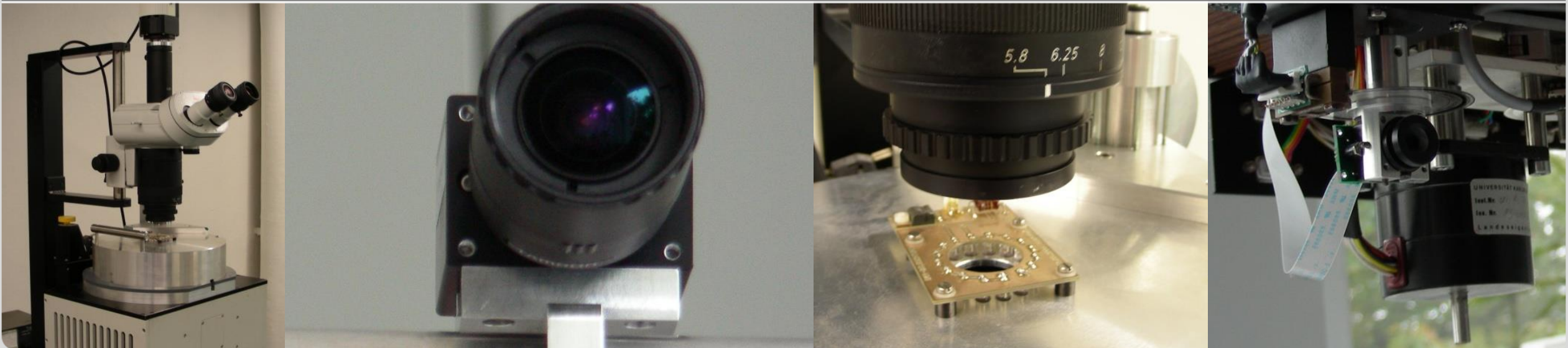


Assignment 07 – Deep Learning

Christian Kinzig

INSTITUTE OF MEASUREMENT AND CONTROL SYSTEMS, DEPARTMENT OF MECHANICAL ENGINEERING



Multi-Layer Perceptrons (MLP)

- MLPs are highly parameterized, non-linear functions

$$f_{MLP}^{\vec{w}} : \mathbb{R}^n \rightarrow \mathbb{R}^q$$

$$f_{MLP}^{\vec{w}} : \vec{x} \mapsto \vec{y}$$

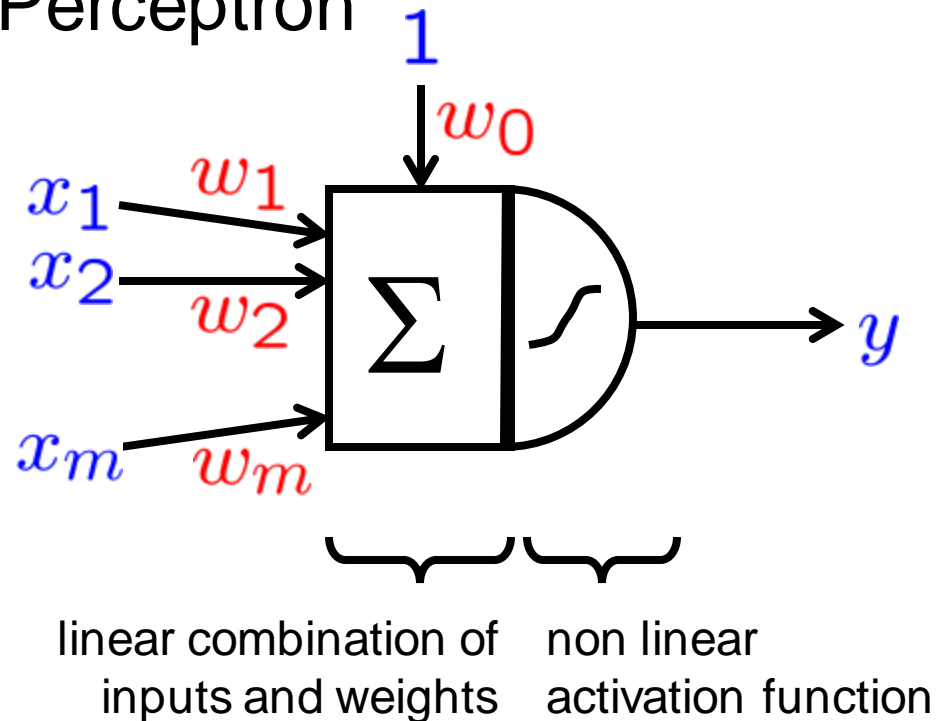
parameter vector,
weight vector (input-)
pattern output



