# Statistical Machine Learning (BE4M33SSU) Lecture 8: Deep Neural Networks

Jan Drchal

Czech Technical University in Prague
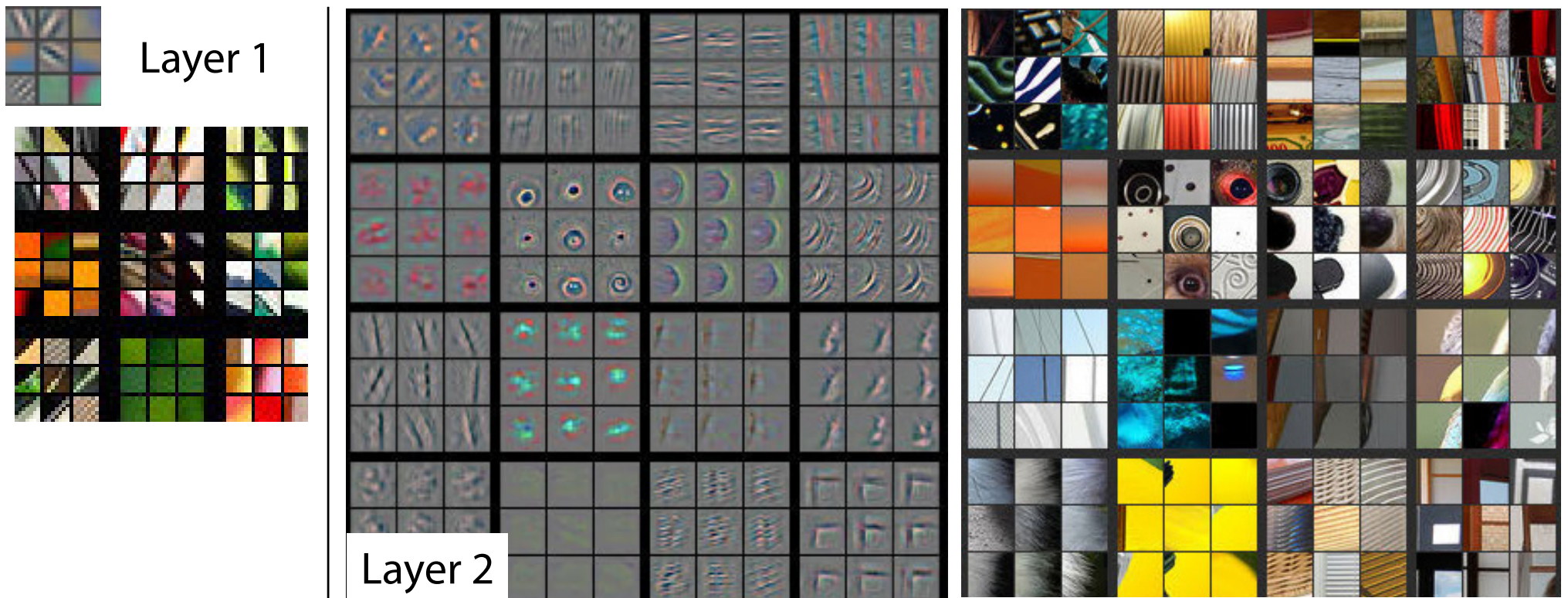Faculty of Electrical Engineering
Department of Computer Science

Topics covered in the lecture:

◆ Deep Architectures

◆ Parameter initialization

◆ Convolutional Neural Networks (CNNs)

◆ Transfer learning

- Is it better to use deep architectures rather than the shallow ones for complex nonlinear mappings?

- We know that deep architectures evolved in Nature (e.g., cortex)

- Universal approximation theorem: one layer is enough so why to bother with more layers?

- Mhaskar et al: *Learning Functions: When Is Deep Better Than Shallow*, 2016:

  - deep neural networks can have exponentially less units than shallow networks for learning the same function

  - functions such as those realized by current deep convolutional neural networks are considered

- Handcrafted features vs. automatic extraction

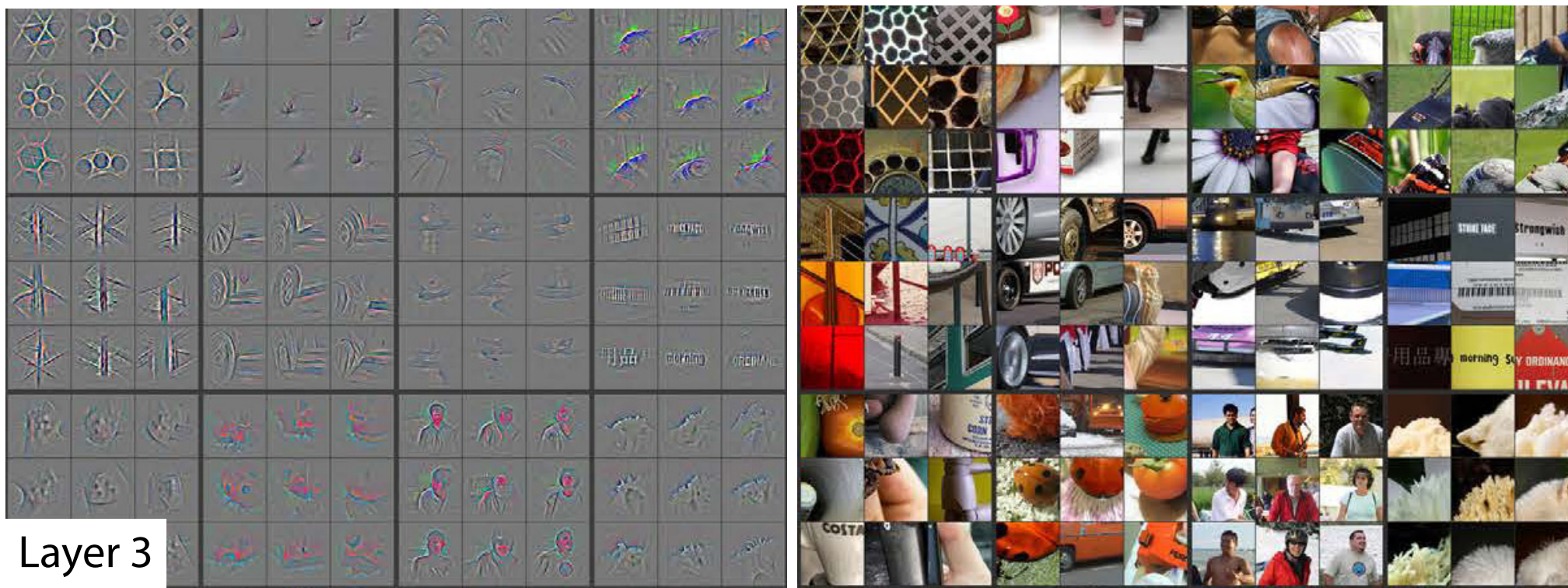- Gradually increasing complexity, intermediate representations: each successive layer brings higher abstraction
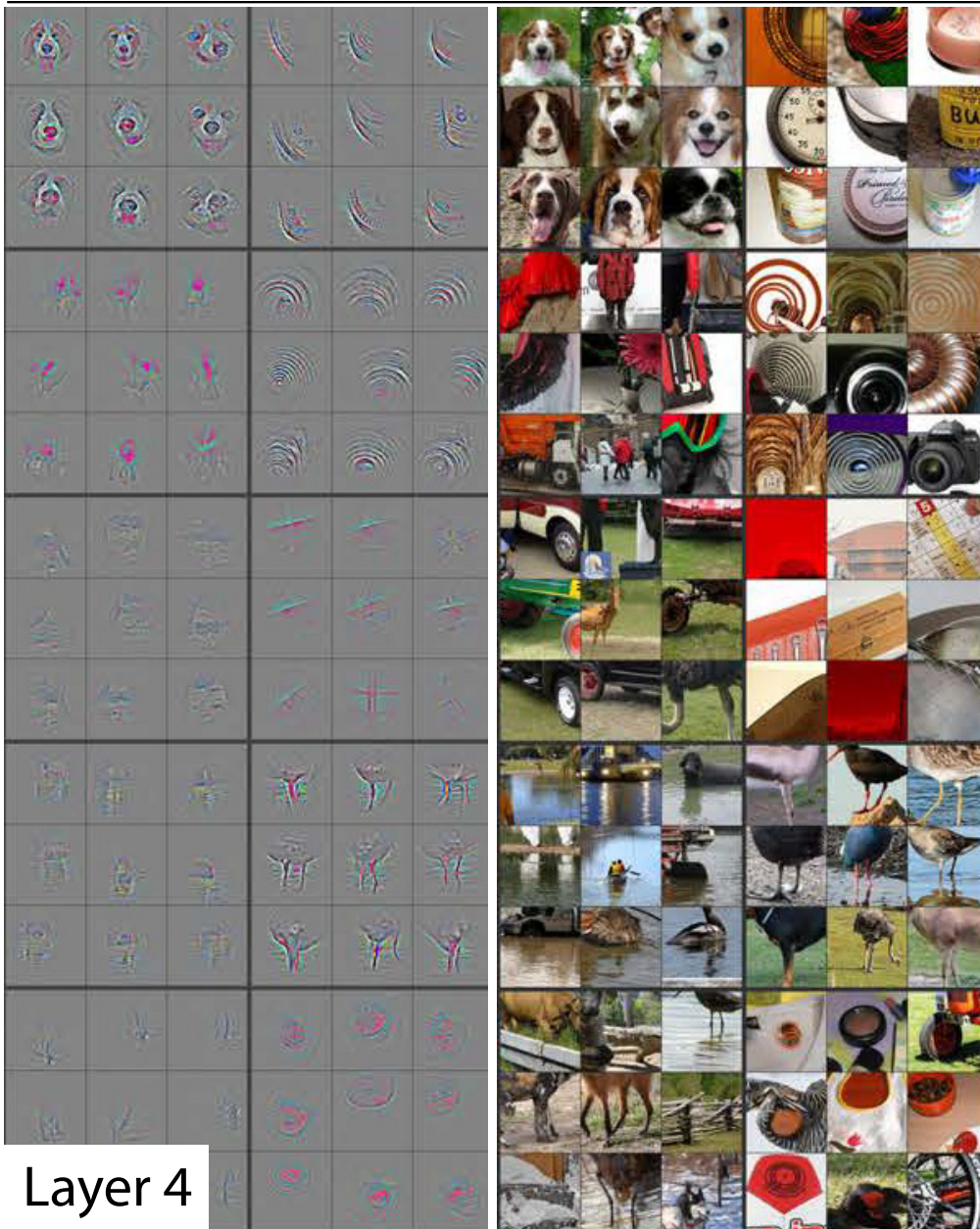
Layer 1

