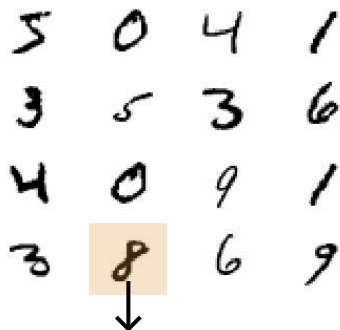# Applied Machine Learning

Convolutional Neural Networks

**Reihaneh Rabbany**

# 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



we can apply an MLP to image data

first vectorize the input $\quad x \to \mathrm{vec}(x) \in \mathbb{R}^{784}$

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

$$\mathrm{softmax} \circ W^{\{L\}} \circ \ldots \circ \mathrm{ReLU} \circ W^{\{1\}}\mathrm{vect}(x)$$

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!

4

# Parameter-sharing

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

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

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

output

$W$

input

**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?

output

$W$

input

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

5

# 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



output

$W$

input

parameter-sharing in W
W is very sparse

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

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

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

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

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

# Convolution (1D)

Cross-correlation is similar to convolution

**Cross-correlation** $\quad y_d = \sum_{k=-\infty}^{\infty} w_k x_{d+k} \qquad w \star x$

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)

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

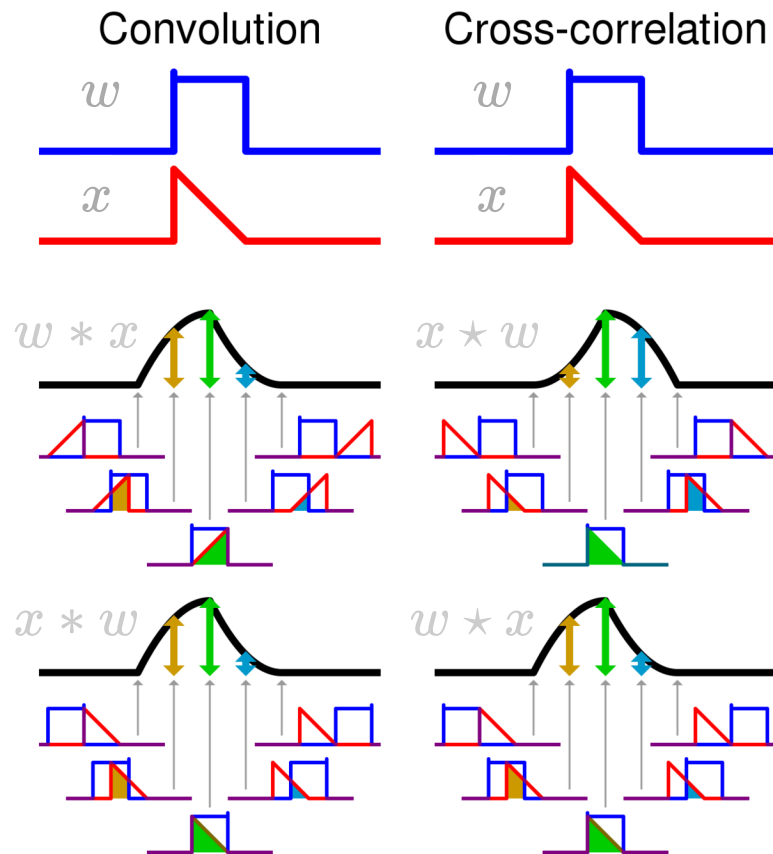$w * x$ $\quad$ *change of variable* $\quad x * w$
$k' = d - k$

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    Cross-correlation

$w$ $\qquad$ $w$

$x$ $\qquad$ $x$

$w * x$ $\qquad$ $x \star w$

$x * w$ $\qquad$ $w \star x$

# Convolution (1D)

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

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

```python
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
```