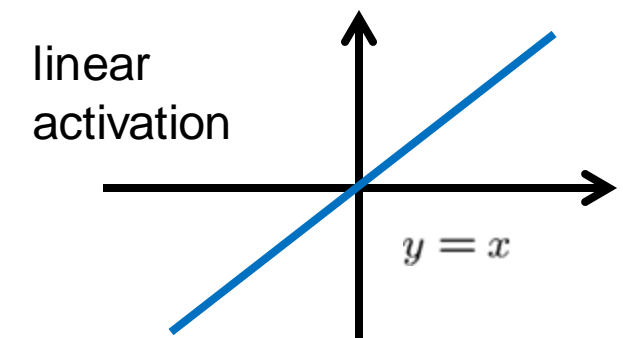
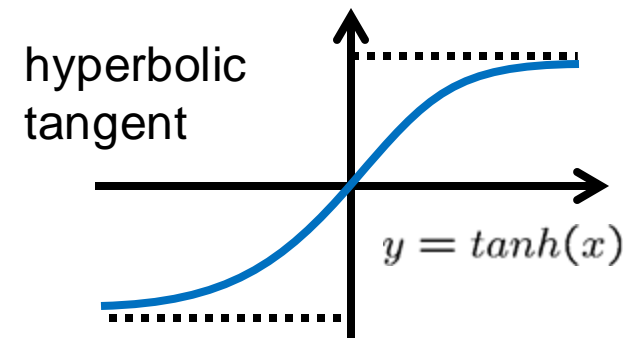
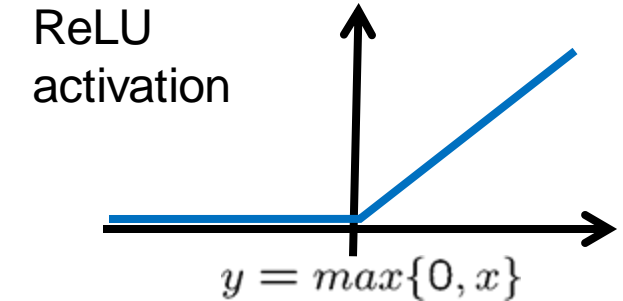
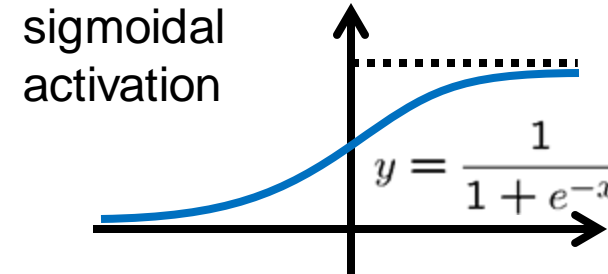
- Example: classification of images
 - \vec{x} : feature vector, e.g. vector of all gray values in image
 - \vec{y} : 1-of-q-vector that models probabilities for each of q possible categories, e.g. smiley is happy/sad/frustrated

