

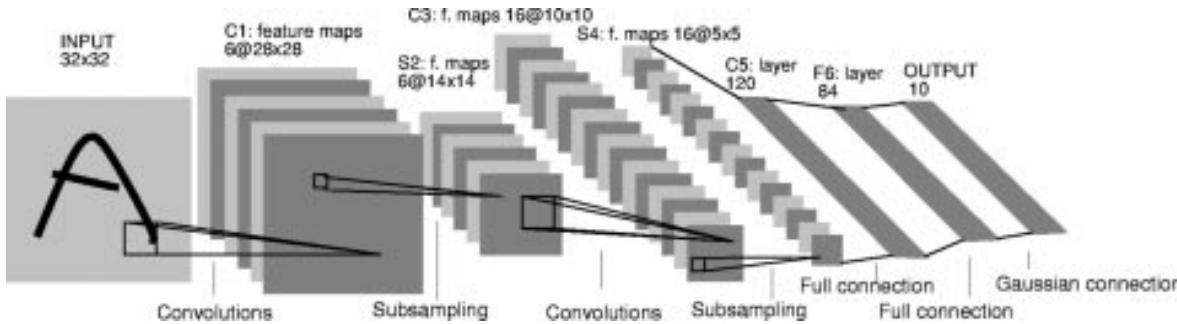
# Deep Learning for Computer Vision

## Part III: Advanced Topics

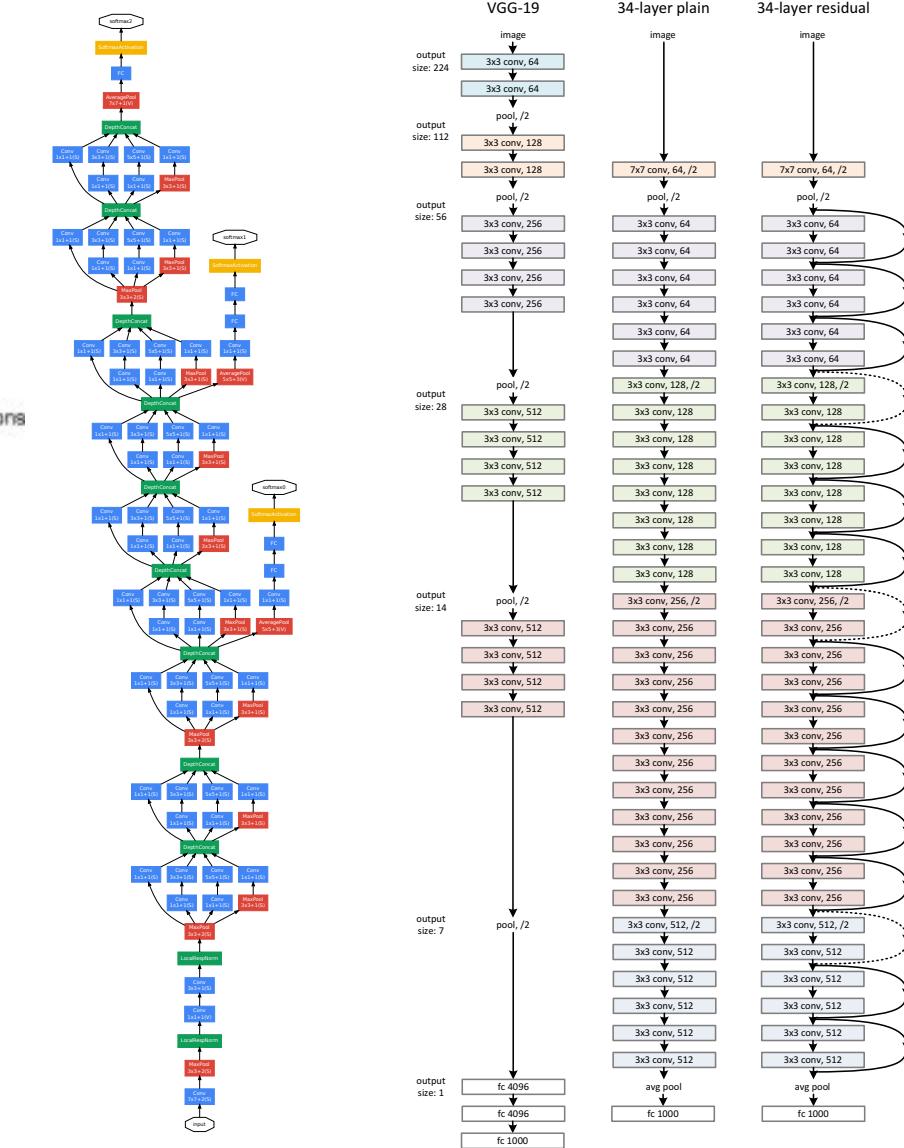
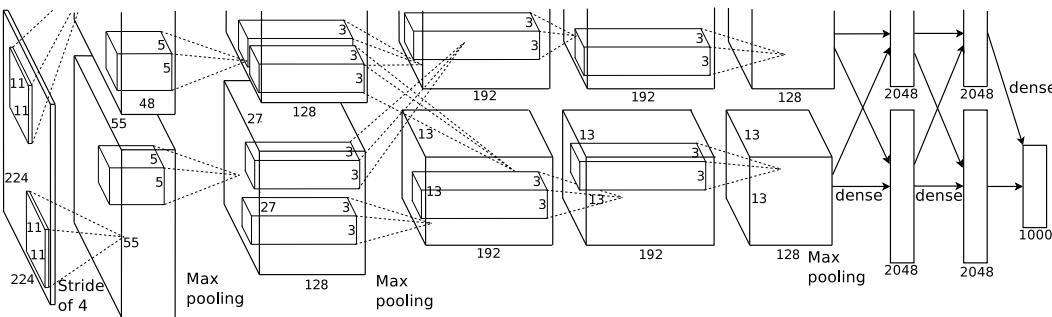
# Outline

1. Introduction to neural networks – this week  
Basics of neural networks
2. Convolutional neural networks
3. Advanced topics and applications

## So far only image-wide classification



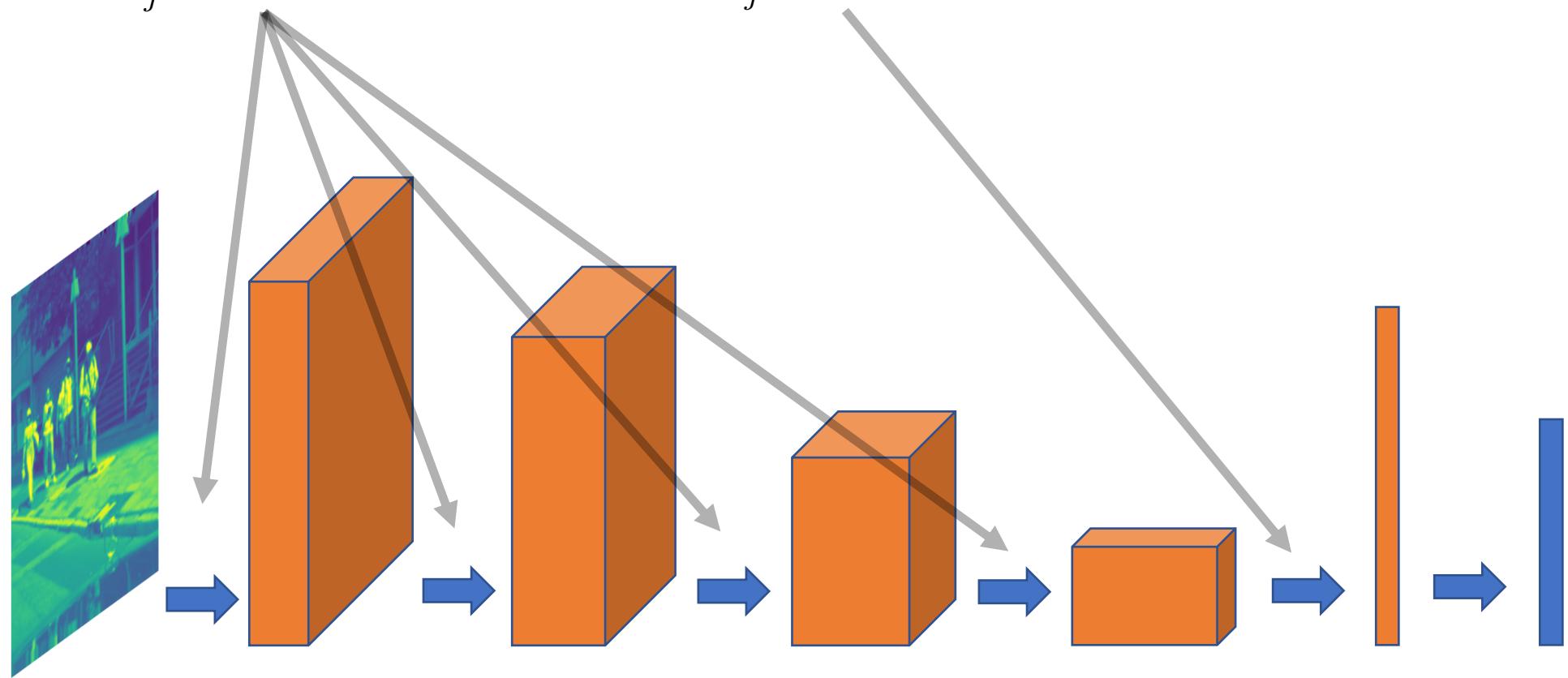
**Fig. 2.** Architecture of LeNet-5, a convolutional NN, here used for digits recognition. Each plane is a feature map, i.e., a set of units whose weights are constrained to be identical.



# Convolutional neural network (CNN) for classification

$h_{l,k} = \sigma(a_{l,k})$  : element-wise non-linear function all across the layers

$$a_{l,k} = \sum_j w_{l,kj} * h_{l-1,j} + b_{l,k} \quad a_{l,k} = \sum_j w_{l,kj} h_{l-1,j} + b_{l,k}$$



# Convolutional neural network (CNN) for classification

$h_{l,k} = \sigma(a_{l,k})$  : element-wise non-linear function all across the layers

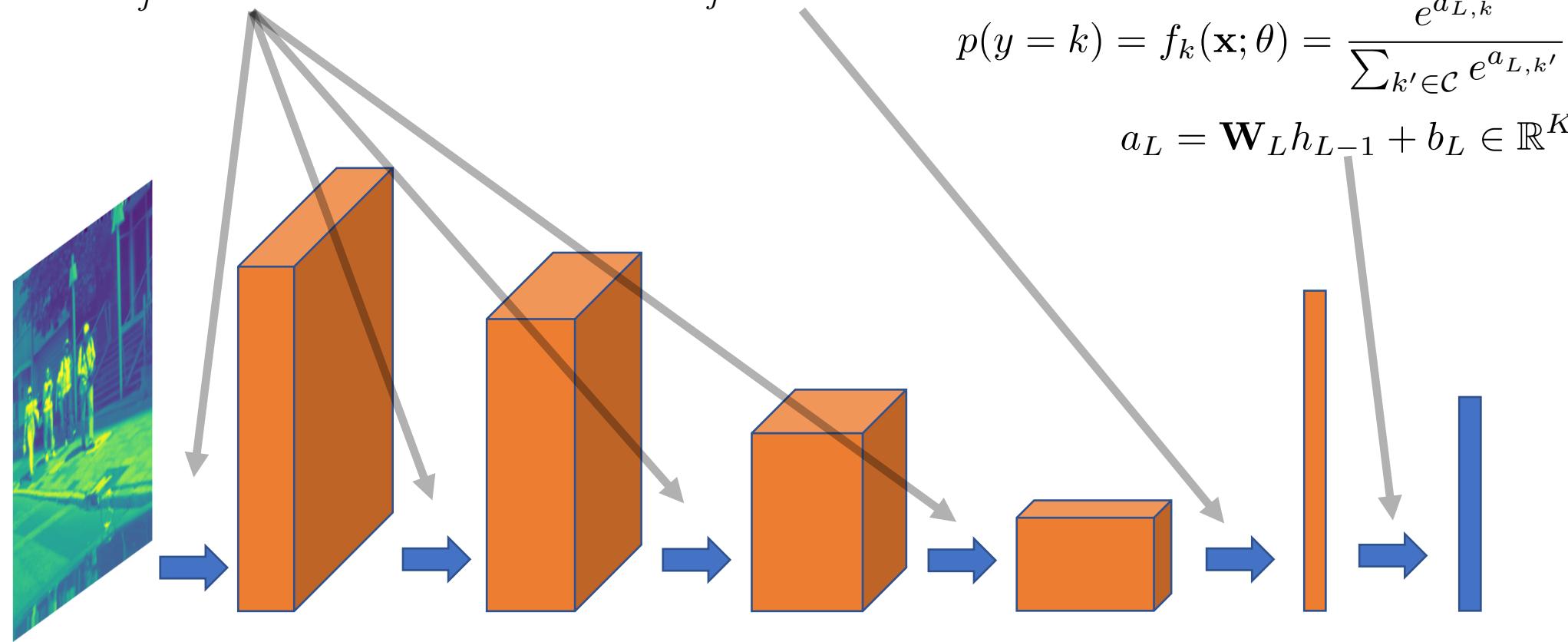
$$a_{l,k} = \sum_j w_{l,kj} * h_{l-1,j} + b_{l,k}$$

$$a_{l,k} = \sum_j w_{l,kj} h_{l-1,j} + b_{l,k}$$

$$\sum_{k \in \mathcal{C}} p(y = k) = 1$$

$$p(y = k) = f_k(\mathbf{x}; \theta) = \frac{e^{a_{L,k}}}{\sum_{k' \in \mathcal{C}} e^{a_{L,k'}}}$$

$$a_L = \mathbf{W}_L h_{L-1} + b_L \in \mathbb{R}^K$$



# Convolutional neural network (CNN) for classification

Cost function

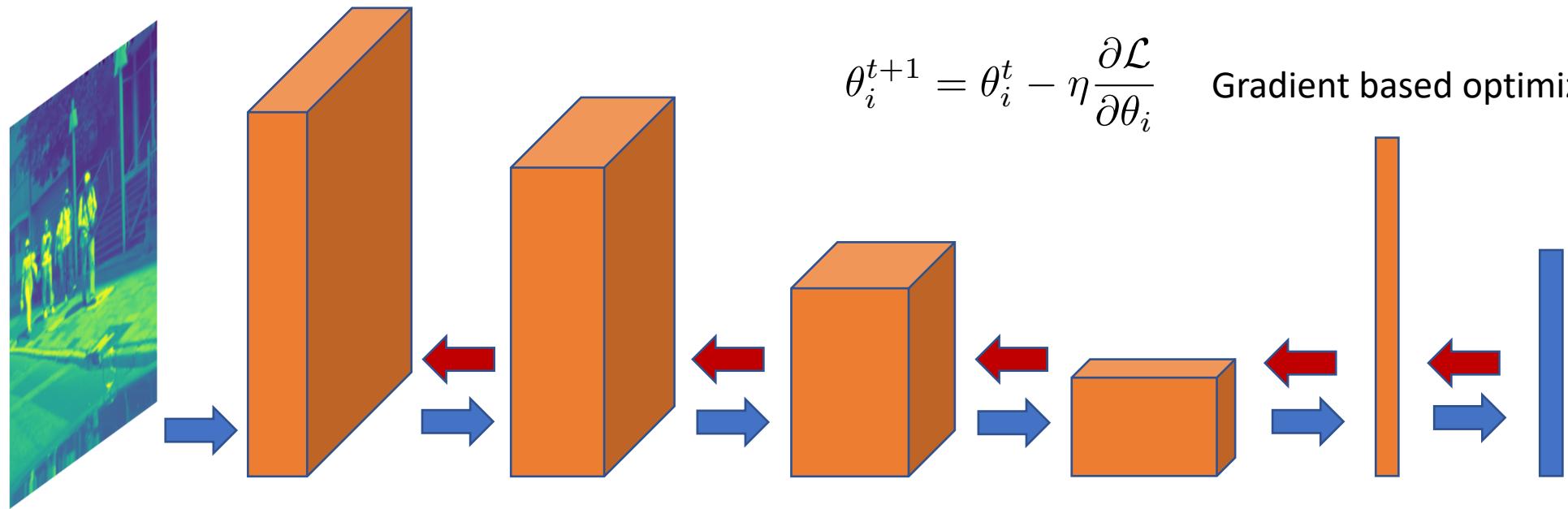
$$\mathcal{L}(y_n, f(\mathbf{x}_n; \theta)) = - \sum_{k \in \mathcal{C}} \log(f_k(\mathbf{x}_n; \theta)) \mathbf{1}(y_n = k)$$

Optimization

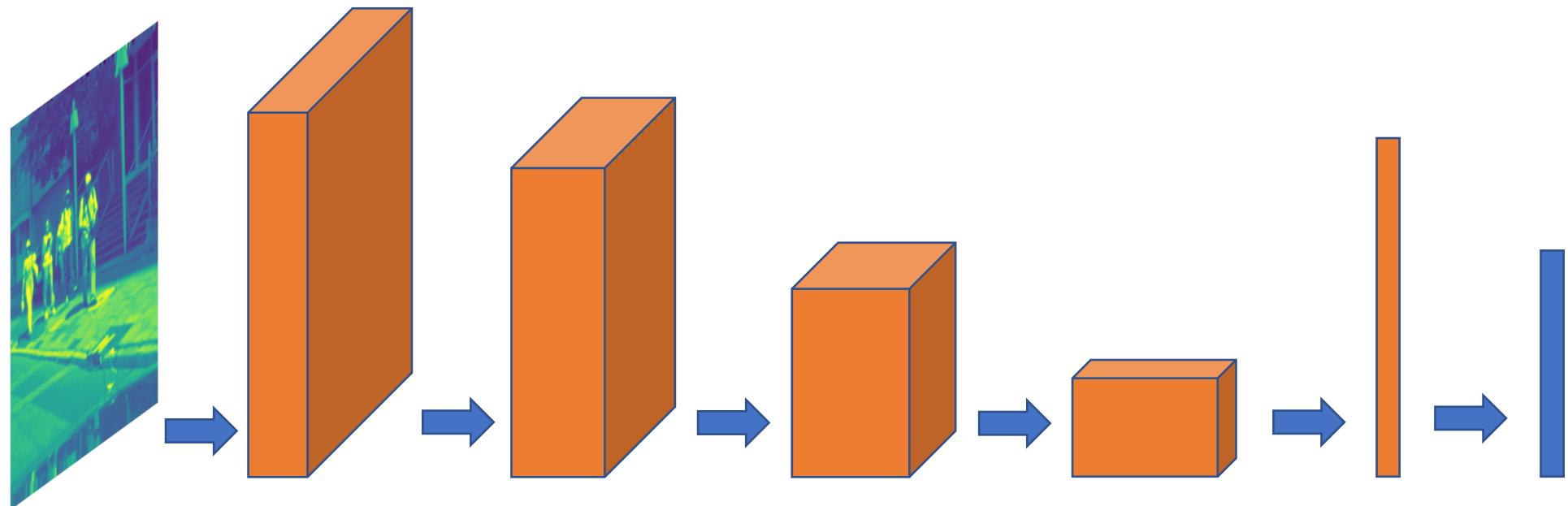
$$\theta^* = \arg_{\theta} \min \sum_{n=1}^N \mathcal{L}(\mathbf{y}_n, f(\mathbf{x}_n; \theta))$$

$$\theta_i^{t+1} = \theta_i^t - \eta \frac{\partial \mathcal{L}}{\partial \theta_i}$$

Gradient based optimization



## How do we use this block to perform other tasks?

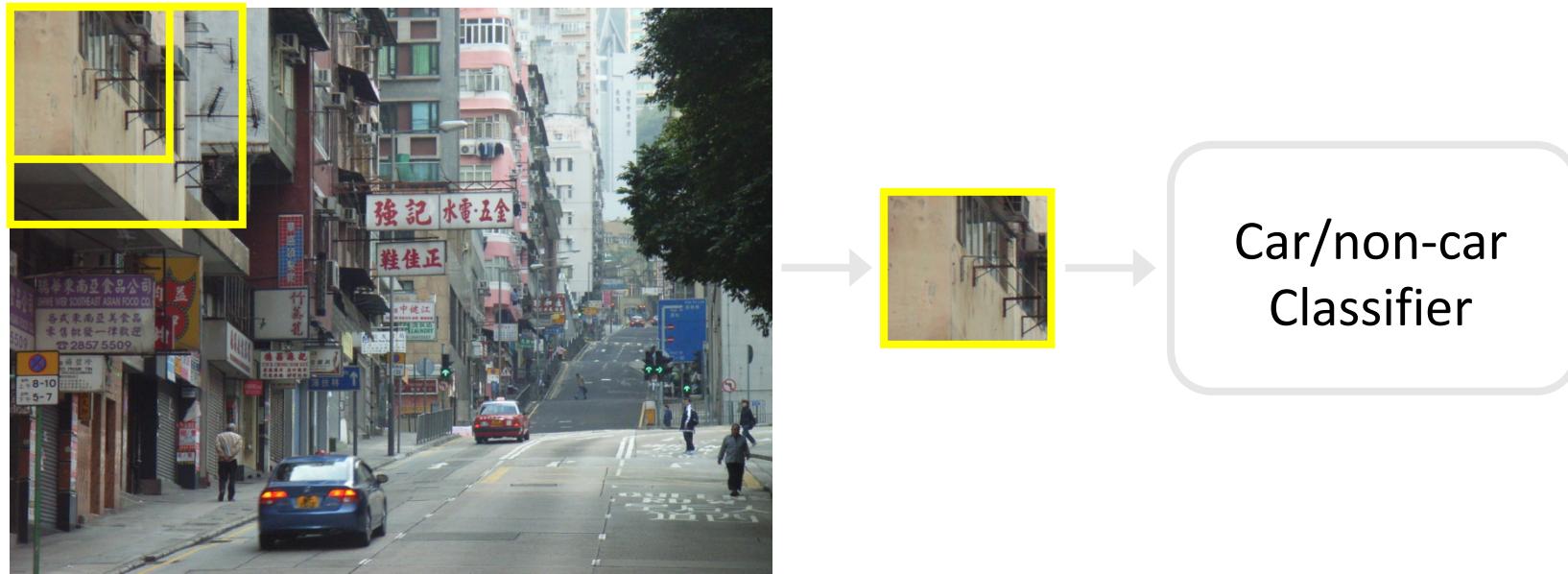


## A bit more detailed outline

3. Advanced topics and applications
  - a) Detection / Localization
  - b) Semantic segmentation
  - c) Visualization and diagnostics
  - d) Unsupervised localization and classification
  - e) Unsupervised learning

# Detection via classification: recall

- If object may be in a cluttered scene, slide a window around looking for it.



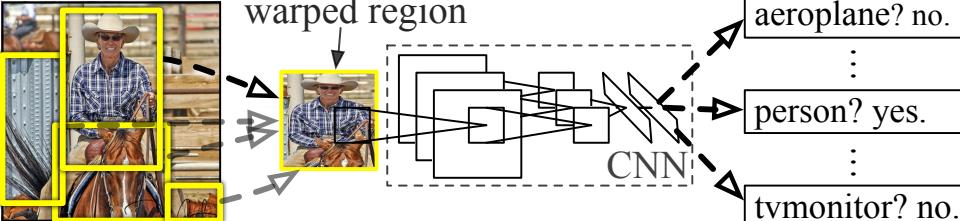
- a brute-force approach with many local decisions.

