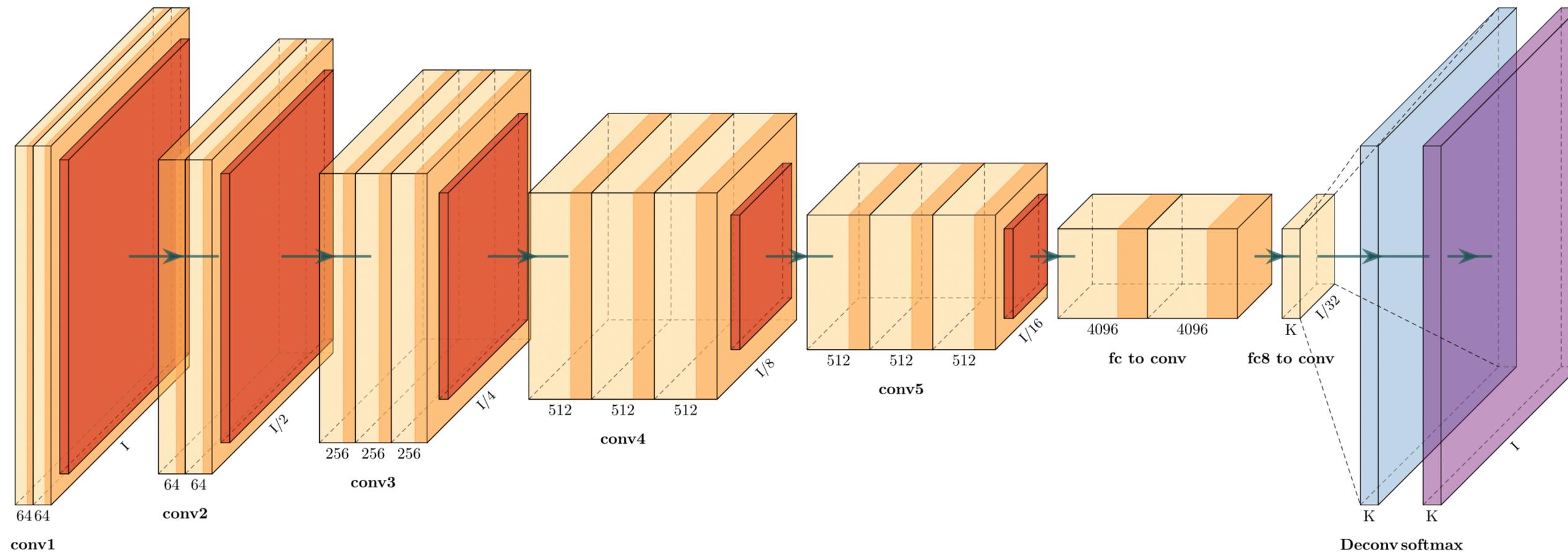


Convolutional neural networks (CNNs)

Stefan Funk - Erlangen Centre for astroparticle physics (ECAP)

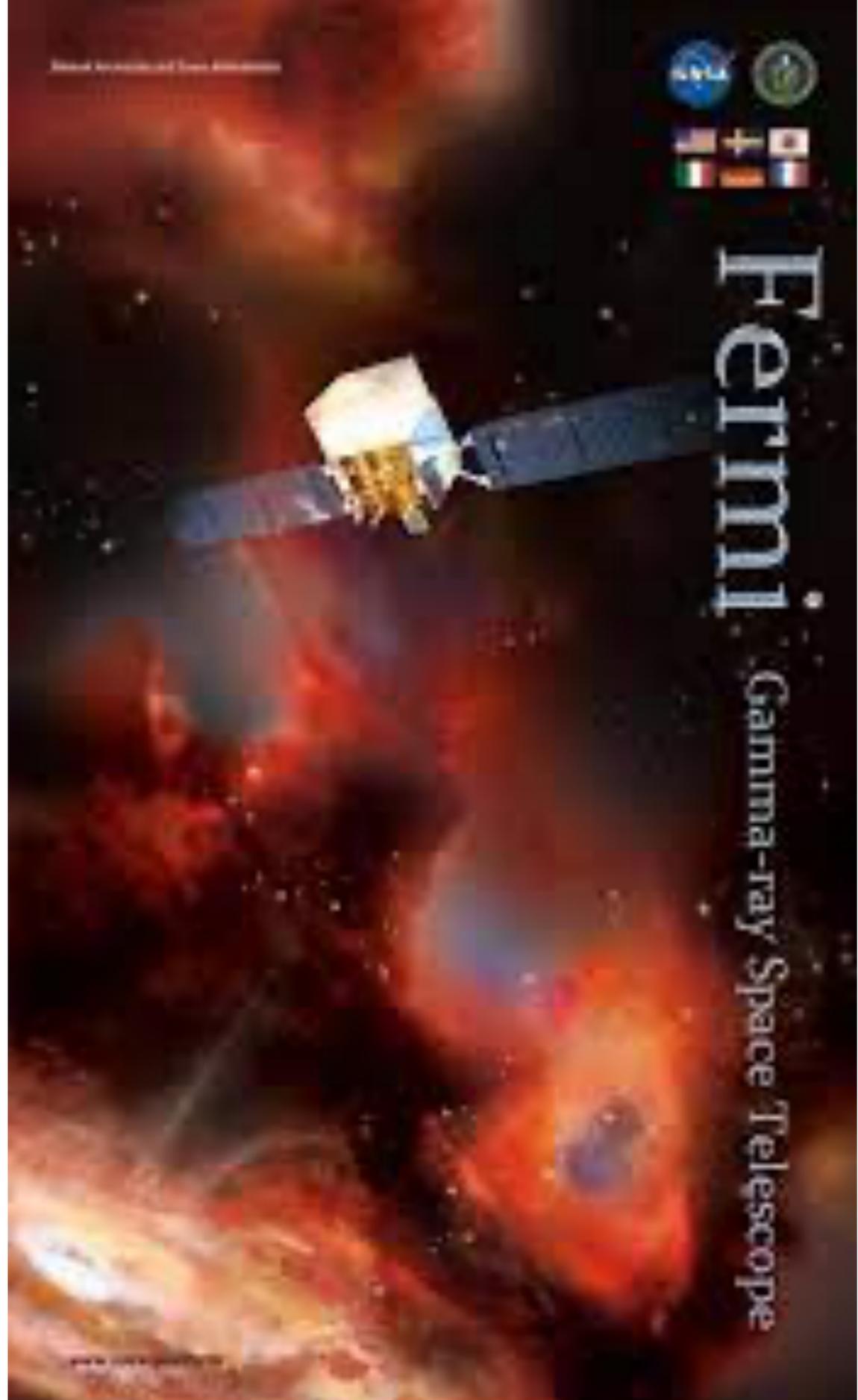
With input from a lot of people, in particular M. Meyer (U. Hamburg) and D. Malyshev (FAU)



me



me





me
and
you



CNNs IN THE PHYSICS DEPARTMENT AT FAU



Friedrich-Alexander-Universität
Faculty of Sciences

- At FAU (Physics department):
 - 1st semester: Datenverarbeitung
 - Later Wahlfächer like: “Machine learning and data analysis in science” and “Machine learning for physicists”
 - https://florianmarquardt.github.io/deep_learning_basics_linkmap.html
 - Currently no single course in physics focussed on Deep Learning
- But of course also courses in other departments (e.g. Data Sciences, ...)
- What I will describe is one or two lectures on CNNs in the context of the above lectures
 - Will assume that everything about machine learning, ANN, back propagation, ... has been covered

GOAL OF THE LECTURE

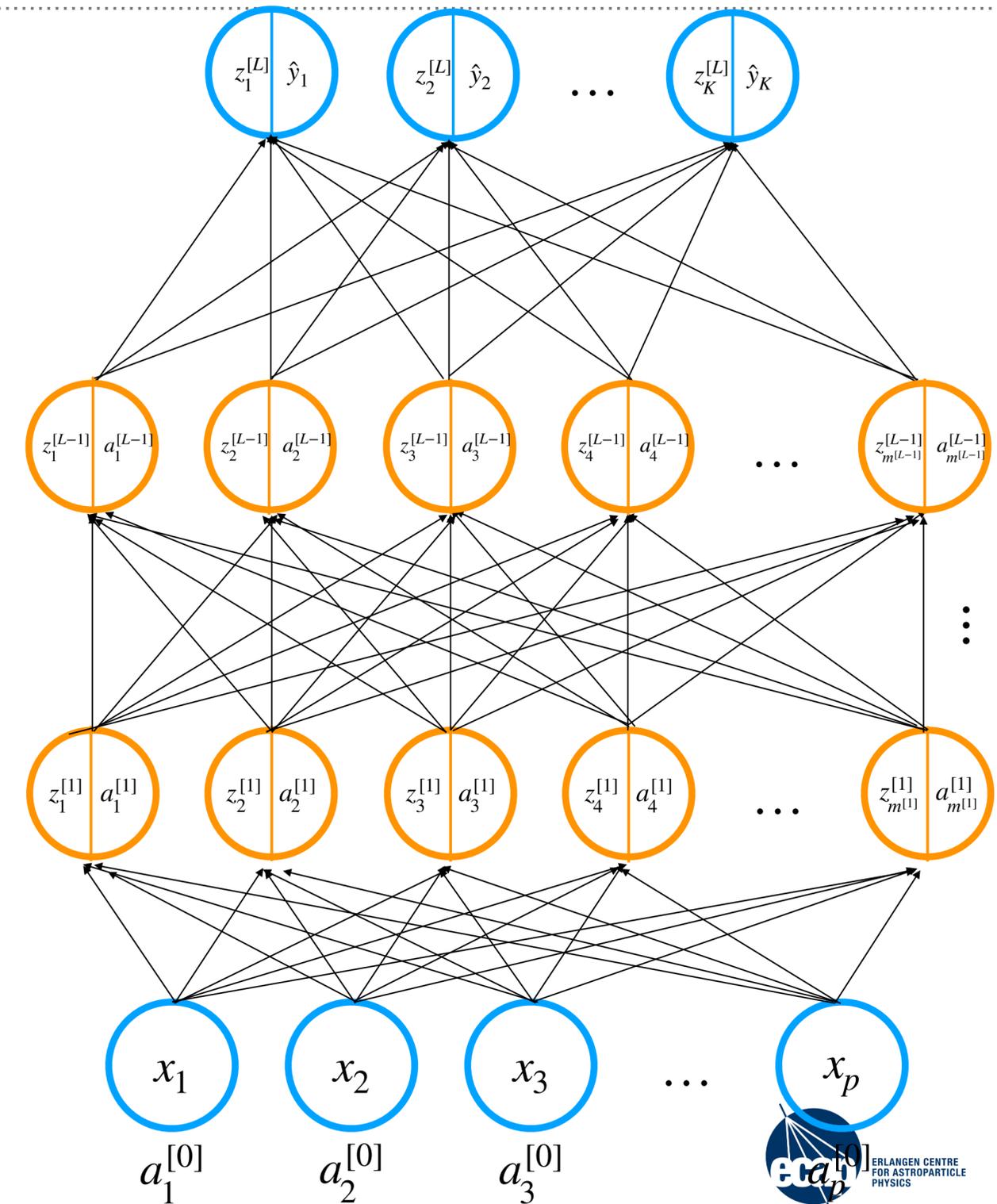
- At the end of this lecture, the students should
 - Understand the basic building blocks of a convolutional neural network
 - Understand why CNNs are very powerful in extracting features
 - Understand how modern CNNs are built
 - Be able to build their own CNN for a classification task



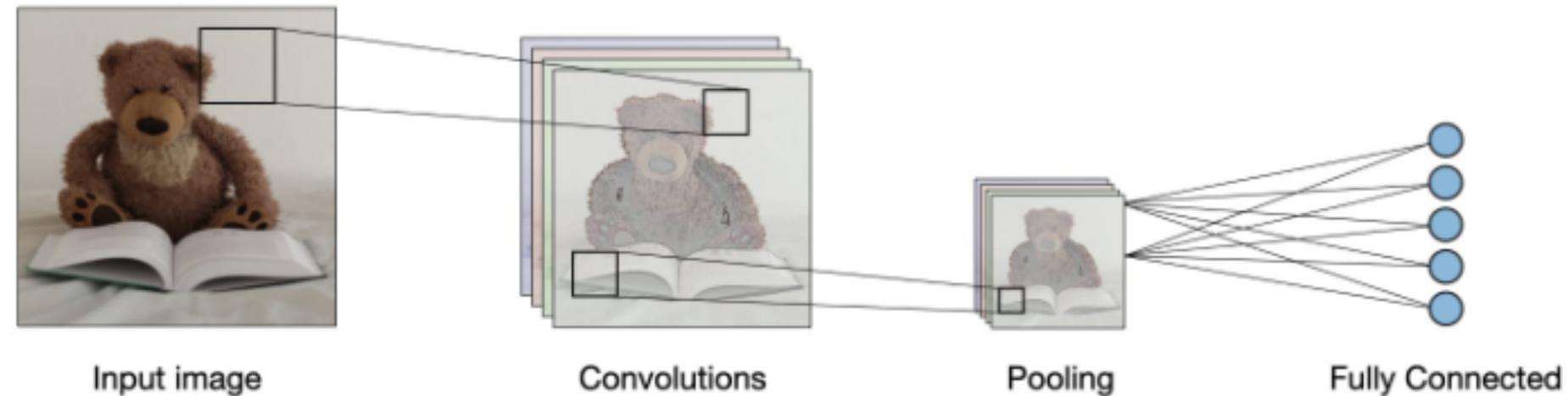
<https://www.istockphoto.com/de/vektor/gruppe-von-studenten-gm910098436-250650793>

REMINDER LAST LECTURE

- We introduced fully connected neural networks (aka multi-layer perceptrons) for 2-class and K -class classification and regression
- Learned about weights, learned about activation function (e.g. ReLU)
- We learned how to train a network for a given cost function J , using **back propagation**
- Neural networks are prone to overfitting, often necessary to include **regularization**



CNNs IN A NUT-SHELL



- Inspired by visual cortex (small region of cells sensitive to specific regions of the visual field)
 - Some neurons fire when they see vertical edges, others for diagonal ones (e.g. <https://arxiv.org/pdf/1406.3284.pdf>)
- CNN learns to look for features (=values of the filters) on its own through learning
 - Technically: slide convolution 'filter' over input volume
 - Learning part: determine optimal parameters of filters
- Able to derive patterns in a highly-complex input space

THE CHALLENGE



What We See

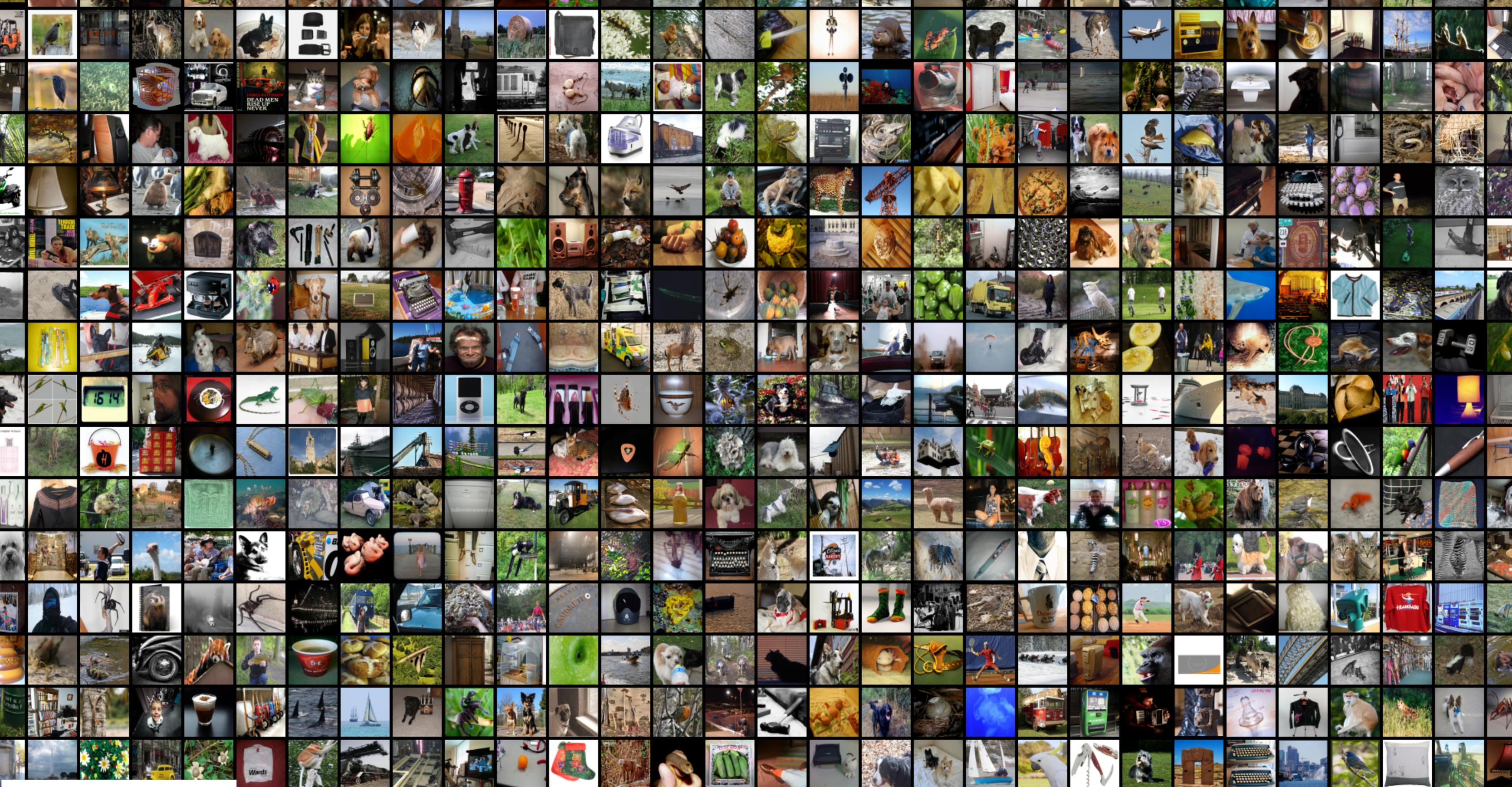
```
08 02 22 97 38 15 00 40 00 75 04 05 07 78 52 12 50 77 91 08
49 49 99 40 17 81 18 57 60 87 17 40 98 43 69 48 04 56 62 00
81 49 31 73 55 79 14 29 93 71 40 67 53 88 30 03 49 13 36 65
52 70 95 23 04 60 11 42 69 24 68 56 01 32 56 71 37 02 36 91
22 31 16 71 51 67 63 89 41 92 36 54 22 40 40 28 66 33 13 80
24 47 32 60 99 03 45 02 44 75 33 53 78 36 84 20 35 17 12 50
32 98 81 28 64 23 67 10 26 38 40 67 59 54 70 66 18 38 64 70
67 26 20 68 02 62 12 20 95 63 94 39 63 08 40 91 66 49 94 21
24 55 58 05 66 73 99 26 97 17 78 78 96 83 14 88 34 89 63 72
21 36 23 09 75 00 76 44 20 45 35 14 00 61 33 97 34 31 33 95
78 17 53 28 22 75 31 67 15 94 03 80 04 62 16 14 09 53 56 92
16 39 05 42 96 35 31 47 55 58 88 24 00 17 54 24 36 29 85 57
86 56 00 48 35 71 89 07 05 44 44 37 44 60 21 58 51 54 17 58
19 80 81 68 05 94 47 69 28 73 92 13 86 52 17 77 04 89 55 40
04 52 08 83 97 35 99 16 07 97 57 32 16 26 26 79 33 27 98 66
88 36 68 87 57 62 20 72 03 46 33 67 46 55 12 32 63 93 53 69
04 42 16 73 38 25 39 11 24 94 72 18 08 46 29 32 40 62 76 36
20 69 36 41 72 30 23 88 34 62 99 69 82 67 59 85 74 04 36 16
20 73 35 29 78 31 90 01 74 31 49 71 48 86 81 16 23 57 05 54
01 70 54 71 83 51 54 69 16 92 33 48 61 43 52 01 89 19 67 48
```

