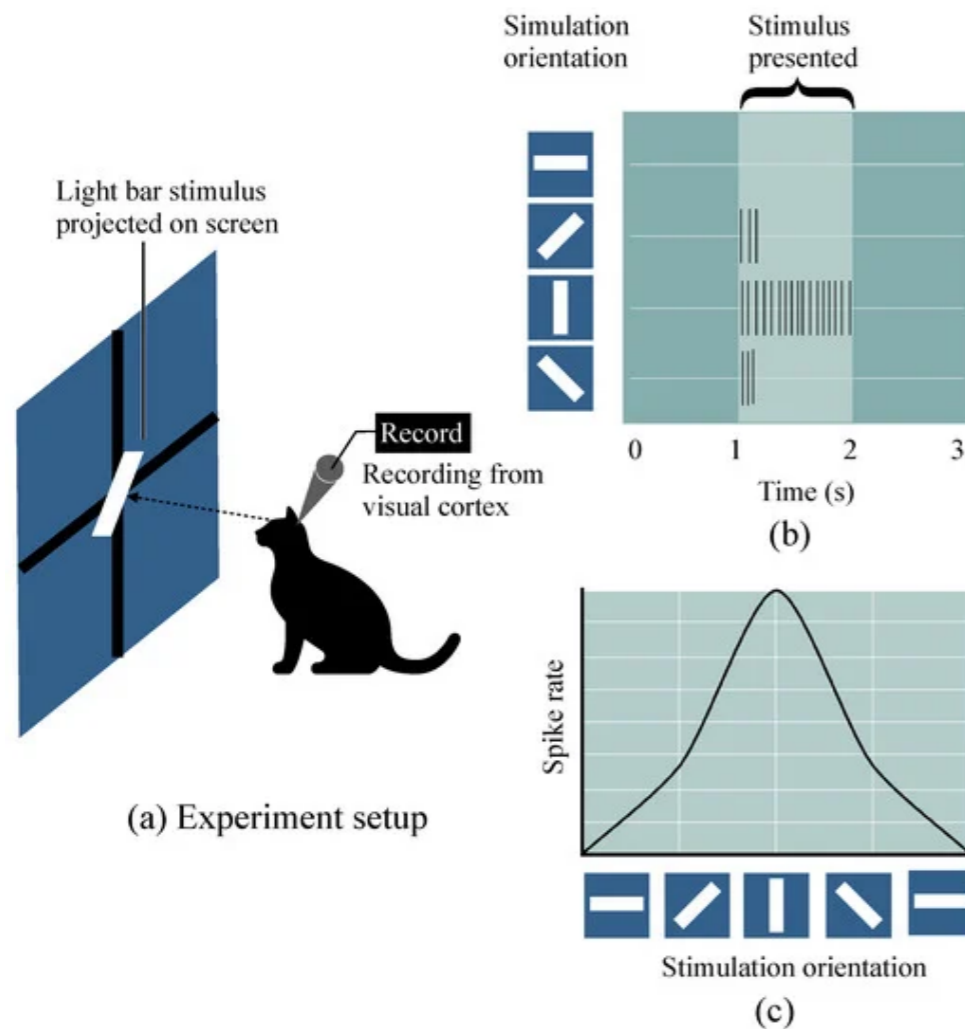# Deep Learning (BEV033DLE)
# Lecture 5
# Convolutional Neural Networks
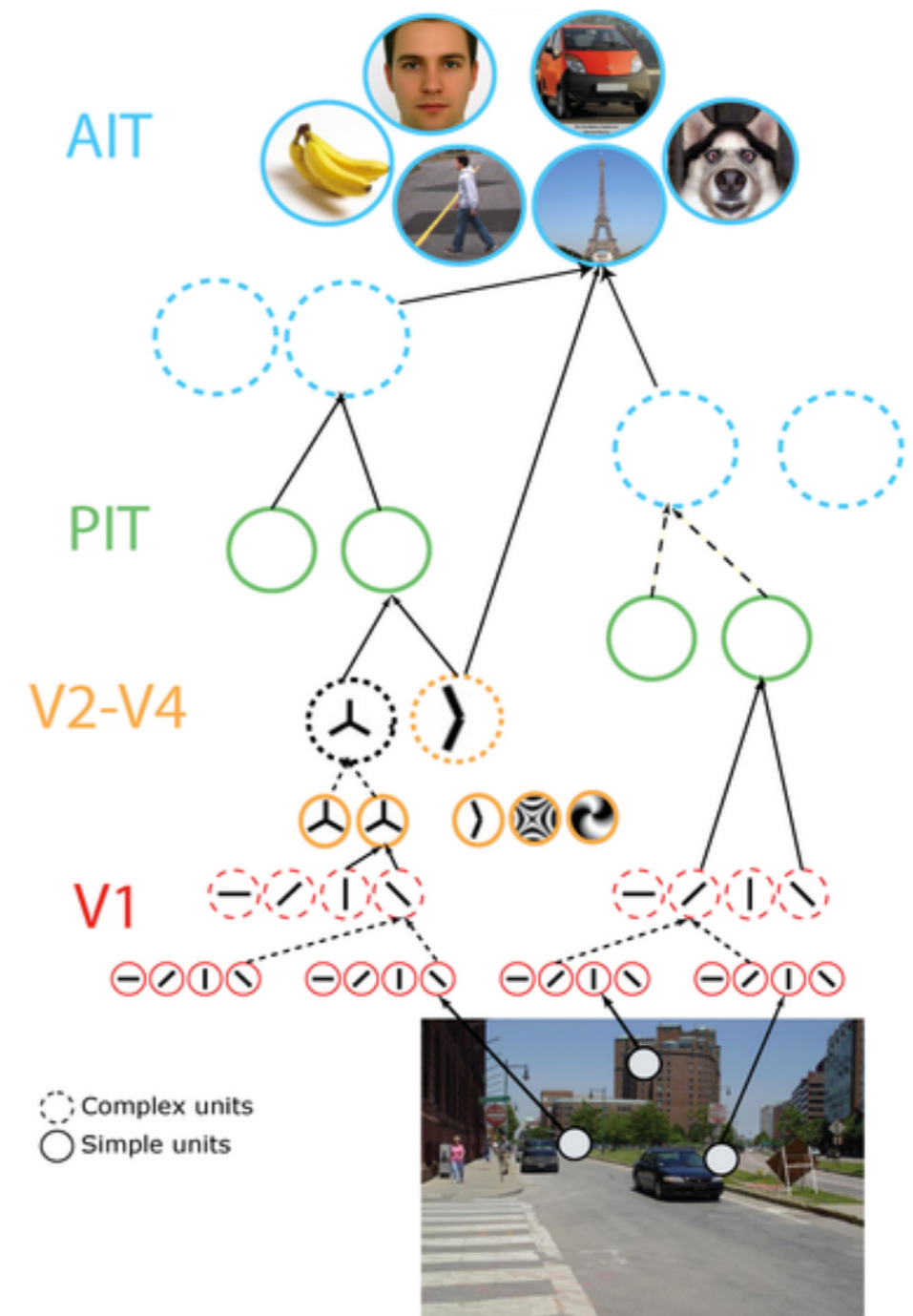
Czech Technical University in Prague

- Overview and Rationales of CNN architecture design

# Inspiration from Neuroscience

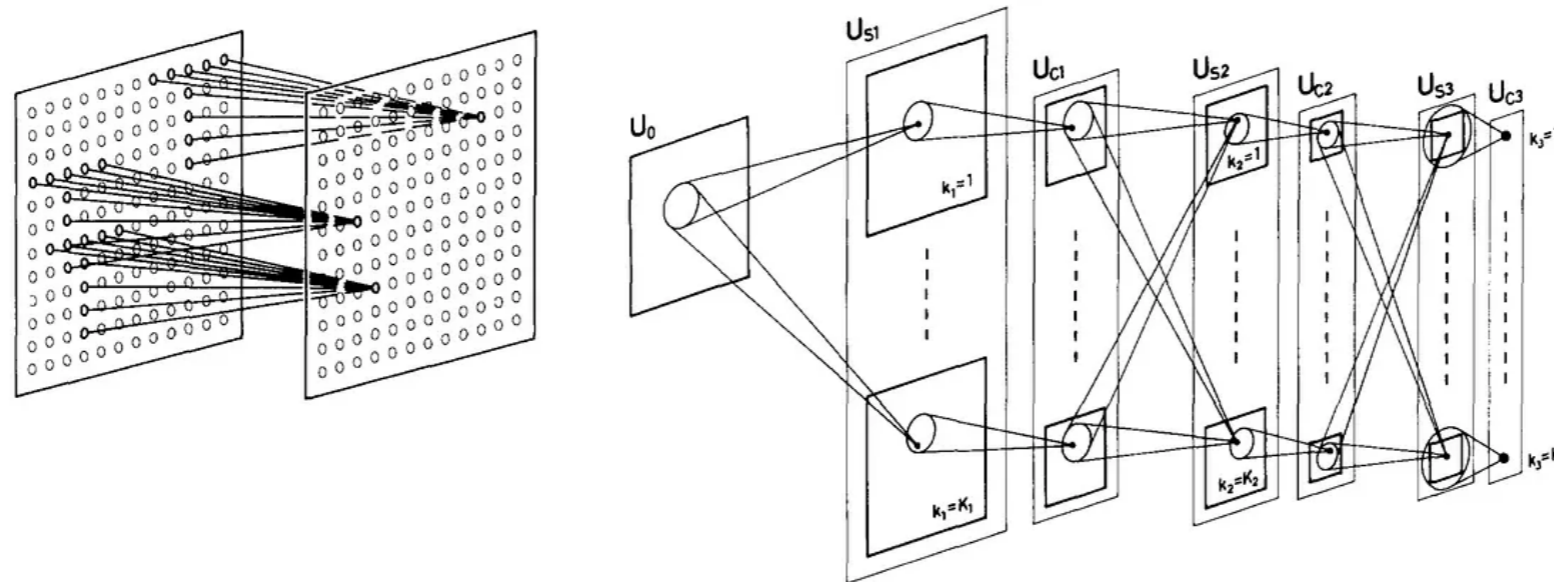✦ Hube and Wiesel (1959): Receptive fields of single neurones in the cat's striate cortex

✦ the organisation appeared to be hierarchical: responses of '*simple cells*' were aggregated by '*complex cells*,'

Inspiration e.g. for SIFT descriptors

# Neocognitron

K. Fukushima, Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position (1980)



K. Fukushima

Inspired by neuroscience (Hubel and Wiesel's observation of response to local patterns and idea of hierarchical organization, excitation-inhibition mechanism)

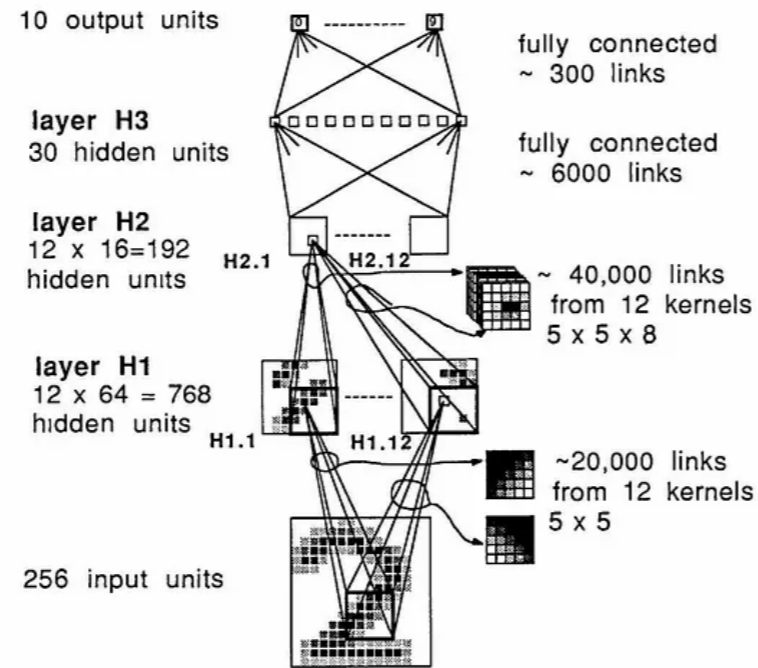A 7-layer network! Local receptive files with shared weights, pooling layers, ReLU activations.