Layer 2

Zeiler and Fergus: *Visualizing and Understanding Convolutional Networks*, 2013
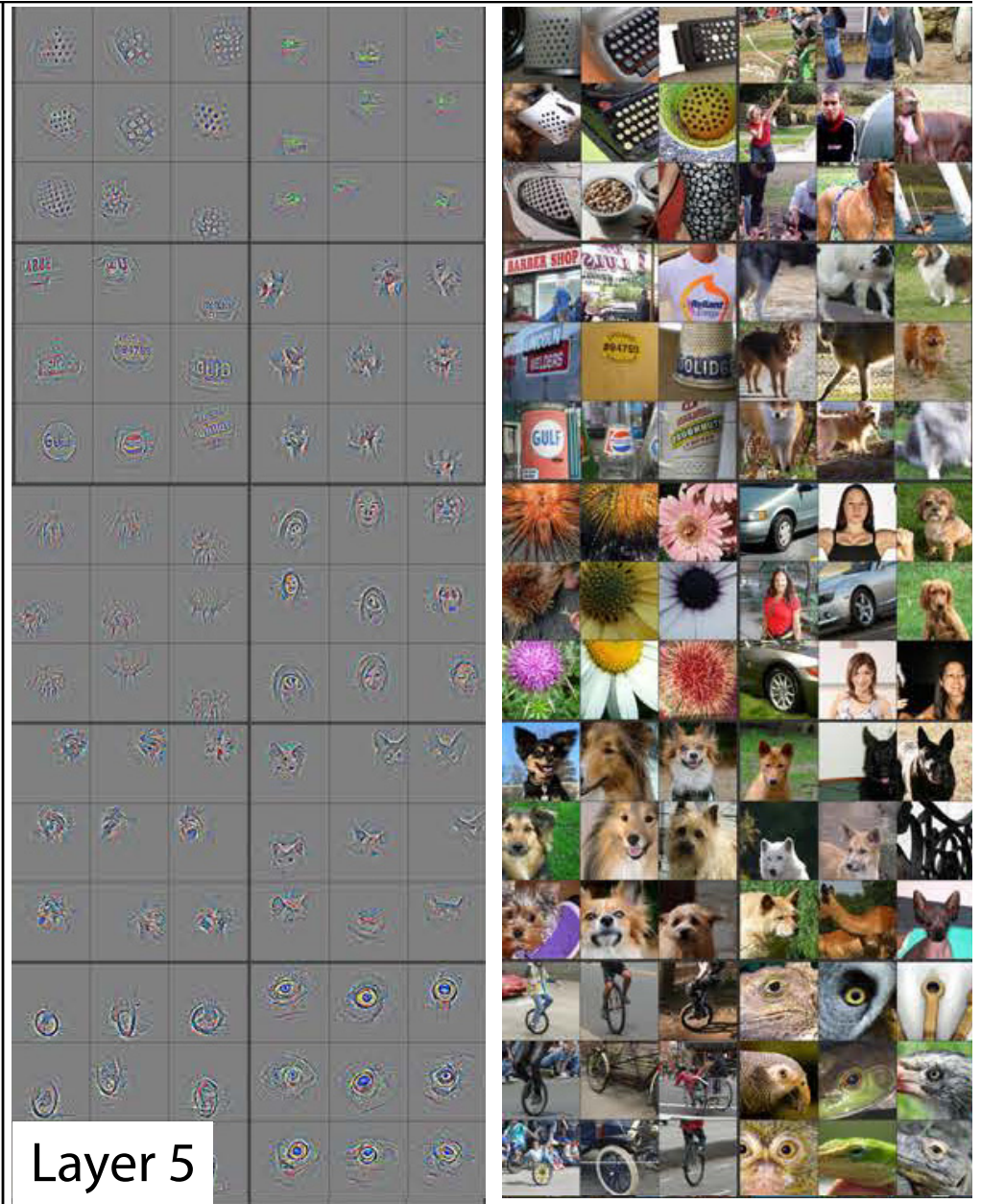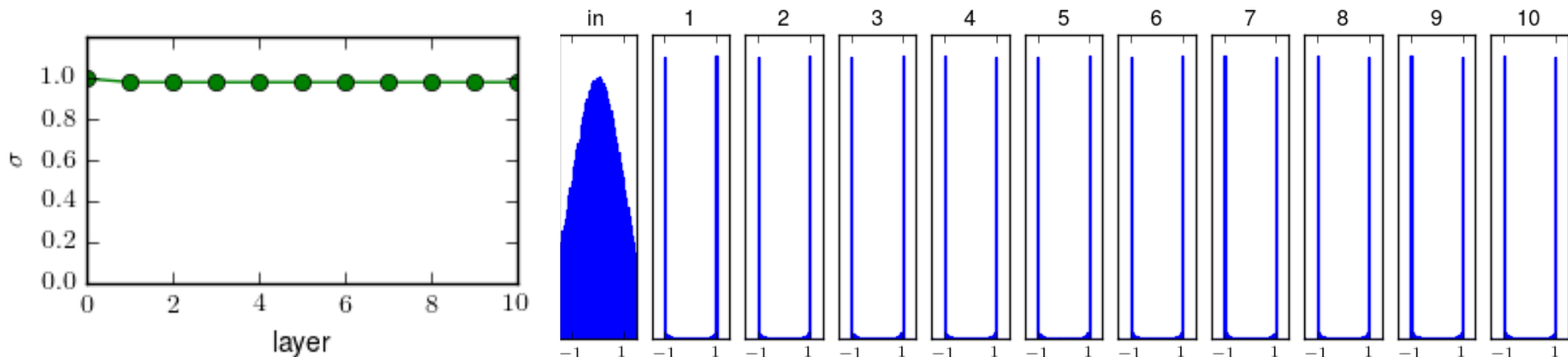
# Features in Deep Neural Networks



Layer 3

Zeiler and Fergus: *Visualizing and Understanding Convolutional Networks*, 2013

Layer 4

Layer 5

Zeiler and Fergus: *Visualizing and Understanding Convolutional Networks*, 2013

◆ Is it a good idea to set all weights to zero?

