

Deep Learning

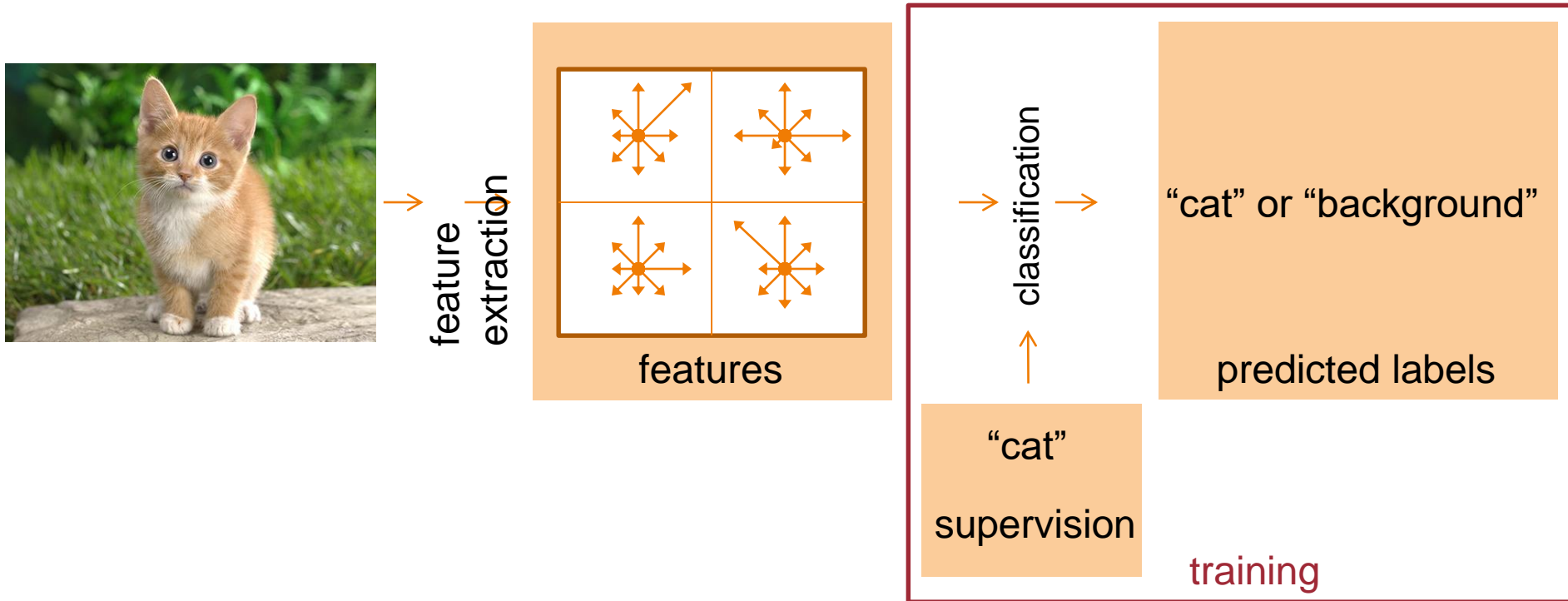
CMSC 422

MARINE CARPUAT

marine@cs.umd.edu

Based on slides by Vlad Morariu

Standard Application of Machine Learning to Computer Vision



- Features: e.g., Scale Invariant Feature Transform(SIFT)
- Classifiers: SVM, Random Forests, KNN, ...
- Features are hand-crafted, not trained
 - eventually limited by feature quality

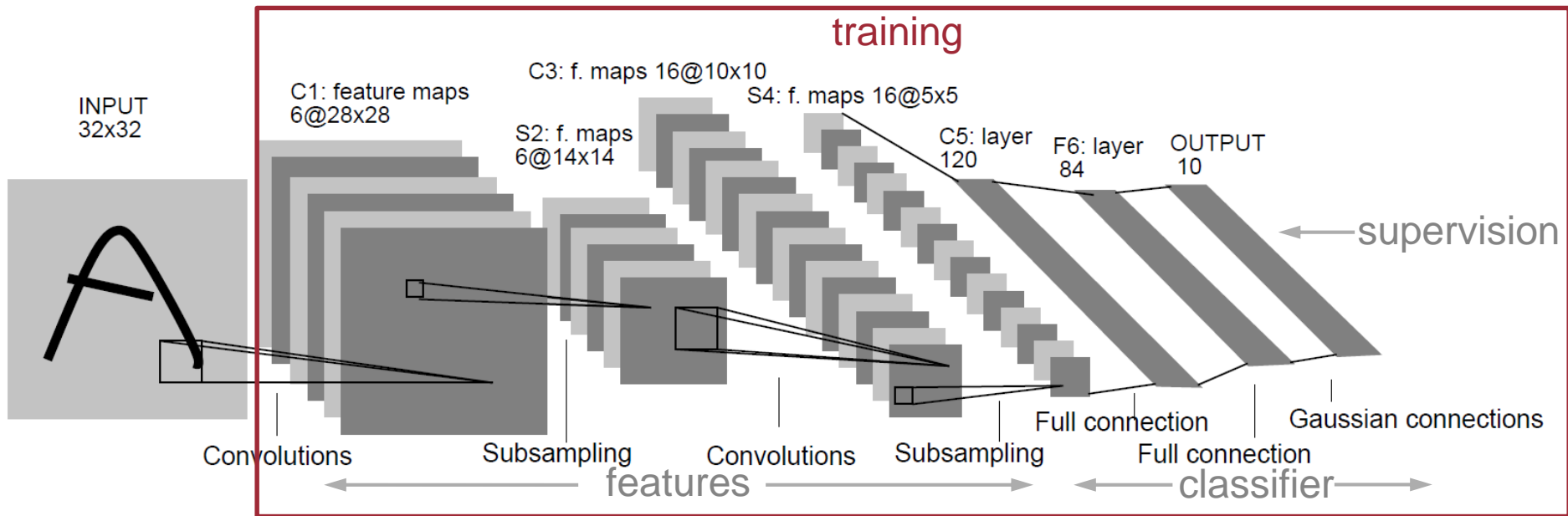


Image credit: LeCun, Y., Bottou, L., Bengio, Y., Haffner, P. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 1998.

- **Deep learning**

- multiple layer neural networks
- learn features and classifiers directly ("end-to-end" training)
- breakthrough in Computer Vision, now in other AI areas

Speech Recognition

According to Microsoft's
speech group:

Using DL

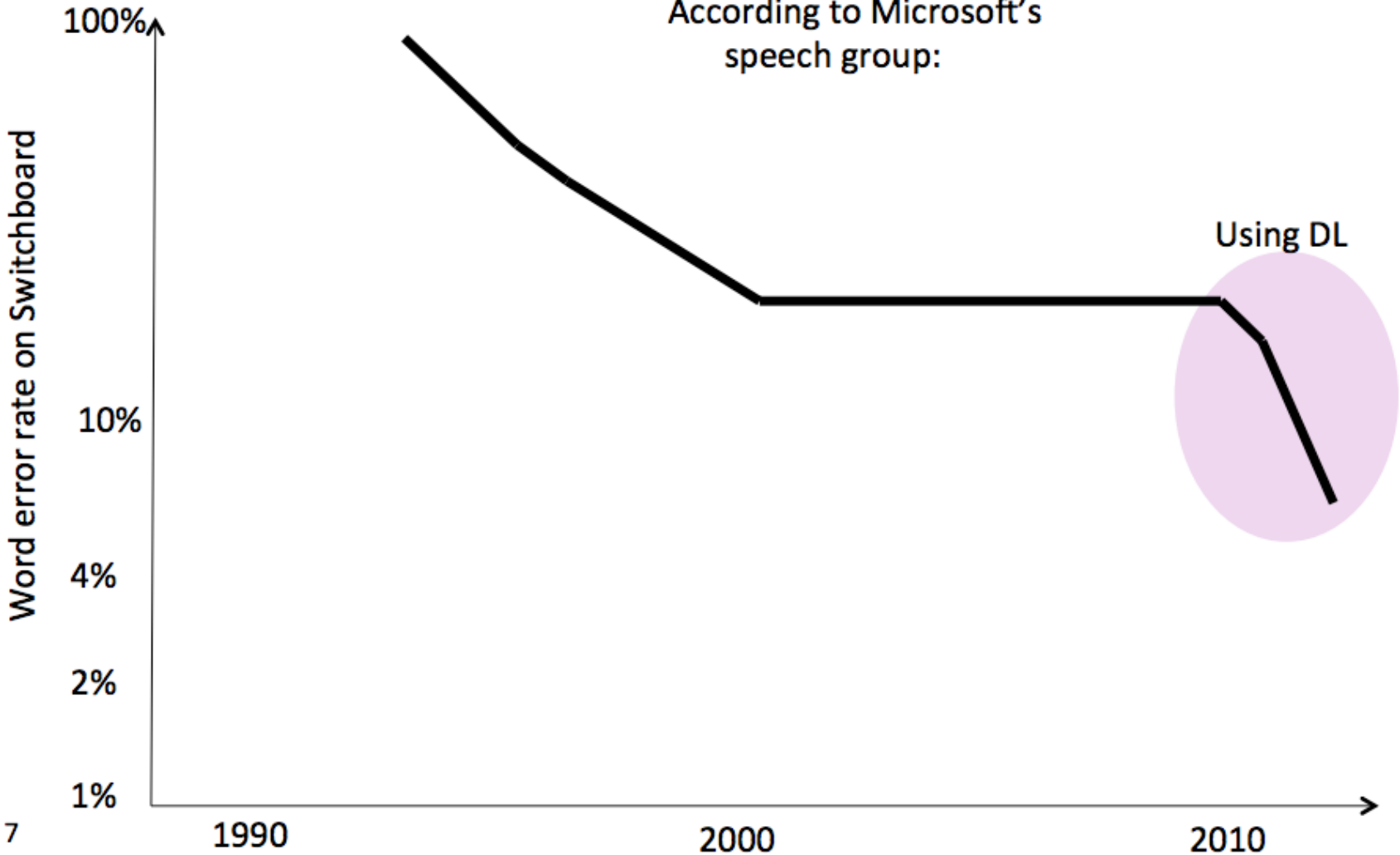


Image Classification Performance

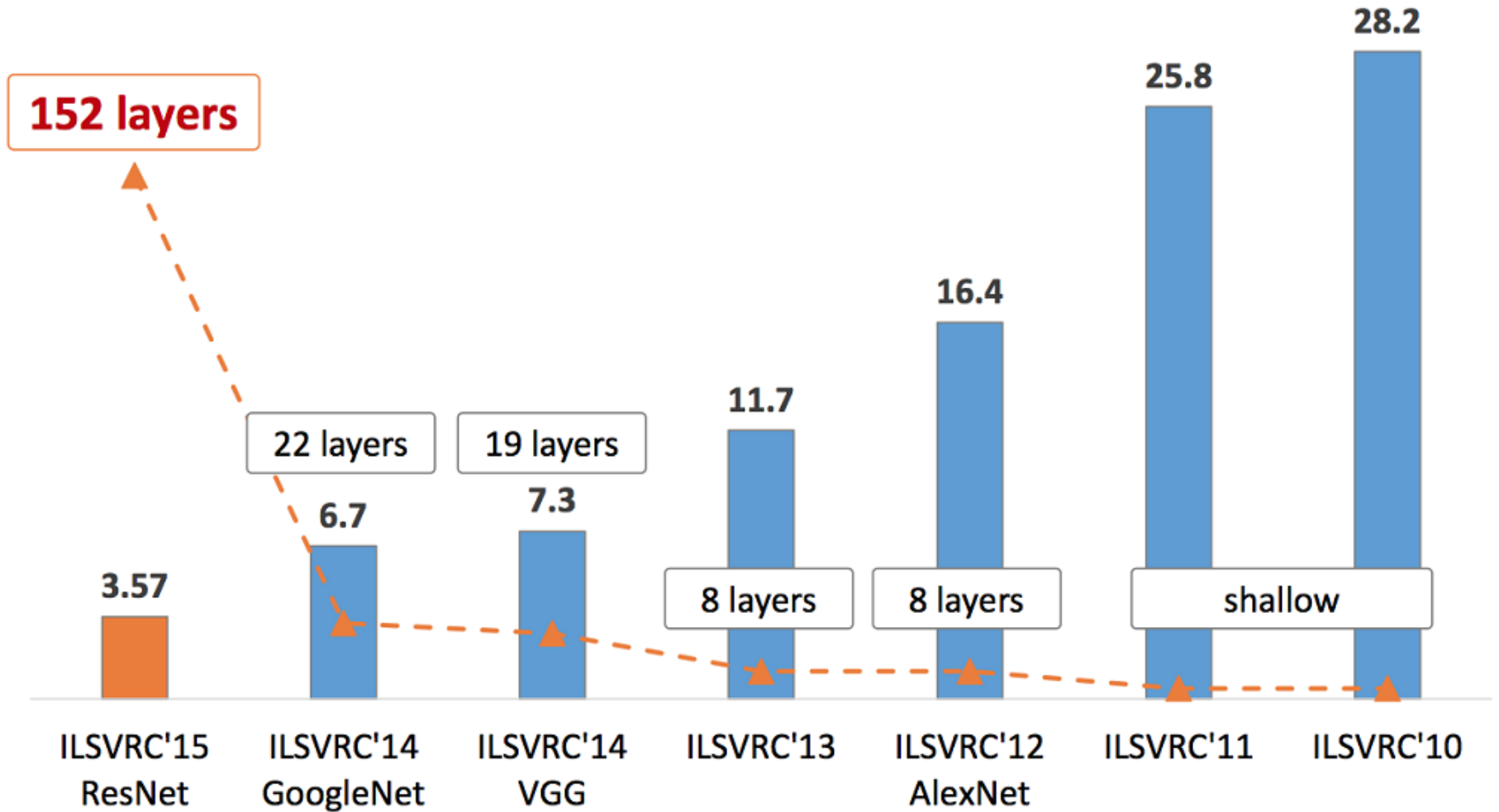


Image Classification Top-5 Errors (%)

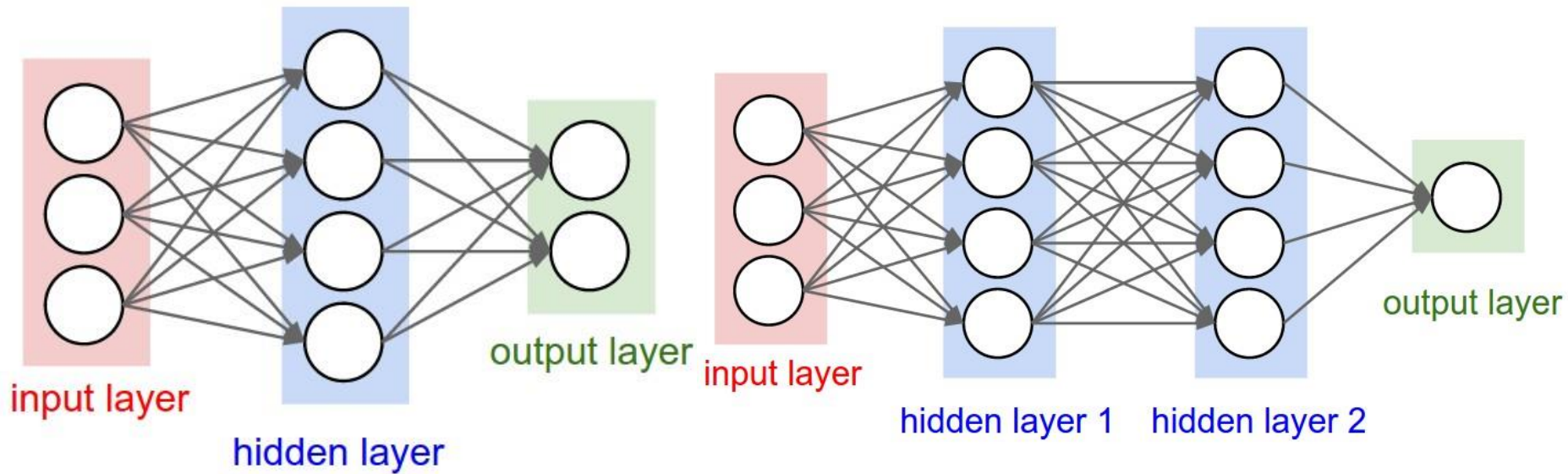
Figure from: K. He, X. Zhang, S. Ren, J. Sun. "Deep Residual Learning for Image Recognition". arXiv 2015. (slides)

Slide credit: Bohyung Han

Today's lecture: key concepts

- Convolutional Neural Networks
- Revisiting Backpropagation and Gradient Descent for Deep Networks

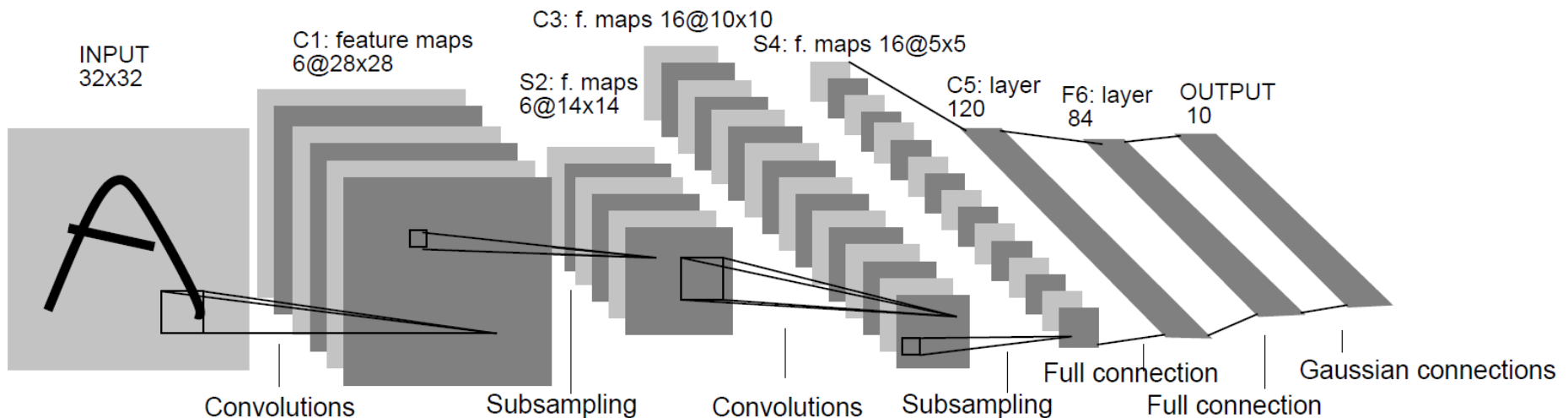
Multi-Layer Perceptron (MLP)



Neural Networks Applied to Vision

LeCun, Y; Boser, B; Denker, J; Henderson, D; Howard, R; Hubbard, W; Jackel, L, "Backpropagation Applied to Handwritten Zip Code Recognition," in Neural Computation, 1989

- USPS digit recognition, later check reading
- Convolution, pooling ("weight sharing"), fully connected layers



Architecture overview

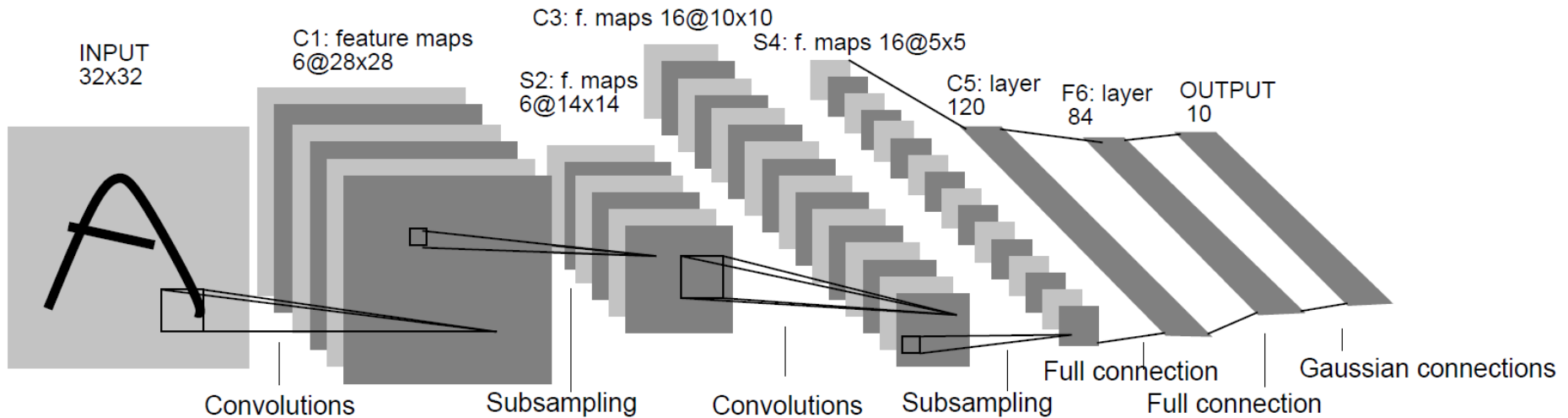


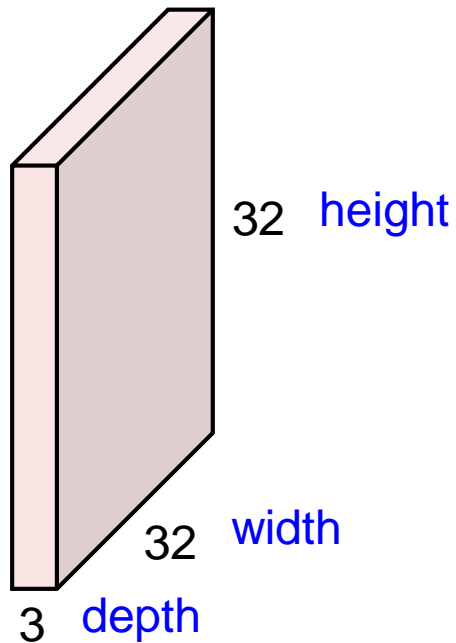
Image credit: LeCun, Y., Bottou, L., Bengio, Y., Haffner, P. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 1998.