# Regions with CNN features

## R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

- [Images adapted from the article]
- Very similar idea as sliding windows
- However, instead of sliding windows, they use **region proposal** schemes to make the algorithm very fast.
- Three modules:
  1. Region proposals by **selective search** [Uijlings et al. 2012] and warp regions to a fixed size
  2. CNNs to extract features
  3. SVM to classify

Rich feature hierarchies for accurate object detection and semantic segmentation

Ross Girshick<sup>1</sup> Jeff Donahue<sup>1,2</sup> Trevor Darrell<sup>1,2</sup> Jitendra Malik<sup>1</sup>

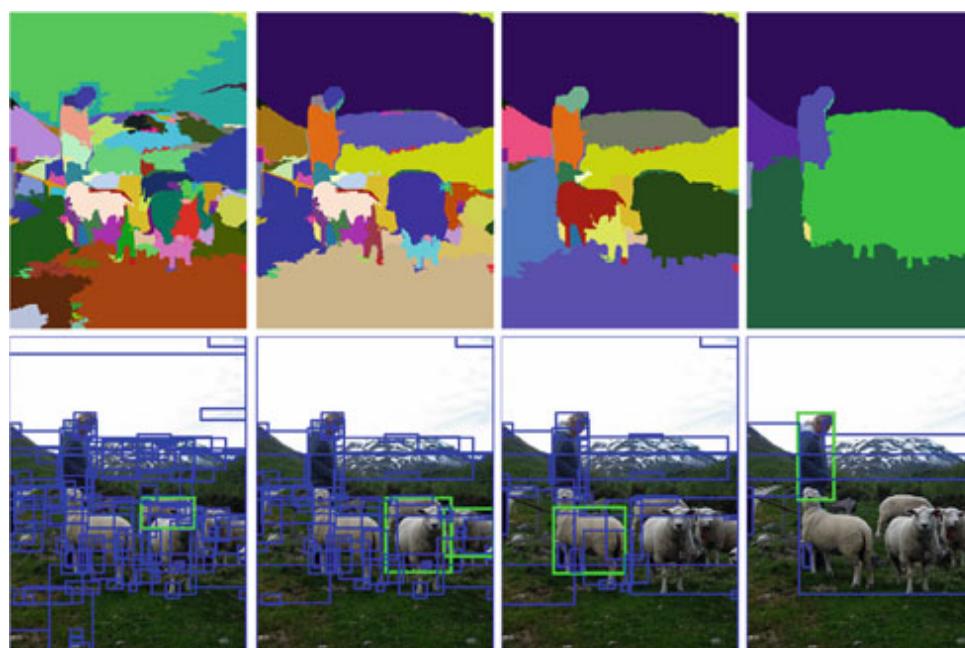
<sup>1</sup>UC Berkeley and <sup>2</sup>ICSI

{rbg, jdonahue, trevor, malik}@eecs.berkeley.edu

# Selective search for region proposals



- Graph-based segmentation
- Clustering pixels based on their intensities in the RGB channels.
- Felzenszwalb and Huttenlocher, IJCV 2004
- [Image taken from the article]



- Hierarchical clustering of the regions.
- Regions at multiple scales
- Aggregate all regions at all scales
- Capturing objects at different scales
- [Uijlings et al., IJCV 2013]
- [Image taken from the article]

# Converting regions into features with a CNN

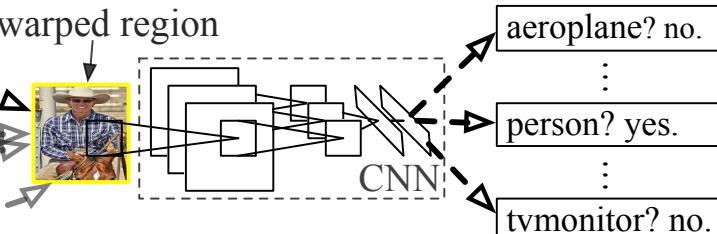
## R-CNN: *Regions with CNN features*



1. Input  
image



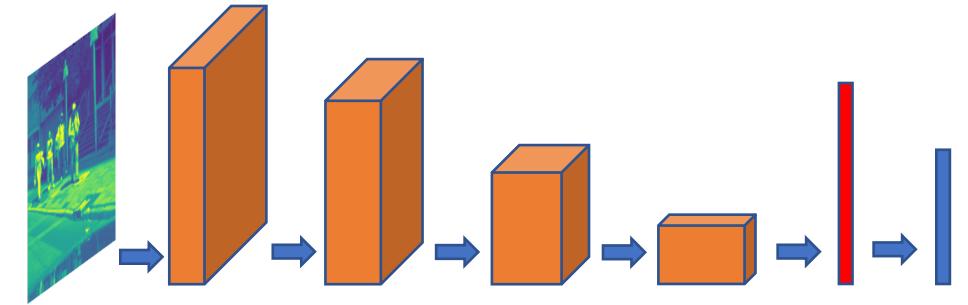
2. Extract region  
proposals (~2k)



3. Compute  
CNN features

4. Classify  
regions

- [Images adapted from the article]
- Regions are reshaped / warped
- Fed into a CNN and converted into a code of 4096 dimensions



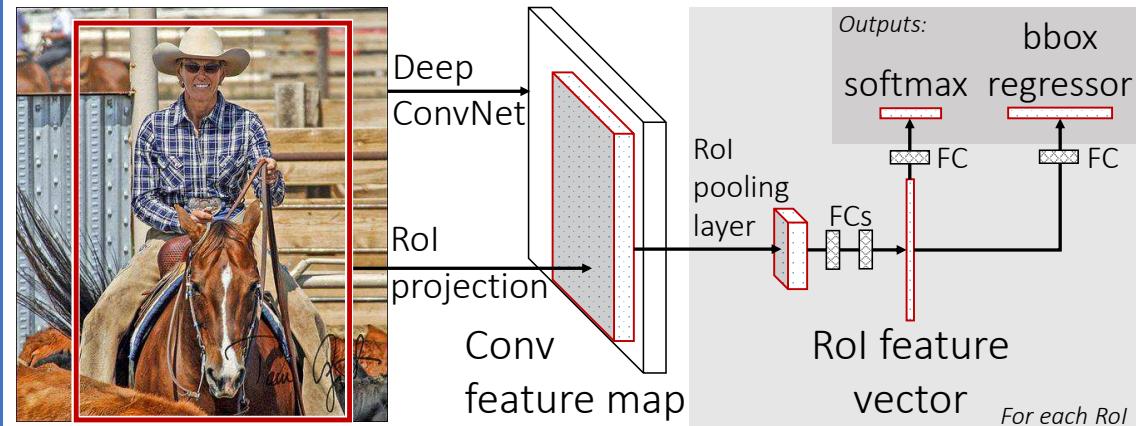
- Codes are used in an SVM for classification : Binary classification
  - Repeat it for all the proposed regions
- Multi-stage method and can be slow

Rich feature hierarchies for accurate object detection and semantic segmentation

Ross Girshick<sup>1</sup> Jeff Donahue<sup>1,2</sup> Trevor Darrell<sup>1,2</sup> Jitendra Malik<sup>1</sup>

<sup>1</sup>UC Berkeley and <sup>2</sup>ICSI

{rbg, jdonahue, trevor, malik}@eecs.berkeley.edu



- [Images adapted from publication]
- Replaced warping with a "ROI pooling" layer – ROI : region of interest
  - 1. Region proposals as before
  - 2. A CNN takes input the image and region proposals and classifies each region.
- ROI pooling converts arbitrary size ROI in any channel to a  $H \times W$  window.
- Defines the back-propagation for the ROI-pooling layer
- Multi-class classification
- Faster at test time

# How do we extend this to semantic segmentation?

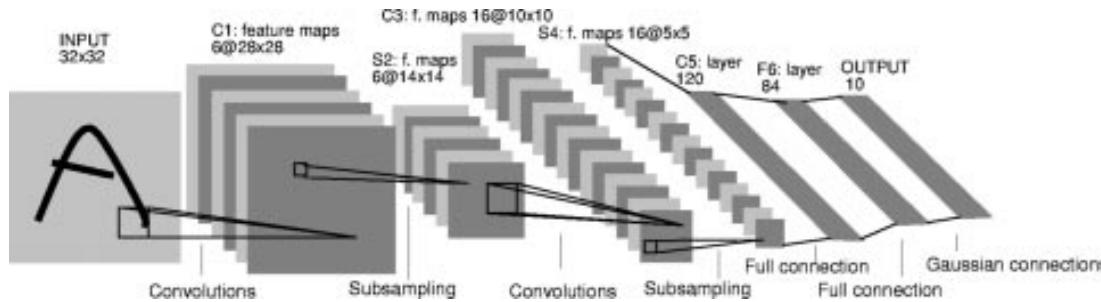
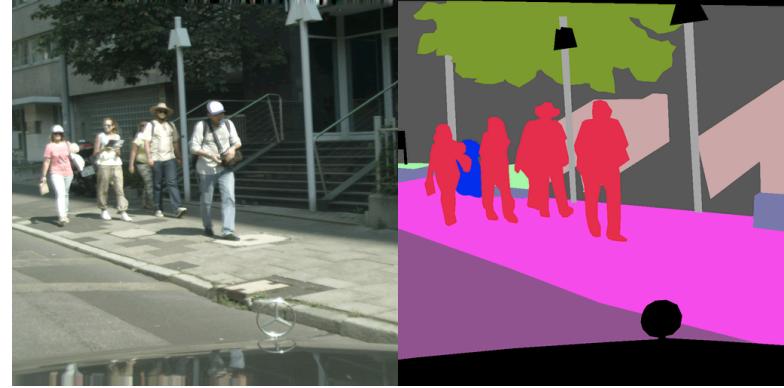


Fig. 2. Architecture of LeNet-5, a convolutional NN, here used for digits recognition. Each plane is a feature map, i.e., a set of units whose weights are constrained to be identical.



- Classification CNN progressively reduces the size of the network
- Segmentation needs to keep the size of the image constant
- Pixel-wise classification
- Commonly used training loss is pixel-wise cross entropy

# Two competing factors



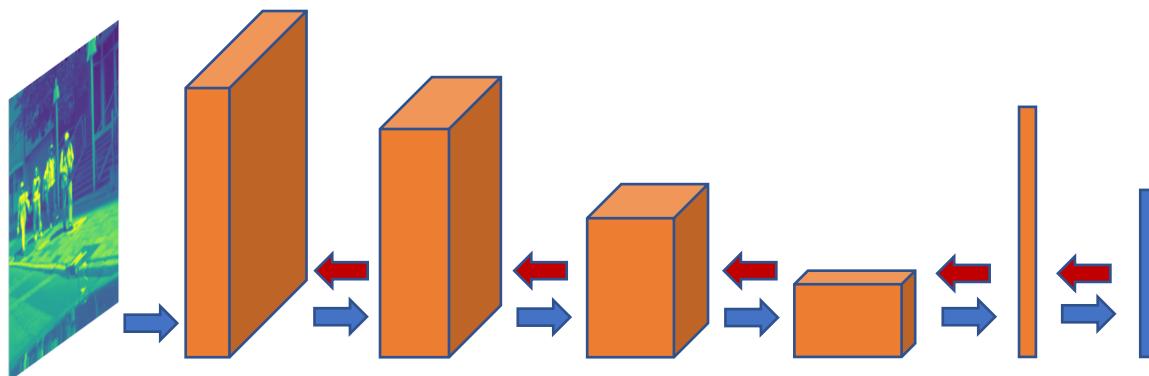
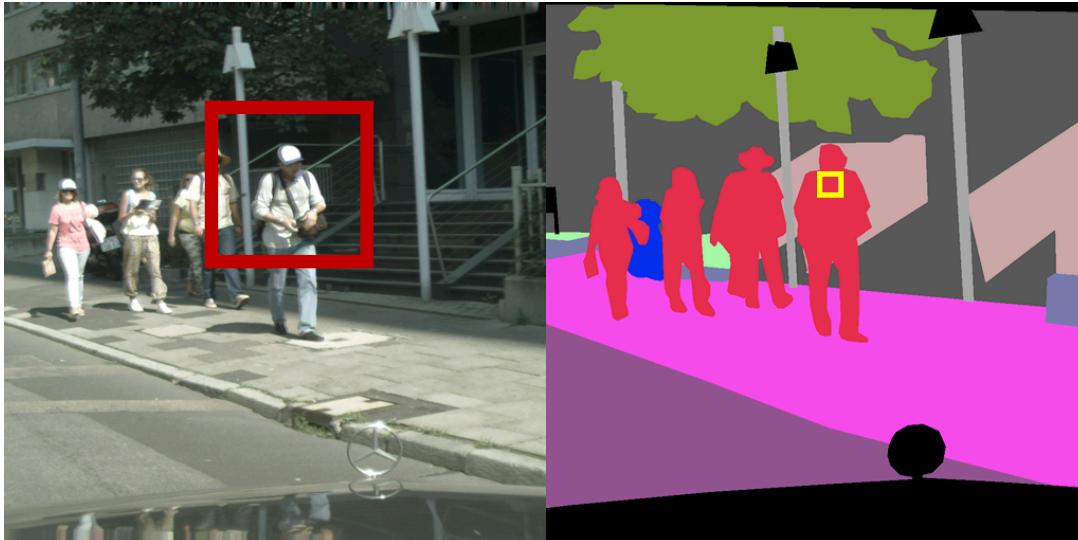
## Integrating larger context

- Decision for each pixel should look at a large portion of the image to decide the semantic label.
- But this means information for neighboring pixels are very similar – difficult to disambiguate boundaries.

## Distinguishing neighboring boundary pixels

- Network should be able to distinguish neighboring pixels – focusing on fine level details.
- Only focusing on fine level details cannot integrate the larger image context and have trouble determining the semantic labels.

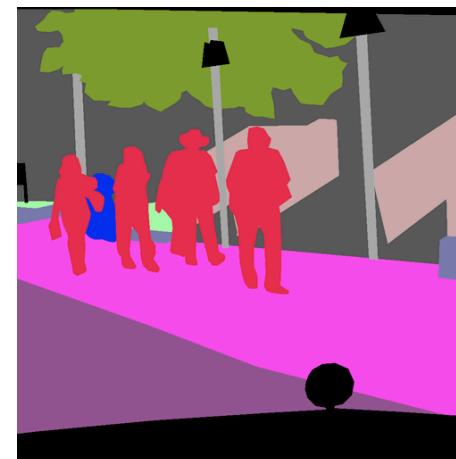
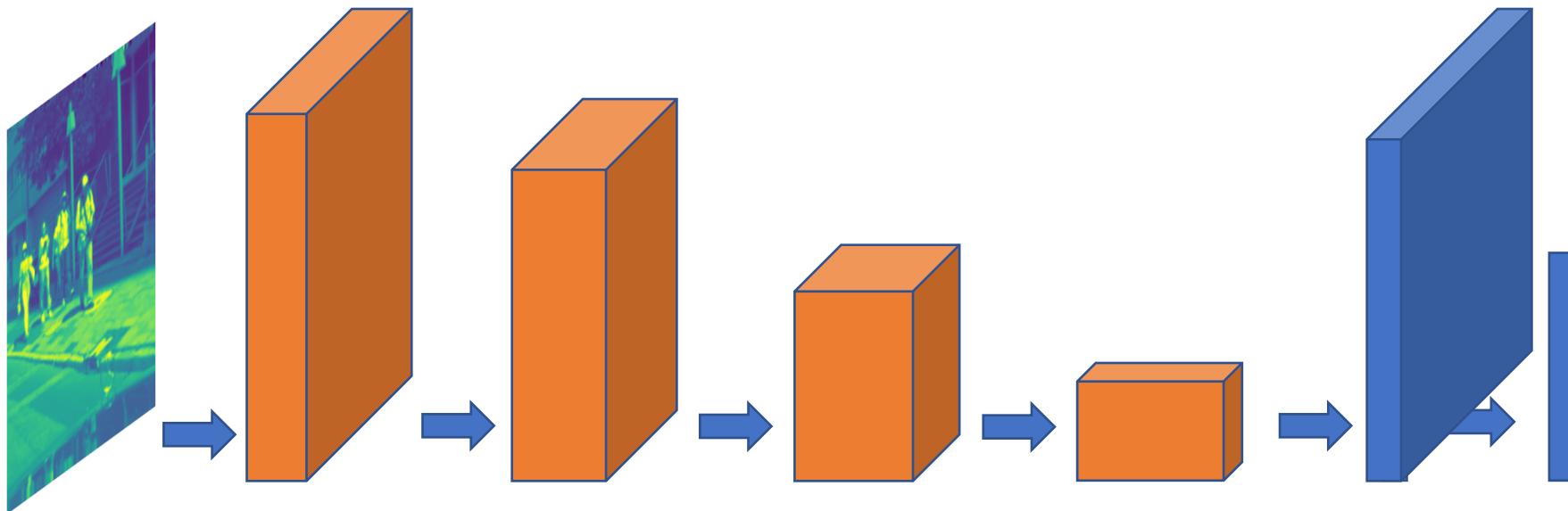
# Naïve approach



- Classification network for each pixel
- Sliding window approach
- Each window only predicts the central pixel's class
- Conceptually easy
- Straightforward training
- Can be slow in test time
- Good in context not very good in boundaries

# Fully convolutional network

Depth equals the  
number of possible  
semantic classes



Pixel-wise  
Softmax

## Fully Convolutional Networks for Semantic Segmentation

Jonathan Long\*

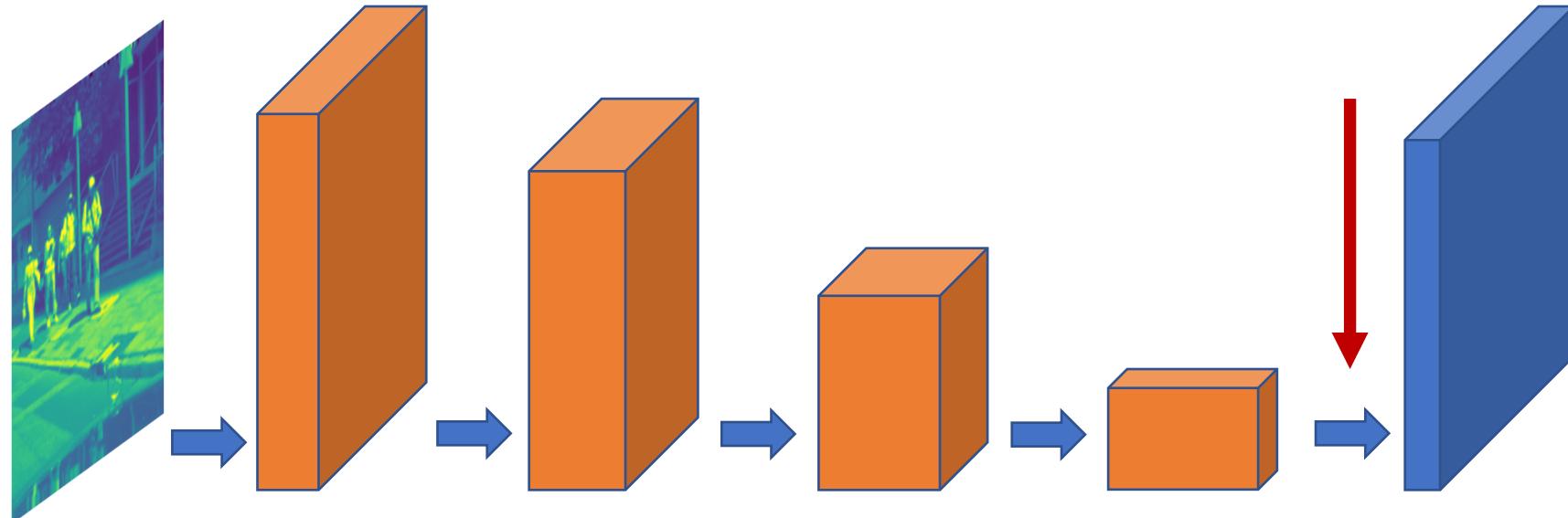
Evan Shelhamer\*

Trevor Darrell

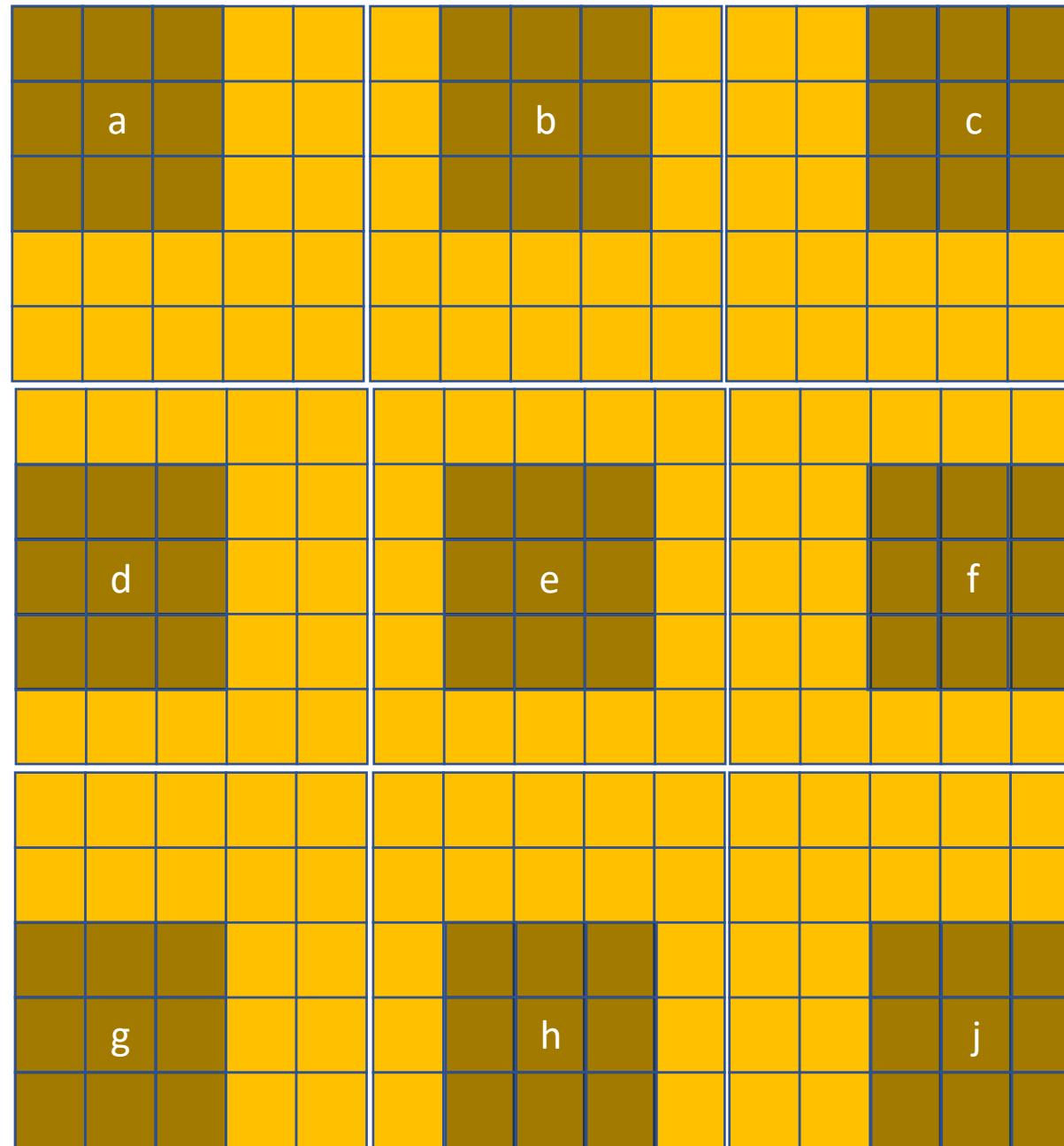
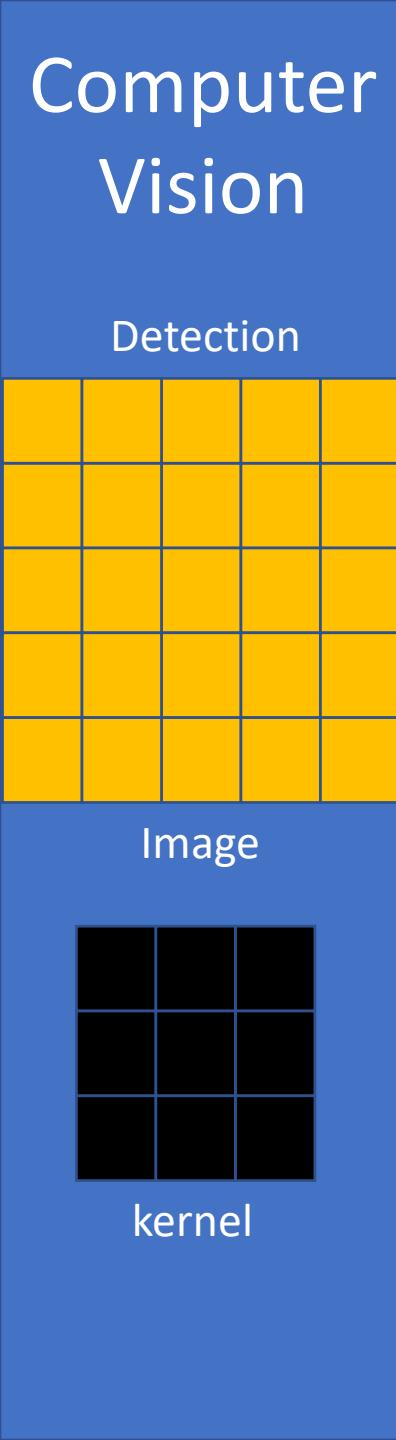
UC Berkeley

{jonlong, shelhamer, trevor}@cs.berkeley.edu

# How did that happen?



- The penultimate layer is deep but has a very small channel size
- The last, output, layer has a very large channel size – as large as the input image
- There are two basic ways of doing this: transposed convolutions and upsampling



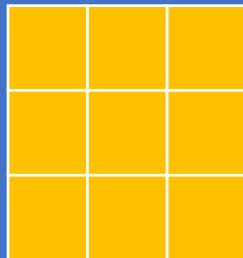
**Convolution**

a	b	c
d	e	f
g	h	j

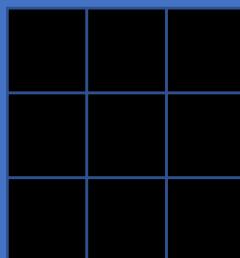
# Computer Vision

Detection

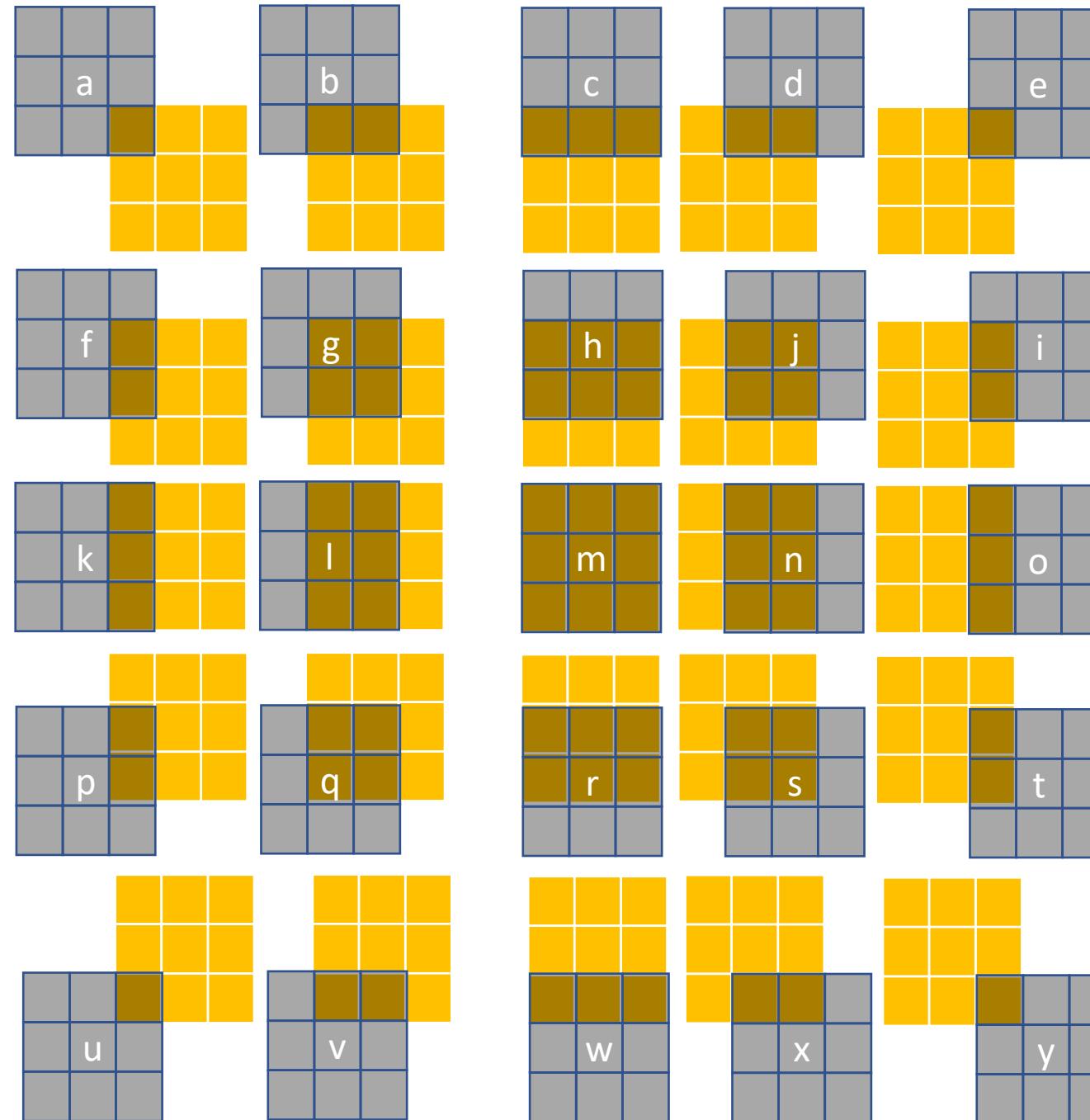
Segmentation



Image



kernel



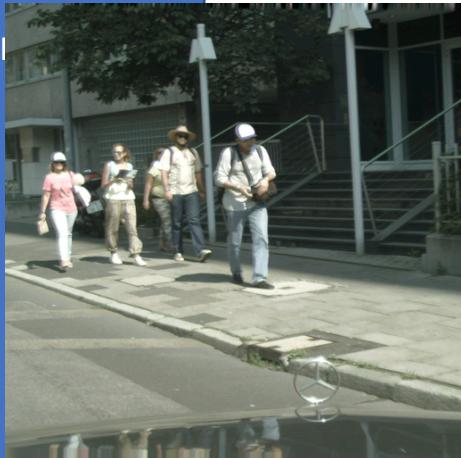
## Transposed Convolution

Sometimes referred  
to as deconvolution  
but that is not  
correct terminology.

