

# Introduction to Machine Learning

## Neural networks (continued)

Dr. Kfir Levy  
Learning and Adaptive Systems ([las.ethz.ch](http://las.ethz.ch))

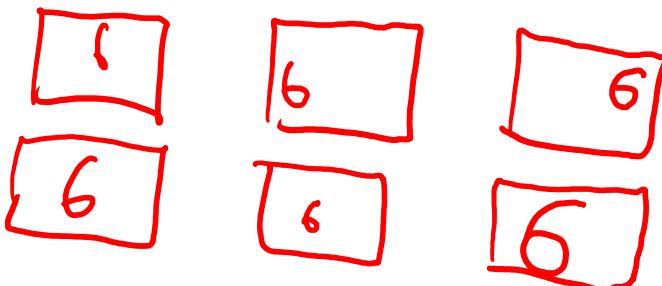
# Invariances

- Predictions should be unchanged under some transformations of the data, e.g.
  - Classification of handwritten-digits

Shift invariance



Rotation invariance



scale invariance

- Speech recognition

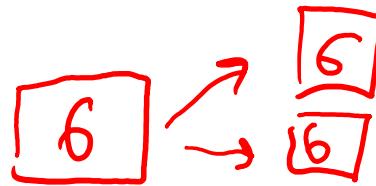
pitch invariant features

spect of speech

# Invariances (contd.)

- How to encourage a model to learn specific invariances?

*SIFT - shift invariant feature transformation  
Ceptrum (speech recognition features)*

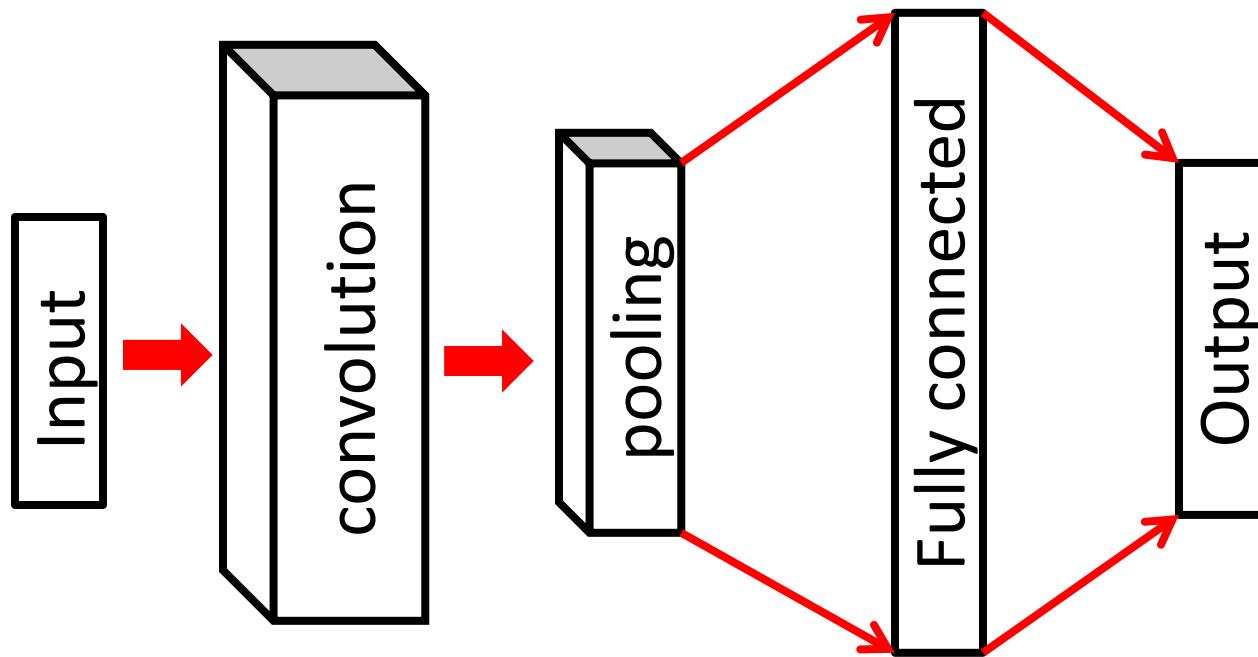


- Augmentation of the training set
- Special regularization terms
- Invariance built into pre-processing
- Implement invariance into structure of ANN  
(→ e.g. using convolutional neural networks)

