

Applied Machine Learning

Convolutional Neural Networks

Oumar Kaba



Admin

- Exam grading
- Quizz
- Assignment 3

Learning objectives

understand the convolution layer and the architecture of Conv-net

- its inductive bias
- its derivation from fully connected layer
- variations of convolution layer

MLP and image data

5	0	4	1
3	5	3	6
4	0	9	1
3	8	6	9

we can apply an MLP to image data

3	5	3	6
4	0	9	1
3	8	6	9

| first vectorize the input $x \rightarrow \text{vec}(x) \in \mathbb{R}^{784}$

4 0 9 1
3 8 6 9

feed it to the MLP (with L layers) and predict the labels

$$\text{softmax} \circ W^{\{L\}} \circ \dots \circ \text{ReLU} \circ W^{\{1\}} \text{vect}(x)$$

the model knows nothing about the image structure

we could **shuffle all pixels** and learn an MLP with similar performance

how to **bias** the model, so that it "knows" its input is image?

image is like 2D version of sequence data



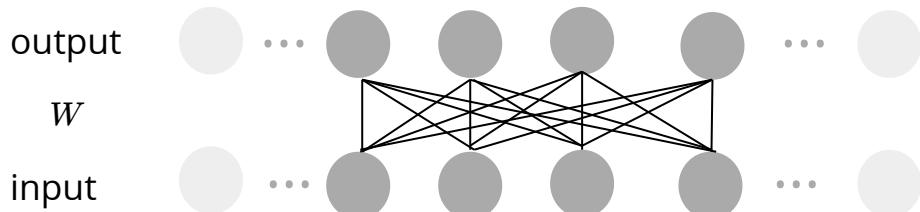
let's find the right model for sequence first!

Parameter-sharing

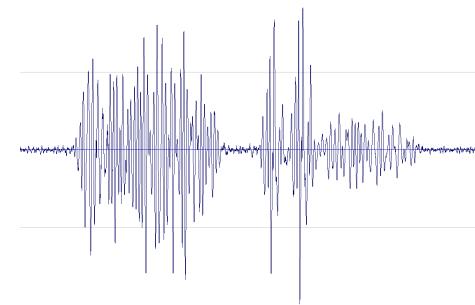
suppose we want to convert one sequence to another $\mathbb{R}^D \rightarrow \mathbb{R}^D$

suppose we have a dataset of input-output pairs $\{(x^{(n)}, y^{(n)})\}_n$

consider only a single layer $y = g(Wx)$

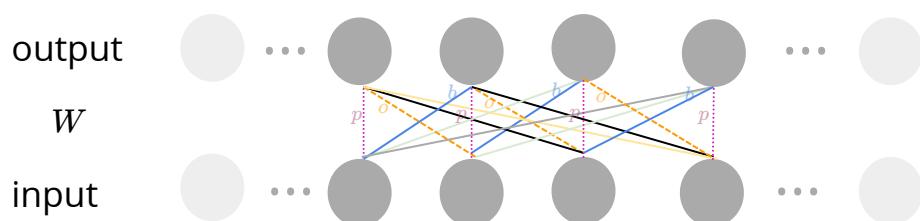


example: remove background noise from audio signal



we may assume, each output unit is the same function shifted along the sequence

when is this a good assumption?



elements of w of the same color are tied together
(parameter-sharing)

Locality & sparse weight

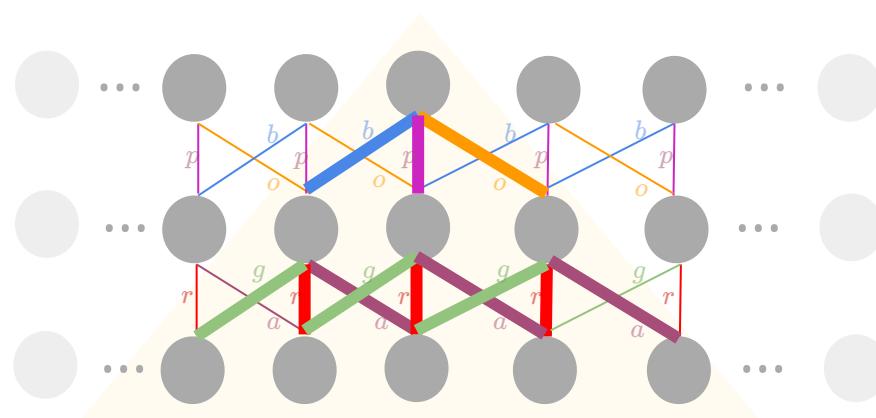
we may further assume each output is a **local** function of input

larger **receptive field** with multiple layers

one layer: the output units "see" 3 neighbouring inputs

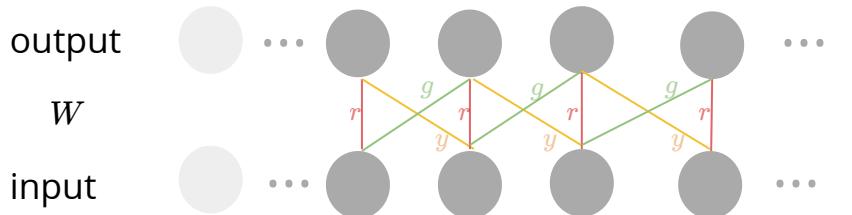
two layers: the output units "see" 5 neighbouring inputs

...



Cross-correlation (1D)

Let's look at the parameter matrix W



r	y			
g	r	y		
g	r	y		
g	r	y		
g	r	y		

parameter-sharing in W
W is very sparse

instead of the whole matrix we can keep the one set of nonzero values

$$w = [w_1, \dots, w_K] = [W_{c,c-\lfloor \frac{K}{2} \rfloor}, \dots, W_{c,c+\lfloor \frac{K}{2} \rfloor}] \longrightarrow$$

r	y			
g	r	y		
g	r	y		
g	r	y		
g	r			

we can write matrix multiplication as **cross-correlation** of w and x

$$y_c = g\left(\sum_{d=1}^D W_{c,d} x_d\right) = g\left(\sum_{k=1}^K w_k x_{c-\lfloor \frac{K}{2} \rfloor + k}\right)$$

feedforward layer: slide on the input, calculate inner product and apply the nonlinearity

Convolution (1D)

Cross-correlation is similar to convolution

Cross-correlation $y_d = \sum_{k=-\infty}^{\infty} w_k x_{d+k}$

$w \star x$

w is called the filter or kernel

ignoring the activation (for simpler notation)