Trained in an unsupervised manner...

By using local receptive fields at each level we achieve more flexibility to geometric variations
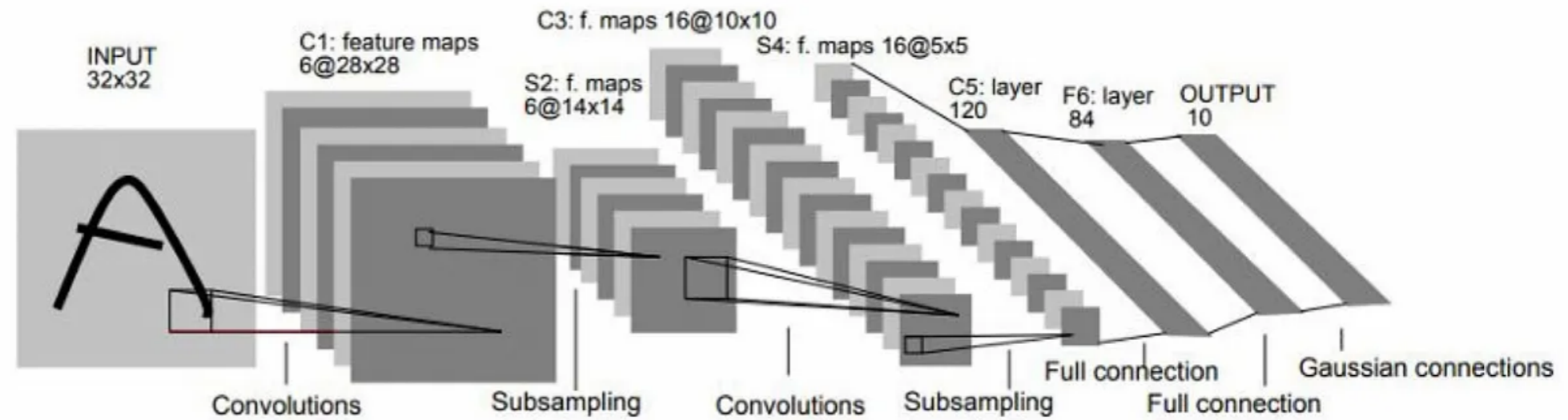
# Convolutional Neural Networks
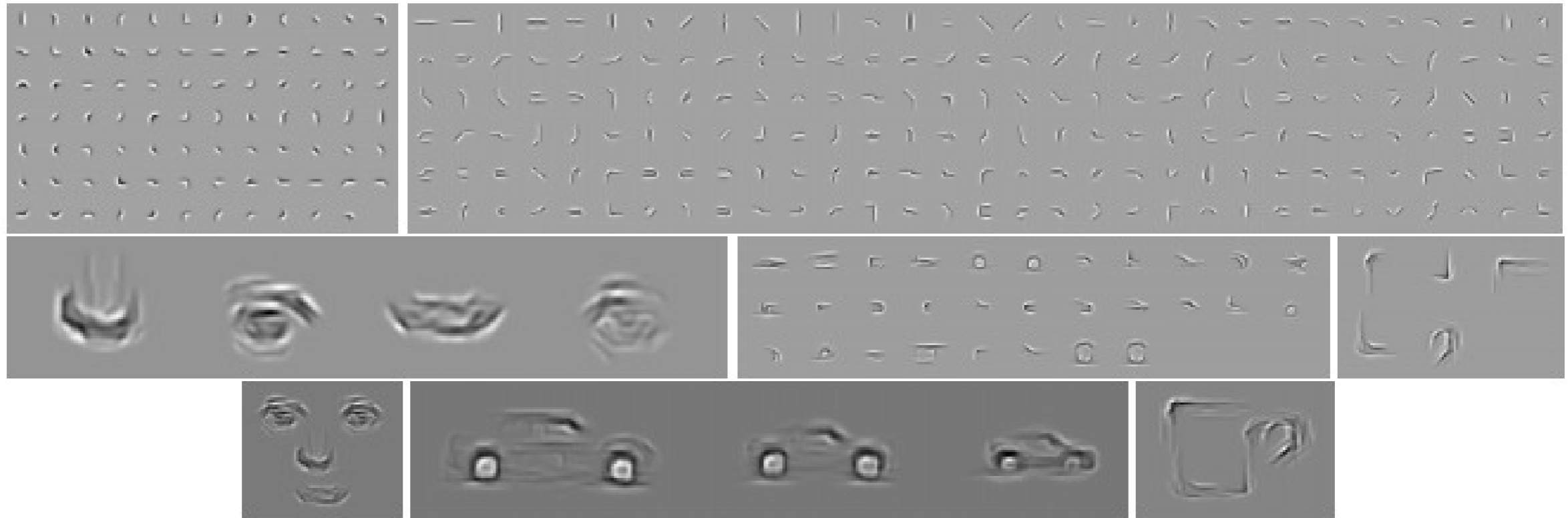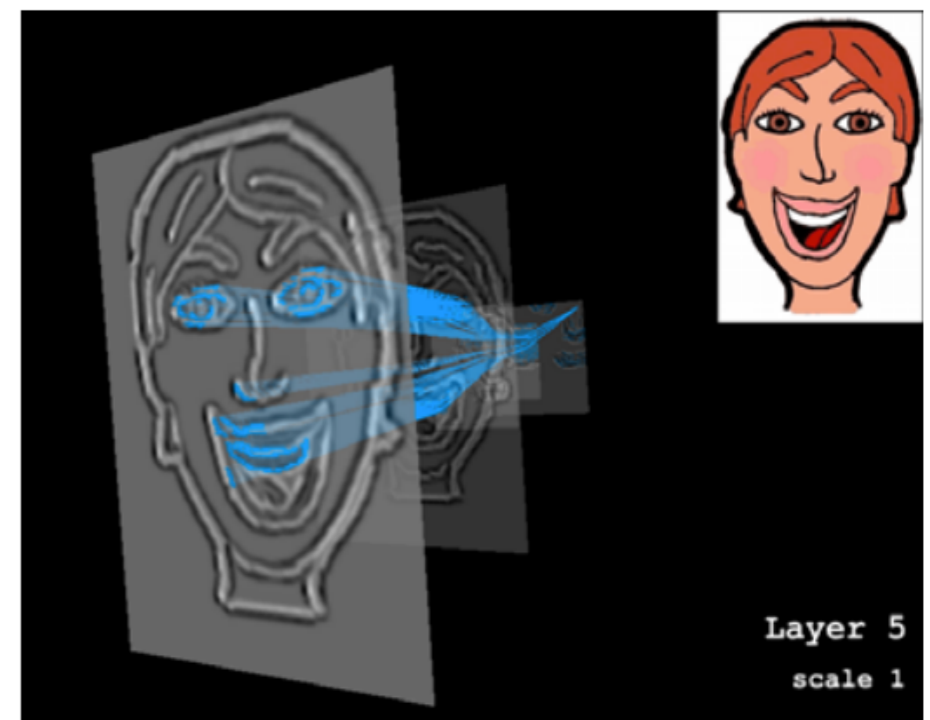
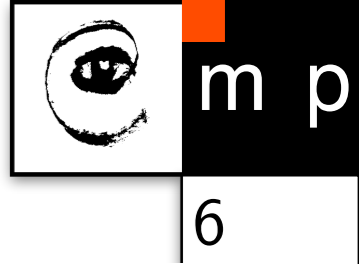LeNet 3 (1989)

Linear filters, backprop

LeNet 5 (1998)

✦ [Fidler and Leonardis (2007): "Towards Scalable Representations of Object Categories: Learning a Hierarchy of Parts"]



Learning layer-by-layer based on statistics and selection
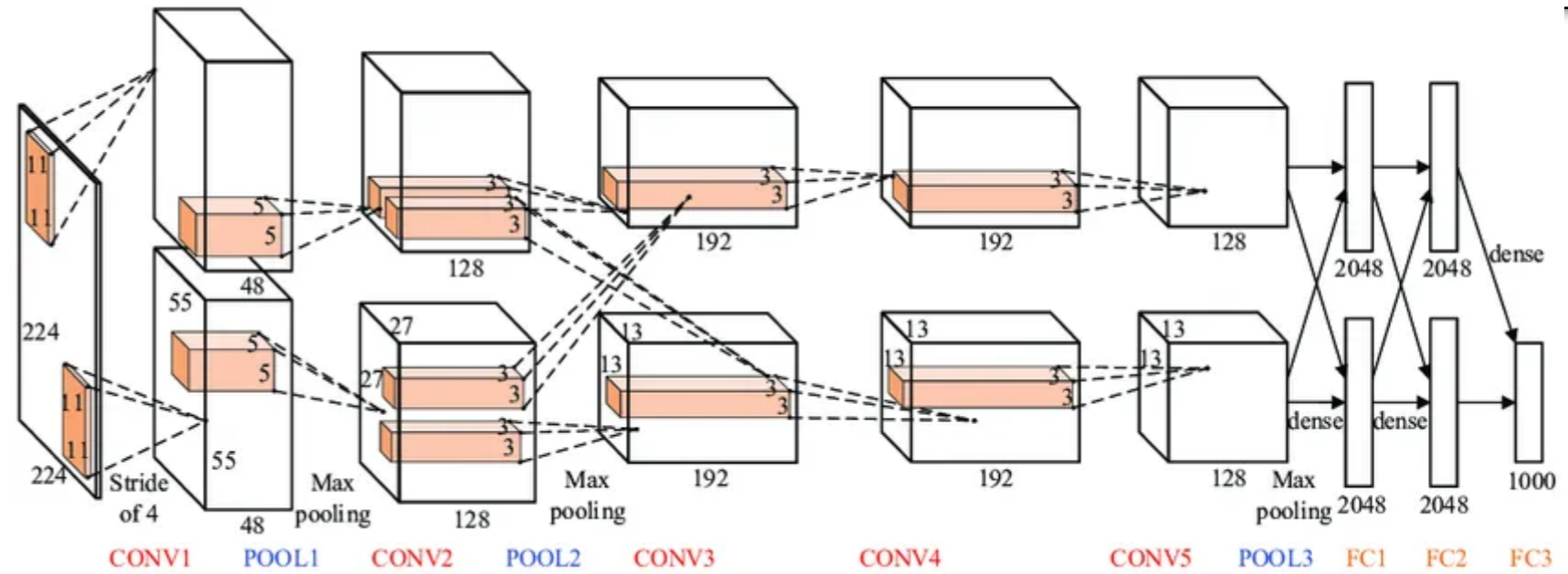