Components:

- Convolution layers
- Pooling/Subsampling layers
- Fully connected layers

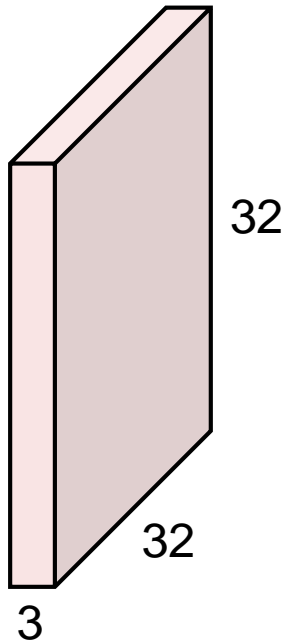
Convolutional Layer

32x32x3 image



Convolutional Layer

32x32x3 image

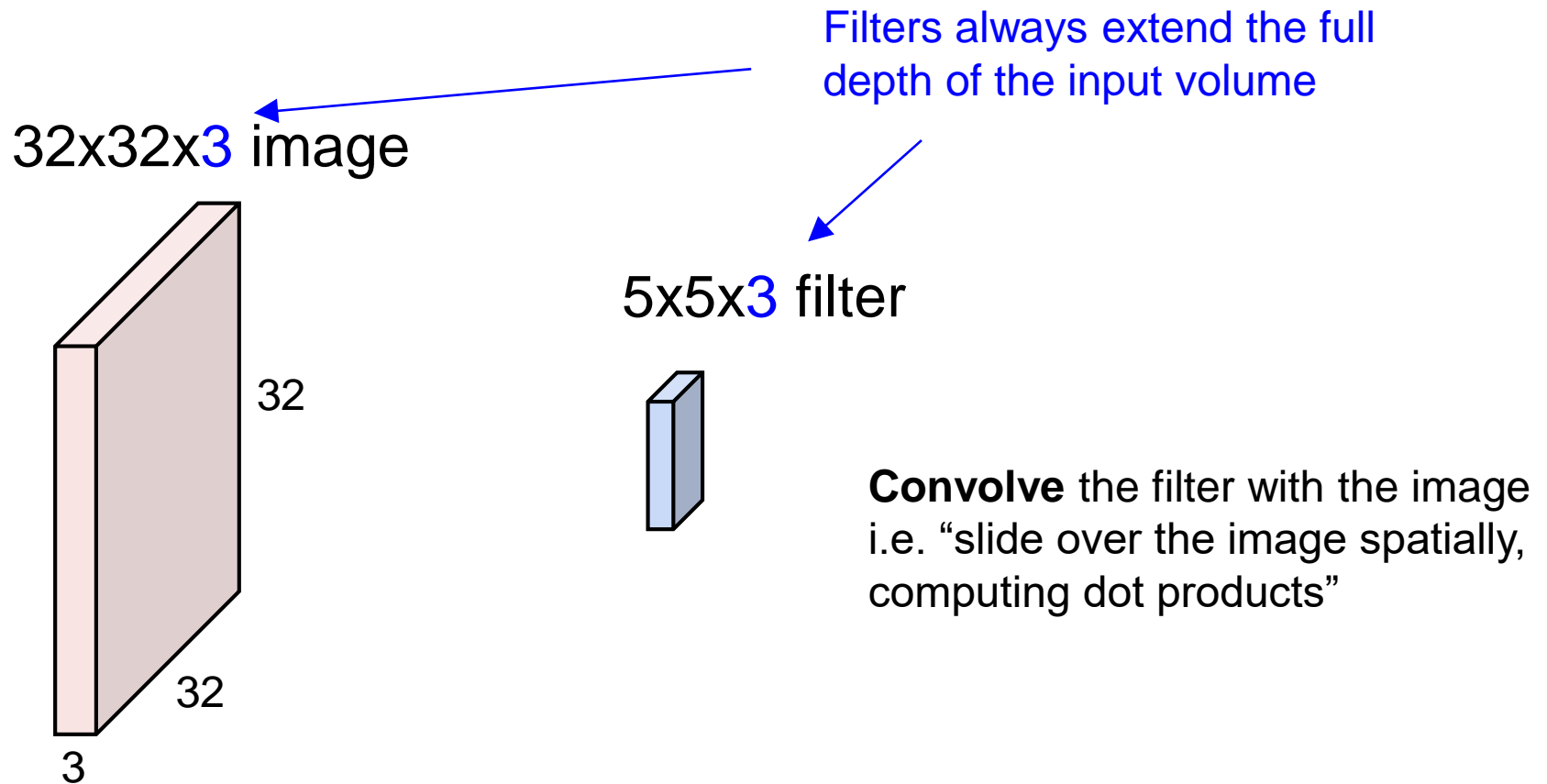


5x5x3 filter

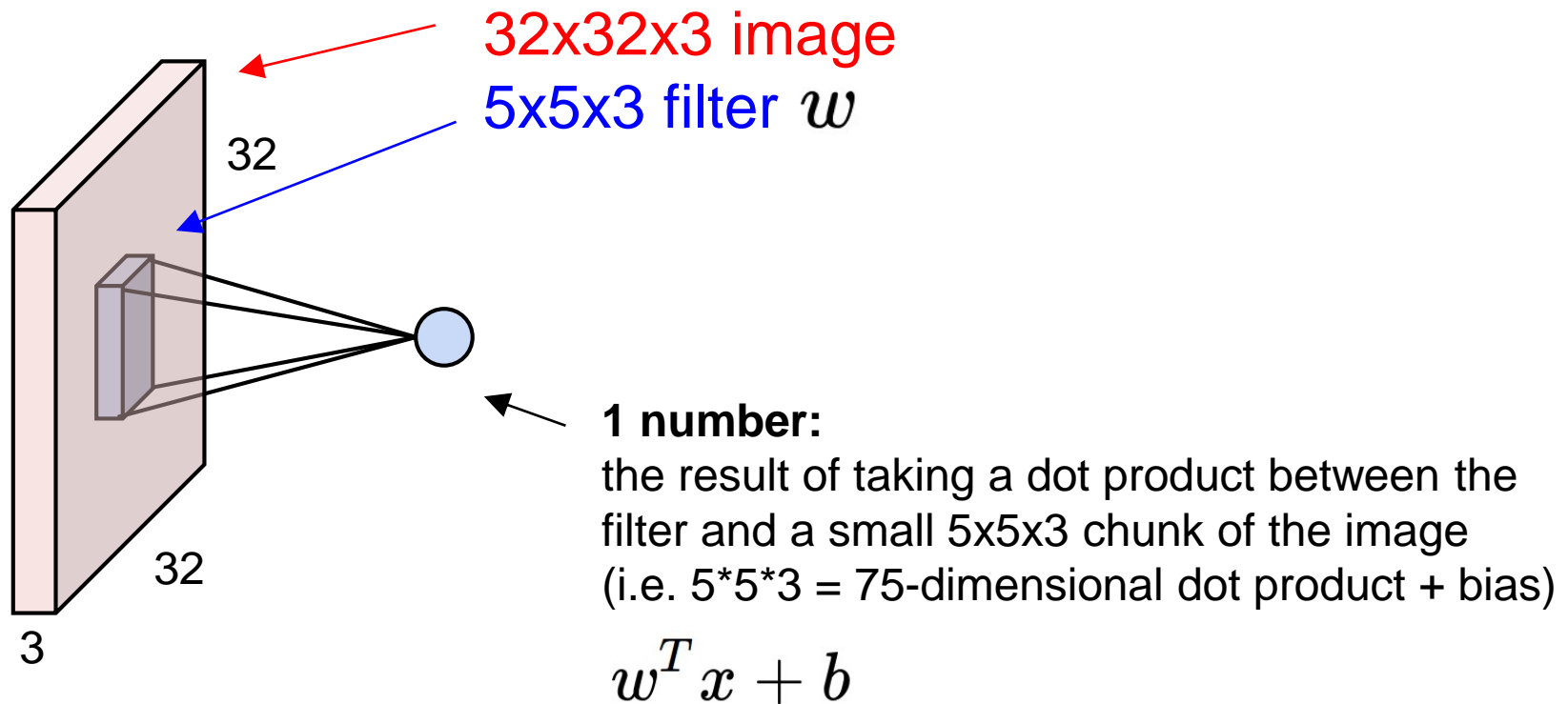


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

