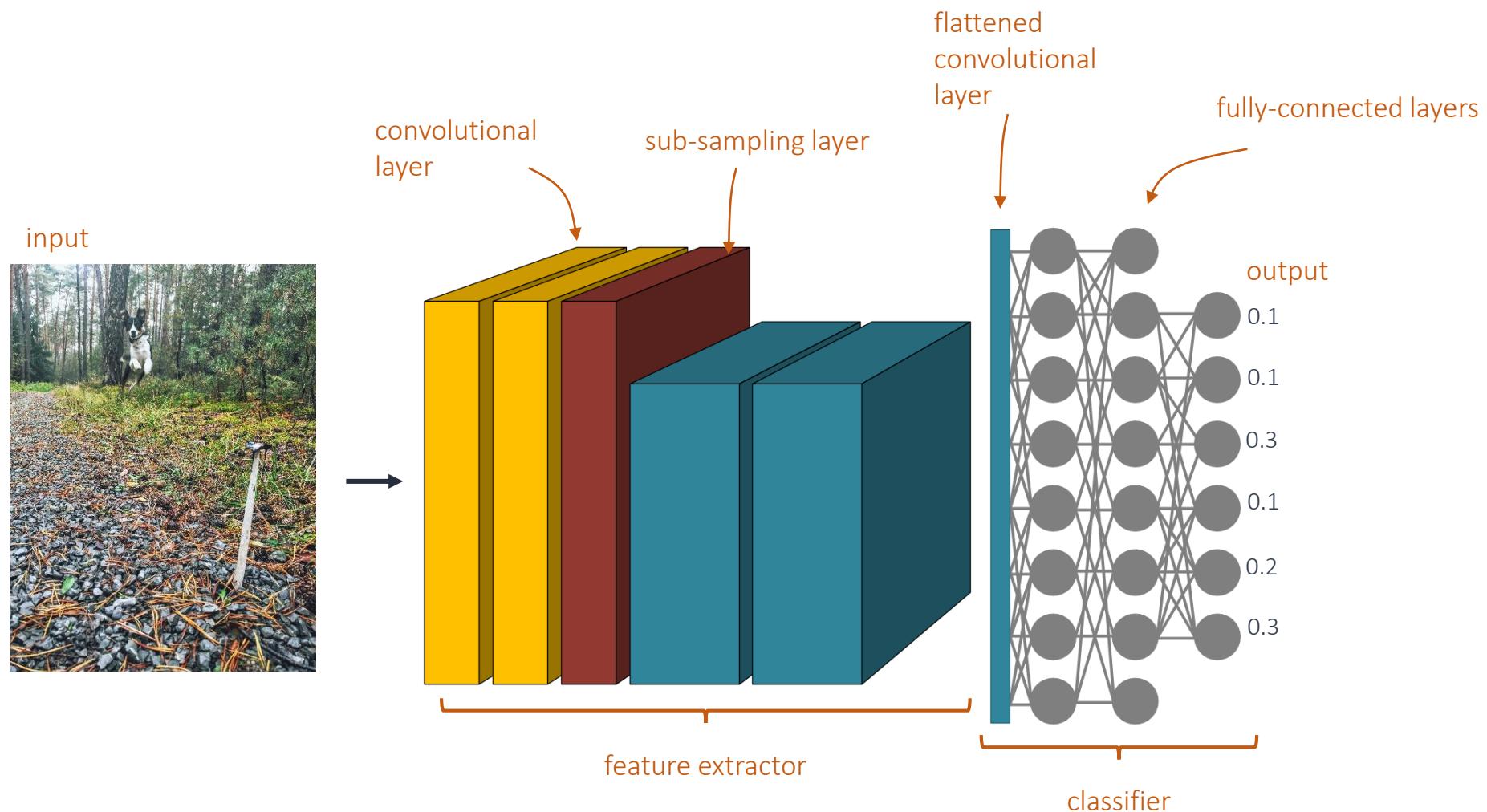


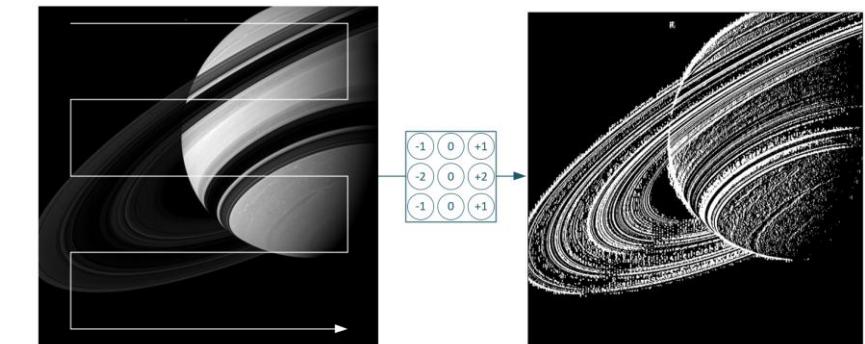