◆ **No.** All neurons would behave the same: the same $\delta$s are backpropagated. We need to *break the symmetry*

◆ Use small numbers, e.g., sample from a Gaussian distribution with zero mean:

  • works well for shallow networks,

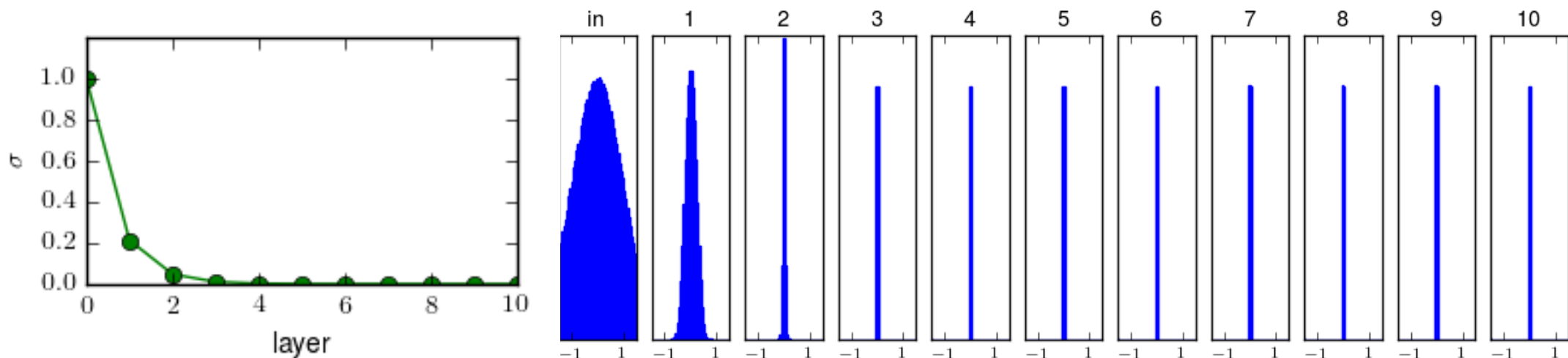  • for deep networks we might get into trouble

# Gaussian Initialization Example

◆ MLP, ten $\mathrm{tanh}$ layers, 500 units each. Each input fed with $\mathcal{N}(0,1)$

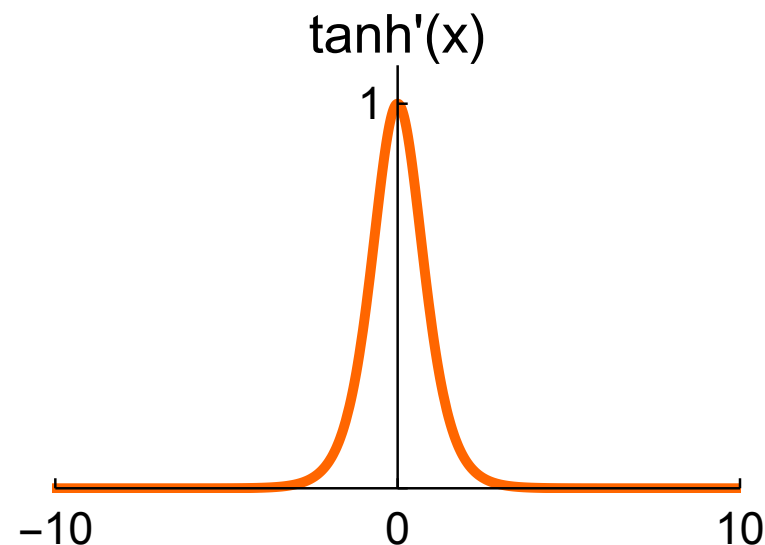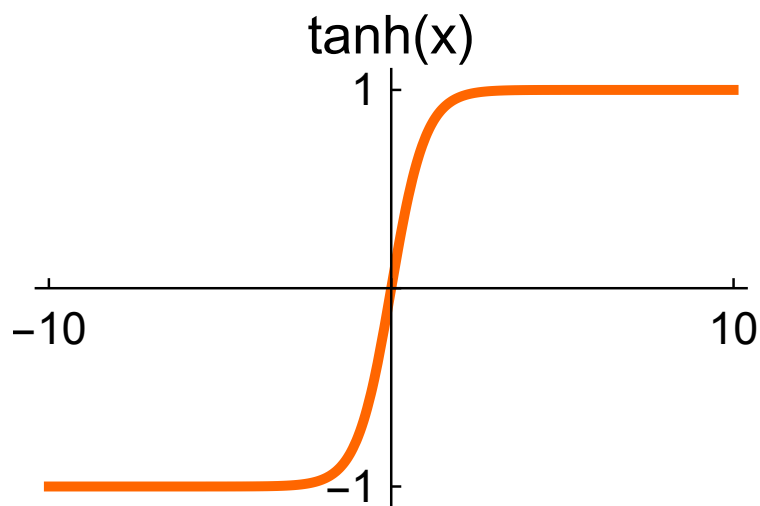◆ Weights initialized to $\mathcal{N}(0,\sigma^2)$

$$\sigma = 1$$



$$\sigma = 0.01$$

◆ For large $\sigma$ (saturation) the $\tanh$ derivative is almost zero

◆ For small $\sigma$ (output close to zero):

  ● the derivative is at most 1,

  ● the weights are very small and $\frac{\partial z_j^{l+1}}{\partial z_i^l} = w_{ij}$ holds for the preceding linear layer

◆ In both cases: $\delta \to 0$ as the number of layers increases

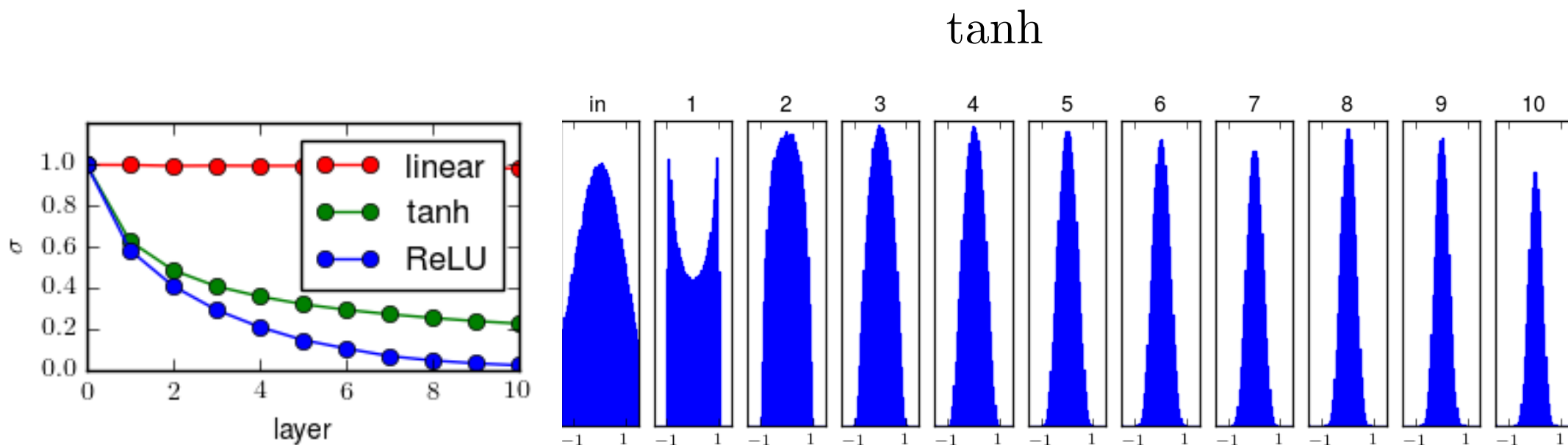◆ Glorot and Bengio: *Understanding the difficulty of training deep feedforward neural networks*, 2010