What Computers See

THE POWER OF CNNs



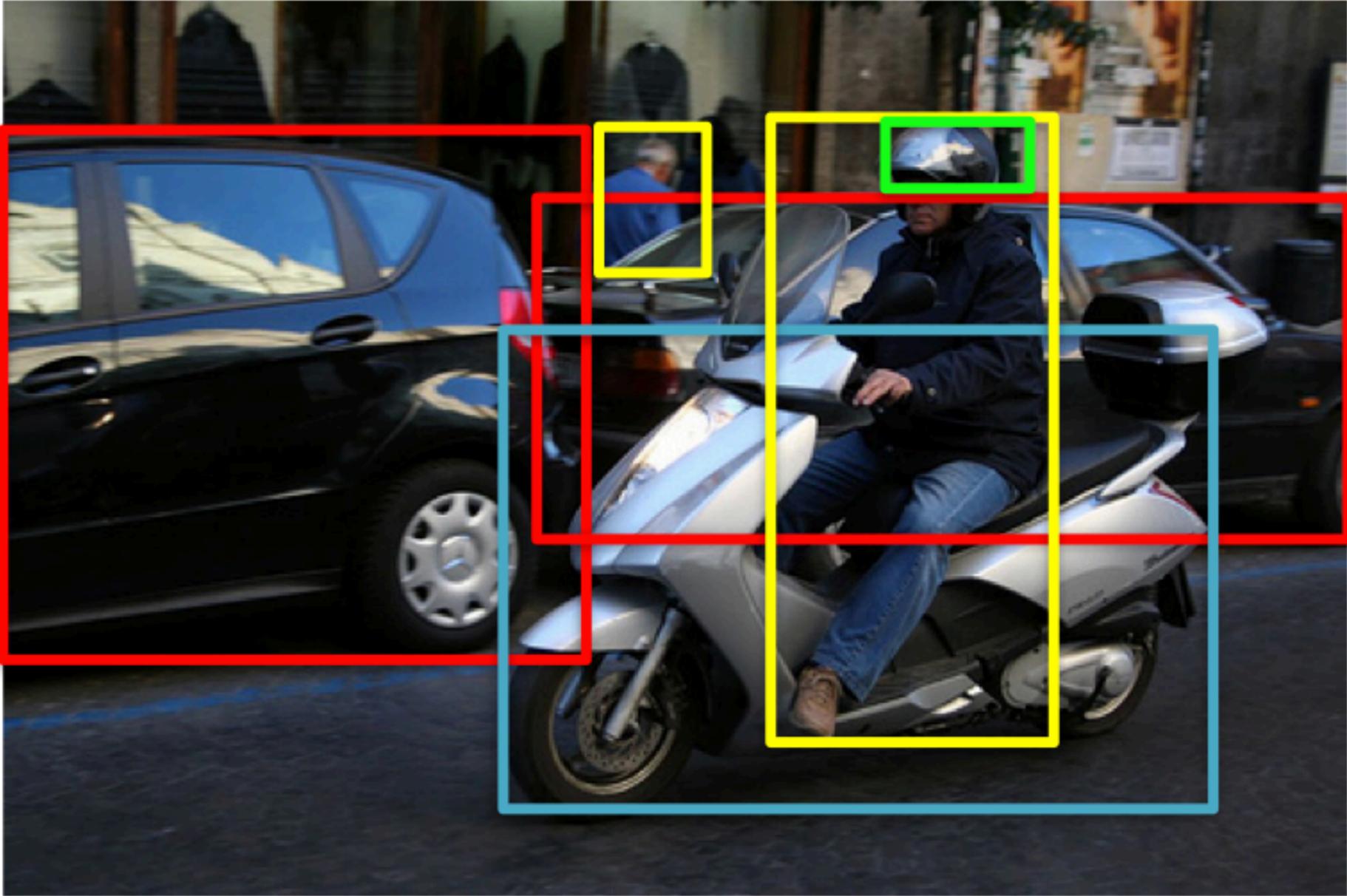
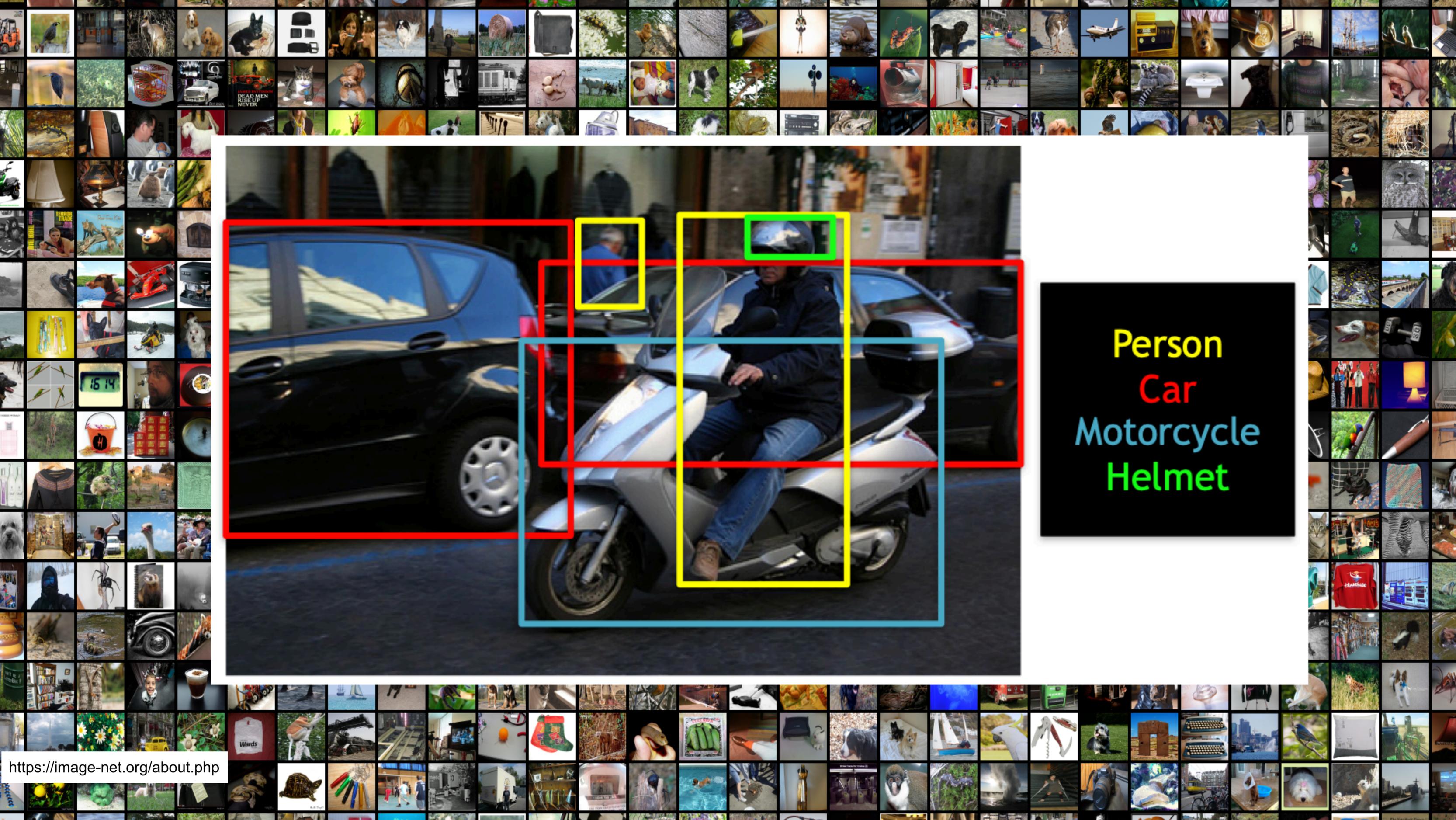
➤ Different image - same meaning



<https://image-net.org/about.php>

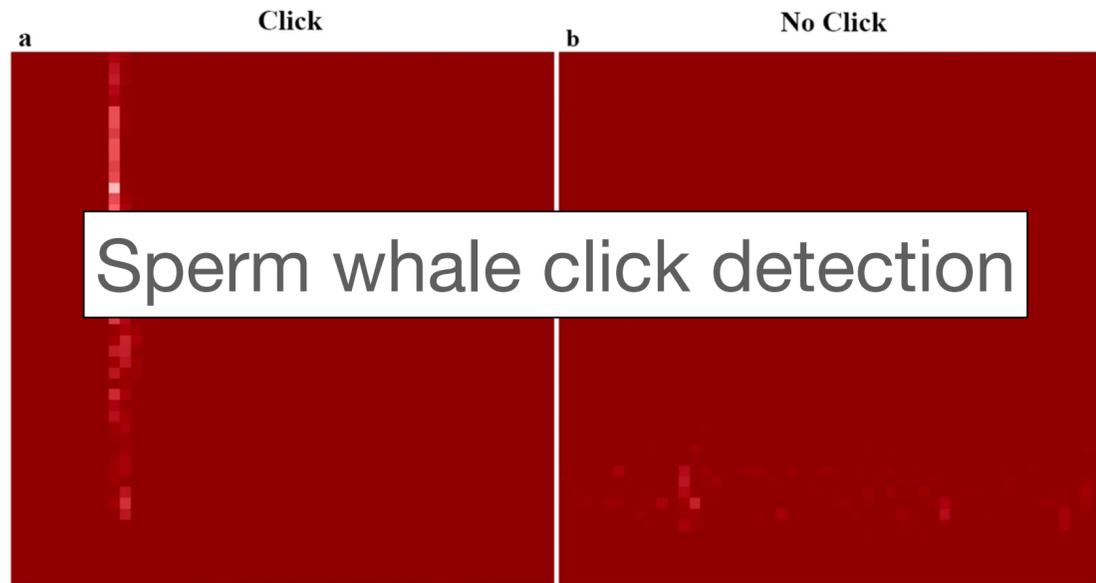


<https://image-net.org/about.php>



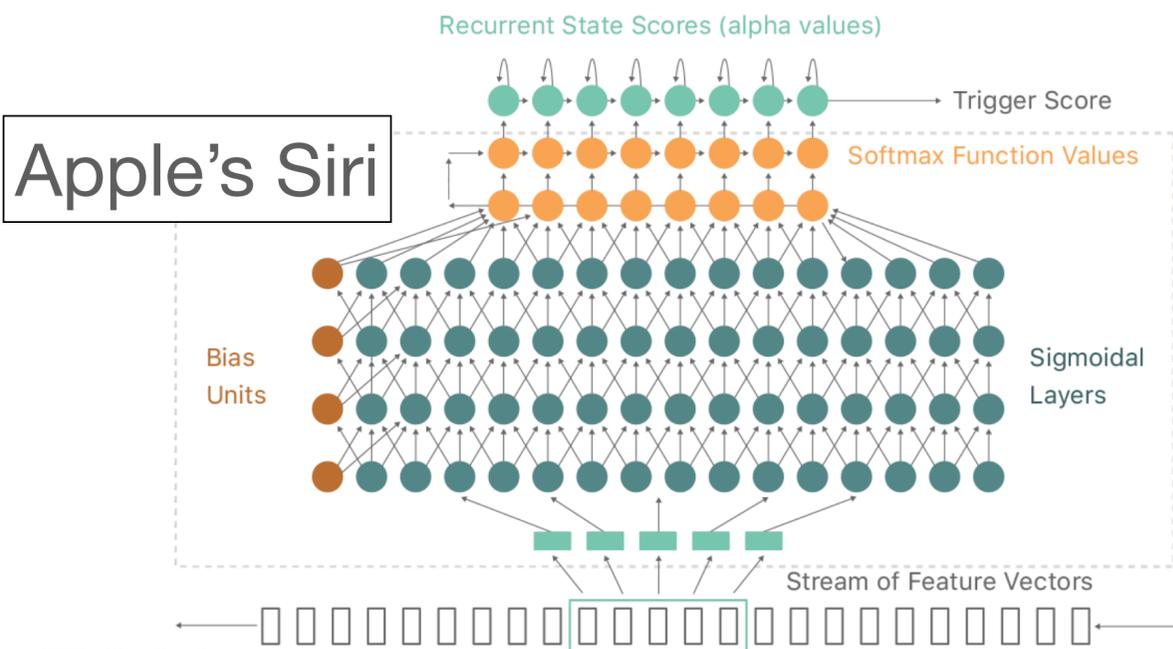
Person
Car
Motorcycle
Helmet

CNNs ARE EVERYWHERE

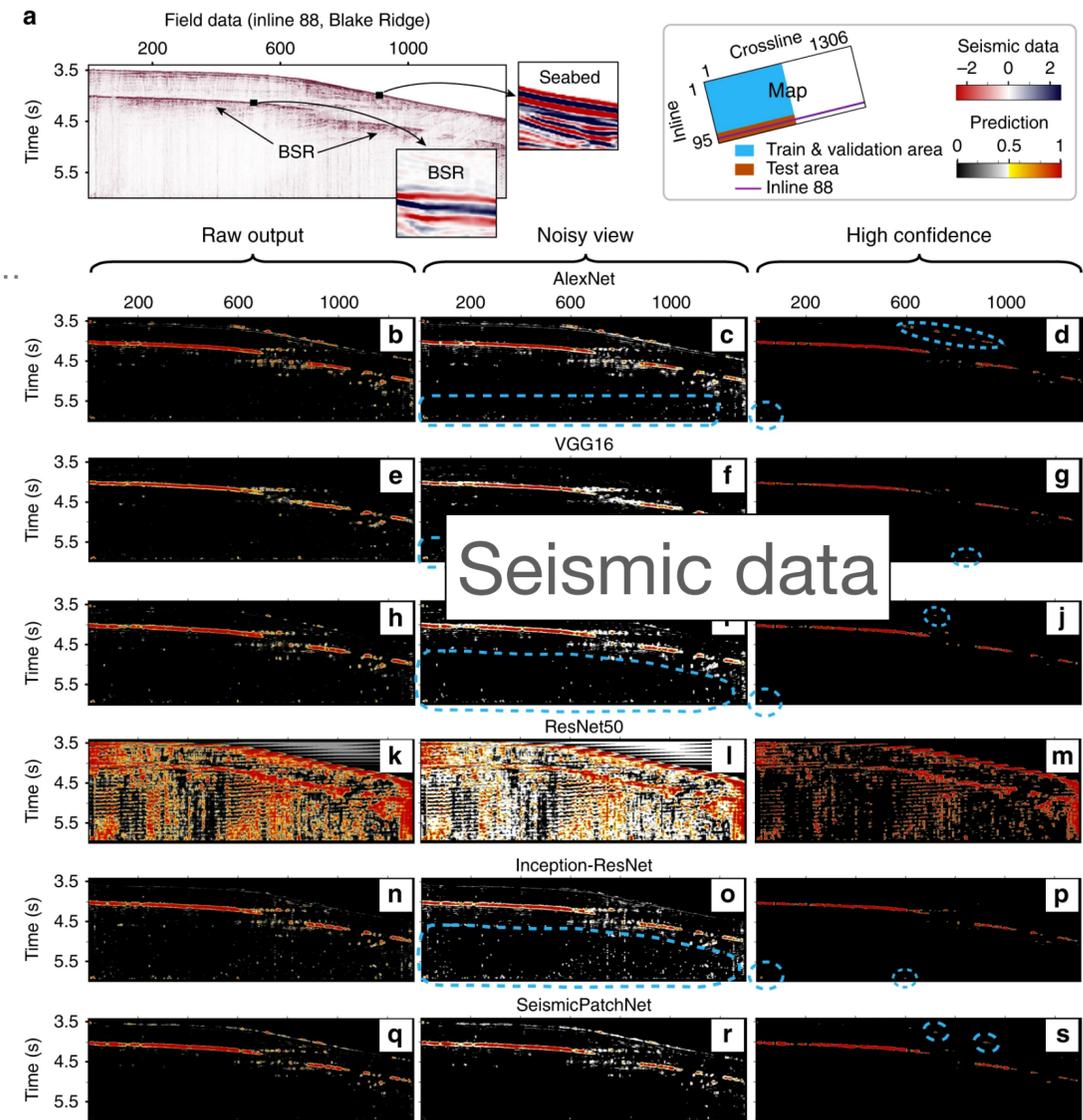


<https://www.nature.com/articles/s41598-019-48909-4>

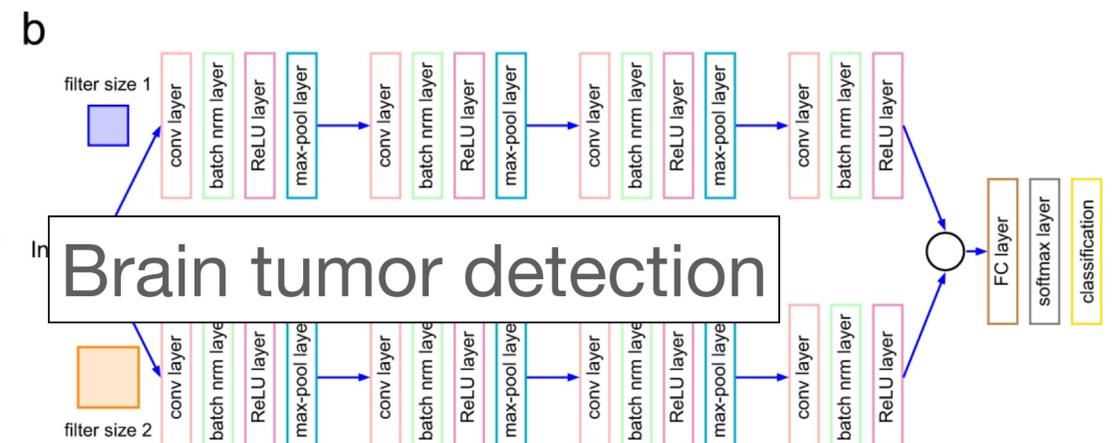
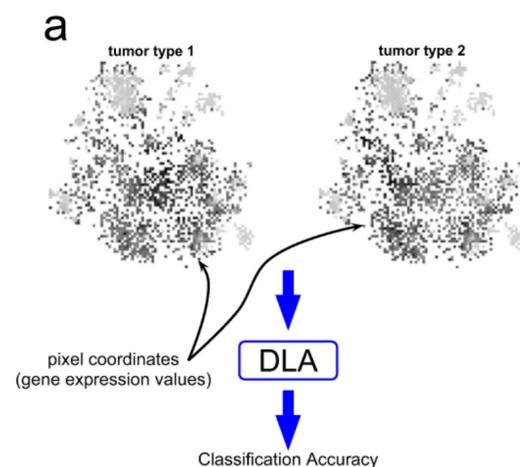
Apple's Siri



<https://machinelearning.apple.com/research/hey-siri>



<https://www.nature.com/articles/s41467-020-17123-6>



<https://www.nature.com/articles/s41598-019-47765-6>

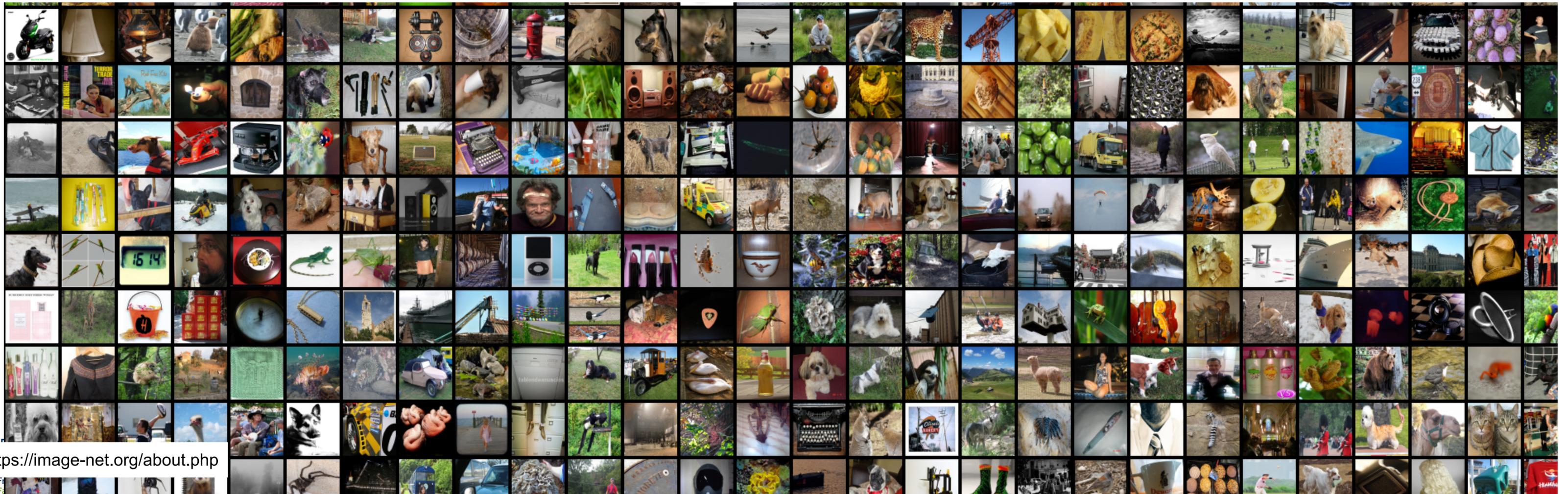
OK, let's get started understanding CNNs

STEPS THAT SHOULD BE COVERED IN A LECTURE ON CNNs

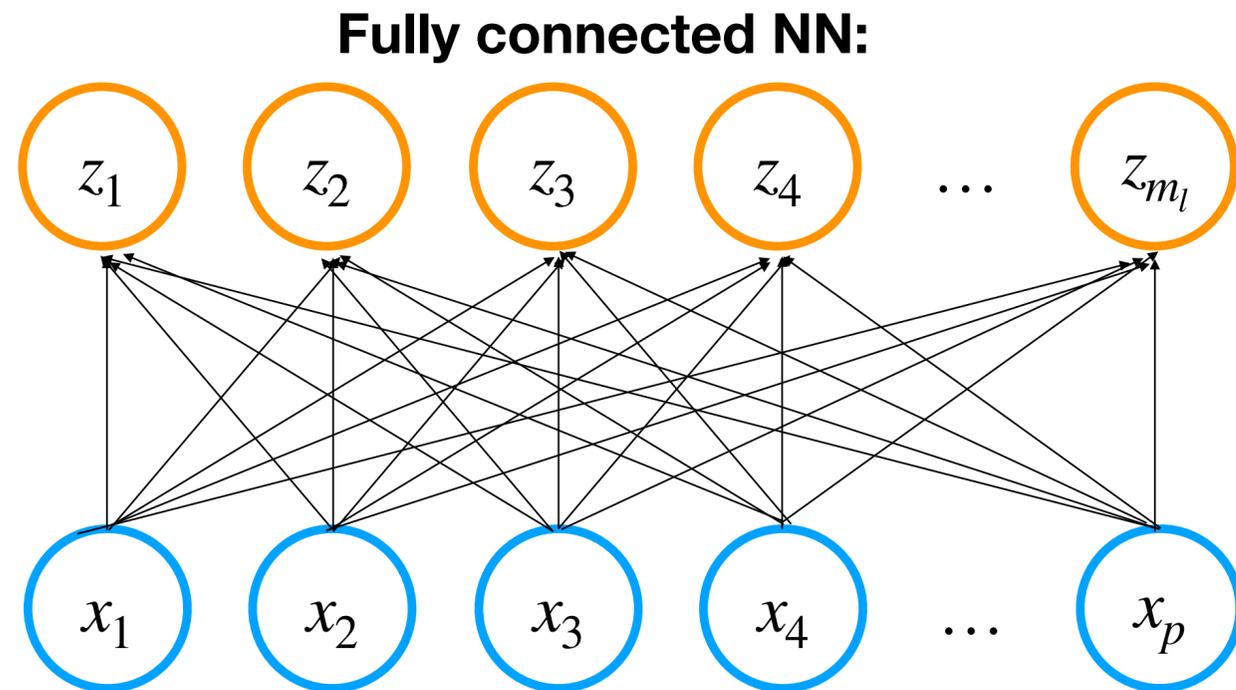
- The basic blocks of machine learning and CNNs
 1. The core of CNNs: convolution in artificial neural networks
 2. Learning in CNNs
 3. (Useful tools: Striding, zero padding and pooling)
 4. Typical CNN architectures
 5. Data augmentation and pre-processing

1 - THE CORE OF CNNs: THE TECHNICAL DETAILS THAT MATTER

- Basic assumption I: input is one or more images (2-dimensional data on a square grid, e.g. RGB)
- Basic assumption II: fully connected networks don't scale well (as the images get larger)



1 - THE CORE OF CNNs: FULLY CONNECTED NNS DON'T SCALE WELL

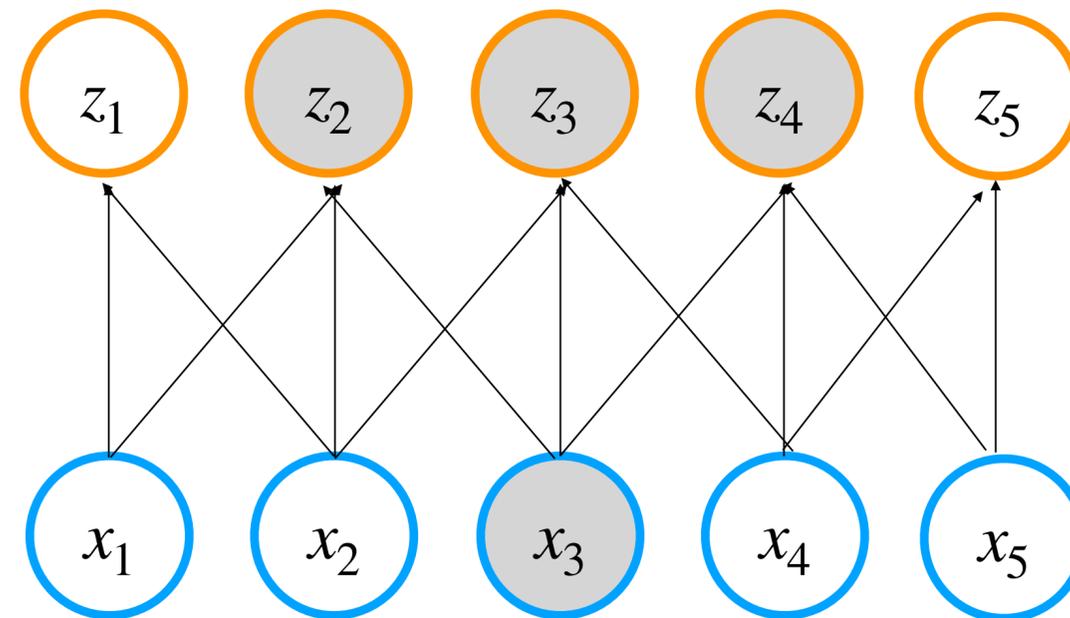


- Imagine an RGB image 32x32 pixels (i.e. $32 \times 32 \times 3 = 3072$ input values)
- Imagine a set of six 5x5 filters to learn the features of the image.
- Output: 28x28 for each of the six filters (=4704 output values).
- Fully connected network: $3072 \times 4704 = 14\text{M}$ weights.