Layer 5
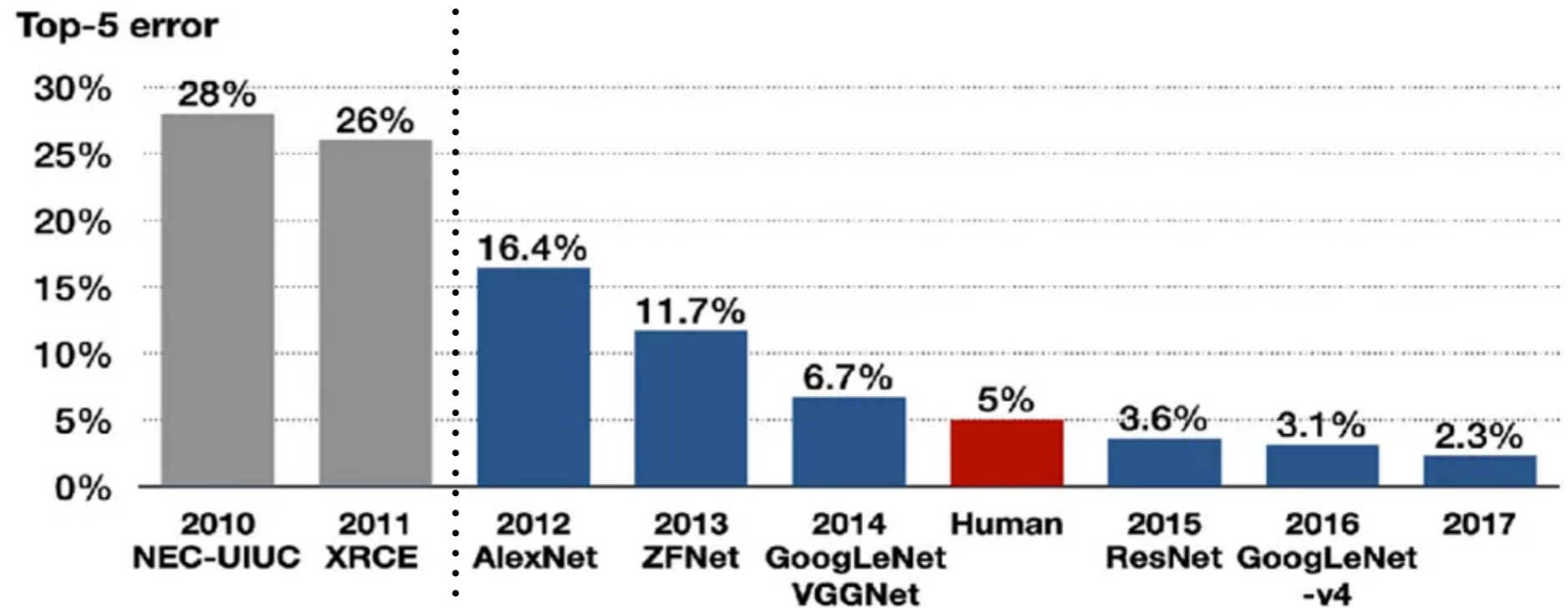scale 1

# Convolutional Neural Networks

AlexNet (2012)

More data & weights, GPU
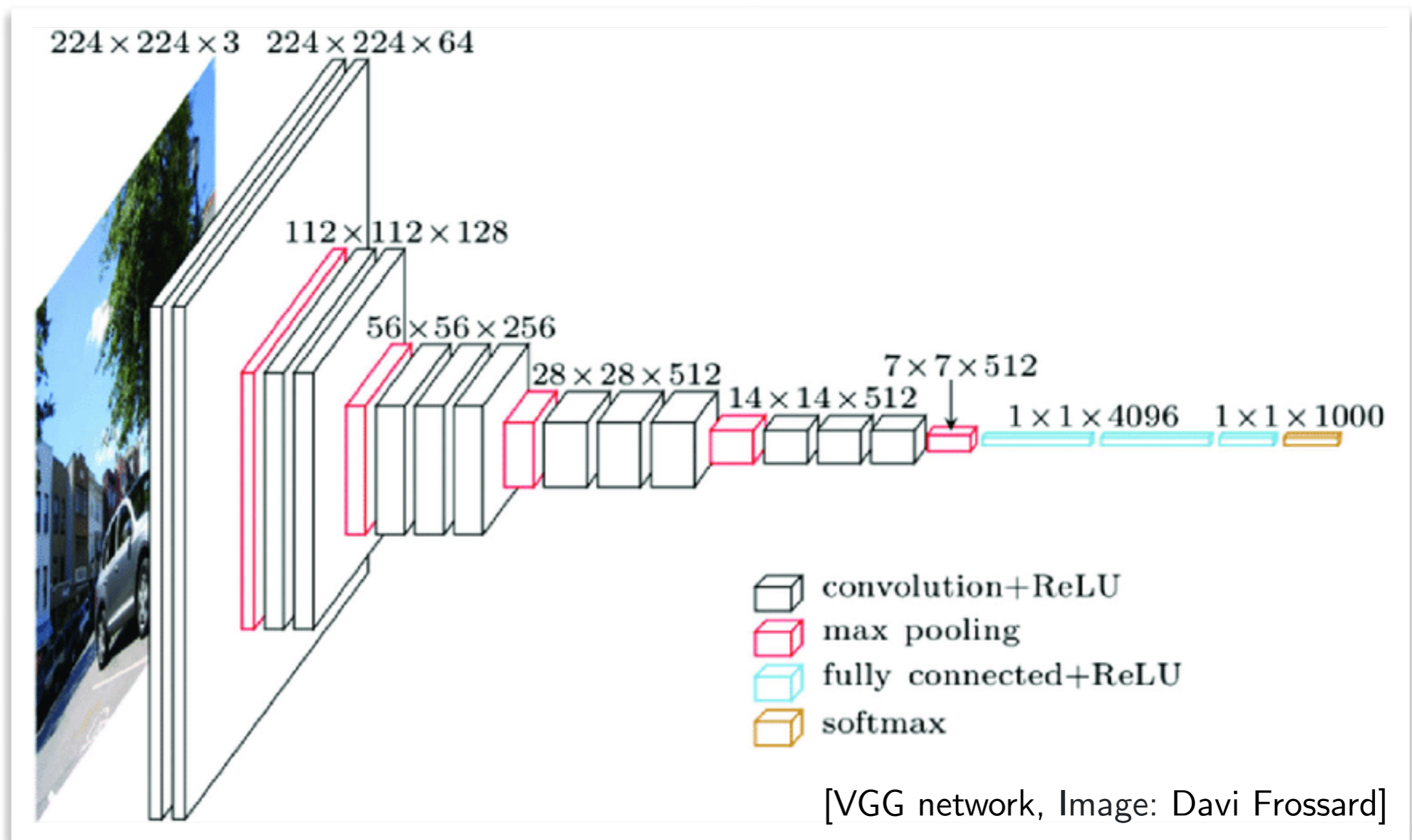


ImageNet classification challange



Engineered features
(incl. hierarchical) + SVM

Deep Learning

# VGG

[Simonyan and Zisserman (2014): Very Deep Convolutional Networks for
Large-Scale Image Recognition]



[VGG network, Image: Davi Frossard]