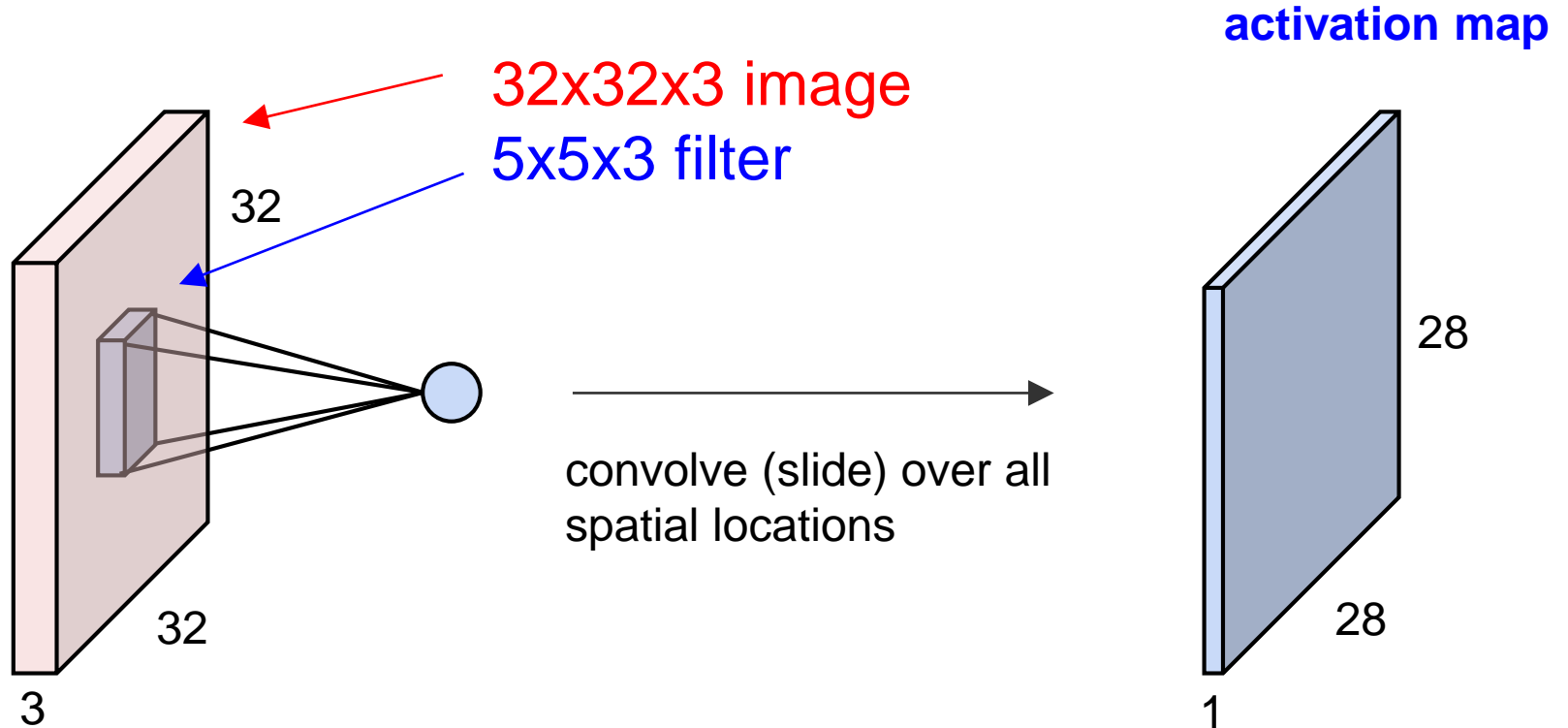
Convolutional Layer



Convolutional Layer

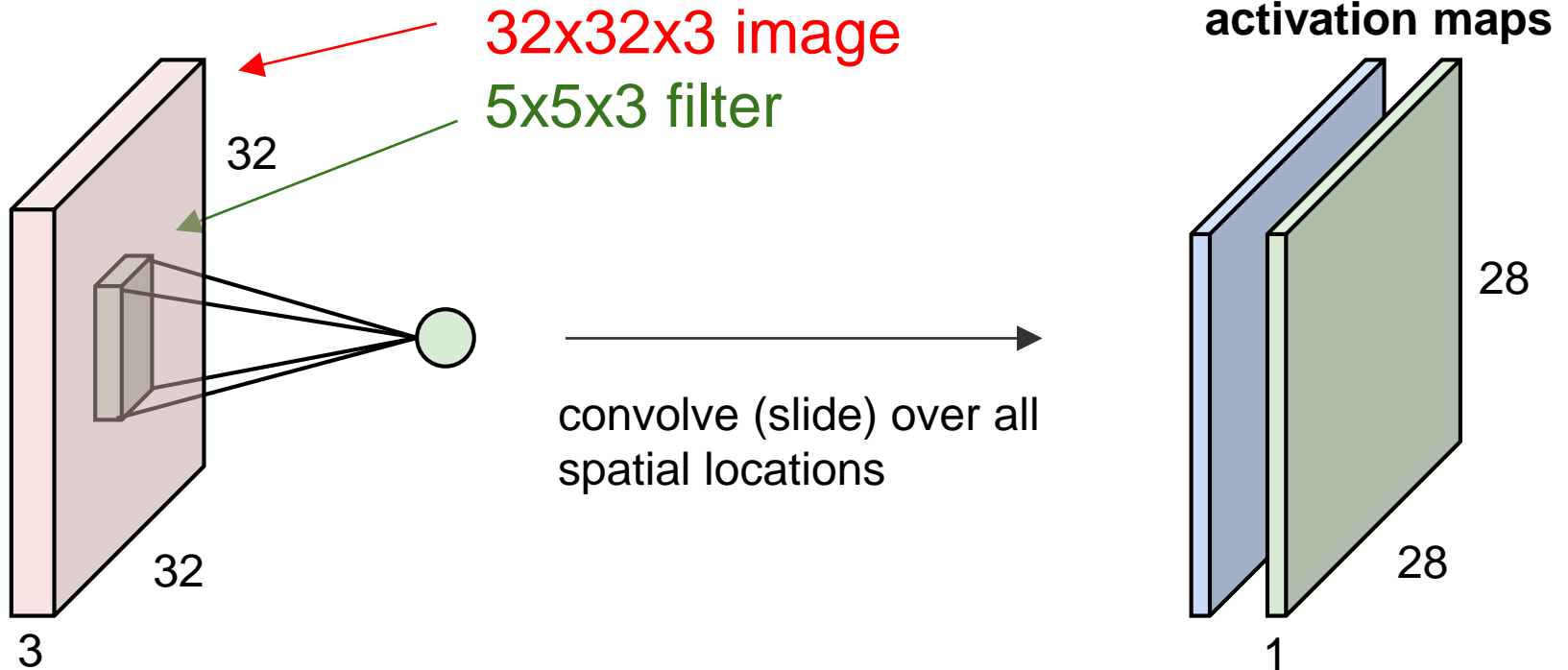


Convolutional Layer



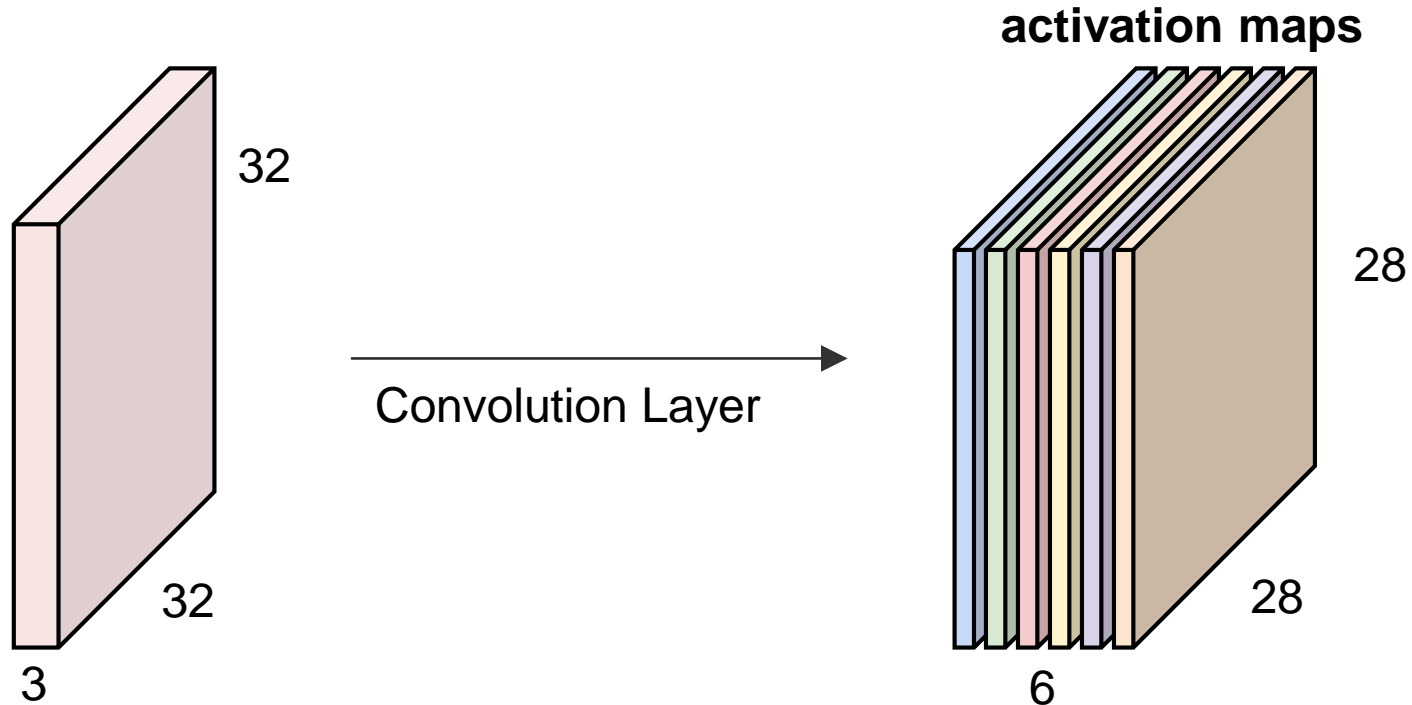
Convolutional Layer

consider a second, **green** filter



Convolutional Layer

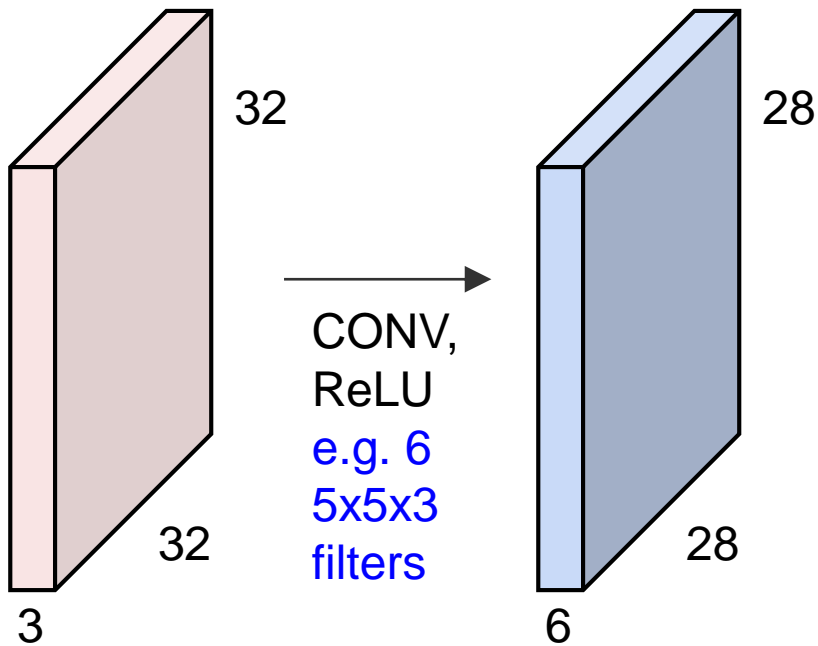
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