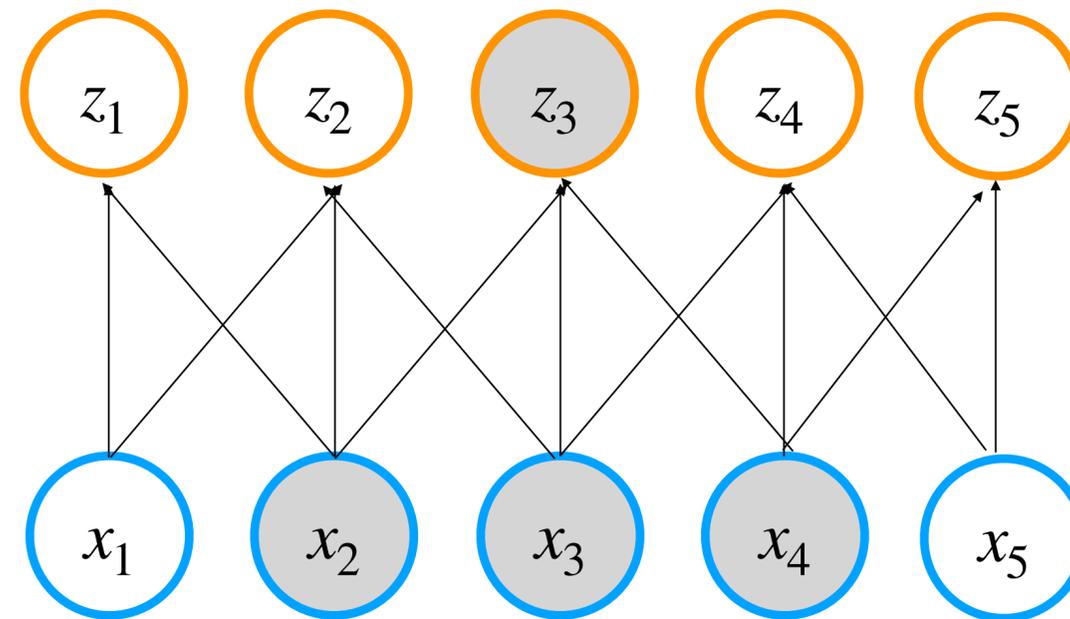
1 - THE CORE OF CNNs: THE TECHNICAL DETAILS THAT MATTER

- Basic assumption I: input is one or more images (2-dimensional data on a square grid, e.g. RGB)
- Basic assumption II: fully connected networks don't scale well (as the images get larger)
- Therefore create a network that contains the following features:
 - **Sparse interactions** (Kernel much smaller than input). Nodes are connected locally (convolution) but not fully connected (exact location of feature in image does not matter)
 - **Parameter sharing** (rather than learning a separate set of parameters for every location, only one set is learned - the Kernel weights - over the whole image)
 - **Translational equivariance** (if we shift the input, we also shift the output)

1 - THE CORE OF CNNs - SPARSE INTERACTIONS



1 – THE CORE OF CNNs: SPARSE INTERACTIONS

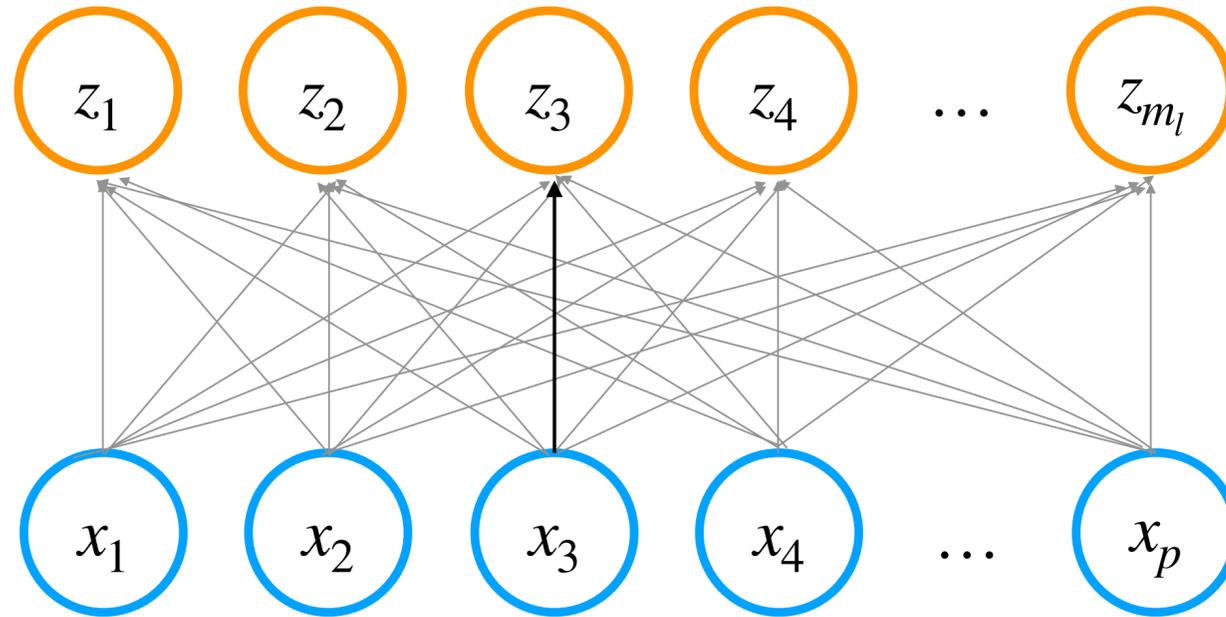


Receptive field

The receptive field in Convolutional Neural Networks (CNN) is the region of the input space that affects a particular unit of the network

1 – THE CORE OF CNNs: WEIGHT SHARING

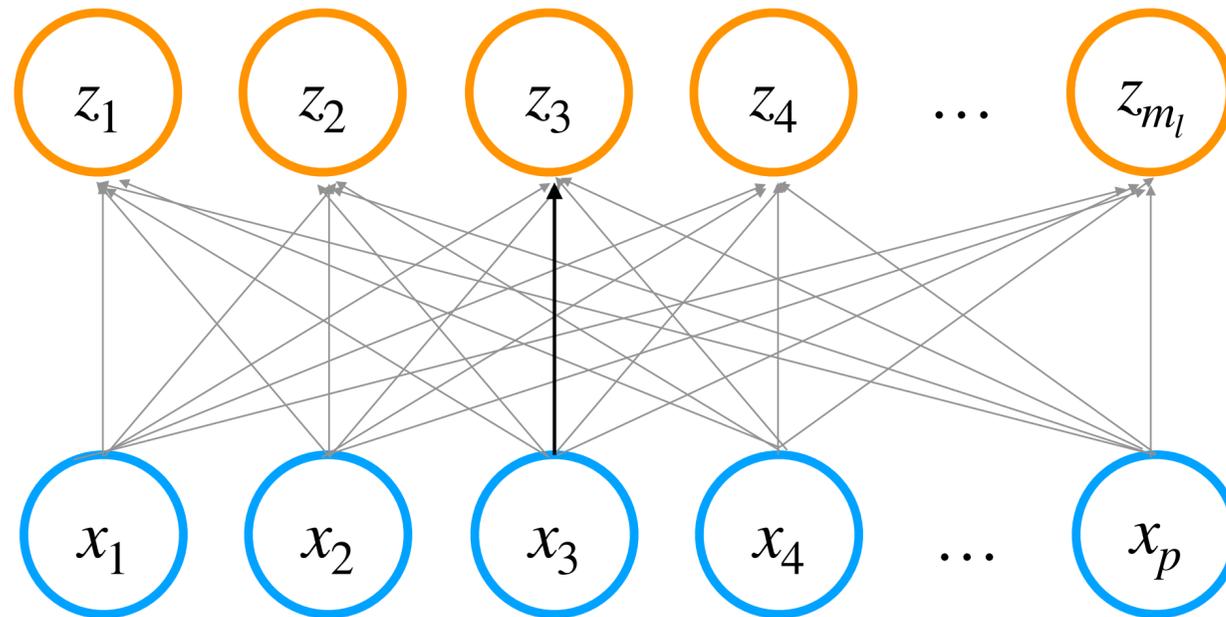
Fully connected NN: each weight only used once



- What is important in one part of the image, is likely also important in another part of the image

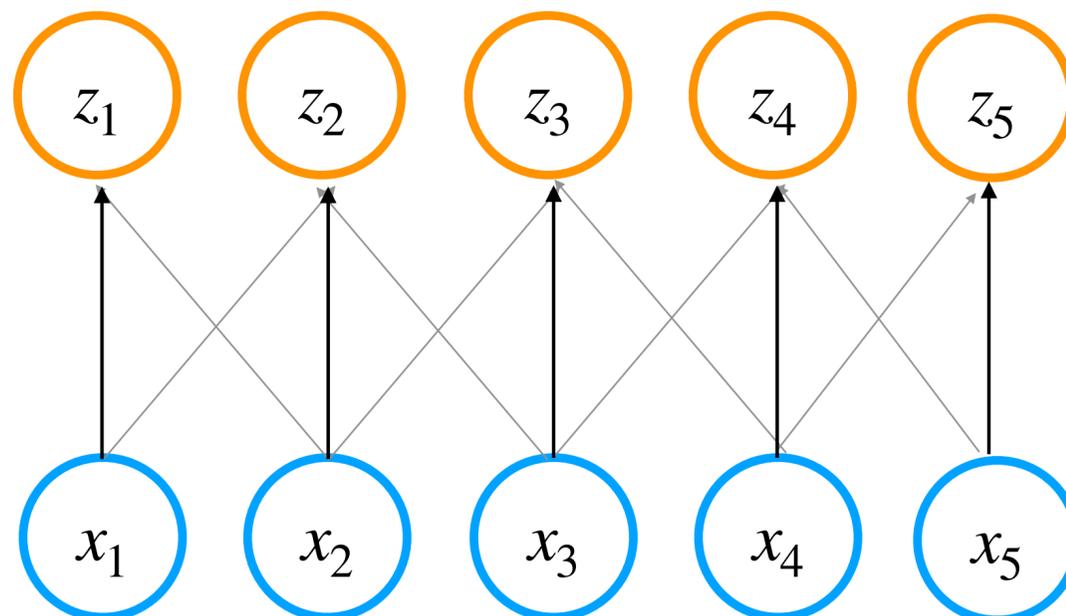
1 - THE CORE OF CNNs: WEIGHT SHARING

Fully connected NN: each weight only used once

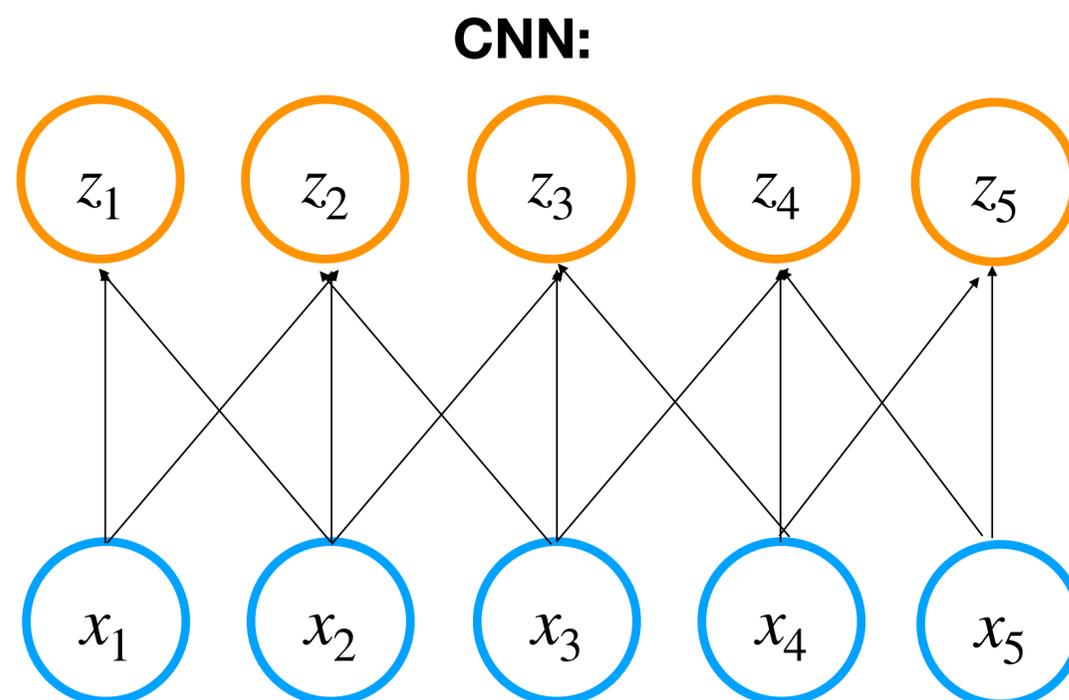
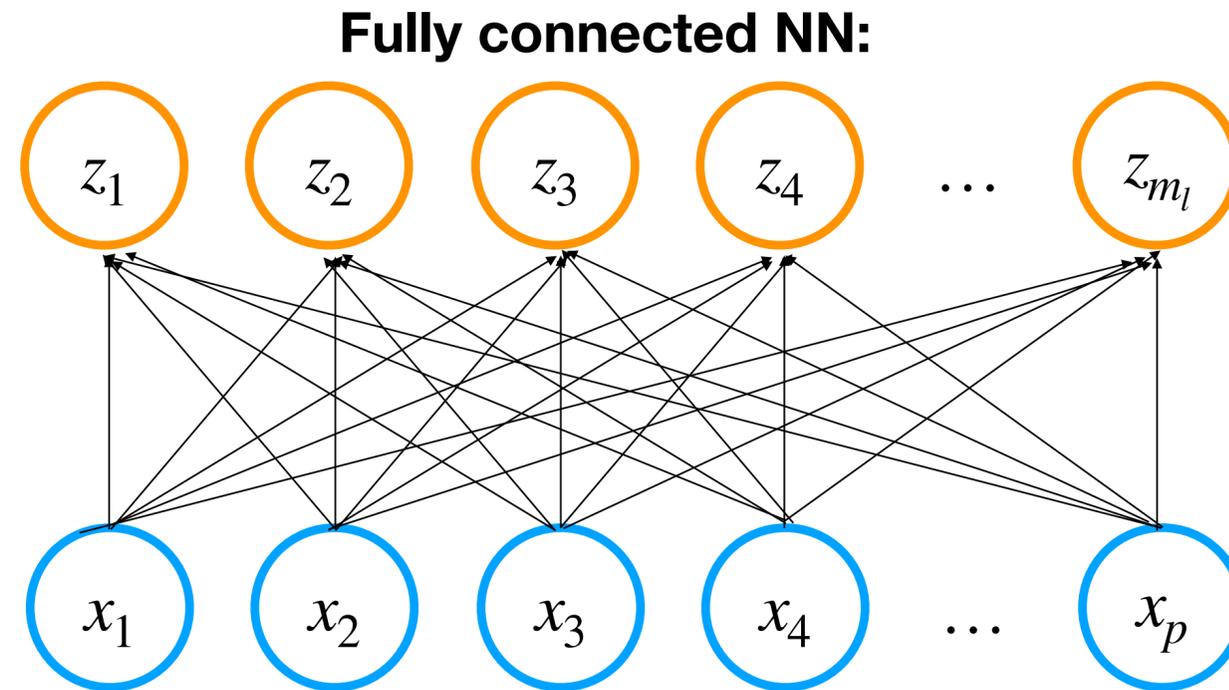


- What is important in one part of the image, is likely also important in another part of the image

CNN: weights shared among neurons



1 - THE CORE OF CNNs: SOME NUMBERS (SPARSE INTERACTIONS AND PARAMETER SHARING)



- Imagine an RGB image 32×32 pixels (i.e. $32 \times 32 \times 3 = 3072$ input values)
- Imagine a set of six 5×5 filters to learn the features of the image.
- Output: 28×28 for each of the six filters (=4704 output values).
- Fully connected network: $3072 \times 4704 = 14\text{M}$ weights.
- For CNN: Each filter has $(5 \times 5 + 1)$ parameters, for the three layers, i.e. $(5 \times 5 + 1) \times 3 \times 6 = 468$ weights

1 – THE CORE OF CNNs: UNDERSTANDING CONVOLUTION – MATHEMATICAL BASICS

- In mathematics, convolution defined as operation $(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau$ (f, g real-valued functions)

- Convolution is commutative:

$$\begin{aligned}(f * g)(t) &= \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau, \quad \text{Substitution: } x = t - \tau, \frac{dx}{d\tau} = -1 \\ &= - \int_{\infty}^{-\infty} f(t - x)g(x)dx = \int_{-\infty}^{\infty} g(\tau)f(t - \tau)dx = (g * f)(t)\end{aligned}$$

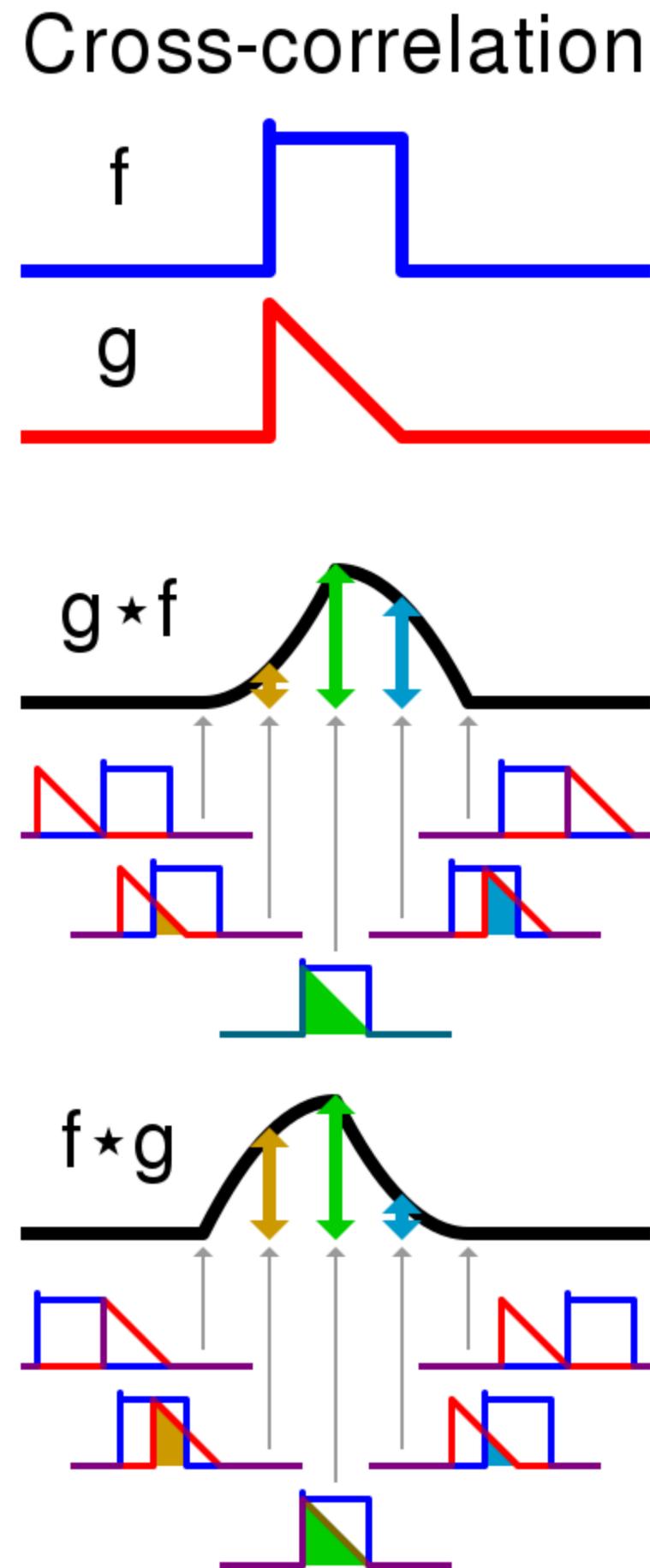
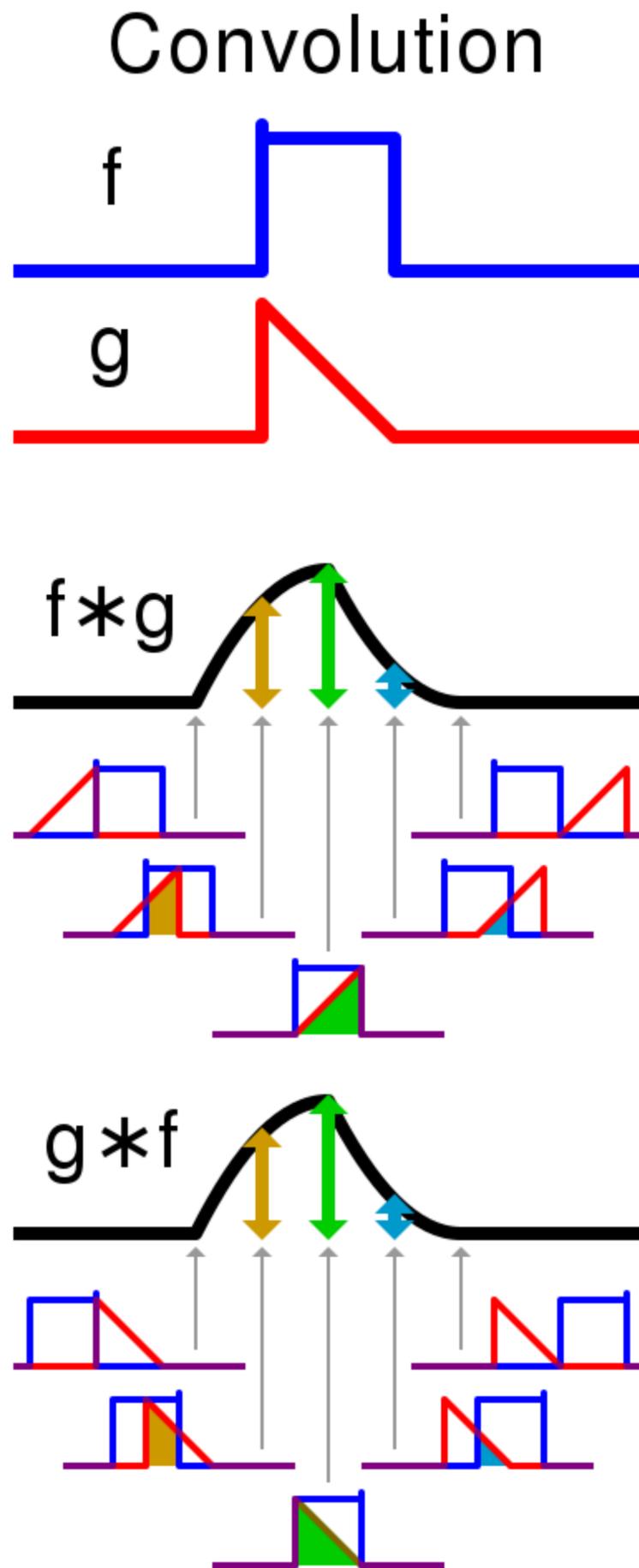
- In discrete form convolution can be written as: $(f * g)(n) = \sum_{m=-\infty}^{\infty} f(m)g(n - m)$

- In the CNN: two-dim image $I(i, j)$ and kernel (or filter) functions $K(i, j)$, for which a convolution can be written as

$$(I * K)(i, j) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} I(m, n)K(i - m, j - n) = (K * I)(i, j) = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} I(i - m, j - n)K(m, n)$$