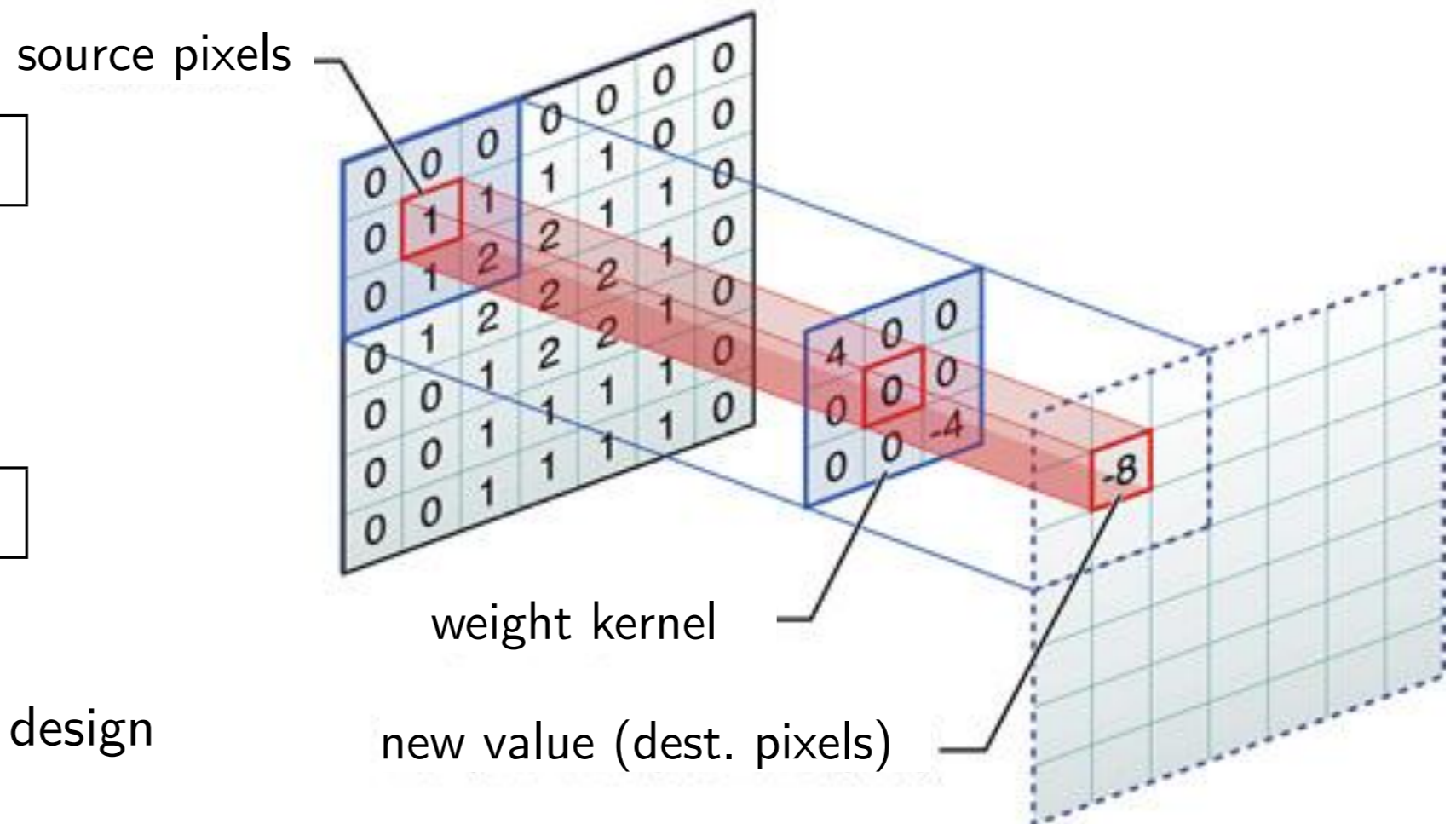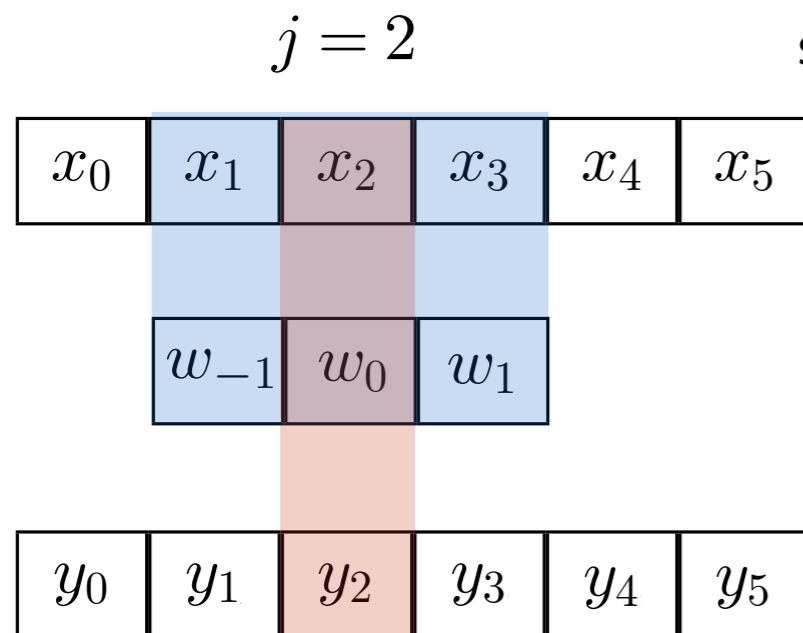✦ Goal: understand building blocks and design principles

◆ Convolution and Correlation (1D)

- Convolution
$$y = w * x: \quad y_j = \sum_{k=-h}^{h} w_k x_{j-k} = \sum_{k=-h}^{h} w_{-k} x_{j+k}$$

- Cross-correlation $y = w \star x: \quad y_j = \sum_{k=-h}^{h} w_k x_{j+k}$

flip of the weight matrix

output

weight kernel

input

Easily convertible, more convenient to consider cross-correlation in Deep Learning

$j = 2$

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

| $w_{-1}$ | $w_0$ | $w_1$ |

| $y_0$ | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ |

source pixels

weight kernel

new value (dest. pixels)
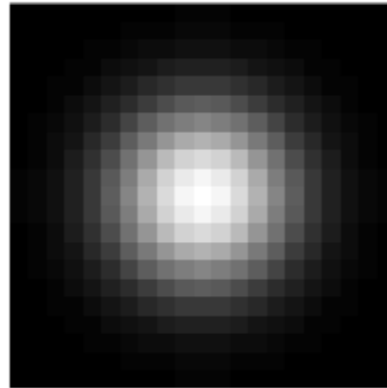
◆ Translation equivariance by design

# Examples (Cross-Correlation)

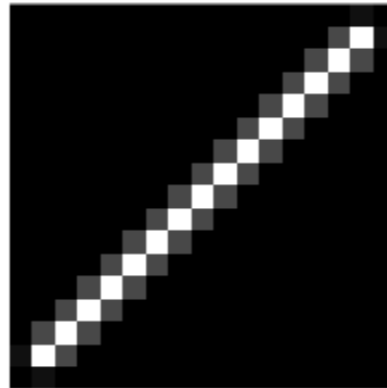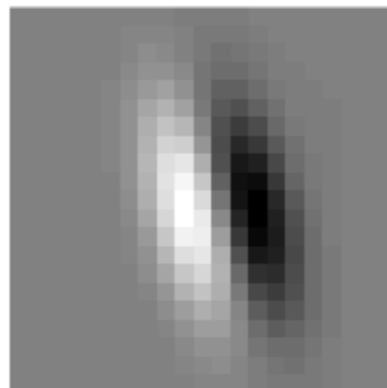| Input | Kernel | Output | |
|---|---|---|---|



Blur

Motion Blur

Input          Kernel          Output



Blur

Motion Blur

Shift

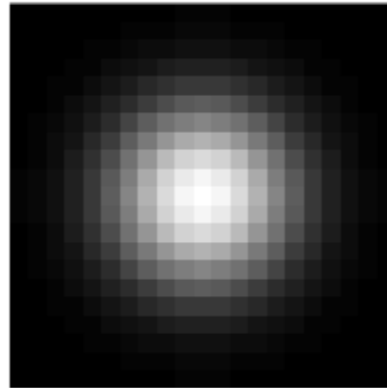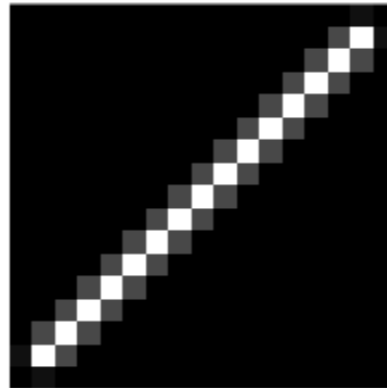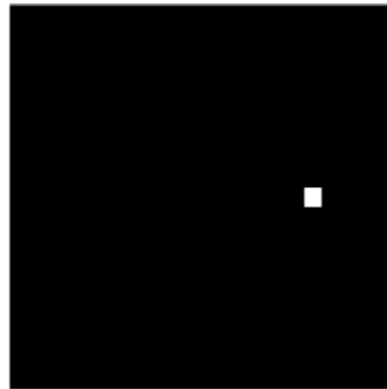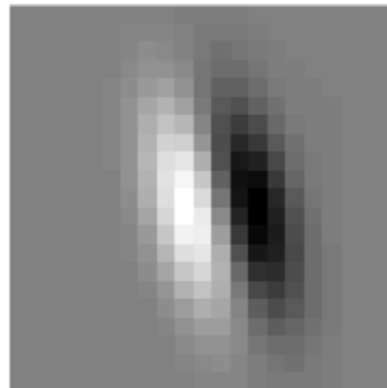# Examples (Cross-Correlation)

| Input | Kernel | Output | |
|-------|--------|--------|--|
| | | | Blur |
| | | | Motion Blur |
| | | | Shift |
| | | | Edge detector |

# Multi-Channel Convolution

✦ Extension:

- color input images -> convolution kernel needs to have 3 channels

- stack of filters -> channels of the output

feature map



✦ Multi-channel cross-correlation:

$$\sum_{c}\sum_{\Delta i}\sum_{\Delta j} x_{c,i+\Delta_i,j+\Delta j}\ w_{o,c,\Delta i,\Delta j} = y_{o,i,j}$$

input channel      filter spatial dimensions      output channel

- input is 3D tensor, weight is 4D tensor, output is 3D tensor

- Essentially: a cross-correlation on spatial dims and fully-connected on channel dims

# Invariance and Equivariance

shift

Classification

→ "Baloon"

Invariant to shift

Segmentation

Segmentation

balloon       sky

shift

sky       balloon

Equivariant to shift (commutes with shift)

Would be hard to achieve if the image was given as a general vector — we are using 2D grid structure and require that all locations are treated equally

[M. Bronstein et al.: "Geometric Deep Learning"]

✦ Concept: systematization of geometries as study of equivariances

• Apply this principle to systematize the zoo of NN architectures



Grids          Groups          Graphs          Geodesics & Gauges

✦ Key message: neural networks for processing geometric data should respect the structure of the domain

◆ Convolution:

- $y_j = \sum_{k=-h}^{h} w_k x_{j-k}$

◆ As matrix-vector product:

- Denote: $i = j - k$, then $k = j - i$ and $y_j = \sum_i w_{j-i} x_i$
- Denote $W_{i,j} = w_{i-j}$
- Then $y = Wx$

✦ Convolution is a linear transform of a special structure:

Example: $w_k = k$

| -1 | 0 | 1 |
|----|---|---|

$W$

| 1 | 0 | -1 |
|---|---|----|

$x$

◆ Cross-correlation has flipped $w$, resulting in transposed $W$

- Backprop of convolution is cross-correlation and vice-versa

Grid



Shift operator (1D)



$\mathbf{y}$     $\mathbf{S}$     $\mathbf{x}$         $\mathbf{y}$     $\mathbf{S}^{\mathsf{T}}$     $\mathbf{x}$

◆ Characterize linear transforms which are equivariant to shift:

- let $S \colon \mathbb{R}^n \to \mathbb{R}^n$ be a shift matrix: $S_{i,j} = [\![ j = i+1 ]\!]$
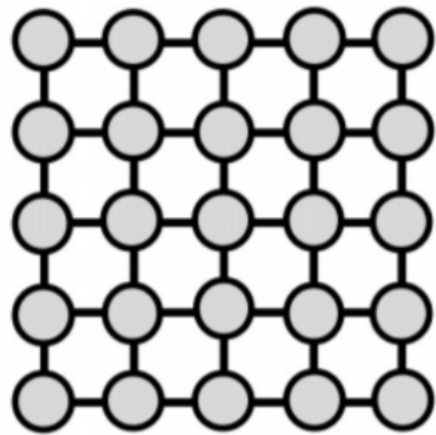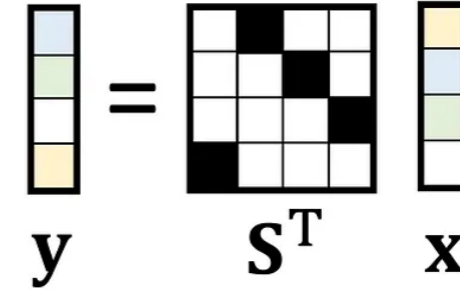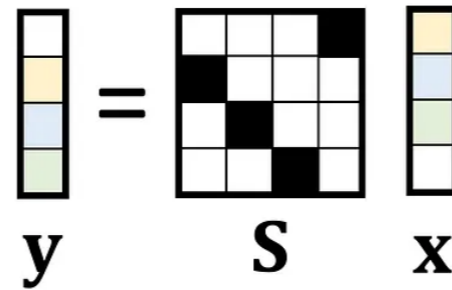
- let $A \colon \mathbb{R}^n \to \mathbb{R}^n$

