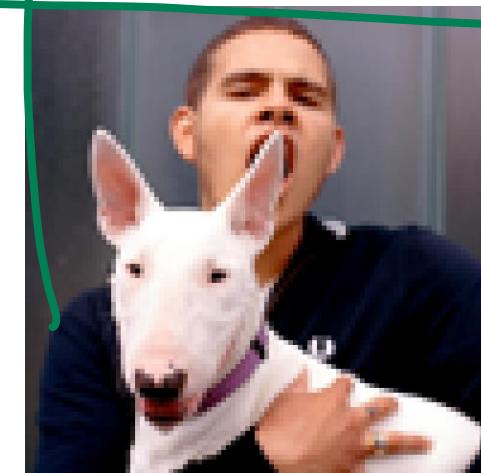
$$\odot \mathcal{G}_{\theta_1} = \Gamma$$

23 cheetah

$$\odot \mathcal{G}_{\theta_1} = -1$$



$$\odot \mathcal{G}_{\theta_2} = 1$$



$$\odot \mathcal{G}_{\theta_2} = -1$$

Example solution

\mathbb{E}
Expectation

$$\odot g_{\theta,1} = 1$$

$$\odot g_{\theta,2} = 1$$

-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

$$\mathbb{E}[\phi(\begin{array}{cc} \bullet & \bullet \\ \vdots & \vdots \\ \bullet & \bullet \end{array} \odot g_\theta)] \neq \phi(\mathbb{E}[\begin{array}{cc} \bullet & \bullet \\ \vdots & \vdots \\ \bullet & \bullet \end{array} \odot g_\theta])$$

Non-linearity

Non-linearity

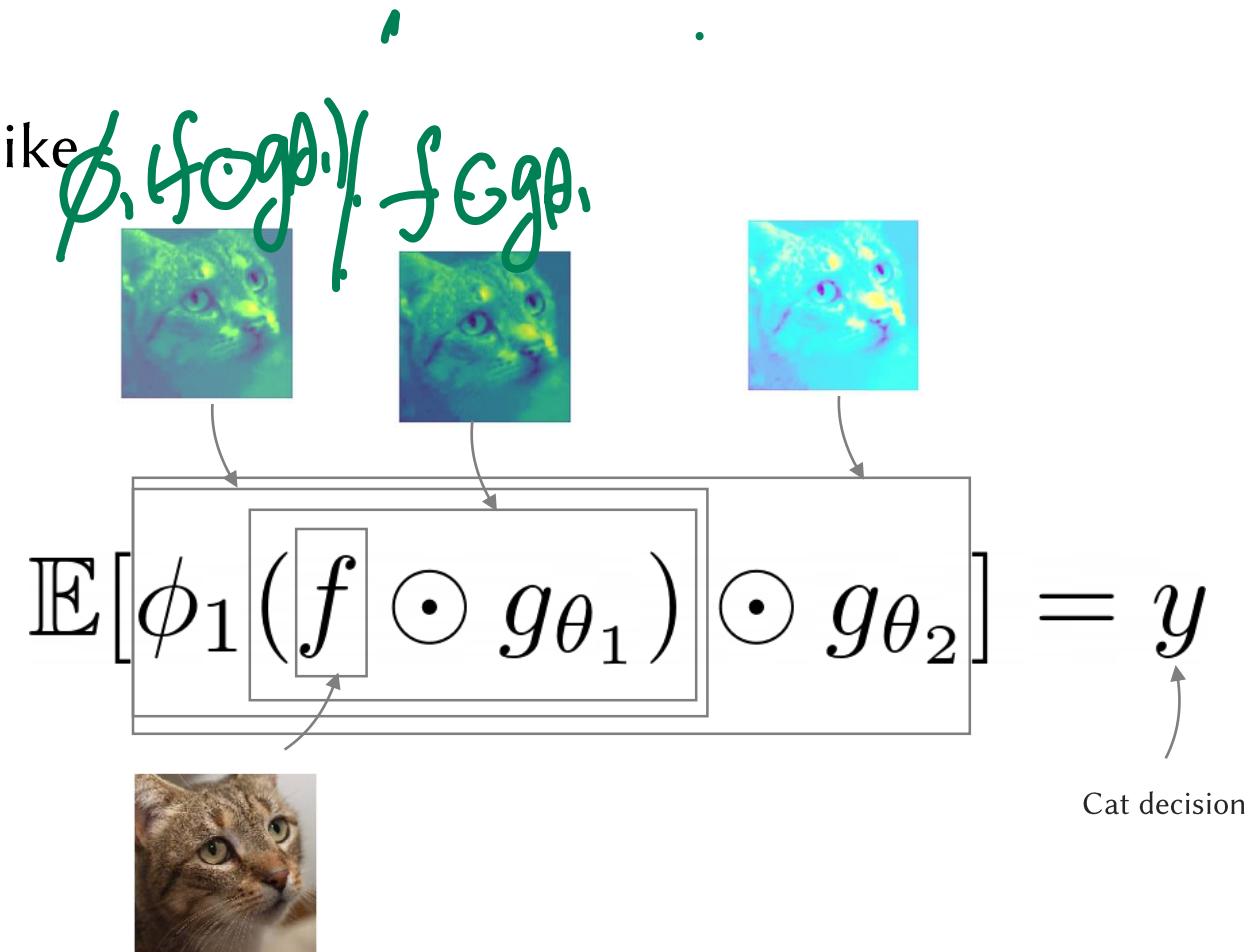
- And after that no-linearity, we could again do a convolution
- We call this **layers**