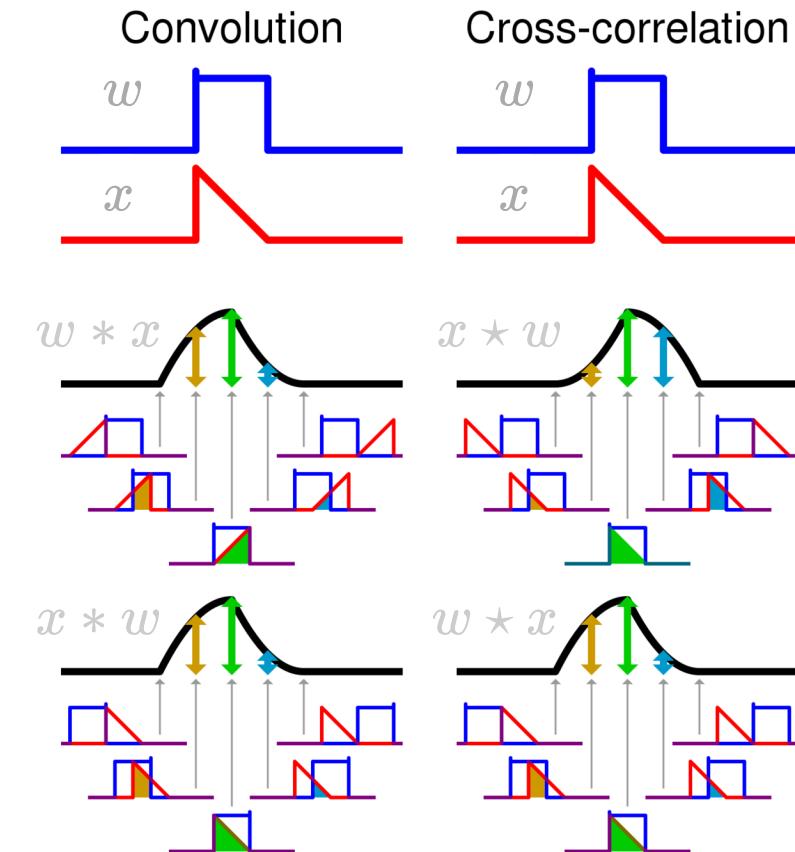
assuming w and x are zero for any index outside the input and filter bound

Convolution flips w or x (to be **commutative**)

$$y_d = \sum_{k=-\infty}^{\infty} w_k x_{d-k} = \sum_{k'=-\infty}^{\infty} w_{d-k'} x_{k'}$$

$w * x$ change of variable
 $k' = d - k$ $x * w$

since we **learn** w, flipping it makes no difference
 in practice, we use cross correlation rather than convolution
 convolution is **equivariant** wrt translation
 -- i.e., shifting x, shifts $w * x$



Convolution (1D)

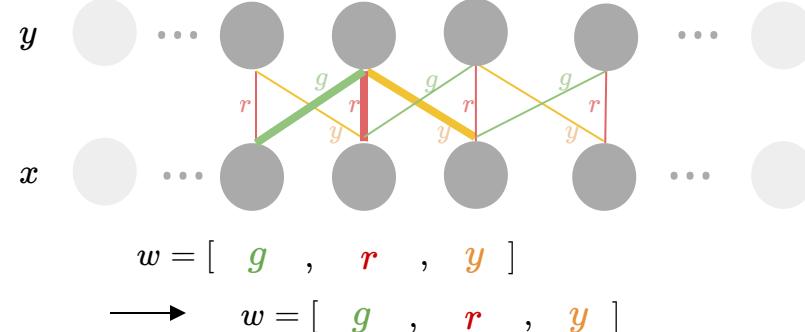
1D convolution layer **so far...**

$$y_d = \sum_{k=1}^K w_k x_{d+k-1}$$

-1 is because the indexing starts from 1
for d=1,k=1 we index the first element of x



```
1 def Conv1D(  
2     x, # D (length)  
3     w, # K (filter length)  
4 ):  
5  
6     D, = x.shape  
7     K, = w.shape  
8     Dp = D - K + 1 #output length  
9     y = np.zeros(Dp)  
10    for dp in range(Dp):  
11        y[dp] = np.sum(x[dp:dp+K] * w)  
12    return y
```



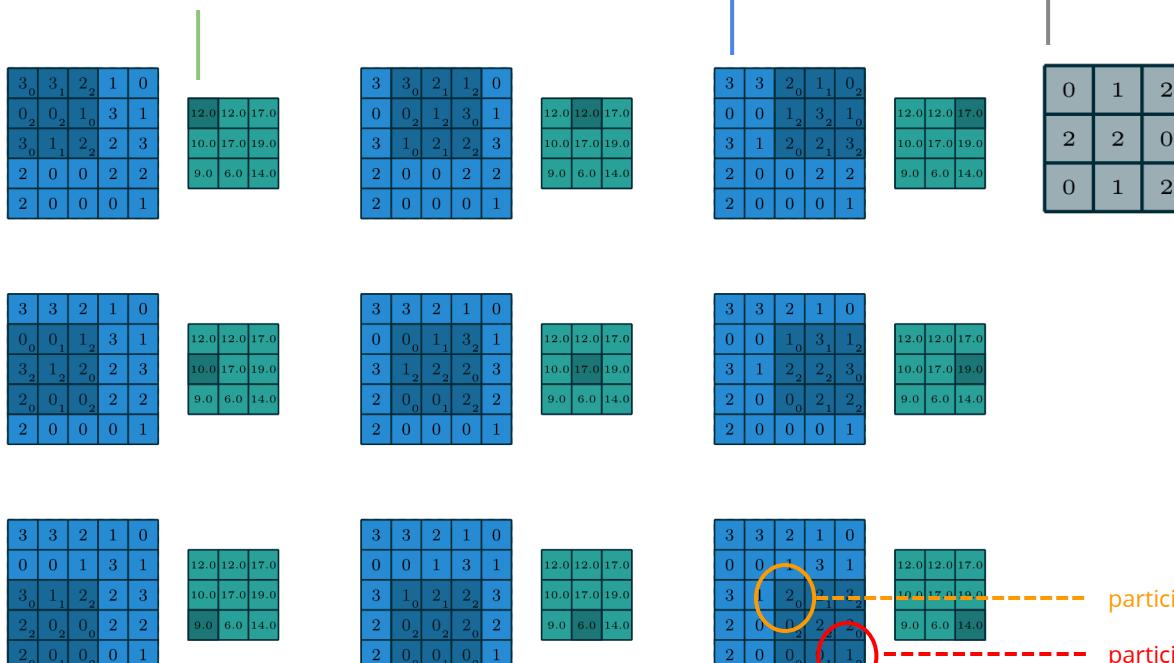
Example:

Input	Kernel	Output															
<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td></tr></table>	0	1	2	3	4	5	6	<table border="1"><tr><td>1</td><td>2</td></tr></table>	1	2	<table border="1"><tr><td>2</td><td>5</td><td>8</td><td>11</td><td>14</td><td>17</td></tr></table>	2	5	8	11	14	17
0	1	2	3	4	5	6											
1	2																
2	5	8	11	14	17												

Convolution (2D)

similar idea of parameter-sharing and locality extends to 2 dimension (*i.e. image data*)

$$y_{d_1, d_2} = \sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} x_{d_1+k_1-1, d_2+k_2-1} w_{k_1, k_2}$$



Convolution (2D)

there are different ways of handling the borders

zero-pad the input, and produce all non-zero outputs (**full**)

the output is larger than the input

output length (for one dimension)

$$D + 2 \times \text{padding} - K + 1$$

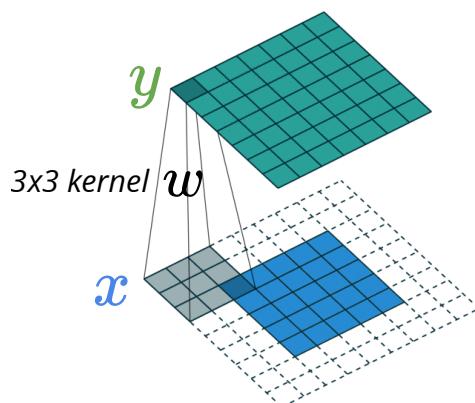
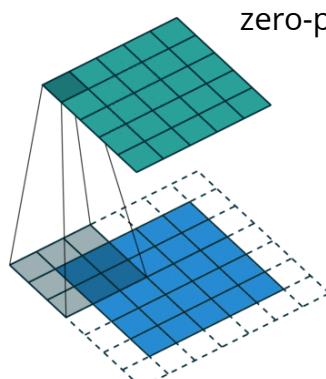
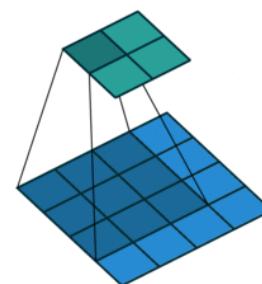