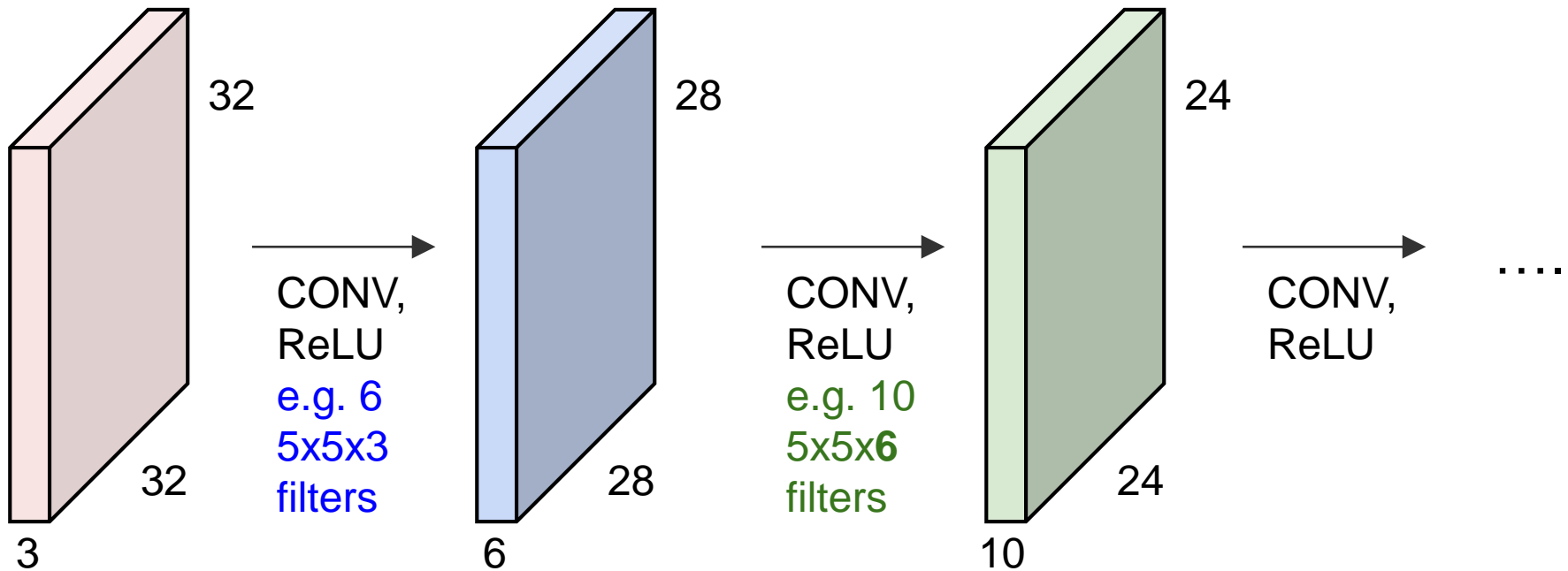
Convolutional Layer

ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

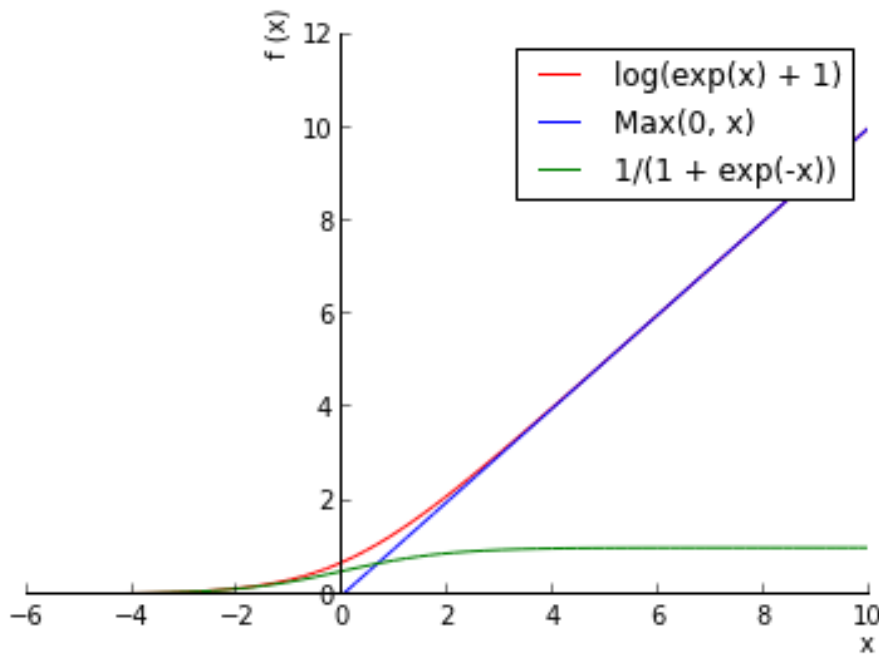


Convolutional Layer

ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



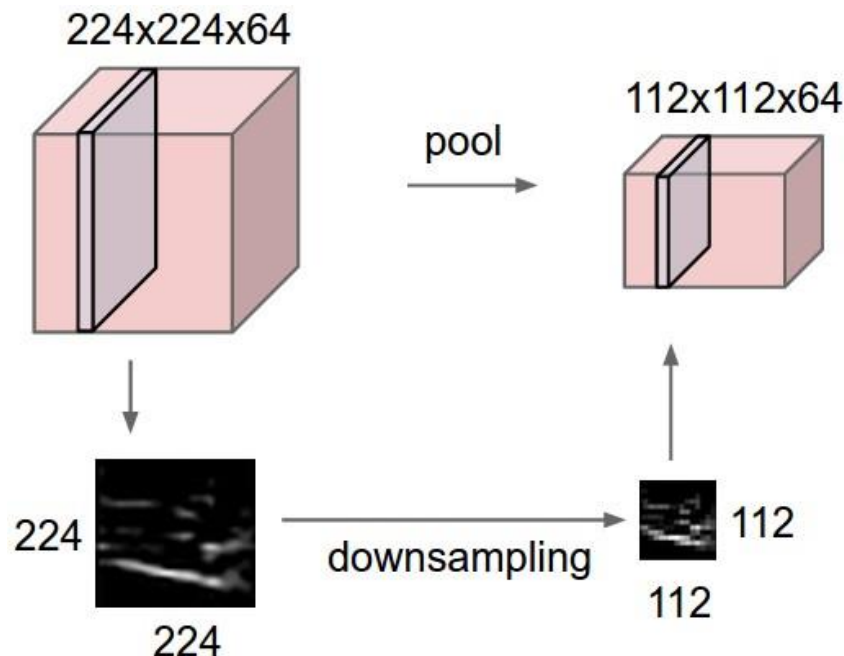
Rectified Linear Units (ReLU)



- Use rectified linear function instead of sigmoid
 $\text{ReL}(x) = \max(0, x)$
- Advantages
 - Fast
 - No vanishing gradients

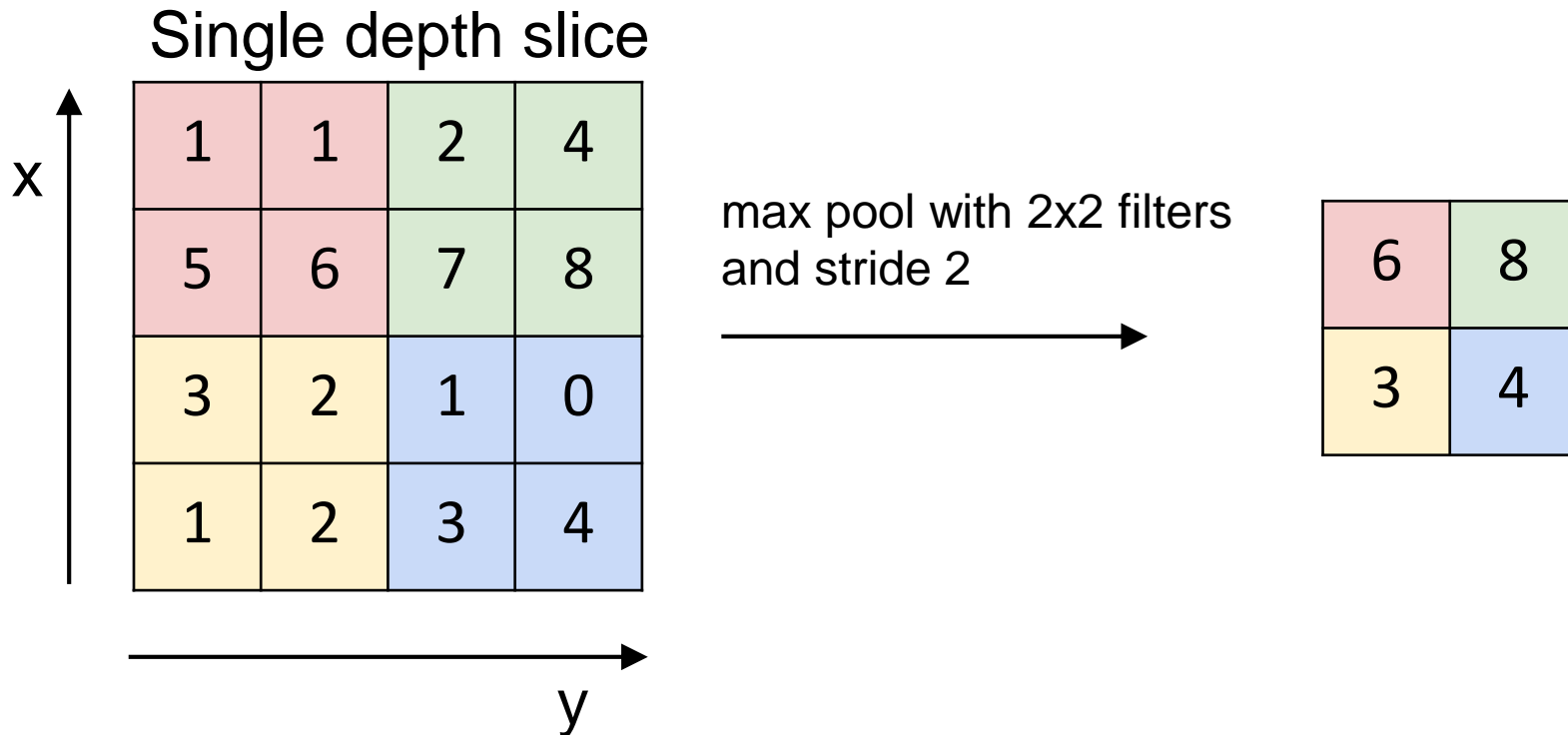
Pooling Layer

- makes the representations smaller and more manageable
- operates over each activation map independently



