

Module 5: Convolutional neural networks

5LSM0: Convolutional neural networks for computer vision

Joost van der Putten

Electrical Engineering / VCA research group



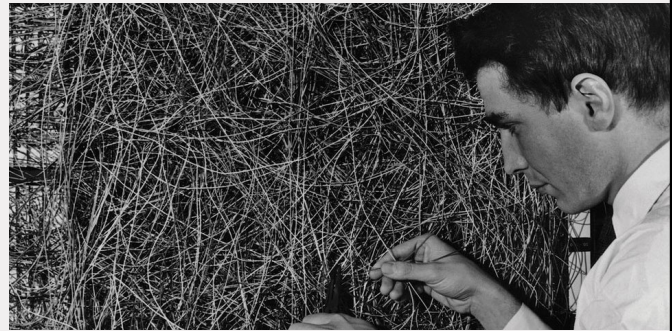
Module outline

- History of CNNs
- Another cat video
- CNNs today
- The convolution layer
 - Convolutions
 - Output dimensions of convolution process
 - Padding
- Activation functions
- Pooling



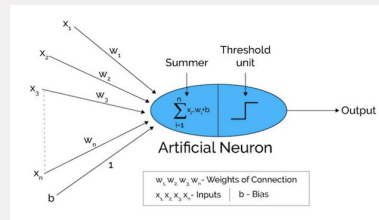
History of CNNs

- The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.
- The machine was connected to a camera that used 20x20 cadmium sulfide photocells to produce a 400-pixel image.

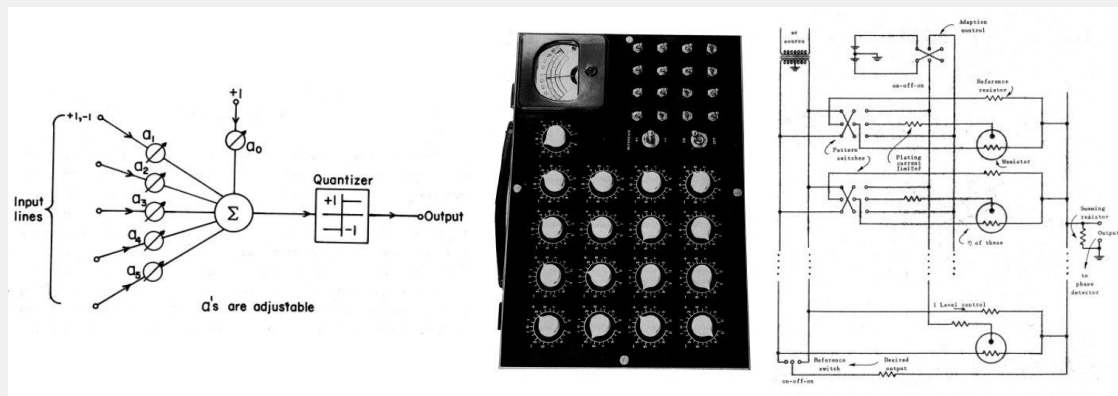


$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

Frank Rosenblatt, 1957: Perceptron



History of CNNs

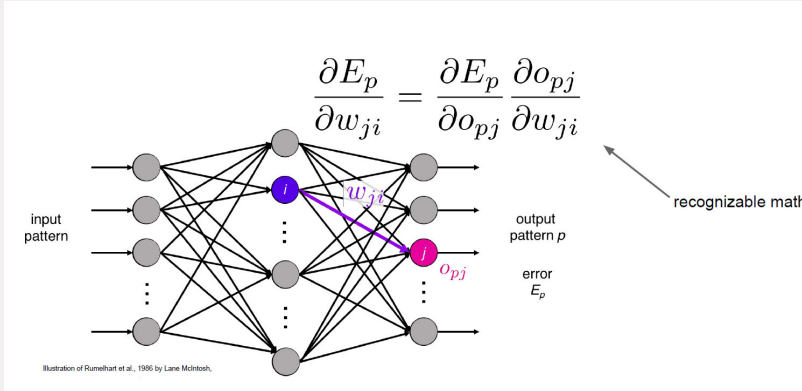


Widrow and Hoff, ~1960: Adaline/Madaline

- Still no backprop



History of CNNs



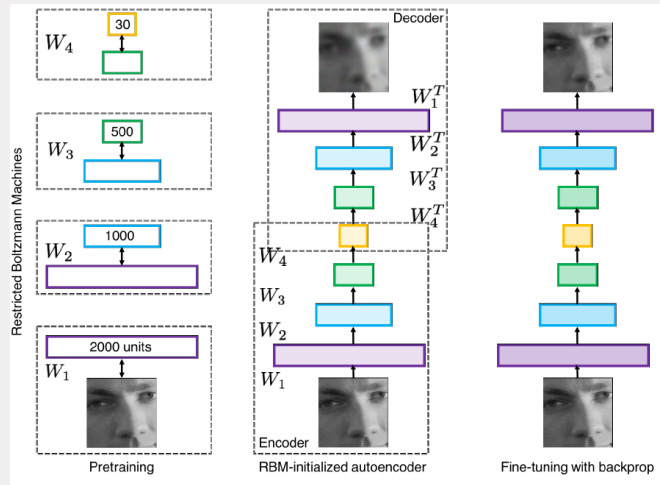
Rumelhart et al., 1986: First time back-propagation was used



History of CNNs

And then Nothing.
Until ...

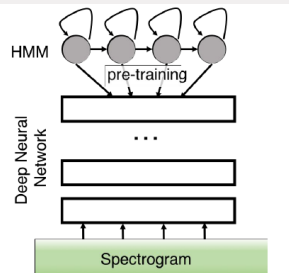
Hinton and Salakhutdinov 2006



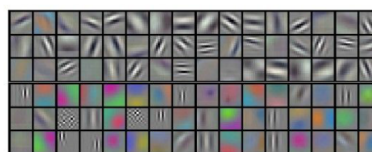
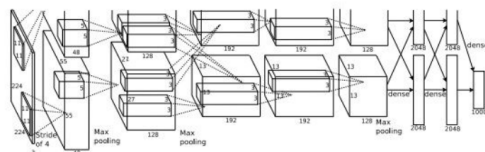
History of CNNs

First strong results

Acoustic Modeling using Deep Belief Networks
 Abdel-rahman Mohamed, George Dahl, Geoffrey Hinton, 2010
Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition
 George Dahl, Dong Yu, Li Deng, Alex Acero, 2012



Imagenet classification with deep convolutional neural networks
 Alex Krizhevsky, Ilya Sutskever, Geoffrey E Hinton, 2012



Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012.