Internal Structure of MLPs

■ Perceptron

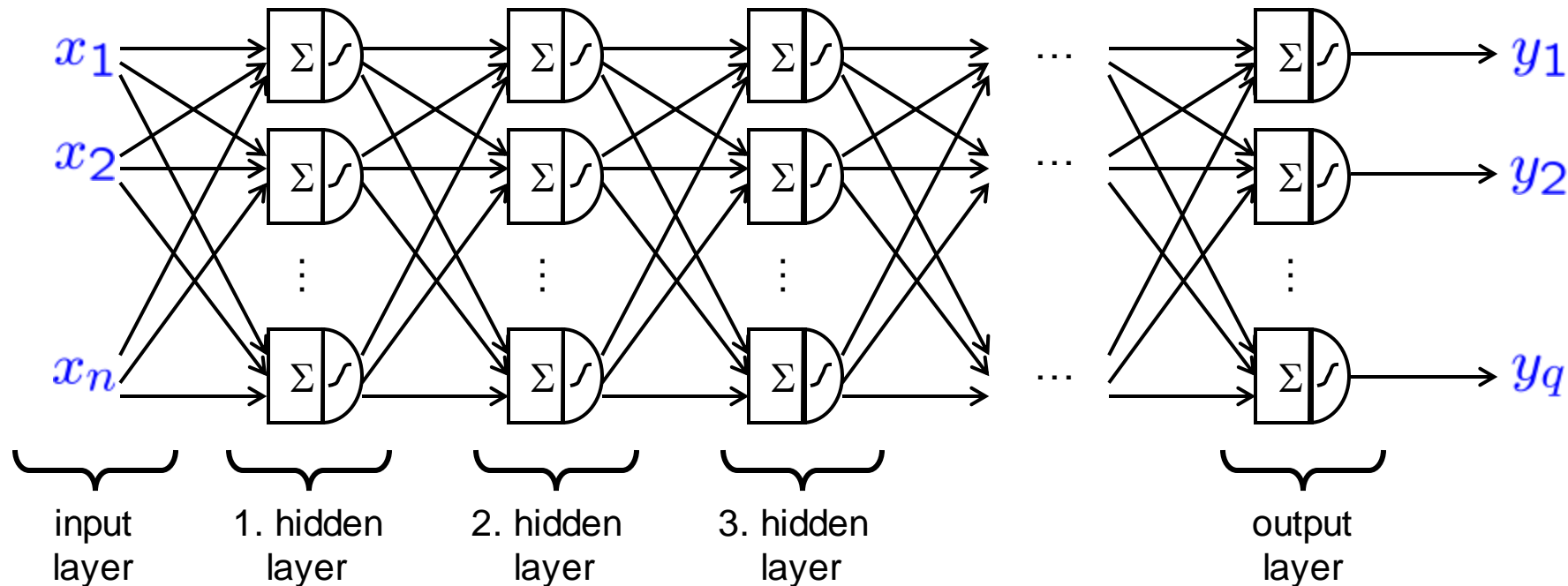


$$y = f_{act}\left(w_0 + \sum_{i=1}^m w_i x_i\right)$$



Internal Structure of MLPs

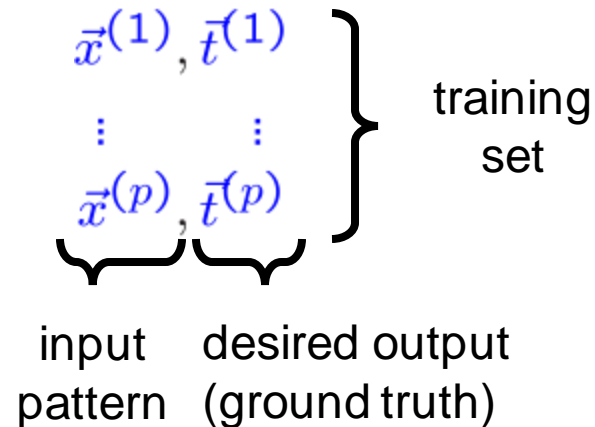
- Layered arrangement of many perceptrons



- Network structure creates set of highly nonlinear function
- Many weights
- Deep architectures: typically >5 hidden layers

Training of MLPs

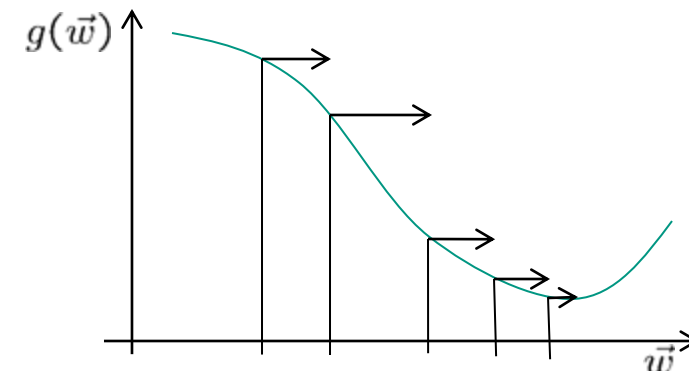
- How do we determine weights of MLP?
 - Basic idea: minimize error for training examples



- Solve $\underset{\vec{w}}{\text{minimize}} \sum_{j=1}^p \text{err}(f_{MLP}^{\vec{w}}(\vec{x}^{(j)}), \vec{t}^{(j)})$ for appropriate error measure err
- Algorithm: gradient descent (backpropagation)

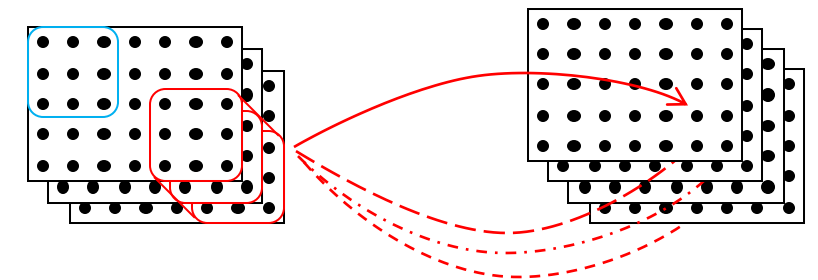
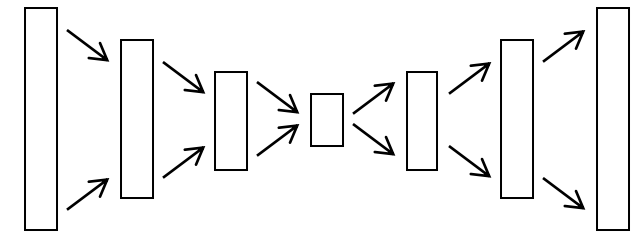
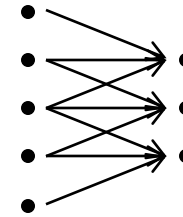
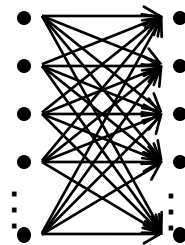
Gradient Descent (Backpropagation)

- Goal: *minimize* $g(\vec{w})$ with $g(\vec{w}) := \sum_{j=1}^p \text{err}(f_{MLP}^{\vec{w}}(\vec{x}^{(j)}), \vec{t}^{(j)})$
- Algorithm:
 1. Initialize weights \vec{w} randomly with small numbers
 2. Calculate gradient $\frac{\partial g(\vec{w})}{\partial \vec{w}}$
 3. Update weights $\vec{w} \leftarrow \vec{w} - \varepsilon \frac{\partial g(\vec{w})}{\partial \vec{w}}$ with small learning rate $\varepsilon > 0$
 4. GoTo 2. until stopping criterion reached