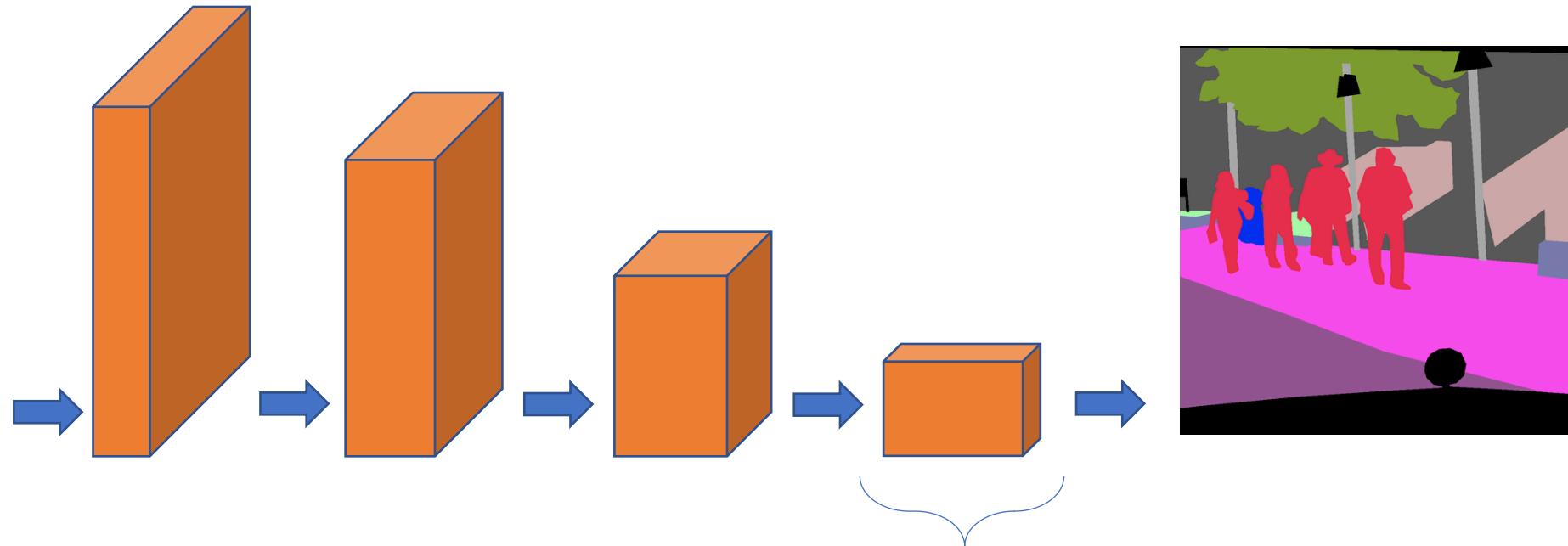
a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y

# Up-sampling followed by convolution



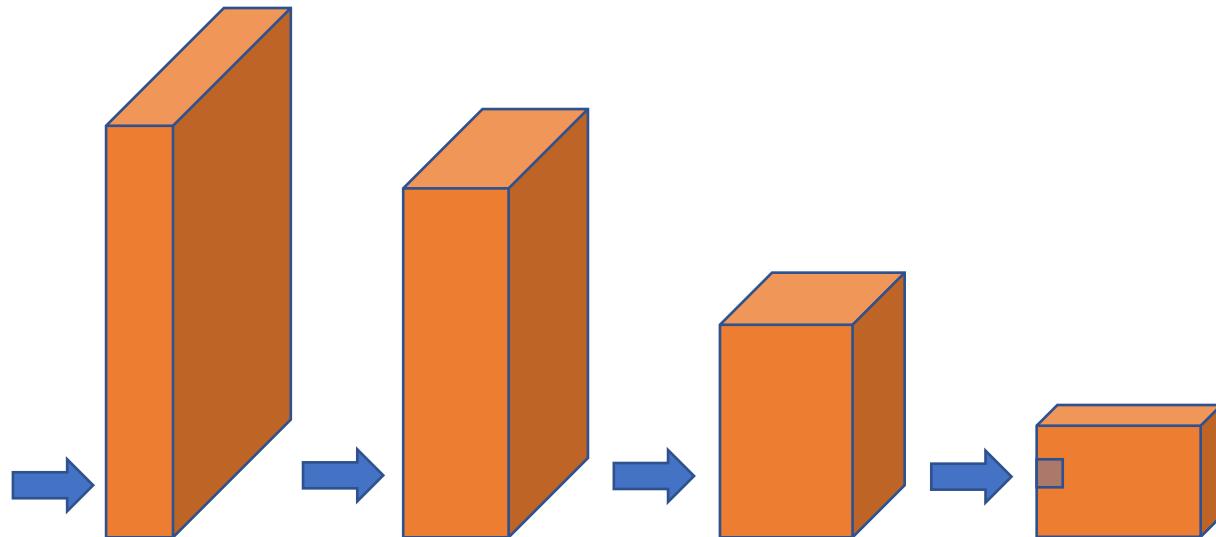


# Receptive fields



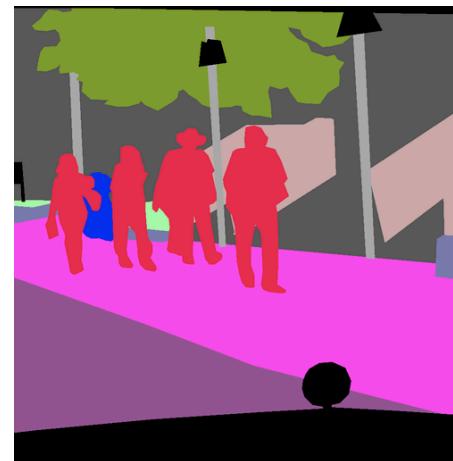
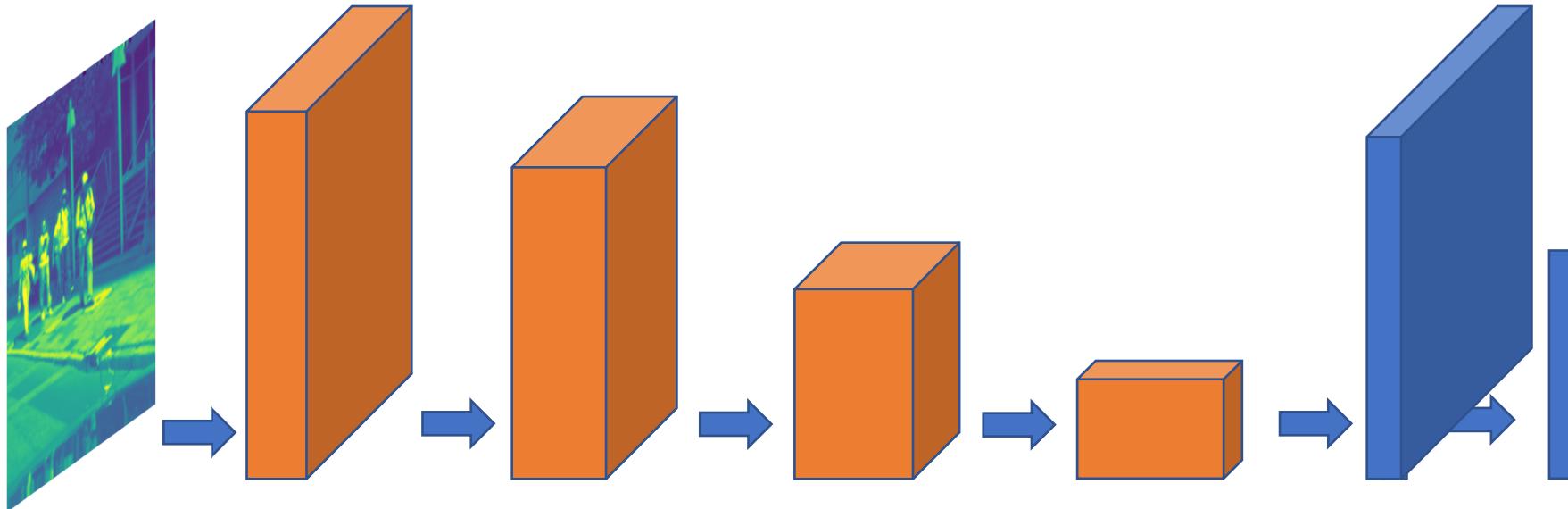
- This channel must see large areas of the image so we can determine semantic labels for the pixels.
- Due to the preceding convolutions and pooling each neuron in this channel sees a large area of the image – it has a **receptive field**

# Receptive fields



- Receptive field grows progressively as we go backwards through the network
- Normally, to grow the receptive field the options are:
  - Convolution, pooling, strides
  - Pooling and strides reduce the size of the channels by a large margin
  - Leads to resolution and coverage loss

# Solves the context integration



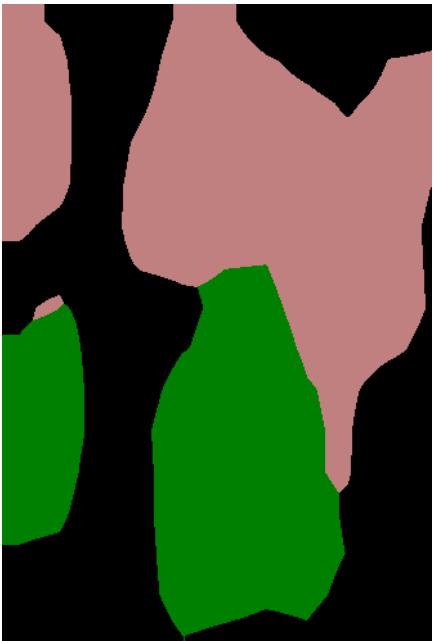
## Fully Convolutional Networks for Semantic Segmentation

Jonathan Long\*   Evan Shelhamer\*   Trevor Darrell  
UC Berkeley

[{jonlong,shelhamer,trevor}@cs.berkeley.edu](mailto:{jonlong,shelhamer,trevor}@cs.berkeley.edu)

# Trouble with distinguishing between boundary pixels

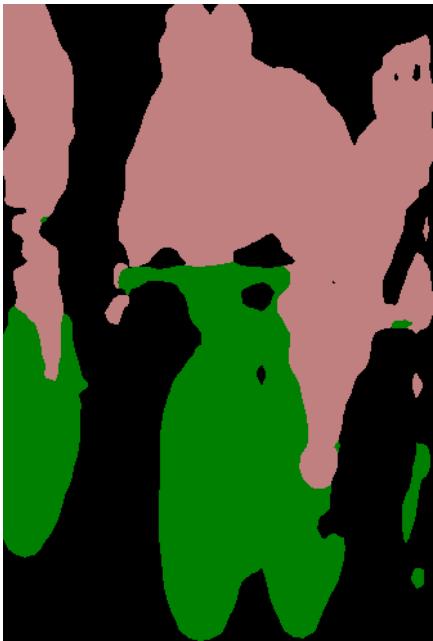
FCN-32s



FCN-16s



FCN-8s



Ground truth



## Fully Convolutional Networks for Semantic Segmentation

Jonathan Long\*

Evan Shelhamer\*

Trevor Darrell

UC Berkeley

{jonlong, shelhamer, trevor}@cs.berkeley.edu

Some modifications can mitigate  
the issue to some extent

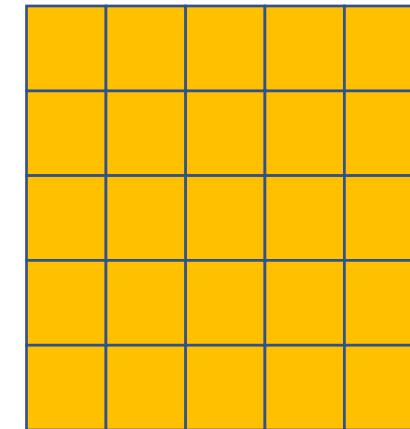
# Alternative I: dilated convolutions

$$J(p) = \sum_t I(p - t)k(t) = k * I$$

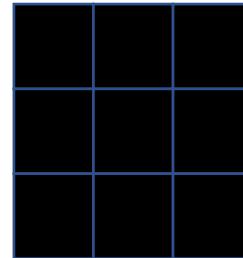
Normal convolution

$$J(p) = \sum_t I(p - lt)k(t) = k_l * I$$

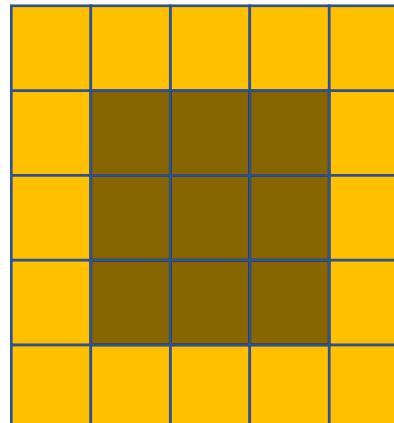
Dilated convolution



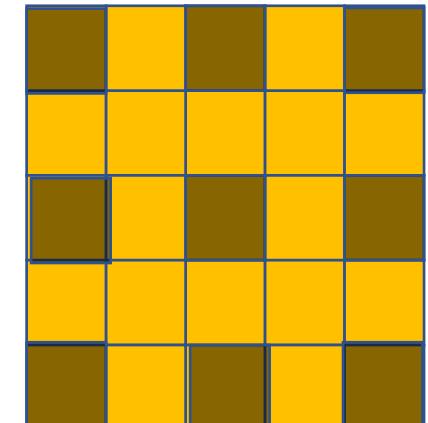
Image



kernel

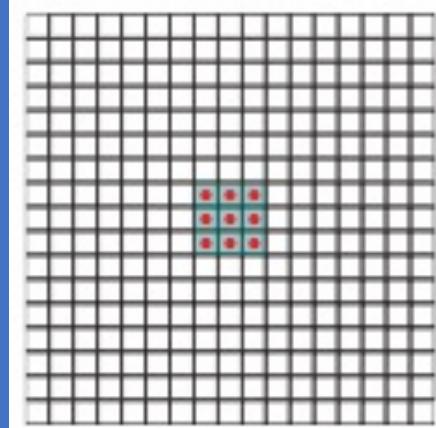


Normal convolution

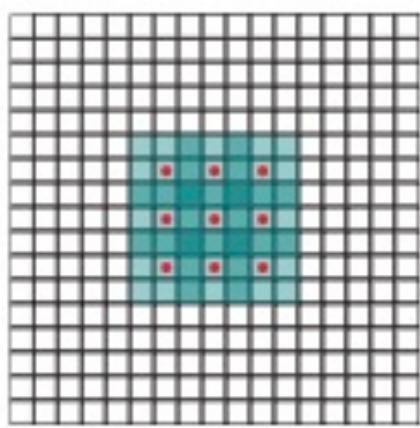


Dilated convolution

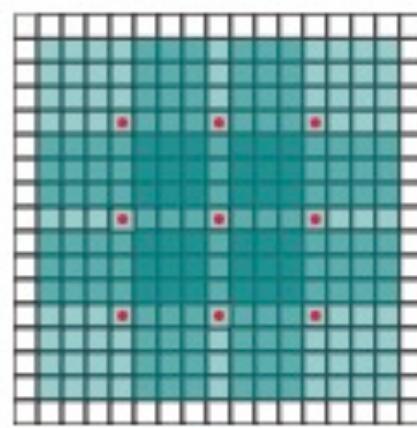
# Alternative I: dilated convolutions



$$k * I$$



$$k'_2 * k * I$$



$$k''_4 * k'_2 * k * I$$

MULTI-SCALE CONTEXT AGGREGATION BY  
DILATED CONVOLUTIONS

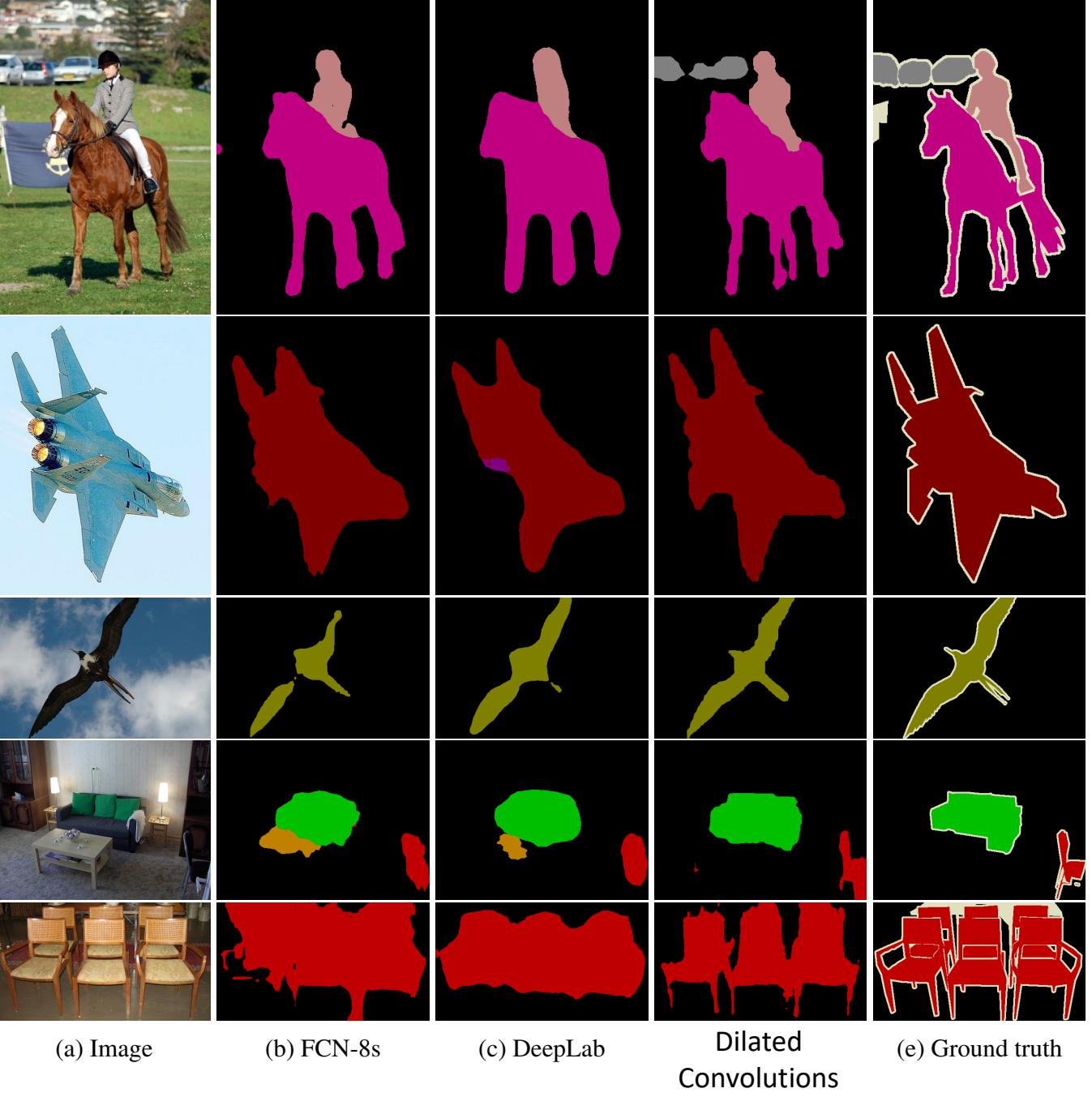
- Convolution with dilated filters
- Can achieve substantial increase in receptive field
- Without substantial increase in parameters

## MULTI-SCALE CONTEXT AGGREGATION BY DILATED CONVOLUTIONS

Fisher Yu  
Princeton University

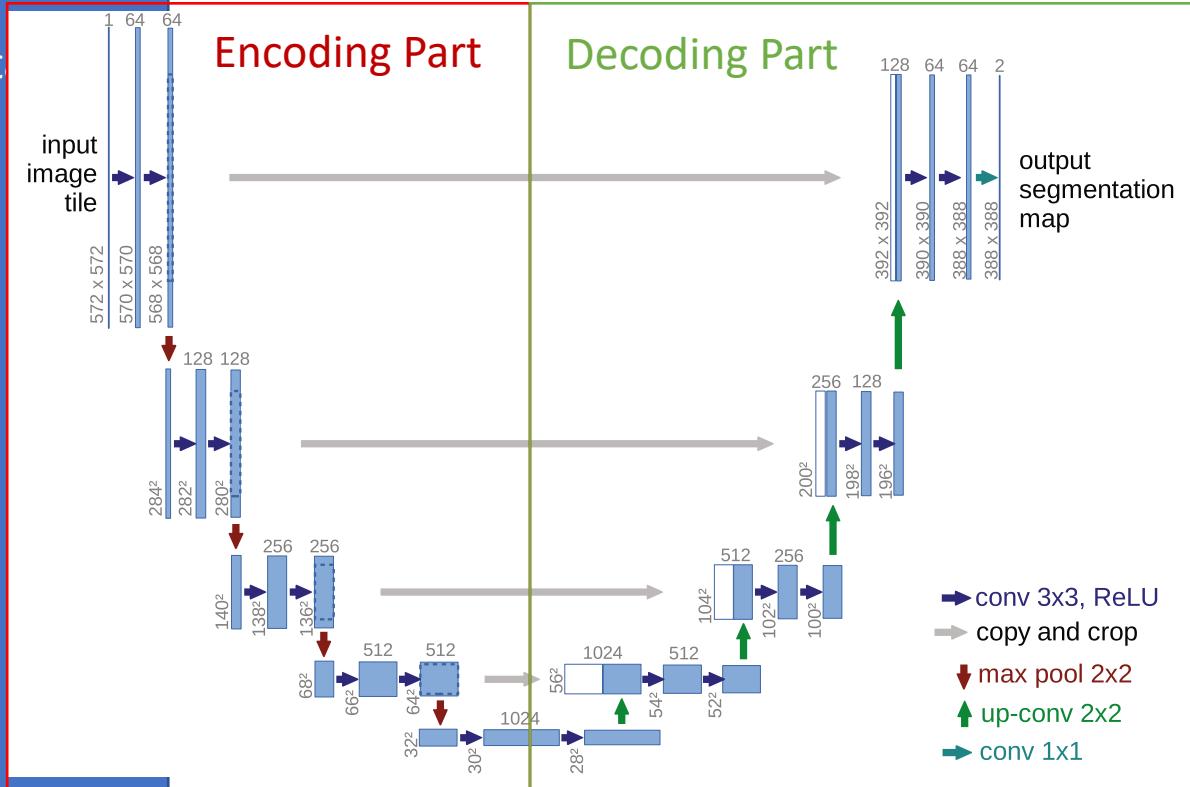
Vladlen Koltun  
Intel Labs

Leads to segmentation  
results that do not lose  
resolution, maintaining high  
level details



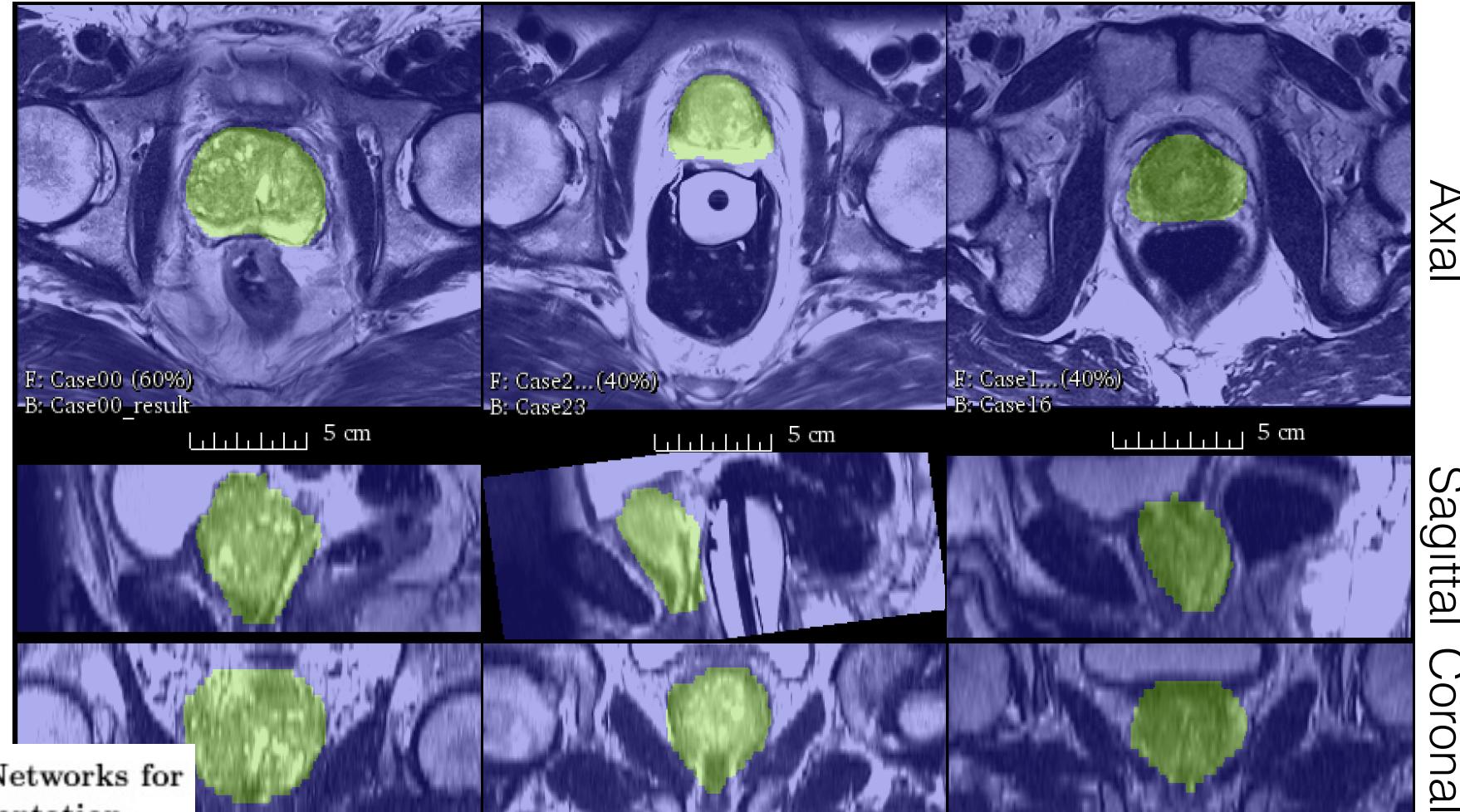
[Image adapted from the publication]

# Alternative II: Skip Connections – U-Net



- The **encoding part** of the network is like the fully convolutional network
- It pools local information and achieve global representation at the bottom layers
- The skip connections passes high-resolution information to retain information necessary to distinguish neighboring boundary pixels.