- equivariance: $ASx = SAx$ for all $x$

- implies $AS = SA$

- implies $A_{i,j+1} = A_{i+1,j}$ for all $i, j$ – "circulant" matrix (convolution)

✦ **Conclusion**: a matrix is circulant if and only if it commutes with shift

◆ Further properties:

- Matrices satisfying $AS = SA$ will have same eigenvectors

- Eigenvectors of shift $S$ are the Fourier basis functions $\Phi$

- All convolutions can be represented as $A = \Phi \Lambda(w) \Phi^{\mathsf{T}}$, where $\Lambda(w) = \mathrm{diag}(\Phi^{\mathsf{T}} w)$

- Convolution Theorem: $Ax = \Phi \Big( (\Phi^{\mathsf{T}} w) \odot (\Phi^{\mathsf{T}} x) \Big)$

# Learned vs Engineered Filters

✦ Gabor Filters (and generalizations) - mathematical model for V1 cells:



Equivariance to transformations + more design principles w.r.t. scale-space

Active research on symmetry as a guiding principle in artificial and brain neural networks
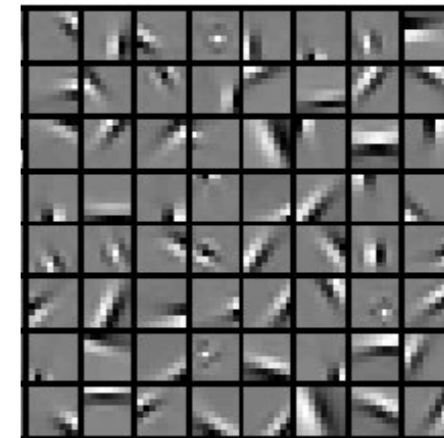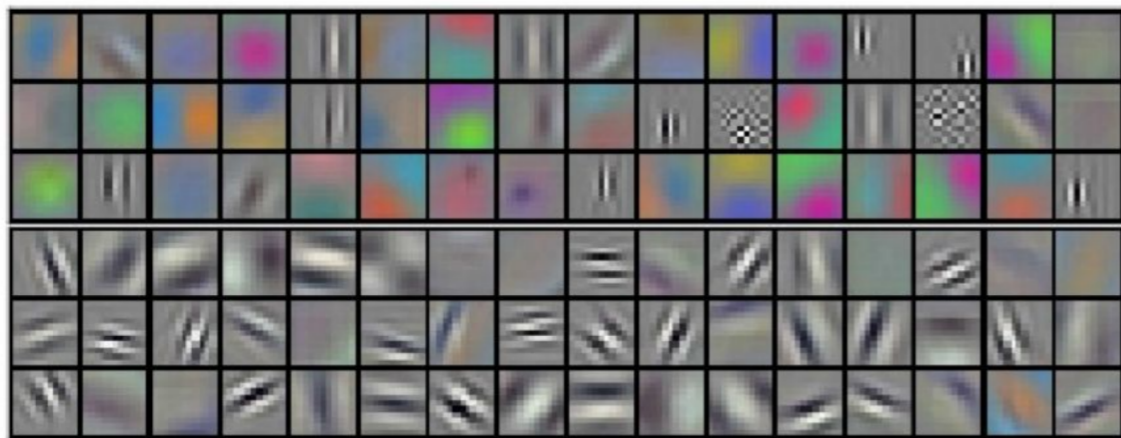
✦ CNN first layer filters (learned)

✦ PCA of Image Patches



equivariance + learning

Data statistics + regularization

# Geometric Deep Learning on other Domains

✦ Popular architectures as instances of GDL blueprint

| Architecture | Domain $\Omega$ | Symmetry Group $\mathfrak{G}$ |
|---|---|---|
| CNN | Grid | Translation |
| Spherical CNN | Sphere / SO(3) | Rotation SO(3) |
| Intrinsic / Mesh CNN | Manifold | Isometry Iso($\Omega$) / Gauge Symmetry SO(2) |
| GNN | Graph | Permutation $\Sigma_n$ |
| Deep Sets | Set | Permutation $\Sigma_n$ |
| Transformer | Complete Graph | Permutation $\Sigma_n$ |
| LSTM | 1D Grid | Time warping |

[M. Bronstein et al.: "Geometric Deep Learning"]

Graph $G = (V, E)$

Node features $\mathcal{X}(G)$

functions $\mathcal{F}(\mathcal{X}(\Omega))$

Permutation group $\Sigma_n$

**Permutation matrix P**

$$\mathbf{PX} = \left(x_{\pi^{-1}(i),j}\right)$$

**Message passing**

$$\mathbf{F}(\mathbf{PX}, \mathbf{PAP}^\top) = \mathbf{PF}(\mathbf{X}, \mathbf{A})$$

*Convolutional*

$$\mathbf{h}_i = \phi\left(\mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} c_{ij}\psi(\mathbf{x}_j)\right)$$
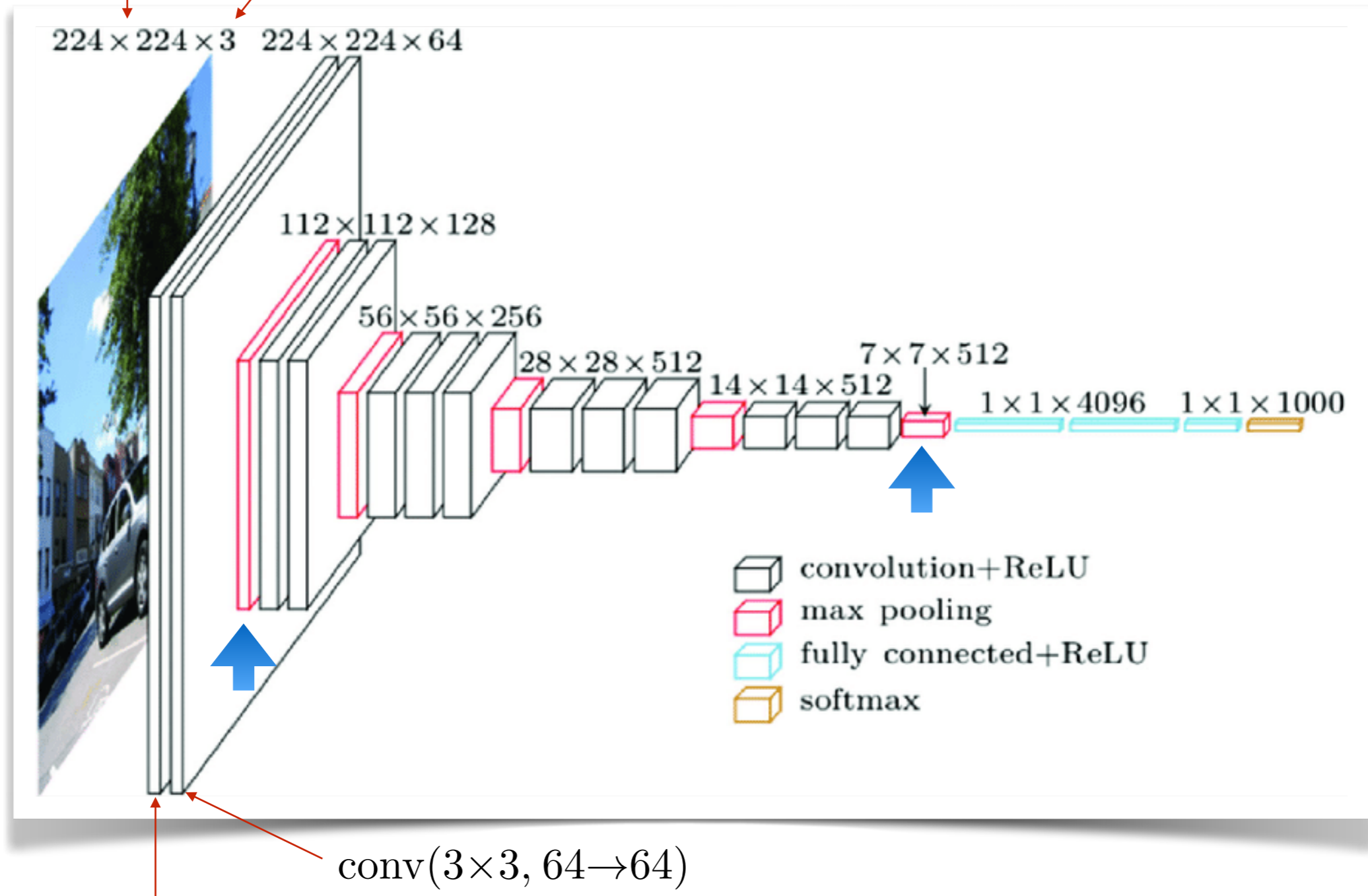
*Attentional*

$$\mathbf{h}_i = \phi\left(\mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} a(\mathbf{x}_i, \mathbf{x}_j)\psi(\mathbf{x}_j)\right)$$

*Message-passing*

$$\mathbf{h}_i = \phi\left(\mathbf{x}_i, \bigoplus_{j \in \mathcal{N}_i} \psi(\mathbf{x}_i, \mathbf{x}_j)\right)$$

[M. Bronstein et al.: "Geometric Deep Learning"]

Spatial size of the input image

channels

224 × 224 × 3    224 × 224 × 64

112 × 112 × 128

56 × 56 × 256

28 × 28 × 512

14 × 14 × 512

7 × 7 × 512

1 × 1 × 4096    1 × 1 × 1000

- convolution+ReLU
- max pooling
- fully connected+ReLU
- softmax

$\mathrm{conv}(3{\times}3, 64{\to}64)$

Result of $\mathrm{conv}(K{\times}K, 3{\to}64)$ followed by ReLU
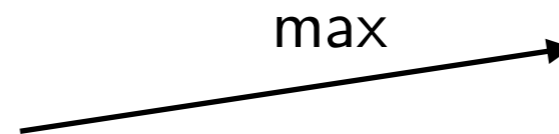
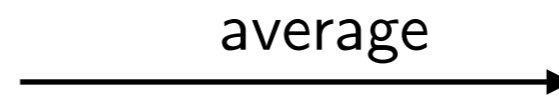✦ Eventually want to classify -> need to reduce spatial dimensions

# Pooling

✦ Following approaches are used to reduce the spatial resolution:

- **max pooling**

- **average pooling**
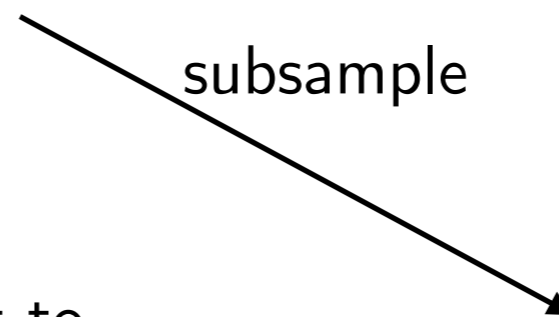
- subsampling -> **convolution with stride**
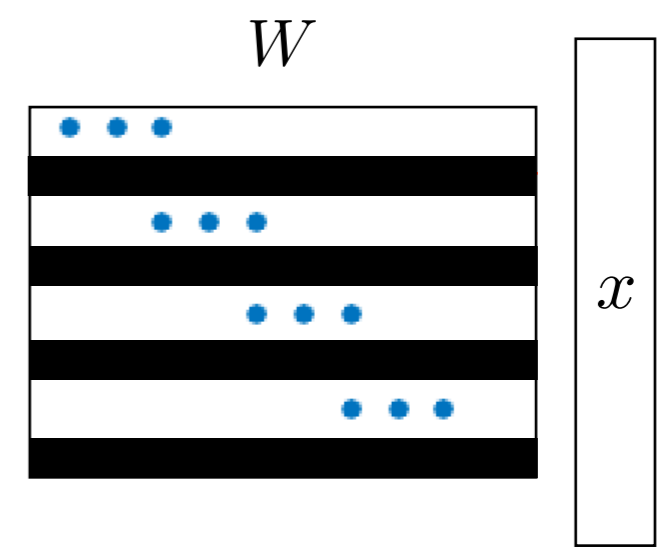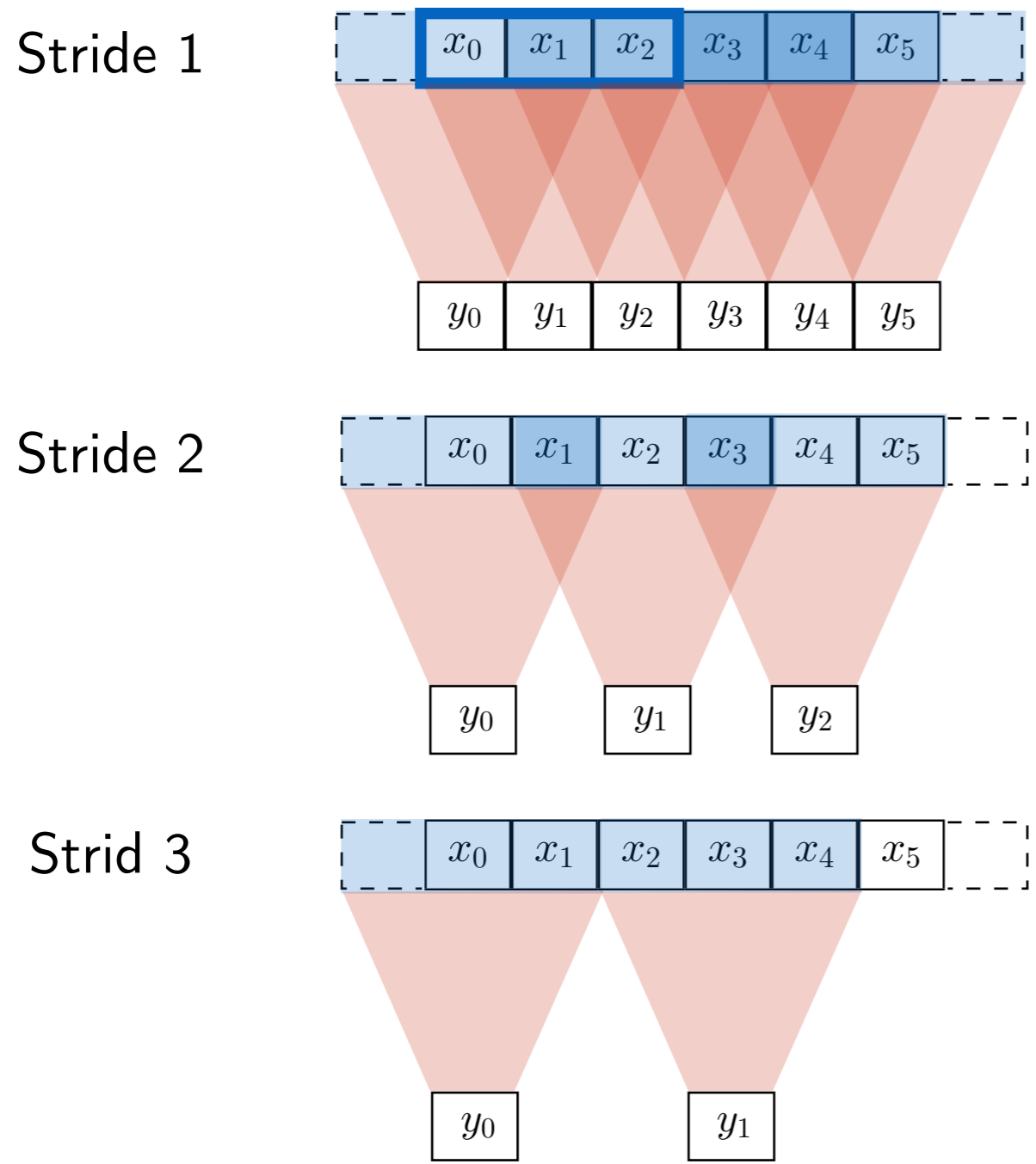


✦ max and average pooling are invariant to permutations of responses within a cell

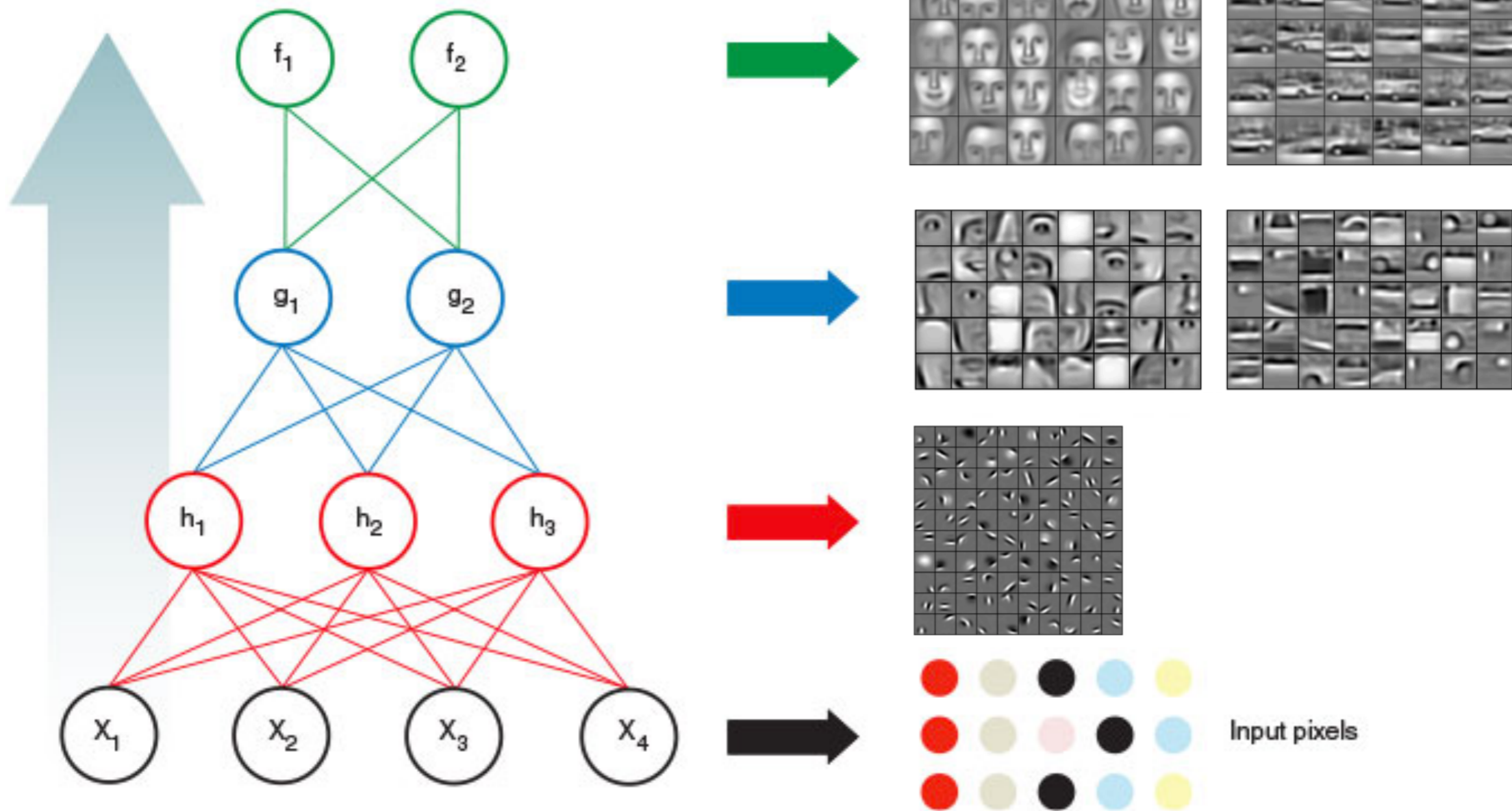✦ Once spacial resolution has been decreased, we can afford to increase the number of channels

✦ Full convolution + subsampling is equivalent to calculating the result at the required locations only, stepping with a **stride**

Stride 1

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

| $y_0$ | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ |

Stride 2

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

| $y_0$ | $y_1$ | $y_2$ |

Strid 3

| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

| $y_0$ | $y_1$ |

$W$

$x$

All variants and detail: [Dumoulin, Visin (2018): A guide to convolution arithmetic for deep learning]

✦ In networks trained for different complex problems many intermediate layers activations correspond object parts