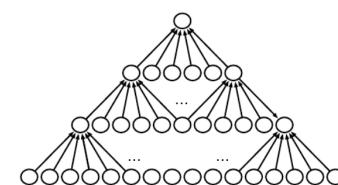
image credit: Vincent Dumoulin, Francesco Visin



zero-pad the input, to keep the output dims similar to input (**same**)



no padding at all (**valid**)
the output is smaller than the input
we can't stack many layers
sometimes we need to maintain the width



Pooling

sometimes we would like to reduce the size of output e.g., from $D \times D$ to $D/2 \times D/2$

a combination of pooling and downsampling is used

1. calculate the output $\tilde{y}_d = g\left(\sum_{k=1}^K x_{d+k-1} w_k\right)$

2. aggregate the output over different regions

$$y_d = \text{pool}\{\tilde{y}_d, \dots, \tilde{y}_{d+p}\}$$

two common aggregation functions are **max** and **mean**

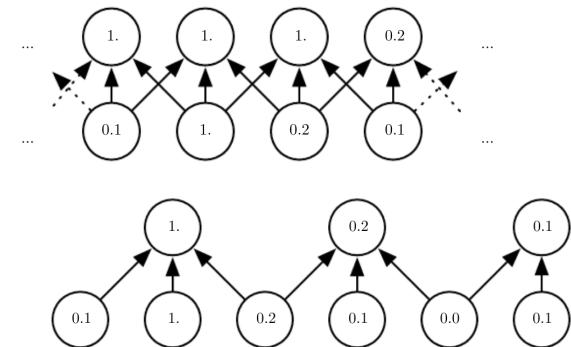
3. often this is followed by subsampling using the same step size

the same idea extends to higher dimensions

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

$\xrightarrow{2 \times 2 \text{ Max-Pool}}$

20	30
112	37



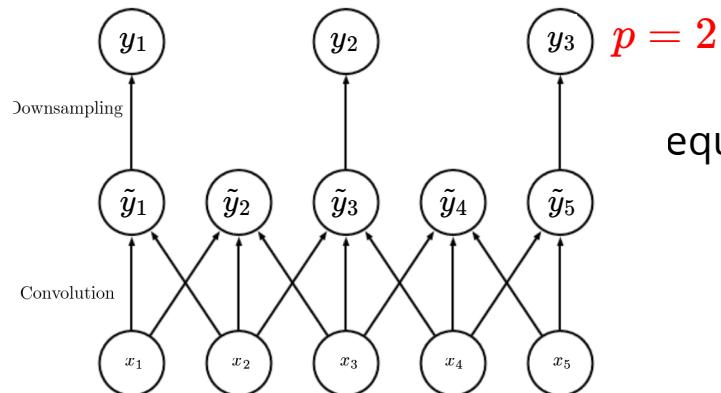
Strided convolution

alternatively we can directly subsample the output

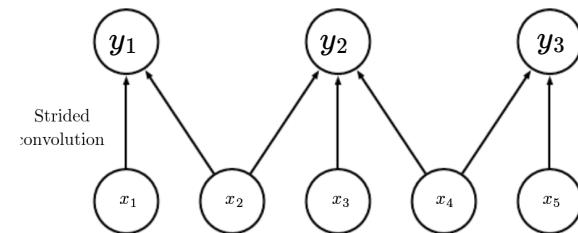
$$\tilde{y}_d = g\left(\sum_{k=1}^K x_{(d-1)+k} w_k\right)$$

$$y_d = \tilde{y}_{\textcolor{red}{p}(d-1)+1}$$

$$\tilde{y}_d = g\left(\sum_{k=1}^K x_{\textcolor{red}{p}(d-1)+k} w_k\right)$$



equivalent to

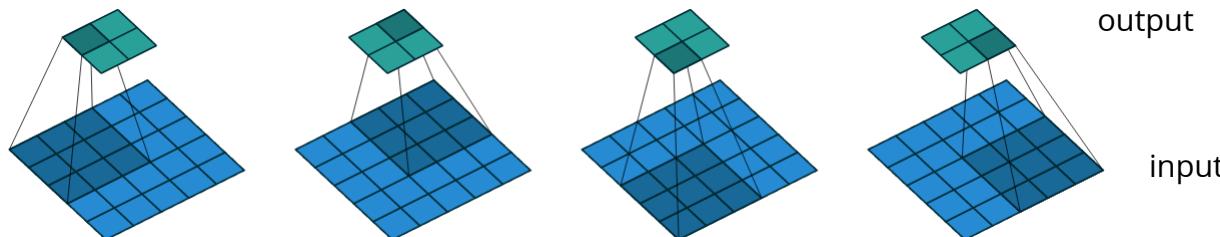


Strided convolution

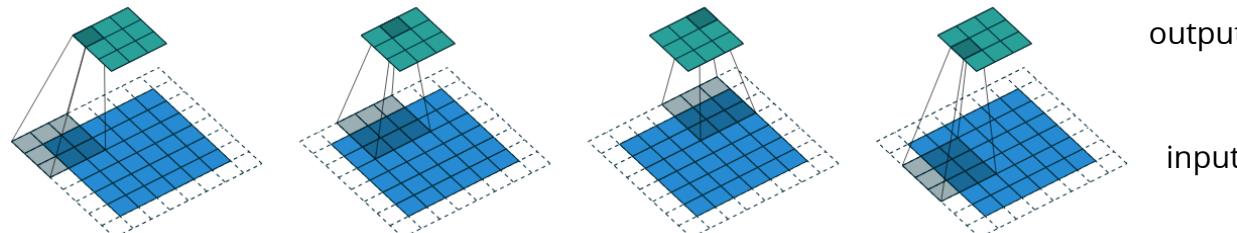
the same idea extends to higher dimensions

$$y_{d_1, d_2} = \sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} x_{p_1(d_1-1)+k_1, p_2(d_2-1)+k_2} w_{k_1, k_2}$$

different strides for different dimensions



with padding

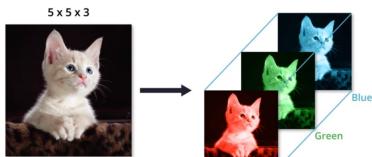


$$\left\lfloor \frac{D+2 \times \text{padding} - K}{\text{stride}} + 1 \right\rfloor$$

output length (for one dimension)

Channels

so far we assumed a single input and output sequence or image

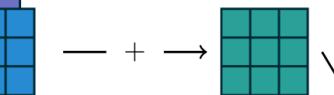
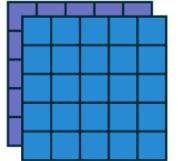


we have one $K_1 \times K_2$ filters per input-output channel combination $w \in \mathbb{R}^{M \times M' \times K_1 \times K_2}$