# Context and fine level



V-Net: Fully Convolutional Neural Networks for  
Volumetric Medical Image Segmentation

# Visualization and diagnostics

# Understanding how a network works

- Challenging task
  - Multivariate interactions, information in different areas of the image are used in interaction with each other
  - Nonlinear mapping between features and labels
  - Hierarchical mapping, information gathered in multiple layers
- Definition of ‘understanding’ is crucial
  - What do you exactly want to get out of the system?
  - There are different approaches with different definitions
  - We will see one particular example: Visualizing features

# Visualizing features

- Discussion based on [Zeiler and Fergus 2013]

---

## Visualizing and Understanding Convolutional Networks

---

**Matthew D. Zeiler**

Dept. of Computer Science, Courant Institute, New York University

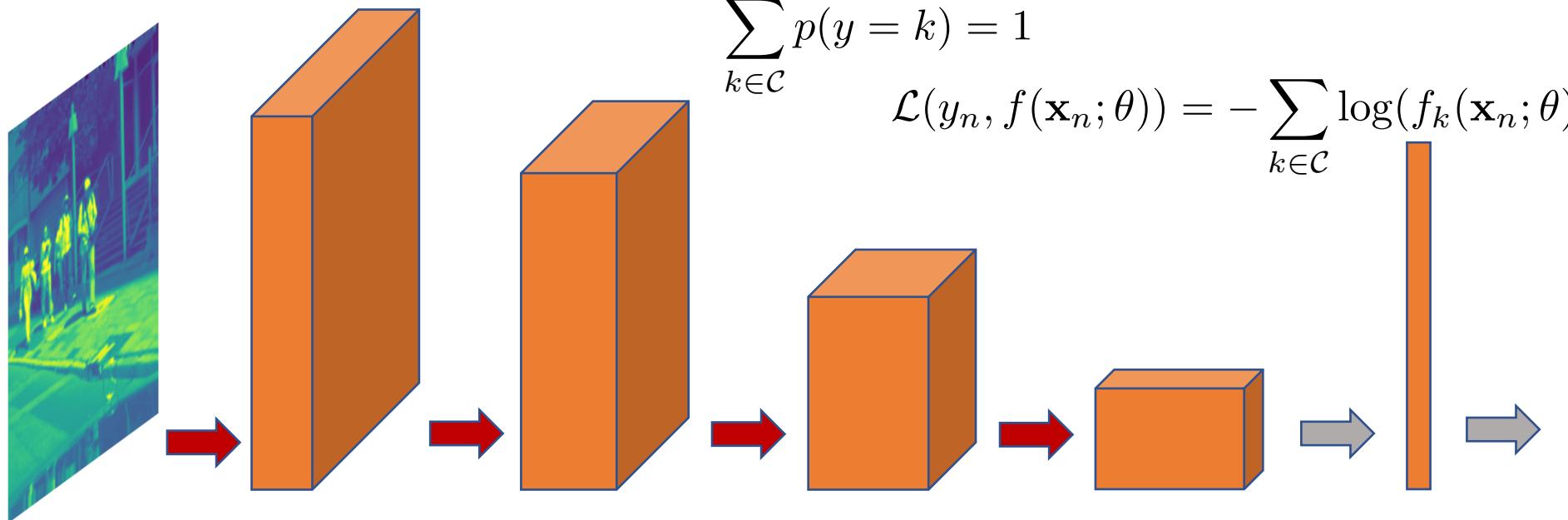
ZEILER@CS.NYU.EDU

**Rob Fergus**

Dept. of Computer Science, Courant Institute, New York University

FERGUS@CS.NYU.EDU

- Visualizing the input that activates a neuron in any layer
- ``Deconvolutional'' network



Convolution followed by ReLu nonlinearity followed by max-pooling [optionally]



Fully connected layer: transformation followed by non-linearity

Output

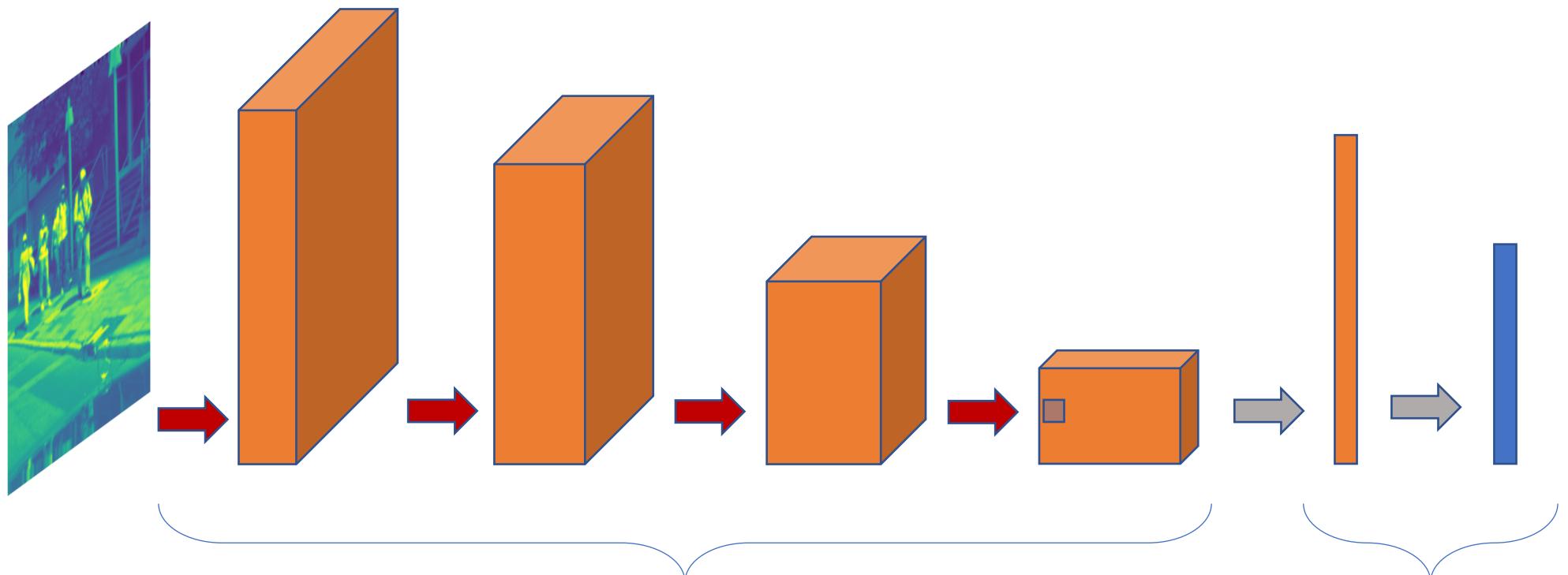
Number of neurons  
equal number of classes

Similar to LeCun et al. 1998 and Krizhevsky et al. 2012

$$a_L = \mathbf{W}_L h_{L-1} + b_L \in \mathbb{R}^K \quad p(y = k) = f_k(\mathbf{x}; \theta) = \frac{e^{a_{L,k}}}{\sum_{k' \in \mathcal{C}} e^{a_{L,k'}}$$

$$\sum_{k \in \mathcal{C}} p(y = k) = 1$$

$$\mathcal{L}(y_n, f(\mathbf{x}_n; \theta)) = - \sum_{k \in \mathcal{C}} \log(f_k(\mathbf{x}_n; \theta)) \mathbf{1}(y_n = k)$$

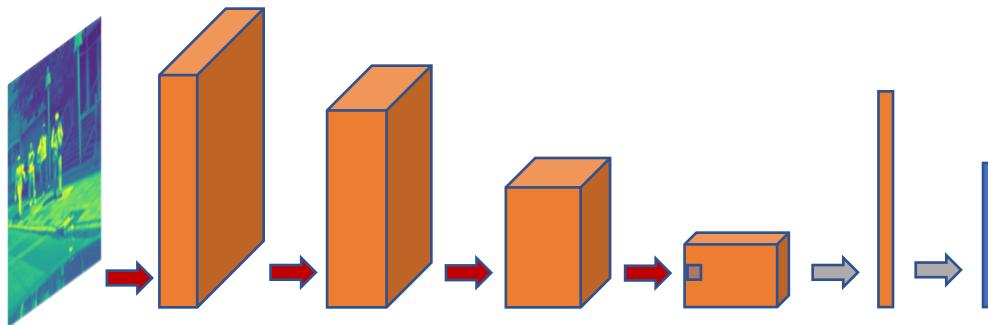


Convolutional layers

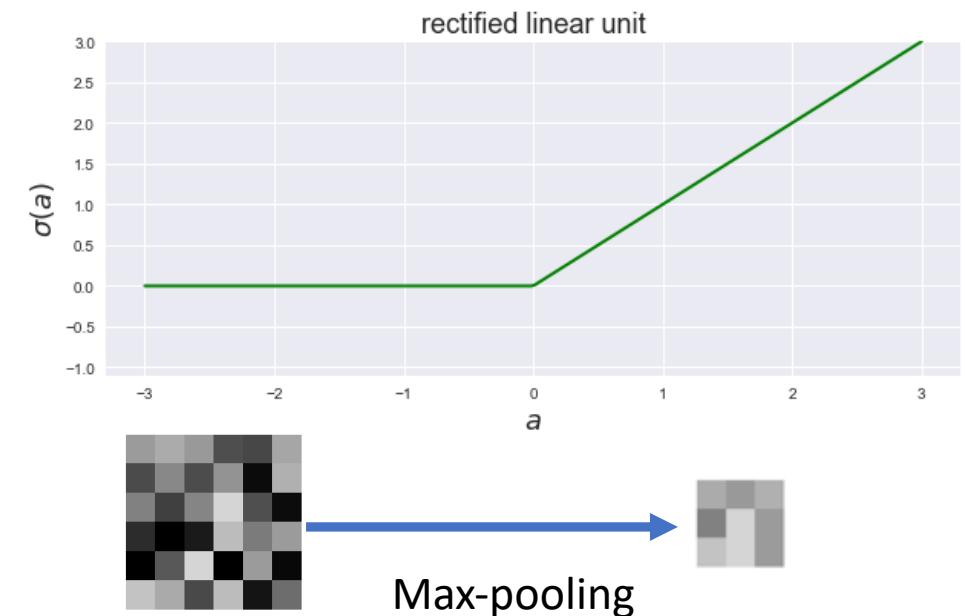
What input pattern activates a given neuron in an  
intermediate layer?

Fully  
connected  
layers

# Image dependent visualization

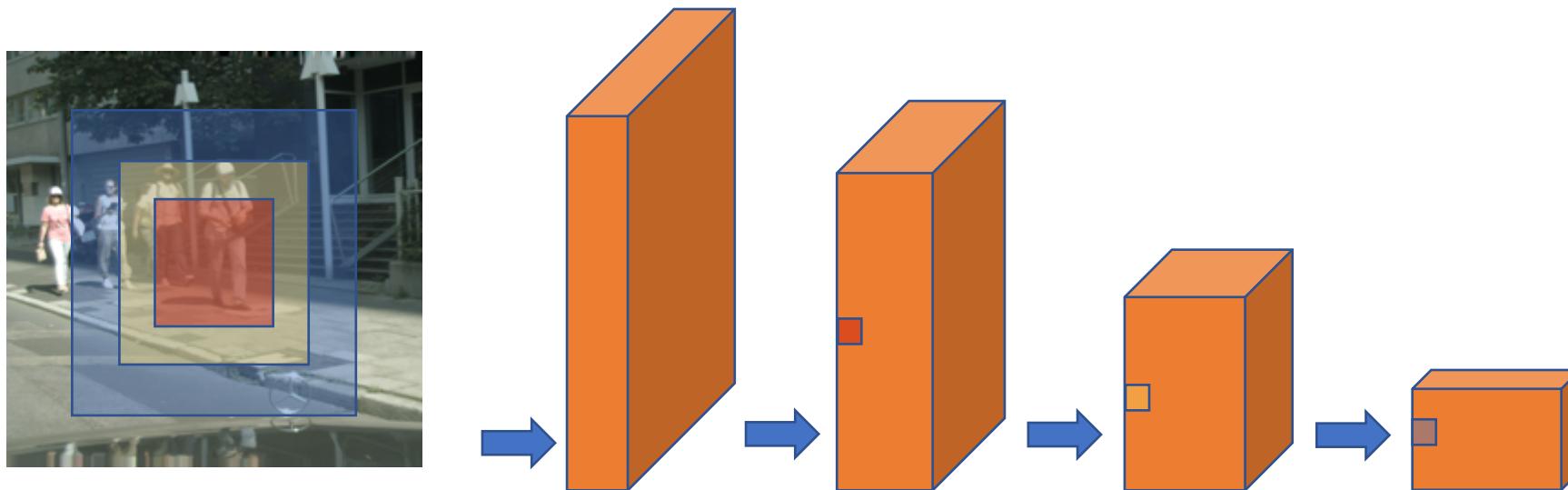


- The activation level in the neuron depends on the input image
- For different inputs it will be activated at different levels
- The difference is due to the non-linearity of the network (activation function + pooling)
- If it was linear, neuron's activation would be based on the respective linear projection



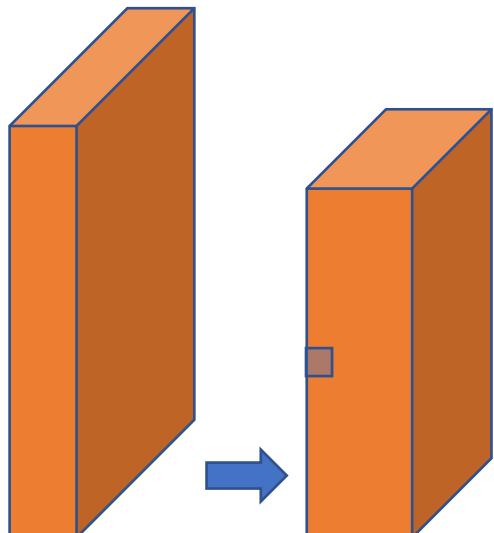
- Analysis should be based on the input image.
- The new question  
“In the input image which pattern caused the activation in a given neuron in an intermediate layer?”

# Size of the pattern in the input image



- Size of the input pattern changes with respect to the receptive field
- Depending on the layer the neuron sits, its receptive field changes
- The size of the input pattern changes as well.
- The image patch that activates the blue neuron is larger than the one that activates the red neuron

# Operations in the forward pass



- Consider the operations we need to do in order to compute the activation in the blue neuron from the layer below
- Three operations

- Convolutions with a set of filters

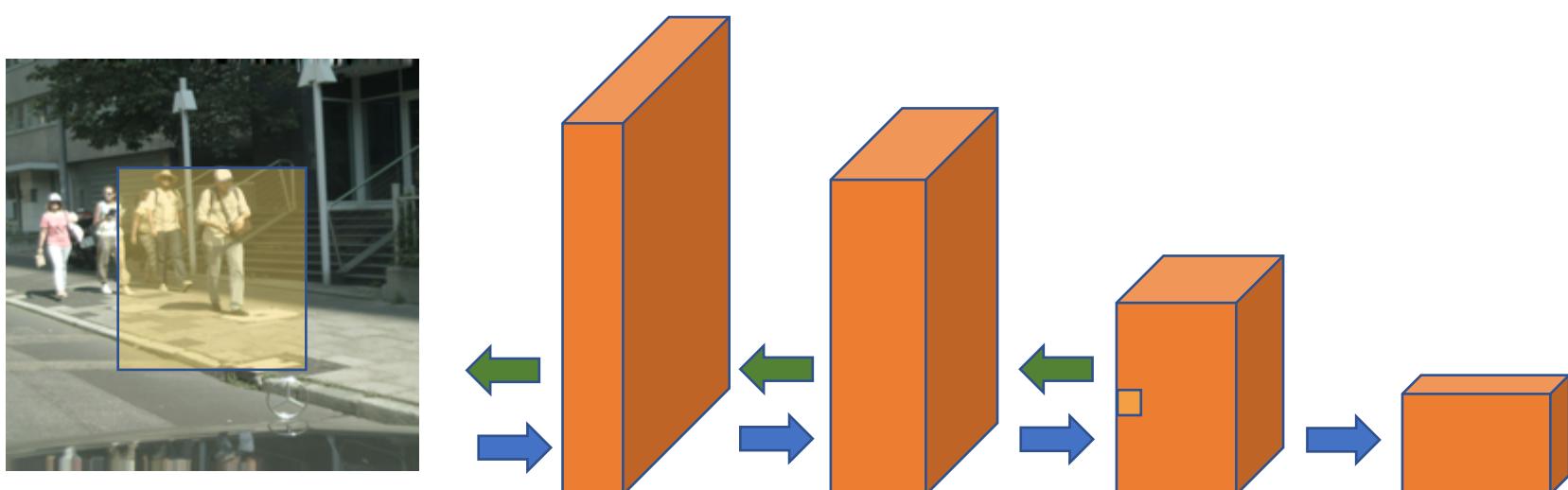
$$a_{l,k}^{(x,y)} = \left[ \sum_j w_{l,kj} * h_{l-1,j} + b_{l,k} \right]_{(x,y)}$$

- Non-linearity with ReLu function

$$h_{l,k}^{(x,y)} = \sigma \left( a_{l,k}^{(x,y)} \right)$$

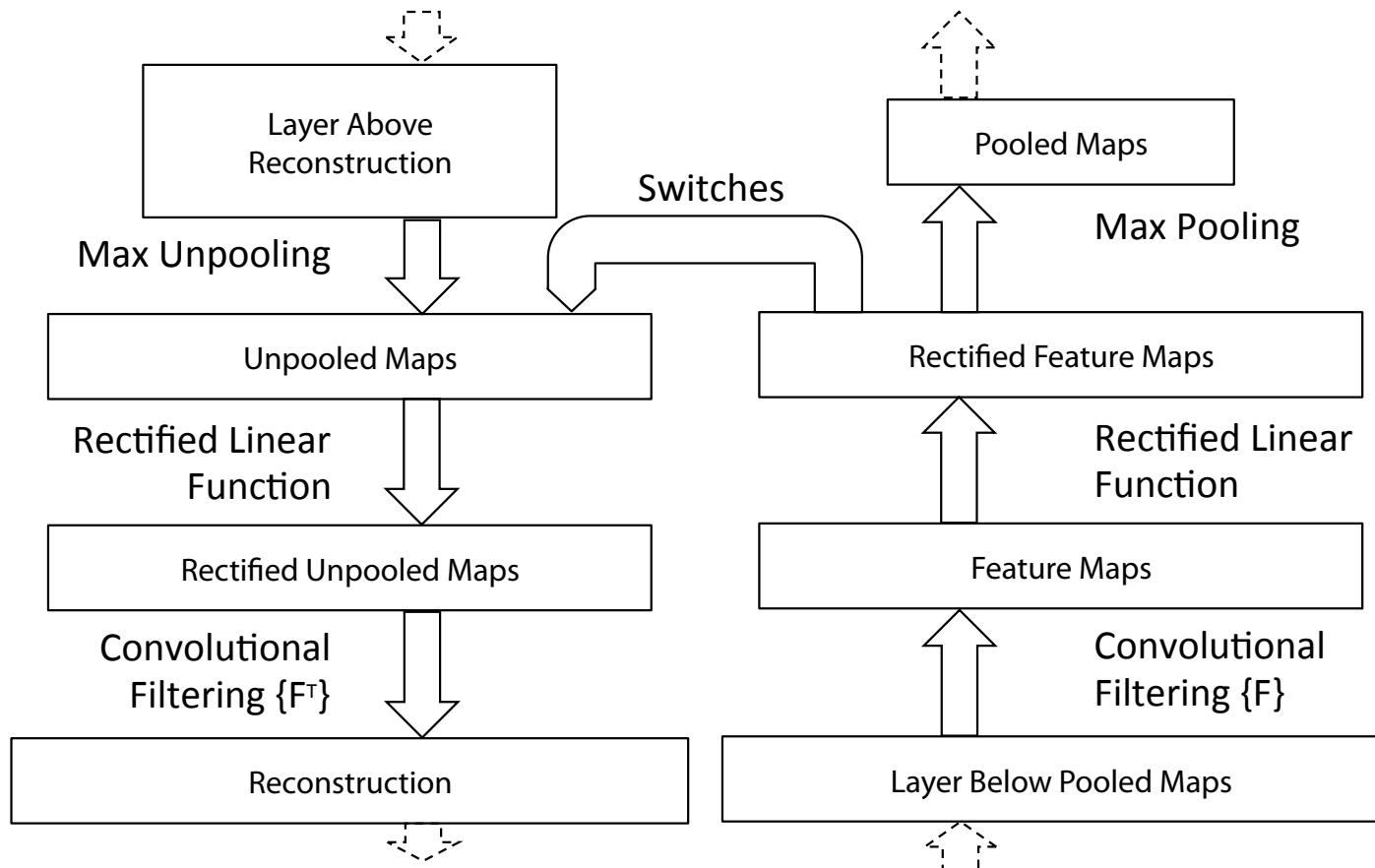
- Max-pooling

# The idea



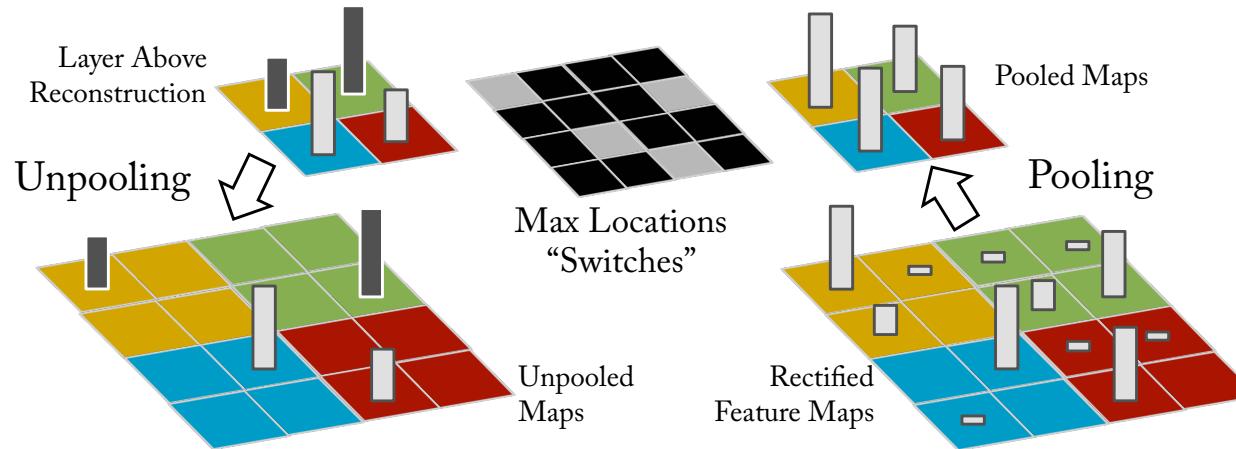
- Run an input image forward and compute all the features
- Keep the activation of the neuron you want and set the rest to 0.
- Starting from the same layer run the operations in reverse order.
- Unpool
- Rectify
- Transposed convolution
- A linked reverse "deconvolutional" network
- Modified layers are the inputs

# The idea



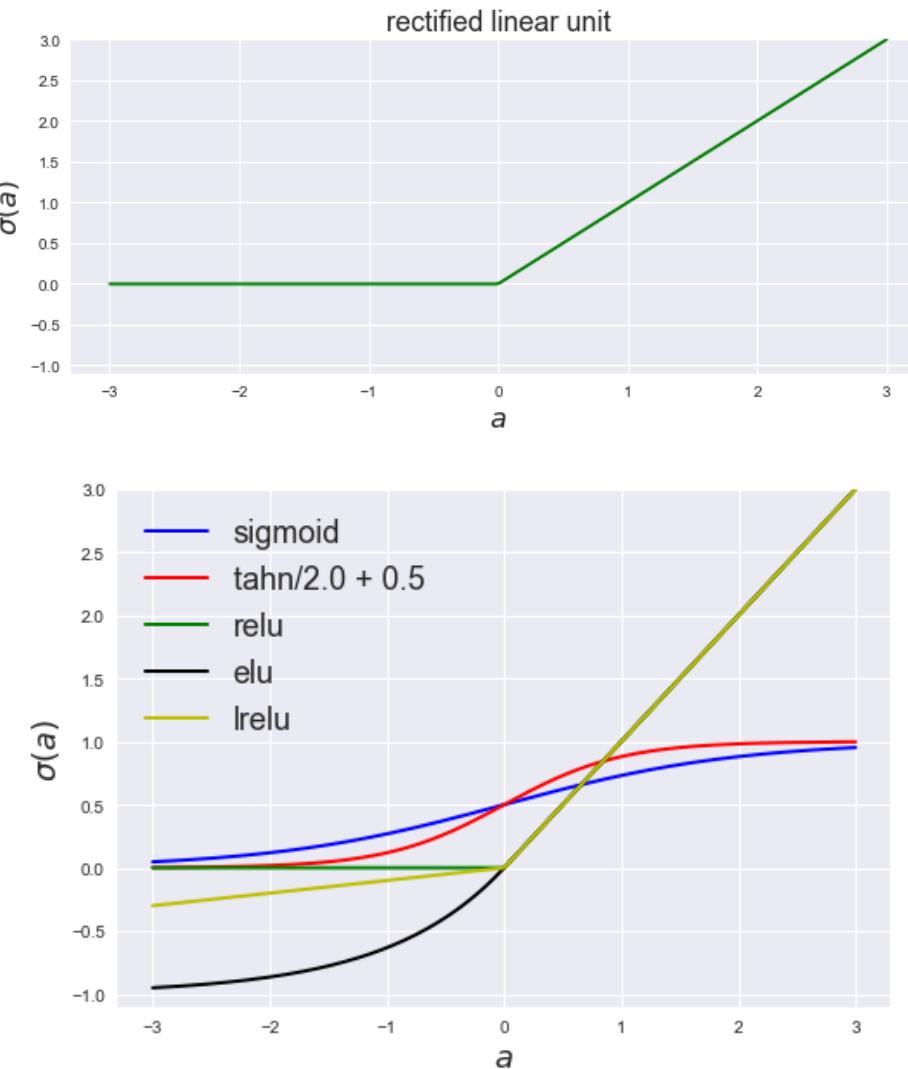
[Image from Zeiler and Fergus 2013]