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

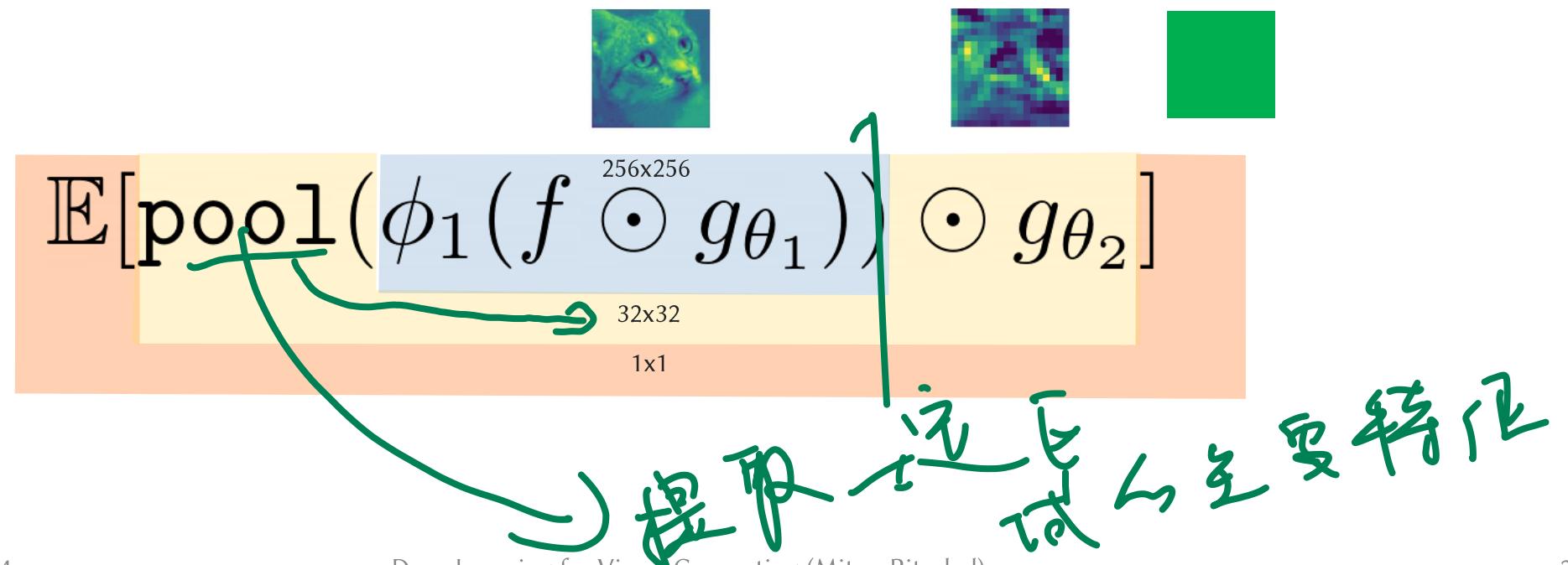
Depth + non-linearity ϕ

- What this looks like



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



Fully-connected layer

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

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

Map c-times-n vector to number

32x32

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

Handwriting recognition: LeNet

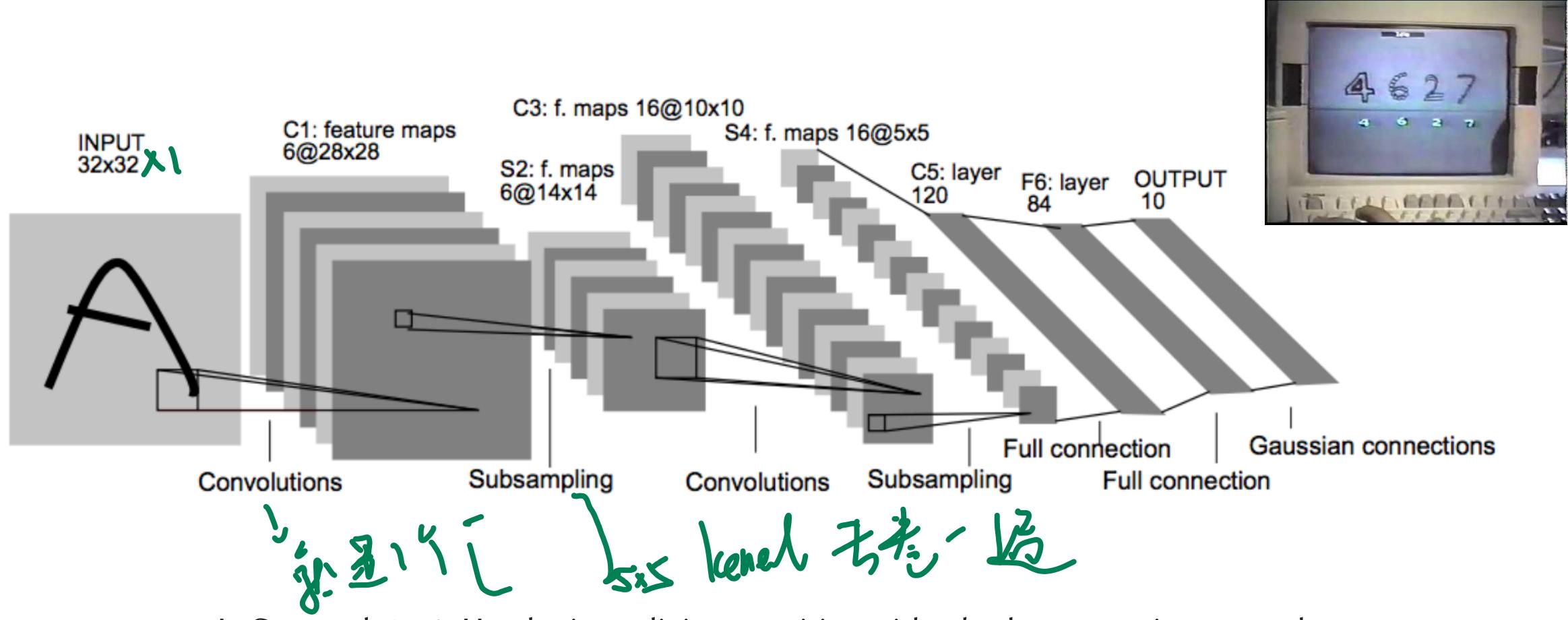
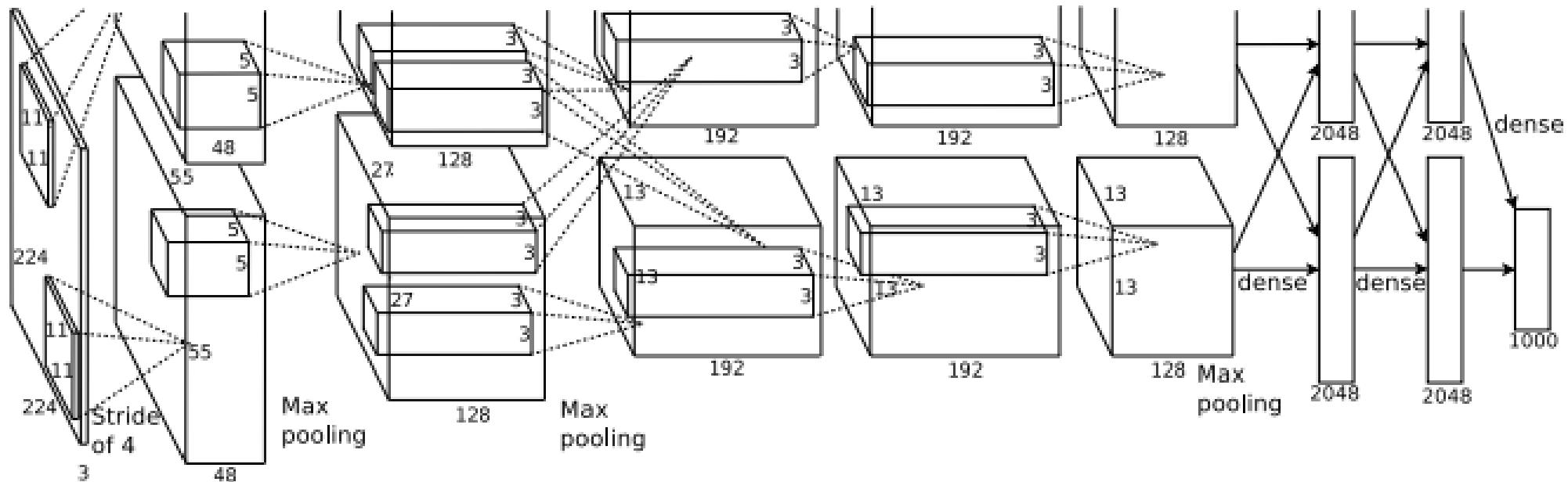