, with $K(i, j), I(i, j) = 0$ if i, j outside the image or kernel

CONVOLUTION AND CROSS CORRELATION



1 - THE CORE OF CNNs: IMAGE CONVO

1	0	1
0	1	0
1	0	1

Kernel

No interaction between all neurons

Weights in kernel stay constant as we move across the image

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

1 - THE CORE OF CNNs: IMAGE CONVO

1	0	1
0	1	0
1	0	1

Kernel

No interaction between all neurons

Weights in kernel stay constant as we move across the image

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

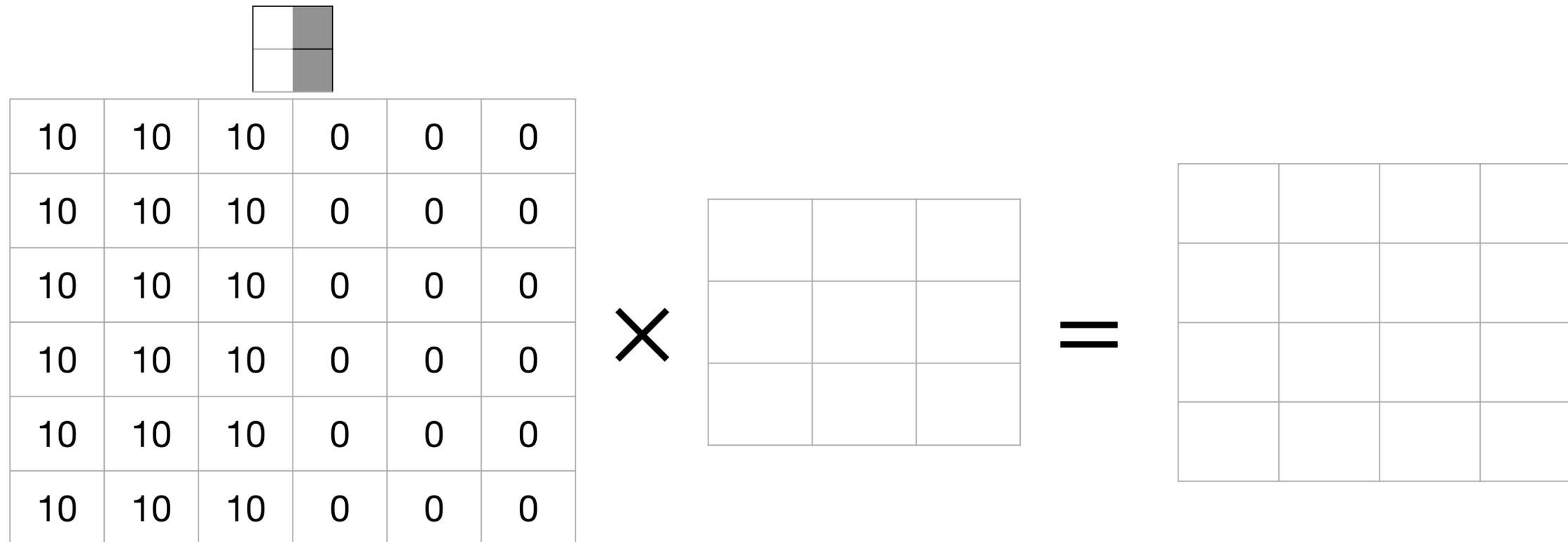
Image

4		

Convolved Feature

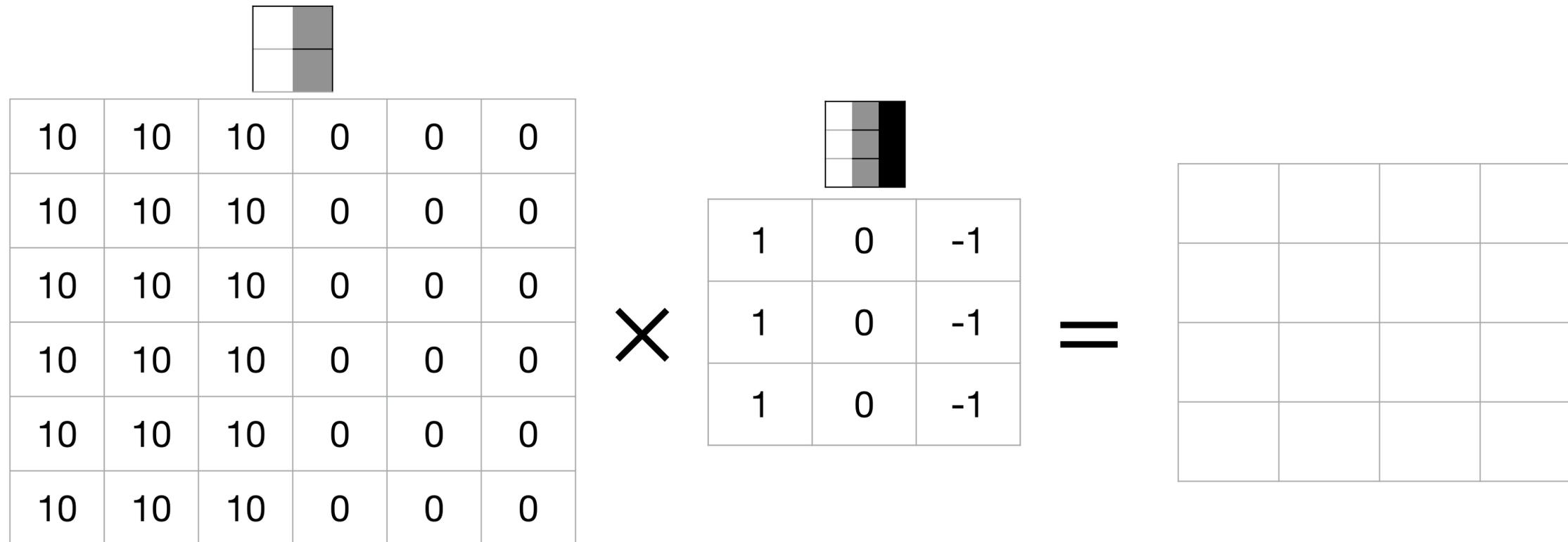
1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Vertical edge detection



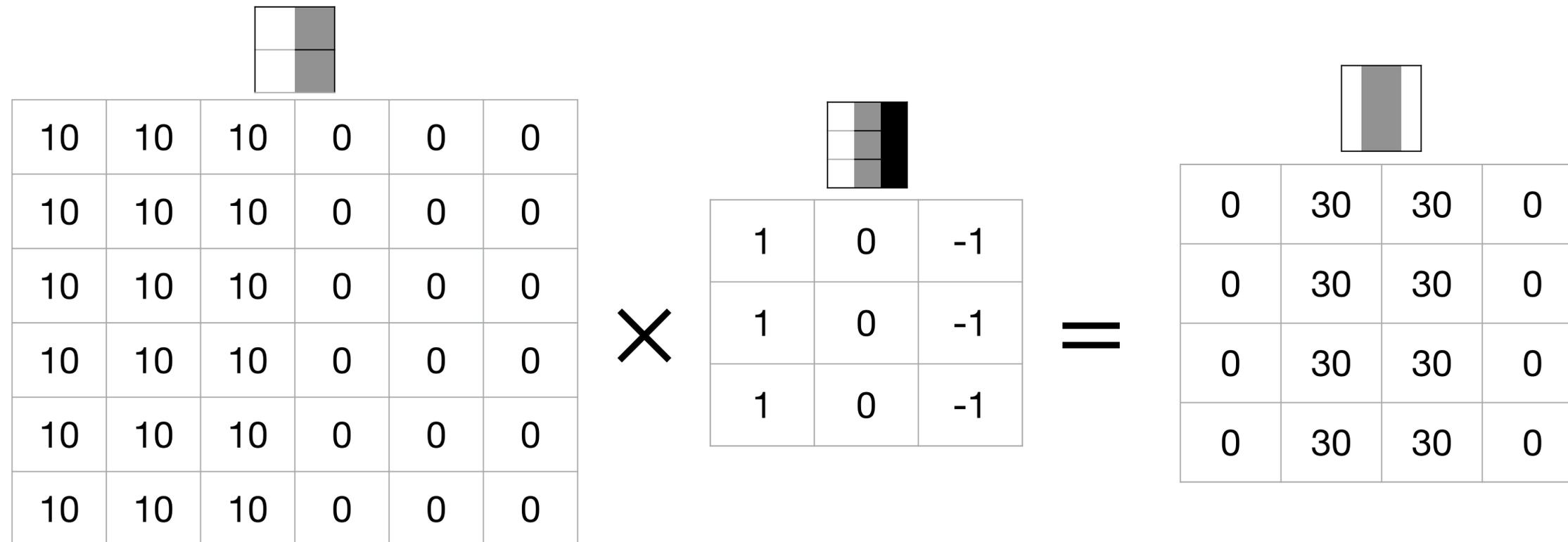
1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Vertical edge detection



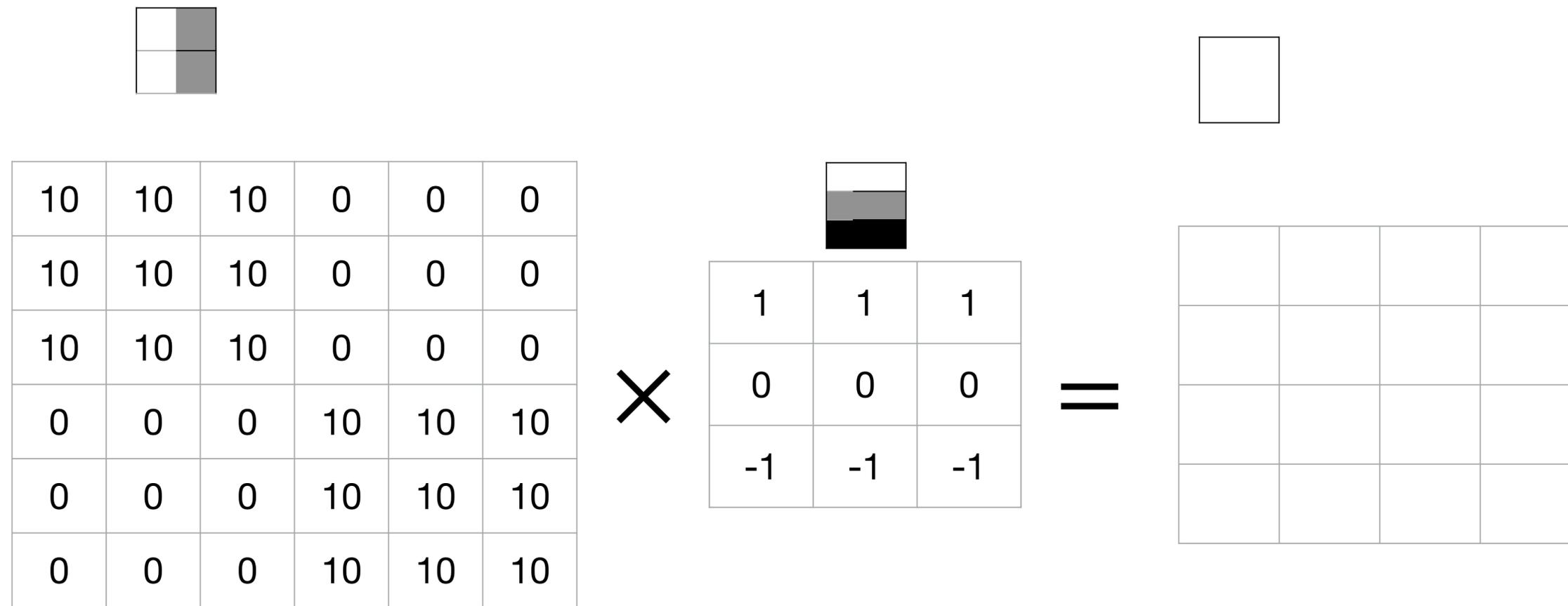
1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Vertical edge detection



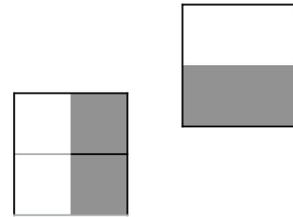
1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Horizontal edge detection



1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Horizontal edge detection



10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

×

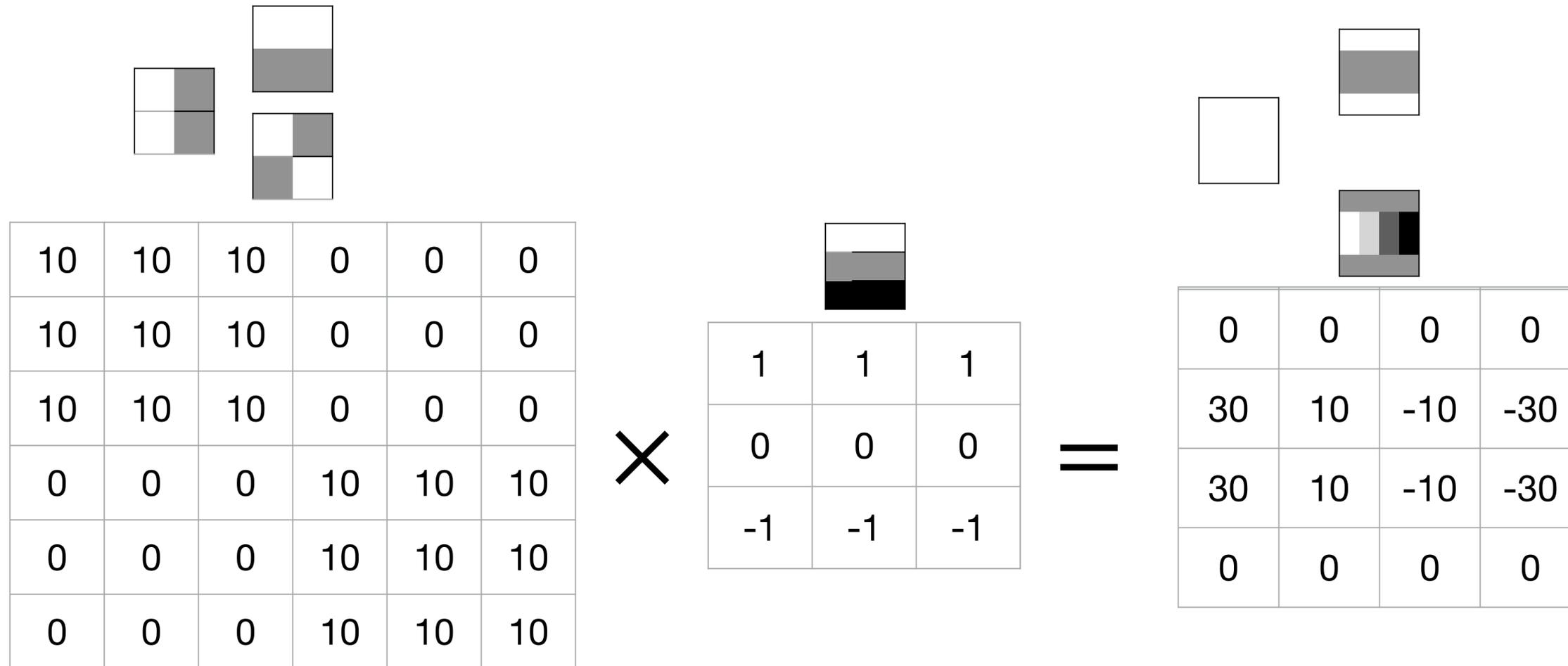
1	1	1
0	0	0
-1	-1	-1

=

0	0	0	0
30	30	30	30
30	30	30	30
0	0	0	0

1 - THE CORE OF CNNs: EXAMPLE CONVOLUTION FOR EDGE DETECTION

Horizontal edge detection



1 – THE CORE OF CNNs: EDGE DETECTION AT WORK



1 - THE CORE OF CNNs: WHY ARE CNNs SO USEFUL?

- Reduce number of computations (sparse interactions, parameter sharing) and memory requirements
- Intuitively: CNN will learn filters that activate for some kind of feature in the image like an edge or blotch of some color; often useful to learn that regardless of position in image **translational equivariance**:
- Equivariance not always useful: e.g. when different features are present at different positions in an image, like a face in the upper part. Then don't use weight sharing

CONVOLUTIONAL NEURAL NETWORKS

Lecture 4

80 Fukushima

89 LeCun et al

12 Krizhevsky et al

KERNEL, CHANNELS, POOLING, ...

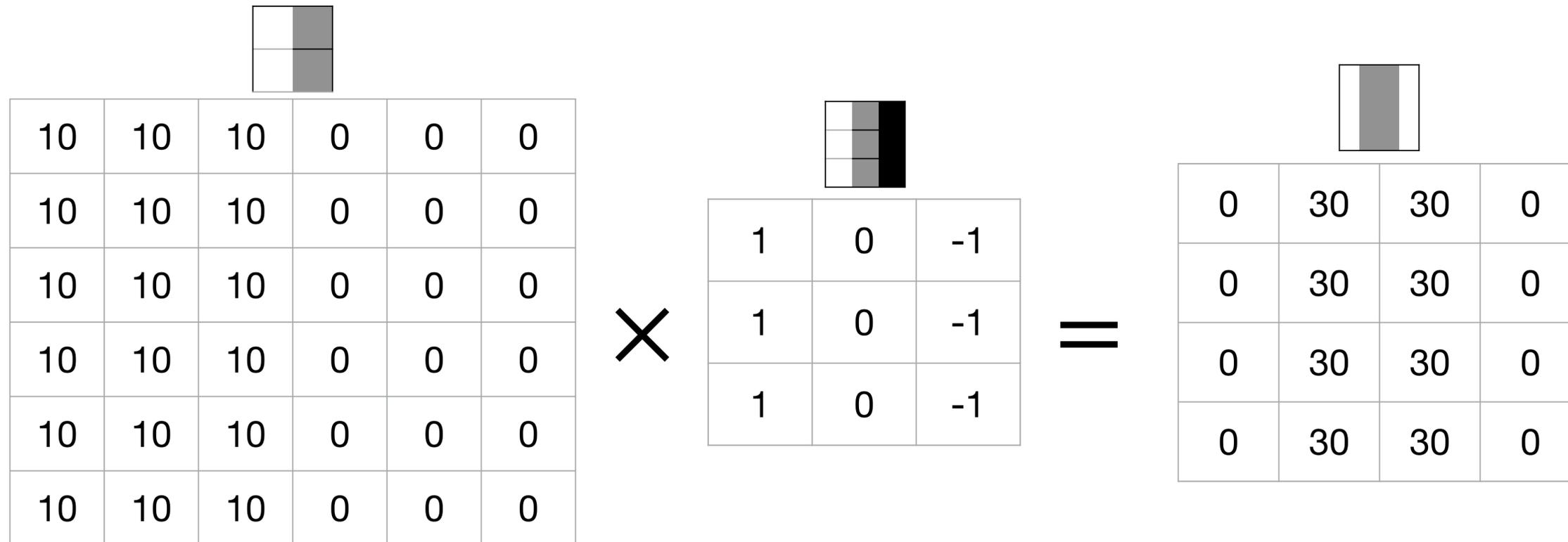
The diagram features a green background with a blue grid representing a kernel being applied to an input plane. The text is written in a hand-drawn, white font. The statistics are presented in orange rounded rectangles.

https://florianmarquardt.github.io/deep_learning_basics_linkmap.html

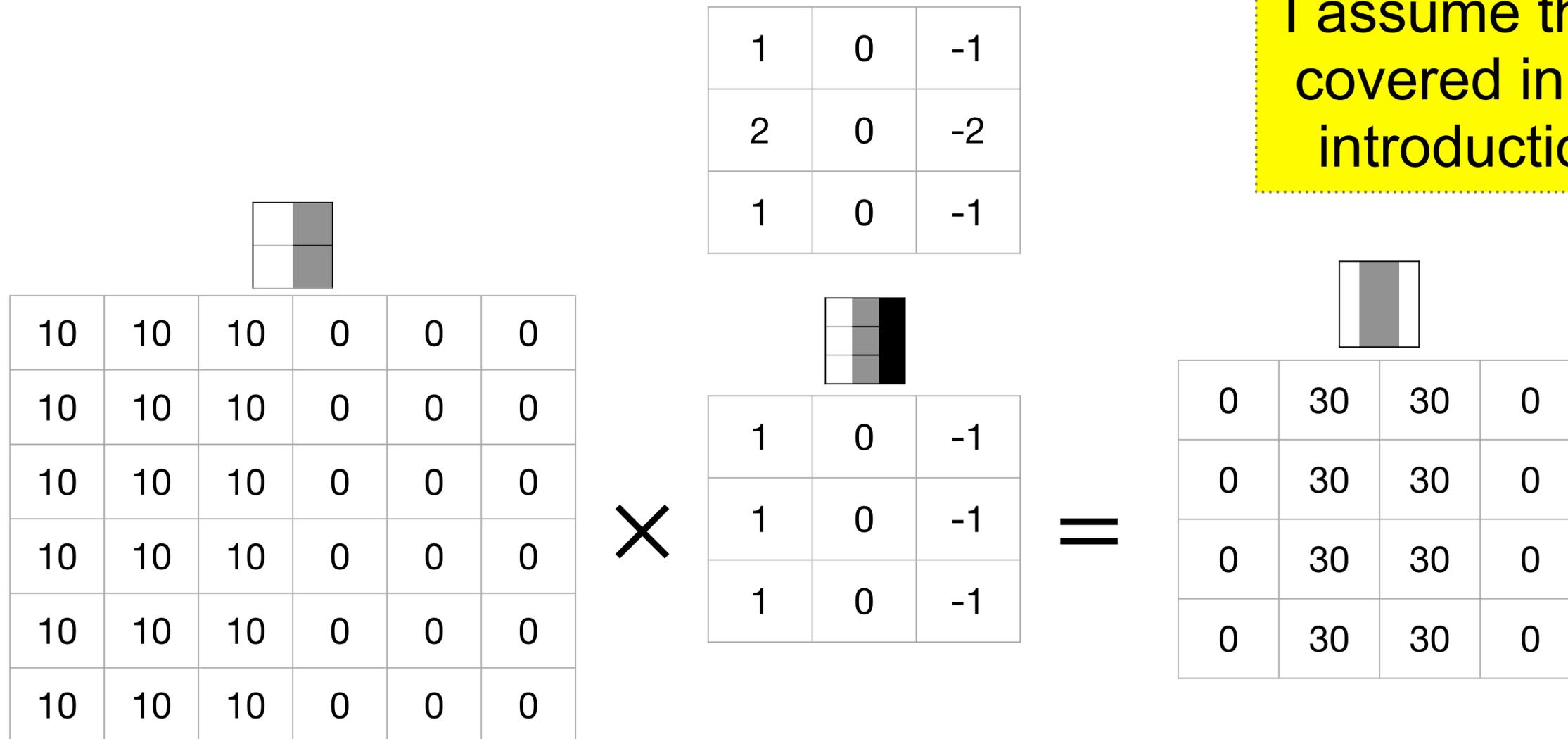
Now that we understand the basic mechanics of a CNN, how does the learning work?