# Inverting the operations – max pooling



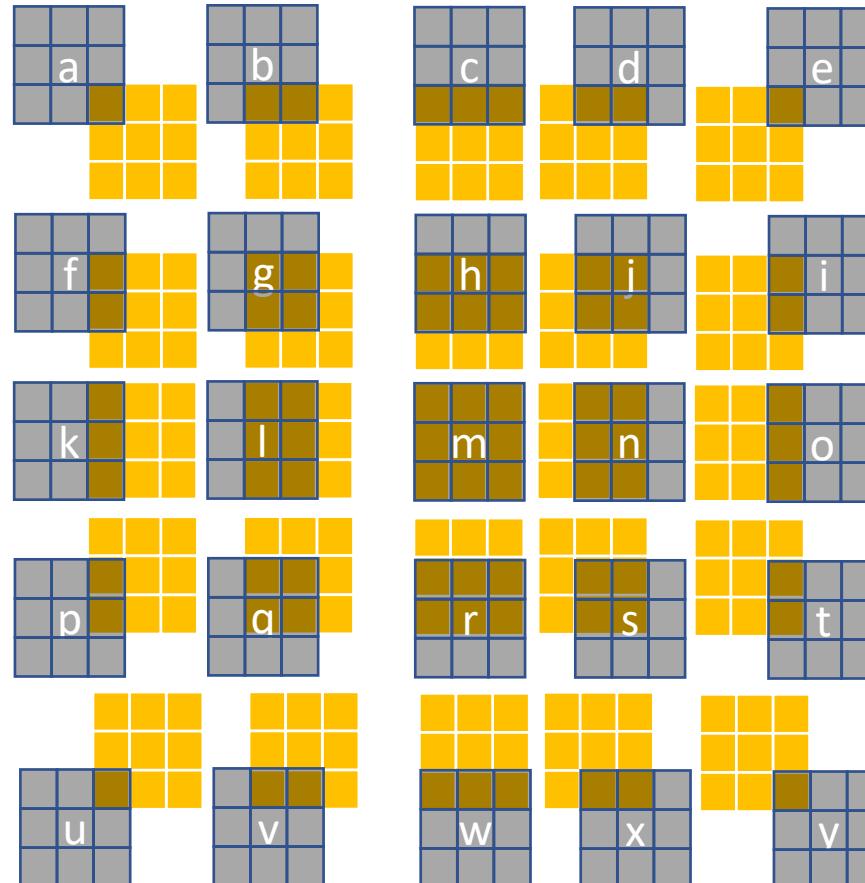
[Image from Zeiler and Fergus 2013]

- Keep the max locations while forward passing the image
- While *unpooling* place the values to the respective positions
- Pooling leads to information loss
- It is not possible to regenerate this information
- Instead zero values are placed for the locations where activations are discarded during forward pass

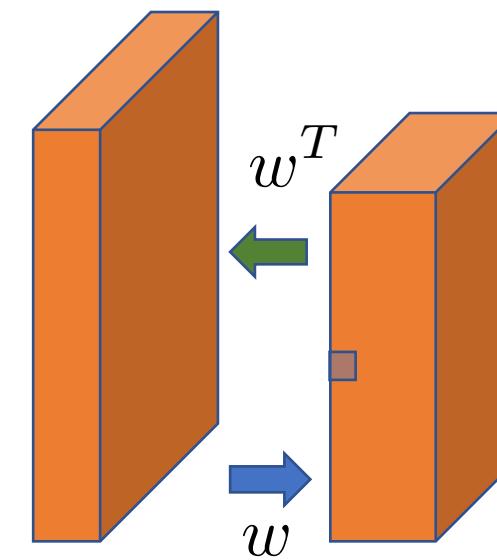


- Only keeps the positive layers
- The same function is used for the reverse operation
- ReLu yields only positive activation maps
- To keep the activation maps the same, ReLu is used again to keep the activations during reconstruction positive
- You can in theory, also use the inverse of a function

# Inverting the operations - Convolution



- Transposed convolution
- The kernel is the kernel used in the forward pass, flipped horizontally and vertically

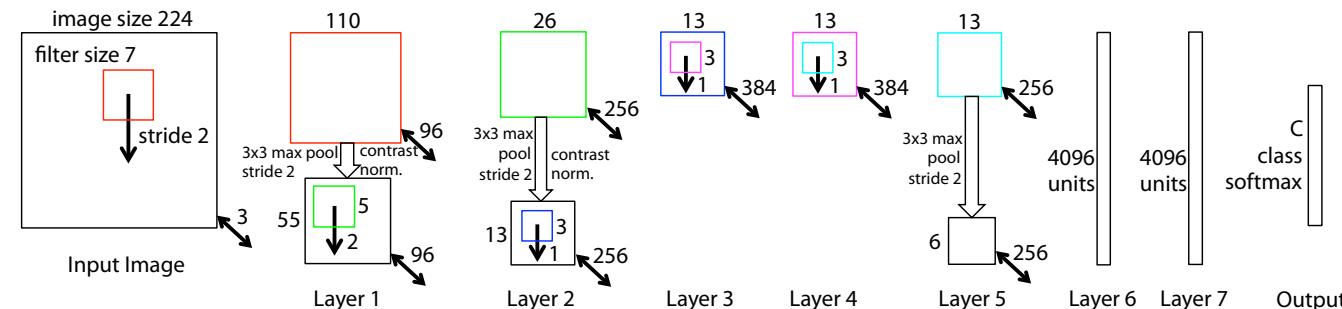
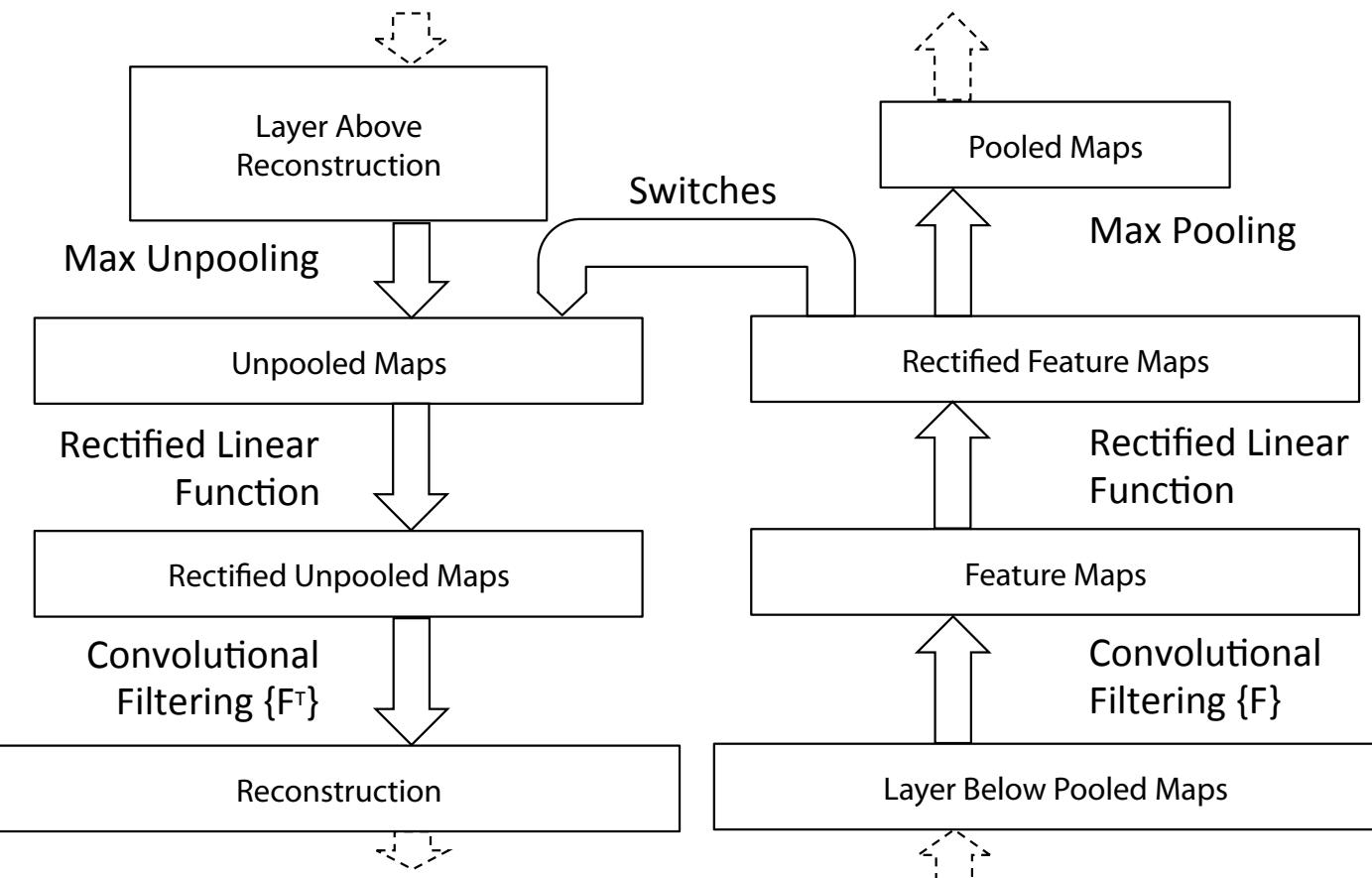


# Computer Vision

Detection

Segmentation

Visualization



[Images from Zeiler and Fergus 2013]

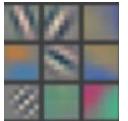
# Computer Vision

Detection

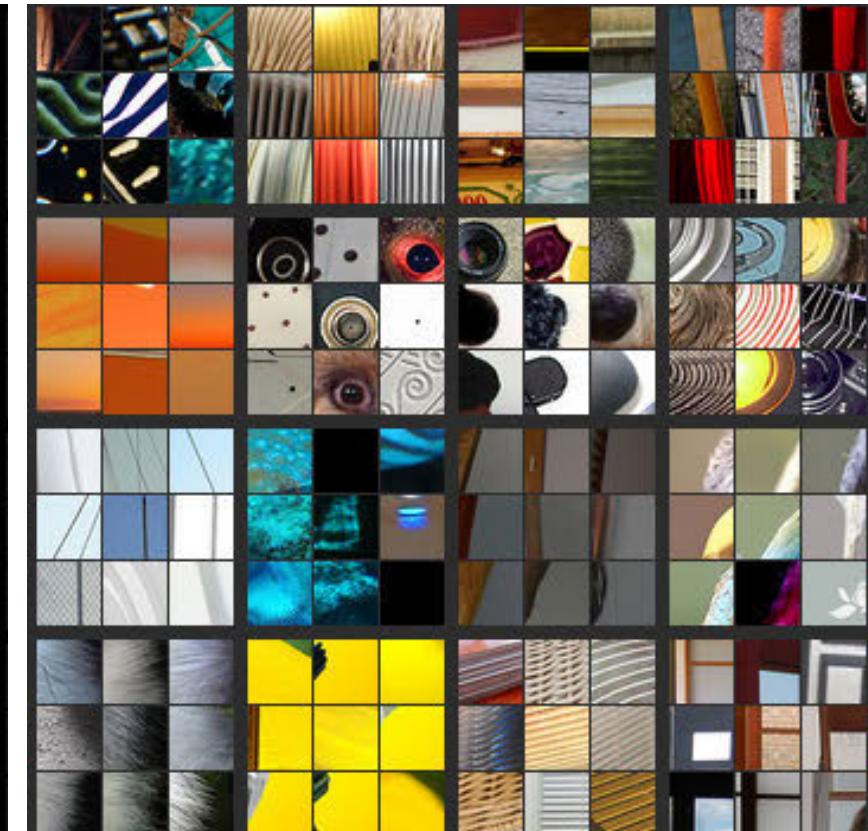
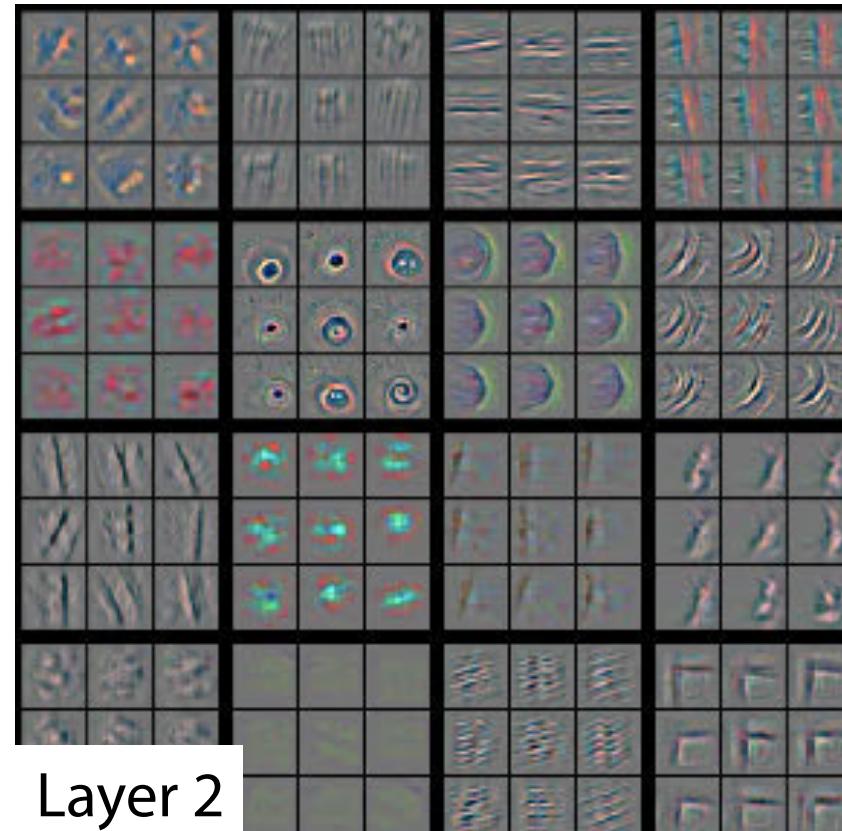
Segmentation

Visualization

## Feature maps

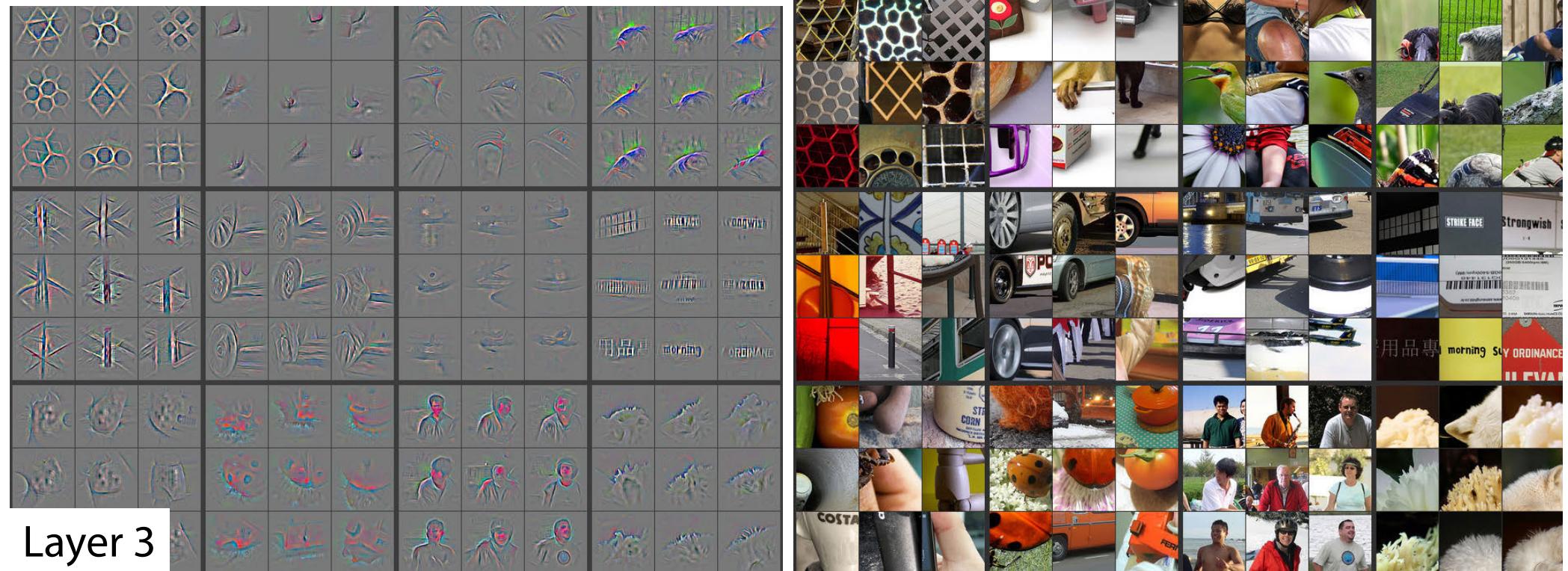


Layer 1



[Images from Zeiler and Fergus 2013]

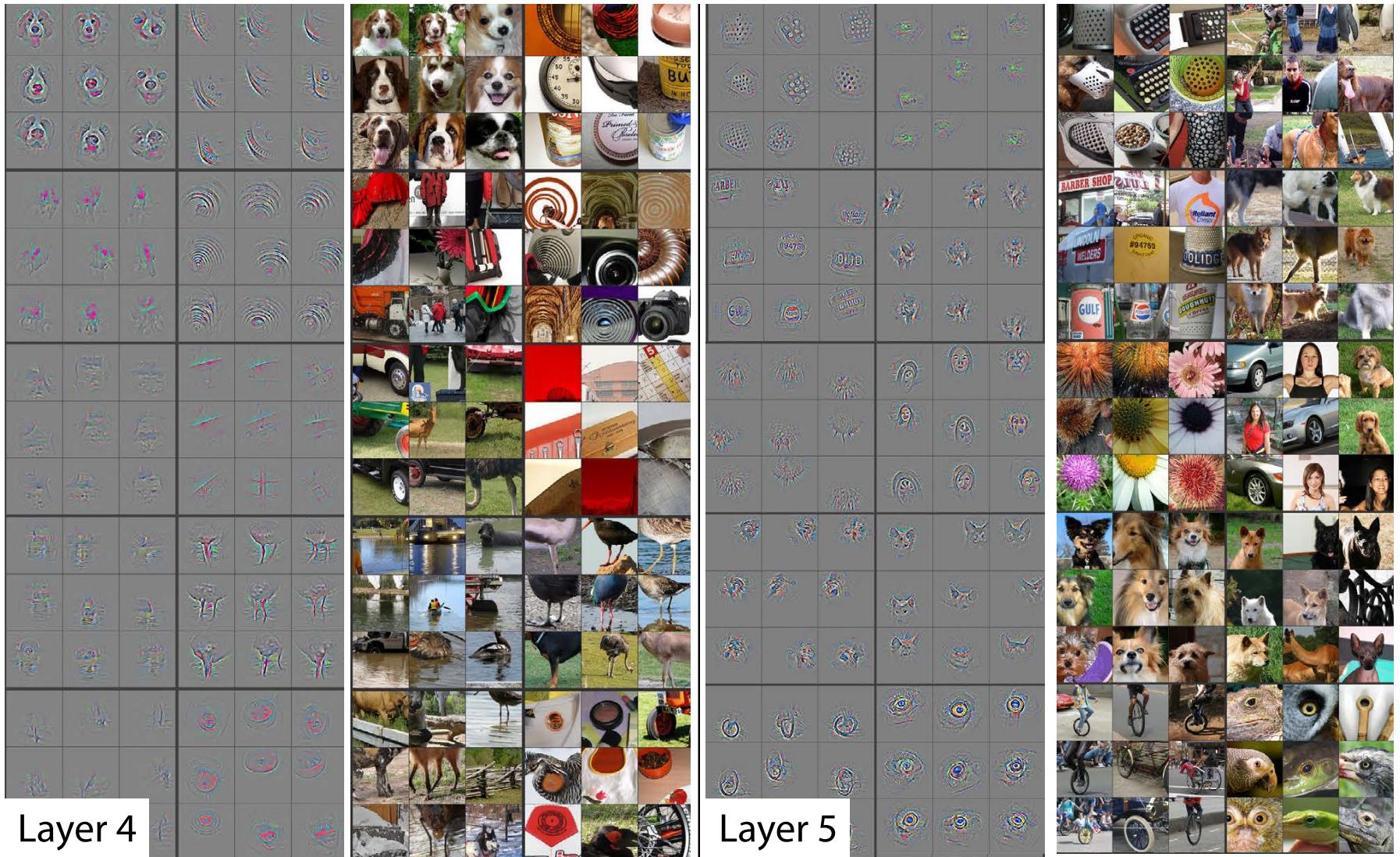
# Further deeper features



[Images from Zeiler and Fergus 2013]

[Images from Zeiler and Fergus 2013]

# Even further deeper





# Localization and classification

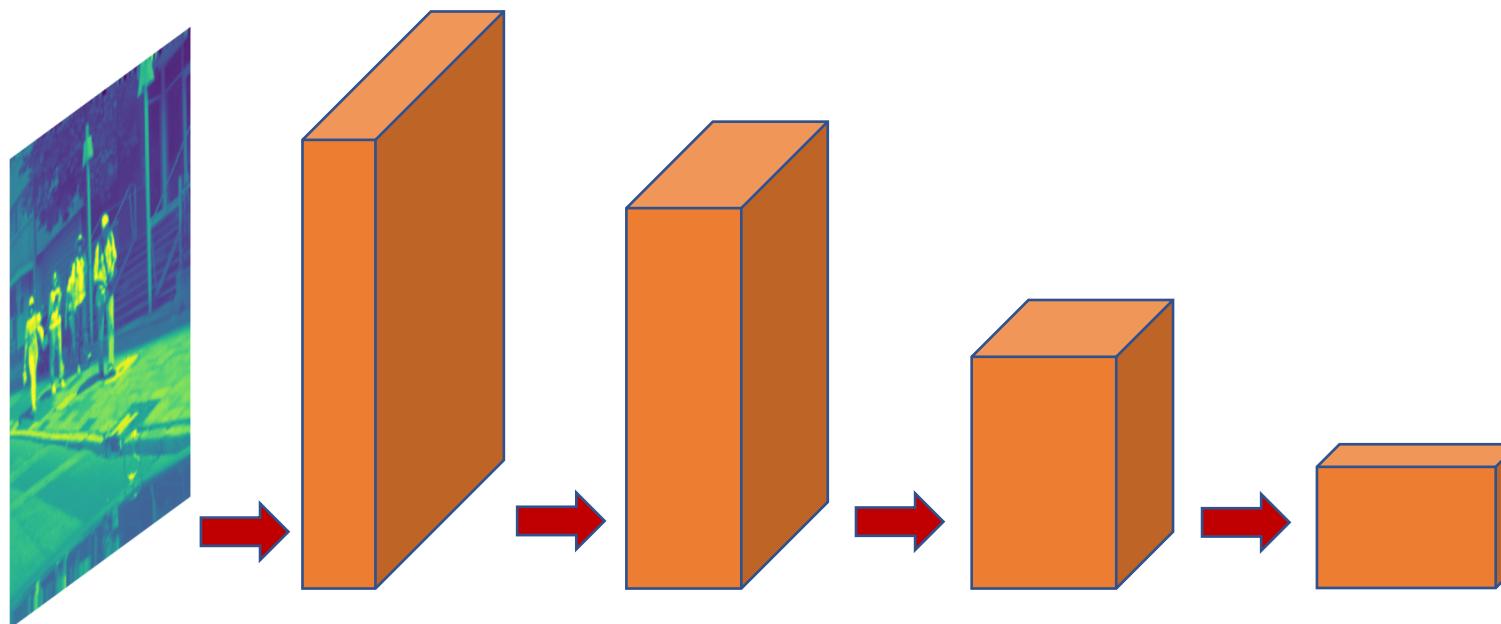
# Class activation maps

- So far for classification we were only interested in determining the class assignment
- We also had a separate localization network that relied on separate classification tasks at proposal regions
- With slight modifications classification networks can identify approximate locations
- Based on global average pooling idea
- Discussion based on [Zhou et al. 2016]

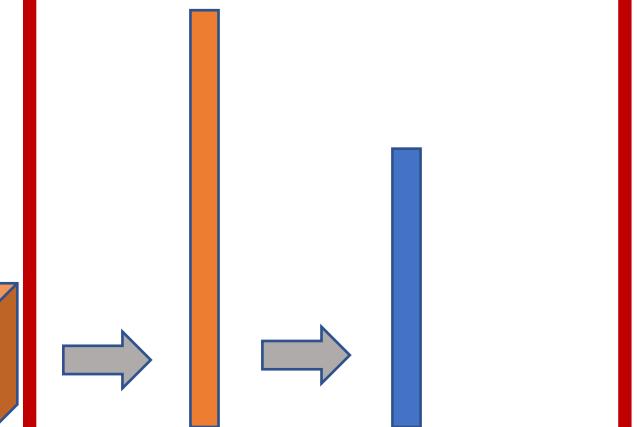
## **Learning Deep Features for Discriminative Localization**

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba  
Computer Science and Artificial Intelligence Laboratory, MIT  
`{bzhou, khosla, agata, oliva, torralba}@csail.mit.edu`