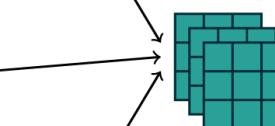
with RGB data, we have 3 input channels ($M = 3$)

this example: 2 input channels

$$x \in \mathbb{R}^{M \times D_1 \times D_2}$$

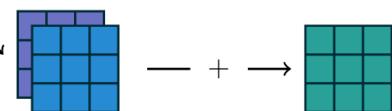


+ add the result of convolution from different input channels



similarly we can produce multiple output channels $M' = 3$

$$y \in \mathbb{R}^{M' \times D'_1 \times D'_2}$$



Channels

we can also add a *bias parameter (b)*, one per each output channel

$$y_{\mathbf{m}', d_1, d_2} = g \left(\sum_{m=1}^M \sum_{k_1} \sum_{k_2} w_{\mathbf{m}, \mathbf{m}', k_1, k_2} x_{\mathbf{m}, d_1+k_1-1, d_2+k_2-1} + b_{\mathbf{m}'} \right)$$

$$y \in \mathbb{R}^{M' \times D'_1 \times D'_2}$$

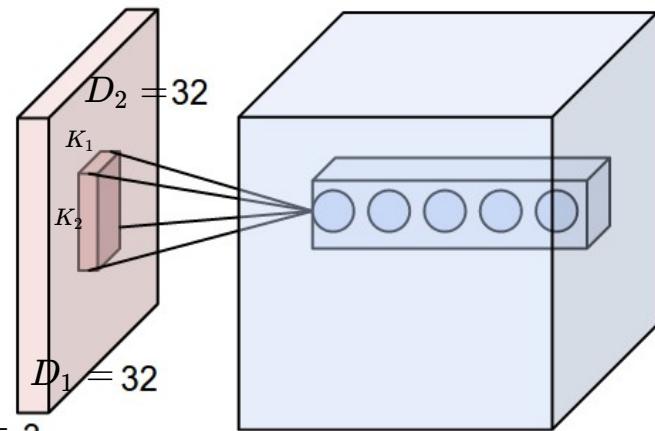
$$w \in \mathbb{R}^{M \times M' \times K_1 \times K_2}$$

$$x \in \mathbb{R}^{M \times D_1 \times D_2}$$

$$b \in \mathbb{R}^{M'}$$

$$\begin{matrix} M = 3 \\ \text{RGB channels} \end{matrix}$$

$$\begin{matrix} M' = 5 \\ \text{ } \end{matrix}$$



Example

<https://cs231n.github.io/assets/conv-demo/>

⋮

demo from: <https://cs231n.github.io/convolutional-networks/>

Convolutional Neural Network (CNN)

CNN or convnet is a neural network with convolutional layers

it could be applied to 1D sequence, 2D image or 3D volumetric data

example: conv-net architecture (LeNet, 1998) for digits recognition

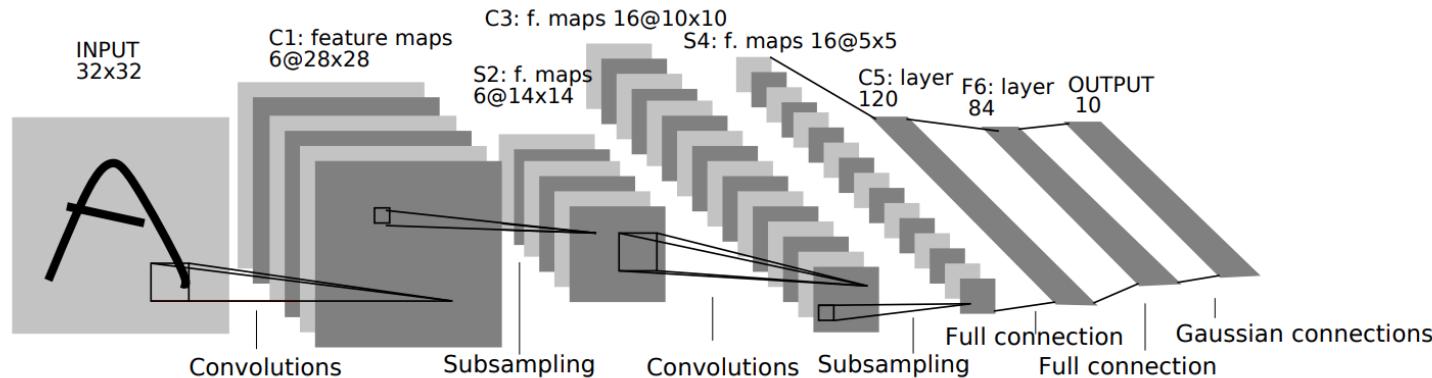
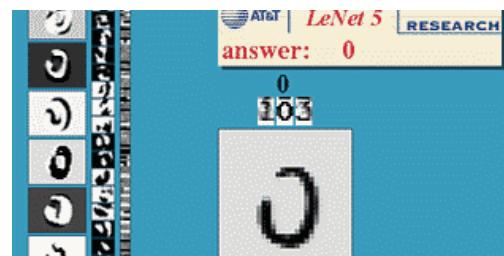


image from [LeNet paper](#)



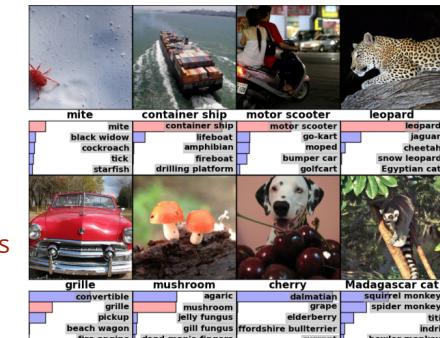
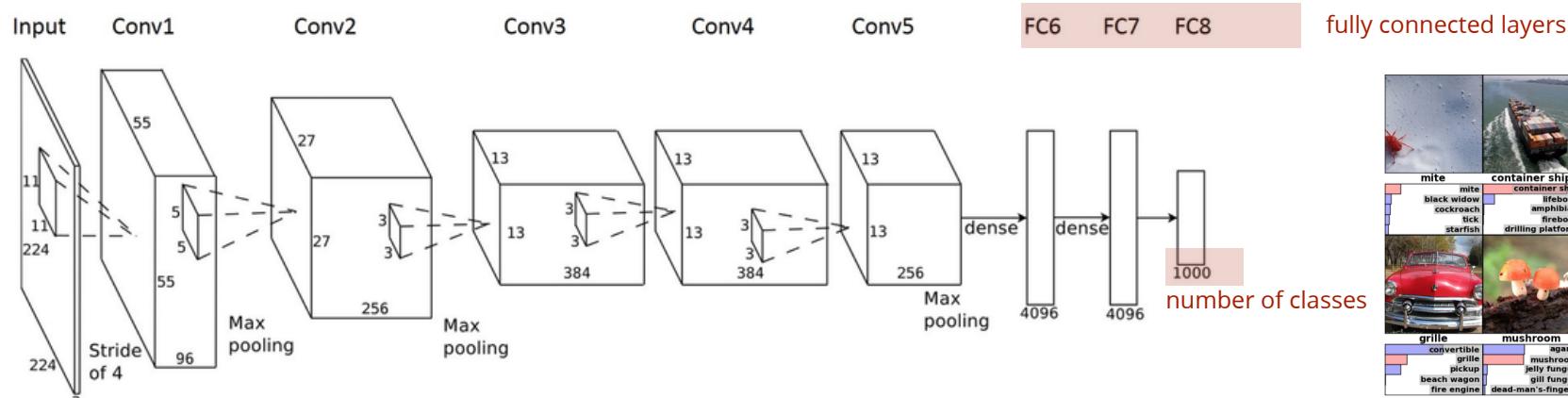
very accurate to be used in large scale in postal services (zip code recognition) and banks (cheques)