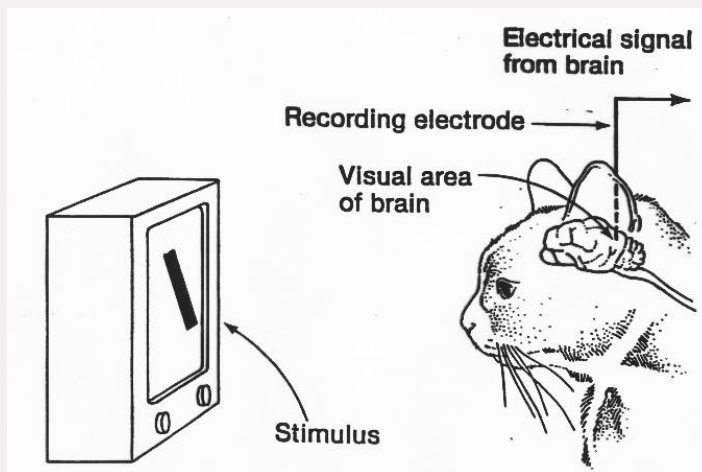
History of CNNs

Hubel & Wiesel, 1959

RECEPTIVE FIELDS OF SINGLE NEURONES IN THE CAT'S STRIATE CORTEX

1962

RECEPTIVE FIELDS, BINOCULAR INTERACTION AND FUNCTIONAL ARCHITECTURE IN THE CAT'S VISUAL CORTEX



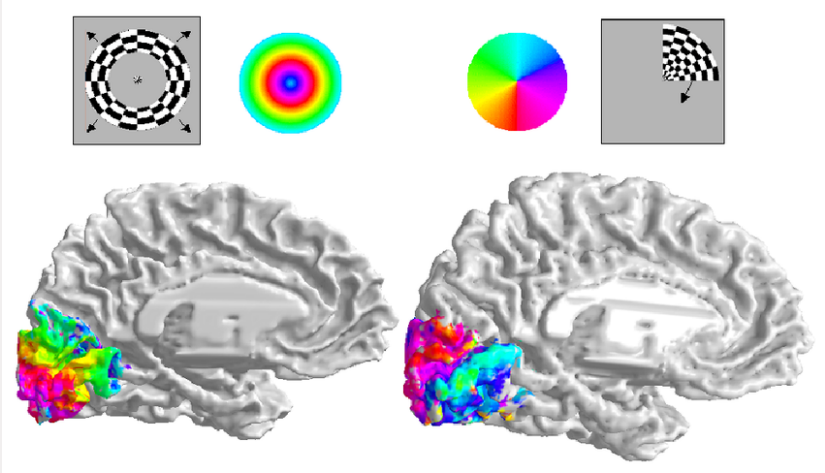
Cat video



9 SLSM0 Module 5: Convolutional neural networks **TU/e**

History of CNNs

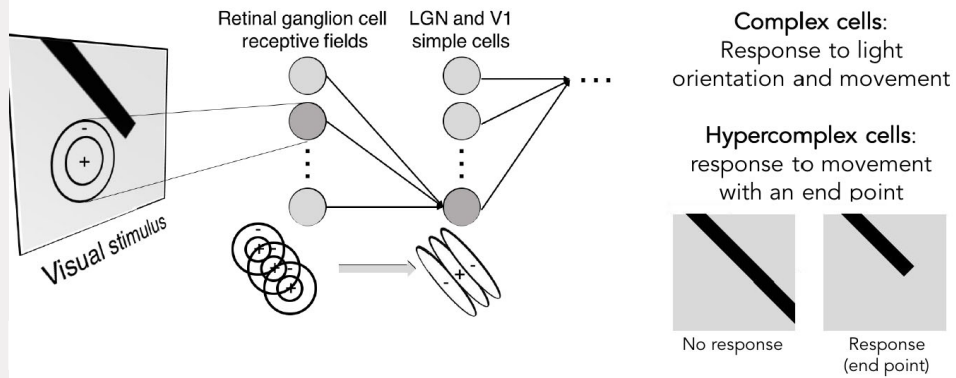
- Nearby cells in visual cortex represent nearby regions in the visual field



10 SLSM0 Module 5: Convolutional neural networks **TU/e**

History of CNNs

Hierarchical organization

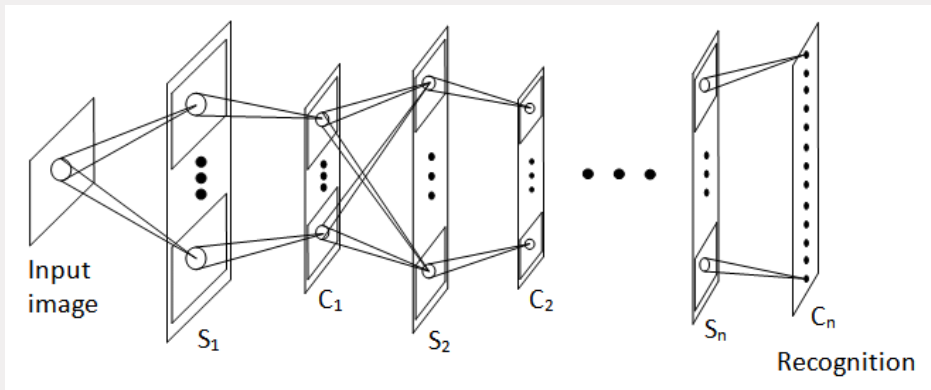


11 SLSM0 Module 5: Convolutional neural networks



History of CNNs

- Neurocognitron (Fukushima 1980)
- Combination of Simple cells (modifiable parameters) and Complex cells (perform pooling)



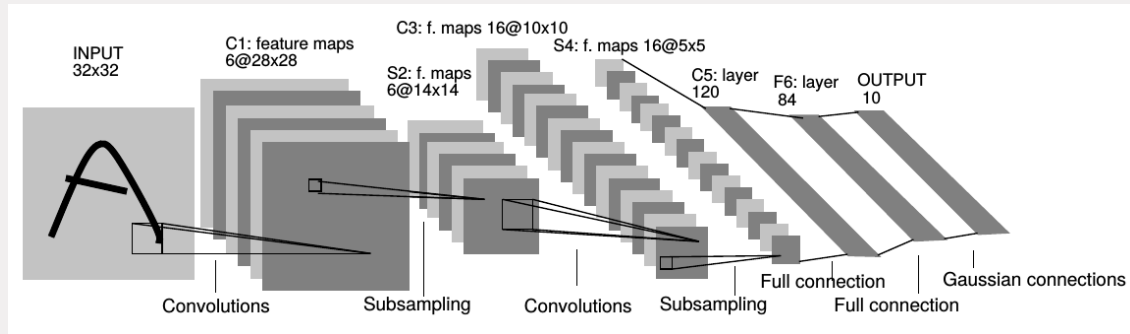
12 SLSM0 Module 5: Convolutional neural networks



History of CNNs

Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner 1998]



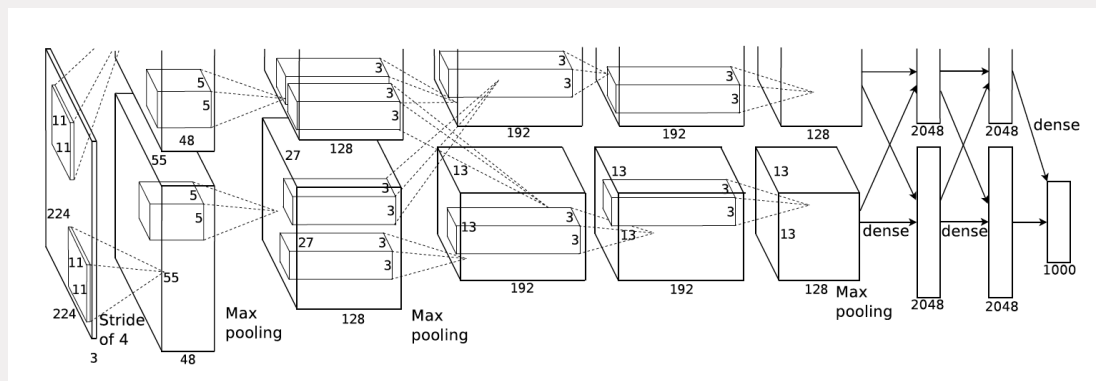
13 SLSM0 Module 5: Convolutional neural networks



History of CNNs

ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]

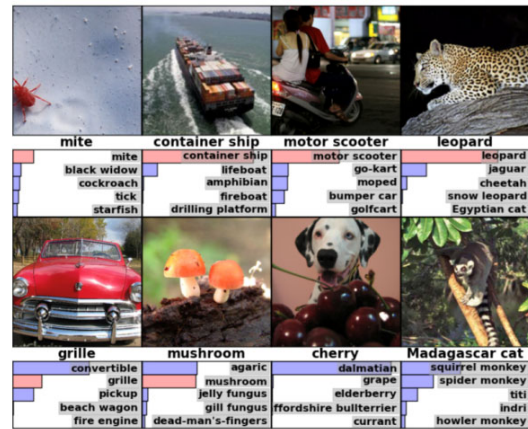


14 SLSM0 Module 5: Convolutional neural networks



CNNs today

Classification

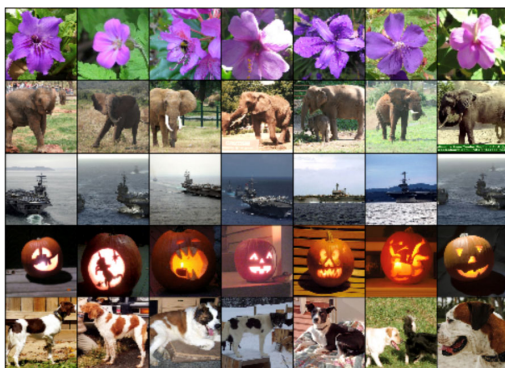


Detection



CNNs today

Image retrieval

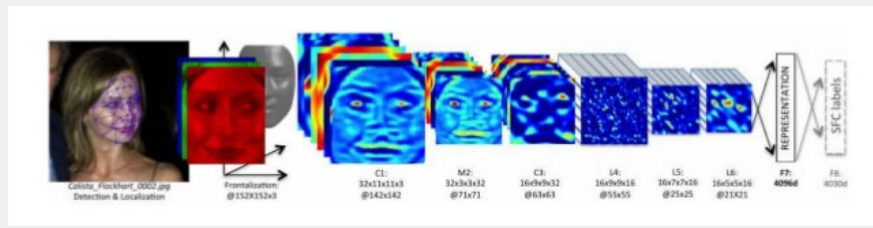


Semantic segmentation

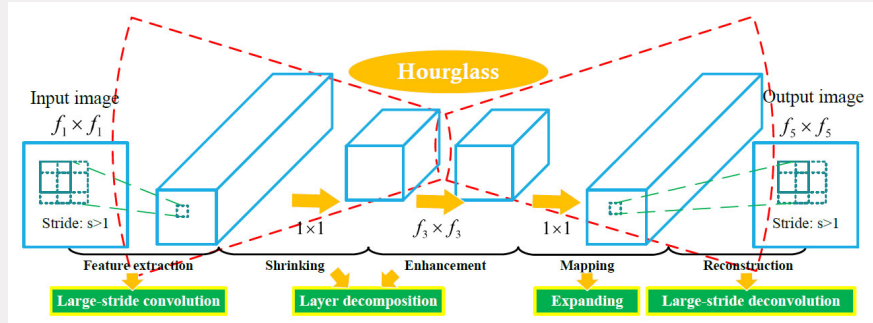


CNNs today

Face verification



Compression



CNNs today

Traffic sign detection

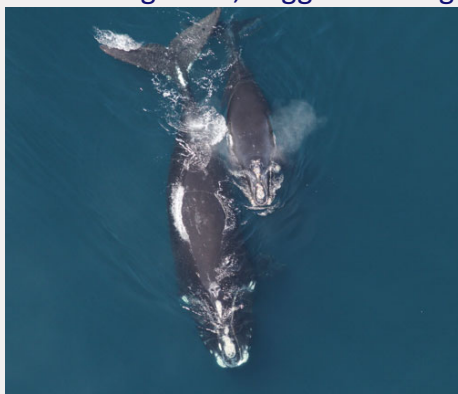
Test images found on the internet

Speed limit (30km/h)	Bumpy road	General caution	Keep right	No entry	Slippery road
✓	✓	✓	✓	✓	✓
94%	99%	100%	100%	99%	100%