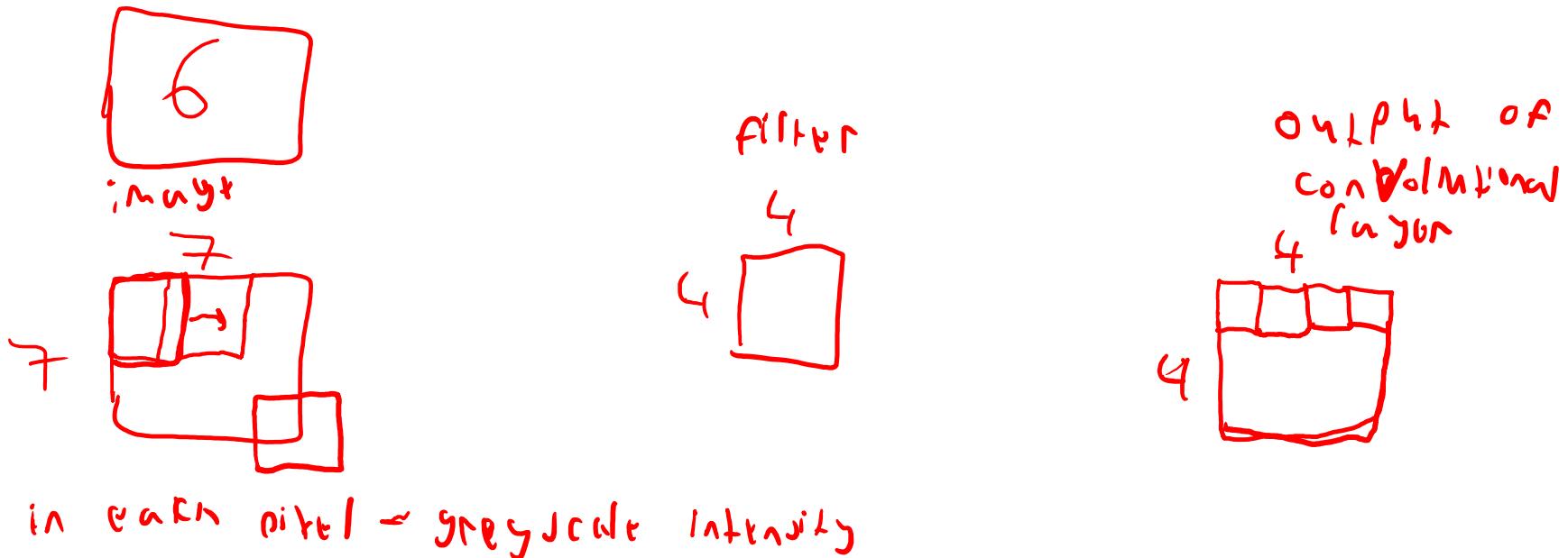
# Convolutional neural networks

- Convolutional neural networks are ANNs for **specialized applications** (e.g., image recognition)
- The hidden layer(s) closest to the input layer **shares parameters**: Each hidden unit only depends on all „closeby“ inputs (e.g., pixels), and weights constrained to be identical across all units on the layer
- This reduces the number of parameters, and encourages robustness against (small amounts of) translation
- The weights can still be optimized via backpropagation

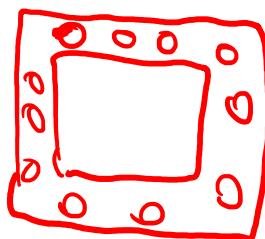
# CNN architecture



# CNN examples

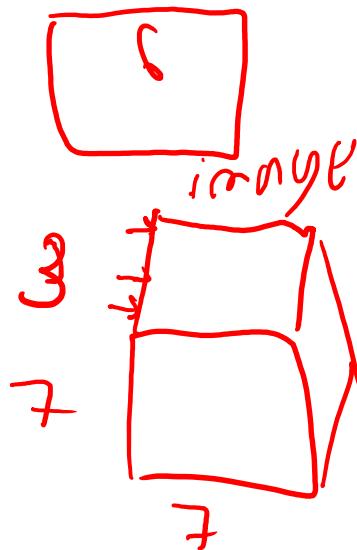


zero padding

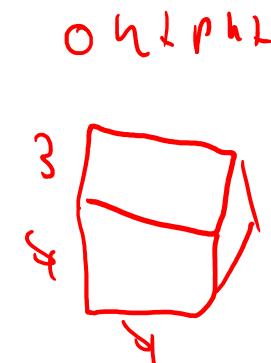
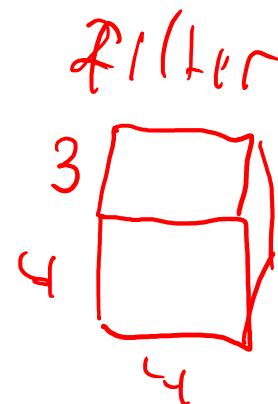