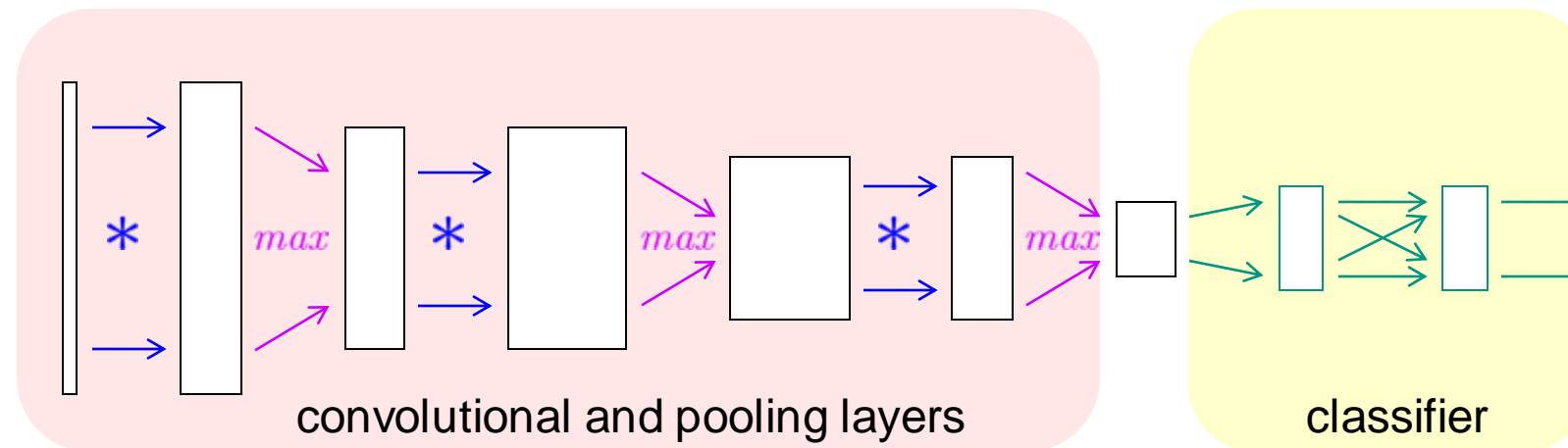
Deep Learning

- Larger training sets (millions instead of hundreds)
- More powerful computers, parallel implementations on multi-core CPUs and GPUs
- Special network structures
 - Autoencoders
 - Convolutional networks
 - ...
- Weight sharing
- Layer-wise learning
- Learning from unlabeled examples
- Learning of useful features



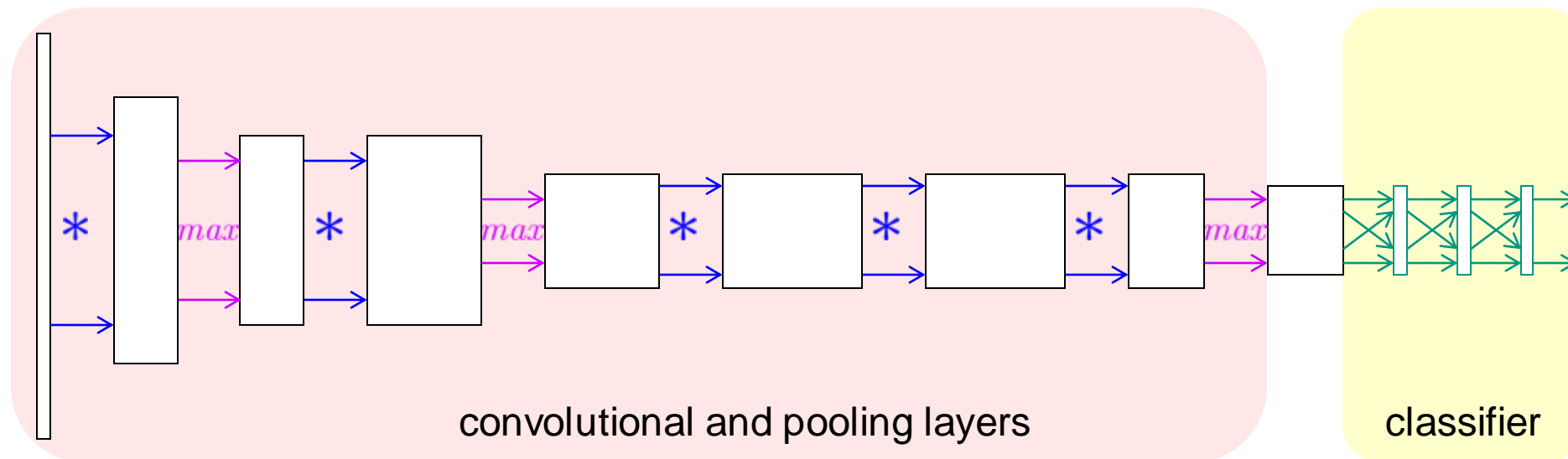
Convolutional Networks

- Convolutional Networks (CNNs) combine
 - Convolutional layers
 - Pooling layers
 - Fully connected classifier network