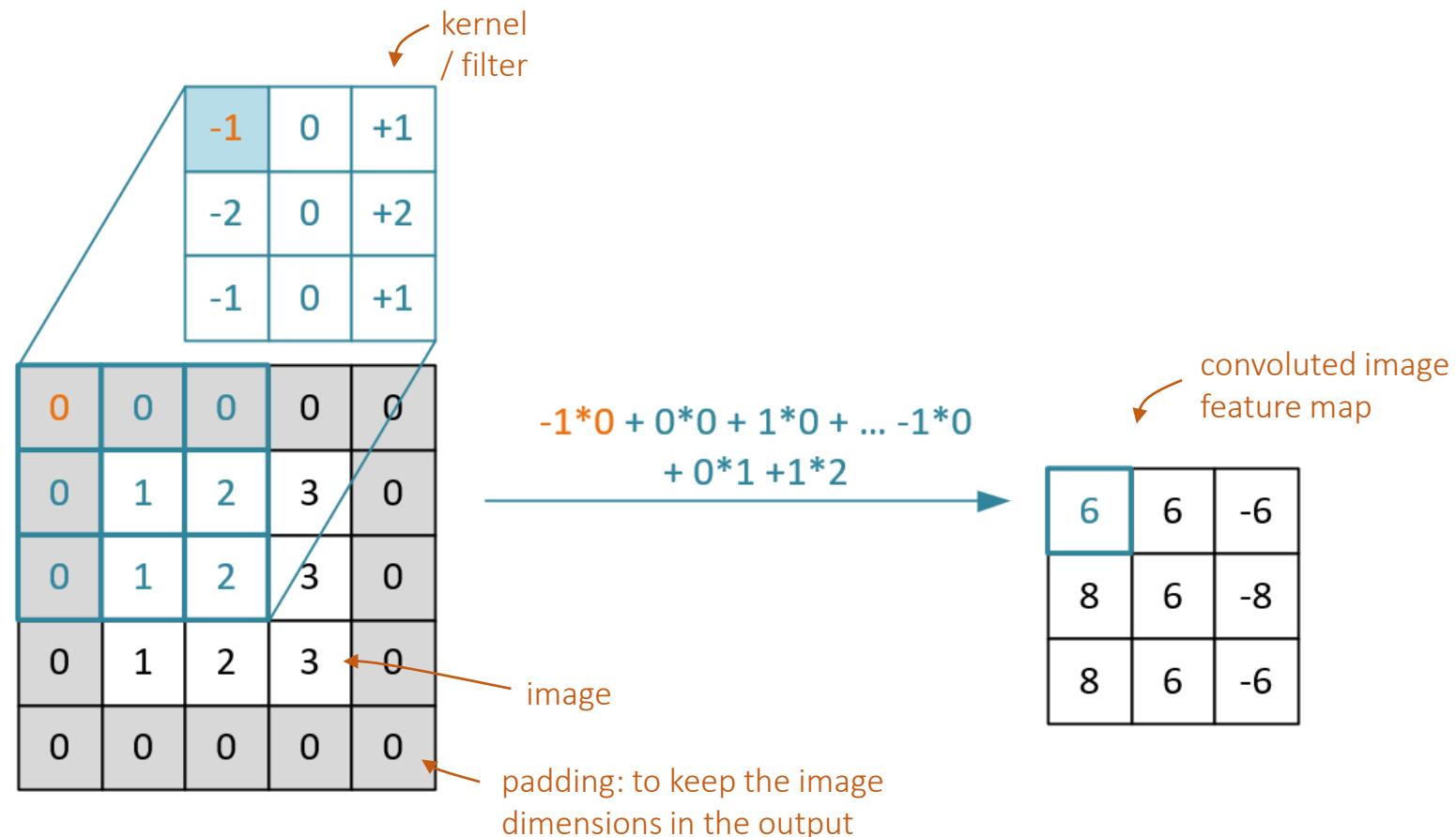
Lecture 6