# Normal classification network



Lastly fully connected  
layers summarizing the  
feature maps



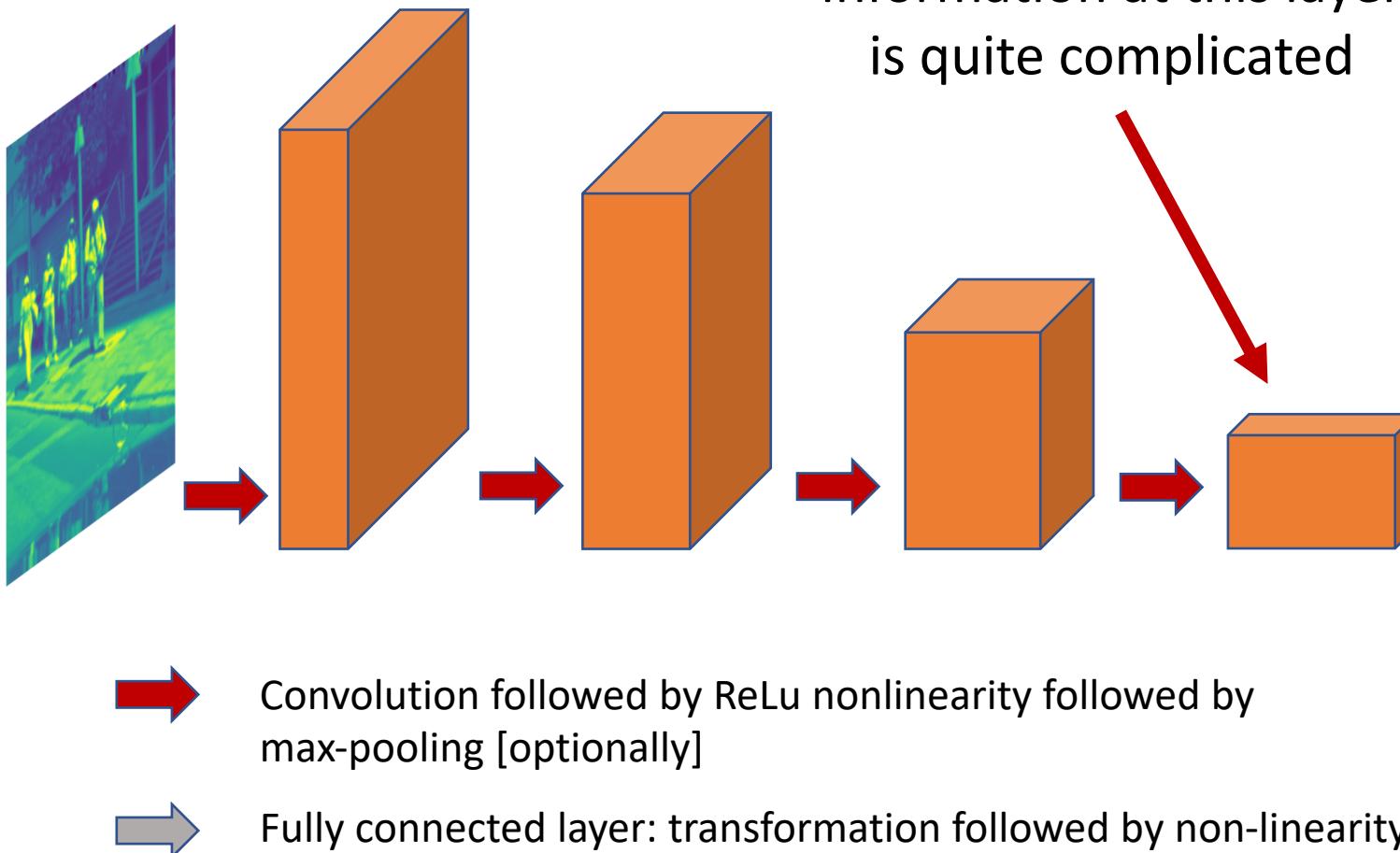
Output

Number of neurons  
equal number of classes

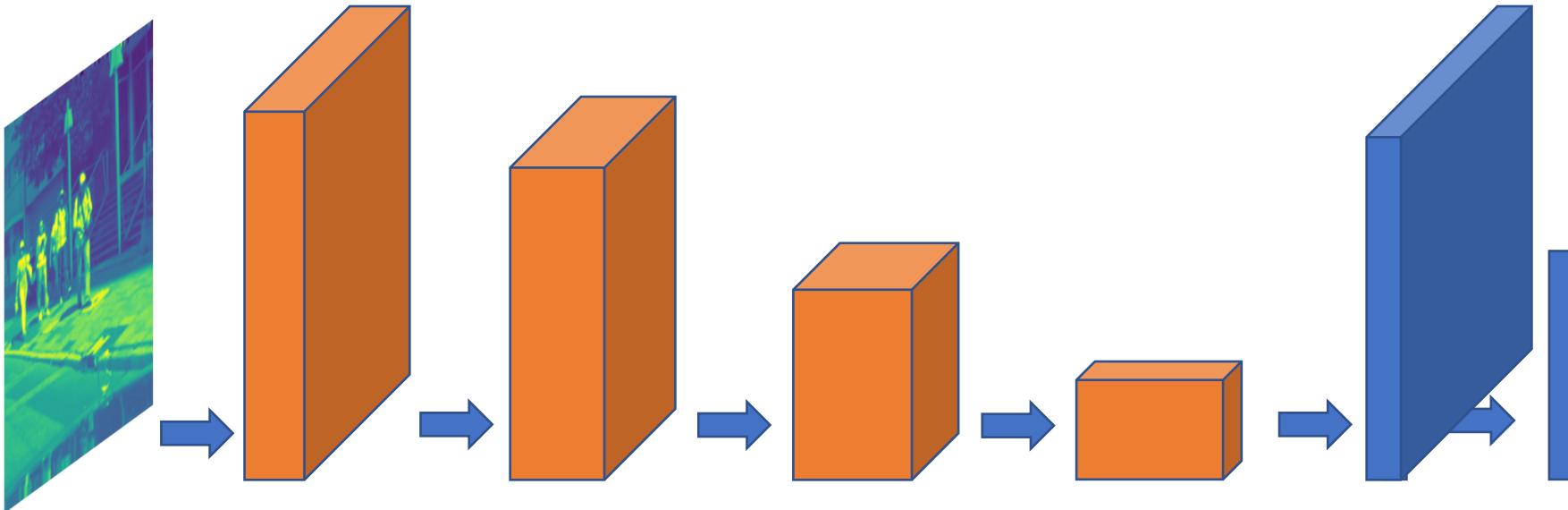
→ Convolution followed by ReLu nonlinearity followed by  
max-pooling [optionally]

→ Fully connected layer: transformation followed by non-linearity

# Normal classification network



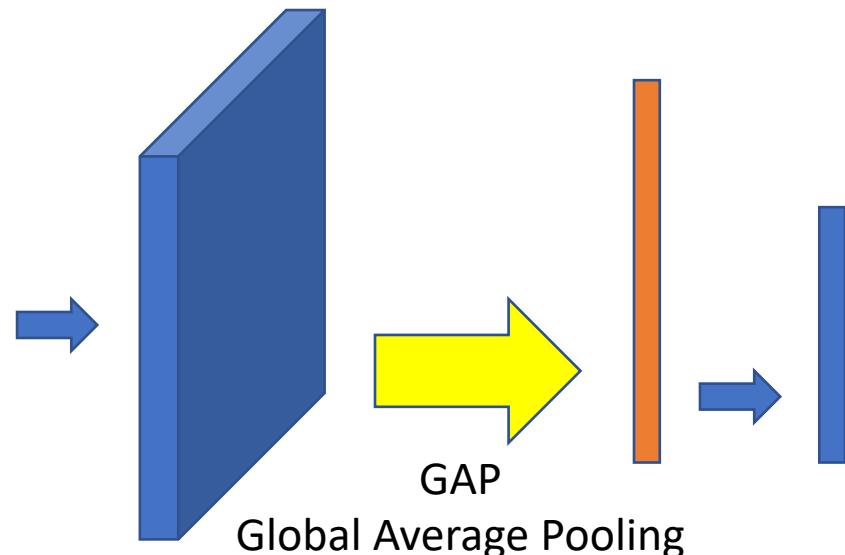
We have also seen fully convolutional networks for segmentation



- Image size outputs
- Replaced the final fully connected layers
- Upsampling using transpose convolutions or bilinear upsampling followed by convolutions

# Combining these ideas

- Class activation maps combines these ideas
- Using Global Average Pooling
  - Like normal pooling, applies to each channel in a layer separately

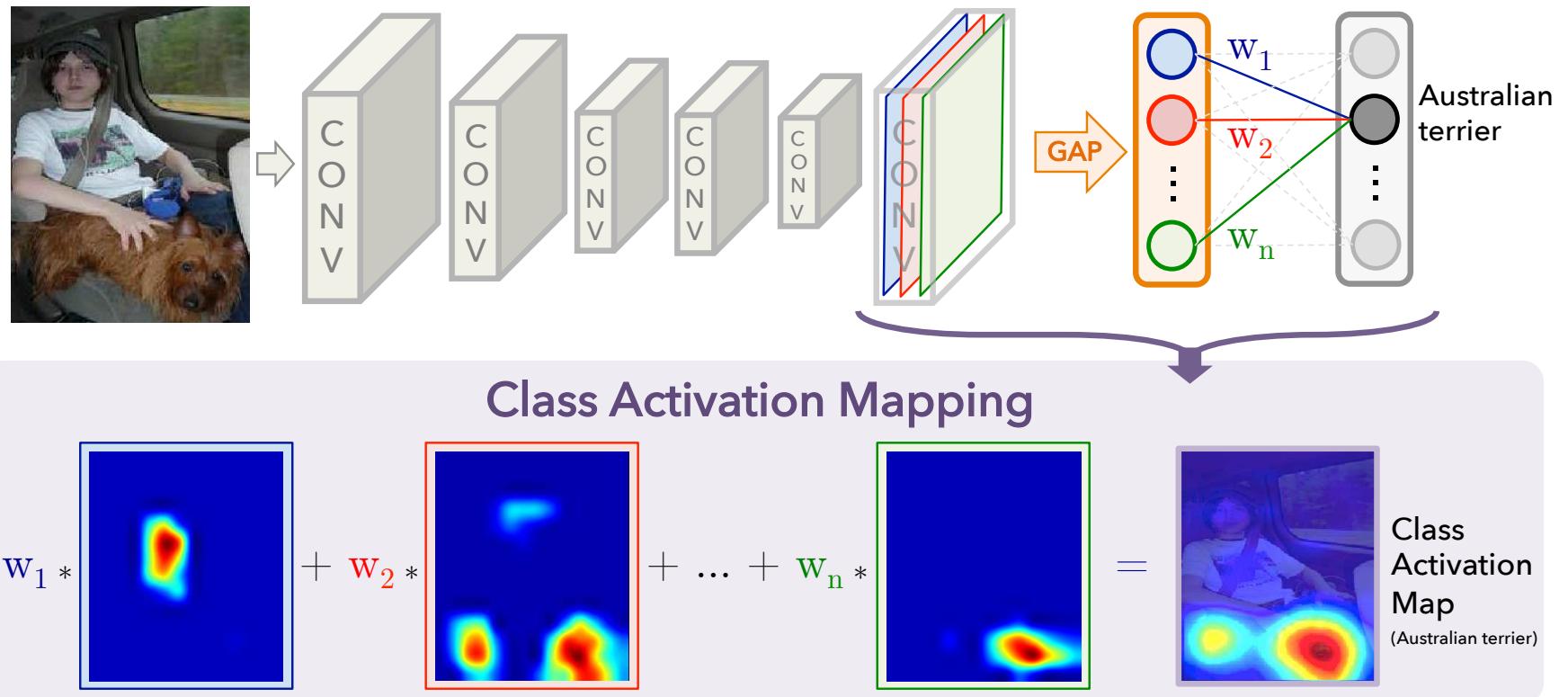


$$h_{l,k} = \frac{1}{M_{l-1}N_{l-1}} \sum_{x,y} h_{l-1,k}^{(x,y)}$$

- Averaging all the information to a single number!
- Then continue as usual

$$a_{L,j} = \sum_k w_{L,jk} h_{l,k}$$

# Activation Maps



- Per-class weighted sum of all the channels before global average pooling yields the class-specific activation map

[Image taken from Zhou et al. 2016]

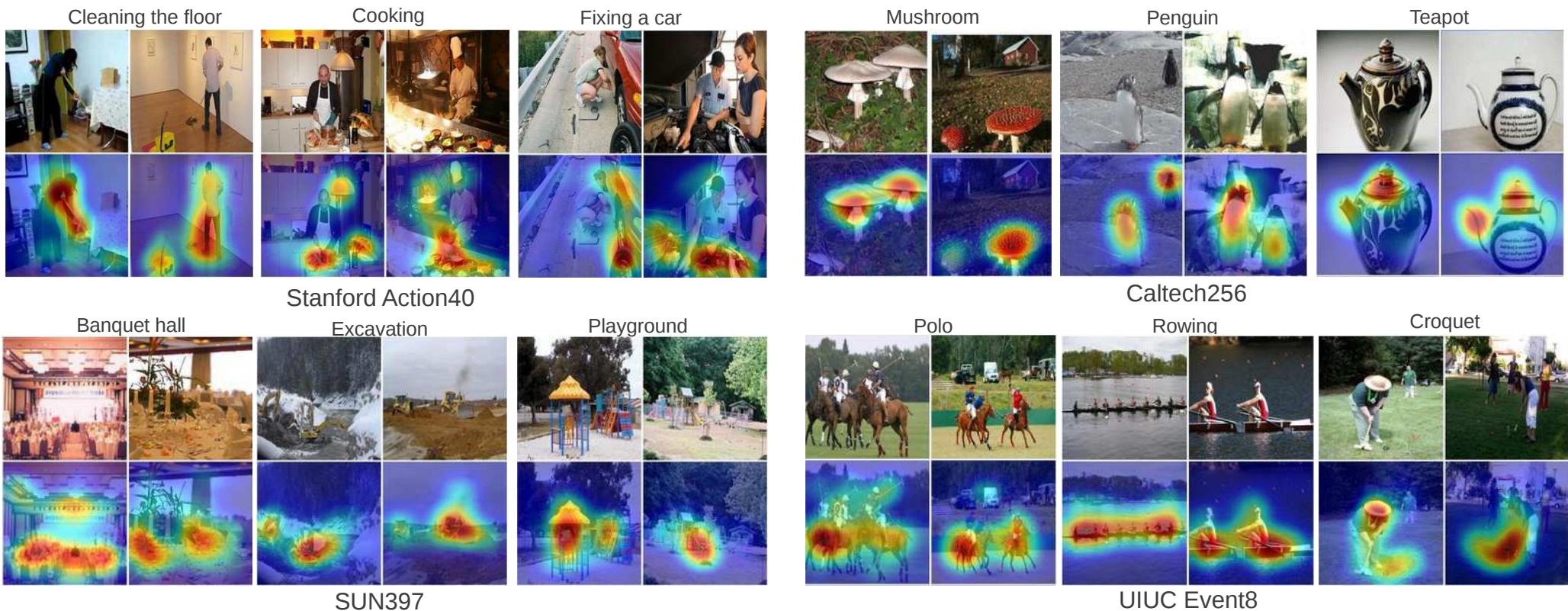
# Computer Vision

Detection

Segmentation

Visualization

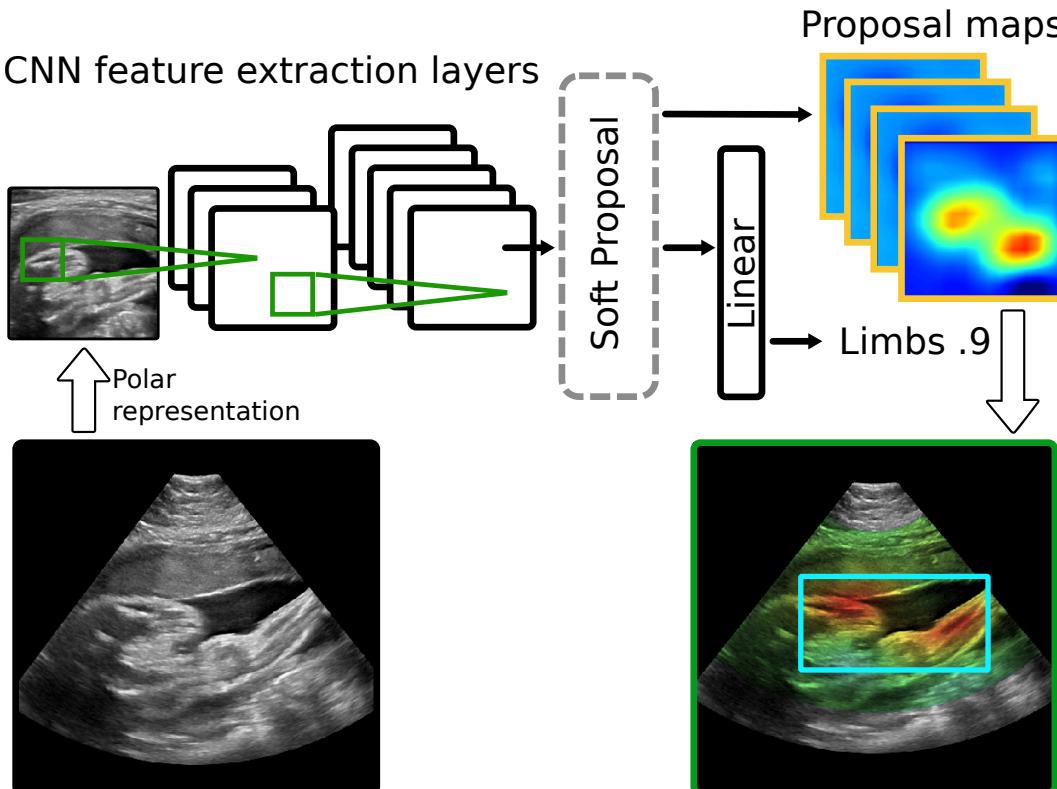
Localization



- Network architecture preceding the GAP layer can change
- Form of weak-supervision for localization

[Image taken from Zhou et al. 2016]

# Various applications



Weakly Supervised Localisation for Fetal  
Ultrasound Images

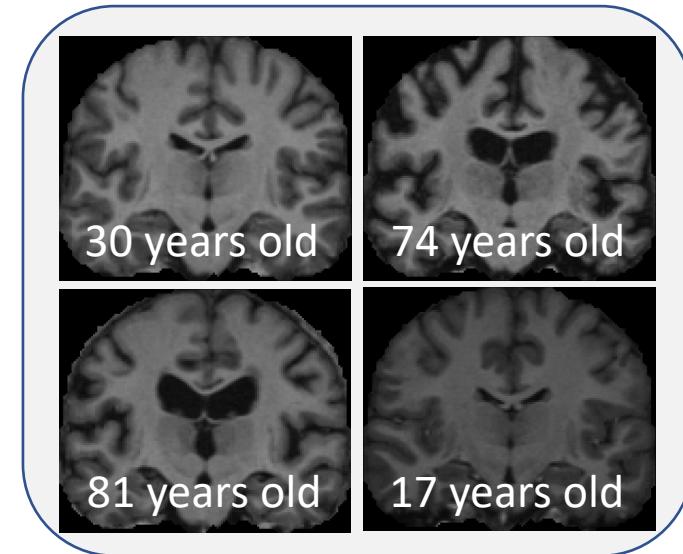
- Especially in medical imaging
- Labels are expensive and difficult to get.
- Approximate localization with CAM allow identifying areas of interest
- Also weak supervision to train stronger localization algorithms

# Unsupervised learning

# Very coarse view on supervised learning

## Supervised learning

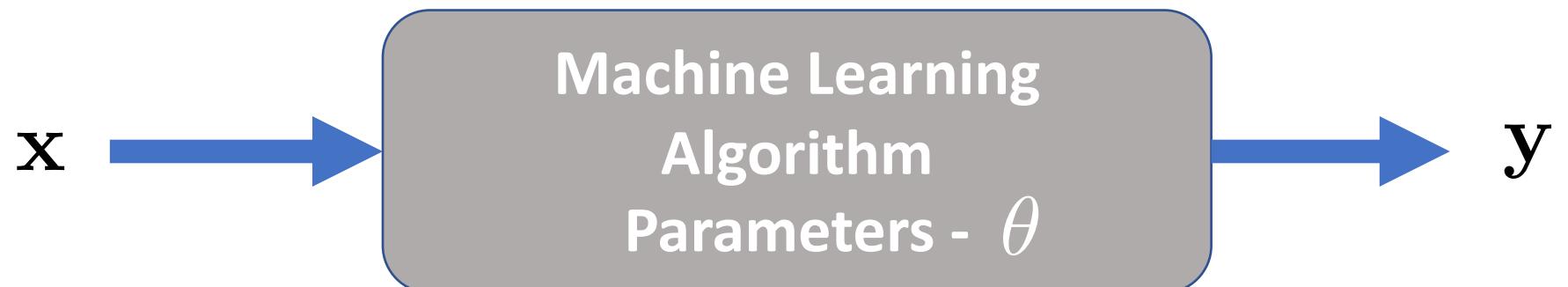
- Patterns between two types of data
- Goal: predicting one from the other
- Examples have both types of data
- At prediction only one exist



30 years old

# General idea in the supervised approach

- Algorithms assume a mathematical model between features and labels



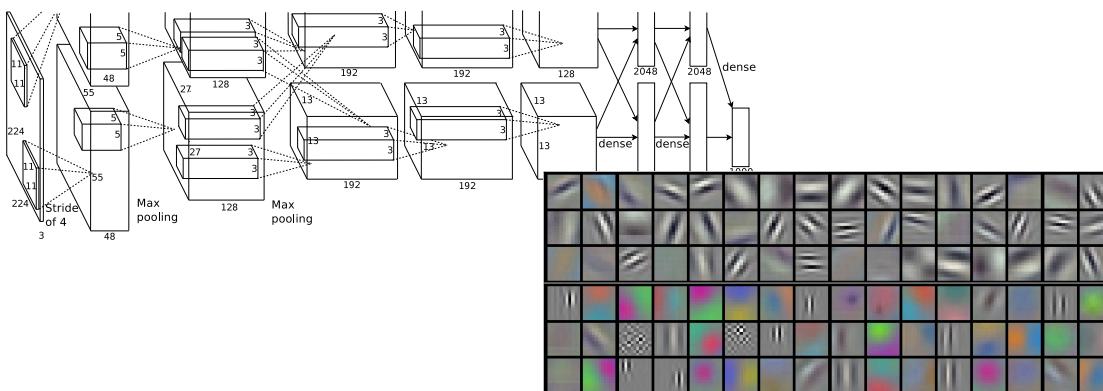
- Estimate the parameters of the model to best predict labels from features in the training examples

$$y = f(\mathbf{x}|\theta)$$

# Unsupervised learning

## Unsupervised learning of features

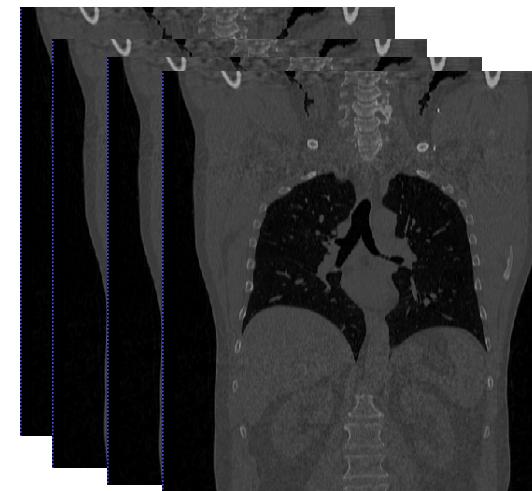
- Filters are important for performing image analysis tasks
- So far, we determine features in a supervised way, task-specific manner



- Determine features in an unsupervised manner
- Examples have only features

## Unsupervised learning of distributions

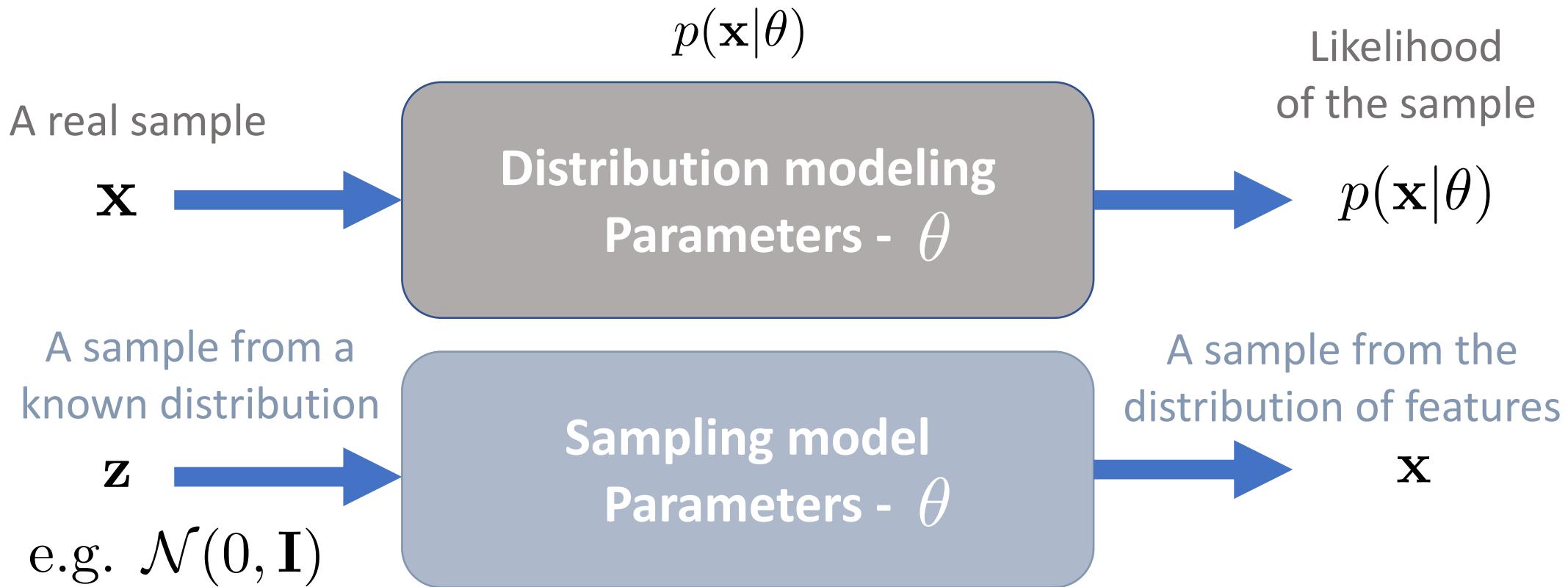
- Patterns within the data
- Goal: describe variability in the data
- Estimate the distribution of the data
- There is still a training dataset
- Examples have only features



Both are unsupervised in the sense that there are no labels!

# General idea in unsupervised distribution learning

- Algorithm assumes a mathematical model for the features
- Ideally this is the probability distribution of features



# Why is this useful?

- Sample from the distribution of image to generate images



Figure 1: Class-conditional samples generated by our model.