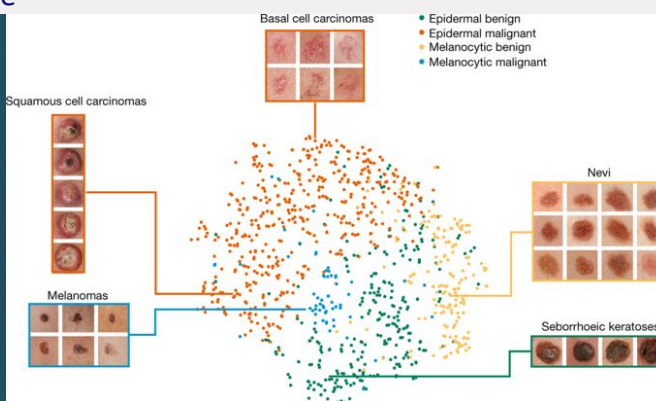


CNNs today

Whale recognition, Kaggle challenge



Skin cancer classification



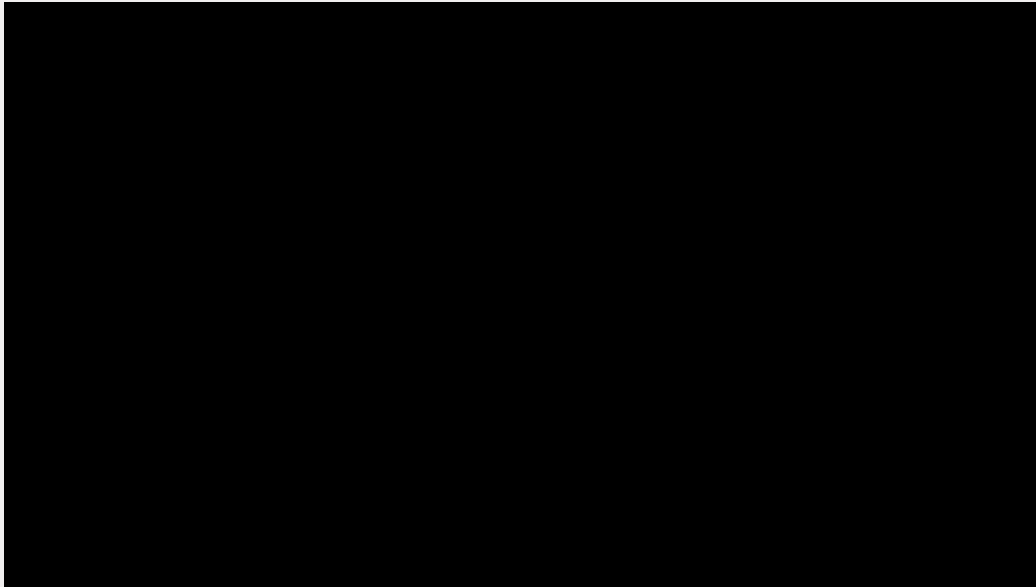
Pose estimation



<https://www.youtube.com/watch?v=pW6nZXeWIGM&t=7s>



AlphaStar



21 SLSM0 Module 5: Convolutional neural networks



CNNs today

Deep dream



22 SLSM0 Module 5: Convolutional neural networks



CNNs today

Road segmentation

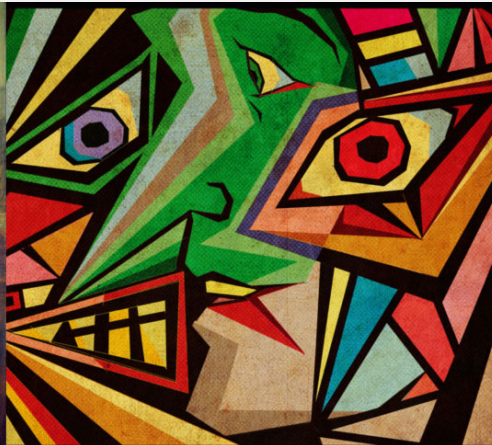


23 5LSM0 Module 5: Convolutional neural networks



CNNs today

Style transfer

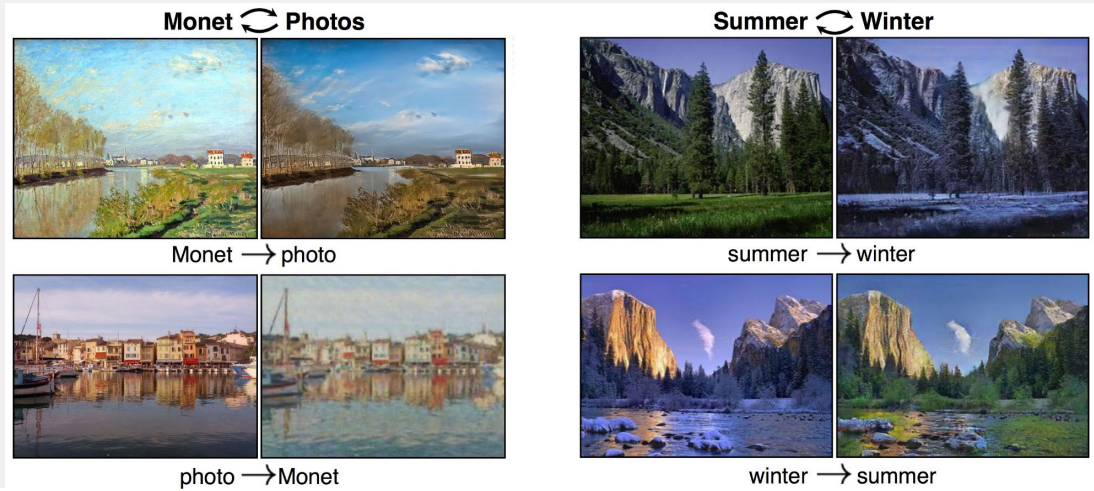


24



CNNs today

Style transfer



Cnns today

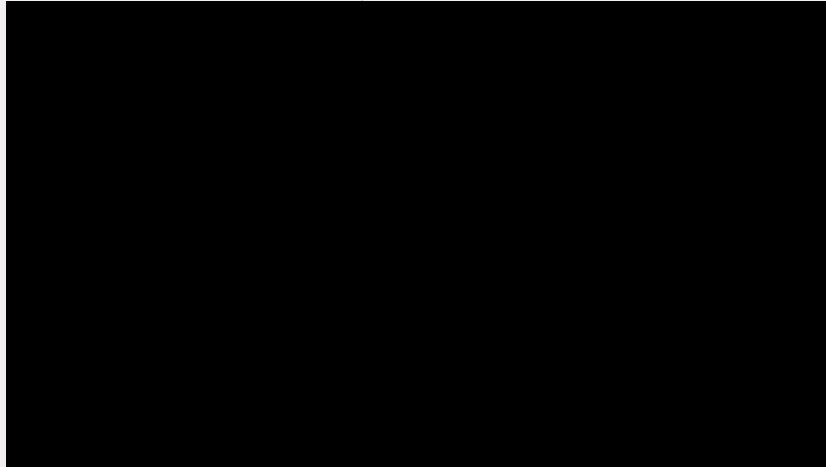


Image captioning



CNNs today

Style transfer



27 5LSM0 Module 5: Convolutional neural networks

TU/e

CNNs today

Image colorization



28 5LSM0 Module 5: Convolutional neural networks

TU/e

CNNs today

Super resolution

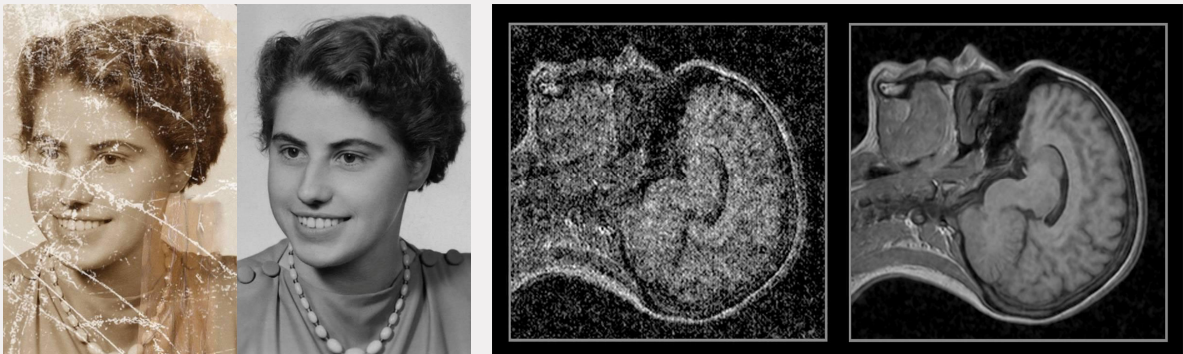


29 SLSM0 Module 5: Convolutional neural networks



CNNs today

Image restoration



30 SLSM0 Module 5: Convolutional neural networks



Convolutional Neural Networks

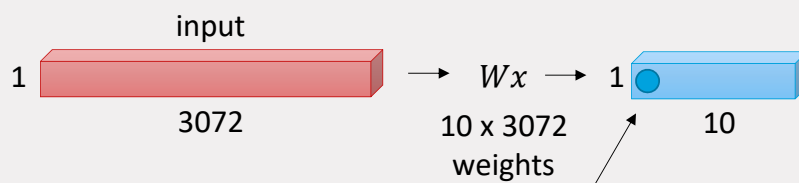
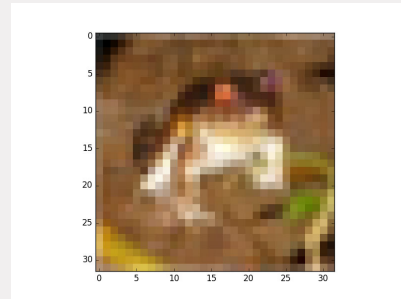


31 SLSM0 Module 5: Convolutional neural networks

TU/e

Fully connected layer

32x32x3 image -> stretch to 3072x1



1 number:
The result of taking a dot product
between a row of W and the input

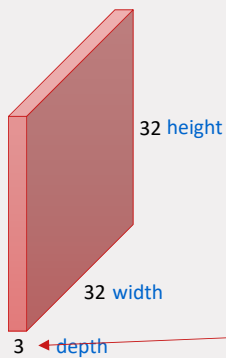


32 SLSM0 Module 5: Convolutional neural networks

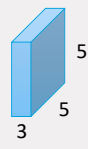
TU/e

Convolutional layer

32x32x3 image



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



Filters always extend the full depth of the input volume

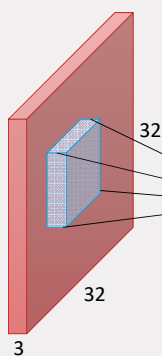


33 SLSM0 Module 5: Convolutional neural networks



Convolutional layer

32x32x3 image
5x5x3 filter w



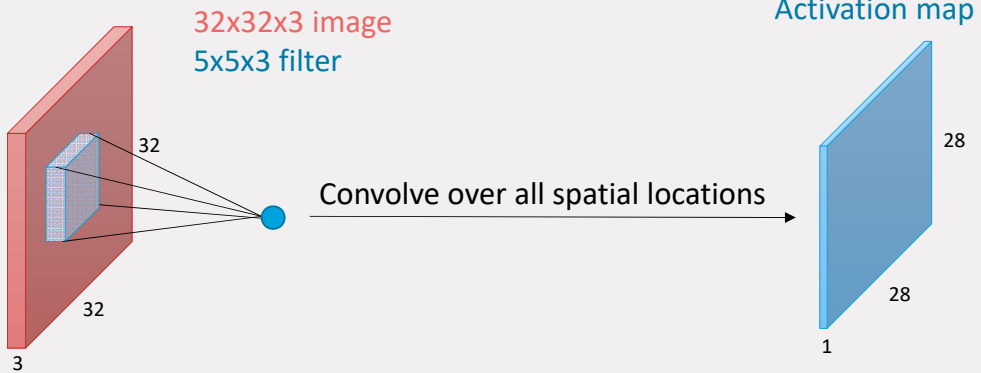
1 number:
The result of taking a dot product between the filter and a small 5x5x3 chunk of the image
 $w^T x + b$



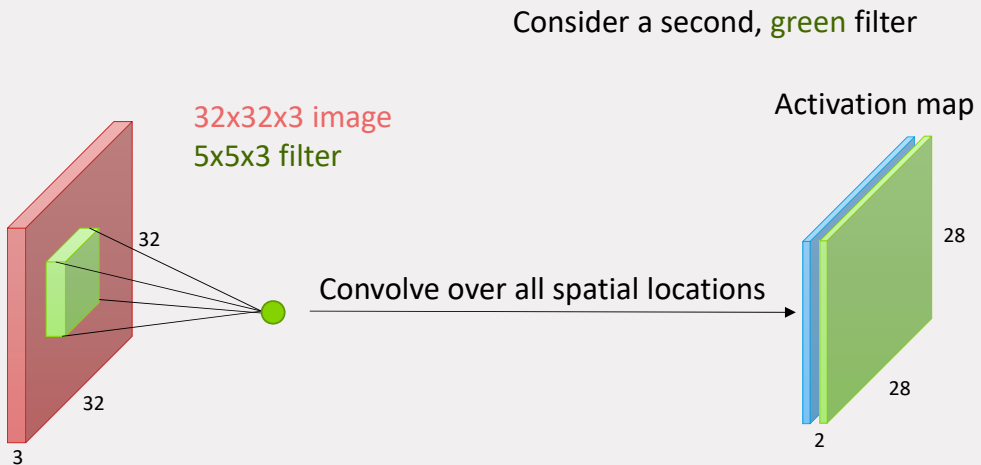
34 SLSM0 Module 5: Convolutional neural networks



Convolutional layer

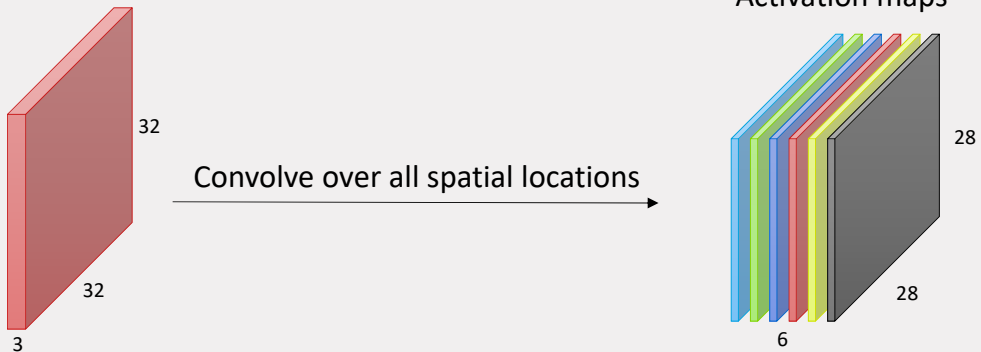


Convolutional layer



Convolutional layer

With 6 5x5 filters, we obtain 6 28x28 activation maps

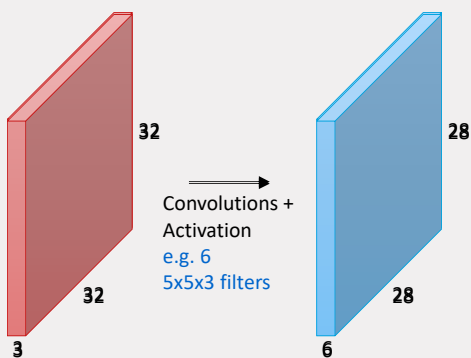


Stack the activation maps to obtain a “new image” of size 28x28x6



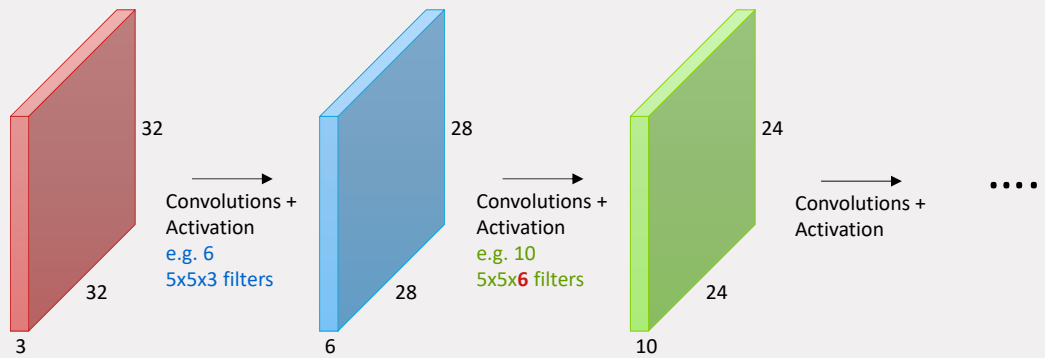
Preview to ConvNets

A **ConvNet** is a sequence of convolutional layers interspersed with activation functions.



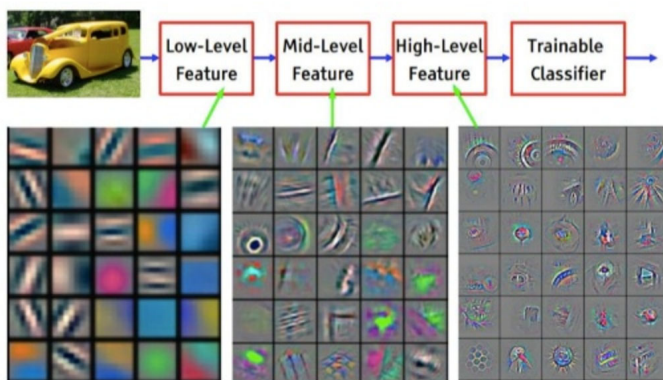
Preview to ConvNets

A **ConvNet** is a sequence of convolutional layers interspersed with activation functions.