Faces                    Cars



Input pixels

# Hierarchy of Parts Phenomenon

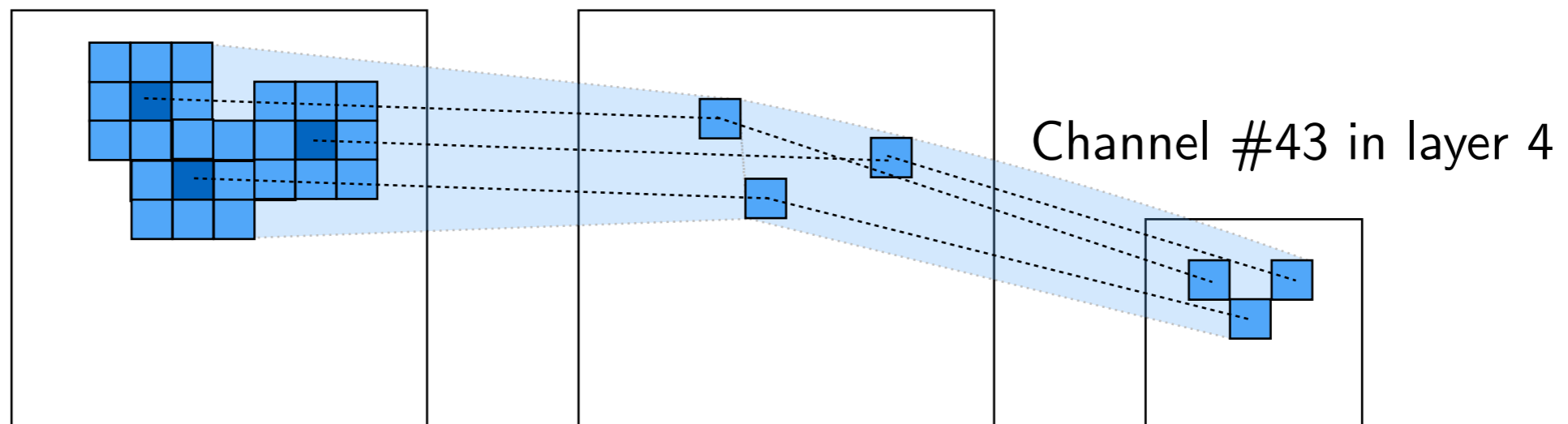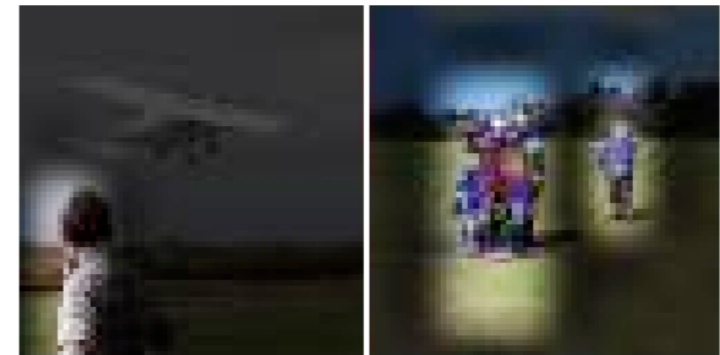✦ In networks trained for different complex problems many intermediate layers activations correspond object parts
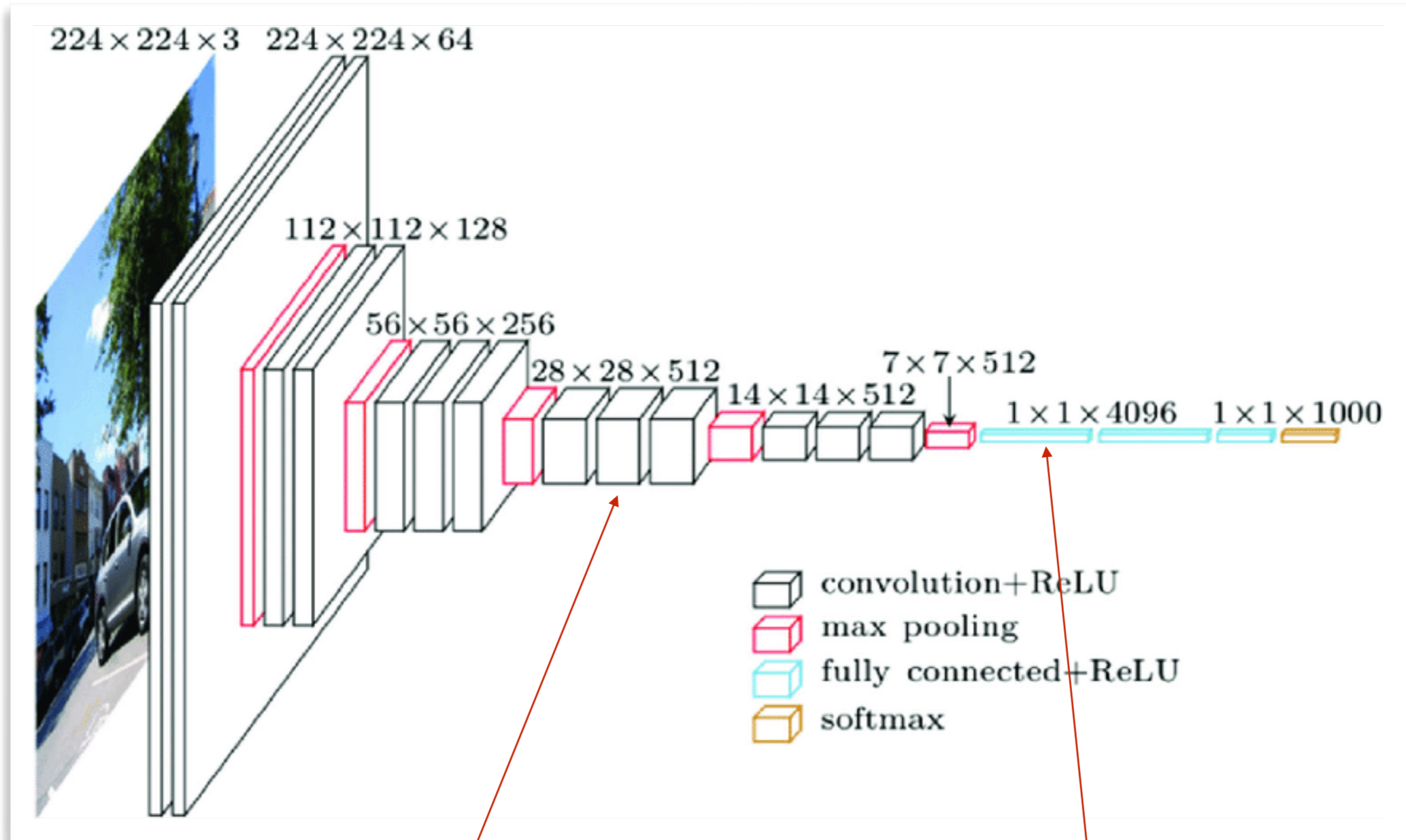
lamps in places net      wheels in object net      people in video net



Channel #43 in layer 4

too many weights here --
could use structured convolution
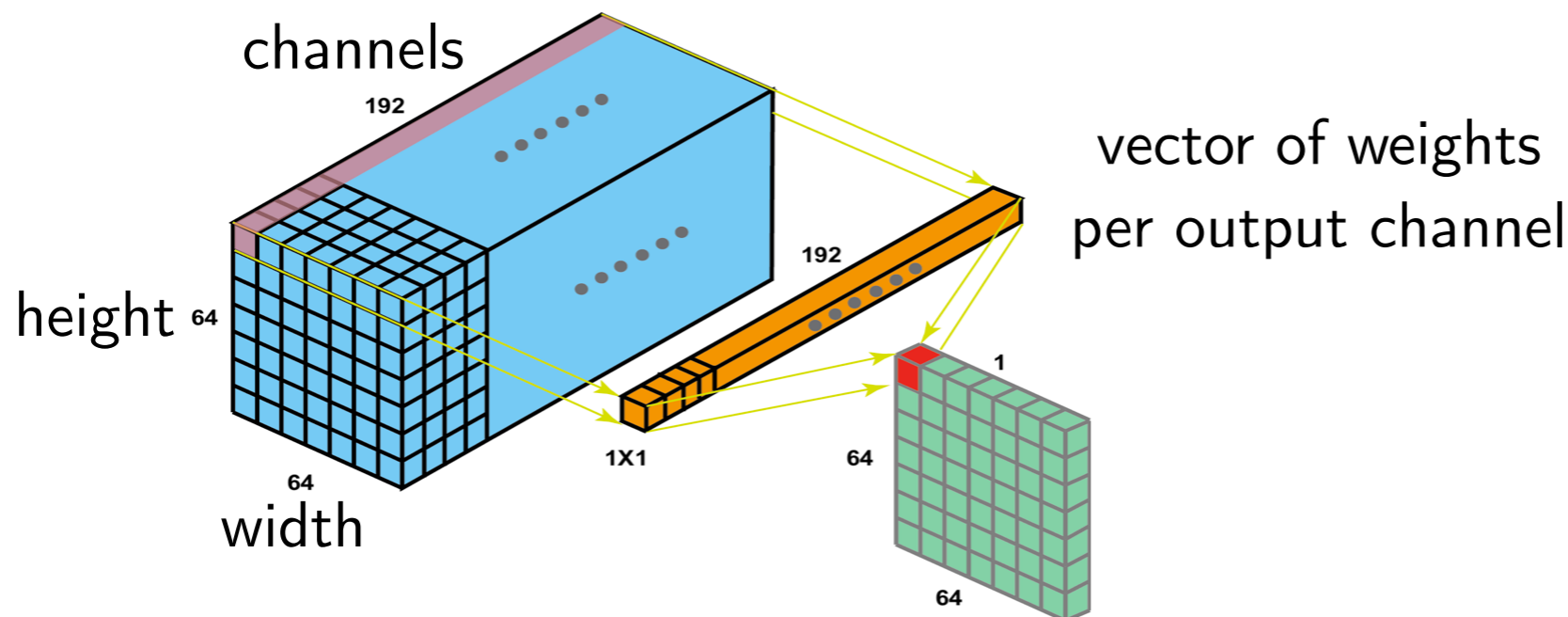
1x1 convolution for input of size 1x1
is equivalent to fully connected
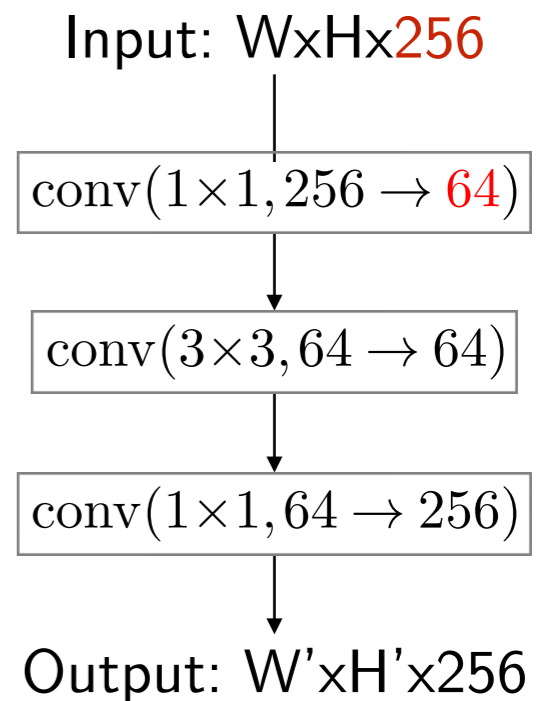
# 1x1 and Structured Convolutions