◆ For the linear neuron $s = \sum_i w_i x_i$, let $w_i$ and $x_i$ be independent random variables, $w_i$ and $x_i$ are i.i.d., $E(x_i) = E(w_i) = 0$:

$$\mathrm{Var}(s) = \mathrm{Var}\left(\sum_i w_i x_i\right) = \sum_i \mathrm{Var}(w_i x_i) =$$

$$= \sum_i \mathbb{E}\left([w_i x_i - \mathbb{E}(w_i x_i)]^2\right) = \sum_i \mathbb{E}\left([w_i x_i - \mathbb{E}(w_i)\mathbb{E}(x_i)]^2\right) =$$

$$= \sum_i \mathbb{E}(w_i^2 x_i^2) = \sum_i \mathbb{E}(w_i^2)E(x_i^2) =$$

$$= \sum_i \mathbb{E}([w_i - E(w_i)]^2)E([x_i - E(x_i)]^2) =$$

$$= \sum_i \mathrm{Var}(x_i)\mathrm{Var}(w_i) = n_{\mathsf{in}}\mathrm{Var}(x)\mathrm{Var}(w)$$

◆ We have $\mathrm{Var}(s) = n_{\mathsf{in}}\mathrm{Var}(x)\mathrm{Var}(w)$

◆ We want $\mathrm{Var}(s) = \mathrm{Var}(x)$

◆ Choose $\mathrm{Var}(w) = \frac{1}{n_{\mathsf{in}}}$

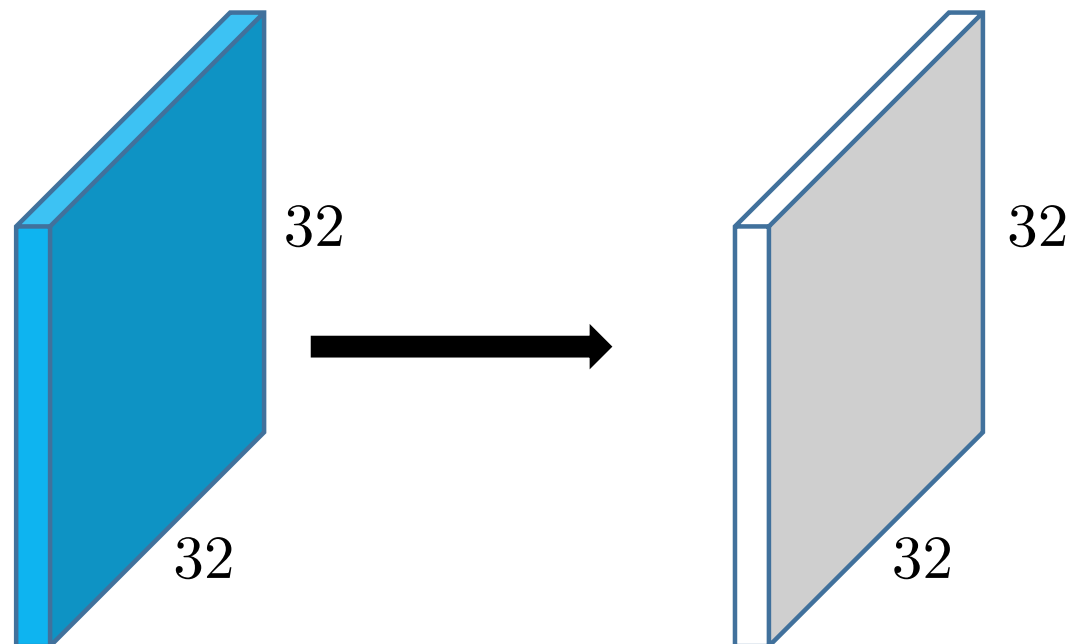◆ Works well for $\tanh$ as it is linear near zero
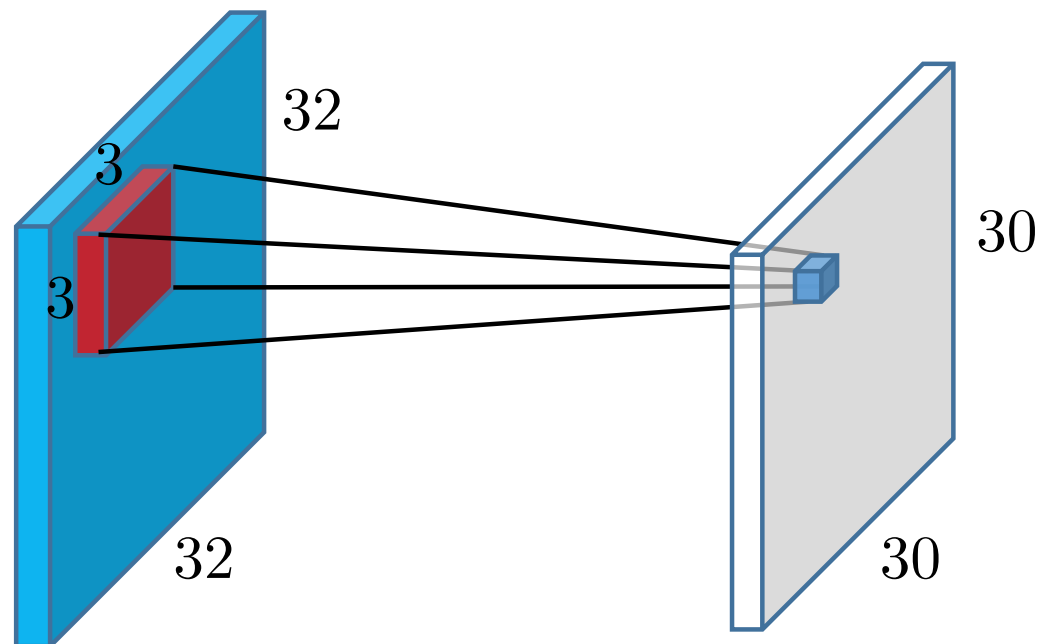
◆ Do not forget to standardize ANN input data

$\tanh$

◆ Topographical mapping in the visual cortex - nearby cells represent nearby regions in the visual field

◆ Input: grayscale image $32 \times 32$ pixels

◆ Output: layer of $32 \times 32$ features

◆ How many parameters do we need when input and output is fully connected?

◆ Topographical mapping in the visual cortex - nearby cells represent nearby regions in the visual field

◆ Input: grayscale image $32 \times 32$ pixels

◆ Output: layer of $32 \times 32$ features

◆ How many parameters do we need when input and output is fully connected?

$$\underbrace{32^2}_{\text{outputs}} \times (\underbrace{32^2}_{\text{inputs}} + \underbrace{1}_{\text{biases}}) \approx 1M$$

◆ Each neuron has a **receptive field** of $3 \times 3$ pixels

◆ It is fully connected only to the corresponding set of 9 inputs

◆ How many parameters do we need now?

# Locally Connected Layer

◆ Each neuron has a **receptive field** of $3 \times 3$ pixels

◆ It is fully connected only to the corresponding set of 9 inputs

◆ How many parameters do we need now?

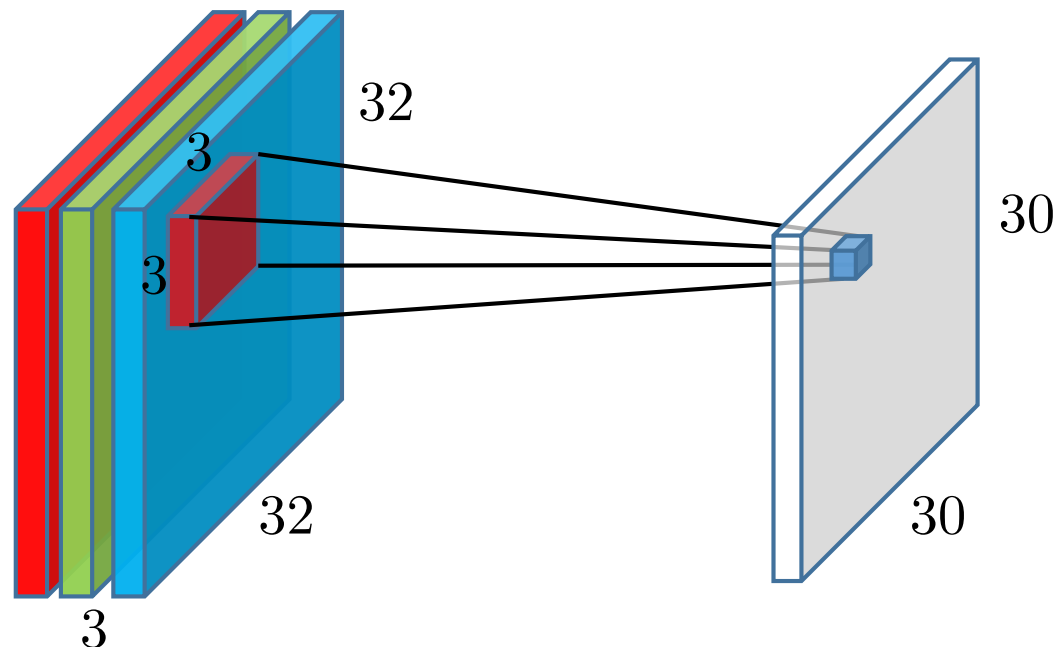$$\underset{\text{outputs}}{30^2} \times ( \underset{\text{inputs}}{3^2} + \underset{\text{bias}}{1} ) = 9\text{k}$$