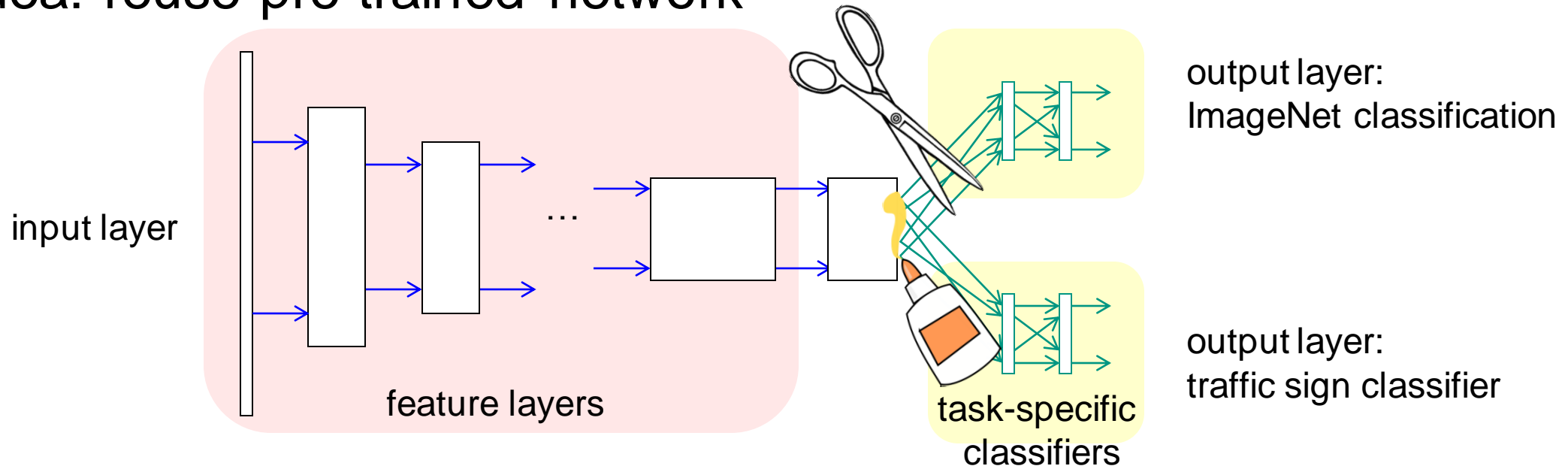
Example: AlexNet

- A. Krizhevsky, I. Sutskever, G. Hinton, „ImageNet Classification with Deep Convolutional Neural Networks“, NIPS, 2012
 - Classification of images, 1000 categories
 - Data set: 1,2 millions of images
 - Approach: convolutional network, 60 millions of weights



Usage of Pre-Trained Feature Networks

- Idea: reuse pre-trained network



1. Train other task with large training set
2. Throw away classification layers of other task
3. Create new classification layers for new task
4. Train weights of new classification layer while preserving feature layers

Convolutional Networks

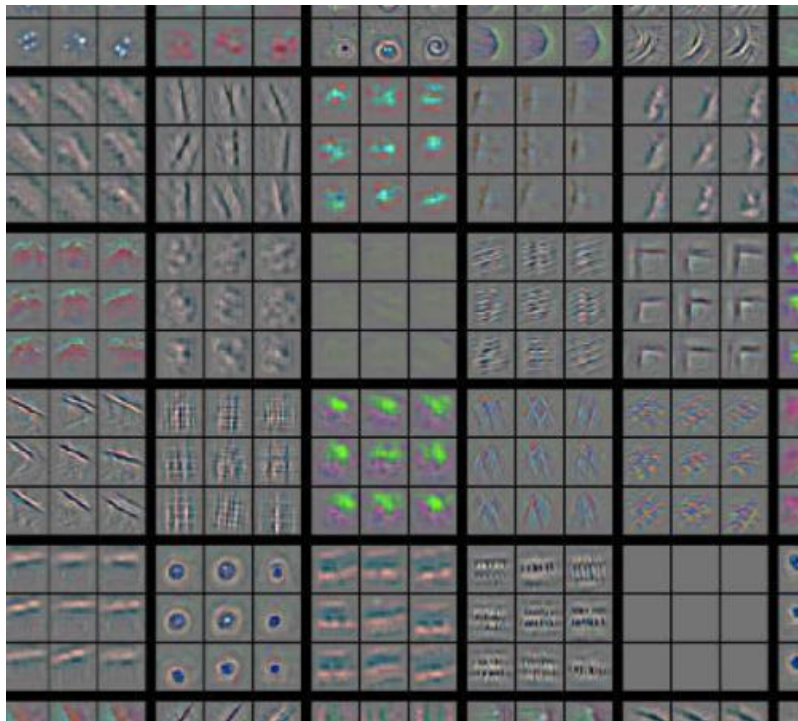
- Which features are learned in hidden layers?
 - 1. layer: gray level edges, color edges, blobs



Taken from:
<http://image-net.org/challenges/LSVRC/2012/supervision.pdf>

Convolutional Networks

- Which features are learned in hidden layers?
 - 2. layer: corners, round structures