◆ Kernel size $1{\times}1$:

$$y_{o,i,j} = \sum_c \sum_{\Delta i=0} \sum_{\Delta j=0} w_{o,c,\Delta i,\Delta j} \; x_{c,i+\Delta i,j+\Delta j}$$

$$= \sum_c w_{o,c,0,0} \; x_{c,i,j}$$

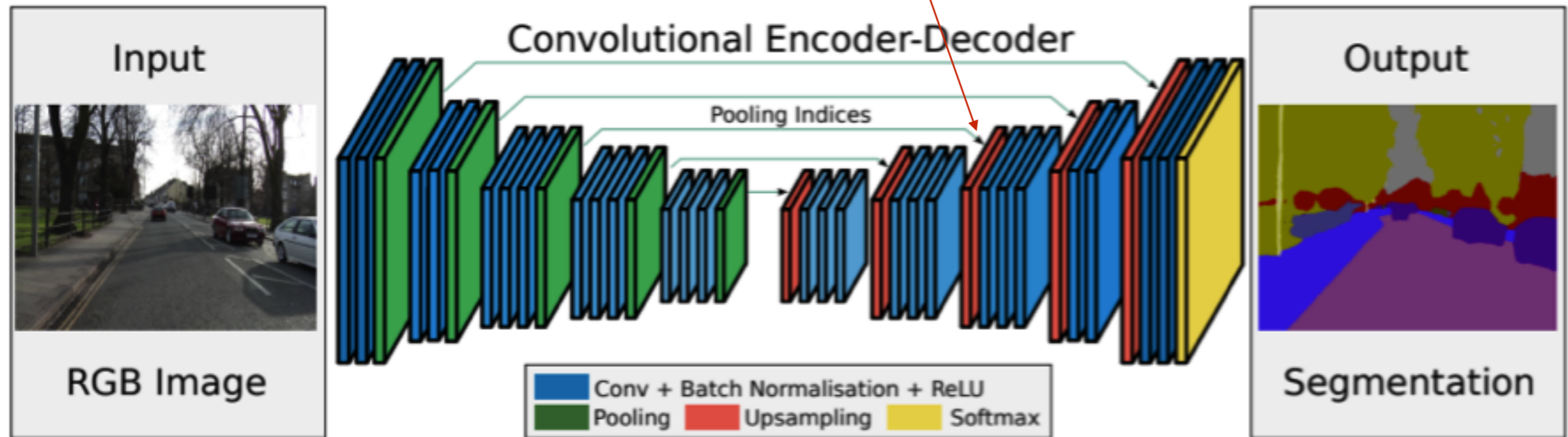◆ For all $i,j$ a linear transformation on channels with a matrix $w_{o,c,0,0}$



vector of weights
per output channel

Example $3{\times}3$, $256{\to}256$,
is too expensive, simplify:

Input: WxHx256

$\mathrm{conv}(1{\times}1, 256 \to 64)$

$\mathrm{conv}(3{\times}3, 64 \to 64)$

$\mathrm{conv}(1{\times}1, 64 \to 256)$

Output: W'xH'x256

◆ Useful to perform operations along channels dimension:
   ● Increase /decrease number of channels
   ● In combination with purely spatial convolution $=$ separable transform

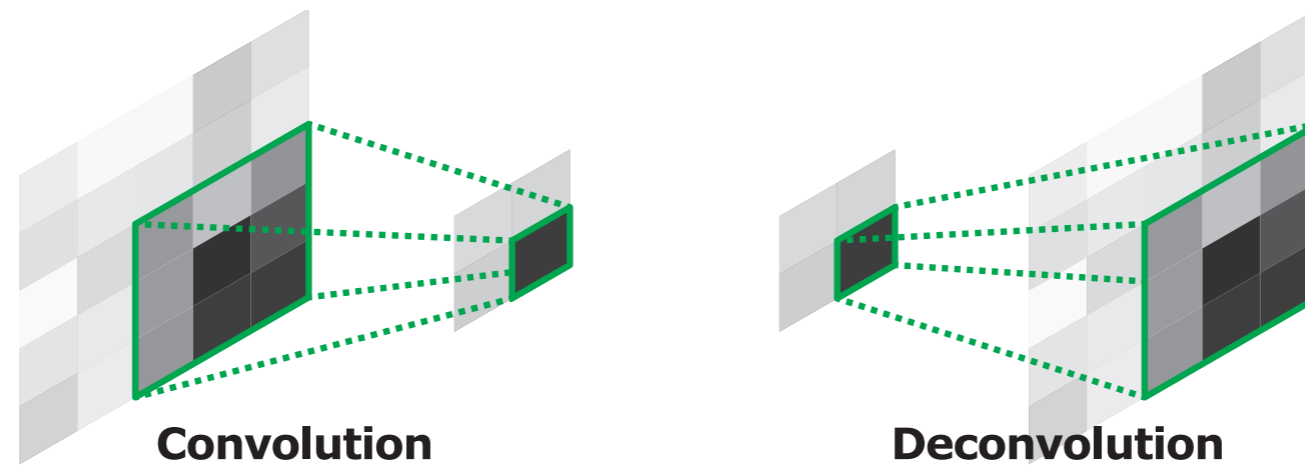Semantic segmentation architectures need unpooling / upsampling



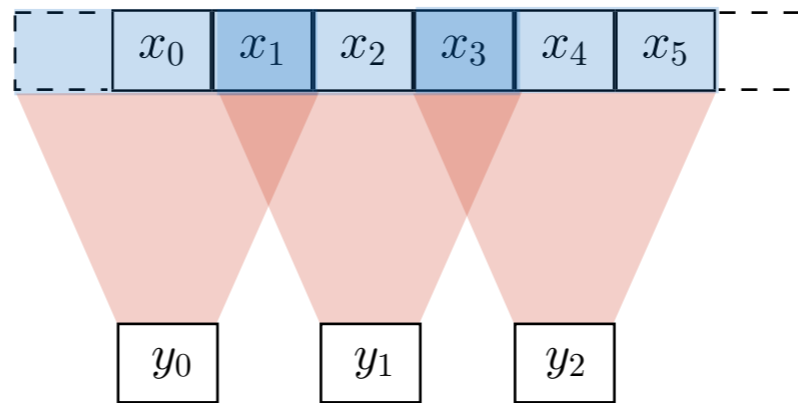We will look at up-sampling with "transposed" convolution ("deconvolution")
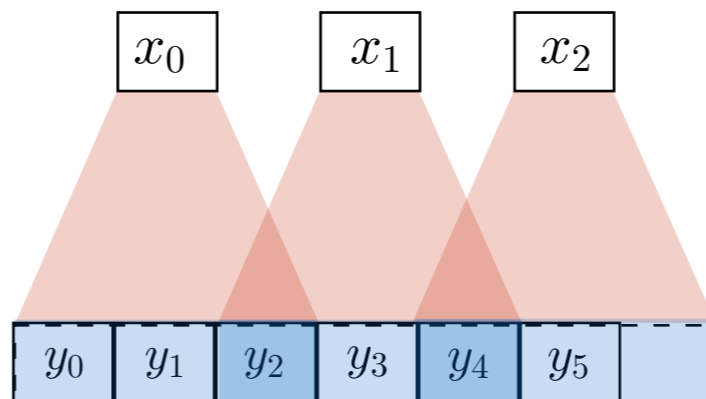
# Transposed Convolution

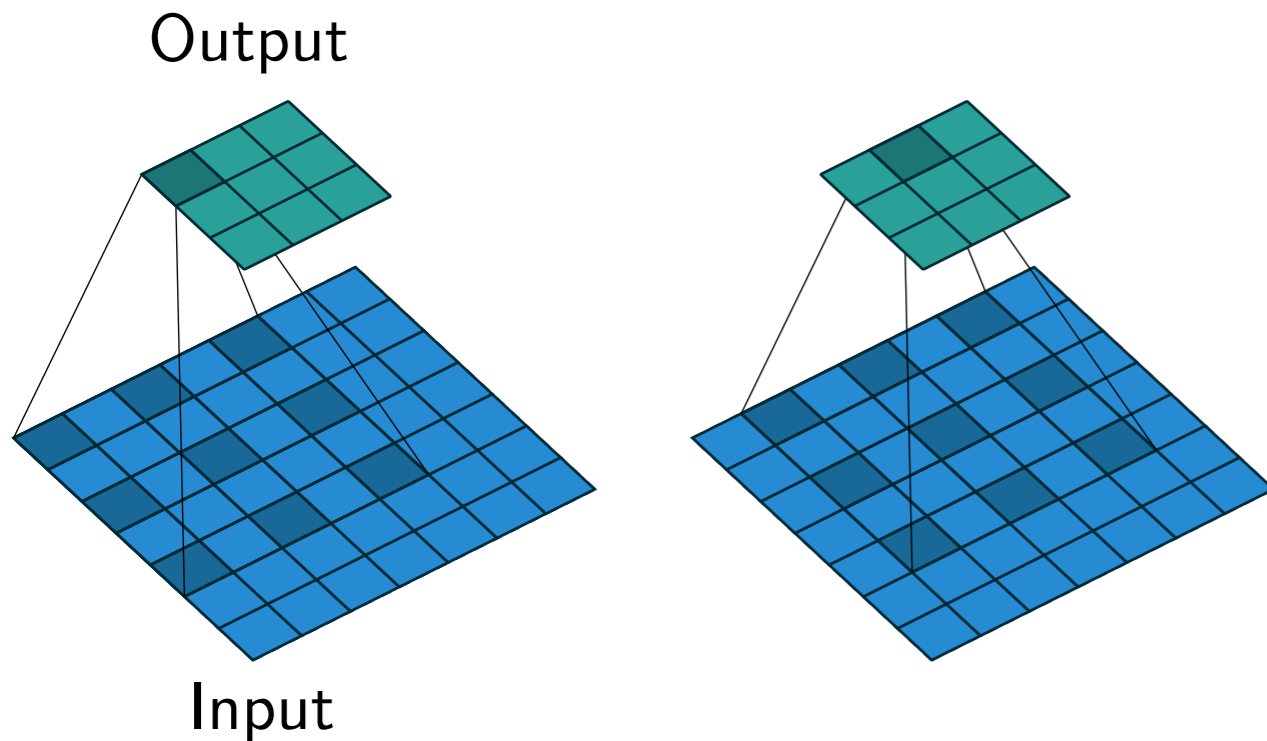✦ Deconvolution = Transposed strided convolution = backprop of strided convolution



**Convolution**   **Deconvolution**

Stride 2 Convolution

$x_0$ $x_1$ $x_2$ $x_3$ $x_4$ $x_5$

$y_0$ $y_1$ $y_2$

Stride 2 Deconvolution

$x_0$ $x_1$ $x_2$

$y_0$ $y_1$ $y_2$ $y_3$ $y_4$ $y_5$

# Sparse & Deformable Convolutions

✦ Want to increase receptive field size

- without decreasing spatial resolution and having too many layers

- Can increase kernel size, but it was also costly

- Can use a sparse mask for the kernel

**Dilated** convolutions

Can even learn sparse locations —

**deformable** convolutions