◆ We can have more input channels, e.g., colors

◆ Now the input is defined by width, height and depth: $32 \times 32 \times 3$

◆ The number of parameters is $\underset{\text{outputs}}{30^2} \times (\underset{\text{channels}}{3} \times \underset{\text{inputs}}{3^2} + \underset{\text{bias}}{1}) \approx 25\text{k}$

- We can further reduce the number of parameters by sharing weights

- Use the same set of weights and bias for all outputs, define a *filter*

- The number of parameters drops to $3 \times 3^2 + 1 = 28$
  $\underbrace{\phantom{3 \times 3^2}}_{\text{inputs}} \quad \underbrace{\phantom{1}}_{\text{bias}}$

- Translation *equivariance*

# Multiple Output Channels

◆ Extract multiple different of features

◆ Use multiple *filters* to get more *feature maps*

◆ For $4$ filters we have $\underbrace{4}_{\text{filters}} \times (\underbrace{3 \times 3^2}_{\text{inputs}} + \underbrace{1}_{\text{bias}}) = 112$ parameters

◆ This is the **convolutional layer**

◆ Processes *volume* into *volume*

| | | | | | | |
|---|---|---|---|---|---|---|
| Identity | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | | Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | | |
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ | | Box blur (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | | |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ | | Gaussian blur (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | | | | | |

https://en.wikipedia.org/wiki/Kernel_(image_processing)

# Convolution in 2D: Forward Message



$$z_{kld} = f_{kld}(\boldsymbol{x}, \boldsymbol{w}, \boldsymbol{b}) = b_d + \sum_{i=1}^{F}\sum_{j=1}^{F}\sum_{c=1}^{C} x_{k+i-1,l+j-1,c}\, w_{ijcd}$$
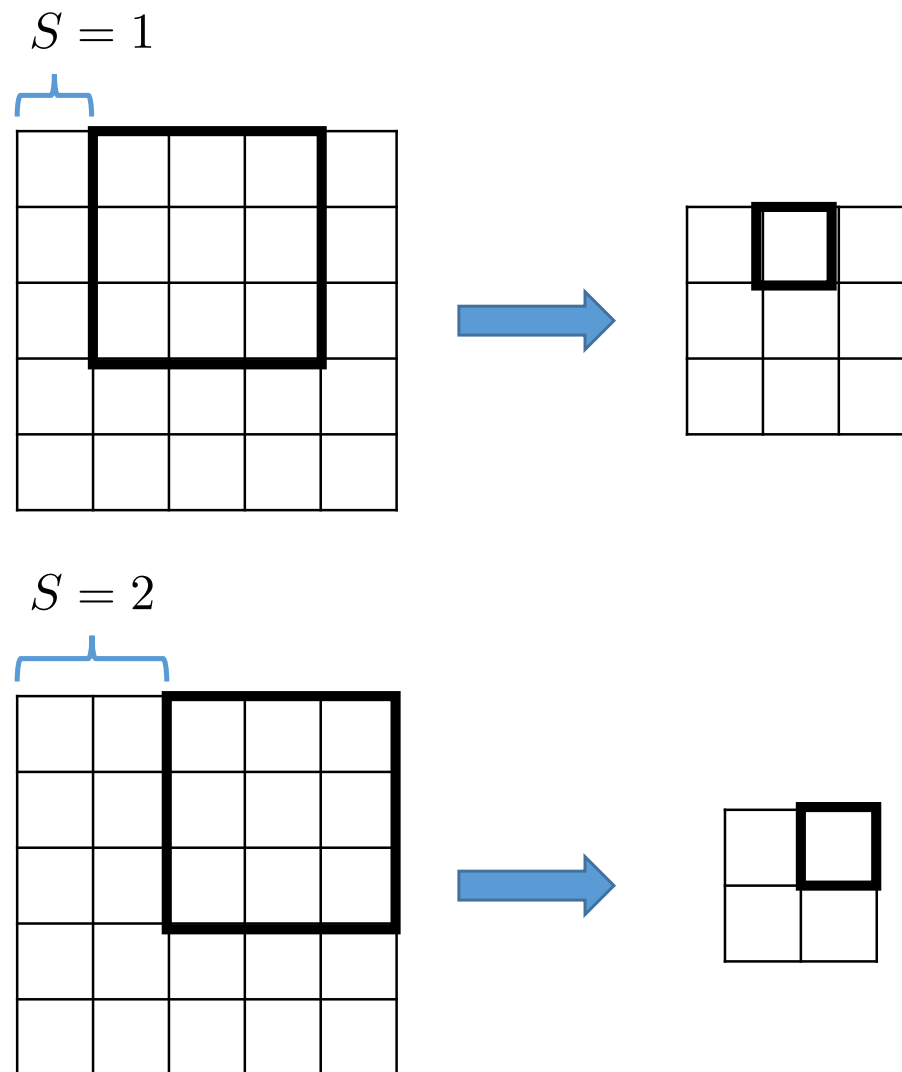
◆ Stride hyper parameter, typically $S \in \{1, 2\}$

◆ Higher stride produces smaller output volumes spatially

$S = 1$



$S = 2$

◆ Stride hyper parameter, typically $S \in \{1, 2\}$

◆ Higher stride produces smaller output volumes spatially



$S = 1$

$S = 2$

# Zero Padding

◆ Convolutional layer reduces the spatial size of the output w.r.t. the input

◆ For many layers this might be a problem

◆ This is often fixed by *zero padding* the input

◆ The size of the zero padding is denoted $P$

$$P = 1, \; S = 1$$

# Convolutional Layer Summary

◆ Input volume: $W_{\text{input}} \times H_{\text{input}} \times C$

◆ Output volume: $W_{\text{output}} \times H_{\text{output}} \times D$

◆ Having $D$ filters:

    • receptive field of $F \times F$ units,

    • stride $S$

    • zero padding $P$

$$W_{\text{output}} = (W_{\text{input}} - F + 2P)/S + 1$$
$$H_{\text{output}} = (H_{\text{input}} - F + 2P)/S + 1$$

◆ Needs $F^2 C D$ weights and $D$ biases

◆ The number of activations and $\delta$s to store: $W_{\text{output}} \times H_{\text{output}} \times D$

# Convolutional Layer: Nonlinearities

◆ In most cases a nonlinearity (sigmoid, tanh, ReLU) is applied to the outputs of the convolutional layer

◆ Example: ReLU units



Input feature map

Output feature map

Black = negative; white = positive values

Only non-negative values

Rob Fergus: MLSS 2015 Summer School

# Max Pooling

◆ Reduces spatial resolution → less parameters → helps with overfitting

◆ Introduces translation invariance and invariance to small rotations

◆ Depth is not affected

$F = 2, S = 2$

| 2 | 2 | 0 | 4 | 3 | 4 |
|---|---|---|---|---|---|
| 0 | 0 | 5 | 0 | 4 | 1 |
| 4 | 5 | 2 | 5 | 1 | 4 |
| 5 | 2 | 1 | 0 | 2 | 1 |
| 2 | 3 | 3 | 3 | 5 | 3 |
| 0 | 3 | 0 | 4 | 0 | 1 |

| 2 | 5 | 4 |
|---|---|---|
| 5 | 5 | 4 |
| 3 | 4 | 5 |

http://cs231n.github.io/convolutional-networks/

# VGGNet 2014

- Simonyan, Zisserman: *Very Deep Convolutional Networks for Large-Scale Image Recognition*, 2014