Convolutional Neural Network (CNN)

CNN or convnet is a neural network with convolutional layers

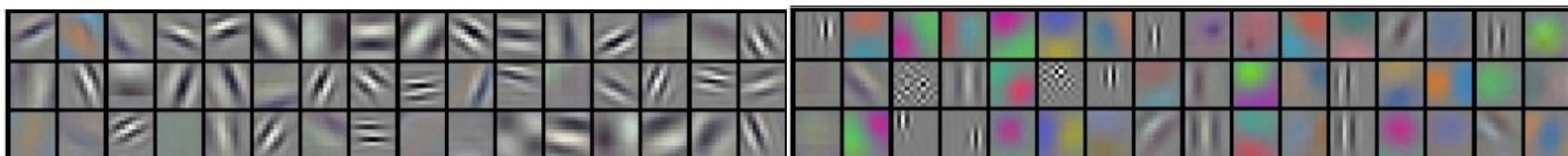
it could be applied to 1D sequence, 2D image or 3D volumetric data

example: conv-net architecture (derived from AlexNet) for image classification



read the paper [here](#)

↙ visualization of the convolution kernel at the first layer 11x11x3x96
96 filters, each one is 11x11x3. each of these is responsible for one of 96 feature maps in the second layer

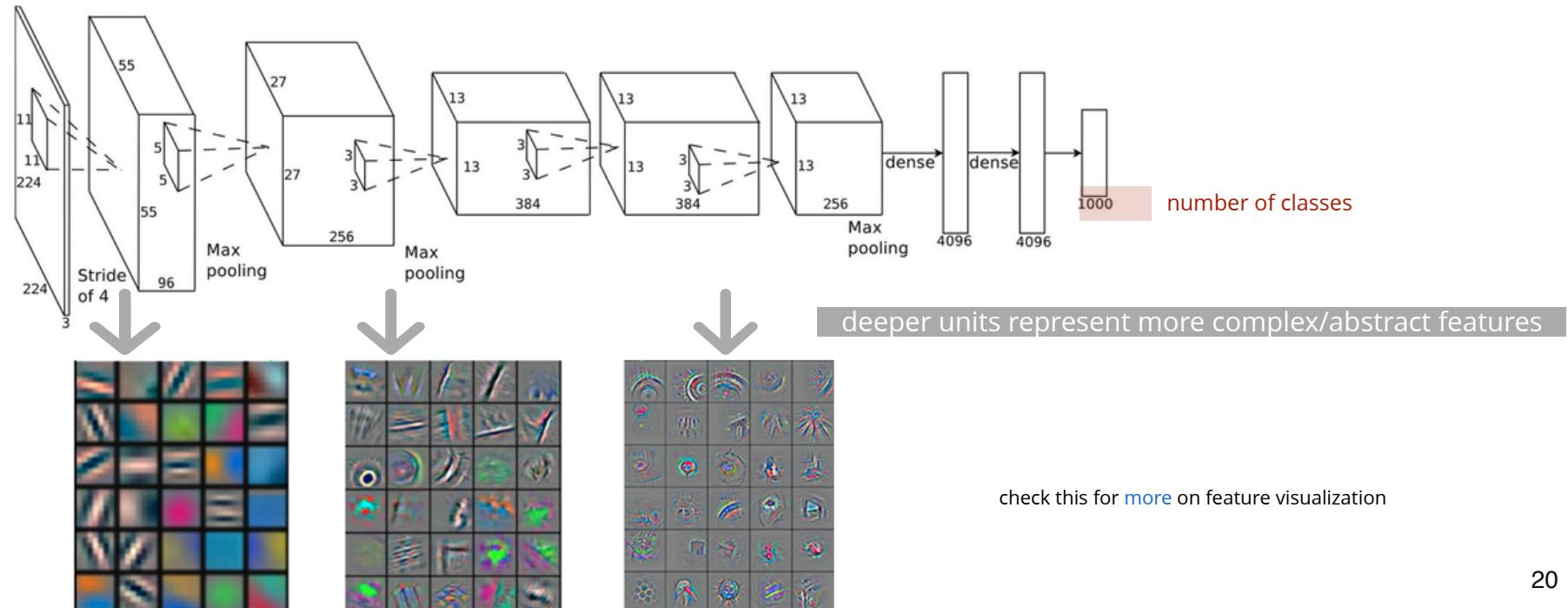


Convolutional Neural Network (CNN)

CNN or convnet is a neural network with convolutional layers (so it's a special type of MLP)
it could be applied to 1D sequence, 2D image or 3D volumetric data

example: conv-net architecture (derived from AlexNet) for image classification

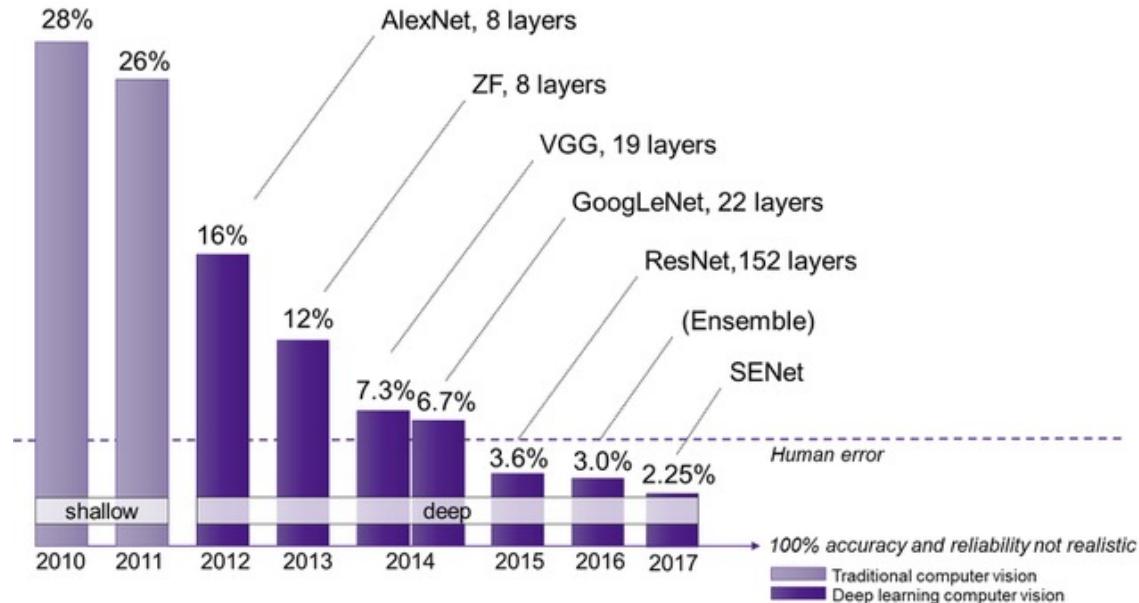
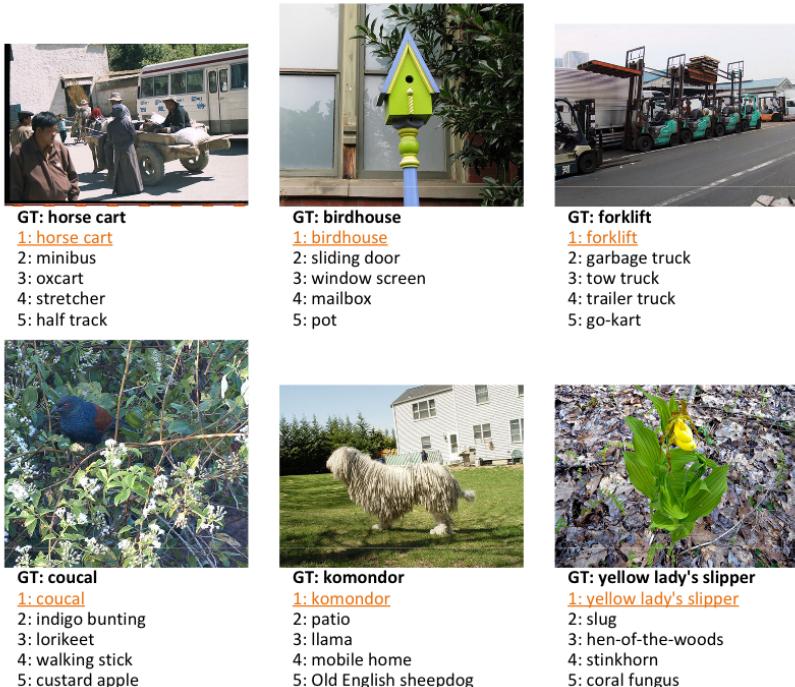
Input Conv1 Conv2 Conv3 Conv4 Conv5 FC6 FC7 FC8 fully connected layers