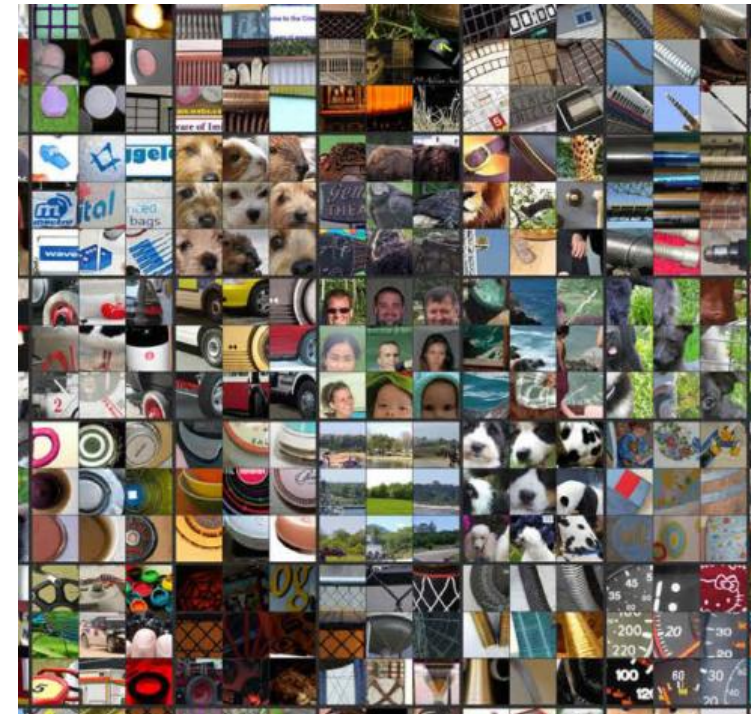
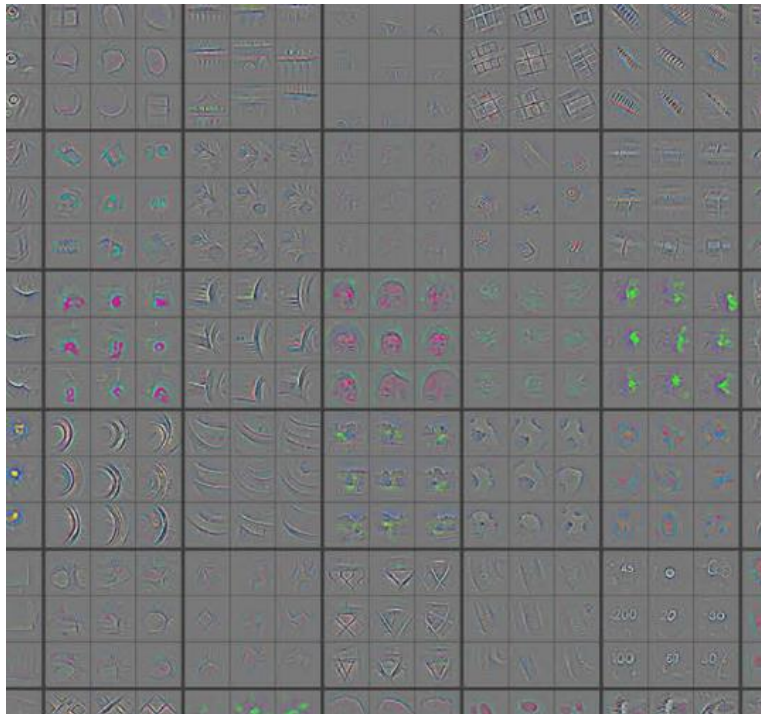


Taken from:

M.D. Zeller, R. Fergus, „Visualizing and Understanding Convolutional Neural Networks“,
arXiv:1311.2901 v2, 13. Nov. 2013

Convolutional Networks

- Which features are learned in hidden layers?
 - 3. layer: shapes, gratings

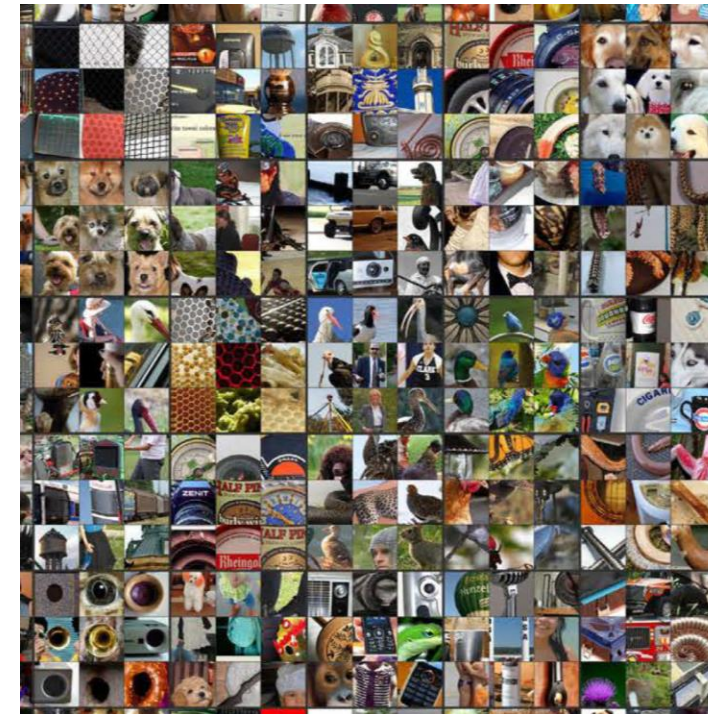
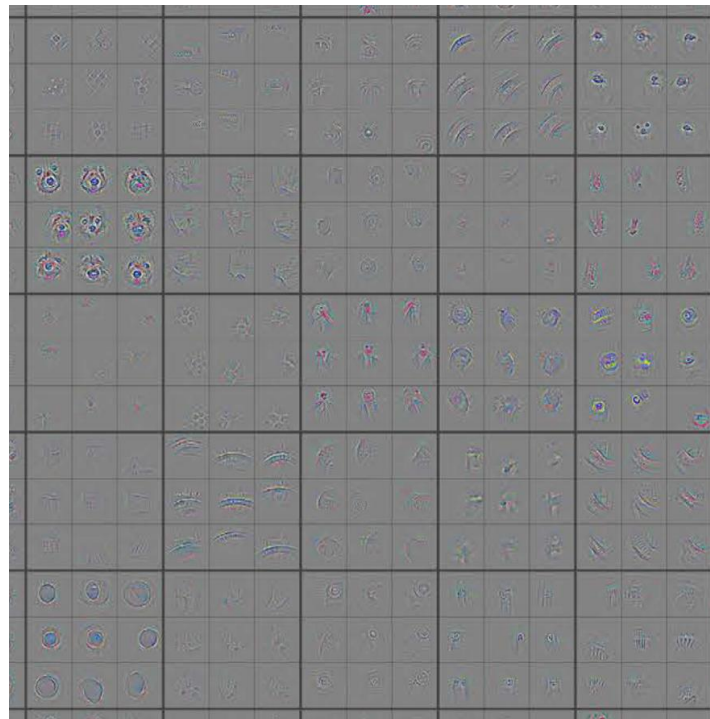