Image classification: AlexNet



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

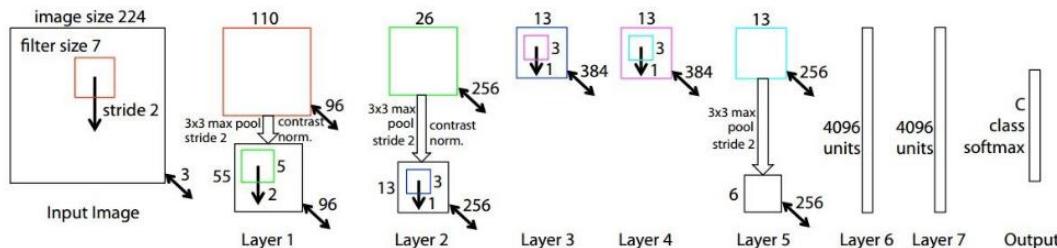


Scripting

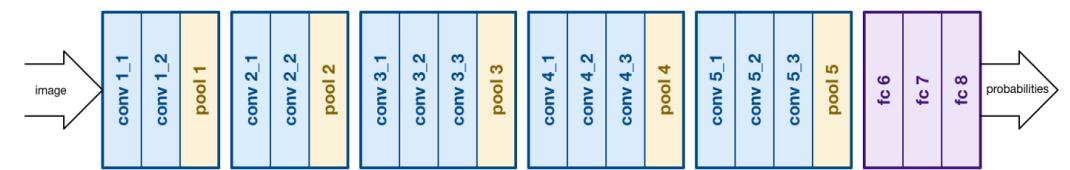
C++ vs
CUDA/Python

More architectures

- Not always relevant or clear what happened



GoogLeNet/Inception (2014)

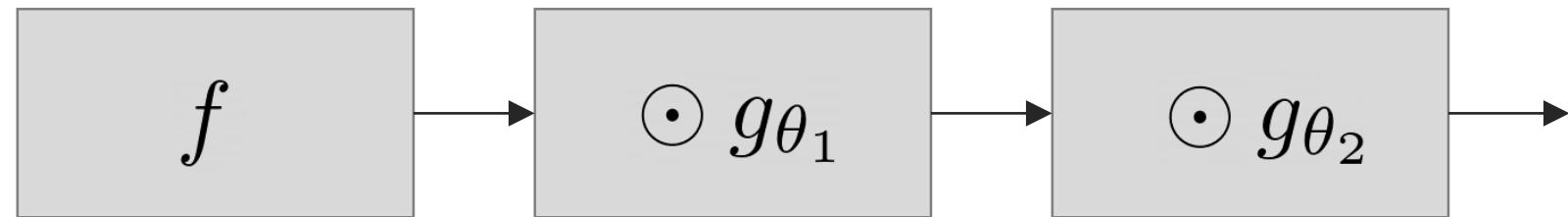


VGG Net (2014)

ResNet, 2015

- Usual sequence makes propagating gradients hard

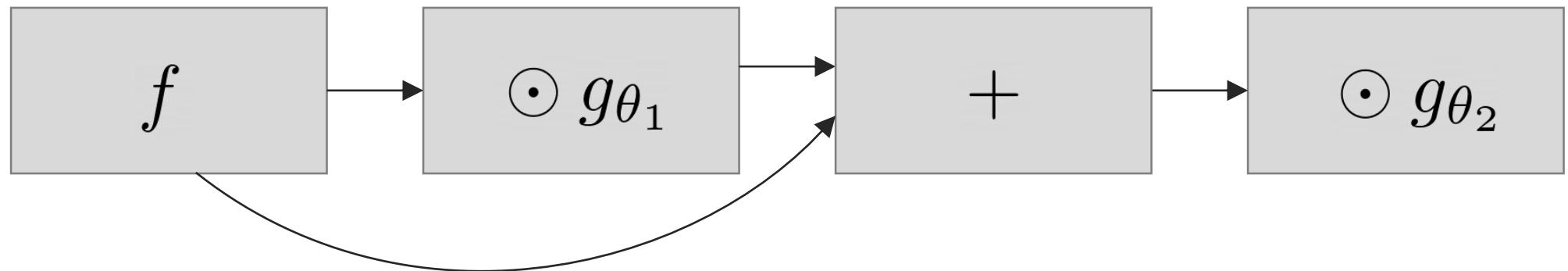
$$\psi(\phi_1(f \odot g_{\theta_1}) \odot g_{\theta_2})$$