- **Terminology:** filters/kernels, stride, padding

# CNN examples



3 RGB channels



- **Terminology:** filters/kernels, stride, padding

Input Volume (+pad 1) (7x7x3)

 $x[:, :, 0]$ 

0	0	0	0	0	0	0
0	2	1	1	2	0	0
0	0	0	2	0	0	0
0	0	2	2	1	2	0
0	0	2	0	1	2	0
0	0	2	1	1	1	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

 $w0[:, :, 0]$ 

1	1	1
1	-1	1
0	1	-1

 $w0[:, :, 1]$ 

0	0	0
1	-1	-1
1	-1	0

 $w0[:, :, 2]$ 

-1	0	0
0	0	1
-1	1	-1

Bias b0 (1x1x1)

 $b0[:, :, 0]$ 

1
---

Bias b1 (1x1x1)

 $b1[:, :, 0]$ 

0
---

Output Volume (3x3x2)

 $o[:, :, 1]$ 

0	-1	-2
4	-1	-2
4	2	-4

$\Sigma = 2$

$\Sigma = -2$

$\Sigma \approx 0$

 $x[:, :, 2]$ 

0	0	0	0	0	0	0
0	0	0	2	1	2	0
0	0	2	1	2	0	0
0	1	2	0	1	1	0
0	2	1	2	2	0	0
0	2	2	2	2	2	0
0	0	0	0	0	0	0

[A. Karpathy, Stanford]

# Computing output dimensions

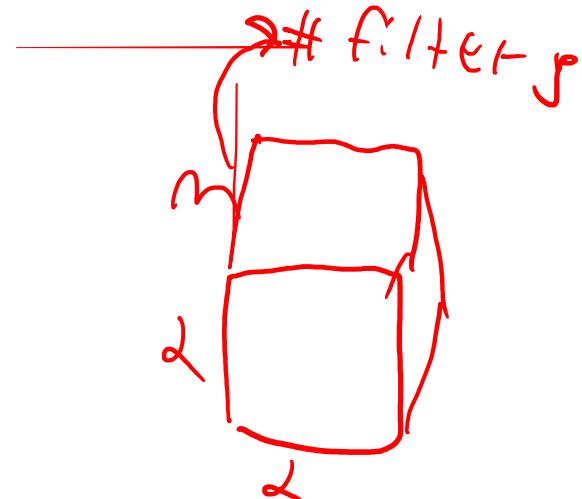
- Applying  $m$  different  $f \times f$  filters
- To a  $n \times n$  image
- Padding  $p$
- Stride  $s$