[Figure from Brock, Donahue and Simonyan 2018 – Class conditional generation of images]

# Computer Vision

Detection

Segmentation

Visualization

Localization

Unsupervised  
learning

## Why is this useful?

- Style transfer



[Figure from Karras, Laine and Aila, 2018]

# Computer Vision

Detection

Segmentation

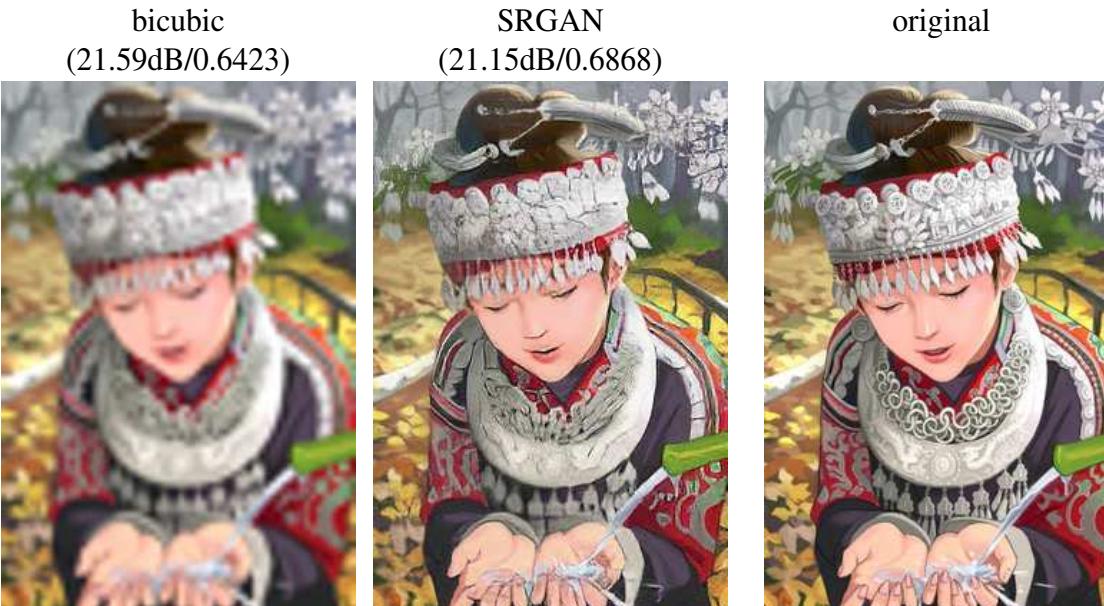
Visualization

Localization

Unsupervised  
learning

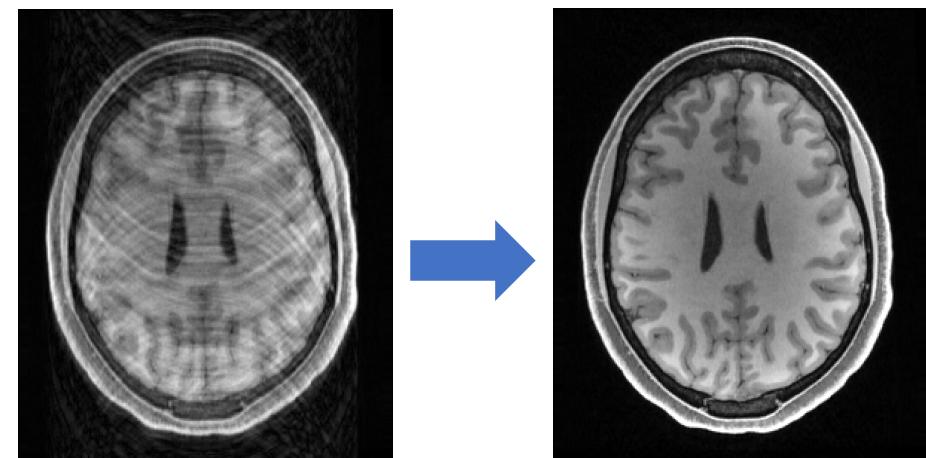
## Why is this useful?

### Improving resolution of an image



[Figure from Ledig et al. 2017]

### Bayesian reconstruction of medical images



[Figure from Tezcan et al. 2018]

## Why is this useful?

- Many more applications:
  - In-painting
  - Realistic video and image editing
  - Video frame prediction
  - Outlier detection
  - ...
- Scientifically
  - Building a model of the visual world
  - Possibly an important component in human learning.
    - We do not see 100s of cups to understand what a cup is
    - We constantly observe around and get visual input to our brains.

# Images are big



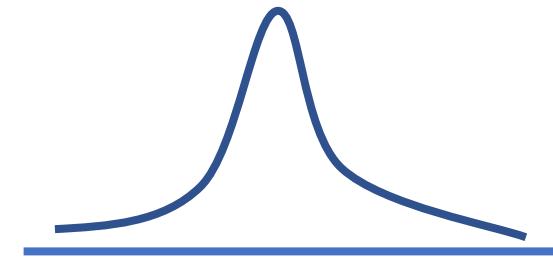
- Images are very high dimensional
- Consider a small image of 64x64
- Even that is 4096 dimensional!
- We need to keep that in mind when we think about unsupervised learning.

# The most straightforward way

- Kernel density estimation (KDE)
- Given a sample set of images the naïve way is

$$p(x|\theta) = \frac{1}{N} \sum_n K_\theta(x, x_n)$$

$$\int K_\theta(x, x_n) dx = 1$$

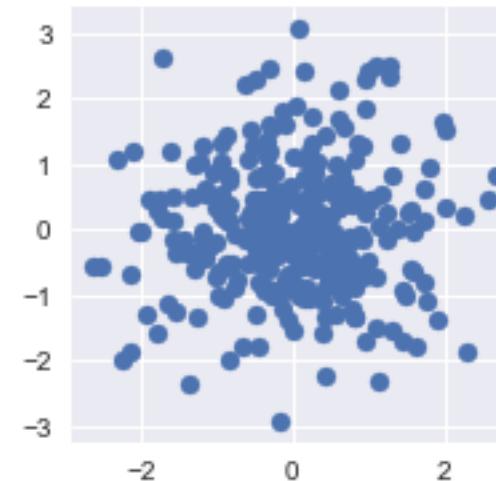


- Place a "kernel" around each training sample
- Determine the likelihood of a new sample based on these kernels
- If kernels depends on Euclidean distance, e.g. Gaussian kernel, then likelihood is related to the distance in Euclidean space.

# Bad idea due to the dimensions



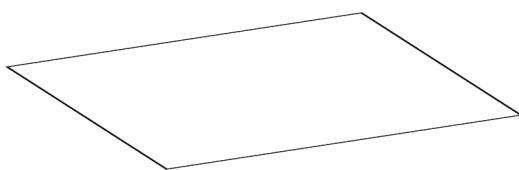
- For the KDE to work, roughly speaking, you need to somehow “fill” the space, e.g.



- To fill a space of 4096 dimensions, you need a lot of samples, we need to find a better solution.

# Latent variable models

- Assume that images live in a lower dimensional sub-space
- We build a mapping between them



latent space  
 $z \in \mathbb{R}^d$

$$p(x|z)$$

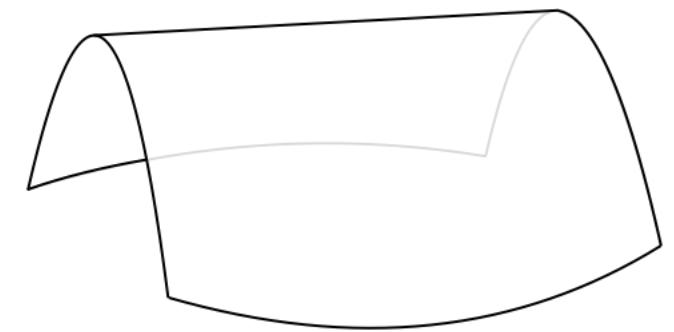


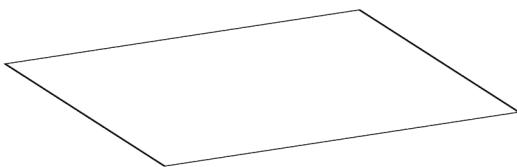
image space  
 $x \in \mathcal{M} \subset \mathbb{R}^D$

$$p(x) = \int p(x|z)p(z)dz$$

$$x \in \mathbb{R}^D, z \in \mathbb{R}^d, d \ll D$$

# Probabilistic principal component analysis

- Assumes the mapping is a linear one
- Probabilistic principal component analysis [Tipping & Bishop 1999]



latent space  
 $z \in \mathbb{R}^d$



$$p(x|z)$$

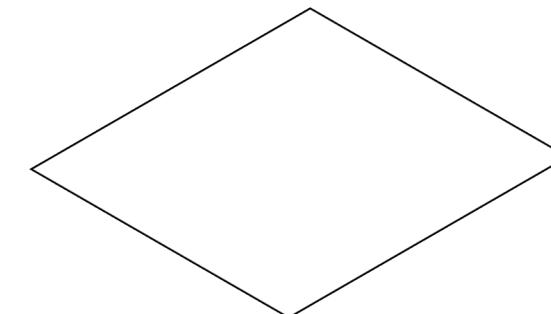


image space  
 $x \in \mathcal{M} \subset \mathbb{R}^D$

$$p(x) = \int p(x|z)p(z)dz$$

$$x \in \mathbb{R}^D, z \in \mathbb{R}^d, d \ll D$$

$$p(x|z) = \mathcal{N}(Wz + \mu_x, \sigma^2)$$

$$p(z) = \mathcal{N}(0, I) \quad \mu_x : \text{mean image}$$
$$W \in \mathbb{R}^{D \times d} \quad \sigma^2 : \text{noise}$$

## Link to PCA

- Maximum likelihood estimate of the parameters yield the PCA
- Eigenvalues and eigenvectors of the sample covariance matrix
- Derivation in [Tipping and Bishop 1999]

$$W_{ML} = U_d (\Lambda_d - \sigma^2 I)^{1/2} R$$

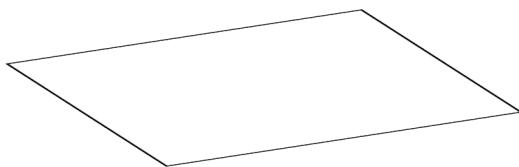
$$W_{ML} = \underbrace{U_d}_{d \text{ eigenvectors}} \begin{pmatrix} \underbrace{\Lambda_d}_{d \text{ eigenvalues}} & -\sigma^2 I \end{pmatrix}^{1/2} \underbrace{R}_{\text{arbitrary rotation}}$$

$\mu_x$  : sample mean image

$$\sigma_{ML}^2 = \frac{1}{D-d} \sum_{q=d+1}^D \lambda_q$$

# Non-linear maps

- In supervised learning linear maps were not enough
- The same idea applies here



latent space  
 $z \in \mathbb{R}^d$

$$p(x|z)$$

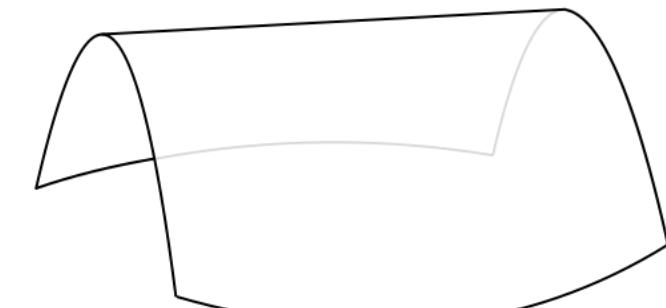


image space  
 $x \in \mathcal{M} \subset \mathbb{R}^D$

$$p(x) = \int p(x|z)p(z)dz$$

$$x \in \mathbb{R}^D, z \in \mathbb{R}^d, d \ll D$$

$$p(x|z) = \mathcal{N}(f(z|\theta) + \mu_x, \sigma^2)$$

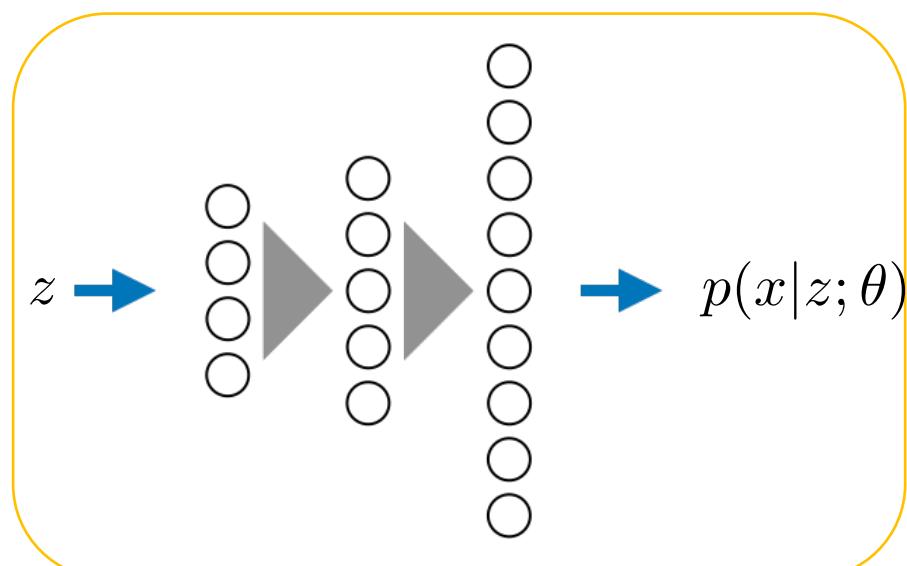
$$p(z) = \mathcal{N}(0, I) \quad \mu_x : \text{mean image}$$
$$\sigma^2 : \text{noise}$$

# Density networks

[MacKay, Nucl. Inst. Met. In Physics Research 1995]

$p(x|z; \theta)$  : Parameterize with a network with parameters  $\theta$

$$p(x; \theta) = \int p(x|z; \theta)p(z)dz$$



For the given samples, maximize with respect to  $\theta$

$$\prod_n p(x_n; \theta) = \prod_n \int p(x_n|z; \theta)p(z)dz$$

using Monte Carlo integration

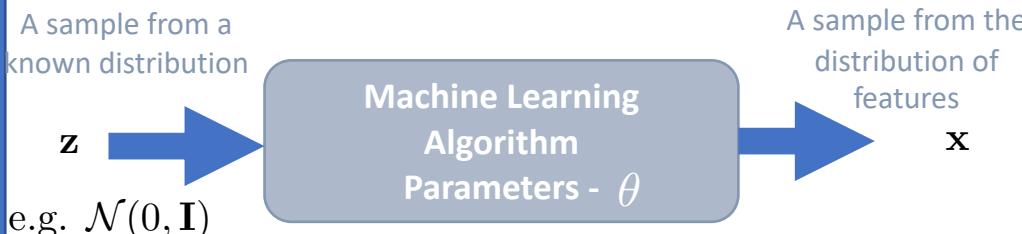
$$\sum_n \ln \frac{1}{R} \sum_r p(x_n|z_r; \theta), \quad z_r \sim p(z)$$

Sampling was not efficient for very large dimensional problems, need too many samples  
MacKay hinted importance sampling

# Two avenues – both end of 2013

## Generative Adversarial Network

### Sampler

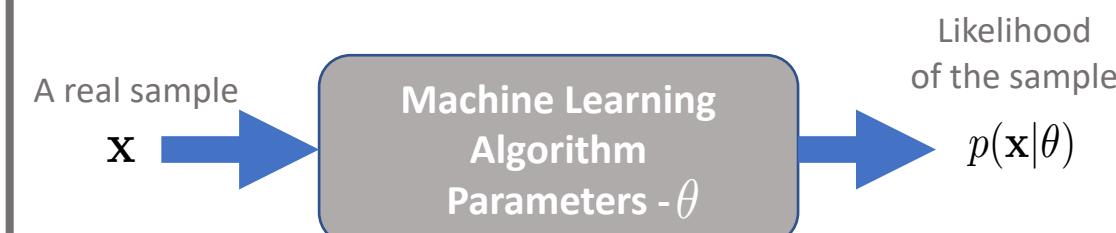


## Generative Adversarial Nets

Ian J. Goodfellow\*, Jean Pouget-Abadie†, Mehdi Mirza, Bing Xu, David Warde-Farley,  
Sherjil Ozair†, Aaron Courville, Yoshua Bengio§  
Département d'informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC H3C 3J7

## Variational Auto-encoders

### Distributional model



## Auto-Encoding Variational Bayes

Diederik P. Kingma  
Machine Learning Group  
Universiteit van Amsterdam  
dpkingma@gmail.com

Max Welling  
Machine Learning Group  
Universiteit van Amsterdam  
welling.max@gmail.com

## Stochastic Backpropagation and Approximate Inference in Deep Generative Models

Danilo J. Rezende, Shakir Mohamed, Daan Wierstra  
{danilor, shakir, daanw}@google.com  
Google DeepMind, London

# Variational auto-encoders

Builds on density networks concept but instead of Monte-Carlo uses variational inference with a network parameterized sampling (approximate) distribution

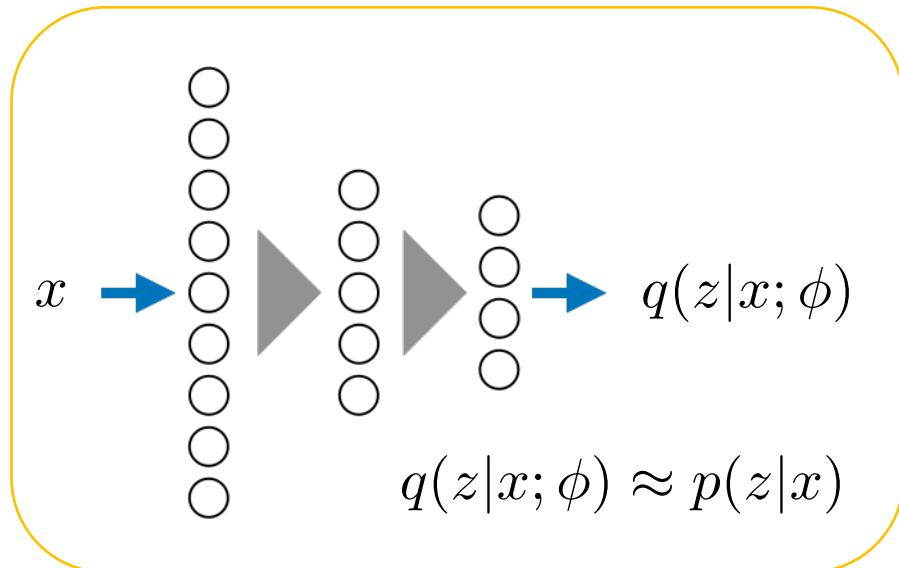
$$\begin{aligned}\ln p(x; \theta) &= \ln \int p(x|z; \theta)p(z)dz \\ &= \ln \int p(x|z; \theta)p(z) \frac{q(z|x; \phi)}{q(z|x; \phi)} dz \\ &\geq \int q(z|x; \phi) \ln p(x|z; \theta) \frac{p(z)}{q(z|x; \phi)} dz \\ &= \mathbb{E} [\ln p(x|z; \theta)] - D_{KL} [q(z|x; \phi) \| p(z)]\end{aligned}$$

Let's find a distribution more focused so I will sample for less for the same approximation  
**Importance Sampling**  
Best option is the posterior  $p(z|x)$

$$\ln p(x; \theta) \geq \mathbb{E}_{q(z|x; \phi)} [\ln p(x|z; \theta)] - D_{KL} [q(z|x; \phi) \| p(z)]$$

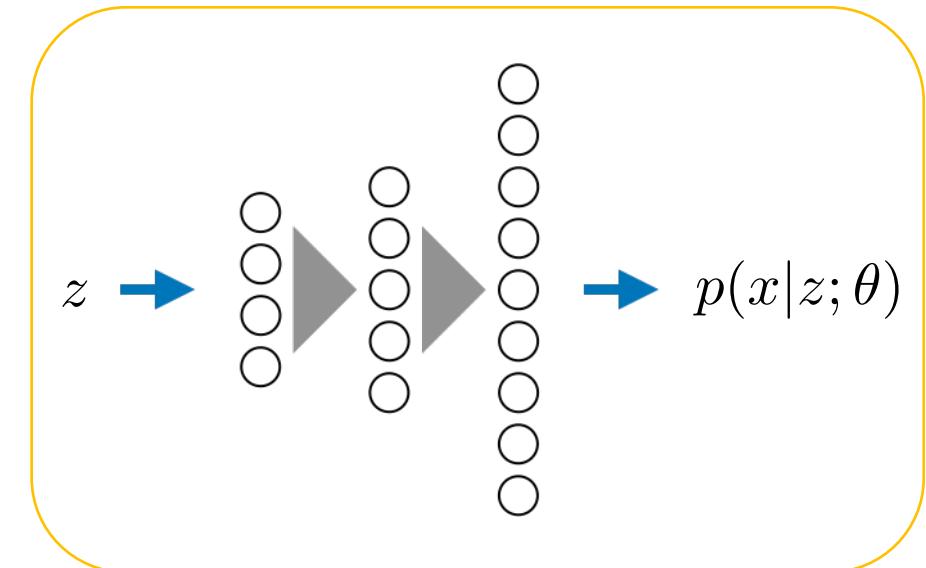
Evidence lower-bound : Maximize this instead of real likelihood

# Variational auto-encoders



Encoding Model

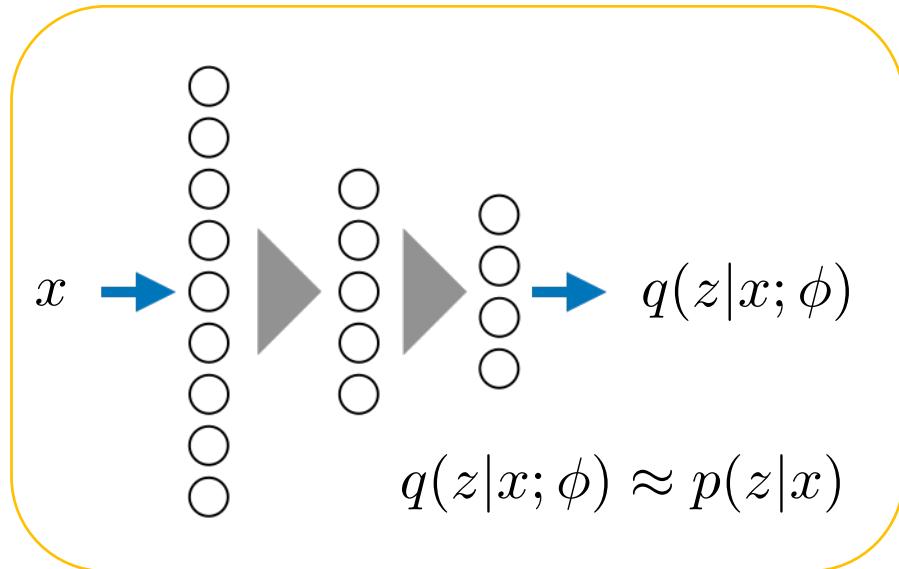
Takes an image and maps it to the posterior distribution in the latent space.  
Encodes to the lower dimensional space



Decoding Model

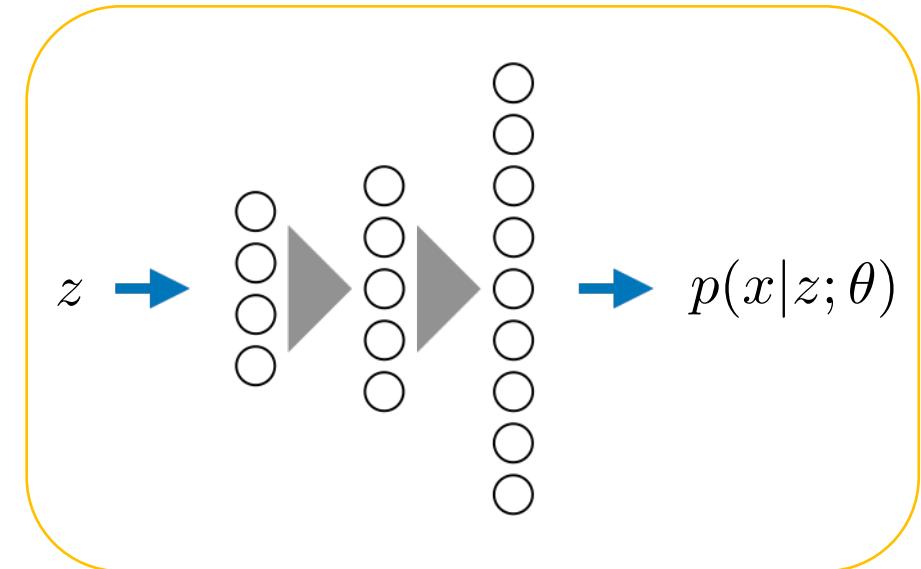
Takes the lower dimensional representation and maps to an image.  
Can be used as a sampler.  
Can be used as a reconstruction tool

# Variational auto-encoders



Encoding Model

$$q(z|x; \phi) = \mathcal{N}(z; \mu_z(x; \phi), \Sigma_z(x; \phi))$$



Decoding Model

$$p(x|z; \theta) = \mathcal{N}(x; \mu_x(z; \theta), \Sigma_x(z; \theta))$$

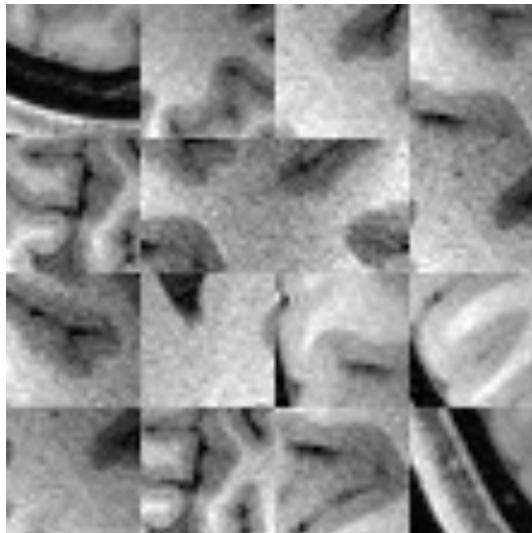
Both Gaussian Models

Homework: Can you determine the link with the probabilistic PCA model?