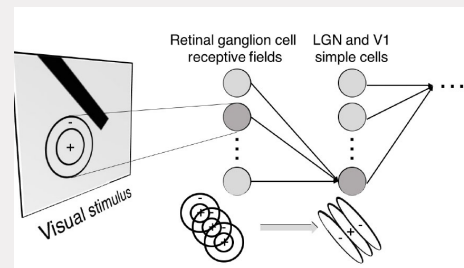


Preview to ConvNets

Convolutional Neural Network



Hubel & Wiesel



Convolutional neural network learns hierarchical structure on its own



Preview to ConvNets

one filter =>
one activation map

example 5x5 filters
(32 total)

Activations:

We call the layer convolutional because it is related to convolution of two signals:

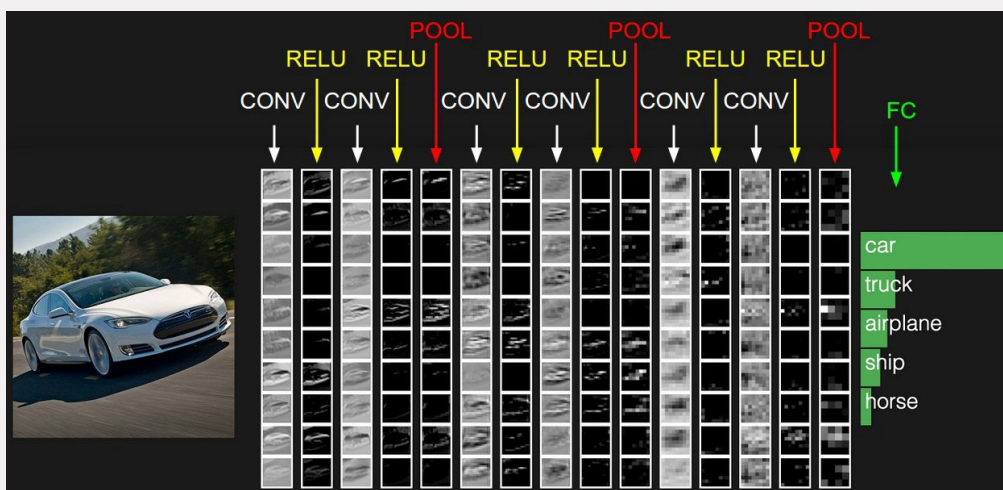
$$f[x,y] * g[x,y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1,y-n_2]$$

↑
elementwise multiplication and sum of a filter and the signal (image)

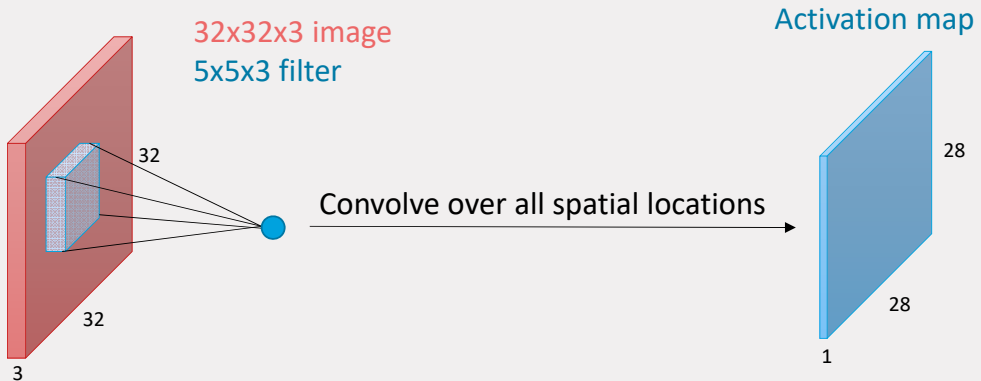
Figure copyright Andrej Karpathy.



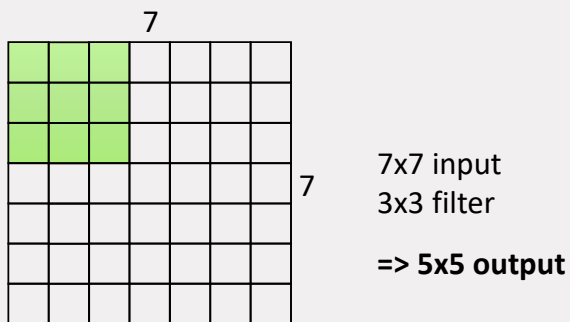
Preview to ConvNets



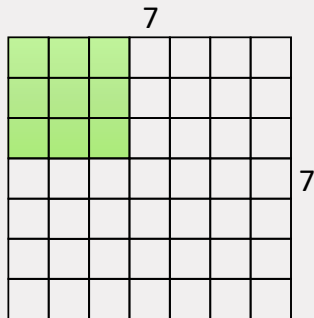
A closer look at spatial dimensions



A closer look at spatial dimensions



A closer look at spatial dimensions



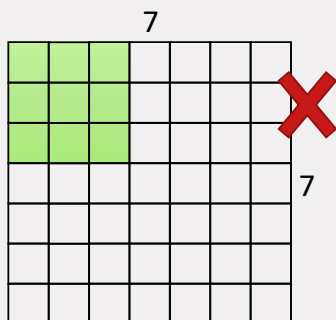
7x7 input
 3x3 filter applied with **Stride 2**
 => **3x3 output**



45 SLSM0 Module 5: Convolutional neural networks

TU/e

A closer look at spatial dimensions



7x7 input
 3x3 filter applied with **Stride 3**

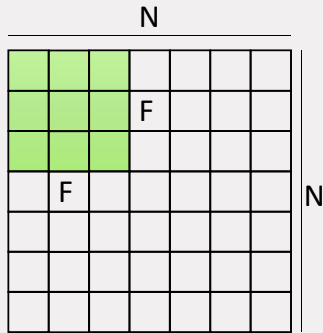
Doesn't fit!
 Cannot apply 3x3 filter on
 7x7 input with stride 3



46 SLSM0 Module 5: Convolutional neural networks

TU/e

A closer look at spatial dimensions



Output size:
 $(N-F)/\text{stride} + 1$

e.g. $N=7, F=3$:

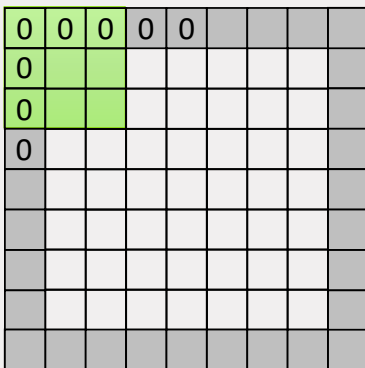
Stride 1 $\rightarrow (7-3)/1 + 1 = 5$

Stride 2 $\rightarrow (7-3)/2 + 1 = 3$

Stride 3 $\rightarrow (7-3)/3 + 1 = 2.33 \text{ ☹}$



In practice: Zero padding



e.g. input **7x7**

3x3 filter, applied with **stride 1**

Pad with 1 pixel border

What are the output dimensions?

7x7 output!

Recall:

$(N-F)/\text{stride} + 1$



In practice: Zero padding

0	0	0	0	0					
0									
0									
0									

e.g. input **7x7**

3x3 filter, applied with **stride 1**

Pad with 1 pixel border

What are the output dimensions?

7x7 output!

In general, it is common to see convolutional layers with stride 1, filters of size $F \times F$, and zero-padding with $(F-1)/2$. (Will preserve size spatially)

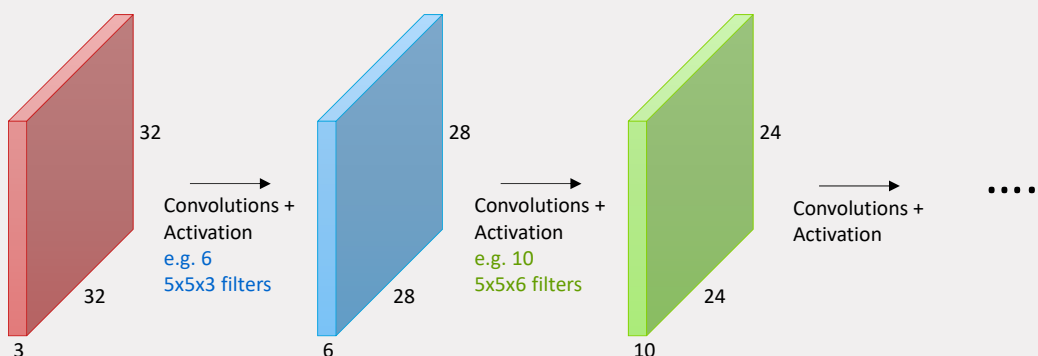
$F = 3 \rightarrow$ zero pad with 1 pixel

$F = 5 \rightarrow$ zero pad with 2 pixels



In practice: Zero padding

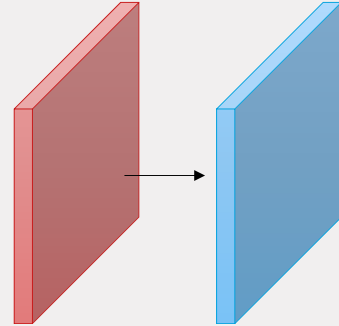
Repeated convolutions shrink volumes spatially! (32 \rightarrow 28 \rightarrow 24 ...). Shrinking too fast does not work well in practice.