ResNet, 2015

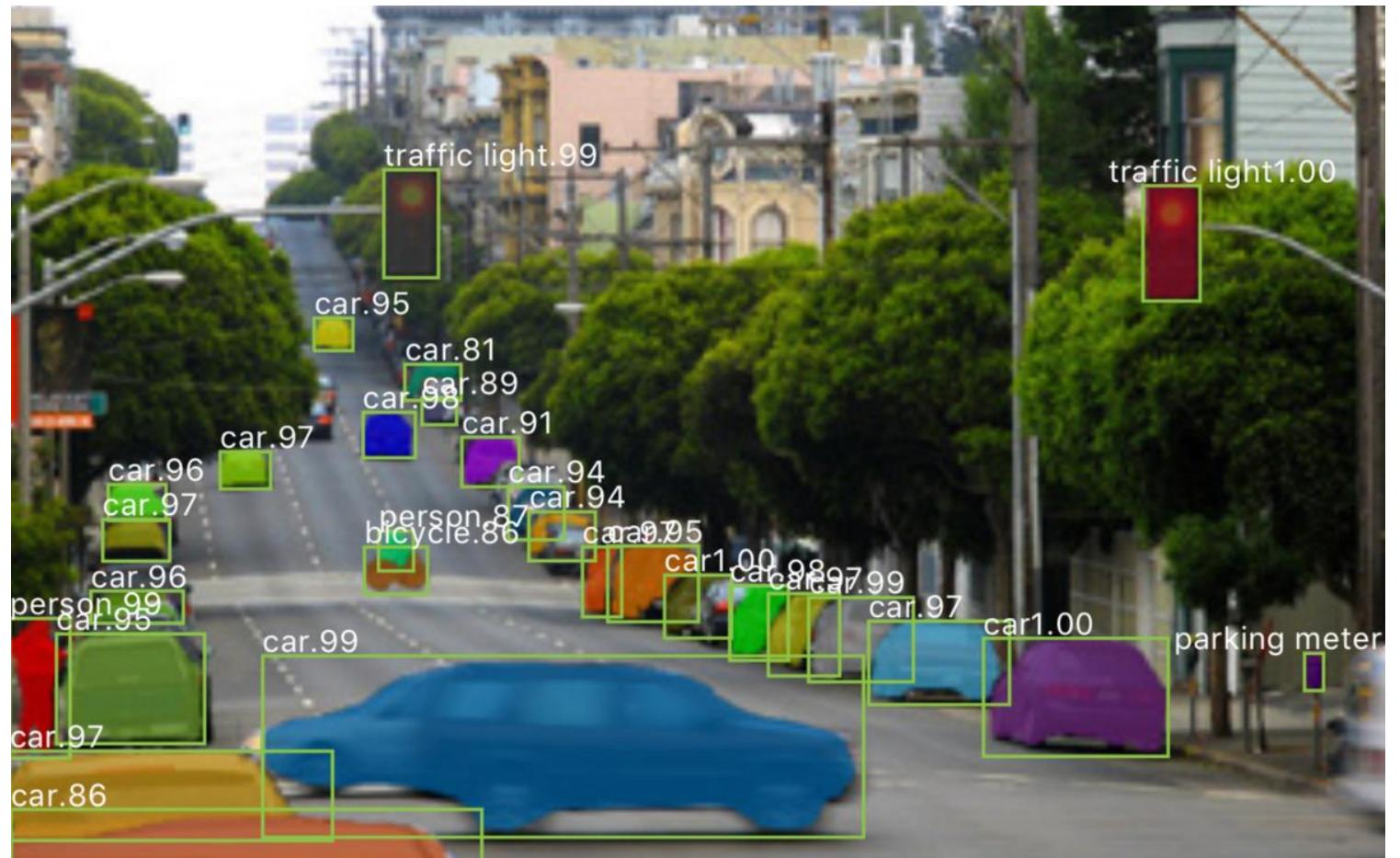
- Usual sequence makes propagating gradients hard