2 UNDERSTAND LEARNING IN CNNs

I assume this has been covered in the general introduction to ANNs

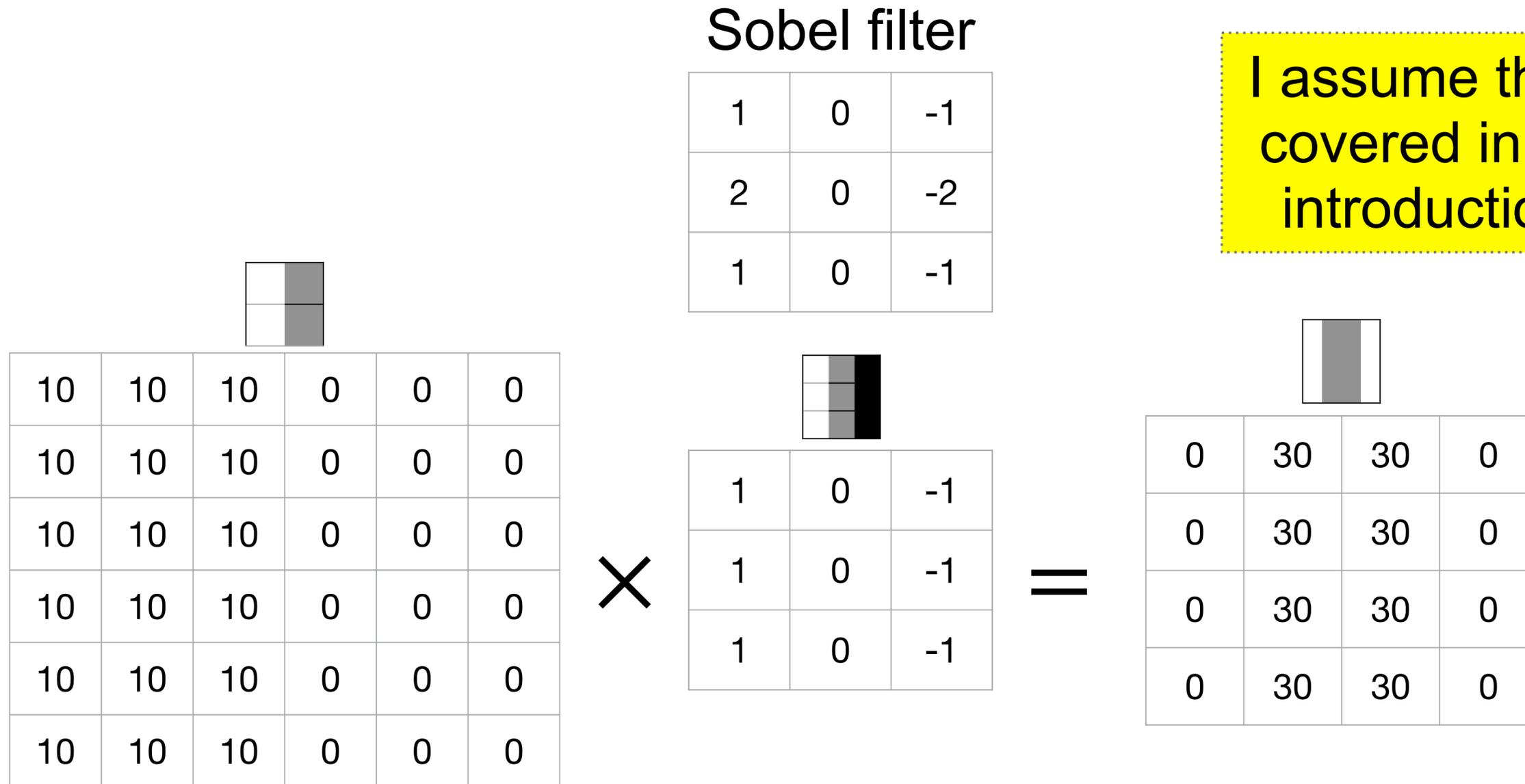


2 UNDERSTAND LEARNING IN CNNs

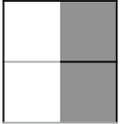


I assume this has been covered in the general introduction to ANNs

2 UNDERSTAND LEARNING IN CNNs



2 UNDERSTAND LEARNING IN CNNs

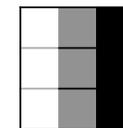


10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

×

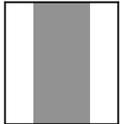
Sobel filter

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



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

=



0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

I assume this has been covered in the general introduction to ANNs

3	0	-3
10	0	-10
3	0	-3

2 UNDERSTAND LEARNING IN CNNs

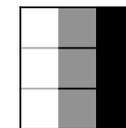


10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

×

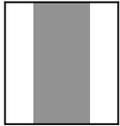
Sobel filter

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



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

=



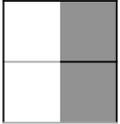
0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

I assume this has been covered in the general introduction to ANNs

Schorr filter

3	0	-3
10	0	-10
3	0	-3

2 UNDERSTAND LEARNING IN CNNs

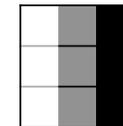


10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

×

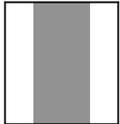
Sobel filter

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



W ₁	W ₂	W ₃
W ₄	W ₅	W ₆
W ₇	W ₈	W ₉

=



0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

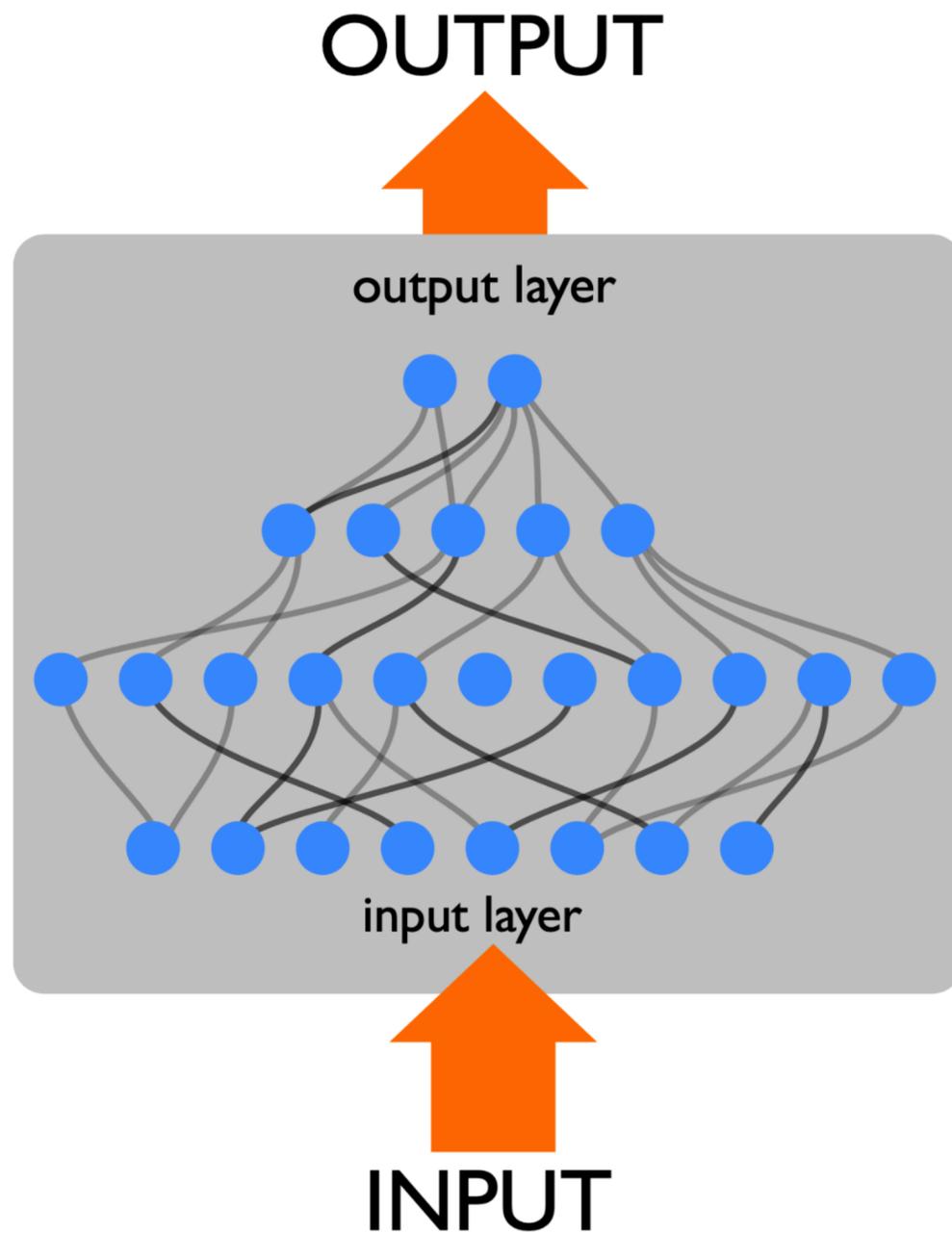
I assume this has been covered in the general introduction to ANNs

Schorr filter

3	0	-3
10	0	-10
3	0	-3

2 UNDERSTAND LEARNING IN CNNs: REMEMBER – TRAINING OF WEIGHTS

A neural network



I assume this has been covered in the general introduction to ANNs

=

Complicated nonlinear function that depends on all the weights and biases

$$y^{\text{out}} = F_w(y^{\text{in}})$$



2 UNDERSTAND LEARNING IN CNNs: REMEMBER – TRAINING OF WEIGHTS

We have:

$$y^{\text{out}} = F_w(y^{\text{in}})$$

neural network
(w here also stands for the weights)

I assume this has been covered in the general introduction to ANNs

We would like: $y^{\text{out}} \approx F(y^{\text{in}})$

desired “target” function

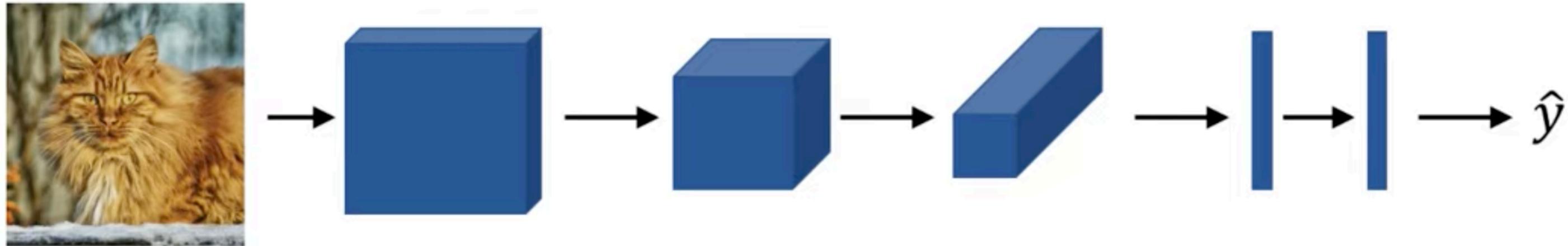
Cost function measures deviation:

$$C(w) = \frac{1}{2} \langle \underbrace{\| F_w(y^{\text{in}}) - F(y^{\text{in}}) \|^2}_{\text{vector norm}} \rangle$$

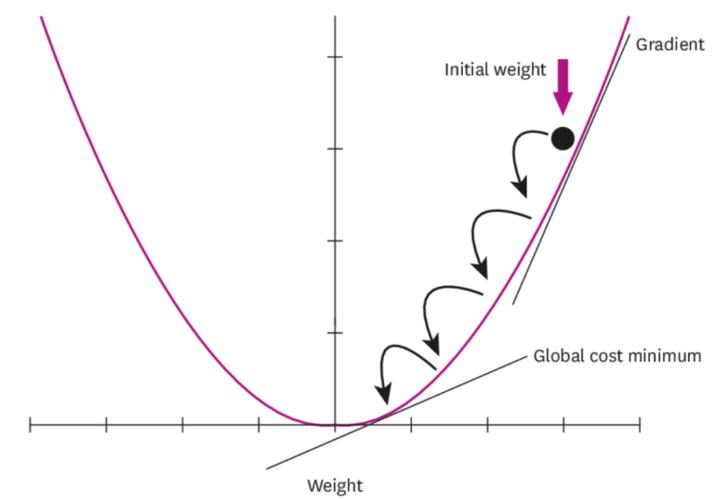
average over all samples

2 UNDERSTAND LEARNING IN CNNs

Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$.



- Define a cost function.
$$J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$
- Use gradient descent to derive the set of W, b , for which the cost is minimised



Before we can build our CNN, we need three more things

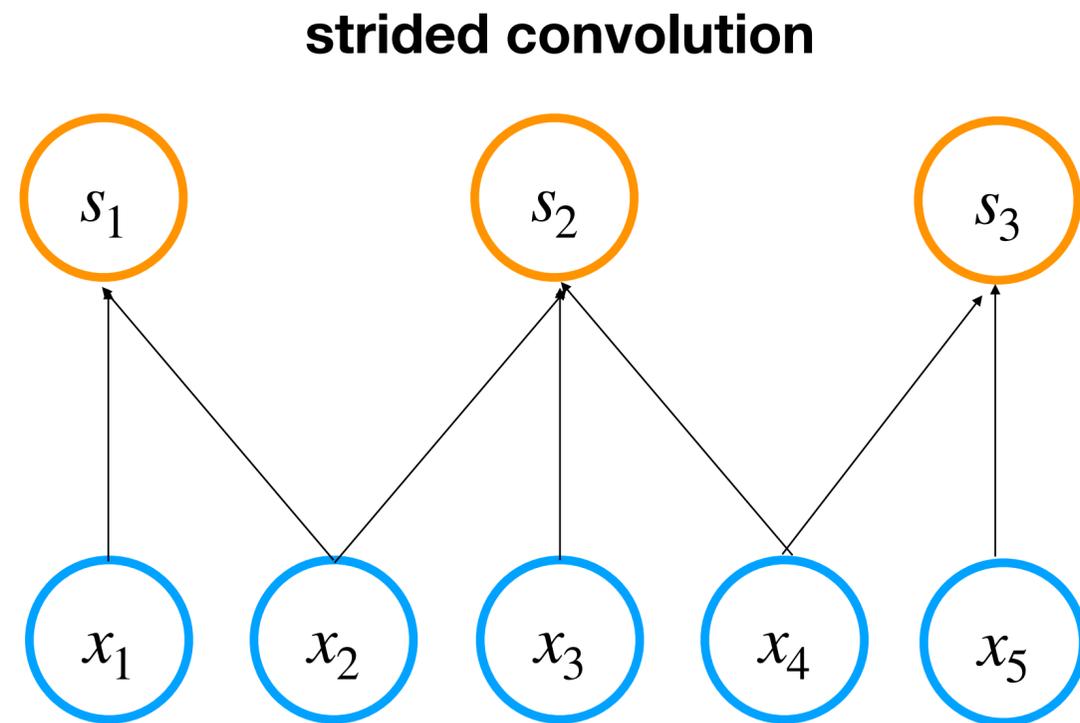
3 USEFUL TOOLS – STRIDING, PADDING, POOLING

- To reduce computational cost we can shift the kernel by more than one pixel, that is, choose a stride $s > 1$
- This is equivalent to a down sampled convolution:

$$z_{i,j,k}^{[l]} = c(K, a^{[l-1]}, s)_{i,j,k} = \sum_m \sum_n \sum_p a_{m,(j-1)\times s+n,(k-1)\times s+p}^{[l-1]} K_{i,m,n,p}^{[l]}$$

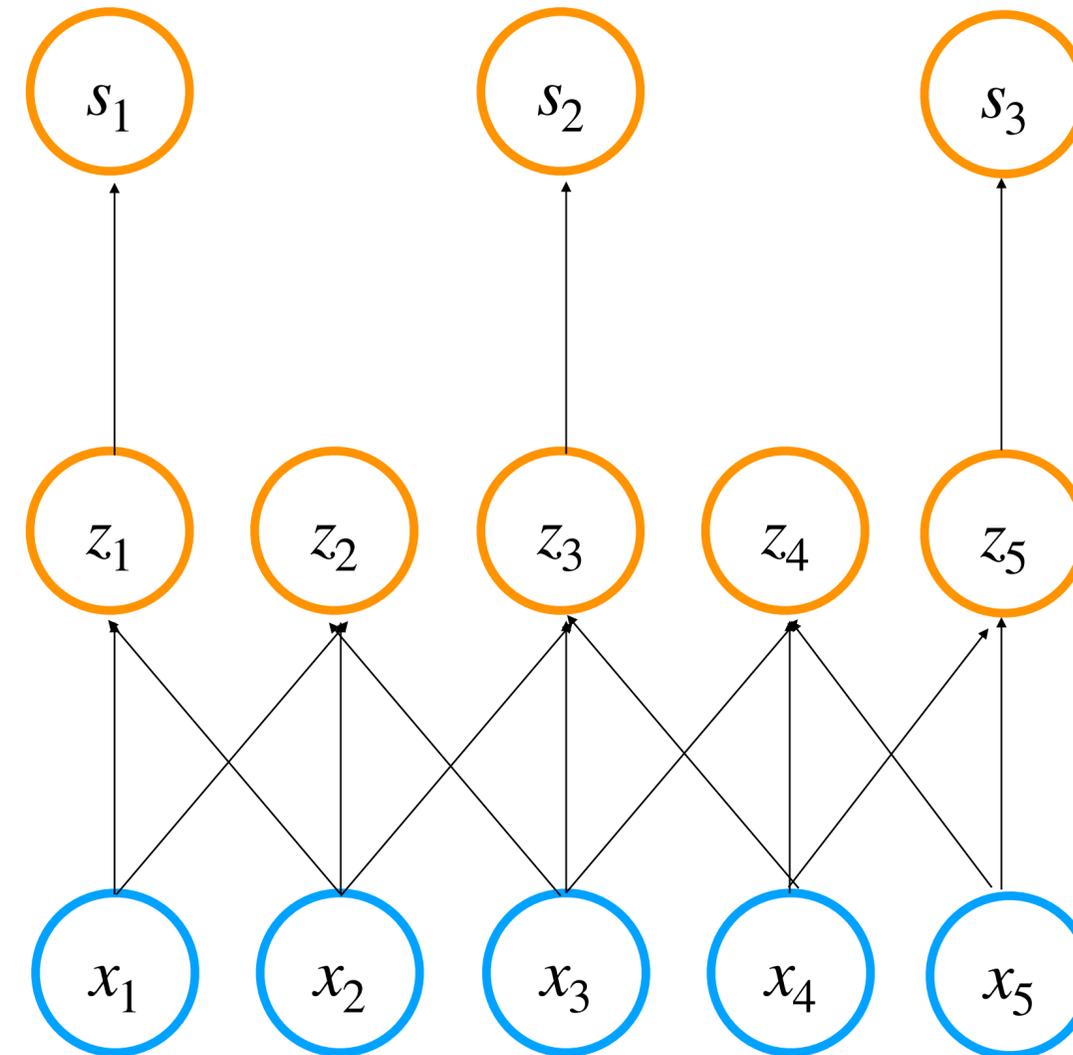
- Which comes at the expense that features can not be extracted as finely

3 USEFUL TOOLS – STRIDING, PADDING, POOLING



=

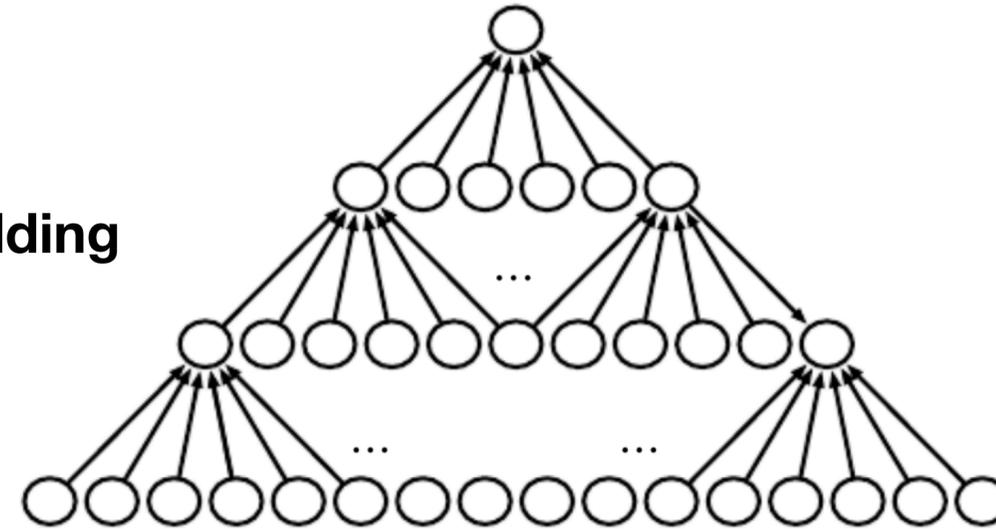
Convolution with downsampling



3 USEFUL TOOLS – STRIDING, PADDING, POOLING

- If only **valid convolutions** (=kernel fully contained in image) are considered, image will quickly shrink in size

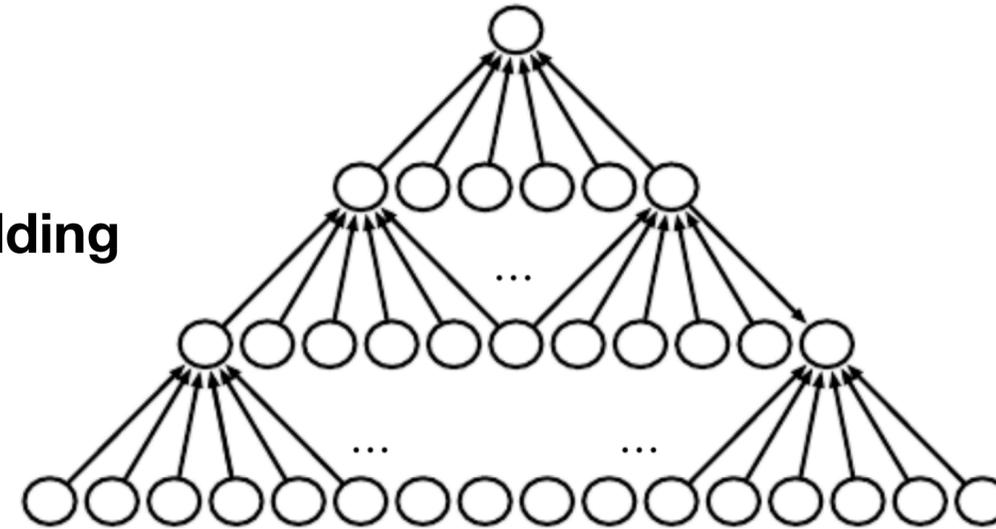
Valid convolution: no zero-padding



3 USEFUL TOOLS – STRIDING, PADDING, POOLING

- If only **valid convolutions** (=kernel fully contained in image) are considered, image will quickly shrink in size

Valid convolution: no zero-padding



Same convolution: zero padding
so that dimensions don't change



- Padding edges with p zeros, allows one to make arbitrarily deep CNNs
- Avoids some edge effects

3 USEFUL TOOLS – STRIDING, PADDING – OUTPUT SIZE OF CNN

- With input of width $W^{[l-1]}$, and kernel with receptive field of size $f^{[l]}$, stride $s^{[l]}$ and zero padding $p^{[l]}$, one output image in the convolutional layer will have width

$$W^{[l]} = (W^{[l-1]} - f^{[l]} + 2p^{[l]})/s^{[l]} + 1$$

- Strides constrained so that $W^{[l]}$ is an integer
- With $s^{[l]} = 1$, zero padding with $p^{[l]} = (f^{[l]} - 1)/2$ will give $W^{[l]} = W^{[l-1]}$

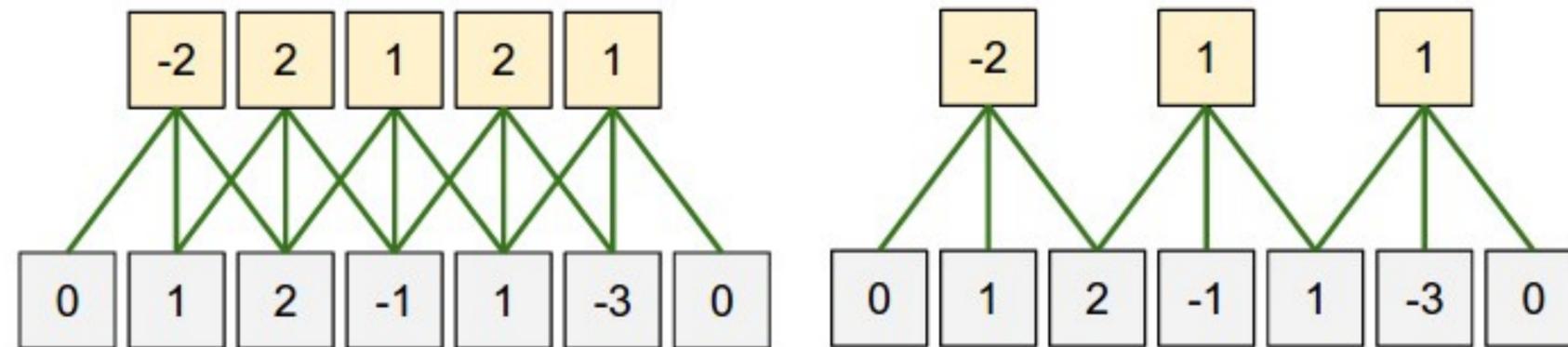
3 USEFUL TOOLS – STRIDING, PADDING, POOLING – OUTPUT SIZE