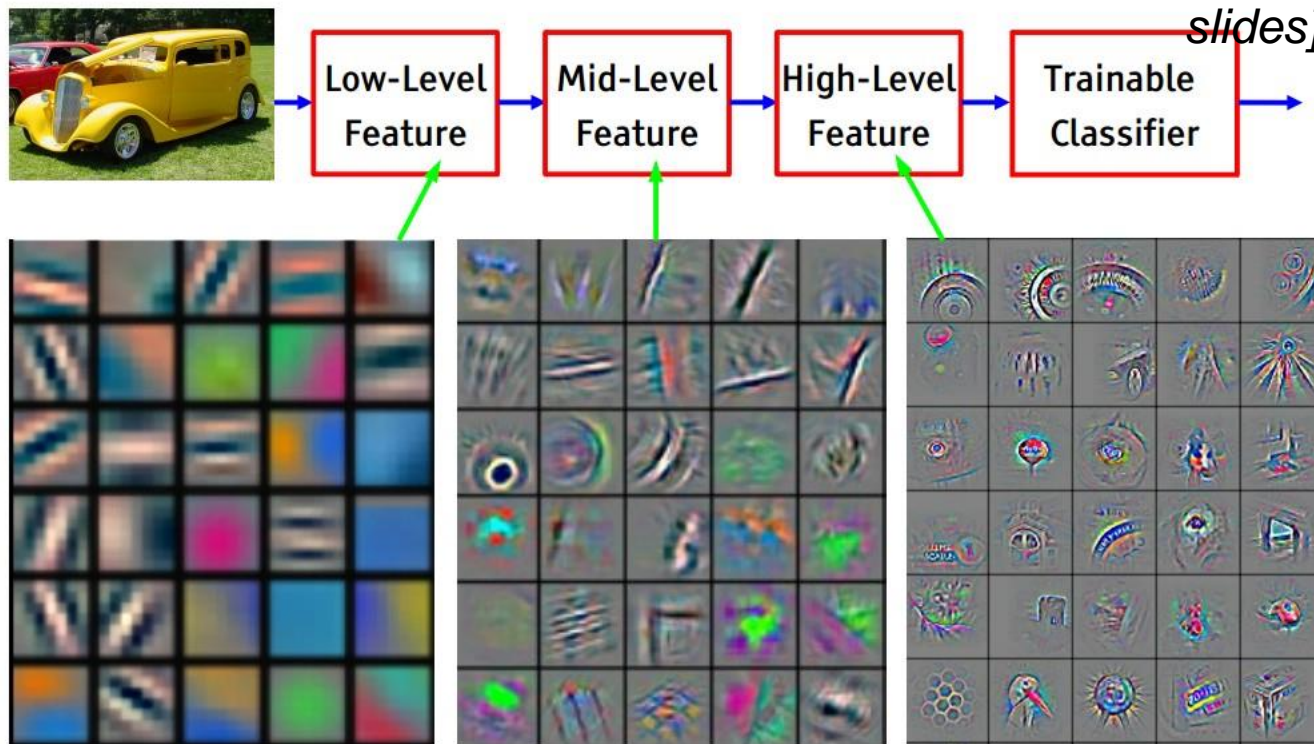
Pooling Layer

MAX POOLING



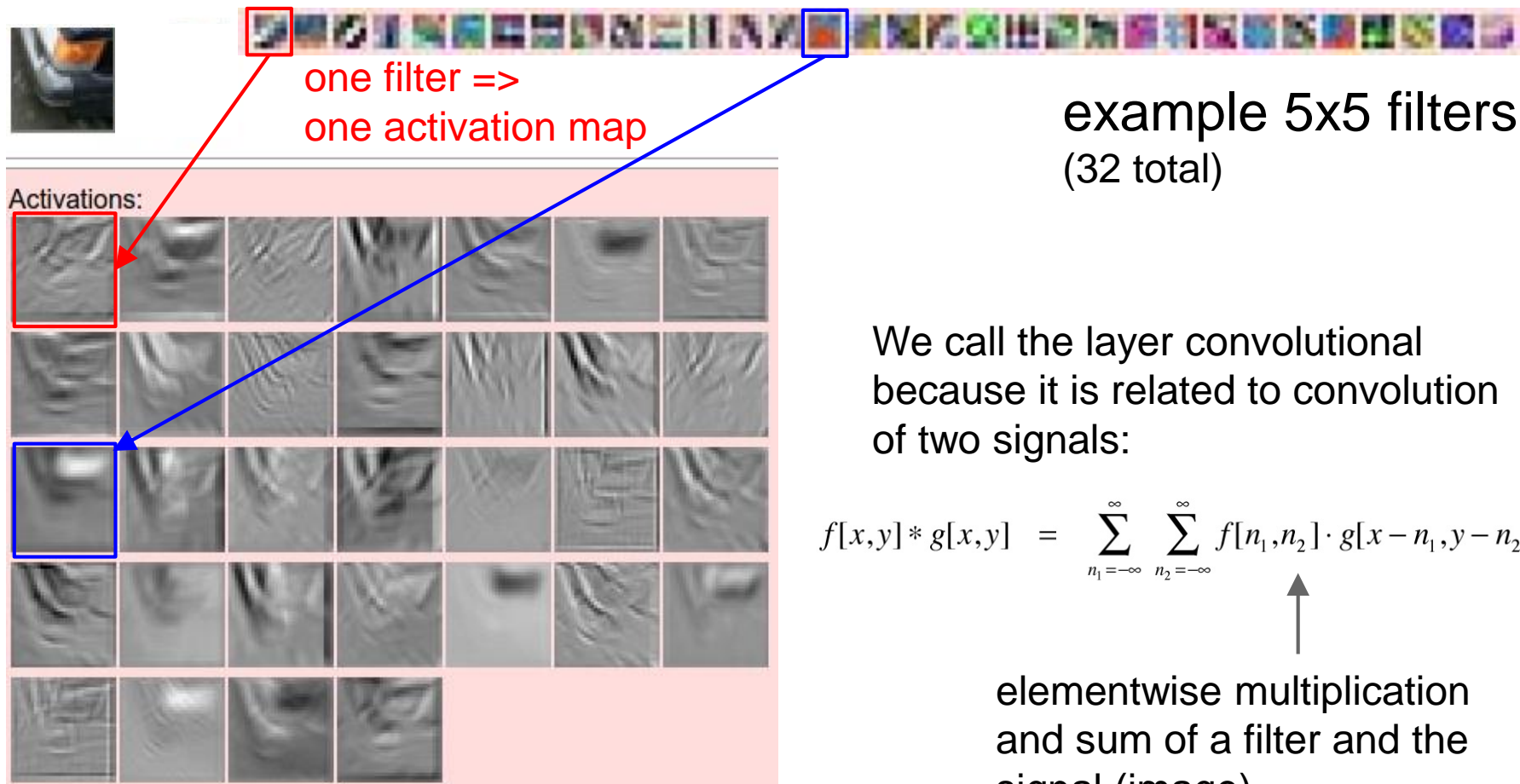
Convolutional filter visualization

*[From recent
Yann LeCun
slides]*



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

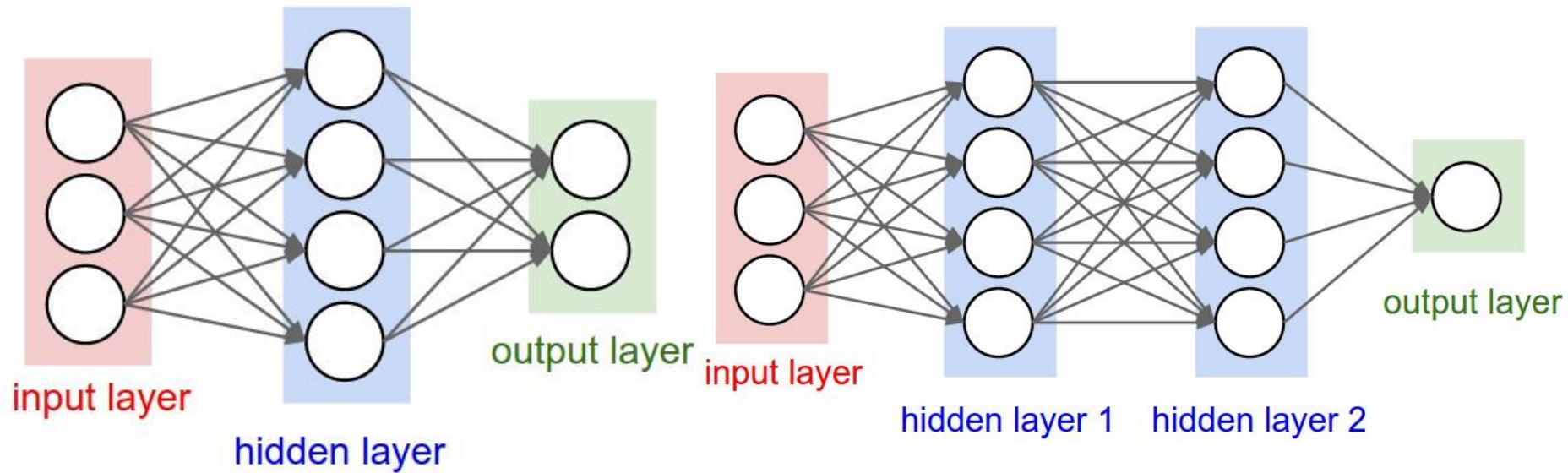
Convolutional filter visualization



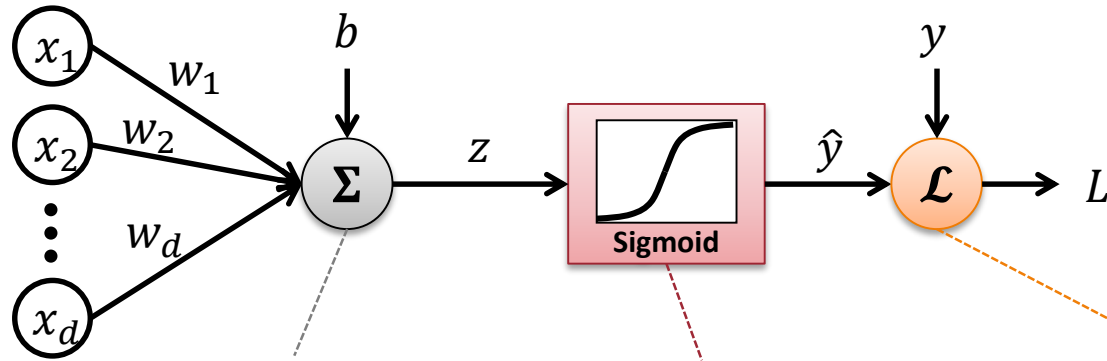
Today's lecture: key concepts

- Convolutional Neural Networks
- Revisiting Backpropagation and Gradient Descent for Deep Networks