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



Object detection task

- Map image to bounding boxes
- Challenge: Number **unknown**



3D Object detection

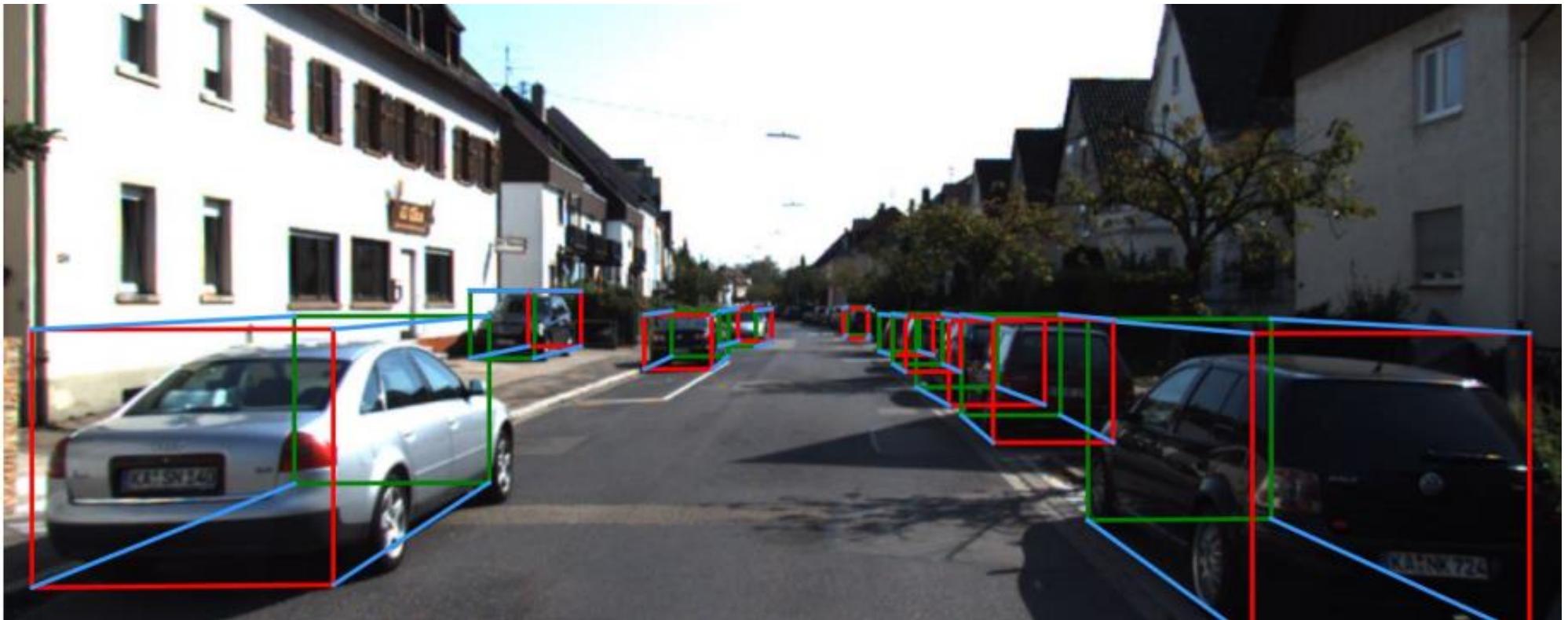


Image-to-image: Segmentation

RGB

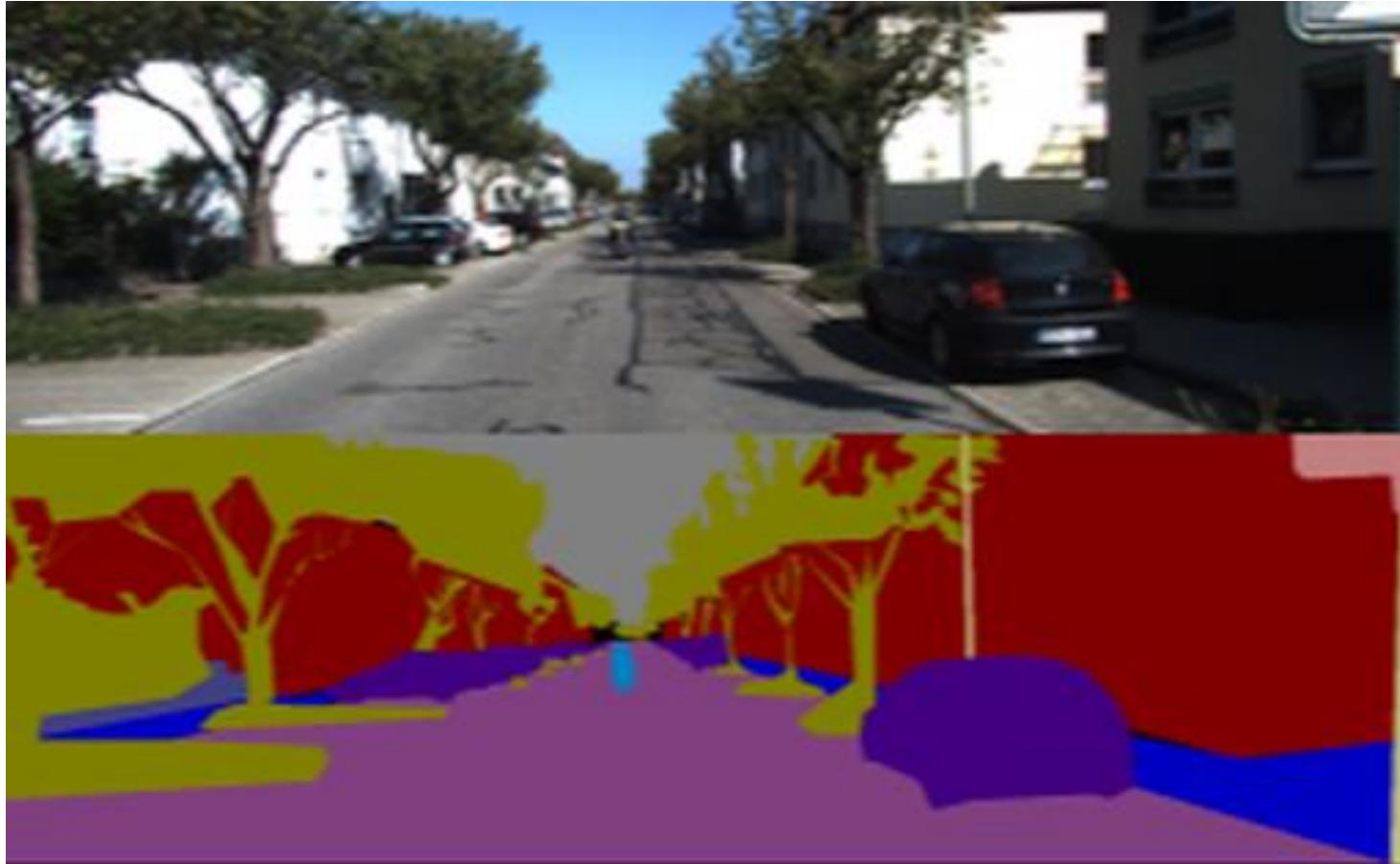
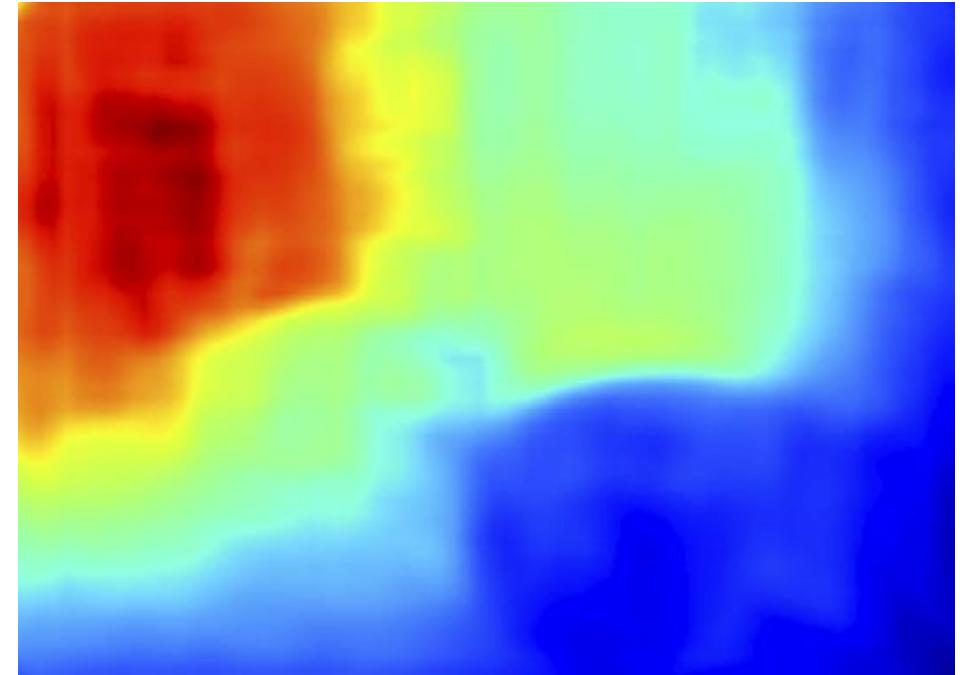
Semantic
segments