- Lowering filter spatial resolution ($F = 3$, $S = 1$, $P = 1$), increasing depth

- A sequence of $3 \times 3$ filters can emulate a single large one

- Top five error $7.3\%$, $6.8\%$ for an ensemble of 2 CNNs

The network layers (left to right): input, conv3-64, conv3-64, MP, conv3-128, conv3-128, MP, conv3-256, conv3-256, conv3-256, MP, conv3-512, conv3-512, conv3-512, MP, conv3-512, conv3-512, conv3-512, MP, FC - 4096, FC - 4096, FC − 1000, softmax

| | input | conv3-64 | conv3-64 | MP | conv3-128 | conv3-128 | MP | conv3-256 | conv3-256 | conv3-256 | MP | conv3-512 | conv3-512 | conv3-512 | MP | conv3-512 | conv3-512 | conv3-512 | MP | FC - 4096 | FC - 4096 | FC − 1000 | softmax |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **parameters** | | 1.7k | 37k | | 74k | 147k | | 295k | 590k | 590k | | 1.2M | 2.4M | 2.4M | | 2.4M | 2.4M | 2.4M | | **103M** | 16.7M | 4M | |
| **activations** | 150k | **3.2M** | **3.2M** | 800k | 1.6M | 1.6M | 400k | 800k | 800k | 800k | 200k | 400k | 400k | 400k | 100k | 100k | 100k | 100k | 25k | 4096 | 4096 | 1000 | 1000 |

Activation dimensions: 224 x 224 x 3 (input), 224 x 224 x 64, 112 x 112 x 64, 112 x 112 x 128, 56 x 56 x 128, 56 x 56 x 256, 28 x 28 x 256, 28 x 28 x 512, 14 x 14  512, 14 x 14 x 512, 7 x 7 x 512, 1 x x1 x 4096, 1 x 1 x 4096, 1 x 1 x 1000

◆ Convolutional layer can be simply transformed to a Fully-connected layer
$\rightarrow$ sparse weight matrix

◆ The other direction is also possible:
FC layer of $D$ units following a $F \times F \times C$ convolutional layer can be replaced by a $1 \times 1 \times D$ convolutional layer using $F \times F$ filters ($P = 0$, $S = 1$)

◆ In both cases you do not have to recompute the weights, you just rearrange them

input

CONV, MP layers

$7 \times 7 \times 512$

$224 \times 224 \times 3$

$1 \times 1 \times 4096$
$F = 7$

$1 \times 1 \times 4096$
$F = 1$

$1 \times 1 \times 1000$
$F = 1$

softmax

input

$384 \times 384 \times 3$

CONV, MP layers

$12 \times 12 \times 512$

$6 \times 6 \times 4096$
$F = 7$

$6 \times 6 \times 4096$
$F = 1$

$6 \times 6 \times 1000$
$F = 1$

softmax

◆ Idea: use an existing model as a base to solve a *similar problem*

◆ Often used when not enough data available to solve the target problem directly

◆ Example: reuse an ImageNet network for object localization

# Transfer Learning

◆ Idea: use an existing model as a base to solve a *similar problem*

◆ Often used when not enough data available to solve the target problem directly

◆ Example: reuse an ImageNet network for object localization

- Idea: use an existing model as a base to solve a *similar problem*

- Often used when not enough data available to solve the target problem directly

- Example: reuse an ImageNet network for object localization

# Transfer Learning

◆ Idea: use an existing model as a base to solve a *similar problem*

◆ Often used when not enough data available to solve the target problem directly

◆ Example: reuse an ImageNet network for object localization

◆ You can:

- cut the original network at various layers,

- fix or not the weights of the original network or use different *learning rates*