CNNs & Advanced Network Architectures

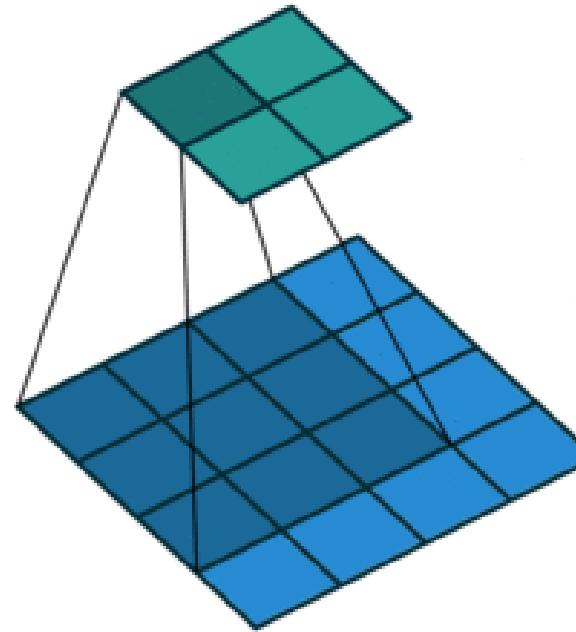
Dr.-Ing. Maike Stern | 19.11.2021

Convolutional neural networks (CNNs)

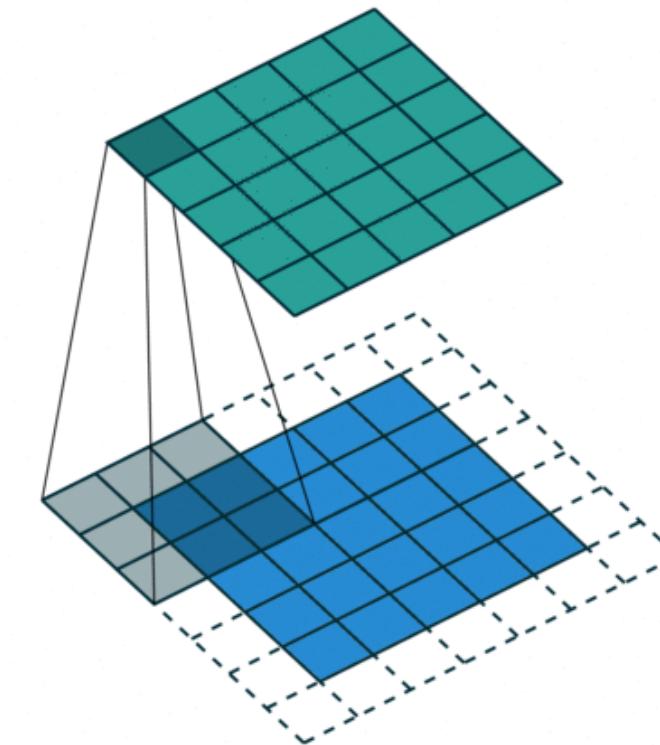




Convolution with a
3x3 filter



Convolution with padding to
keep the dimensions



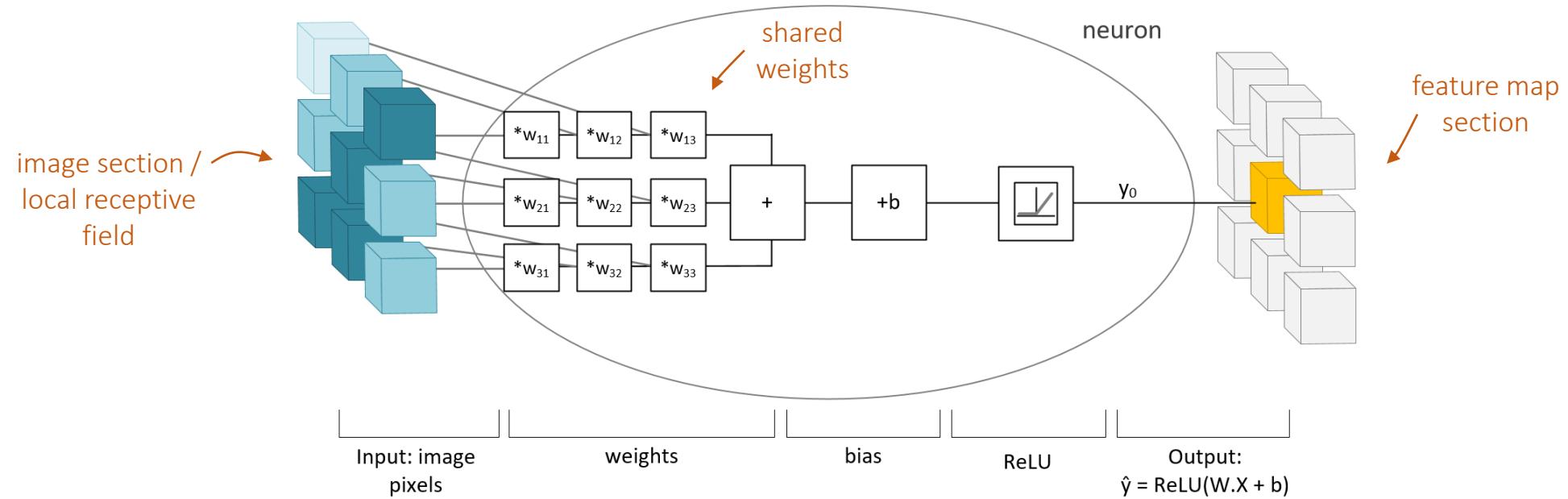
Dumoulin & Visin. A guide to convolution arithmetic for deep learning (2018) https://github.com/vdumoulin/conv_arithmetic