Taken from:

M.D. Zeller, R. Fergus, „Visualizing and Understanding Convolutional Neural Networks“,
arXiv:1311.2901 v2, 13. Nov. 2013

Convolutional Networks

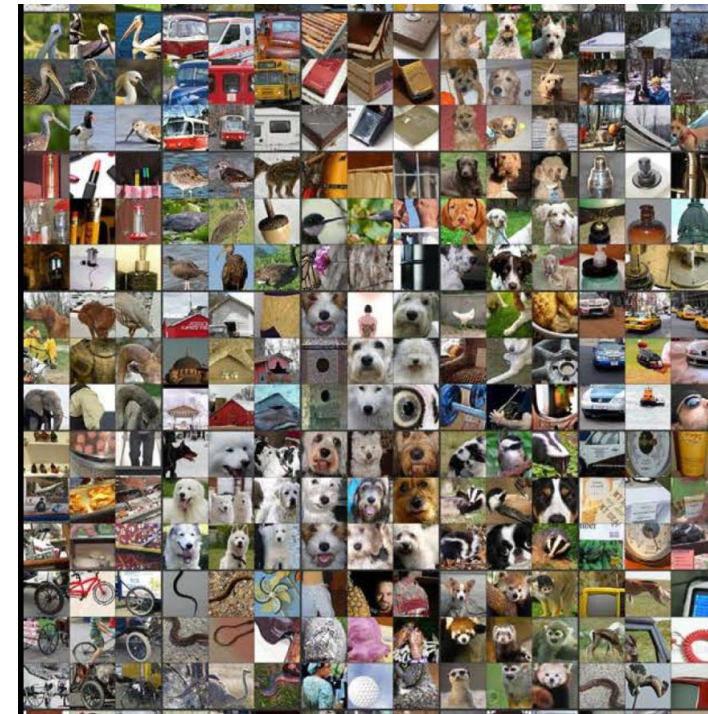
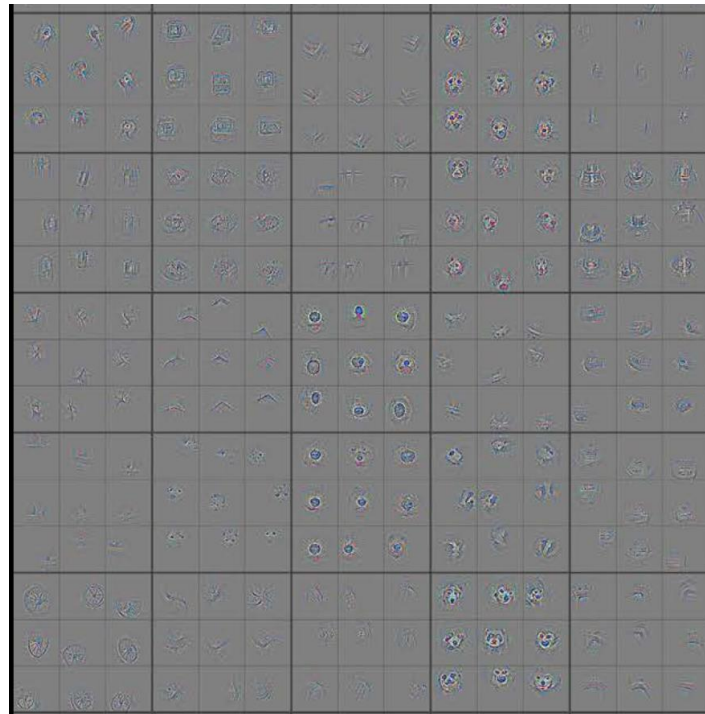
- Which features are learned in hidden layers?
 - 4. layer: textured geometries



Taken from:
M.D. Zeller, R. Fergus, „Visualizing and Understanding Convolutional Neural Networks“,
arXiv:1311.2901 v2, 13. Nov. 2013

Convolutional Networks

- Which features are learned in hidden layers?
 - 5. layer: objects



Taken from:

M.D. Zeller, R. Fergus, „Visualizing and Understanding Convolutional Neural Networks“,
arXiv:1311.2901 v2, 13. Nov. 2013