$$f=3$$

$$n=7$$

$$p=1$$

$$s=2$$

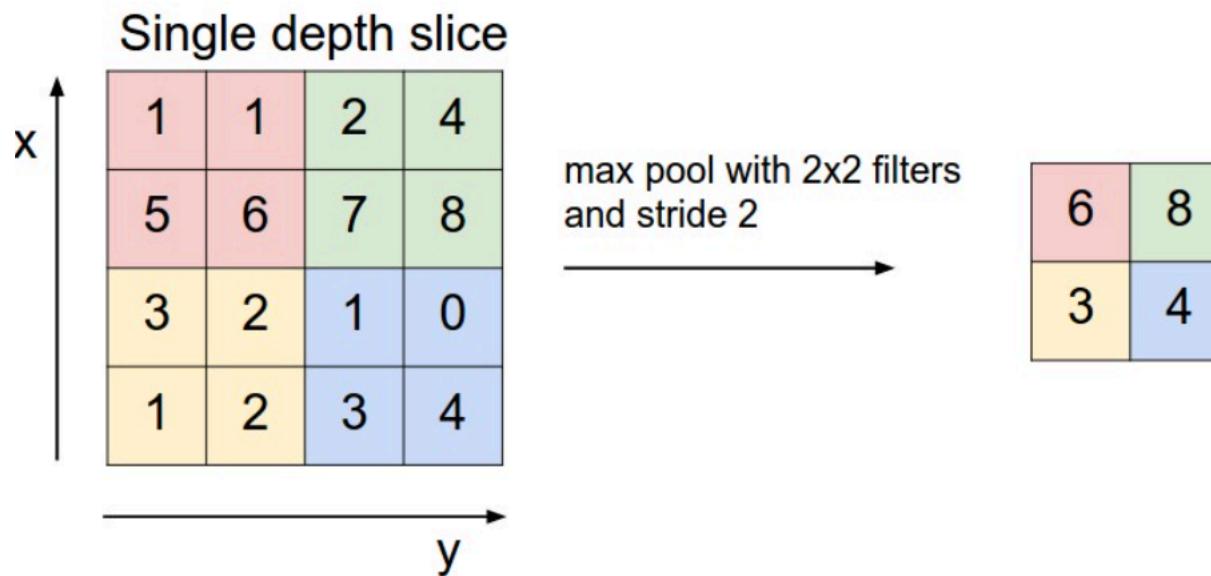


$$d = \frac{n + 2p - f}{s} + 1$$

# Pooling layers (subsampling)

- In some applications (e.g., image classification), it can make sense to aggregate several units to decrease the width of the network (and hence # of parameters)
- Usually, one considers either the average or the maximum value

# Illustration of pooling layers

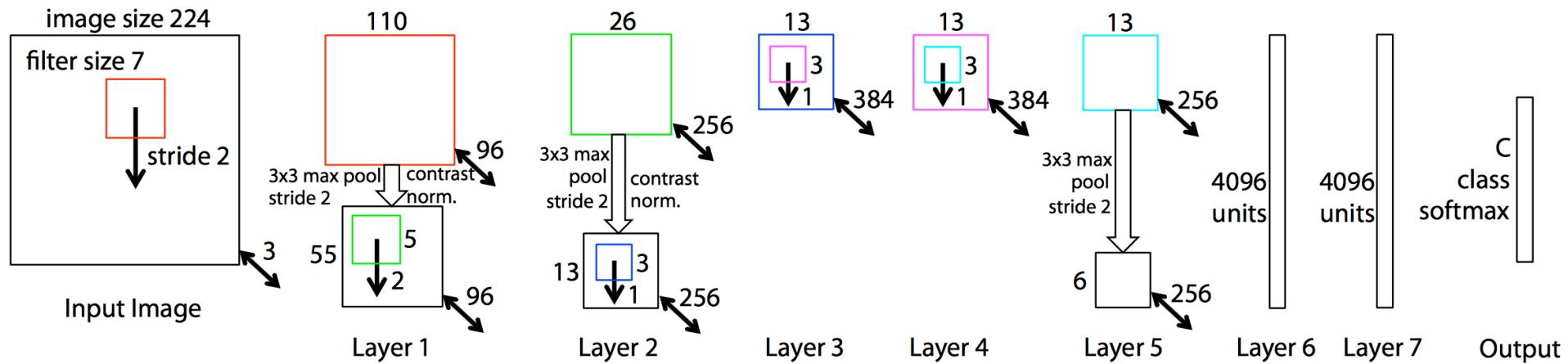


# CNN Demo

---

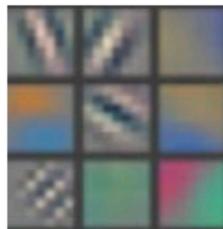
- <http://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

# Example CNN architecture



[Zeiler & Fergus '13]