$w = [\ g\ ,\ r\ ,\ y\ ]$

$\longrightarrow\quad w = [\ g\ ,\ r\ ,\ y\ ]$

# Convolution (2D)

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

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



image credit: Vincent
Dumoulin, Francesco Visin

participates in all outputs

participates in a single output

this is related to the borders

# 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

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

$y$

*3x3 kernel* $w$

$x$

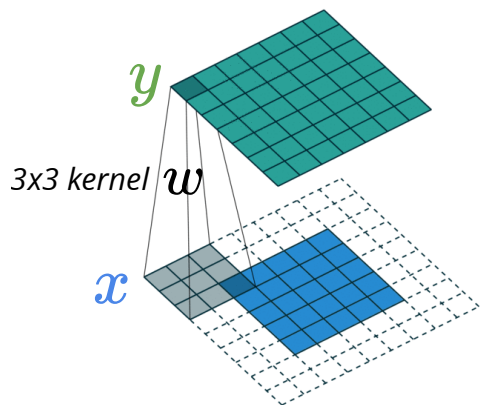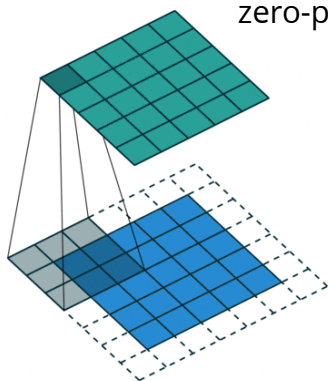image credit: Vincent Dumoulin, Francesco Visin

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

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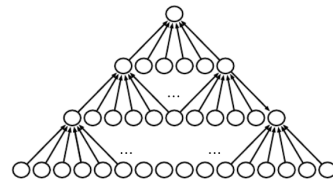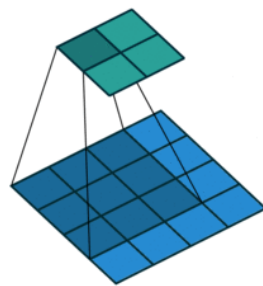the output is smaller than the input
we can't stack many layers
sometimes we need to maintain the width

# Pooling

sometimes we would like to reduce the size of output  *e.g., from D x D to D/2 x 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

$$y_d = \mathrm{pool}\{\tilde{y}_d, \ldots, \tilde{y}_{d+p}\}$$

*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 |

$2 \times 2$ Max-Pool $\longrightarrow$

| 20 | 30 |
|-----|----|
| 112 | 37 |

# Strided convolution

alternatively we can directly subsample the output

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

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

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



$p = 2$

Downsampling

Convolution

equivalent to →

Strided convolution

# Strided convolution
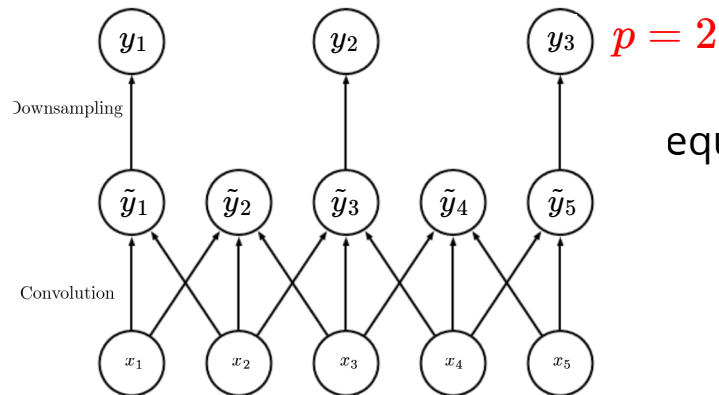
the same idea extends to higher dimensions

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

*different strides for different dimensions*



output

input

with padding



output

input

output length (for one dimension)