- use different type of model for the head, e.g., linear SVM

Layer 1

Layer 2

Layer 3

Layer 4

Layer 5

32

32

32

32

32

32

32

32

32

3

3

32

30

30

32

3

3

3

32

30

30

32

3

3

3

32

30

30

3

32

3

3

32

3

30

30

4

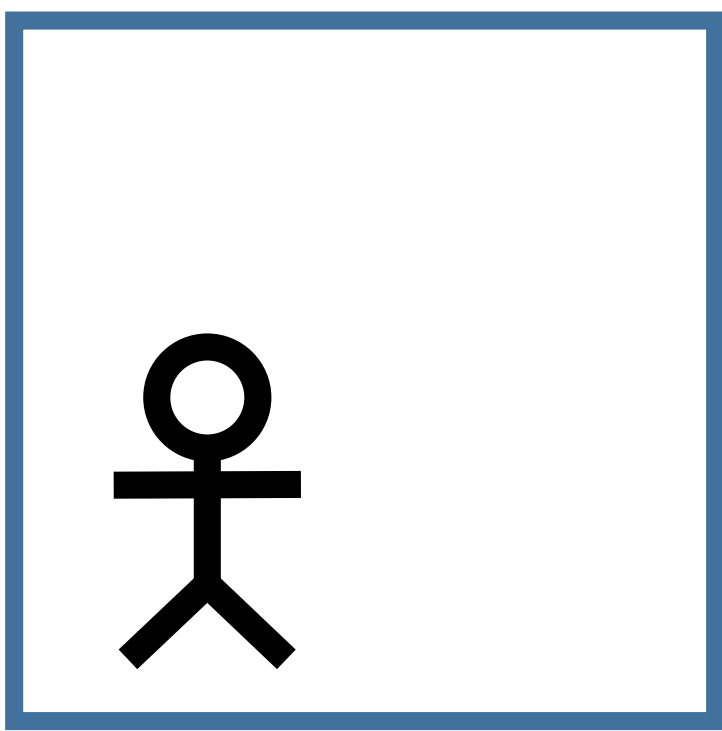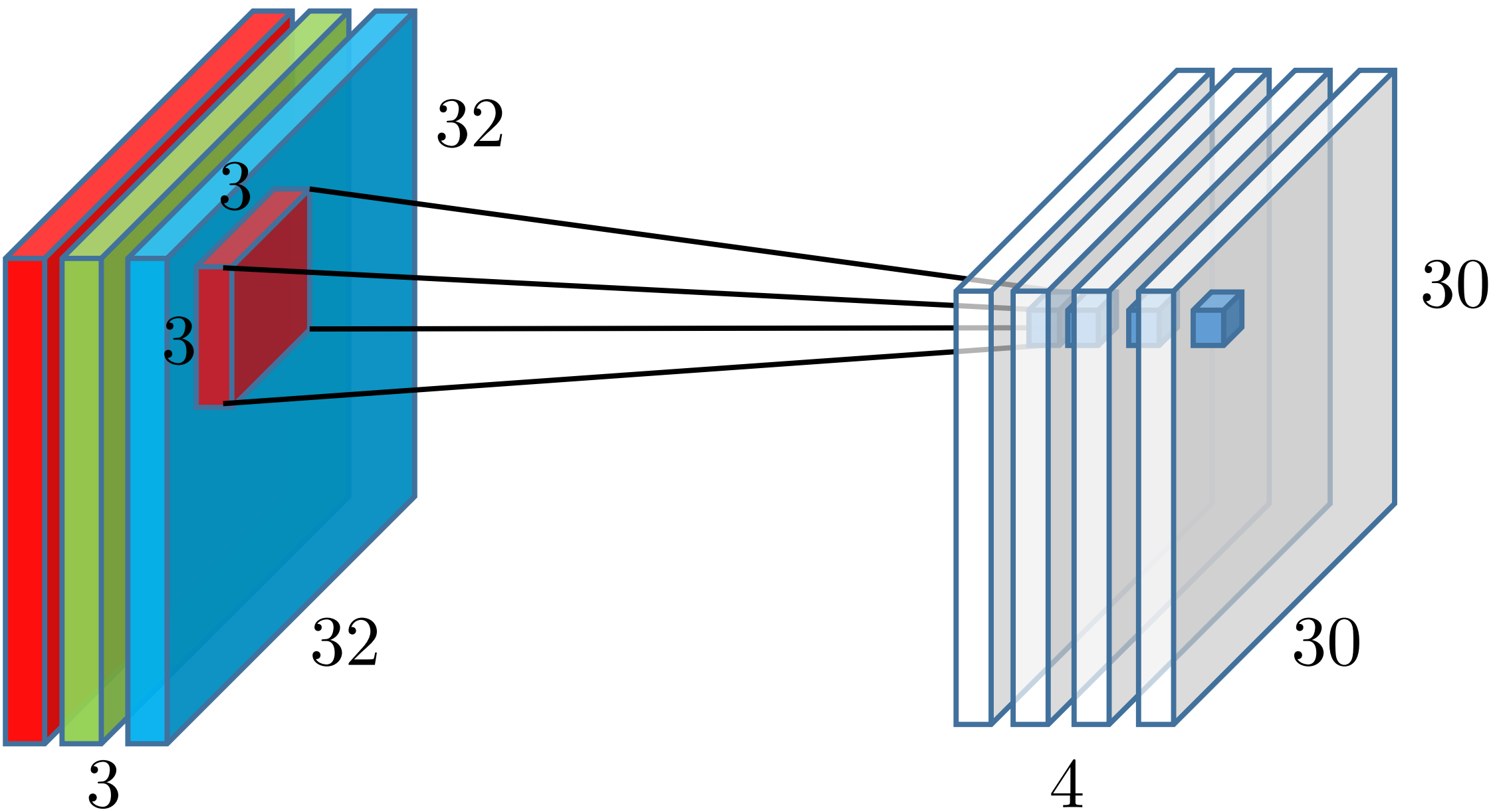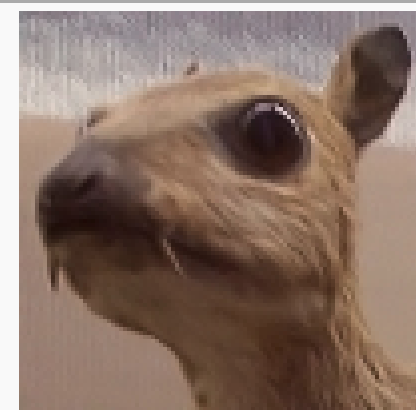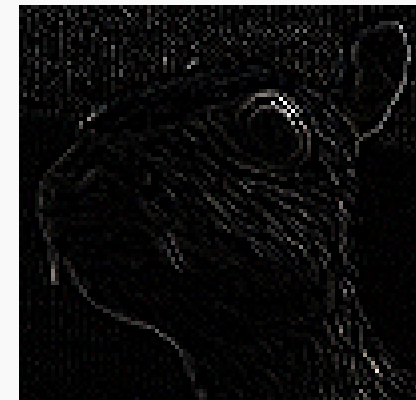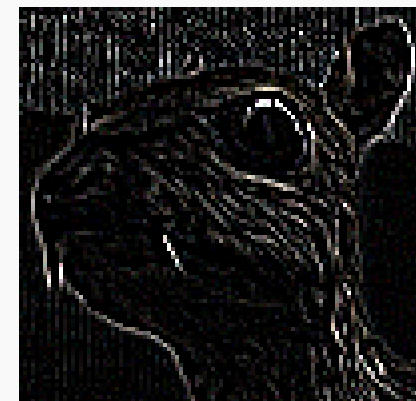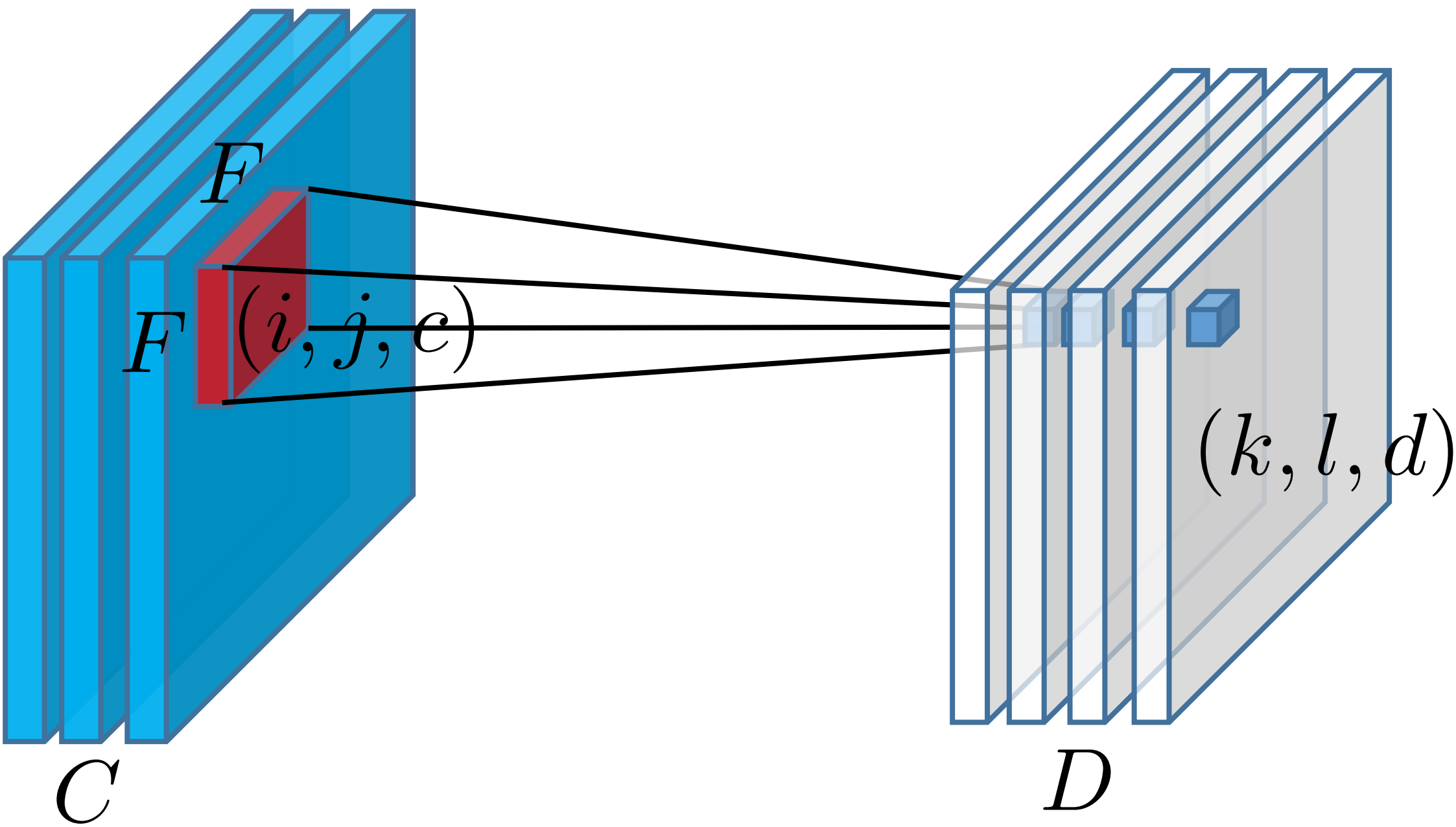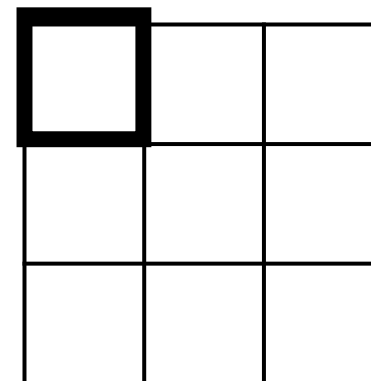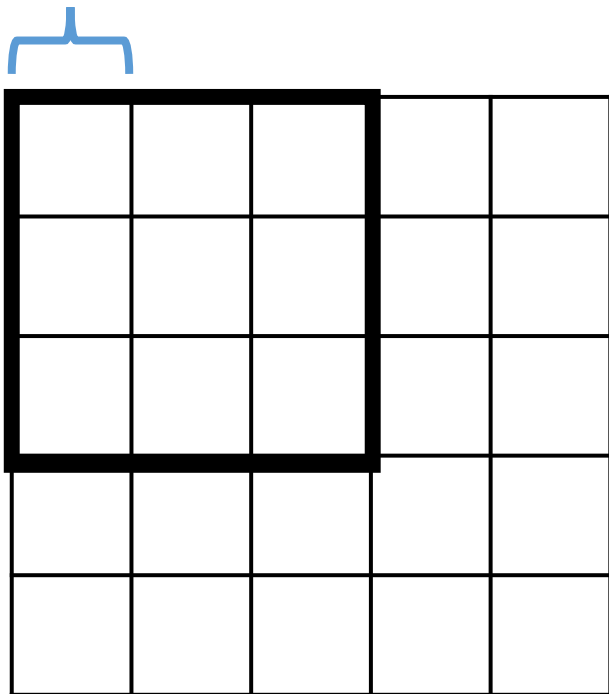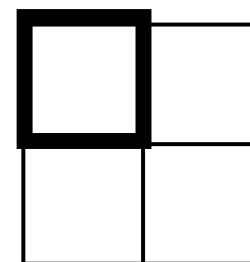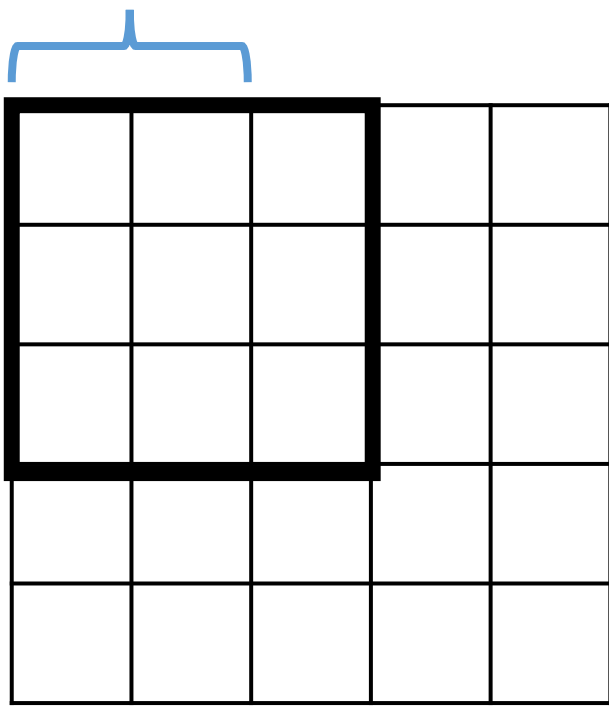| Identity | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ |  |
|---|---|---|
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ |  |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ |  |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ |  |