# Visualizing the filters

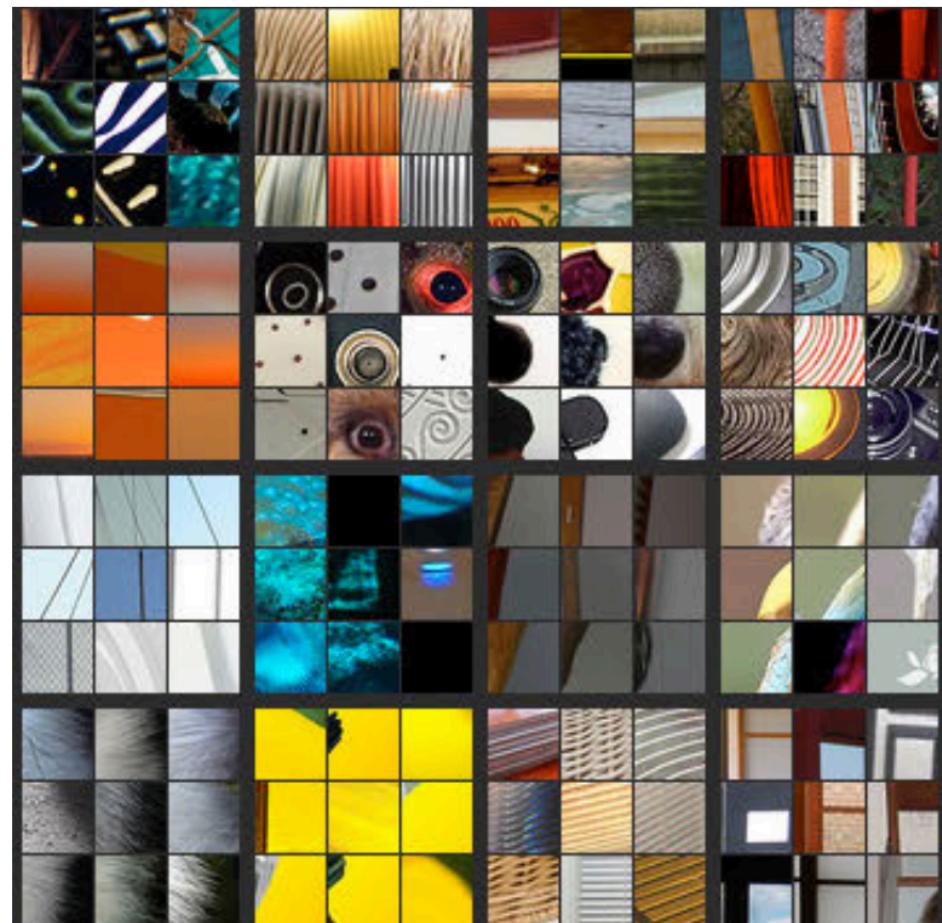
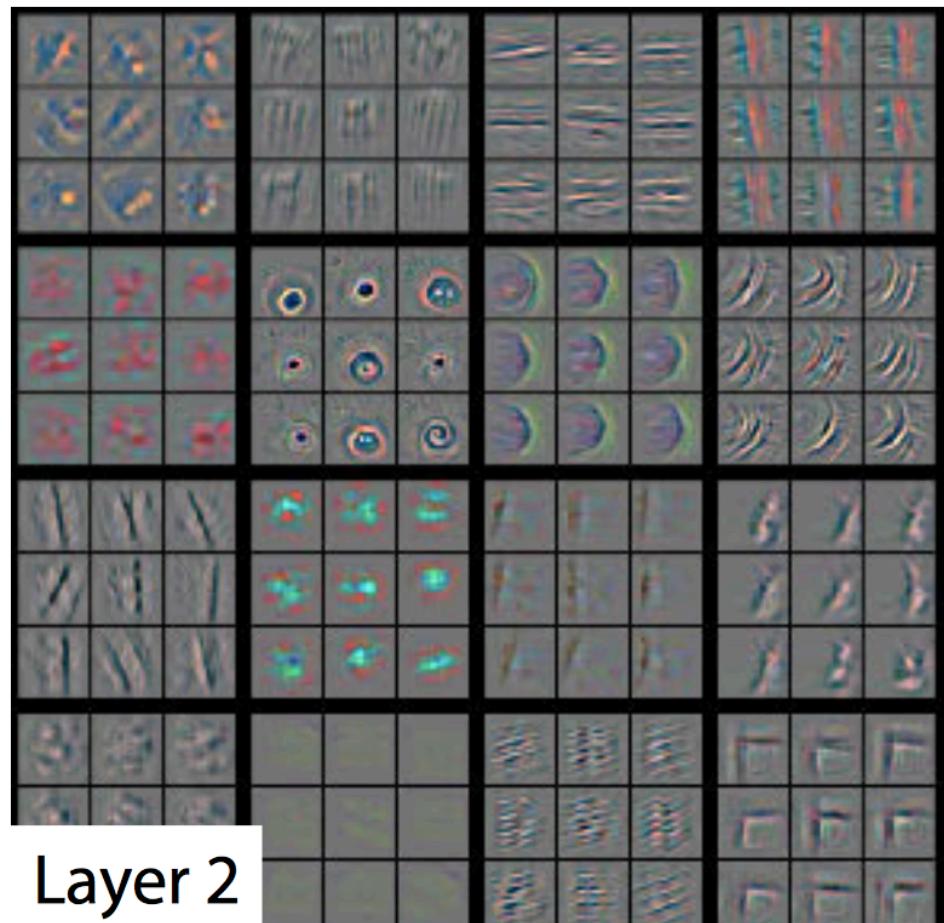