Image-to-image



RGB



Depth

Image-to-image

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

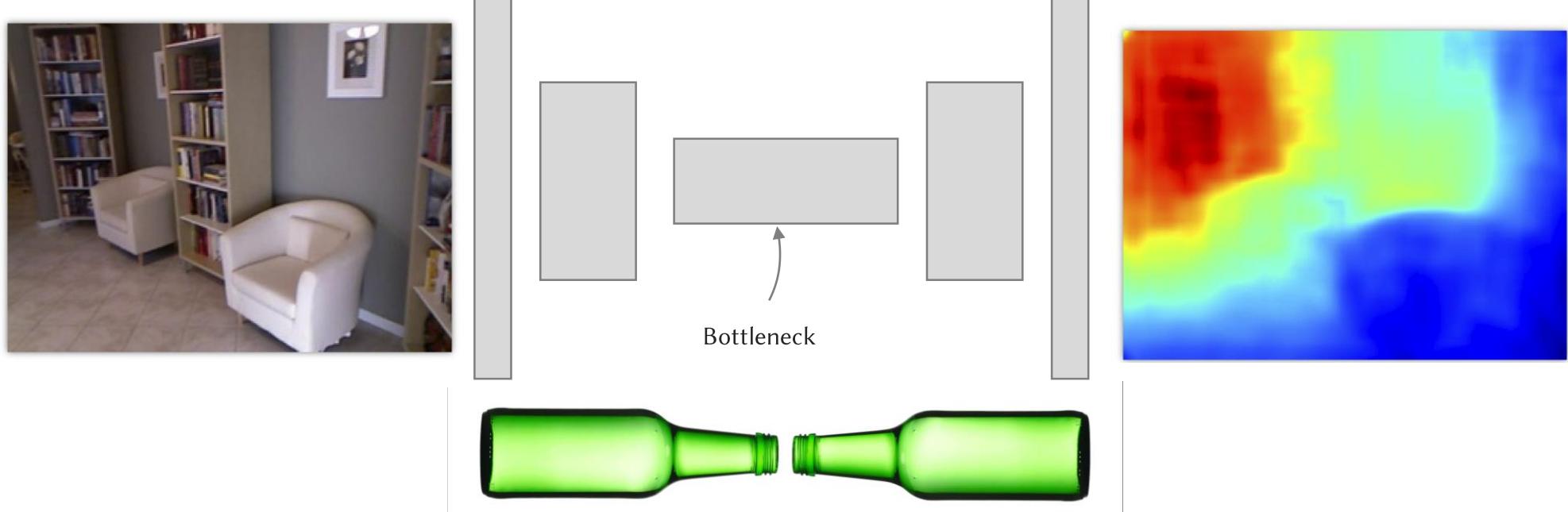
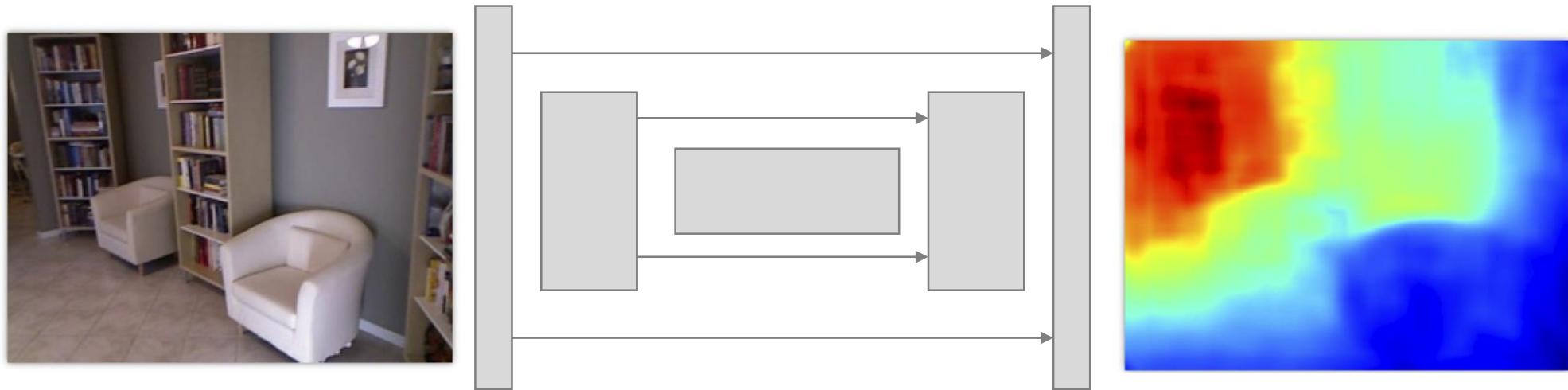
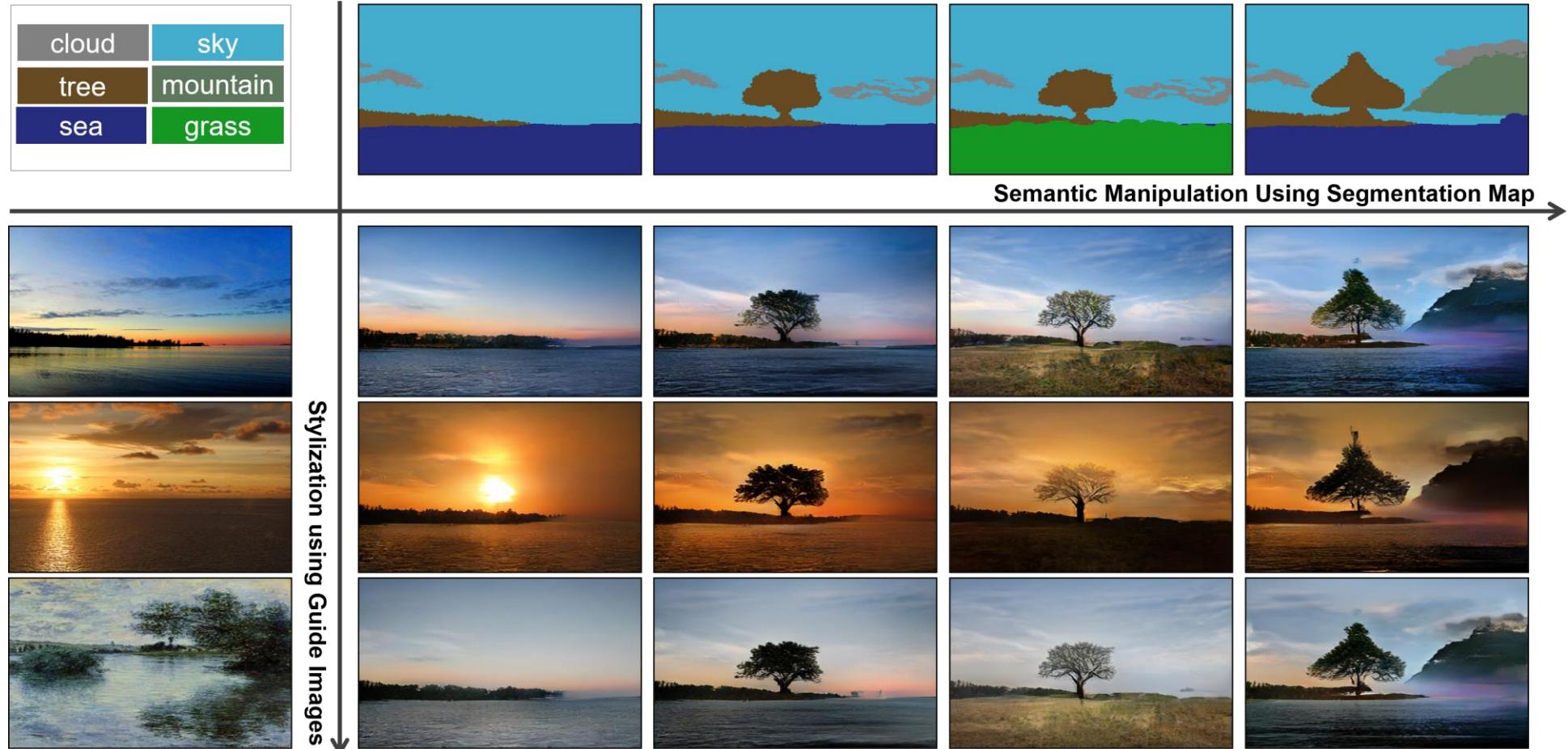