- With input of width w , and kernel with receptive field of size f , stride s and zero padding p , one output image in the convolutional layer will have width

$$w_2 = (w - f + 2p)/s + 1$$

- Strides constrained so that w_2 is an integer

- With $s = 1$, zero padding with $p = (f - 1)/2$ will give $w_2 = w$



$$W^{[1]} = 7$$

$$f^{[2]} = 3$$

$$p^{[2]} = 0$$

$$s^{[2]} = 1$$

$$\Rightarrow W^{[2]} = \frac{7 - 3 + 2 \cdot 0}{1} + 1 = 5$$

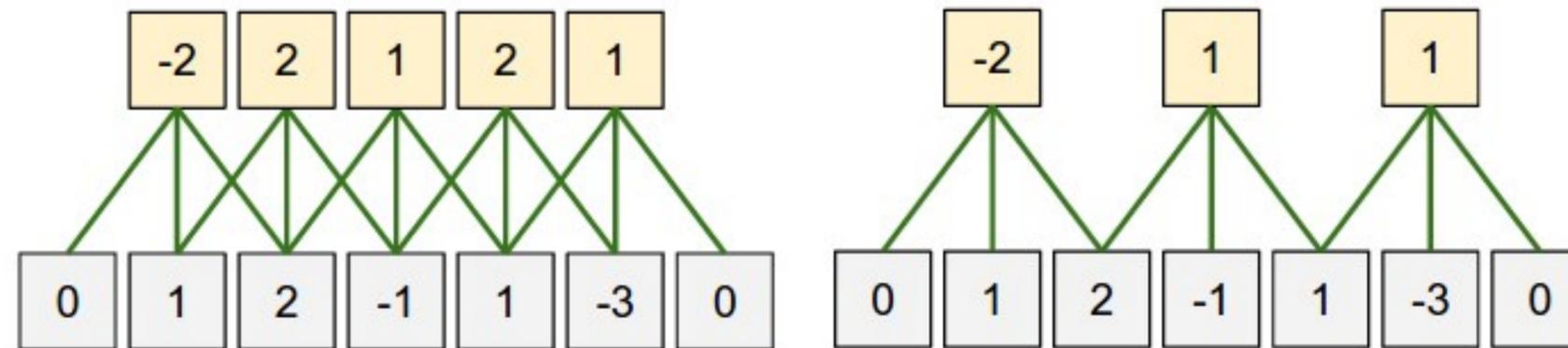
3 USEFUL TOOLS – STRIDING, PADDING, POOLING – OUTPUT SIZE

- With input of width w , and kernel with receptive field of size f , stride s and zero padding p , one output image in the convolutional layer will have width

$$w_2 = (w - f + 2p)/s + 1$$

- Strides constrained so that w_2 is an integer

- With $s = 1$, zero padding with $p = (f - 1)/2$ will give $w_2 = w$



$$W^{[1]} = 7$$

$$f^{[2]} = 3$$

$$p^{[2]} = 0$$

$$s^{[2]} = 1$$

$$\Rightarrow W^{[2]} = \frac{7 - 3 + 2 \cdot 0}{1} + 1 = 5$$

$$W^{[1]} =$$

$$f^{[2]} =$$

$$p^{[2]} =$$

$$s^{[2]} =$$

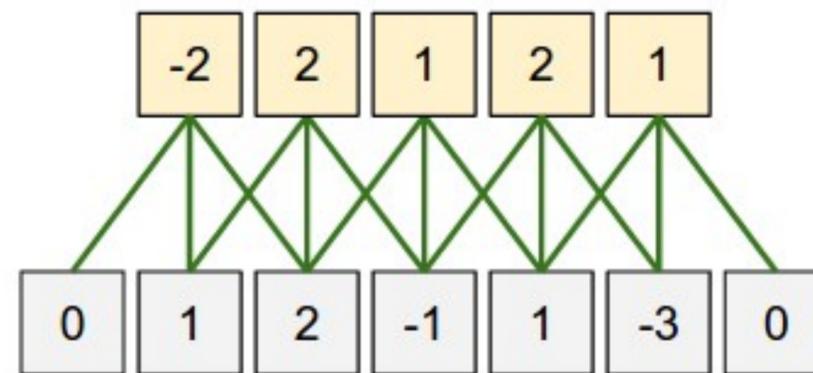
3 USEFUL TOOLS – STRIDING, PADDING, POOLING – OUTPUT SIZE

- With input of width w , and kernel with receptive field of size f , stride s and zero padding p , one output image in the convolutional layer will have width

$$w_2 = (w - f + 2p)/s + 1$$

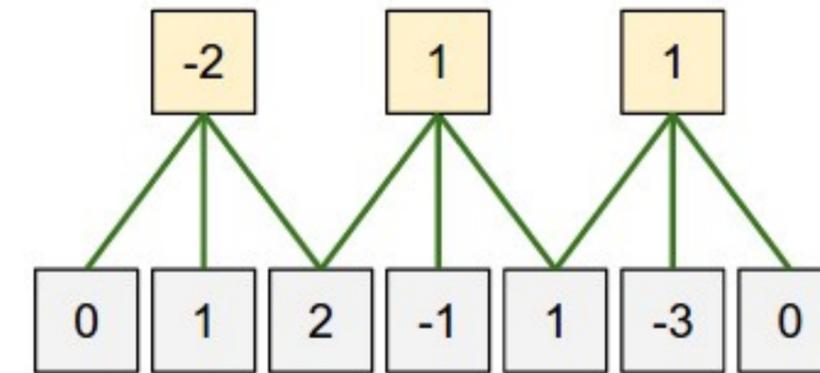
- Strides constrained so that w_2 is an integer

- With $s = 1$, zero padding with $p = (f - 1)/2$ will give $w_2 = w$



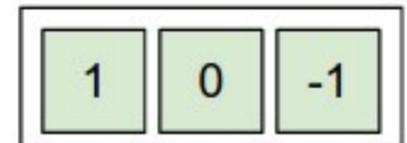
$$\begin{aligned} W^{[1]} &= 7 \\ f^{[2]} &= 3 \\ p^{[2]} &= 0 \\ s^{[2]} &= 1 \end{aligned}$$

$$\Rightarrow W^{[2]} = \frac{7 - 3 + 2 \cdot 0}{1} + 1 = 5$$



$$\begin{aligned} W^{[1]} &= 7 \\ f^{[2]} &= 3 \\ p^{[2]} &= 0 \\ s^{[2]} &= 2 \end{aligned}$$

Weights:



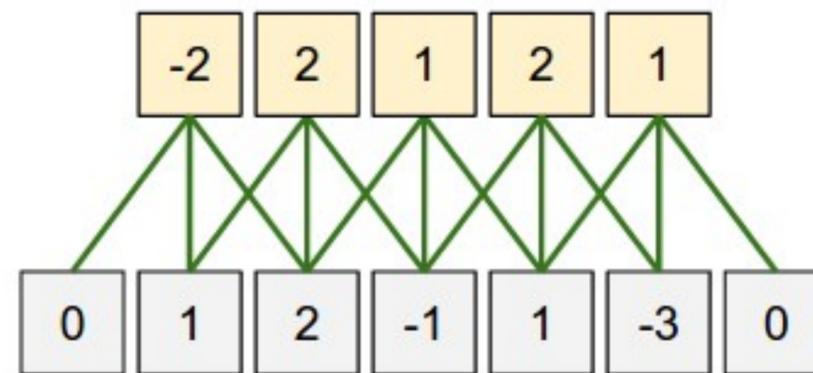
3 USEFUL TOOLS – STRIDING, PADDING, POOLING – OUTPUT SIZE

- With input of width w , and kernel with receptive field of size f , stride s and zero padding p , one output image in the convolutional layer will have width

$$w_2 = (w - f + 2p)/s + 1$$

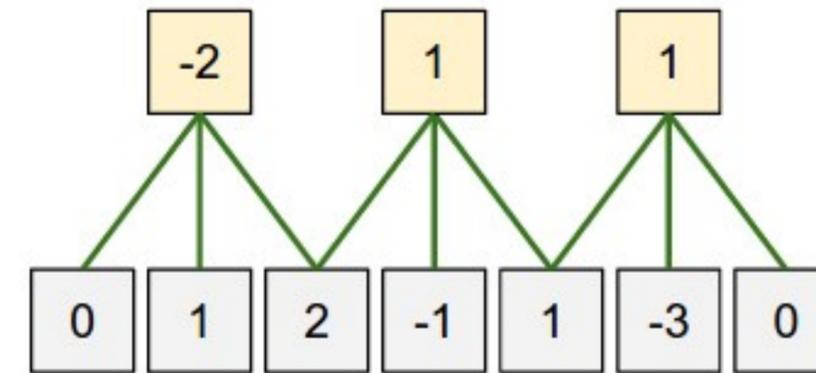
- Strides constrained so that w_2 is an integer

- With $s = 1$, zero padding with $p = (f - 1)/2$ will give $w_2 = w$



$$\begin{aligned} W^{[1]} &= 7 \\ f^{[2]} &= 3 \\ p^{[2]} &= 0 \\ s^{[2]} &= 1 \end{aligned}$$

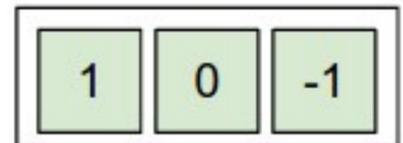
$$\Rightarrow W^{[2]} = \frac{7 - 3 + 2 \cdot 0}{1} + 1 = 5$$



$$\begin{aligned} W^{[1]} &= 7 \\ f^{[2]} &= 3 \\ p^{[2]} &= 0 \\ s^{[2]} &= 2 \end{aligned}$$

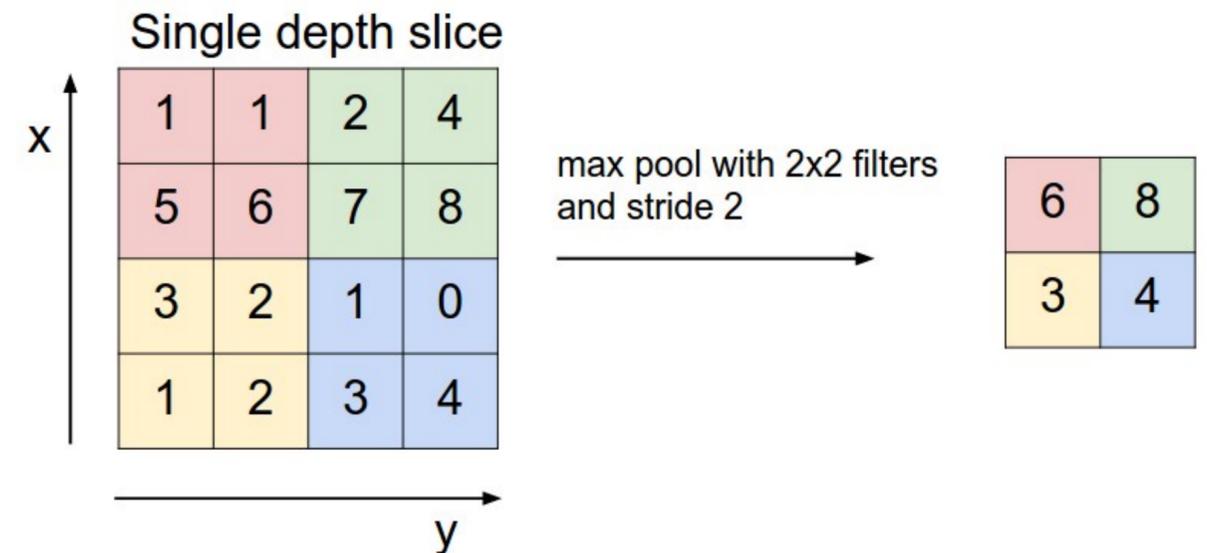
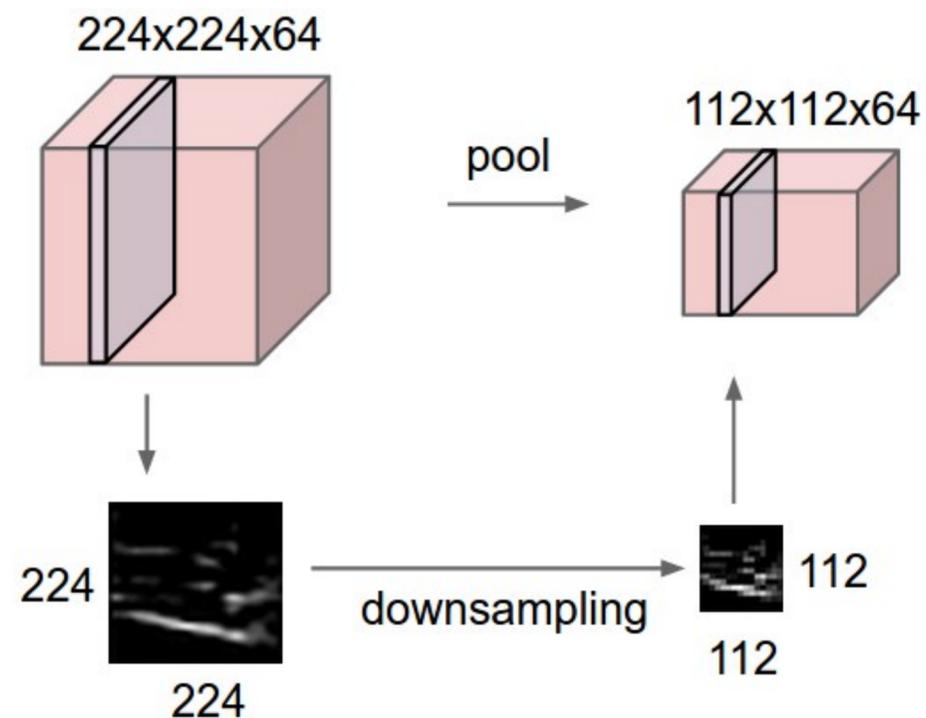
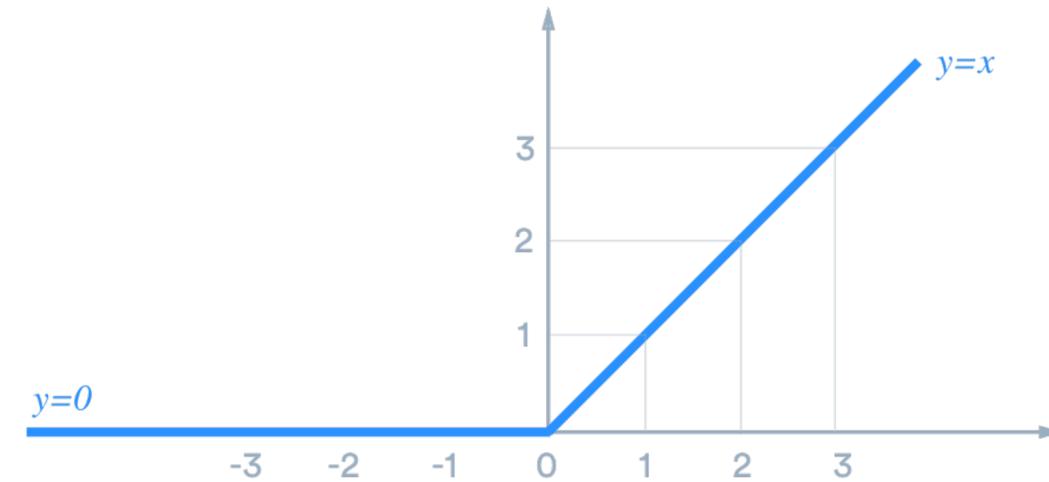
$$\Rightarrow W^{[2]} = \frac{7 - 3 + 2 \cdot 0}{2} + 1 = 3$$

Weights:



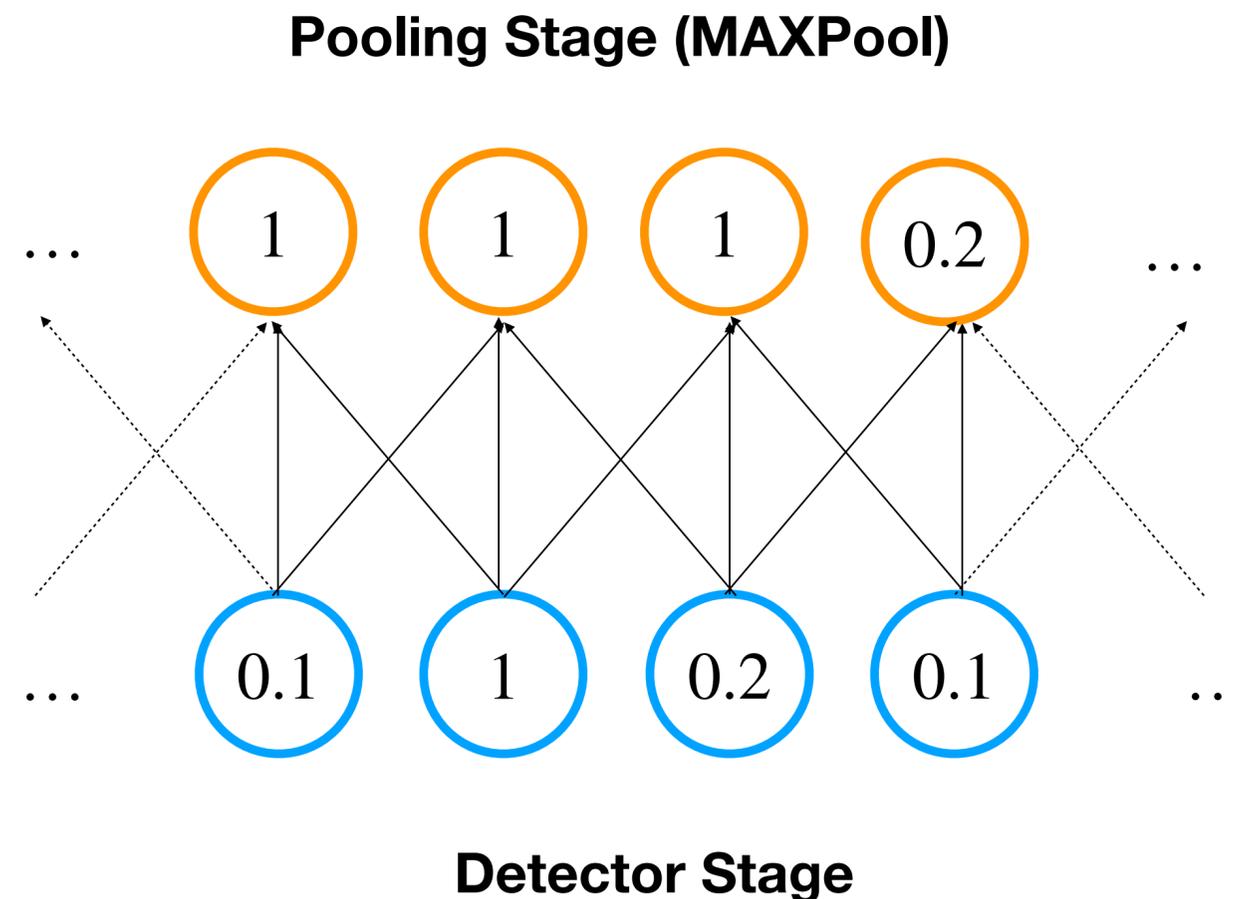
3 USEFUL TOOLS – STRIDING, PADDING, POOLING

- Result of convolution passed to non-linear activation function (e.g. Rectified Linear Unit RELU)
- Result of non-linear activation function usually passed to a ‘pooling layer’ that reduces spatial size (like a downsampling)
- Most common: 2x2 filter with stride 2 that selects the maximum of the input fields (MAX pool)



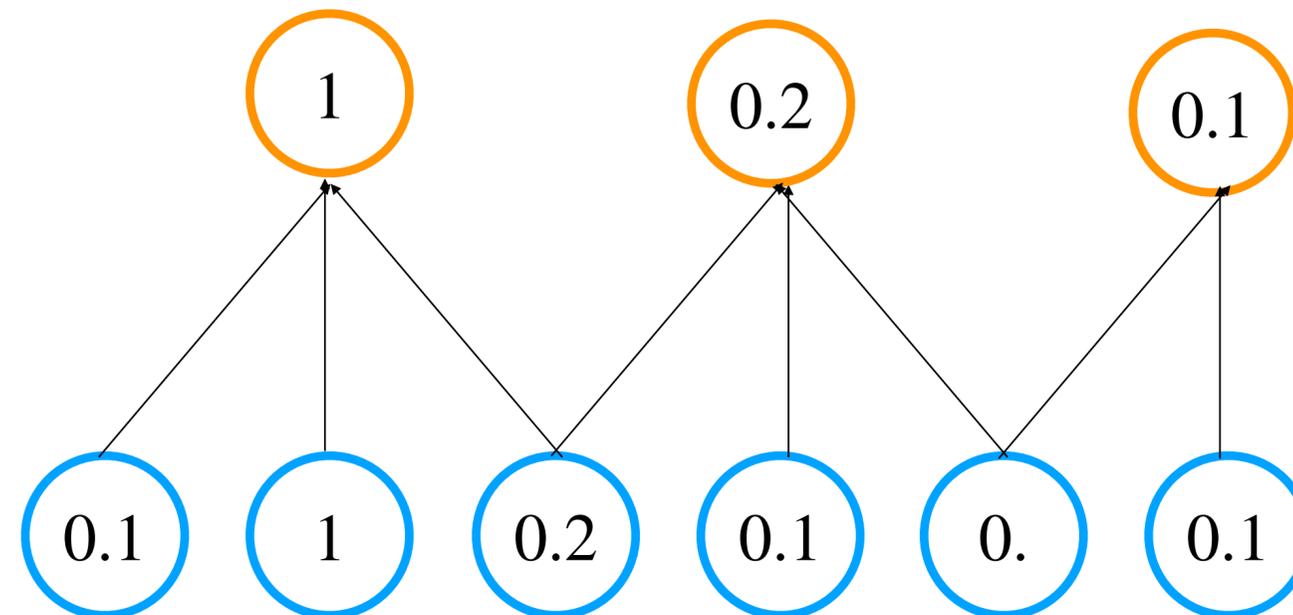
POOLING

- Helps to make representation **invariant** against small translations
- Invariance to local translations useful if we care more about a feature itself than its position in the image