Layer 1



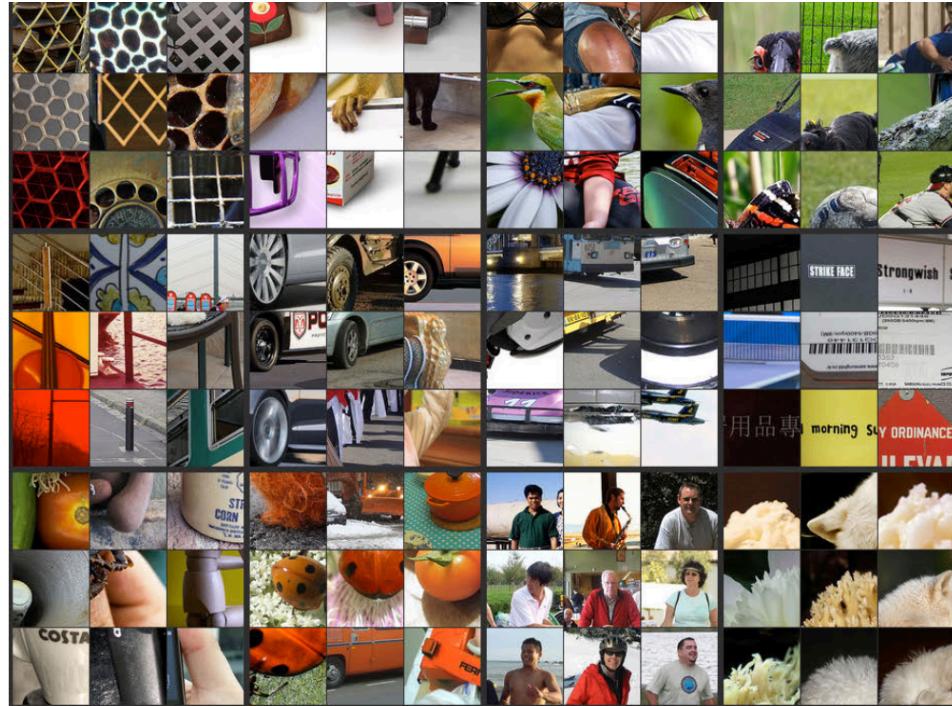
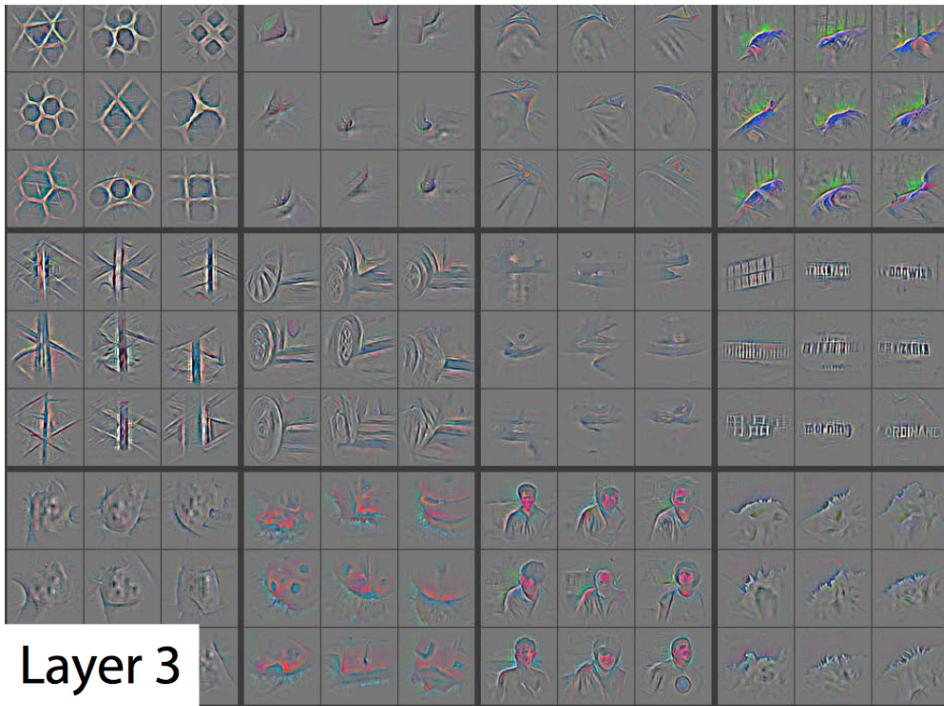
[Zeiler & Fergus '13]