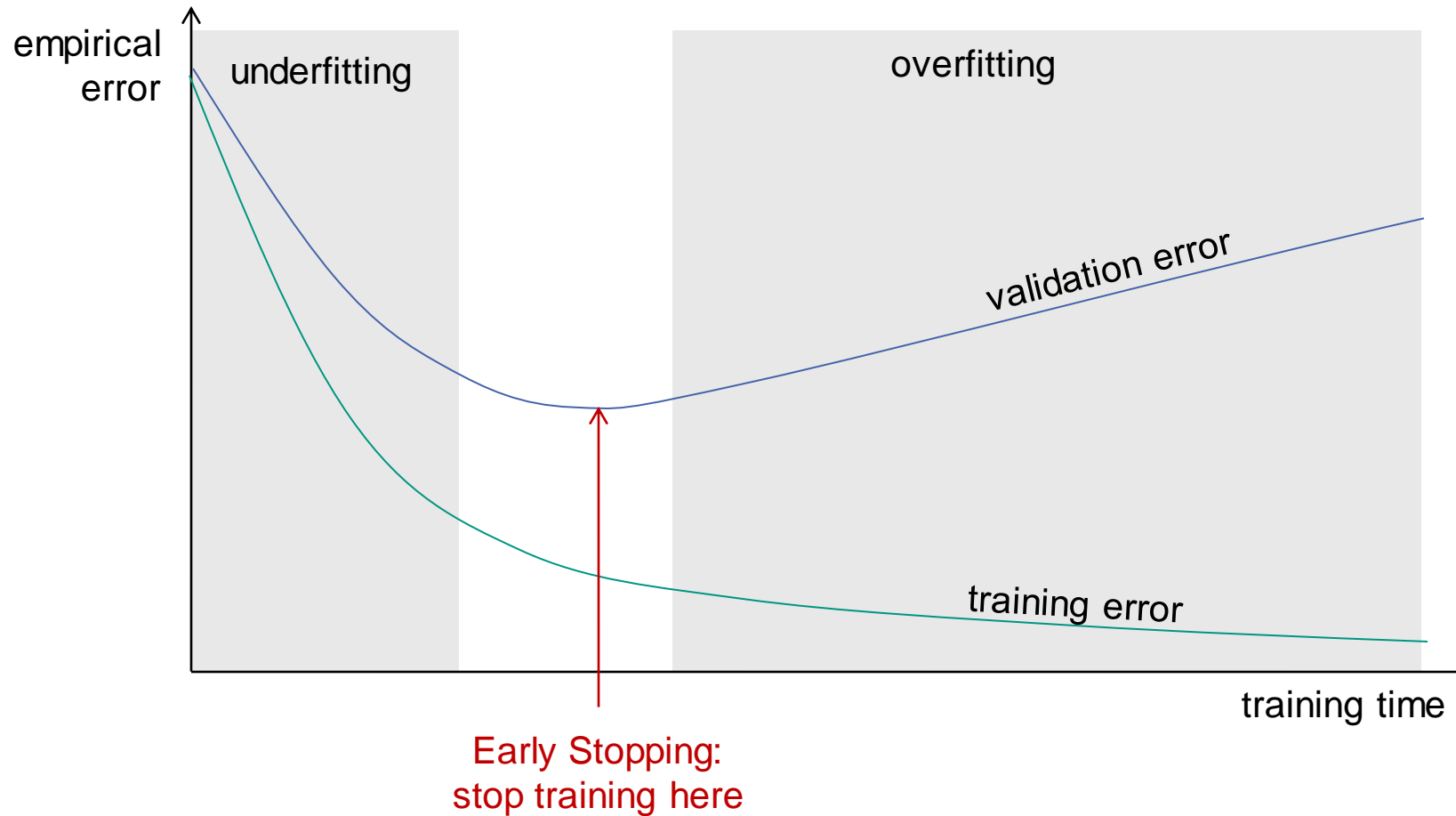
Convolutional Networks

- From layer to layer...
 - Features become more and more geometrically complex
 - Features become more and more independent of position
 - Features become more and more independent of pattern size
 - Features become more and more specific

Principles for Training MLPs

- There's no data like more data!
 - Remember data tuning
- Rigorous validation of training process
- Regularisation of training process
 - Early stopping, Weight decay/L2 regularisation, Dropout, Stochastic gradient descent, Multi task learning, Use pretrained networks
- Reuse of practical knowledge (of others)
 - Successful network structures
 - Successful training processes

Typical Progression of Error during Training



Modifications of Gradient Descent

■ Stochastic gradient descent

$$\vec{w} \leftarrow \vec{w} - \varepsilon \cdot \frac{\partial}{\partial \vec{w}} \sum_{j=1}^p \text{err}(f_{MLP}^{\vec{w}}(\vec{x}^{(j)}), \vec{t}^{(j)}) \quad \left. \vphantom{\sum_{j=1}^p} \right\} \begin{array}{l} \text{Calculate gradient from} \\ \text{all training examples.} \end{array}$$

$$\vec{w} \leftarrow \vec{w} - \varepsilon \cdot \frac{\partial}{\partial \vec{w}} \sum_{j \in S} \text{err}(f_{MLP}^{\vec{w}}(\vec{x}^{(j)}), \vec{t}^{(j)}) \quad \left. \vphantom{\sum_{j \in S}} \right\} \begin{array}{l} \text{Calculate gradient from} \\ \text{subset of all training} \\ \text{examples. Subsets typically} \\ \text{cycle through all examples.} \end{array}$$

with $S \subseteq \{1, \dots, p\}$

Advantages:

- Speed up
- A little bit less overfitting