Image-to-image

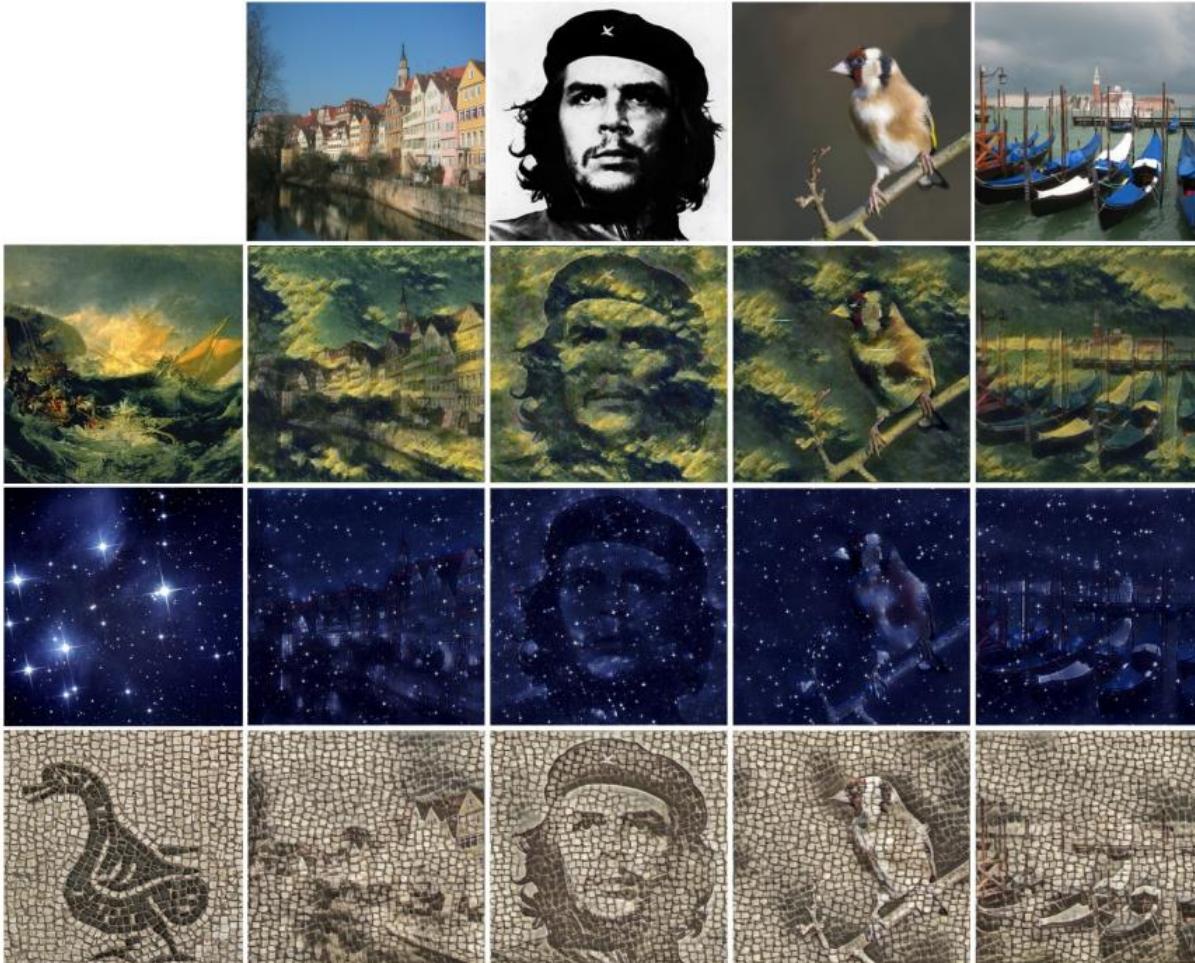
- No chance to get back details
- Idea: **skip connections**, allow decoder to access encoder state (Unet)