Multi-Layer Perceptron (MLP)



Single neuron gradient



$$z = b + \sum_i w_i x_i$$

$$\hat{y} = \frac{1}{1 + e^{-z}}$$

$$L = \frac{1}{2} \sum_n (y^n - \hat{y}^n)^2$$

$$\frac{\partial z}{\partial w_i} = x_i$$

$$\frac{d\hat{y}}{dz} = \hat{y}(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}^n} = -(y^n - \hat{y}^n)$$

Chain rule: If $y = f(x)$, $z = g(y)$, then $\frac{dz}{dx} = \frac{dy}{dx} \frac{dz}{dy}$

$$\frac{\partial L}{\partial w_i} = \sum_n \frac{\partial \hat{y}^n}{\partial w_i} \frac{\partial L}{\partial \hat{y}^n} = \sum_n \frac{\partial z^n}{\partial w_i} \frac{d\hat{y}^n}{dz^n} \frac{\partial L}{\partial \hat{y}^n} = - \sum_n x_i^n \hat{y}^n (1 - \hat{y}^n) (y^n - \hat{y}^n)$$

Single neuron training

for $t = 1, \dots, T$

$$\hat{y}^n = f(\mathbf{x}^n, \mathbf{w}_t) \quad (n = 1, \dots, N)$$

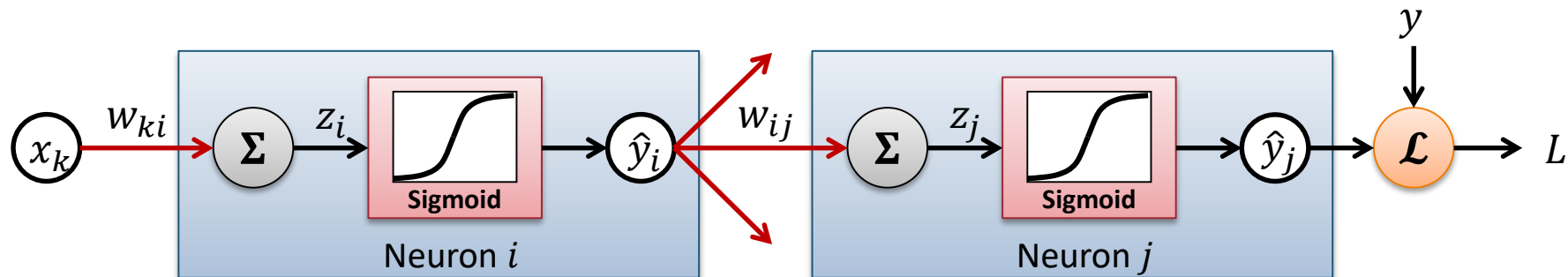
$$\frac{\partial L}{\partial w_i} = - \sum_n x_i^n \hat{y}^n (1 - \hat{y}^n) (y^n - \hat{y}^n) \quad (i = 1, \dots, d)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \Delta \mathbf{w}$$

endfor

an epoch

Multi-Layer: Backpropagation



$$\frac{\partial L}{\partial z_j} = \frac{d\hat{y}_j}{dz_j} \frac{\partial L}{\partial \hat{y}_j}$$

$$\frac{\partial L}{\partial \hat{y}_i} = \sum_j \frac{dz_j}{d\hat{y}_i} \frac{\partial L}{\partial z_j} = \sum_j w_{ij} \frac{\partial L}{\partial z_j} = \sum_j w_{ij} \frac{d\hat{y}_j}{dz_j} \frac{\partial L}{\partial \hat{y}_j}$$

$$\frac{\partial L}{\partial w_{ki}} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \frac{d\hat{y}_i^n}{dz_i^n} \frac{\partial L}{\partial \hat{y}_i^n} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \frac{d\hat{y}_i^n}{dz_i^n} \sum_j w_{ij} \frac{d\hat{y}_j^n}{dz_j^n} \frac{\partial L}{\partial \hat{y}_j^n}$$

Backpropagation in practice

Two passes per iteration:

- **Forward pass:** compute value of loss function (and intermediate neurons) given inputs
- **Backward pass:** propagate gradient of loss (error) backwards through the network using the chain rule

Stochastic Gradient Descent (SGD)

- Update weights for each sample

$$E = \frac{1}{2} (y^n - \hat{y}^n)^2 \quad \mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \epsilon \frac{\partial E^n}{\partial \mathbf{w}_i}$$

+ **Fast, online**

– **Sensitive to noise**

- Minibatch SGD: Update weights for a small set of samples

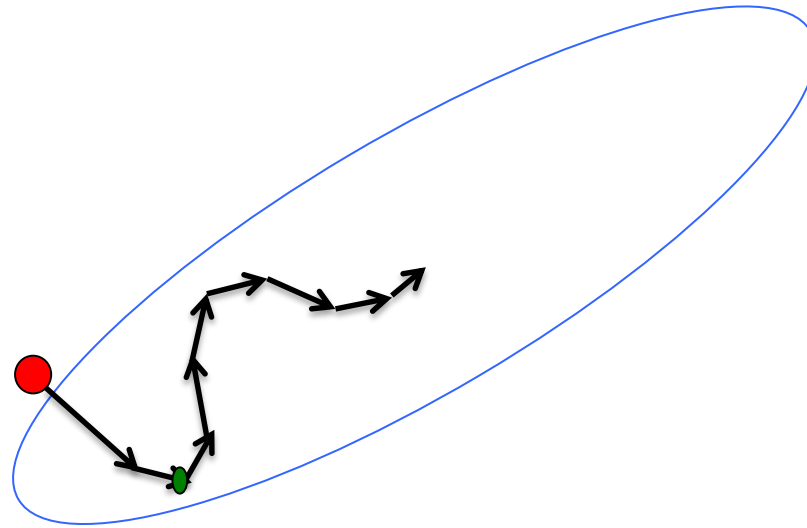
$$E = \frac{1}{2} \sum_{n \in B} (y^n - \hat{y}^n)^2 \quad \mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \epsilon \frac{\partial E^B}{\partial \mathbf{w}_i}$$

+ **Fast, online**

+ **Robust to noise**

SGD improvements: Momentum

- Remember the previous direction



+ Converge faster
+ Avoid oscillation

$$v_i(t) = \alpha v_i(t - 1) - \epsilon \frac{\partial E}{\partial w_i}(t)$$

$$\mathbf{w}(t + 1) = \mathbf{w}(t) + \mathbf{v}(t)$$

SGD improvements: Weight Decay

- Penalize the size of the weights

$$C = E + \frac{1}{2} \sum_i w_i^2$$

$$w_i(t + 1) = w_i(t) - \epsilon \frac{\partial C}{\partial w_i} = w_i(t) - \epsilon \frac{\partial E}{\partial w_i} - \lambda w_i$$

+ Improve generalization a lot!

Key concepts

- Convolutional Neural Networks
- Revisiting Backpropagation and Gradient Descent for Deep Networks

History: NN Revival in the 1980's

Backpropagation discovered in 1970's but popularized in 1986

- David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams. "Learning representations by back-propagating errors." In Nature, 1986.

MLP is a universal approximator

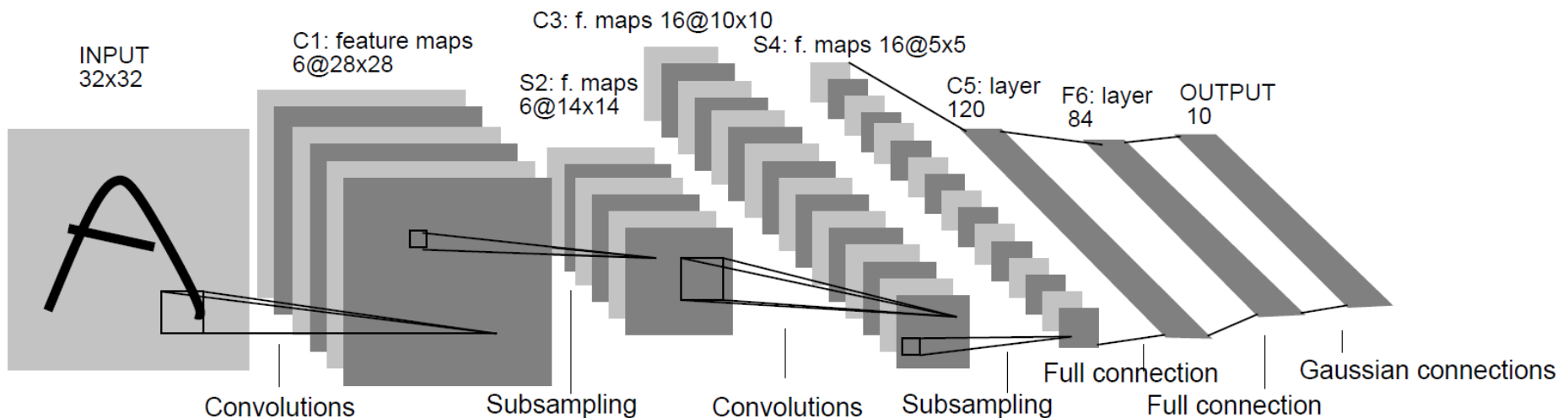
- Can approximate any non-linear function in theory, given enough neurons, data
- Kurt Hornik, Maxwell Stinchcombe, Halbert White. "Multilayer feedforward networks are universal approximators." Neural Networks, 1989

Generated lots of excitement and applications

Neural Networks Applied to Vision

LeNet – vision application

- LeCun, Y; Boser, B; Denker, J; Henderson, D; Howard, R; Hubbard, W; Jackel, L, "Backpropagation Applied to Handwritten Zip Code Recognition," in Neural Computation, 1989
- USPS digit recognition, later check reading
- Convolution, pooling ("weight sharing"), fully connected layers

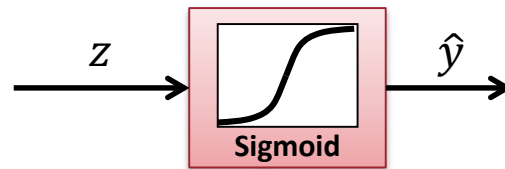


Issues in Deep Neural Networks

- Prohibitive training time
 - Especially with lots of training data
 - Many epochs typically required for optimization
 - Expensive gradient computations
- Overfitting
 - Learned function fits training data well, but performs poorly on new data (high capacity model, not enough training data)

Issues in Deep Neural Networks

Vanishing gradient problem



$$\frac{\partial E}{\partial w_{ki}} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \frac{d\hat{y}_i^n}{dz_i^n} \frac{\partial E}{\partial \hat{y}_i^n} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \boxed{\frac{d\hat{y}_i^n}{dz_i^n}} \sum_j w_{ij} \boxed{\frac{d\hat{y}_j^n}{dz_j^n}} \frac{\partial E}{\partial \hat{y}_j^n}$$

- Gradients in the lower layers are typically extremely small
- Optimizing multi-layer neural networks takes huge amount of time

New “winter” and revival in early 2000’s

New “winter” in the early 2000’s due to

- problems with training NNs
- Support Vector Machines (SVMs), Random Forests (RF) – easy to train, nice theory

Revival again by 2011-2012

- Name change (“neural networks” -> “deep learning”)
- + Algorithmic developments
 - unsupervised layer-wise pre-training
 - ReLU, dropout, layer normalization
- + Big data + GPU computing =
- Large outperformance on many datasets (Vision: ILSVRC’12)

Big Data

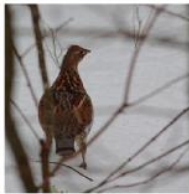
- ImageNet Large Scale Visual Recognition Challenge
 - 1000 categories w/ 1000 images per category
 - 1.2 million training images, 50,000 validation, 150,000 testing



flamingo



cock



ruffed grouse



quail



partridge

...