# Visualizing the filters

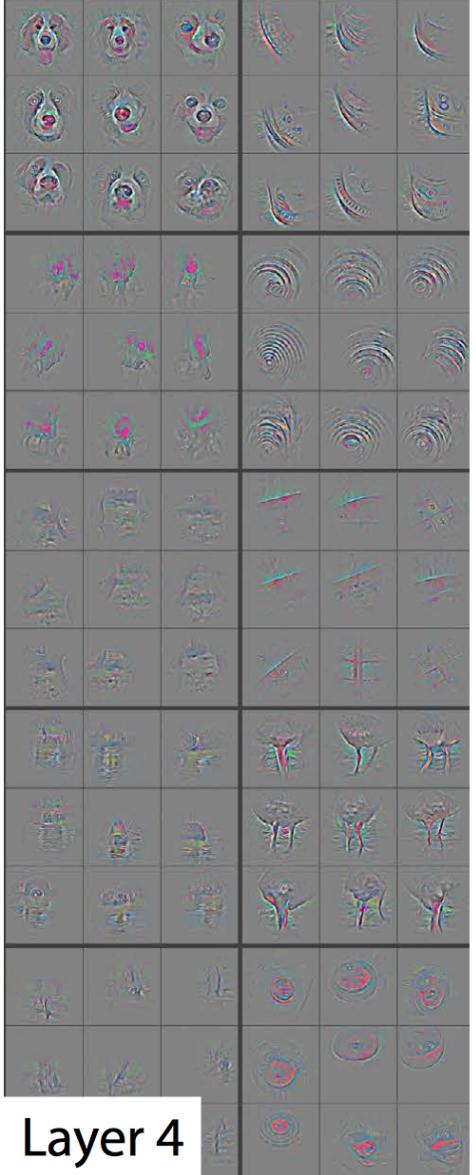


[Zeiler & Fergus '13]

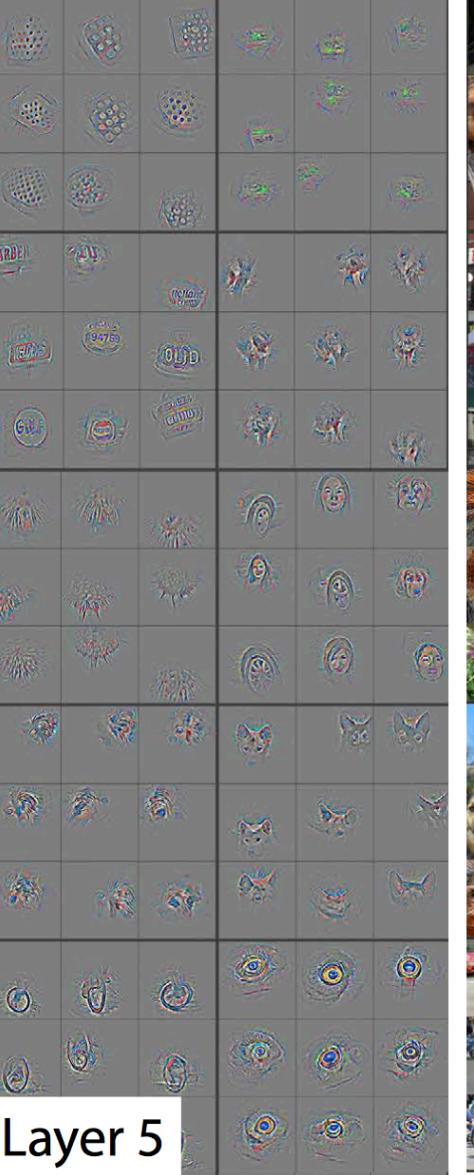
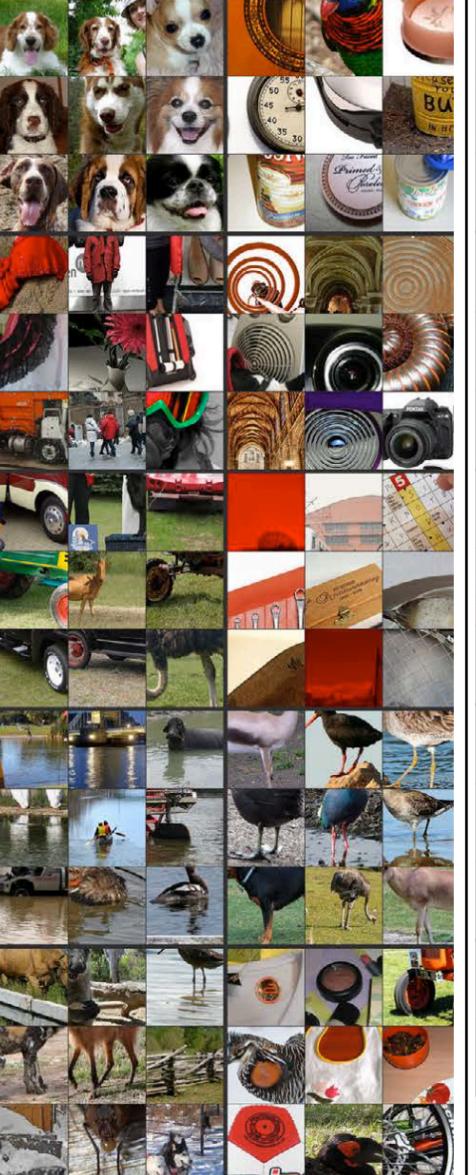
# Visualizing the filters