Application: image classification

Convnets have achieved super-human performance in image classification

ImageNet challenge: > 1M images, 1000 classes

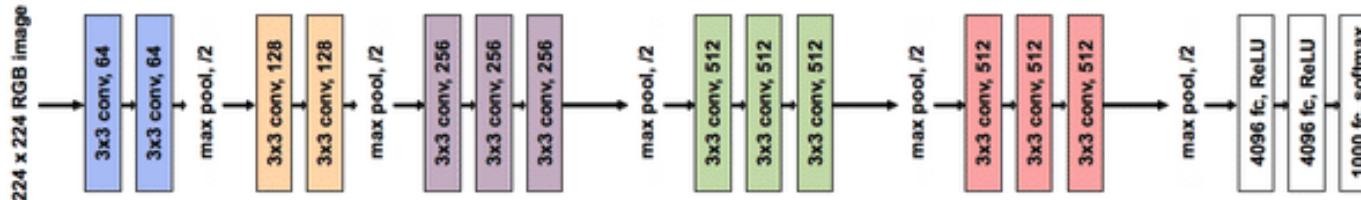


Application: image classification

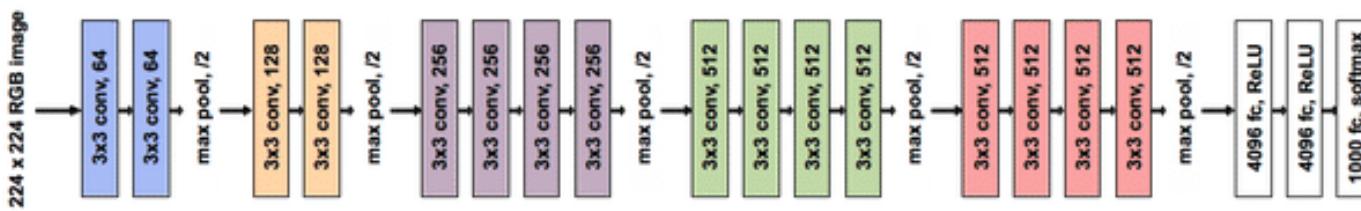
optional

variety of increasingly deeper architectures have been proposed

VGG16



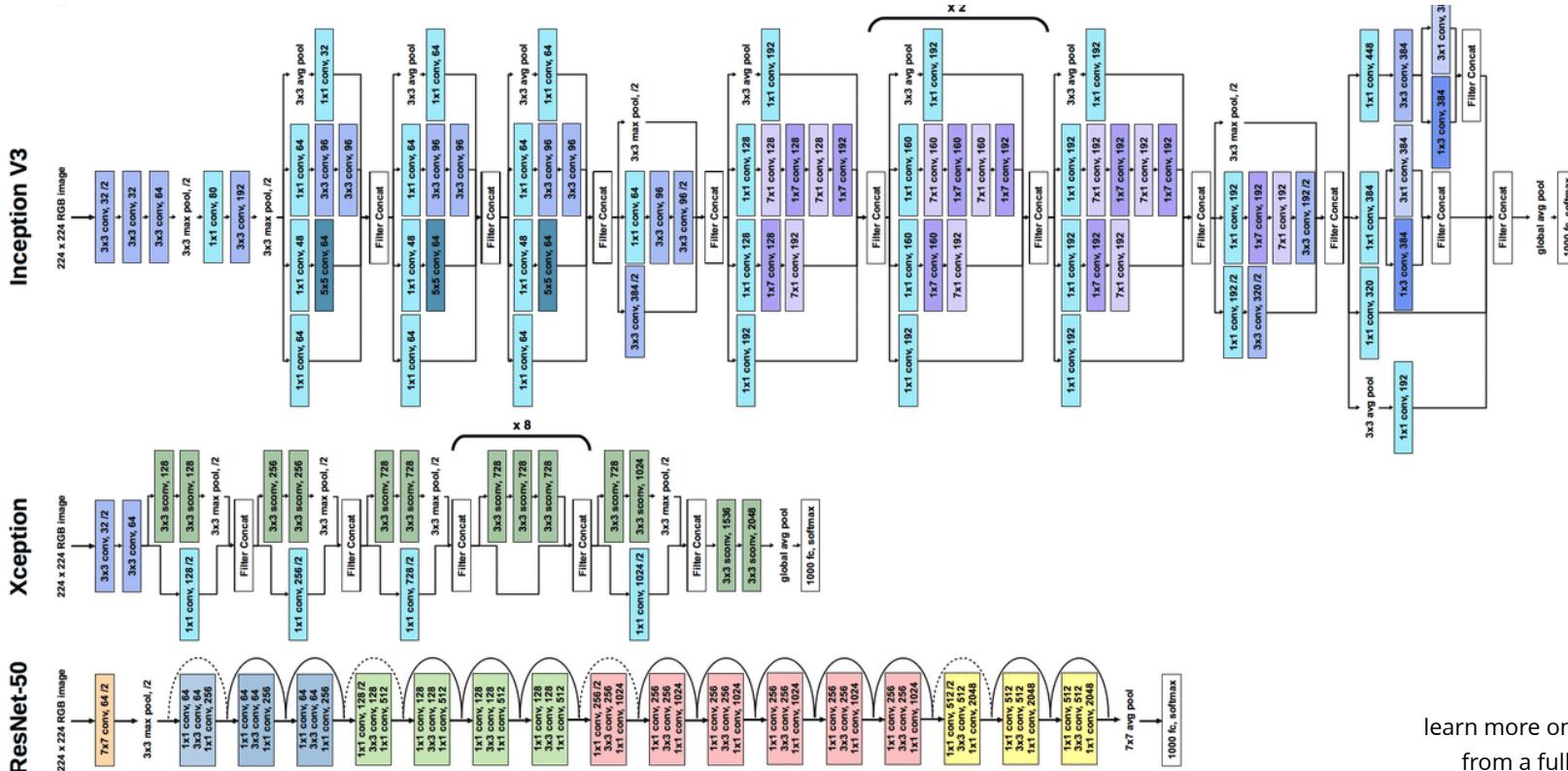
VGG19



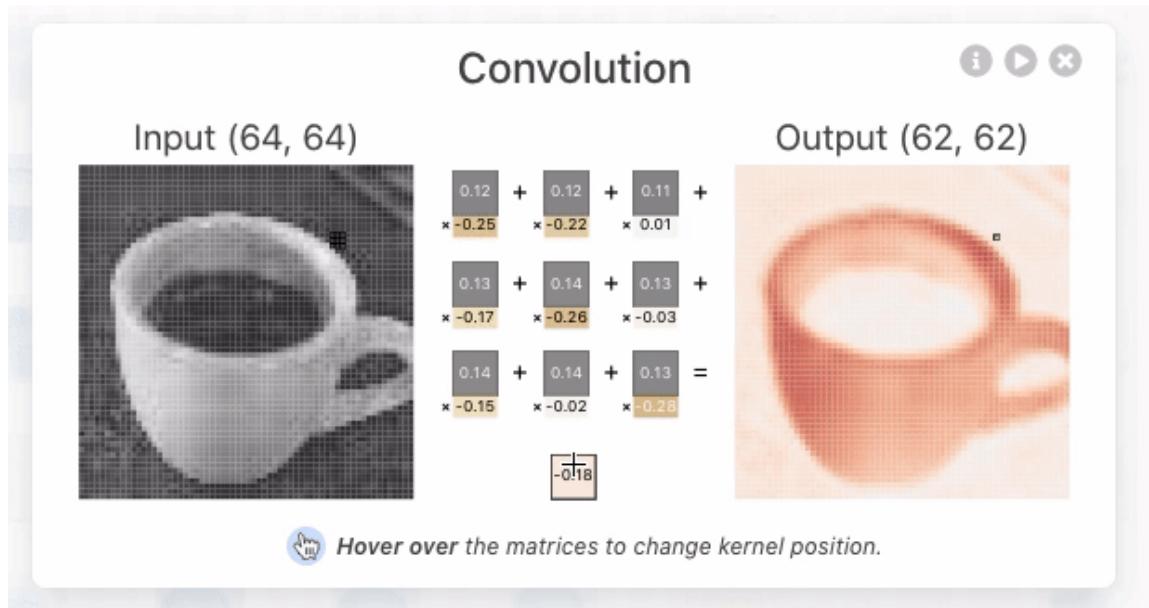
Application: image classification

variety of increasingly deeper architectures have been proposed

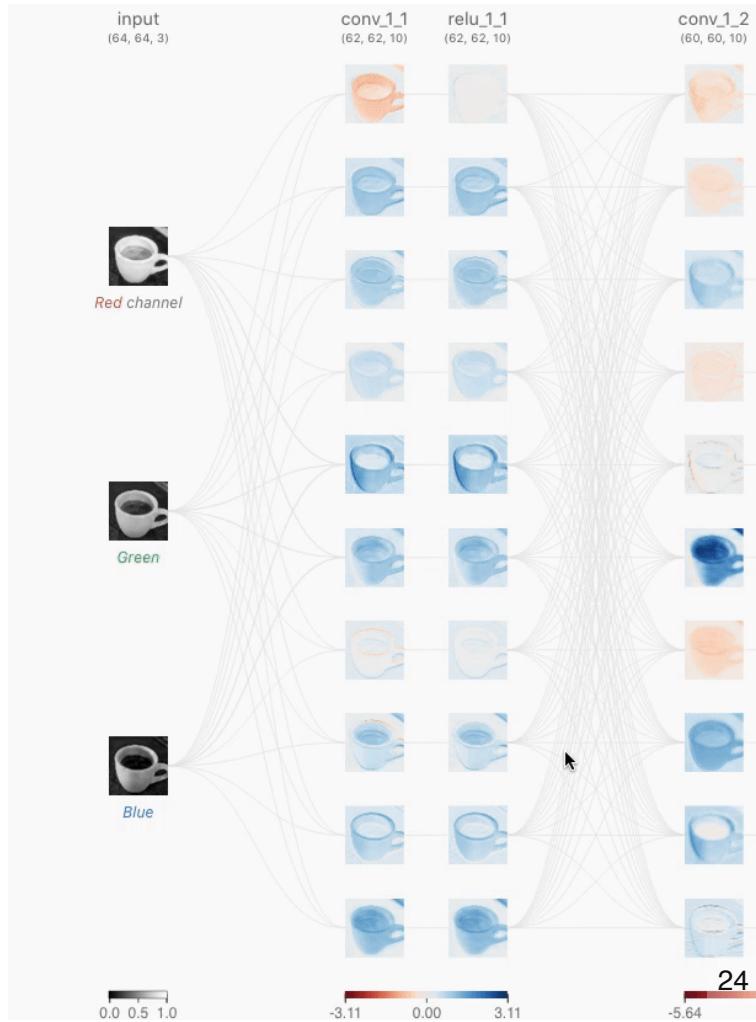
optional



Visual Examples



see the interactive [demo here](#)



Training: backpropagation through convolution

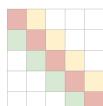
consider the 1D convolution op. $y_d = \sum_k w_k x_{d+k-1}$

using backprop. we have $\frac{\partial J}{\partial y_d}$ so far and we need

$$1) \frac{\partial J}{\partial w_k} = \sum_{d'} \frac{\partial J}{\partial y_{d'}} \frac{\partial y_{d'}}{\partial w_k}$$