Egyptian cat



Persian cat



Siamese cat

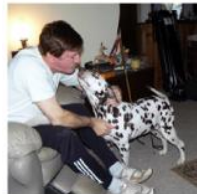


tabby



lynx

...



dalmatian



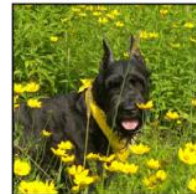
keeshond



miniature schnauzer



standard schnauzer



giant schnauzer

...

AlexNet Architecture

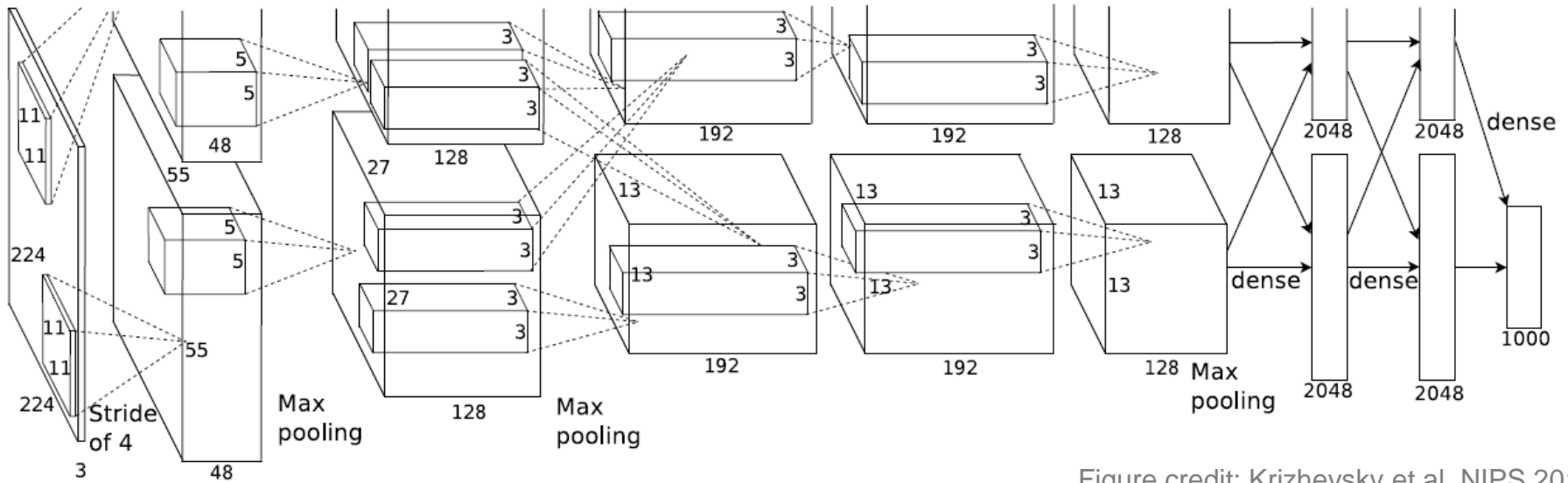


Figure credit: Krizhevsky et al, NIPS 2012.

60 million parameters!

Various tricks

- ReLU nonlinearity
- Overlapping pooling
- Local response normalization
- Dropout – set hidden neuron output to 0 with probability .5
- Data augmentation
- Training on GPUs

GPU Computing

- **Big data** and **big models** require lots of computational power
- GPUs
 - thousands of cores for parallel operations
 - multiple GPUs
 - still took about 5-6 days to train AlexNet on two NVIDIA GTX 580 3GB GPUs (much faster today)

Recurrent Neural Networks

Networks with loops

- The output of a layer is used as input for the same (or lower) layer
- Can model dynamics (e.g. in space or time)

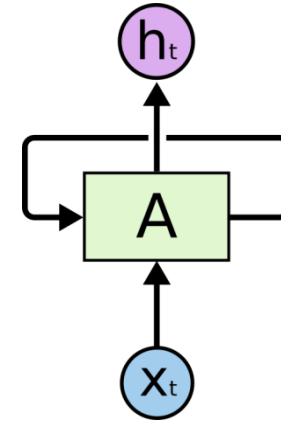


Image credit: Christopher Olah's blog <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
Sepp Hochreiter (1991), Untersuchungen zu dynamischen neuronalen Netzen, Diploma thesis. Institut f. Informatik, Technische Univ. Munich. Advisor: J. Schmidhuber.

Y. Bengio, P. Simard, P. Frasconi. Learning Long-Term Dependencies with Gradient Descent is Difficult. In TNN 1994.

Recurrent Neural Networks

Let's unroll the loops

- Now a standard feed-forward network with many layers
- Suffers from vanishing gradient problem
- In theory, can learn long term memory, in practice not (Bengio et al, 1994)

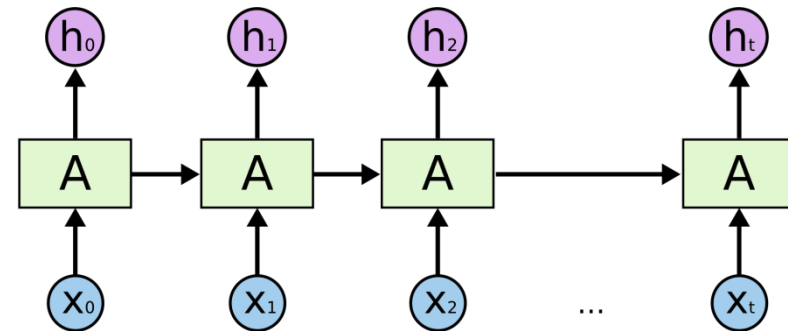
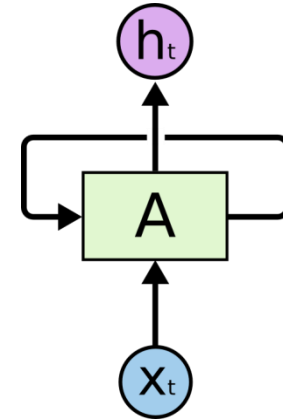
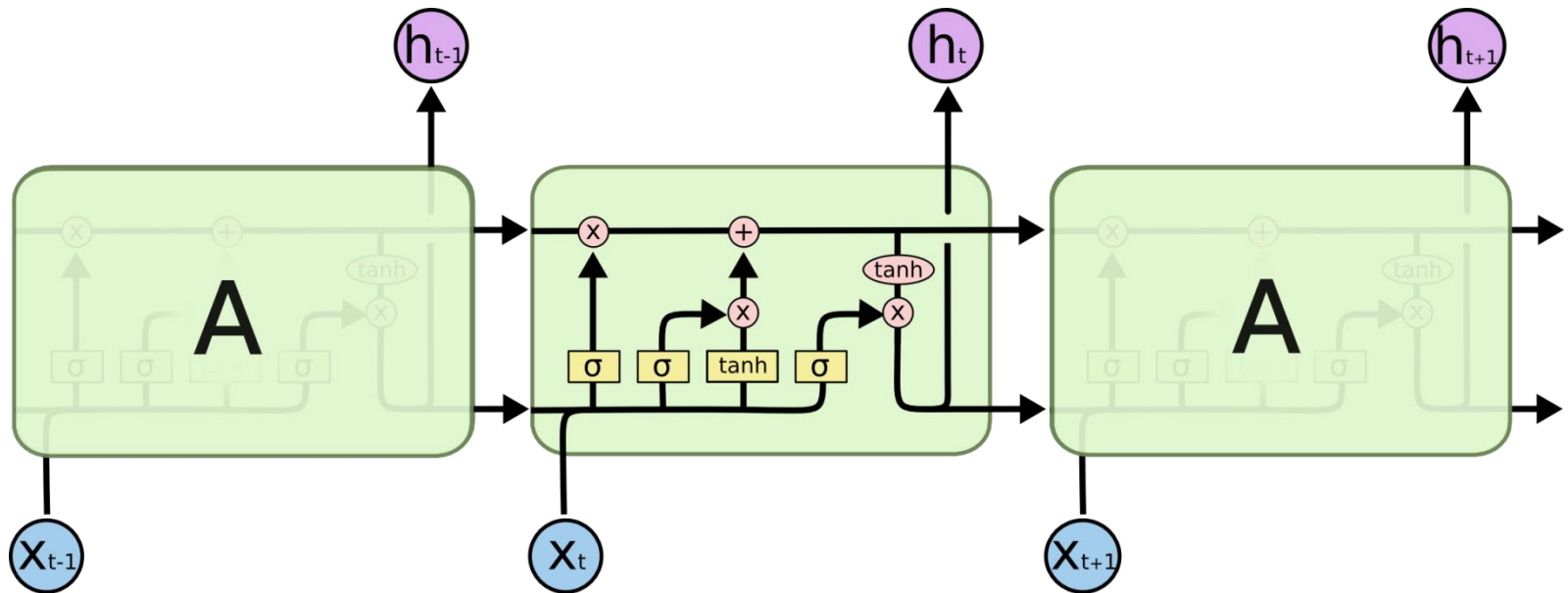


Image credit: Christopher Olah's blog <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
Sepp Hochreiter (1991), Untersuchungen zu dynamischen neuronalen Netzen, Diploma thesis. Institut f. Informatik, Technische Univ. Munich. Advisor: J. Schmidhuber.
Y. Bengio, P. Simard, P. Frasconi. Learning Long-Term Dependencies with Gradient Descent is Difficult. In TNN 1994.

Long Short Term Memory (LSTM)



- A type of RNN explicitly designed not to have the vanishing or exploding gradient problem
- Models long-term dependencies
- Memory is propagated and accessed by gates
- Used for speech recognition, language modeling ...

Unsupervised Neural Networks

Autoencoders

- Encode then decode the same input
- No supervision needed