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

15

image: Dumoulin & Visin'16

# Channels

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



5 x 5 x 3

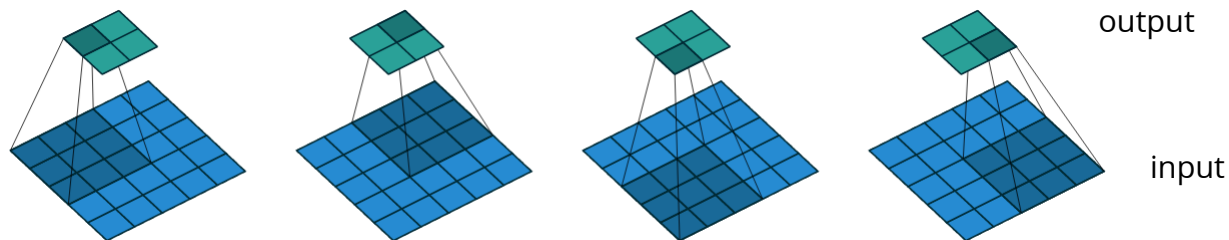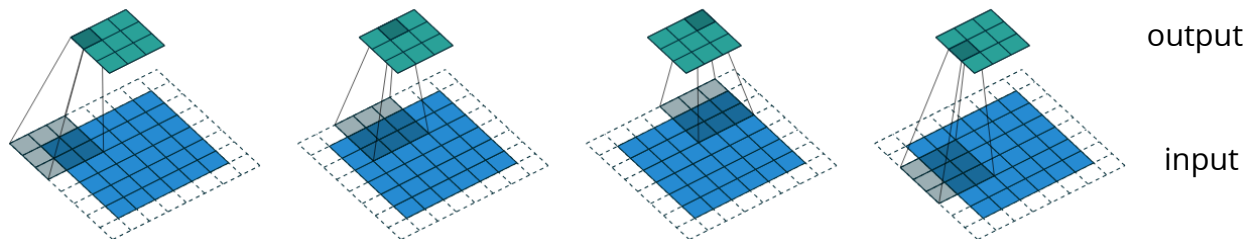we have one $K_1 \times K_2$ filters per input-output channel combination    $w \in \mathbb{R}^{M \times M' \times K_1 \times K_2}$

**+** add the result of convolution from different input channels

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

this example: 2 input channels

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

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

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

17

image: Dumoulin & Visin'16

# Channels

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

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

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

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

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

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



$D_2 = 32$

$K_1$

$K_2$

$D_1 = 32$

$M = 3$

RGB channels

$M' = 5$

18

# Example

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

# 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



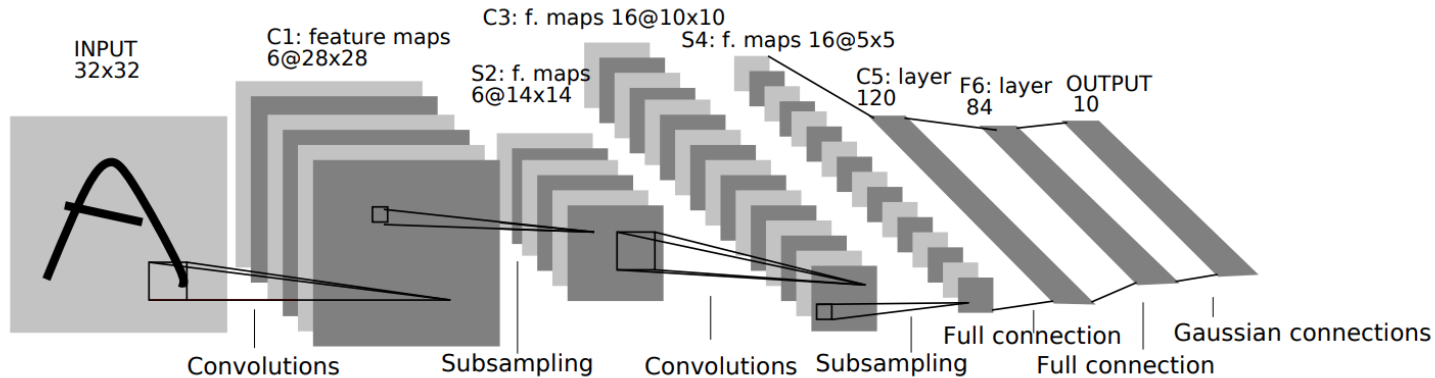very accurate to be used in large scale in postal services (zip code recognition) and banks (cheques)