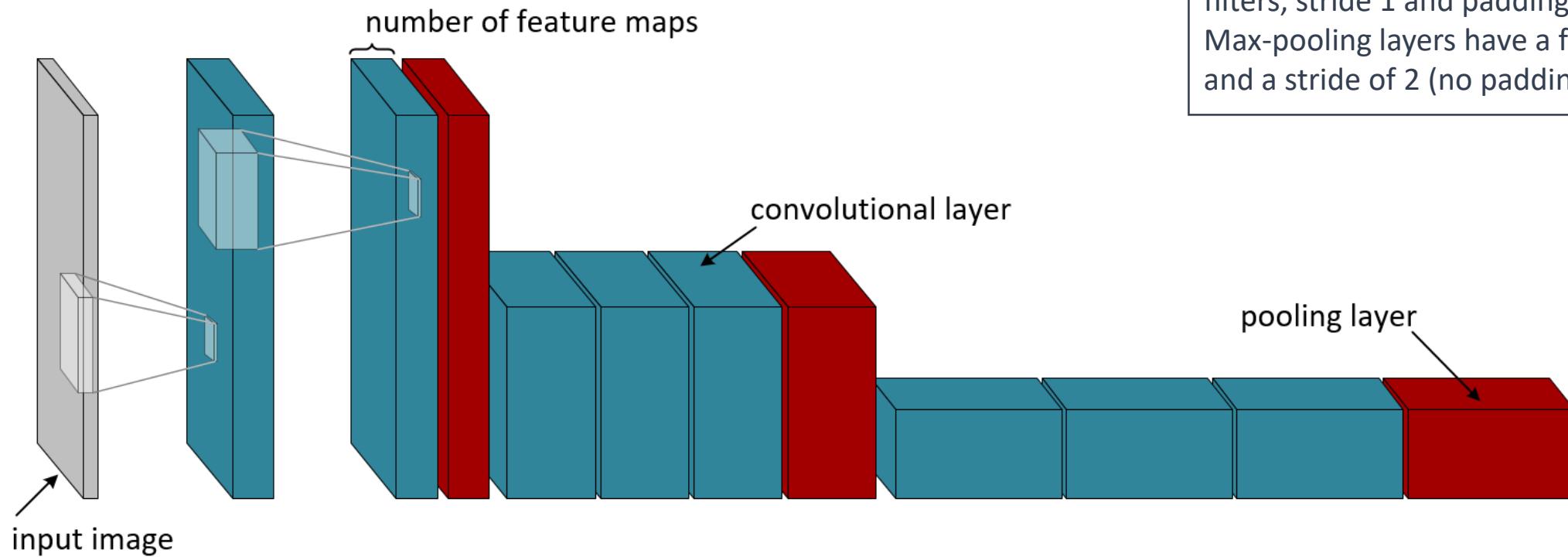
output, located at the
filter's center

input feature
bias kernel height H
maps und width W
input image
pixel kernel weight

$$\text{Convolution: } z_{i',j',f'} = \text{ReLU}\left(b_{f'} + \sum_{f=1}^F \sum_{i=1}^H \sum_{j=1}^W x_{i'+i-1,j'+j-1,f} \cdot w_{ijff'}\right),$$

with i', j' feature map indices, i, j kernel indices, F are the input feature maps and f' is the output feature map