[Zeiler & Fergus '13]



Layer 4



Layer 5



[Zeiler & Fergus '13]

# Tutorial: Keras

---

# Computational efficiency

- Key operation for backpropagation training:
  - Dense matrix vector multiplications
- Very suitable for implementation on  
**General Purpose Graphical Processing Units (GPUs)**
- Much work on special purpose hardware (e.g., TPUs)
- Responsible for large improvements in recent years

# Choice of activation functions

- Several possible activation functions available
- Popular in the past: sigmoid / tanh activations
  - Differentiable!
- Currently used extensively: Rectified linear units (ReLUs)
  - Not differentiable ☹
  - Very fast to compute ☺
  - Gradients do not „vanish“ (important for deep networks) ☺

# Parameter / architecture choice

- A number of parameters have to be chosen
  - Number and width of hidden layers
  - Use convolution / pooling layers? How many?
  - Skip connections
  - Type of activation functions
  - Weight initialization
  - Learning rate schedule
  - Regularization method
- Frameworks for **autodifferentiation** (e.g., Theano, Torch, TensorFlow)
- Can use cross-validation to compare models  
(but training is usually very expensive)
  - ➔ Often use single validation („development“) set
  - ➔ Much work on ML-based experimentation (“AutoML”)

# Connection: ANNs vs. kernels

- When using the *tanh*-activation function, an ANN with a *single hidden layer* learns functions of the form

$$f(\mathbf{x}) = \sum_{i=1} w_i \tanh \theta_i^T \mathbf{x}$$

- This is exactly the same type of functions learned with kernels, when using the *tanh kernel*

$$k(x_i, x) = \langle \phi(x_i), \phi(x) \rangle$$

$$f(\mathbf{x}) = \sum_{i=1} \alpha_i \tanh \mathbf{x}_i^T \mathbf{x}$$

- Difference:** kernels optimize  $\alpha$ 's only → convex  
ANNs optimize both  $w$ 's and  $\theta_i$  → non-convex

# Comparison: Kernels vs. ANNs

<b>Method</b>	<i>Kernels</i>	<i>ANNs</i>
<b>Advantages</b>	Convex optimization, no local minima Robust against noise Models grow with size of data ☺	Flexible nonlinear models, with fixed parameterization Multiple layers discover representations at „multiple levels of abstraction“
<b>Disadvantages</b>	Models grow with size of data ☹ Don't allow „multiple layers“	Many free parameters / architectural choices that need to be tuned Often suffer from very noisy data

# Supervised learning summary so far

Representation/ features	Linear hypotheses; nonlinear hypotheses with nonlinear feature transforms, kernels, <b>learn nonlinear features via neural nets</b>		
Model/ objective:	<b>Loss-function</b>	+	<b>Regularization</b>
	Squared loss, 0/1 loss, Perceptron loss, Hinge loss, cost sensitive losses, multi-class hinge loss		$L^2$ norm, $L^1$ norm, <b>early stopping, dropout</b>
Method:	Exact solution, Gradient Descent, (mini-batch) SGD, Reductions		
Evaluation metric:	Mean squared error, Accuracy, F1 score, AUC, Confusion matrices		
Model selection:	K-fold Cross-Validation, Monte Carlo CV		