using this we can update the convolution kernel at the current layer

$$2) \text{ to backpropagate to previous layer } \frac{\partial J}{\partial x_d} = \sum_{d'} \frac{\partial J}{\partial y_{d'}} \frac{\partial y_{d'}}{\partial x_d}$$



even when we have stride, and padding, this operation is similar to multiplication by transpose of the parameter-sharing matrix (**transposed convolution**)

Transposed convolution

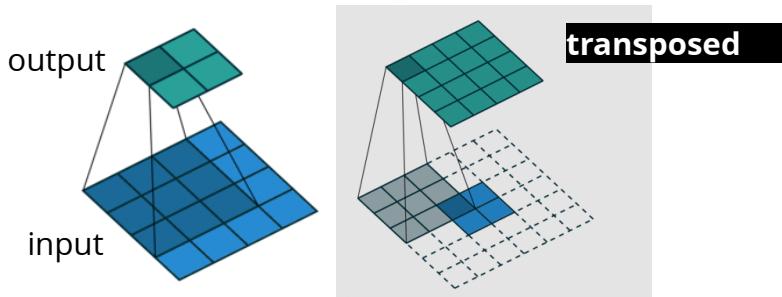
Transposed convolution produces a larger output from a smaller input

Example:

Input	Kernel	$=$					Output
$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$\begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix}$	$+ \quad \begin{matrix} 0 & 1 \\ 2 & 3 \end{matrix}$	$+ \quad \begin{matrix} 0 & 2 \\ 4 & 6 \end{matrix}$	$+ \quad \begin{matrix} 0 \\ 6 \end{matrix}$	$= \quad \begin{matrix} 0 & 0 & 1 \\ 0 & 4 & 6 \\ 4 & 12 & 9 \end{matrix}$	

Transposed convolution can recover the shape of the original input

no padding of the original convolution corresponds to *full* padding of in transposed version