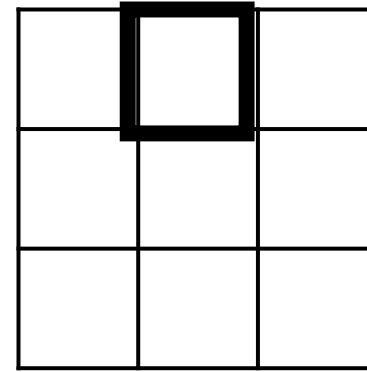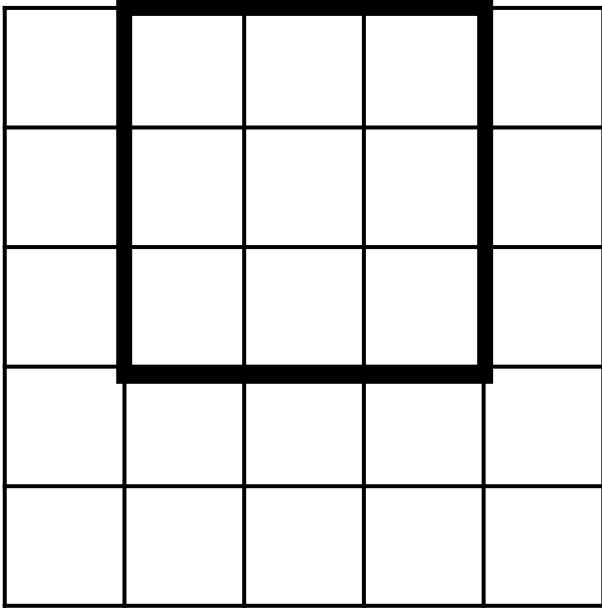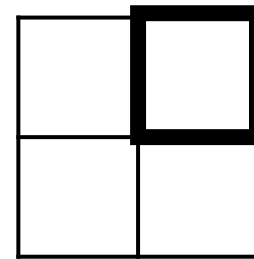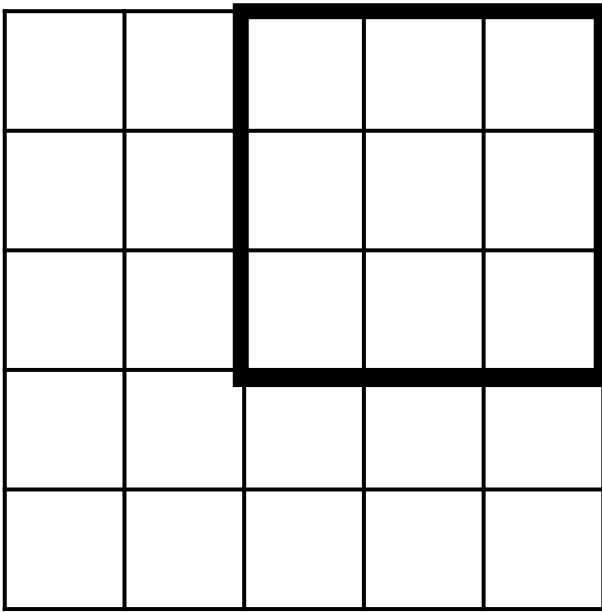| | | |
|---|---|---|
| **Sharpen** | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ |  |
| **Box blur**<br>(normalized) | $\dfrac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ |  |
| **Gaussian blur**<br>(approximation) | $\dfrac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ |  |

$S = 1$

$S = 2$

$S = 1$

$S = 2$

$P = 1, \ S = 1$

Input feature map

Output feature map

Black = negative; white = positive values

Only non-negative values

$F = 2, S = 2$

RELU RELU RELU RELU RELU RELU

POOL POOL POOL

CONV CONV CONV CONV CONV CONV

FC

car
truck
airplane
ship
horse

| | input | conv3-64 | conv3-64 | MP | conv3-128 | conv3-128 | MP | conv3-256 | conv3-256 | conv3-256 | MP | conv3-512 | conv3-512 | conv3-512 | MP | conv3-512 | conv3-512 | conv3-512 | MP | FC - 4096 | FC - 4096 | FC – 1000 | softmax |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **parameters** | | 1.7k | 37k | | 74k | 147k | | 295k | 590k | 590k | | 1.2M | 2.4M | 2.4M | | 2.4M | 2.4M | 2.4M | | **103M** | 16.7M | 4M | |
| **activations** | 150k | **3.2M** | **3.2M** | 800k | 1.6M | 1.6M | 400k | 800k | 800k | 800k | 200k | 400k | 400k | 400k | 100k | 100k | 100k | 100k | 25k | 4096 | 4096 | 1000 | 1000 |
| | $224 \times 224 \times 3$ | $224 \times 224 \times 64$ | | $112 \times 112 \times 64$ | $112 \times 112 \times 128$ | | $56 \times 56 \times 128$ | $56 \times 56 \times 256$ | | | $28 \times 28 \times 256$ | $28 \times 28 \times 512$ | | | $14 \times 14\ 512$ | $14 \times 14 \times 512$ | | | $7 \times 7 \times 512$ | $1 \times 1 \times 4096$ | $1 \times 1 \times 4096$ | $1 \times 1 \times 1000$ | |

input

CONV, MP layers

$224 \times 224 \times 3$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$
$F = 7$

$1 \times 1 \times 4096$
$F = 1$

softmax

$1 \times 1 \times 1000$
$F = 1$

input

CONV, MP
layers

$12 \times 12 \times 512$

$6 \times 6 \times 4096$
$F = 7$

$6 \times 6 \times 4096$
$F = 1$

$6 \times 6 \times 1000$
$F = 1$

softmax

$384 \times 384 \times 3$

input

CONV & MP

class

input

CONV & MP

$(x, y, w, h)$

input

CONV & MP

class

$(x, y, w, h)$