Semantic image synthesis

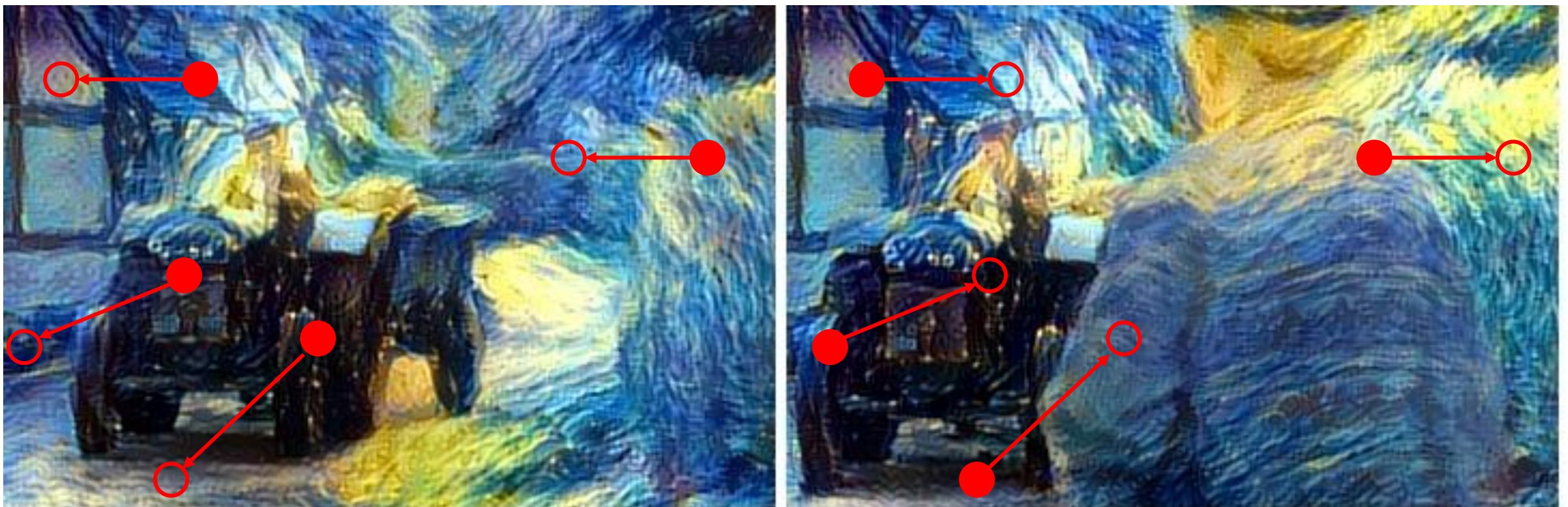


Style transfer



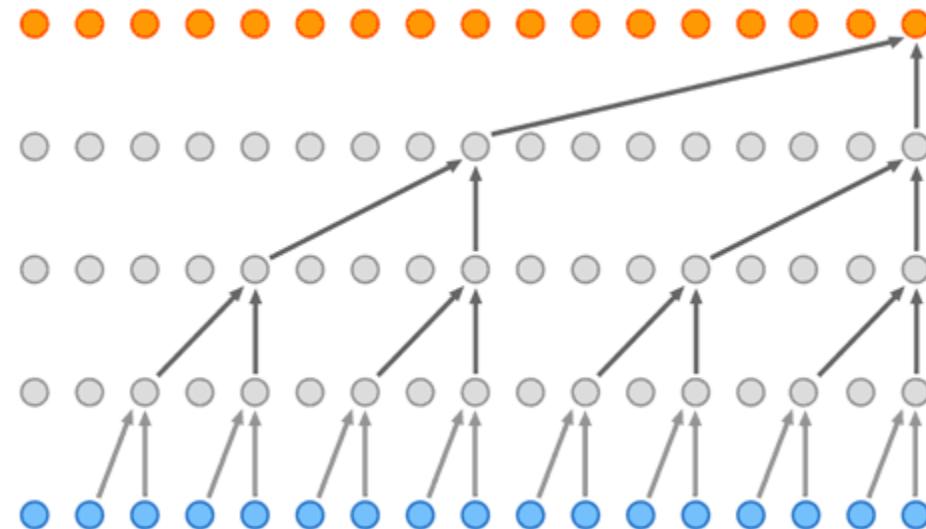
Encoder-decoder for video

- Insert flow-compensated difference into loss



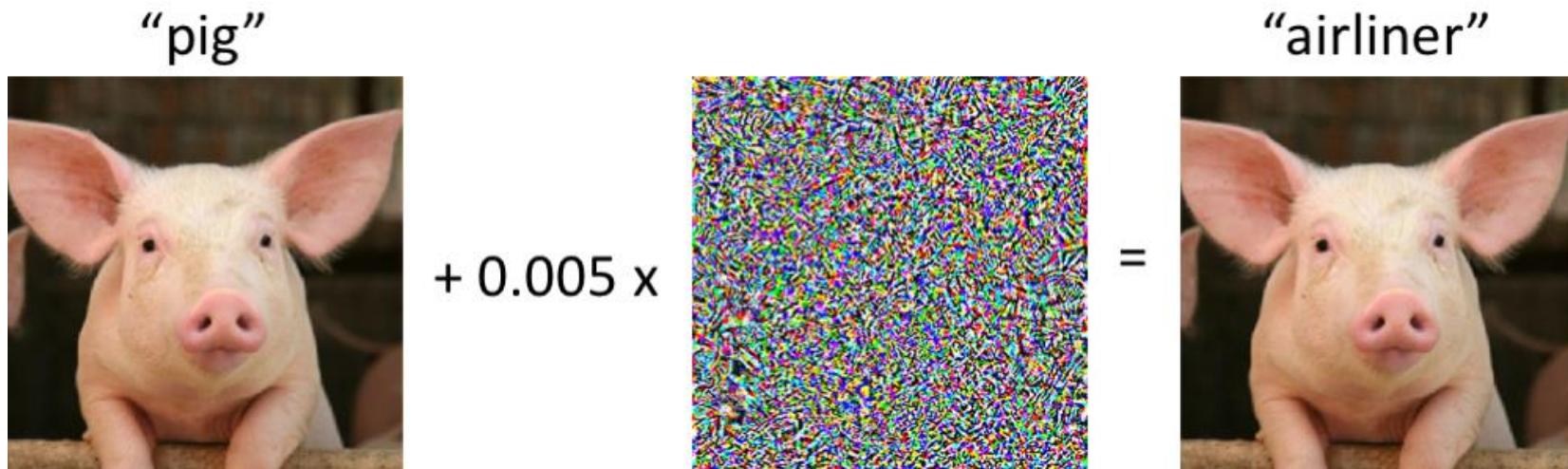
Encoder-Decoder for music

- 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