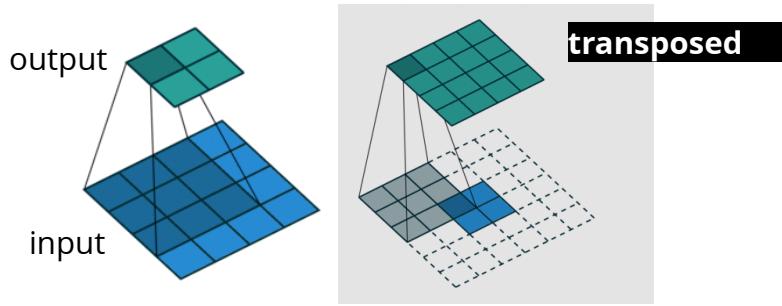


optional

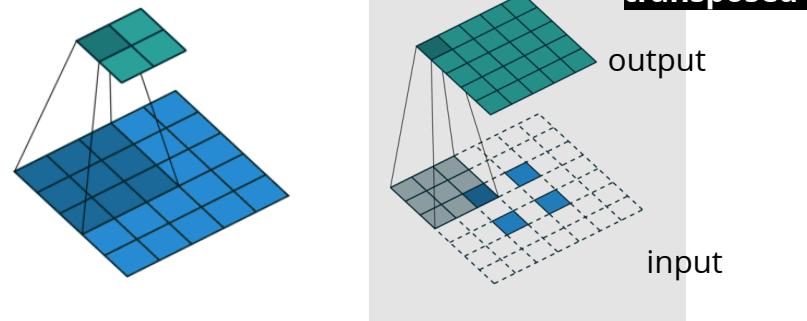
Transposed convolution

Transposed convolution recovers the shape of the original input

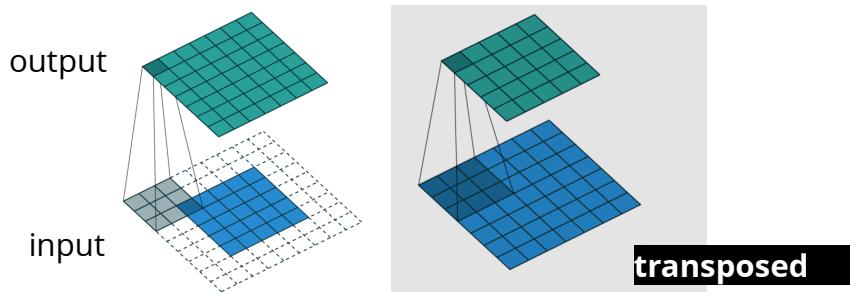
no padding of the original convolution corresponds to *full* padding of in transposed version



Convolution **with stride** and its transpose



full padding of the original convolution corresponds to no padding of in transposed version



this can be used for up-sampling (opposite of stride/pooling)

as expected the transpose of a transposed convolution is the original convolution

Solving other discriminative vision tasks with CNNs

Structured Prediction: the output itself may have (image) structure (e.g., predicting text, audio, image)

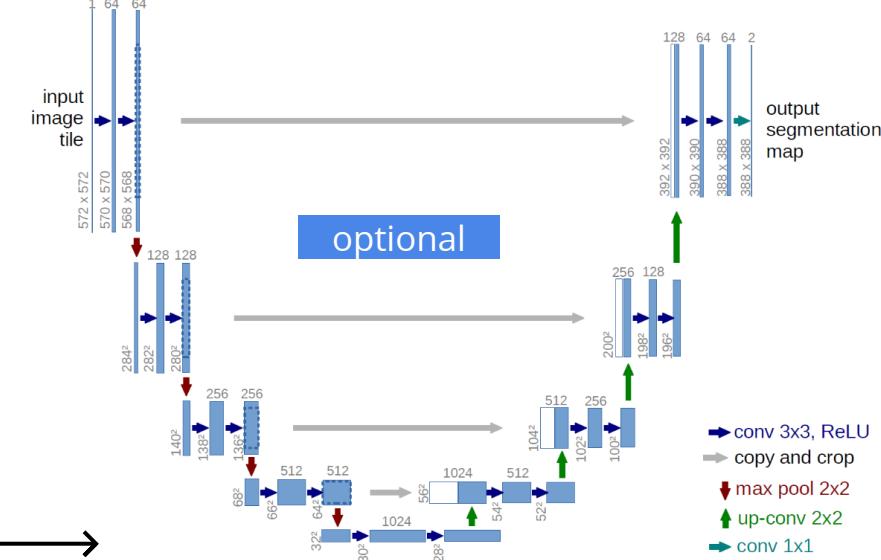
example

in (semantic) segmentation, we classify each pixel
loss is the sum of cross-entropy loss across the whole image



variety of architectures... one that performs well is **U-Net** →

transposed convolution (upconv), concatenation, and skip connection are common in architecture design
architecture search (i.e., combinatorial hyper-parameter search) is an expensive process and an active research area



Generating images by inverting CNNs

generating images which maximize the class label

$$p(x|y) \propto p(x)p(y|x)$$

CNN

$$x_{t+1} = x_t + \epsilon_1 \frac{\partial \log p(x_t)}{\partial x_t} + \epsilon_2 \frac{\partial \log p(y=c|x_t)}{\partial x_t} + \mathcal{N}(0, \epsilon_3^2 \mathbf{I})$$

e.g. using **Gaussian prior** we have:

$$x_{t+1} = (1 - \epsilon_1)x_t + \frac{\partial \log p(y=c|x_t)}{\partial x_t}$$

gradients

Images that maximize the probability of ImageNet classes “goose” and “ostrich”

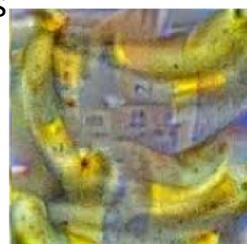


e.g. using **Total variation (TV) prior** gives more realistic

images



Anemone Fish



Banana



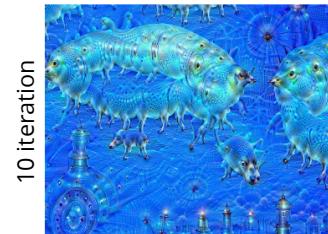
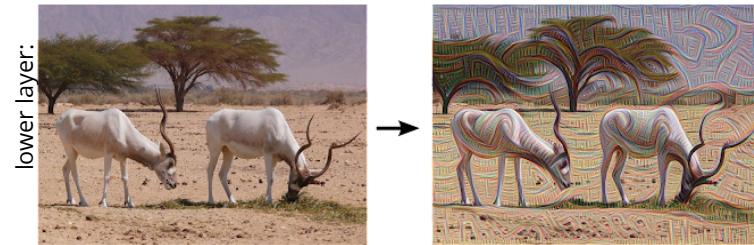
Parachute



Screw

Generating images by inverting CNNs

- Deep Dream
 - *generate versions of an input image that emphasize certain features by picking a layer and ask the network to enhance whatever it detected*



[read more here](#)

- Neural style transfer
 - *specify a reference “style image” xs and “content image” xc .*



Summary

convolution layer introduces an **inductive bias** (equivariance) to MLP

- translation of the same model is applied to produce different outputs (pixels)
- the layer is equivariant to **translation**
- achieved through **parameter-sharing**

conv-nets use combinations of

- convolution layers
- ReLU (or similar) activations
- pooling and/or stride for down-sampling
- skip-connection and/or batch-norm to help with optimization / regularization
- potentially fully connected layers in the end

training

- backpropagation (similar to MLP)
- SGD or its improved variations with adaptive learning rate
- monitor the validation error for early stopping