Spatial dimension and number of parameters example

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Output volume size: ?



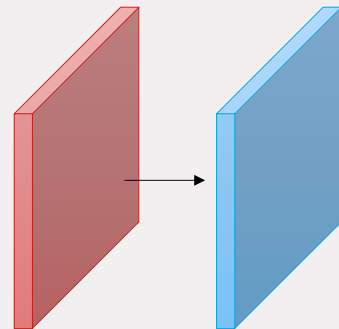
51 SLSM0 Module 5: Convolutional neural networks

TU/e

Spatial dimension and number of parameters example

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Output volume size:
 $(32+2*2-5)/1+1 = 32$ spatially, so
32x32x10



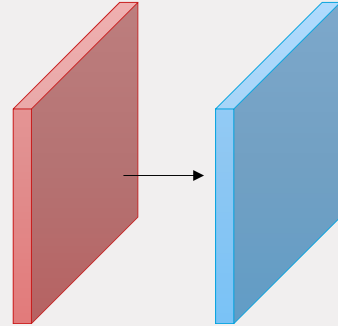
52 SLSM0 Module 5: Convolutional neural networks

TU/e

Spatial dimension and number of parameters example

Input volume: **32x32x3**
 10 5x5 filters with stride 1, pad 2

Number of parameters in this layer: ?



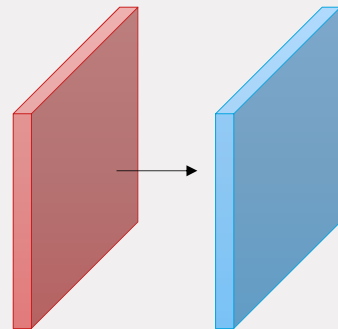
53 SLSM0 Module 5: Convolutional neural networks

TU/e

Spatial dimension and number of parameters example

Input volume: **32x32x3**
 10 5x5 filters with stride 1, pad 2

Number of parameters in this layer:
 Each filter has $5 \times 5 \times 3 + 1 = 76$ params
 $\rightarrow 76 \times 10 = 760$



Each filter has a bias term

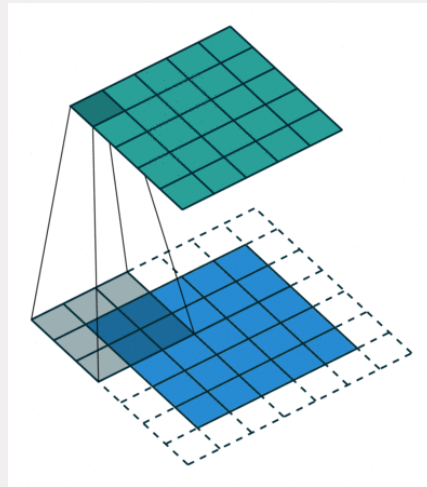


54 SLSM0 Module 5: Convolutional neural networks

TU/e

Summary

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K
 - Their spatial extent F
 - The stride S
 - The amount of padding P
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = \frac{W_1 - F + 2P}{S} + 1$
 - $H_2 = \frac{H_1 - F + 2P}{S} + 1$
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.



Summary

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K
 - Their spatial extent F
 - The stride S
 - The amount of padding P
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = \frac{W_1 - F + 2P}{S} + 1$
 - $H_2 = \frac{H_1 - F + 2P}{S} + 1$
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Common settings:

$K = 32, 64, 128, 512$ (power of 2)

- $F = 3, S=1, P=1$

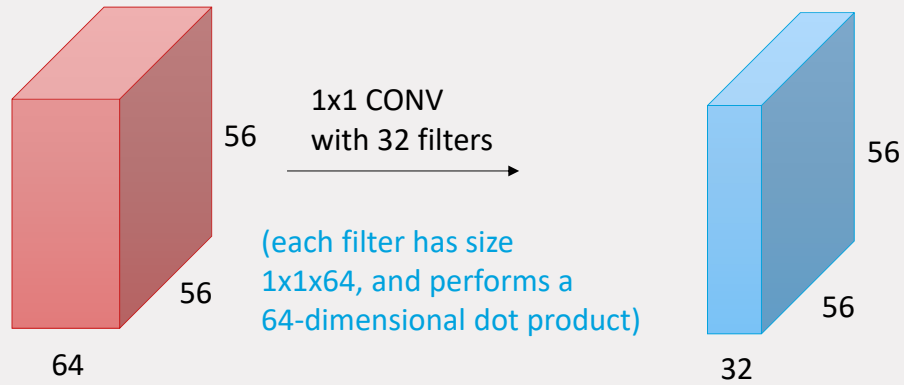
- $F = 5, S=1, P=2$

- $F = 5, S=2, P=?$

- $F = 1, S=1, P=0$



1x1 convolutions



Pytorch example

```

1 nn.Conv2d(in_channels=3,
2           out_channels=16,
3           kernel_size=3,
4           stride=1,
5           padding=1)
    
```

Conv2d

CLASS `torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)` [SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

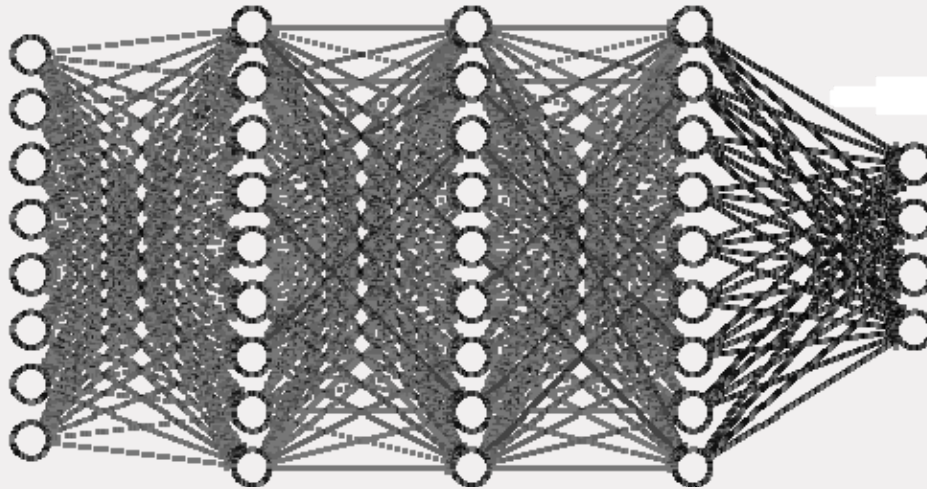
$$out(N_i, C_{out_j}) = bias(C_{out_j}) + \sum_{k=0}^{C_{in}-1} weight(C_{out_j}, k) * input(N_i, k)$$

where $*$ is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

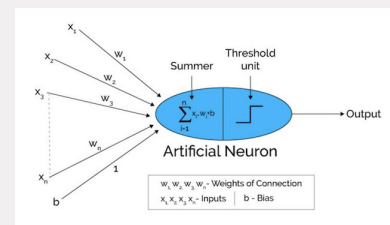
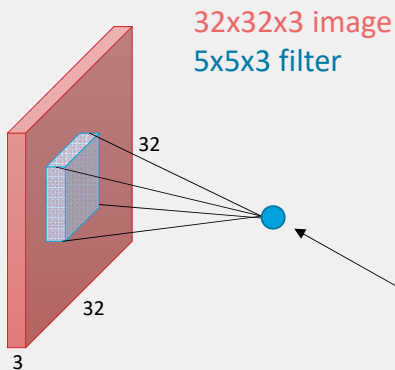
- `stride` controls the stride for the cross-correlation, a single number or a tuple.
- `padding` controls the amount of implicit zero-paddings on both sides for `padding` number of points for each dimension.
- `dilation` controls the spacing between the kernel points; also known as the à trous algorithm. It is harder to describe, but this [link](#) has a nice visualization of what `dilation` does.
- `groups` controls the connections between inputs and outputs. `in_channels` and `out_channels` must both be divisible by `groups`. For example,
 - At `groups=1`, all inputs are convolved to all outputs.
 - At `groups=2`, the operation becomes equivalent to having two conv layers side by side, each seeing half the input channels, and producing half the output channels, and both subsequently concatenated.
 - At `groups= in_channels`, each input channel is convolved with its own set of filters, of size: $\begin{bmatrix} C_{out} \\ C_{in} \end{bmatrix}$.