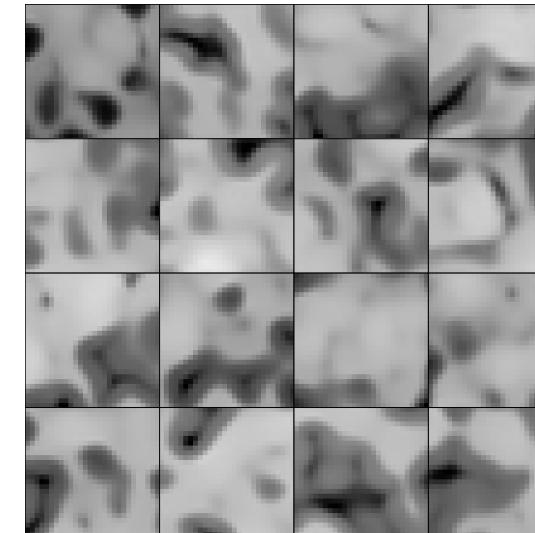
# Difference with PCA

Image patches from Magnetic Resonance Images of the brain

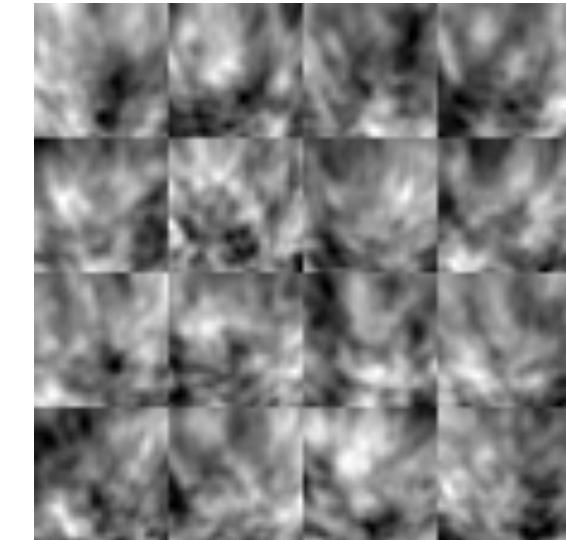
Real patches of 28x28



VAE Generated



PCA

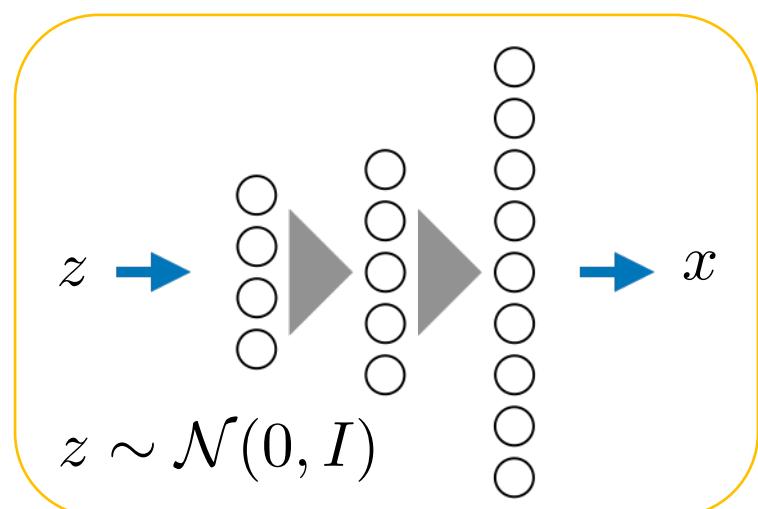


60 components

250 components

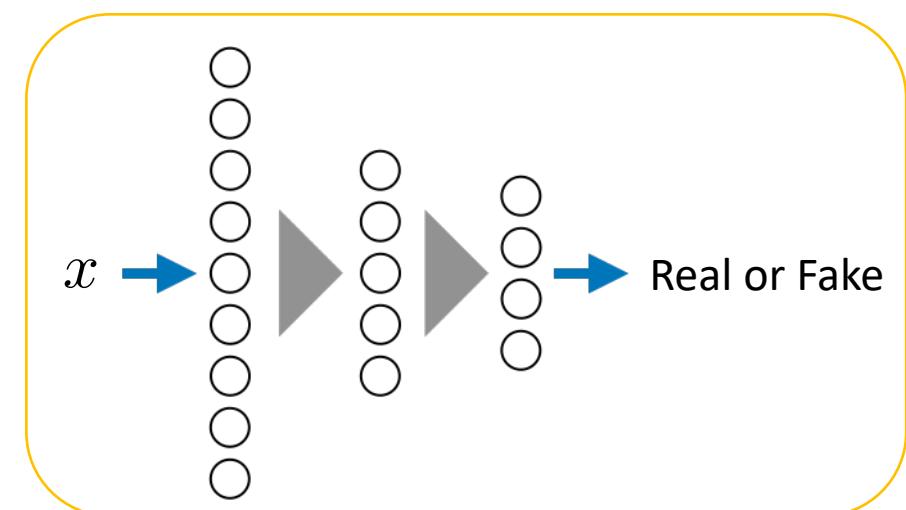
# Generative adversarial networks

Instead of an explicit probabilistic model, a GAN is a sampling tool that generates samples from the data distribution



Generator

Generates realistic looking images  
from random samples in the  
latent space.



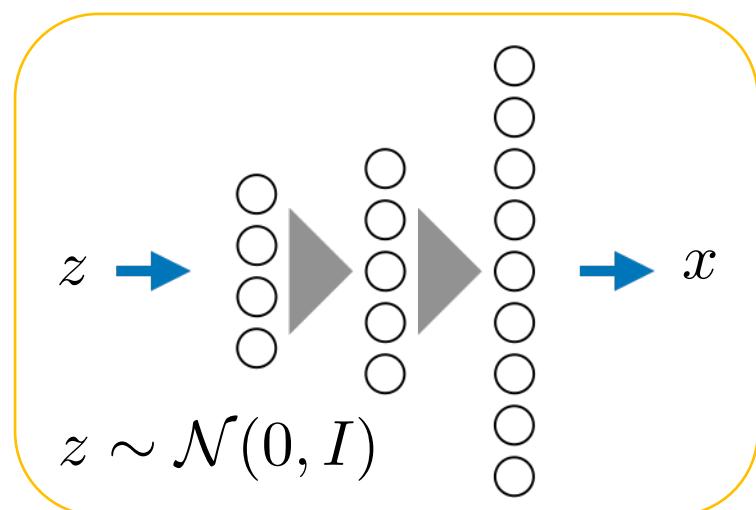
Discriminator

Tries to classify images into two categories:  
Real or generated (Fake)

# During training they compete

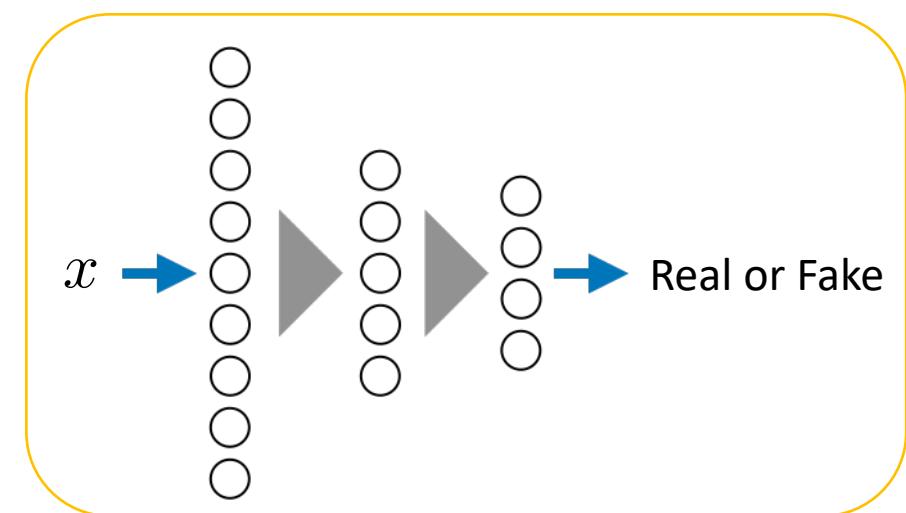
Generator - G

Tries to create samples that can  
fool the discriminator



Discriminator - D

Tries to identify the images  
the generator creates



Solve this problem: Optimize the network weights with a two-player game

$$\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{\text{real}}}(x) [\ln D(x; \phi)] + \mathbb{E}_{z \sim p(z)} [\ln(1 - D(G(z; \theta)))]$$

# Random samples



a)



b)



c)



d)

[Images from Goodfellow et al. 2014]

# Computer Vision

Detection

Segmentation

Visualization

Localization

Unsupervised  
learning

## Very active area of research



### A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras  
NVIDIA

[tkarras@nvidia.com](mailto:tkarras@nvidia.com)

Samuli Laine  
NVIDIA

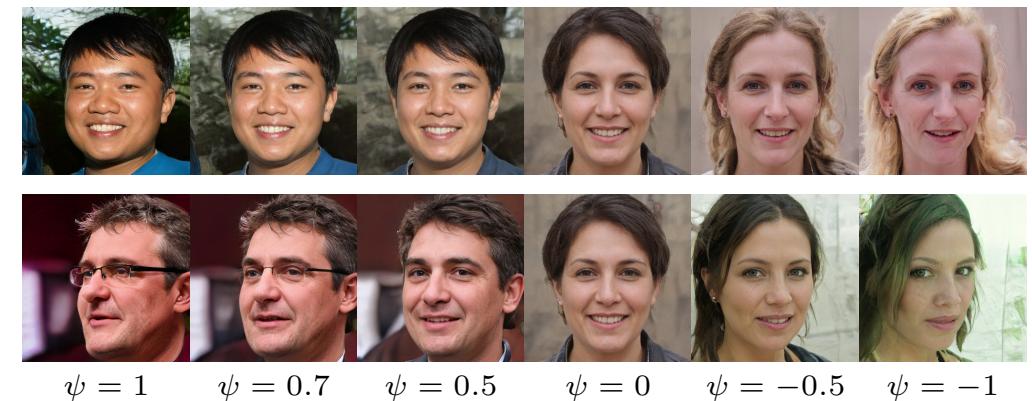
[slaine@nvidia.com](mailto:slaine@nvidia.com)

Timo Aila  
NVIDIA

[taila@nvidia.com](mailto:taila@nvidia.com)

December 12, 2018

- **The model is not yet peer-reviewed**
- However, the samples they claim to generate are remarkable.

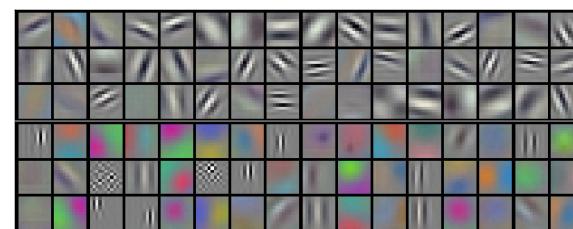
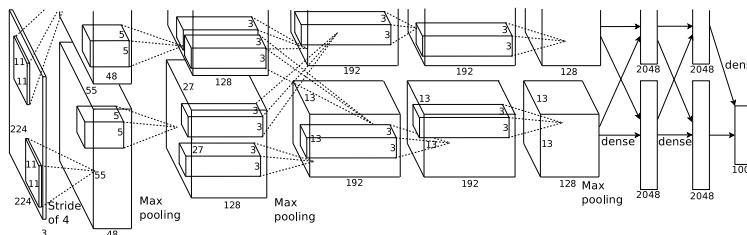


Interpolation between images

# Unsupervised learning

## Unsupervised learning of features

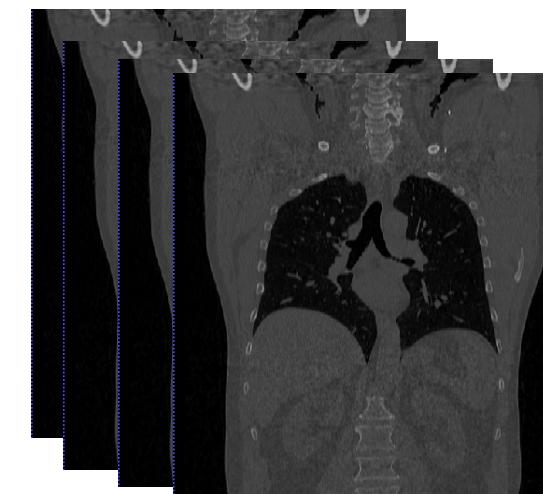
- Filters are important for performing image analysis tasks
- So far, we determine features in a supervised way, task-specific manner



- Determine features in an unsupervised manner
- Examples have only features

## Unsupervised learning of distributions

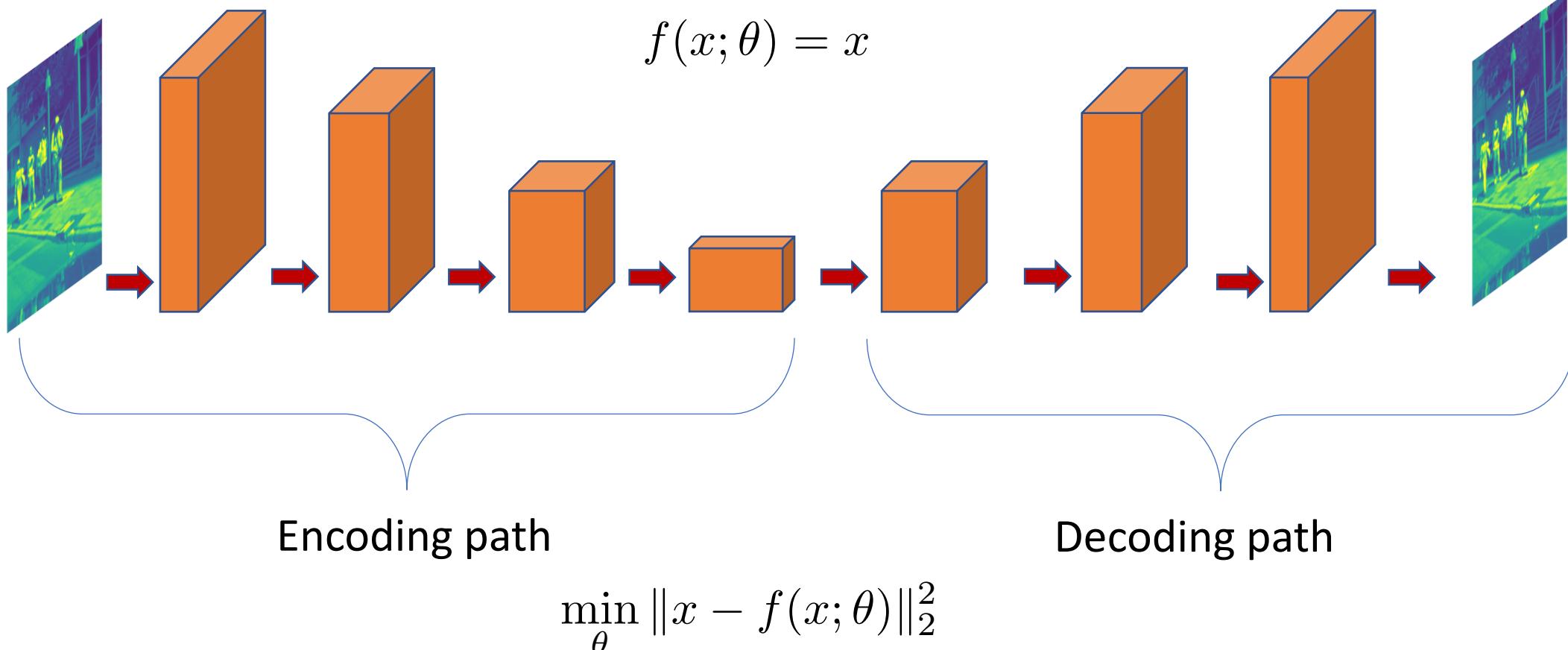
- Patterns within the data
- Goal: describe variability in the data
- Estimate the distribution of the data
- Only features
- Examples have only features



# Unsupervised learning of features

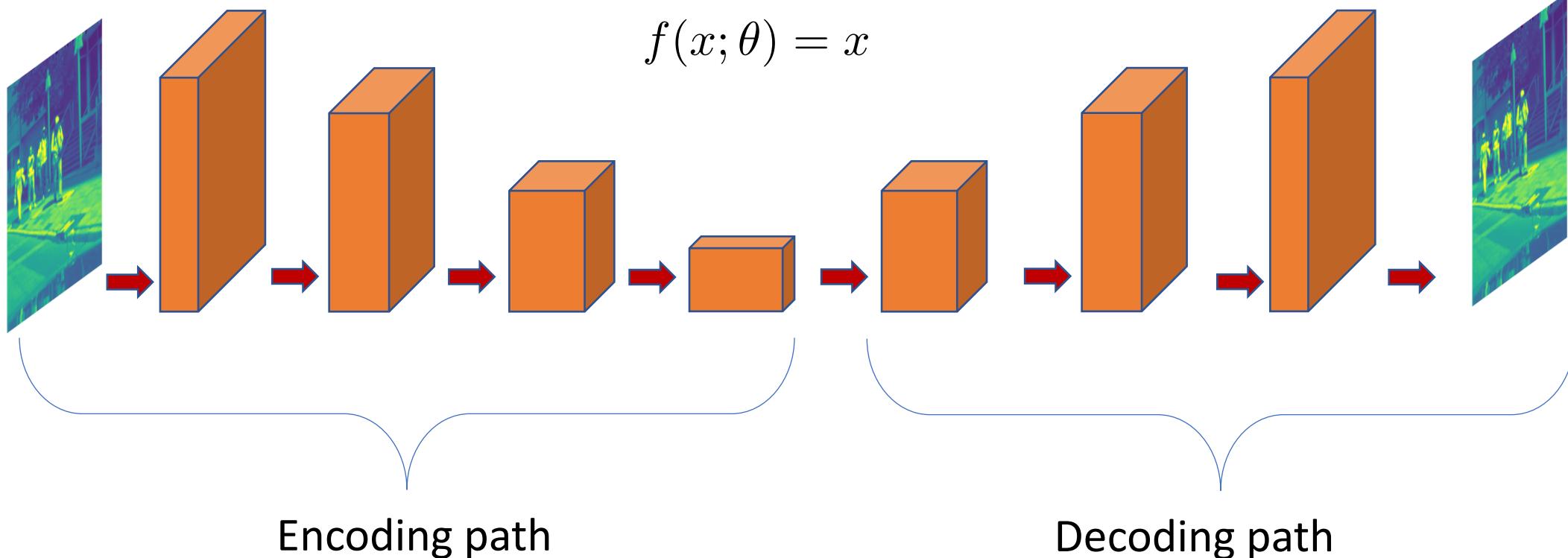
- Features are important, they are the essential building blocks
- For any task it is important to get the right features
- It requires large number of labelled images to do this
  1. It would be wonderful if we could do it with only few images
  2. Humans do not seem to require lots of labelled images for good features, assuming humans do have good features
  3. Are there features that can be used for any visual task?

# Auto-encoding models



The bottleneck layer does not allow the network to learn an identity map  
It learns to summarize the most important information for reconstruction

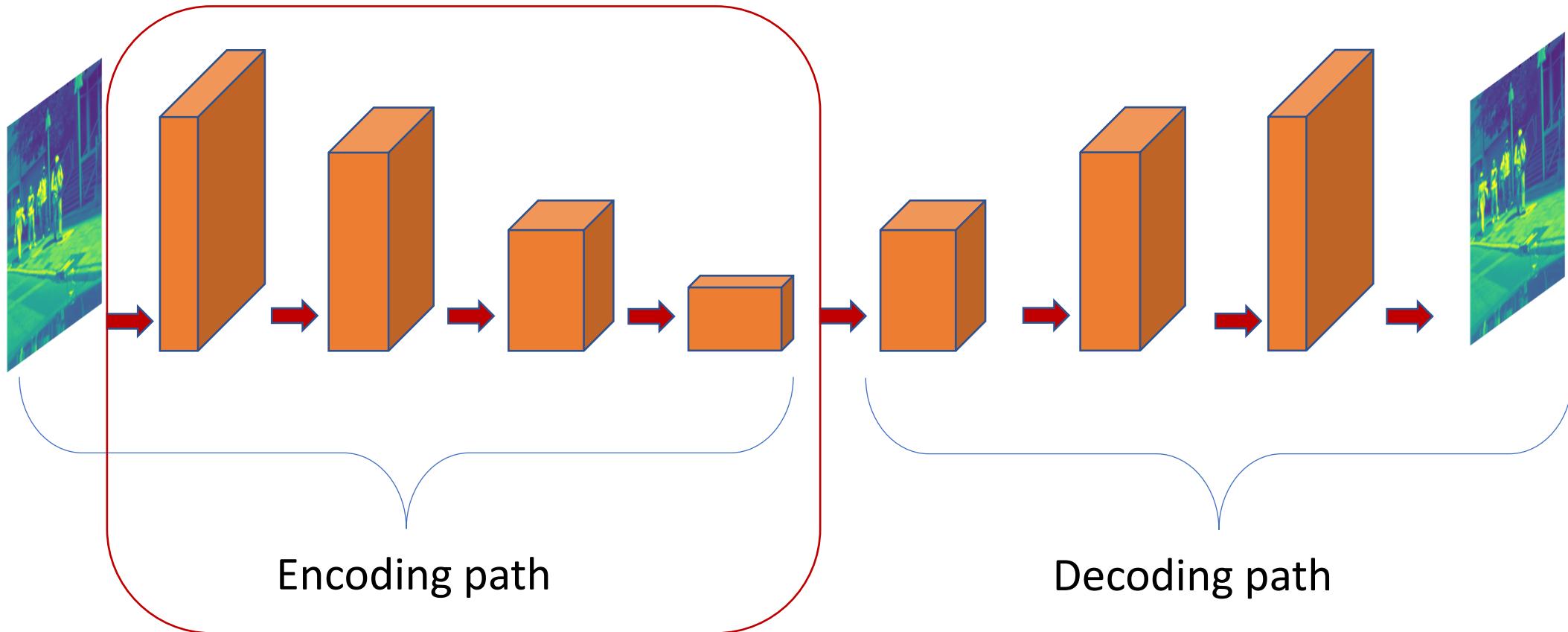
# Auto-encoding models



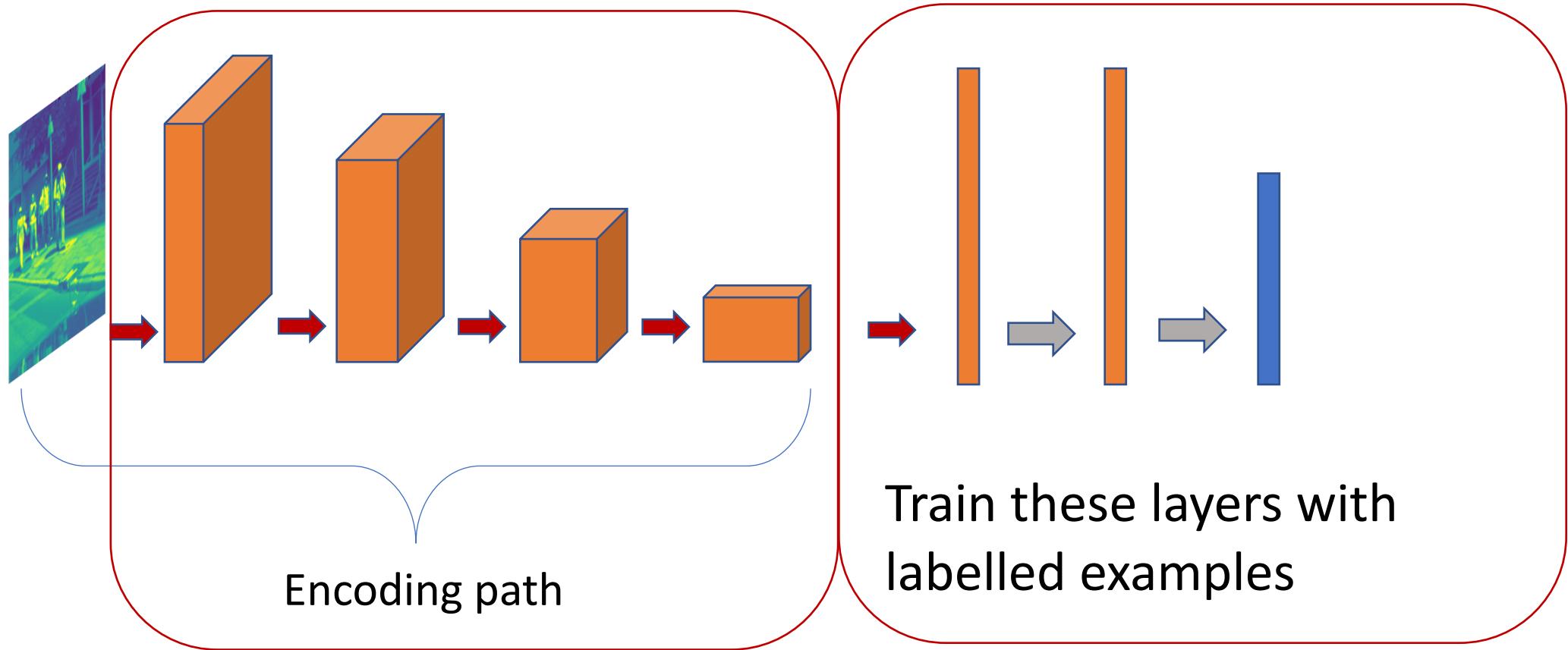
$$\min_{\theta} \|x - f(x; \theta)\|_2^2$$

Minimization only requires the images. The goal is to be able to reconstruct the image with high fidelity.

# Auto-encoding models



# Auto-encoding models



Features learnt here can then be translated to another task either directly or by fine-tuning, i.e. starting the optimization from the pre-learnt weights

## In practice

- The features learnt from a simple auto-encoder are not extremely useful
- In the end, you may still need large number of labelled examples
- Not as large as training from scratch though

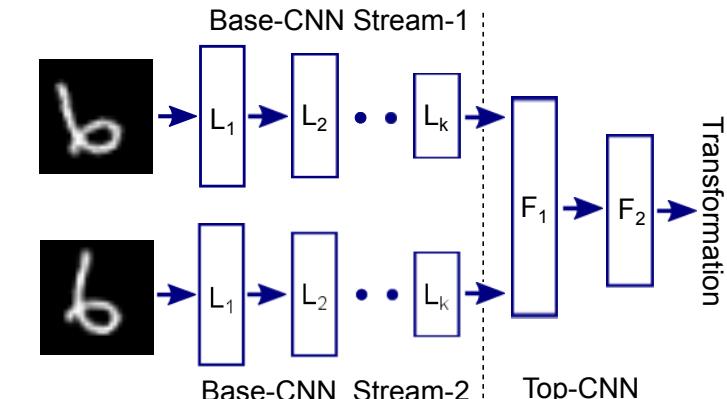
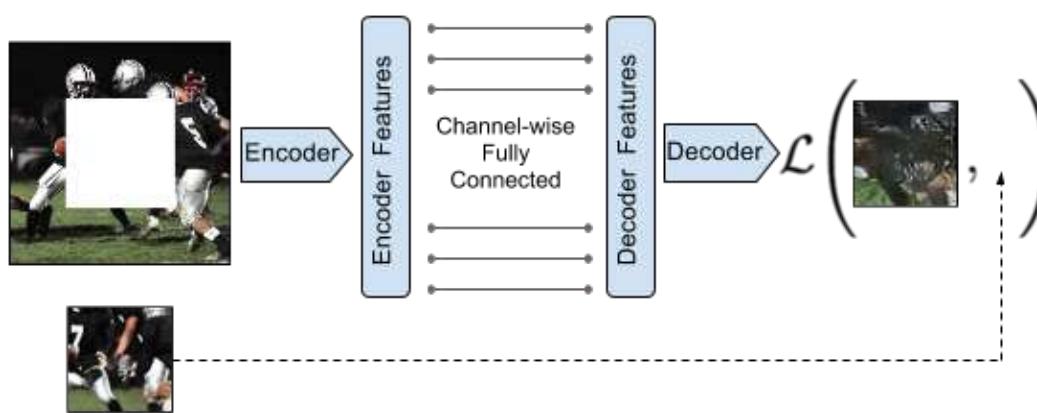
# Other “self supervised” tasks

## Context Encoders: Feature Learning by Inpainting

Deepak Pathak Philipp Krähenbühl Jeff Donahue Trevor Darrell Alexei A. Efros

University of California, Berkeley

{pathak, philkr, jdonahue, trevor, efros}@cs.berkeley.edu



## Learning to See by Moving

Pulkit Agrawal  
UC Berkeley

pulkitag@eecs.berkeley.edu

João Carreira  
UC Berkeley

carreira@eecs.berkeley.edu

Jitendra Malik  
UC Berkeley

malik@eecs.berkeley.edu