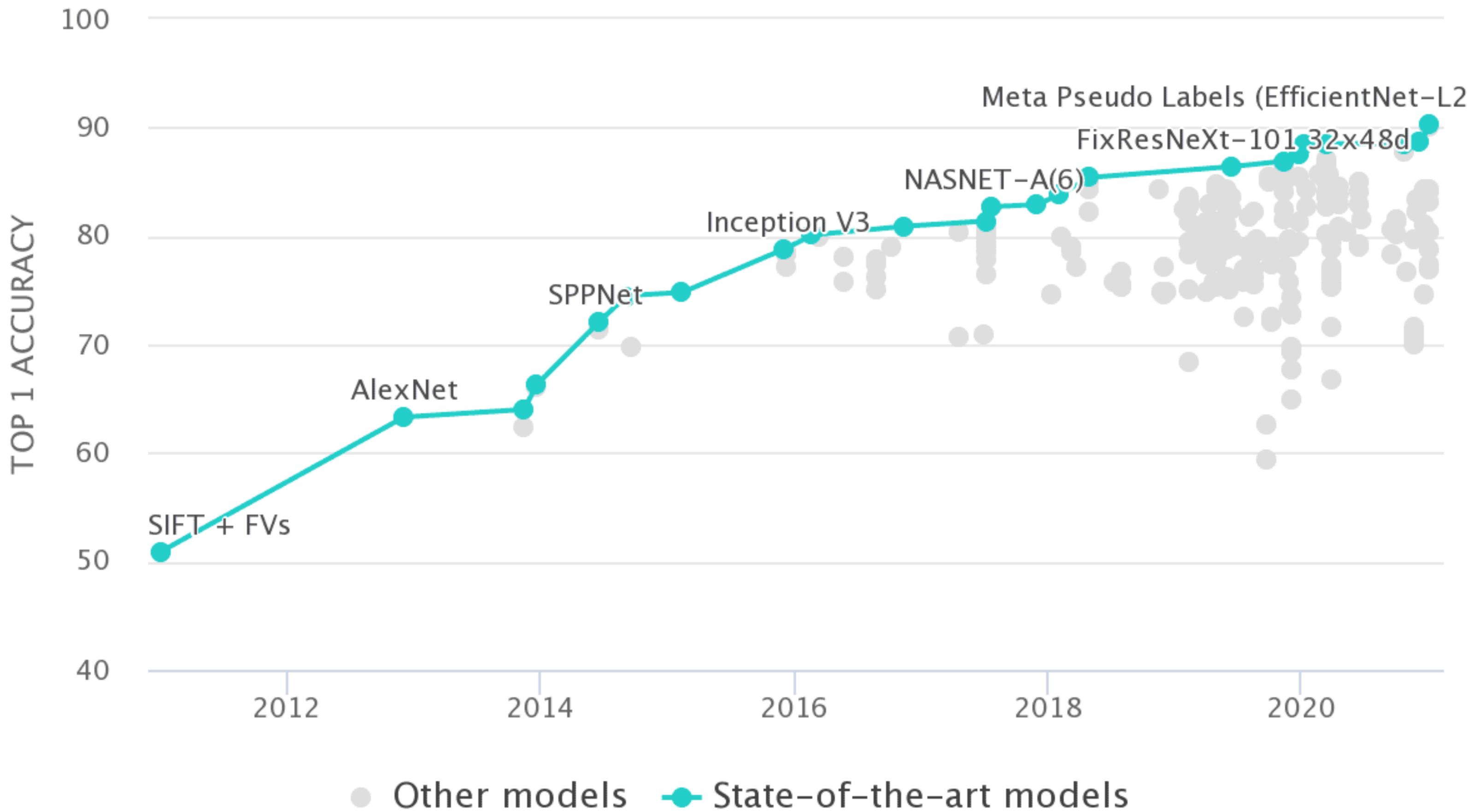


POOLING WITH DOWNSAMPLING

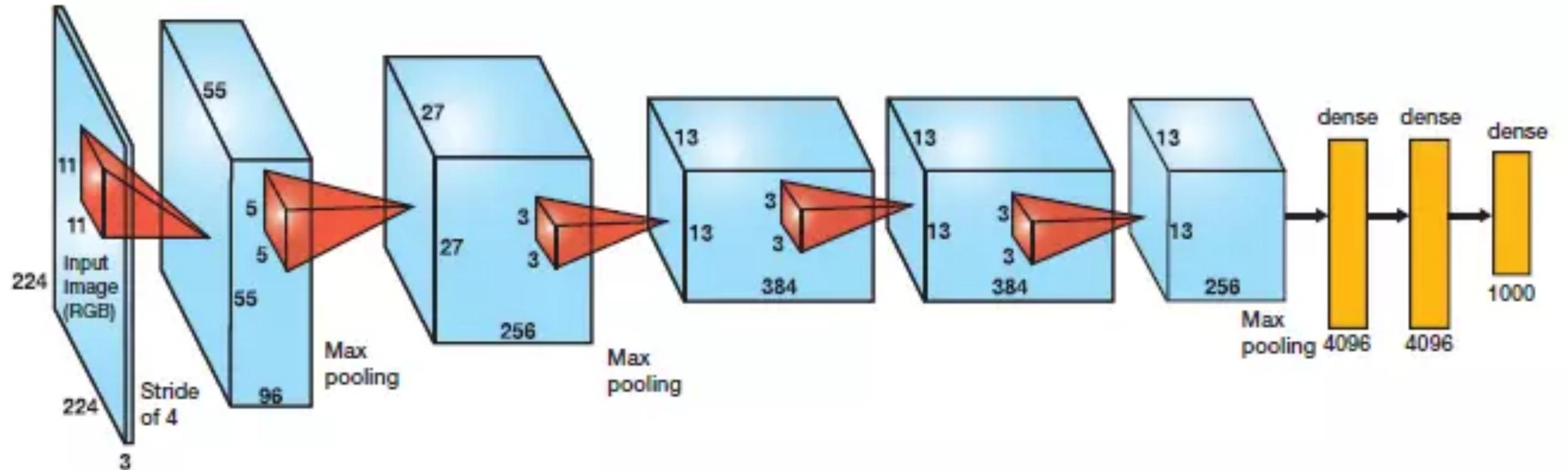
- Pooling summarizes activations of whole neighborhood
- Thus, makes sense to use fewer pooling units than detector units
- Example with stride $s = 2$:



OK, so we have the building blocks. Let's build a CNN.

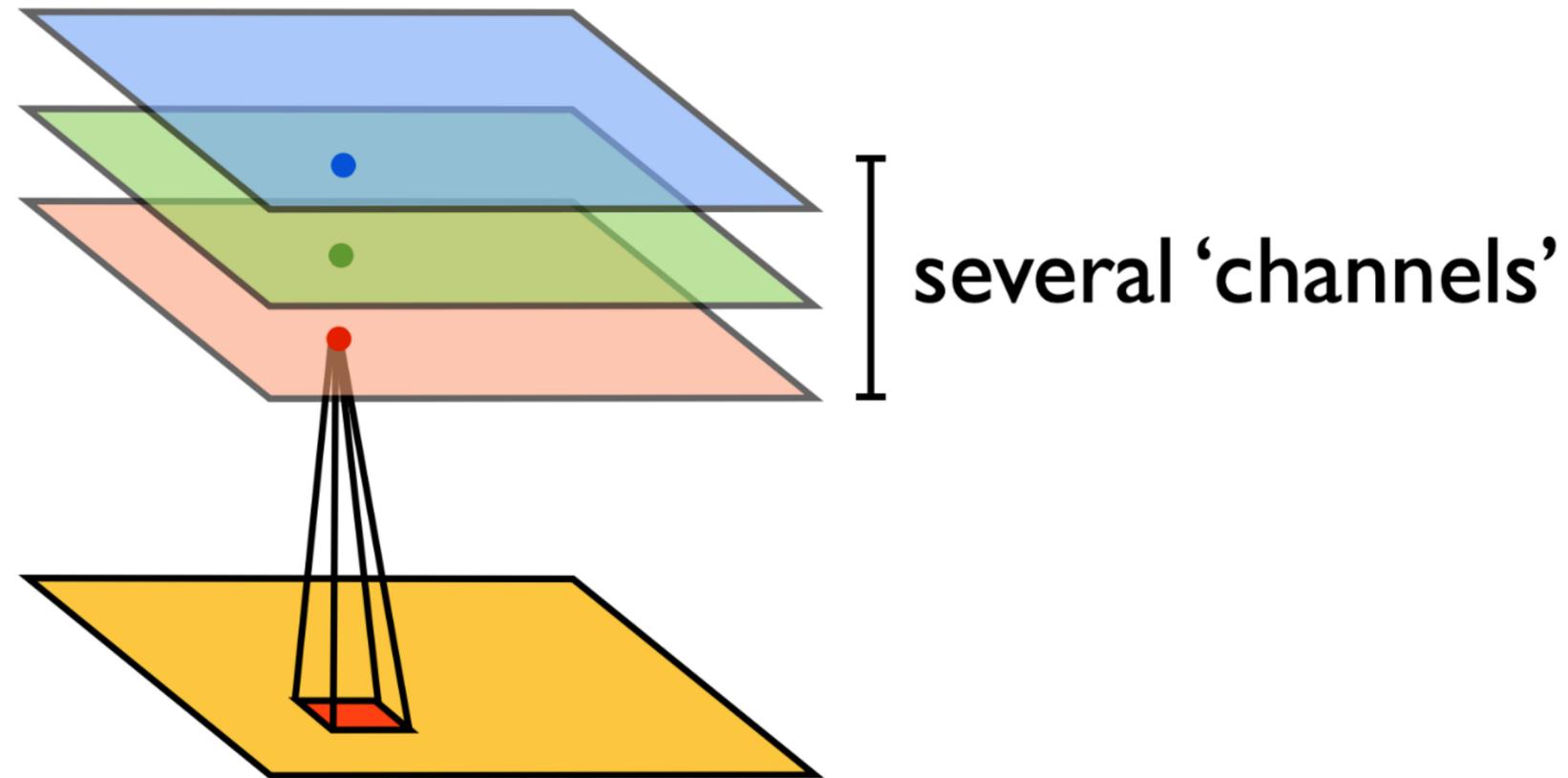


4 TYPICAL CNN ARCHITECTURES: THE FIRST SUCCESSFUL CNN – ALEX NET



Several filters (kernels)

e.g. one for smoothing, one for contours, etc.



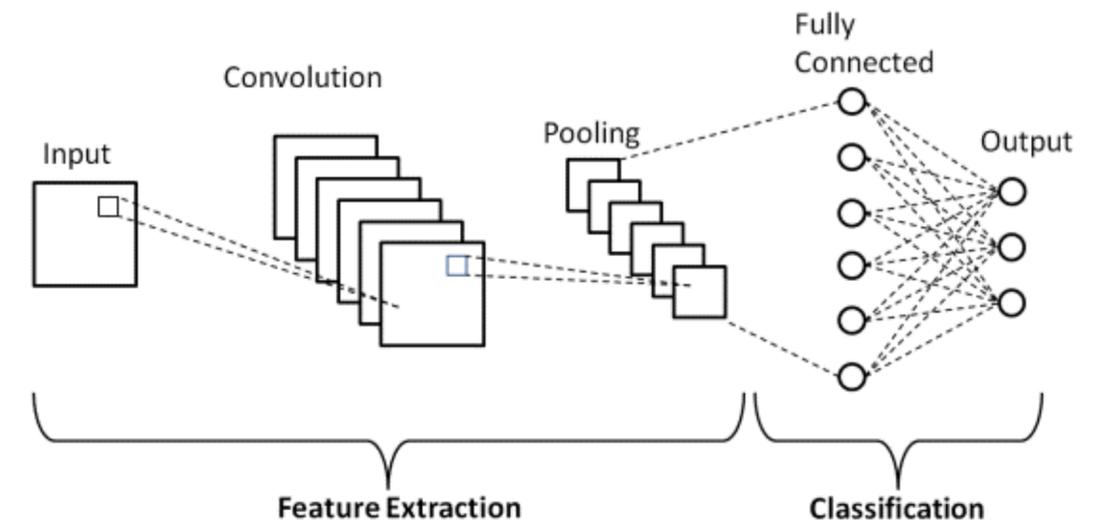
4 TYPICAL CNN ARCHITECTURES: THE FIRST SUCCESSFUL CNN – ALEX NET

- developed by Alex Krizhevsky, Ilya Sutskever and Geoff Hinton. First to popularize CNNs for computer vision, won 2012 ImageNet contest. Significantly outperformed runner up
- First to implement maxpool layers, ReLU activation and dropout layers
- nowadays can be implemented in 35 lines of Torch code
- How was the exact configuration chosen? Trial and error

```
class AlexNet(nn.Module):
    def __init__(self, num_classes: int = 1000) -> None:
        super(AlexNet, self).__init__()
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )
        self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
        self.classifier = nn.Sequential(
            nn.Dropout(),
            nn.Linear(256 * 6 * 6, 4096),
            nn.ReLU(inplace=True),
            nn.Dropout(),
            nn.Linear(4096, 4096),
            nn.ReLU(inplace=True),
            nn.Linear(4096, num_classes),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = self.features(x)
        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.classifier(x)
        return x
```

4 TYPICAL CNN ARCHITECTURES



➤ Typical architecture of CNN generally will look like:

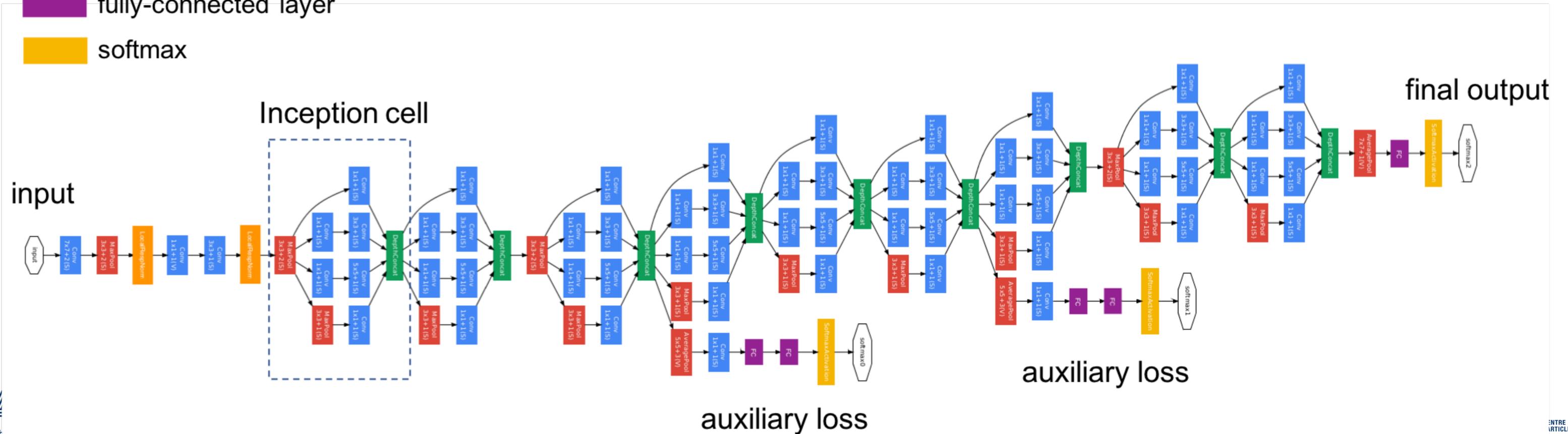


- In general: prefer repeated convolutions with small kernels over one convolution with large kernel
- Use zero padding such that convolution does not change spatial dimensions, only use pooling for that. Otherwise, you have to make sure that striding works
- Fully connected layers to classify based on the high-level features learned in the convolution part

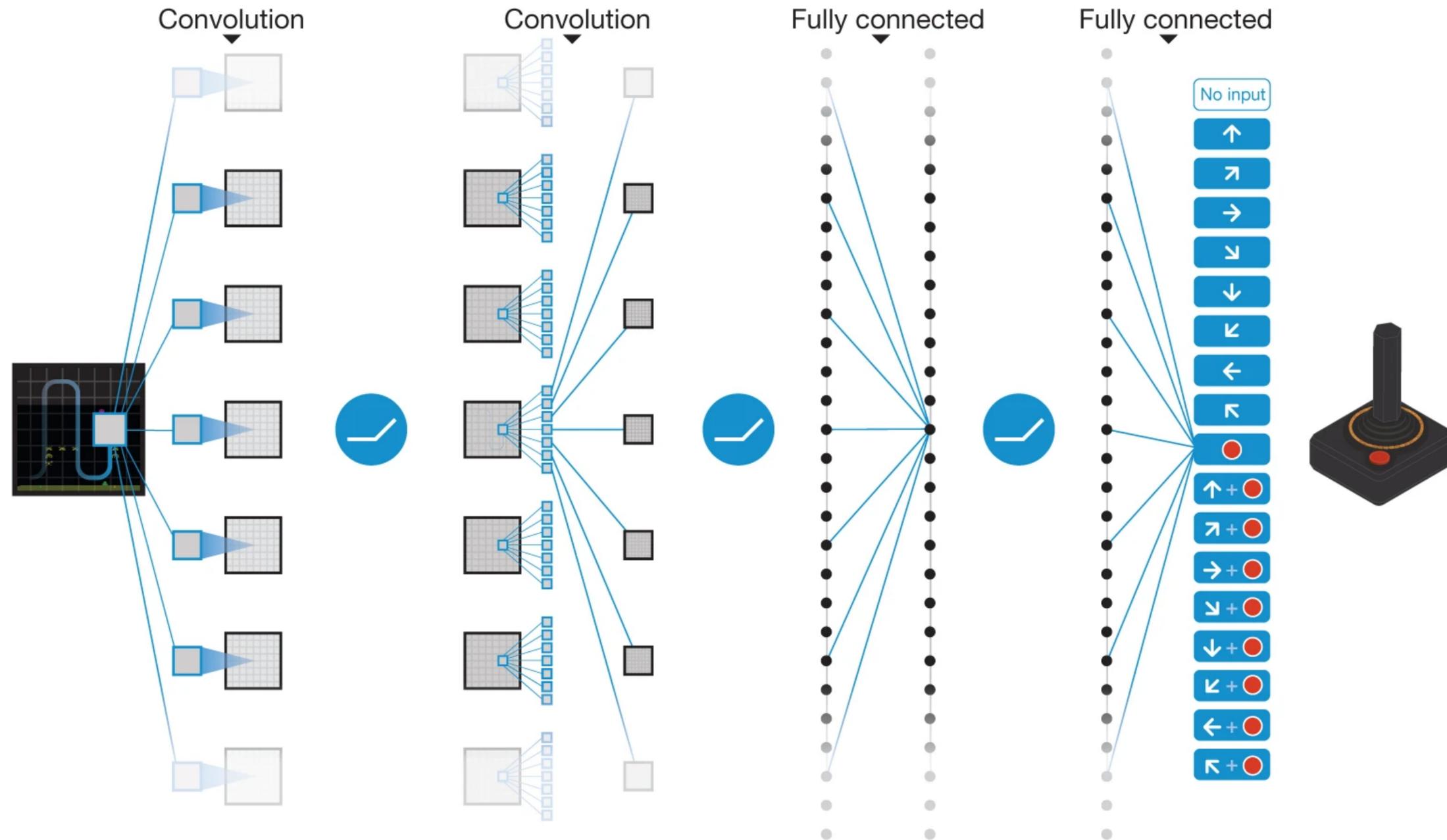
4 TYPICAL CNN ARCHITECTURES: DEVELOPMENT FOLLOWING ALEXNET

► Deeper, e.g. Inception (GoogLeNet)

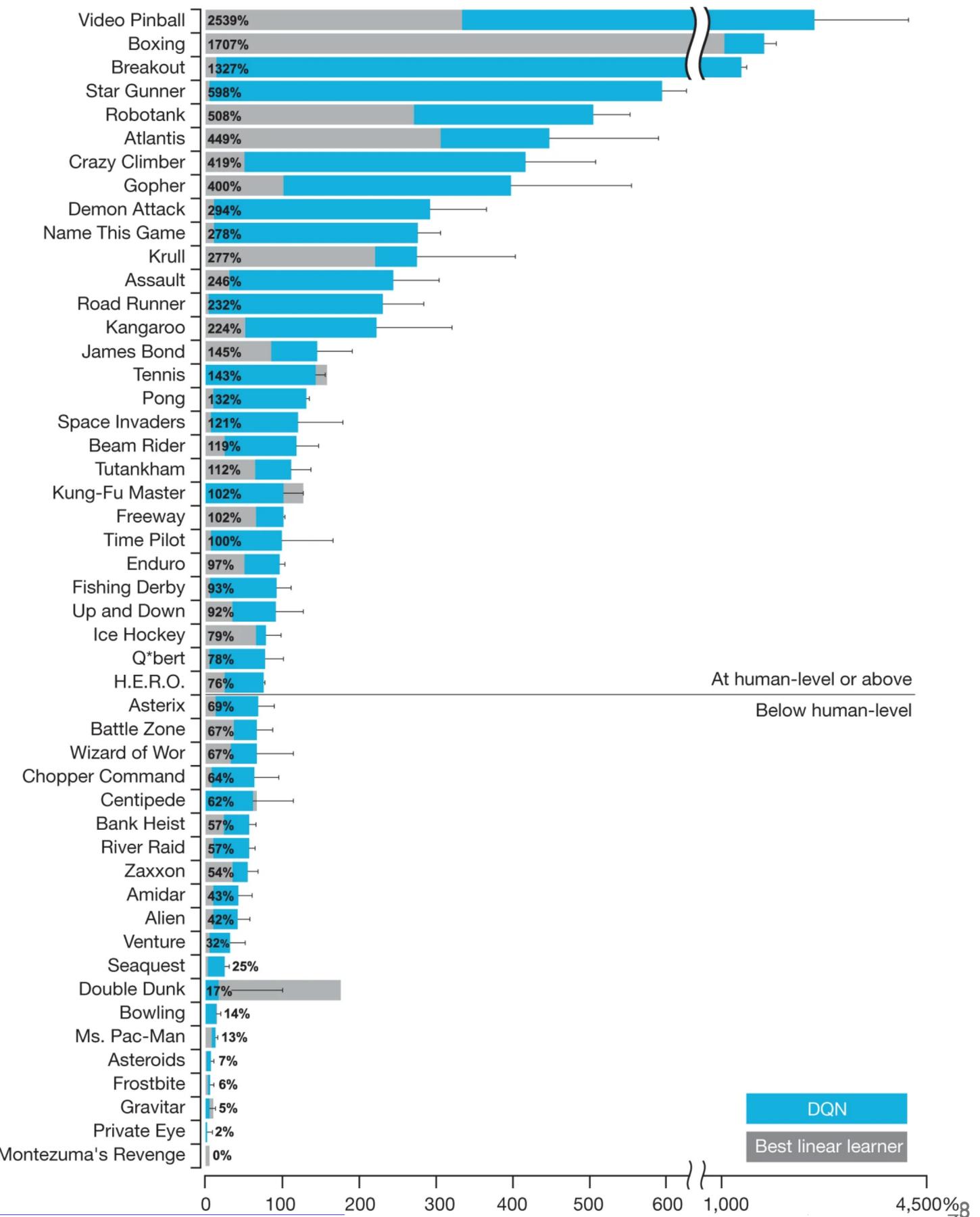
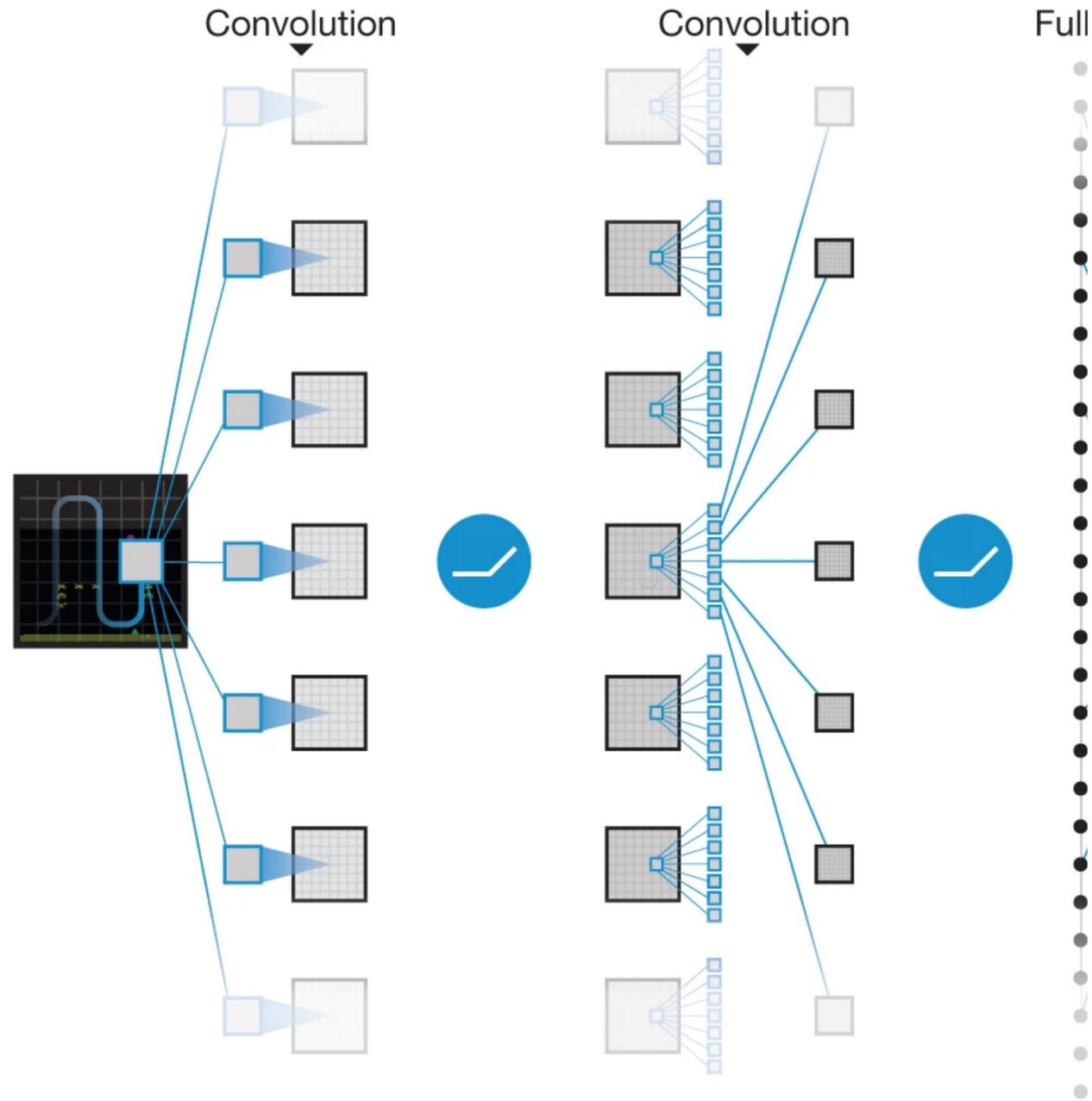
- convolution
- max pooling
- channel concatenation
- channel-wise normalization
- fully-connected layer
- softmax