# 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



fully connected layers

number of classes

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

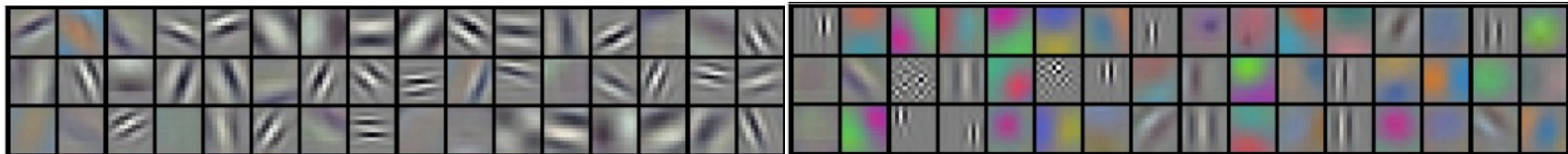**example:** conv-net architecture (derived from AlexNet) for image classification



fully connected layers

number of classes

deeper units represent more complex/abstract features

check this for more on feature visualization

# Application: image classification

Convnets have achieved super-human performance in image classification

ImageNet challenge: > 1M images, 1000 classes



GT: horse cart
1: horse cart
2: minibus
3: oxcart
4: stretcher
5: half track

GT: birdhouse
1: birdhouse
2: sliding door
3: window screen
4: mailbox
5: pot

GT: forklift
1: forklift
2: garbage truck
3: tow truck
4: trailer truck
5: go-kart
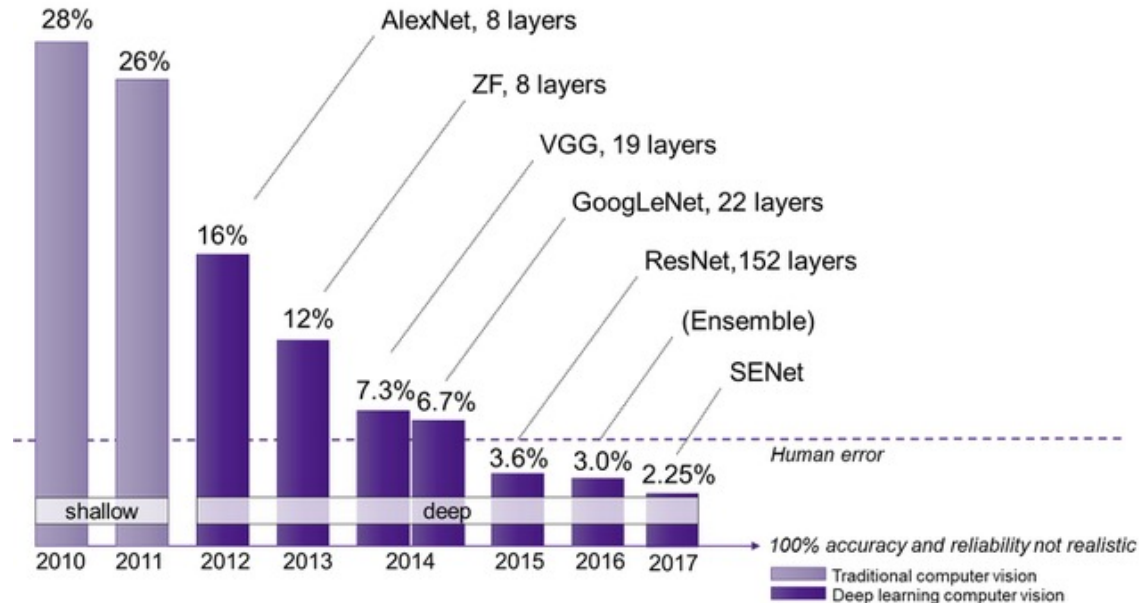
GT: coucal
1: coucal
2: indigo bunting
3: lorikeet
4: walking stick
5: custard apple

GT: komondor
1: komondor
2: patio
3: llama
4: mobile home
5: Old English sheepdog

GT: yellow lady's slipper
1: yellow lady's slipper
2: slug
3: hen-of-the-woods
4: stinkhorn
5: coral fungus



24

image credit: He et al'15, https://semiengineering.com/new-vision-technologies-for-real-world-applications/

# Application: image classification

variety of increasingly deeper architectures have been proposed

image credit: He et al'15, https://semiengineering.com/new-vision-technologies-for-real-world-applications/