Recap: fully connected network



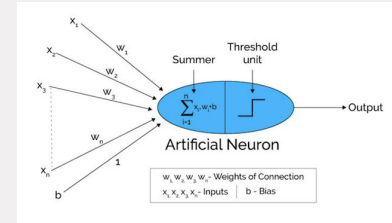
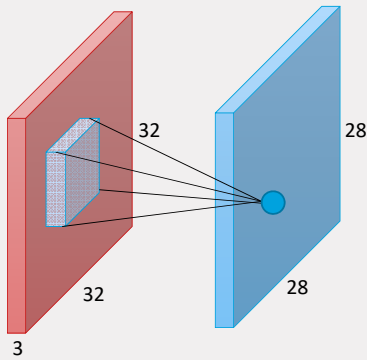
Parallel with the brain



Its just a neuron with local connectivity



Parallel with the brain



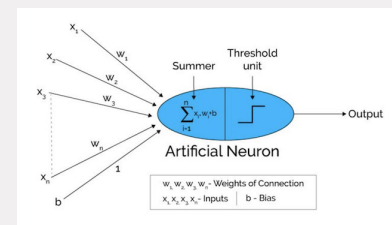
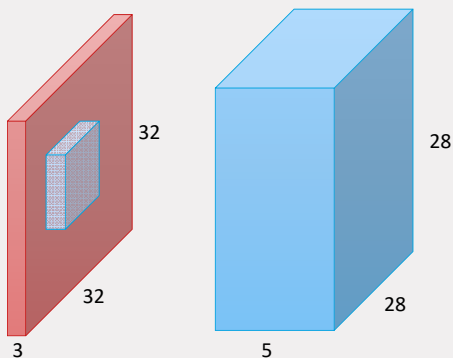
An activation map is a 28x28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share the same parameters

5x5 filter -> 5x5 'receptive field' for each neuron



Parallel with the brain

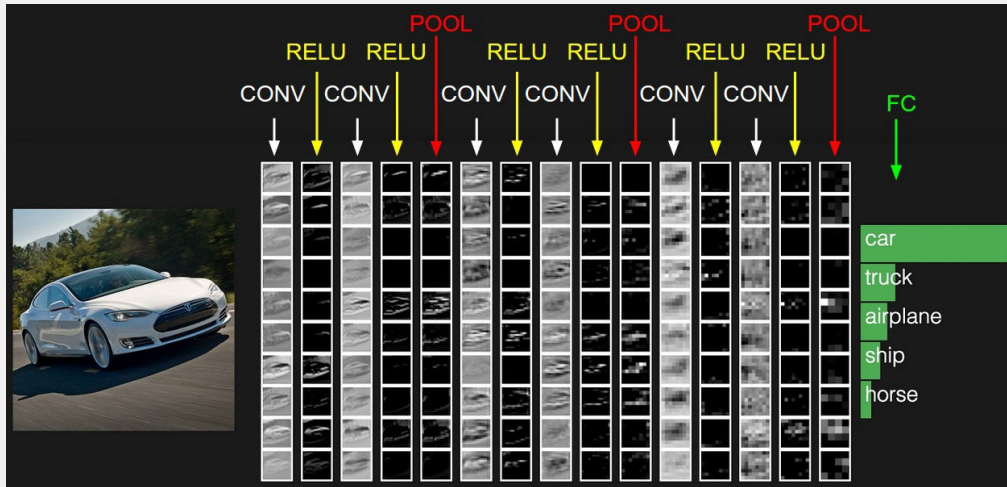


e.g. with 5 filters, CONV layer consists of neurons arranged in a 3D grid (28x28x5)

There will be 5 different neurons all looking at the same region in the input volume



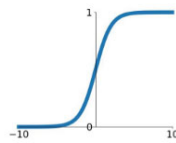
Activation functions



Activation Functions

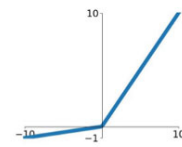
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



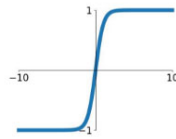
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

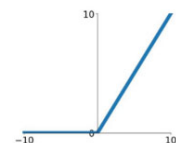


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

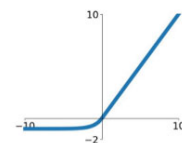
ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Activation functions - Sigmoid

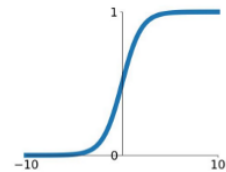
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating 'firing rate' of a neuron

3 problems

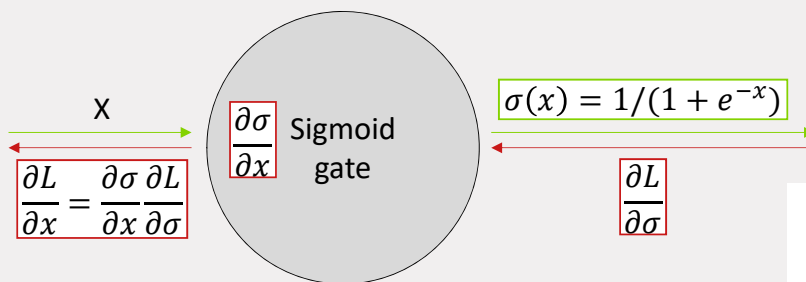
- **Saturated neurons 'kill off' the gradient**

Sigmoid

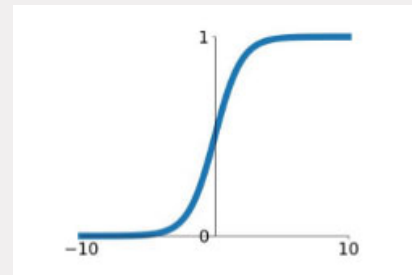
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Activation functions - Sigmoid



What happens when $x = -10$?
 What happens when $x = 0$?
 What happens when $x = 10$?



Activation functions - Sigmoid

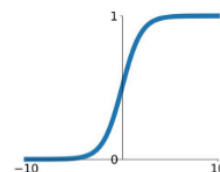
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating 'firing rate' of a neuron

3 problems

- Saturated neurons 'kill off' the gradient
- Sigmoid outputs are not zero-centered

Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



67 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions - Sigmoid

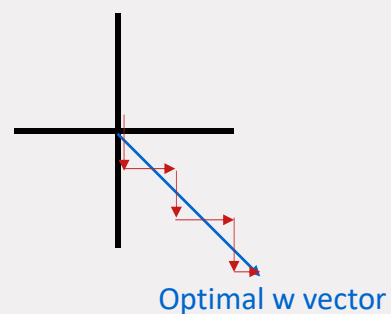
What happens when the input to a neuron is always positive ...

$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on \mathbf{w} ?

Always all positive or all negative ☹️

(This is also why we want zero-centered data)



68 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions - Sigmoid

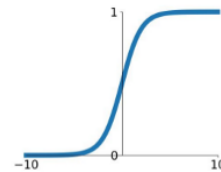
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating 'firing rate' of a neuron

3 problems

- Saturated neurons 'kill off' the gradient
- Sigmoid outputs are not zero-centered
- Exp() is relatively expensive

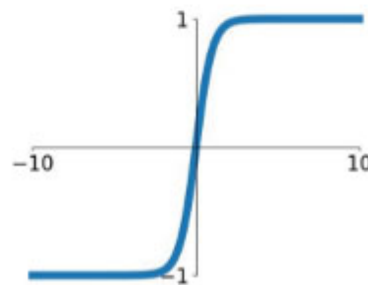
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Activation functions - Tanh

- Squashes numbers to range [-1,1]
- Zero-centered 😊
- Still kills gradients when saturated 😞

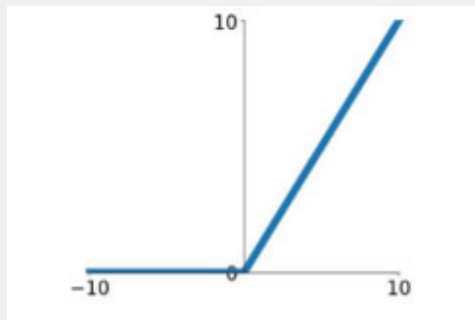


$\tanh(x)$



Activation functions - ReLU

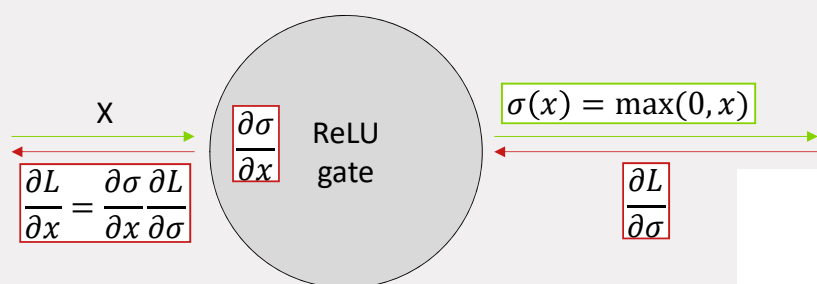
- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice ($\sim 6x$)
- Actually more biologically plausible than sigmoid
- Not zero-centered output
- One small annoyance?



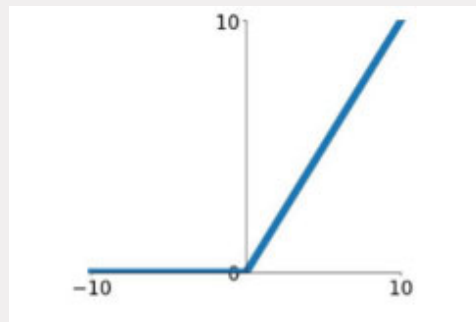
71 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions - ReLU



What happens when $x = -10$?
 What happens when $x = 0$?
 What happens when $x = 10$?