FUN EXAMPLE: PLAYING THE ATARI SUITE

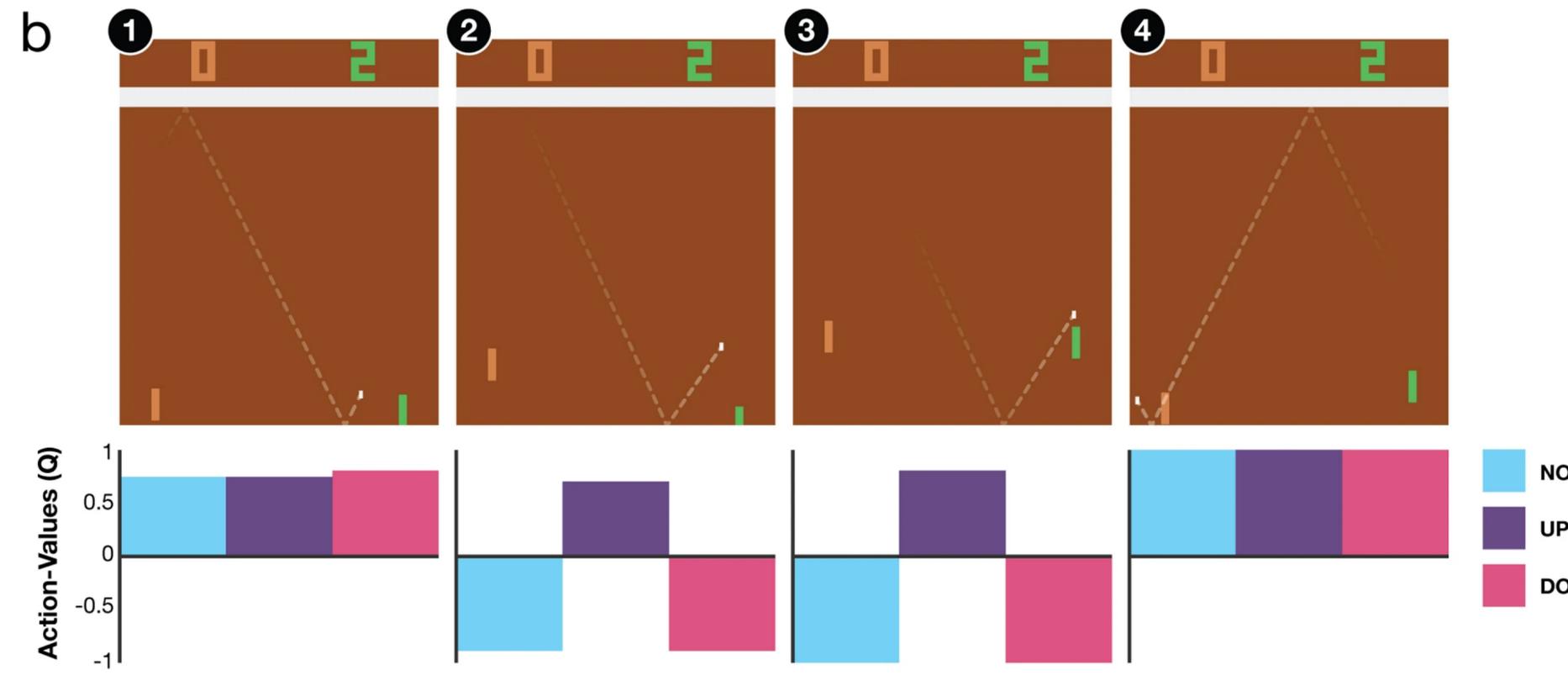
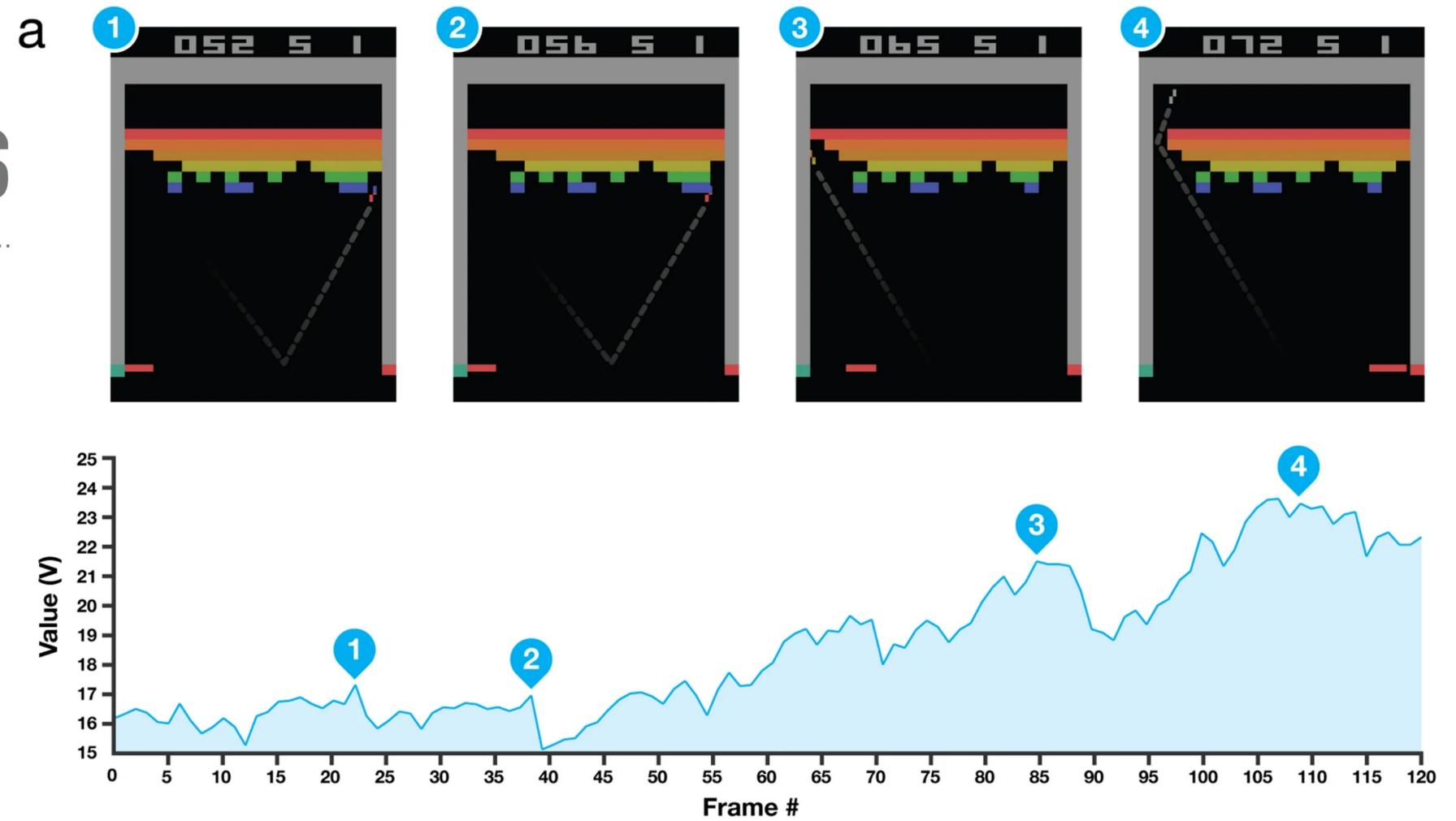


FUN EXAMPLE: PLAYING THE ATARI SUITE



FUN EXAMPLE: PLAYING THE ATARI S

► https://static-content.springer.com/esm/art%3A10.1038%2Fnature14236/MediaObjects/41586_2015_BFnature14236_MOESM124_ESM.mov



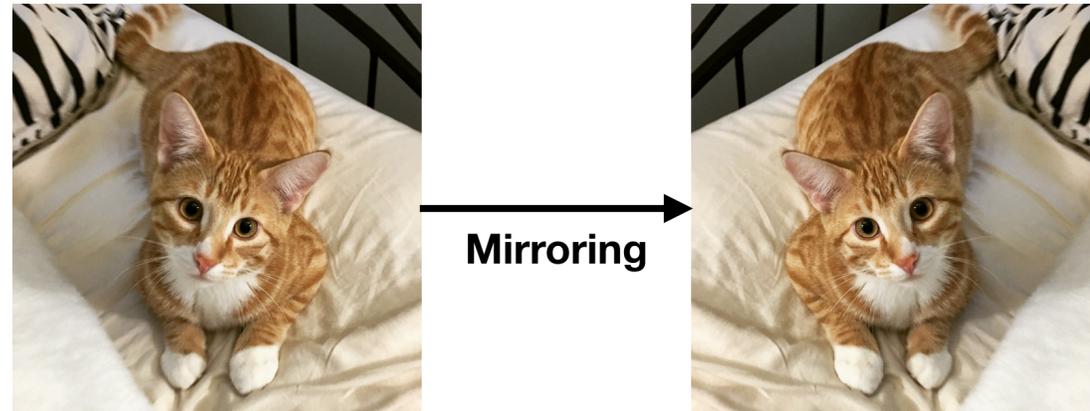
DATA AUGMENTATION

- Improve generalization error of classifier by adding copies of data samples that have been transformed in such a way that class does not change



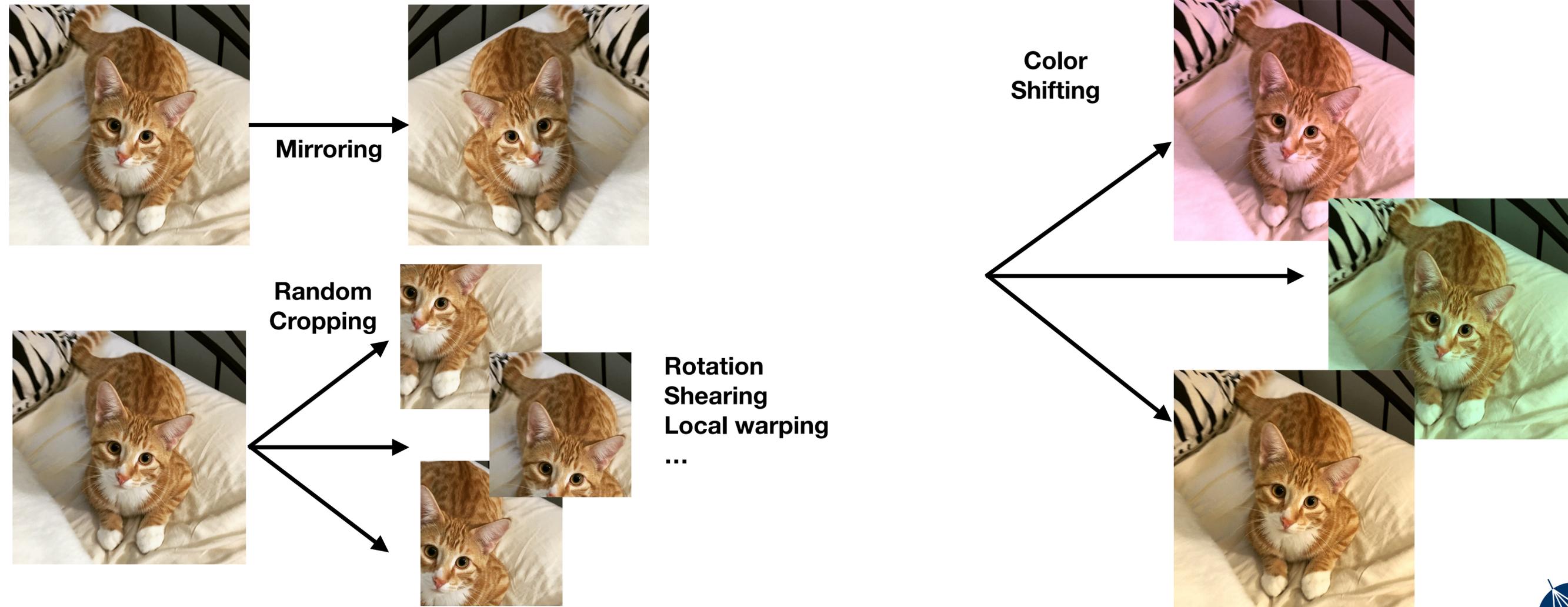
DATA AUGMENTATION

- Improve generalization error of classifier by adding copies of data samples that have been transformed in such a way that class does not change



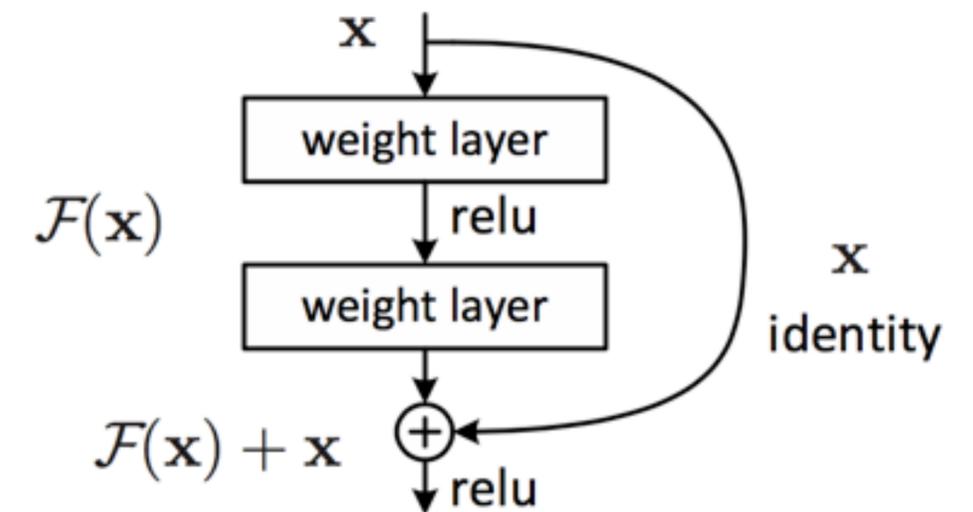
DATA AUGMENTATION

- Improve generalization error of classifier by adding copies of data samples that have been transformed in such a way that class does not change



4 TYPICAL CNN ARCHITECTURES: COULD GO MUCH DEEPER HERE

- Degradation problem, batch normalisation, residual blocks, ResNet
- Dense convolutional networks
- Network efficiency, ...
-



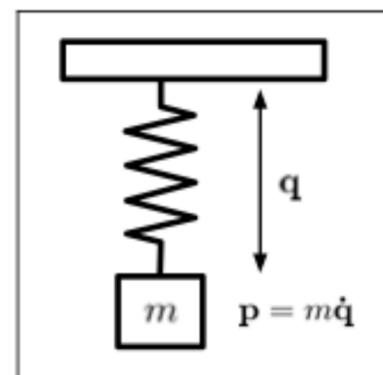
CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. Zhang et al. 2020 or Wu et al., 2018
 - Use physics laws as a guideline for constructing NNs (e.g. Hamiltonian NN)
 - Letting ML find optimal observables (e.g. Datta et al. 2019)
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (typically with GANs)
- Refinement of MC simulations to match data (also with GANs)

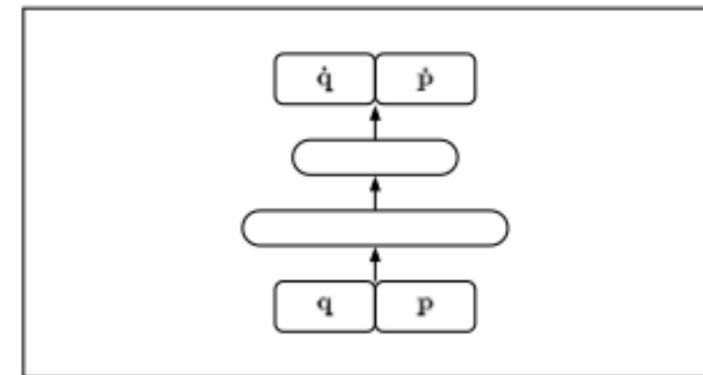
CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. [Zhang et al. 2020](#) or [Wu et al., 2018](#)
 - Use physics laws as a guideline for constructing NNs (e.g. [Hamiltonian NN](#))
 - Letting ML find optimal observables (e.g. [Datta et al. 2019](#))
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, [HLS4ML](#) open source, multi backend)
- Network causality - how are decisions taken? (e.g. [Kindermans et al. 2017](#))
- Uncertainty quantification (typically with GANs)
- Refinement of MC simul

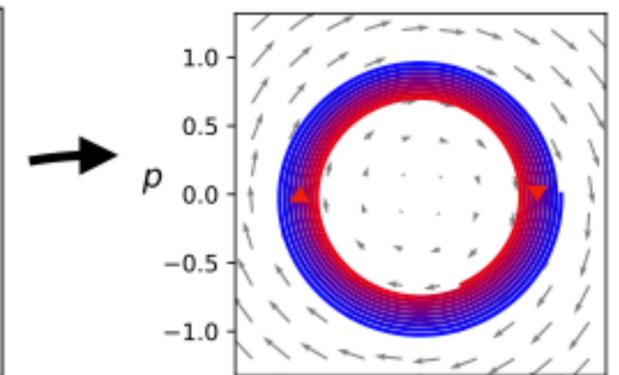
Ideal mass-spring system



Baseline NN



Prediction

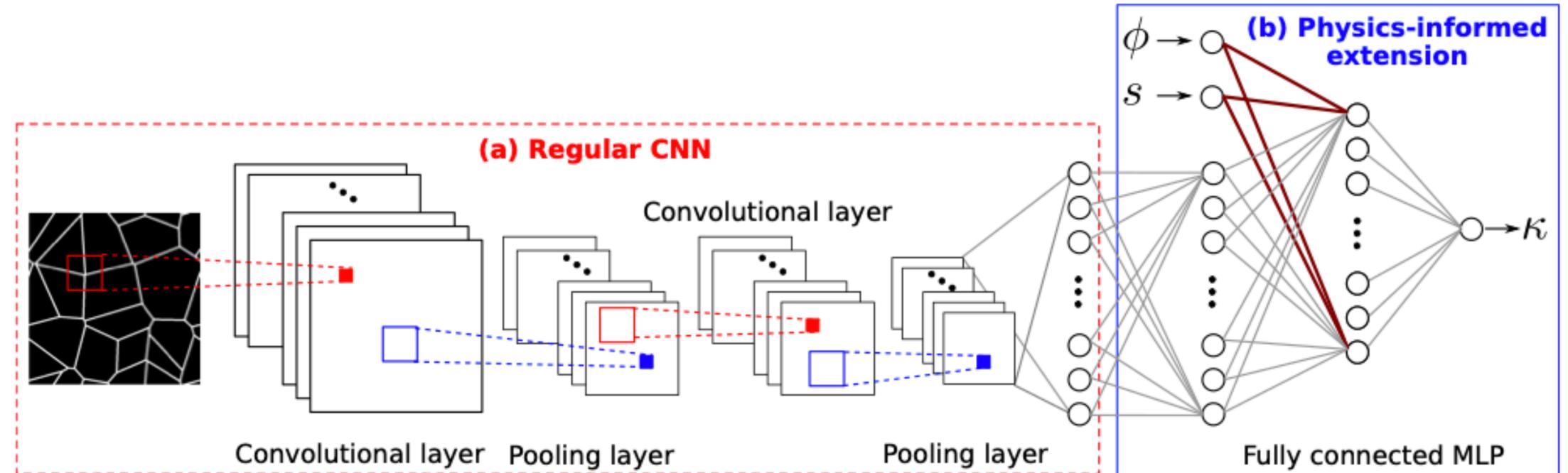


CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. Zhang et al. 2020 or Wu et al., 2018
 - Use physics laws as a guideline for constructing NNs (e.g. Hamiltonian NN)
 - Letting ML find optimal observables (e.g. Datta et al. 2019)
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (typically with GANs)
- Refinement of MC simulations to match data (also with GANs)

CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs
2020 or Wu et al., 2019
- Use physics laws as a constraint
- Letting ML find optimal physics parameters
- Improve computational efficiency
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (typically with GANs)
- Refinement of MC simulations to match data (also with GANs)

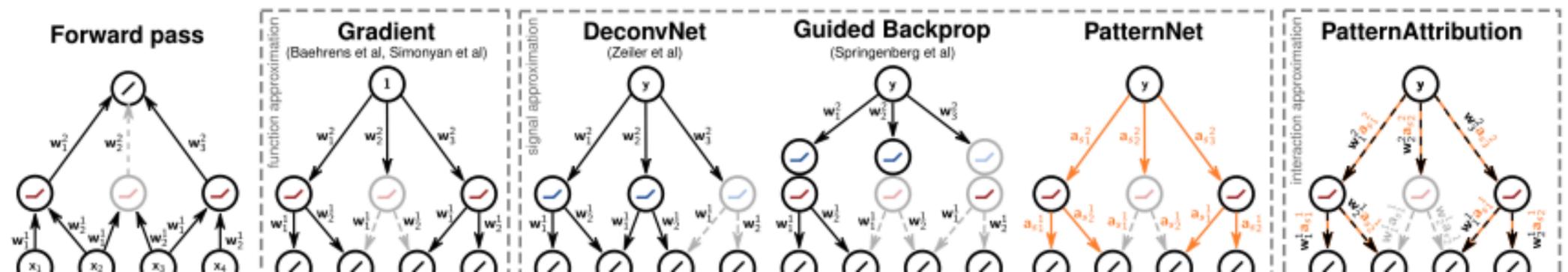


CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. Zhang et al. 2020 or Wu et al., 2018
 - Use physics laws as a guideline for constructing NNs (e.g. Hamiltonian NN)
 - Letting ML find optimal observables (e.g. Datta et al. 2019)
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (typically with GANs)
- Refinement of MC simulations to match data (also with GANs)

CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. Zhang et al. 2020 or Wu et al., 2018
 - Use physics laws as a guideline for constructing NNs (e.g. Hamiltonian NN)
 - Letting ML find optimal observables (e.g. Datta et al. 2019)
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (currently with GANs)
- Refinement of MC simulation



CURRENT EFFORTS TO IMPROVE CNNs

- Physics-guided CNNs (provide physics constraints on the network output) - e.g. Zhang et al. 2020 or Wu et al., 2018
 - Use physics laws as a guideline for constructing NNs (e.g. Hamiltonian NN)
 - Letting ML find optimal observables (e.g. Datta et al. 2019)
- Improve computational efficiency of CNNs, real-time processing
- CNNs on FPGAs (e.g. for L1T at LHC, HLS4ML open source, multi backend)
- Network causality - how are decisions taken? (e.g. Kindermans et al. 2017)
- Uncertainty quantification (typically with GANs)
- Refinement of MC simulations to match data (also with GANs)

NOW LET'S APPLY IT – TENSORFLOW, OR KERAS

- Convenient neural network package for python
- Set up and training of a network in a few lines
- Based on underlying neural network package (also provides run-time compilation to CPU and GPU) either *tensorflow* or *theano*.



SUMMARY

- Machines exploit physics contained in data deeper than before
- CNNs are the workhorse for many of the more advanced applications, as example in investigations of causality, stability, uncertainties