# Application: image classification

variety of increasingly deeper architectures have been proposed

image credit: He et al'15, https://semiengineering.com/new-vision-technologies-for-real-world-applications/

learn more on different CNN here

from a full course on CNN:
http://cs231n.stanford.edu/

26

# Visual Examples
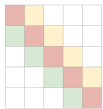


see the interactive demo here

# Training: backpropagation through convolution

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

using backprop. we have $\dfrac{\partial J}{\partial y_{d'}}$ so far and we need

1) $\dfrac{\partial J}{\partial w_k} = \sum_{d'} \dfrac{\partial J}{\partial y_{d'}} \boxed{\dfrac{\partial y_{d'}}{\partial w_k}}^{x_{d'+k-1}}$ using this we can update the convolution kernel at the current layer

2) to backpropagate to previous layer $\dfrac{\partial J}{\partial x_d} = \sum_{d'} \dfrac{\partial J}{\partial y_{d'}} \boxed{\dfrac{\partial y_{d'}}{\partial x_d}}^{w_{d-d'+1}}$
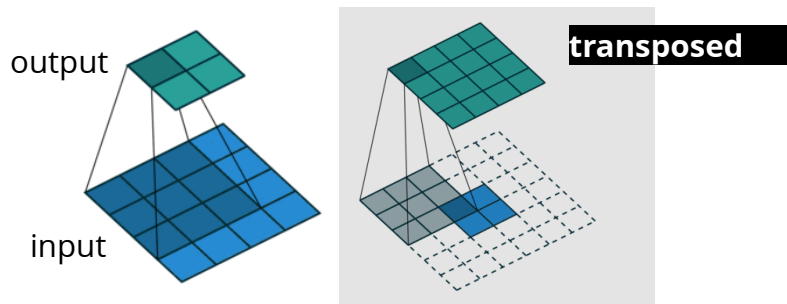
even when we have stride, and padding, this operation is similar to multiplication by transpose of the parameter-sharing matrix **(transposed convolution)**
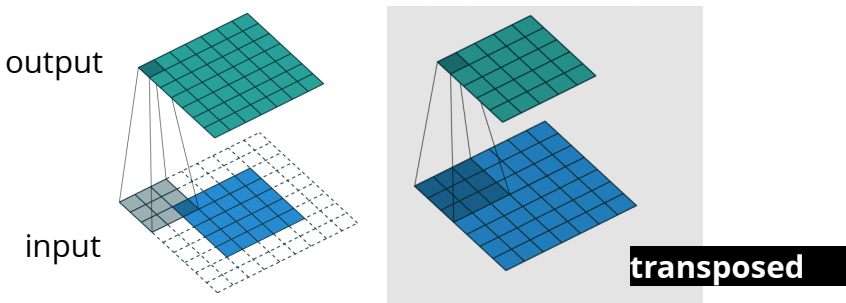
# Transposed convolution

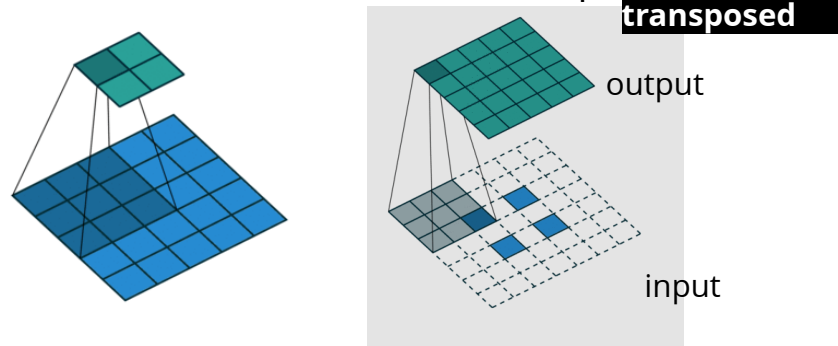Transposed convolution recovers the shape of the original input