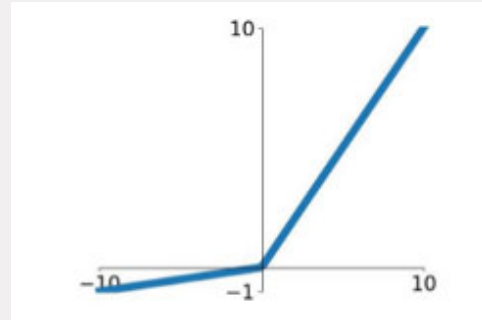


72 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions – Leaky ReLU

- Computes $f(x) = \max(0.01x, x)$
- **Does not saturate**
- **Very computationally efficient**
- **Converges much faster than sigmoid/tanh in practice (~6x)**
- **Will not 'die'**



- **Parametric ReLU** Computes $f(x) = \max(\alpha x, x)$

Learnable parameter α



73 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions - Maxout

- Does not have the basic form of dot product -> nonlinearity
- **Generalises ReLU and Leaky ReLU**
- **Linear regime! Does not saturate! Does not die!**

Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

- **Doubles the number of parameters/neuron ☹**



74 SLSM0 Module 5: Convolutional neural networks

TU/e

Activation functions: TLDR

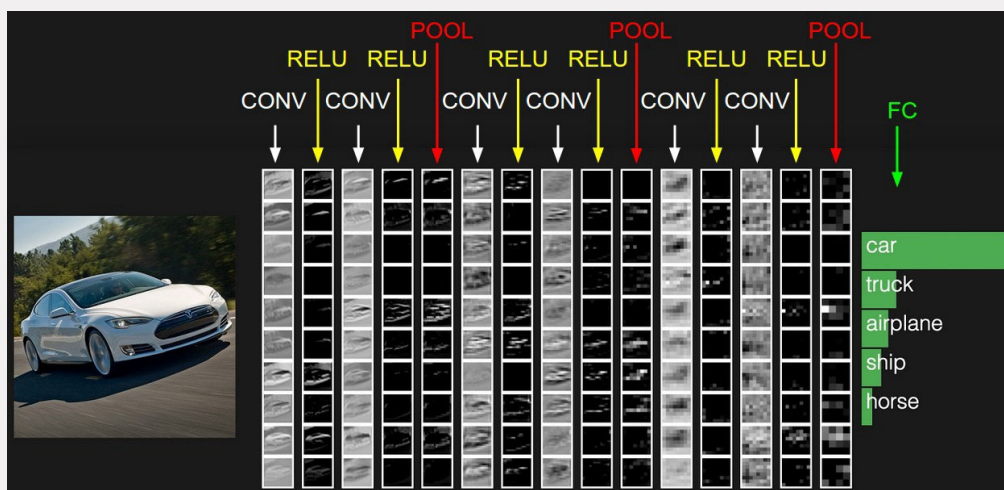
- Use **ReLU**.
- Try out **Leaky ReLU / Maxout**
- **Don't use sigmoid/tanh**



75 SLSM0 Module 5: Convolutional neural networks

TU/e

Pooling

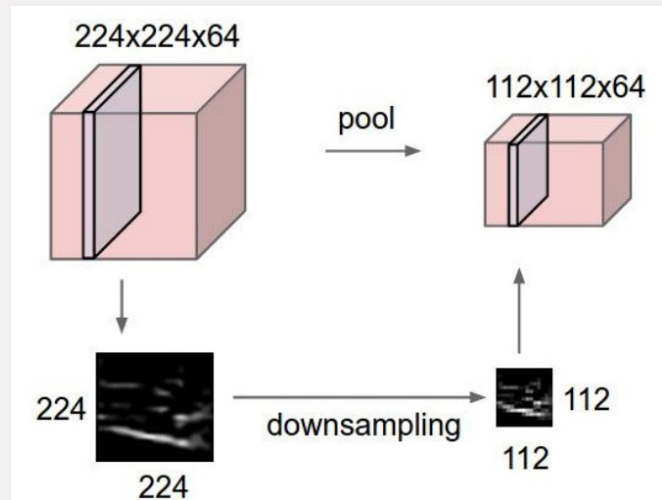


76 SLSM0 Module 5: Convolutional neural networks

TU/e

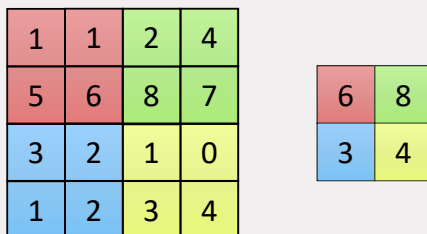
Pooling layer

- Makes the representations smaller and more manageable
- Operates over each activation map independently

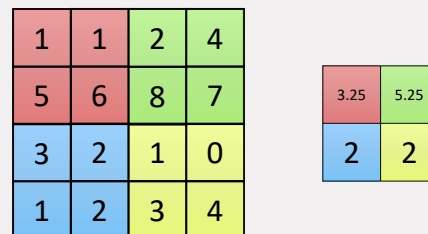


Pooling

Max pool with 2x2 filter and stride 2



Average pool with 2x2 filter and stride 2

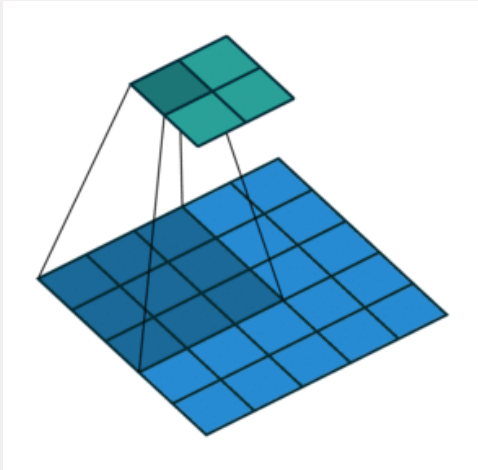


Typically no zero-padding in pooling operations

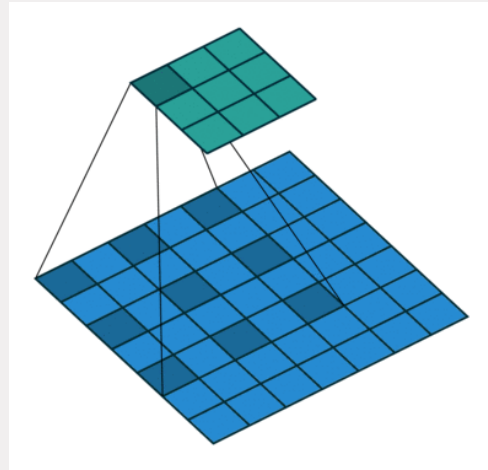


Pooling

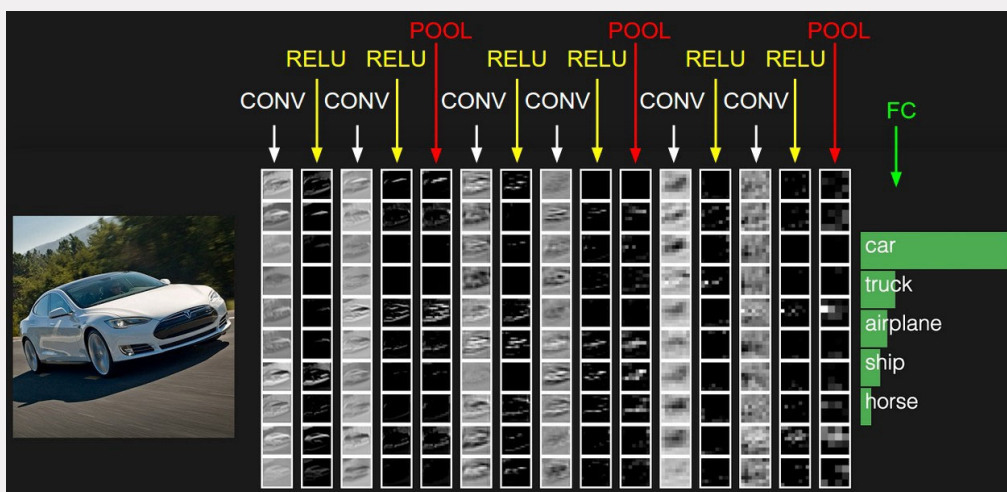
Strided convolutions



dilated convolutions



Fully connected layer (FC layer)



Summary

- ConvNets stack CONV, POOL and FC layers
- Trend towards smaller filters and deeper architectures
- Trend of getting rid of POOL/FC layers, use just convolutional layers instead
- Typical architectures look as follows:
 - $[(\text{CONV} \rightarrow \text{Relu}) * N \rightarrow \text{POOL}] * M \rightarrow (\text{FC} \rightarrow \text{ReLU}) * K \rightarrow \text{softmax}$
where N is usually up to ≈ 3 , M is ≈ 5 , and $0 \leq K \leq 2$.