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

output

**transposed**

input

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

output

input

**transposed**

Convolution **with stride** and its transpose

**transposed**

output

input

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

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

image: Dumoulin & Visin'16

30

# Structured Prediction

the output itself may have (image) structure (e.g., predicting text, audio, image)
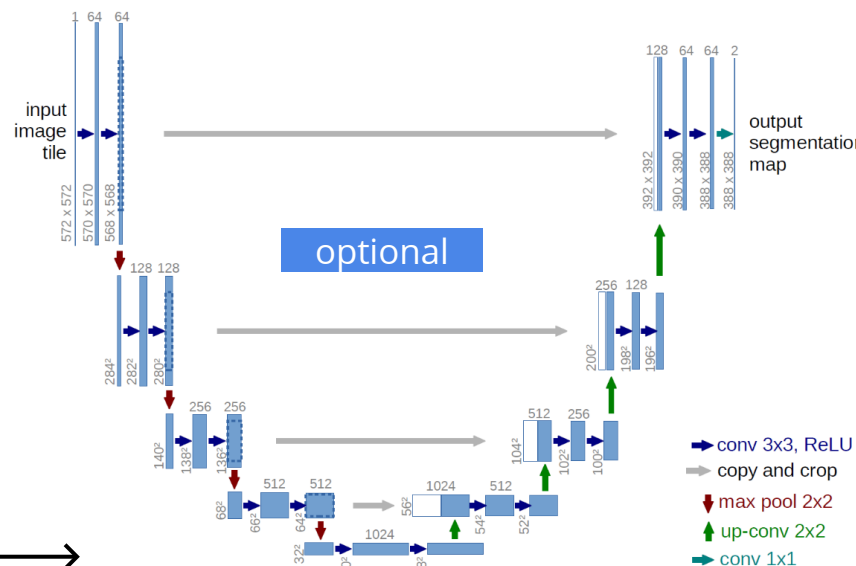
**example**

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** $\longrightarrow$

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

31

image:https://sthalles.github.io/deep_segmentation_network/

# 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**

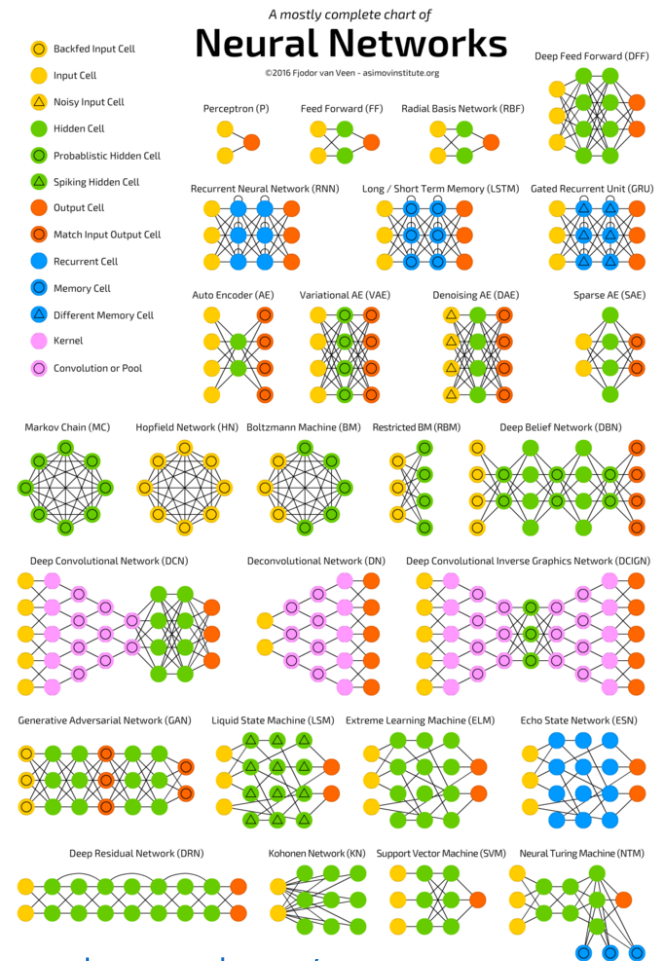**conv-nets** use combinations of

- 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

# Next?

Many other architutures for different task/data



from: https://www.asimovinstitute.org/neural-network-zoo/