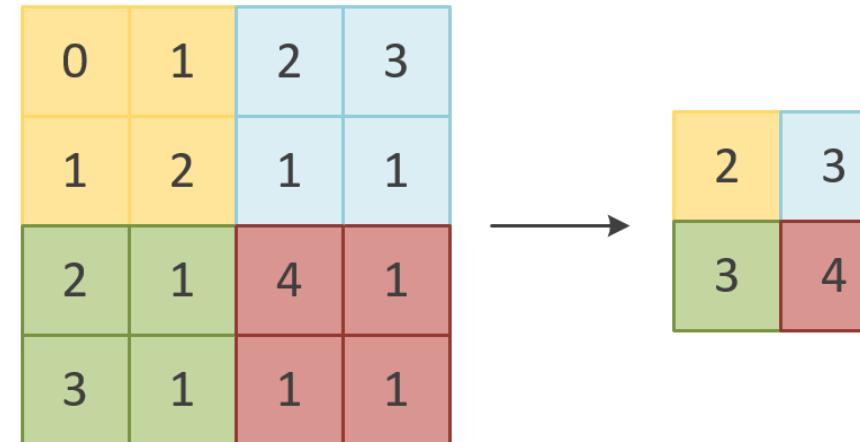




Design rules:
Blocks of two to three convolutional layers followed by a max-pooling layer.
Each layer has 64 to 512 filters, with 3×3 filters, stride 1 and padding 1.
Max-pooling layers have a filter size of 2×2 and a stride of 2 (no padding).

Sub-sampling reduces the feature map size

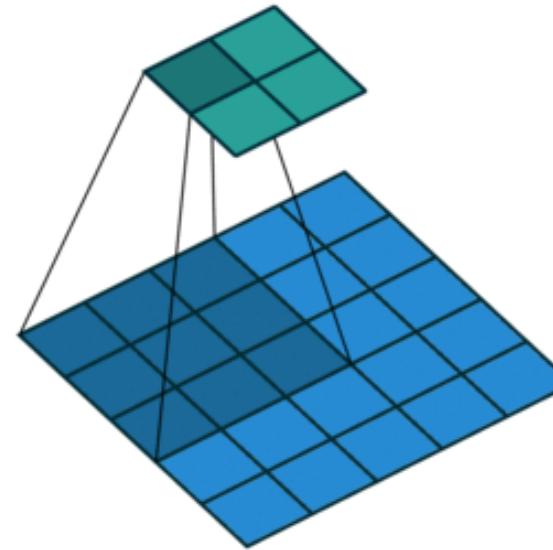
- Reduces the sensitivity of the output to shifts and distortions by introducing invariance to small translations
- Reduces the computational cost and thus allows the typical bi-pyramid:
decreasing feature map sizes & increasing number of filters with increasing network depth
- Mostly: max-pooling, also possible: average pooling, L2 pooling



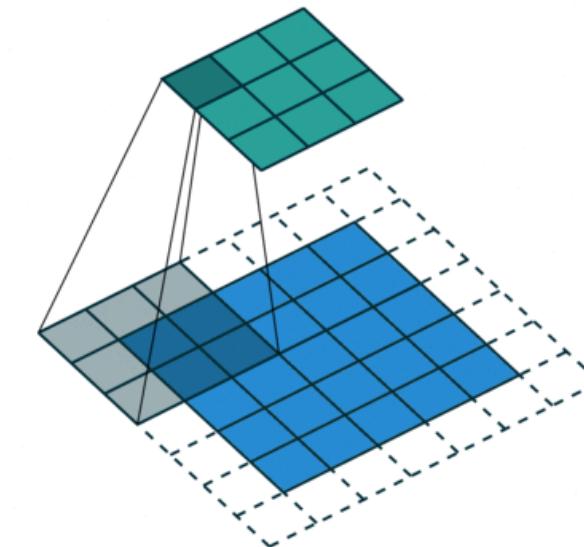
Input size: $W_1 \times H_1 \times D_1$
 New size: $W_2 \times H_2 \times D_2$, with:
 $W_2 = \frac{W_1 - F}{S} + 1$
 $H_2 = \frac{H_1 - F}{S} + 1$
 $D_2 = D_1$
 Where F is the pooling kernel size,
 here 2, and S is the stride, here 2

- Max-pooling may result in a loss of accurate spatial information
- Instead of sub-sampling with max pooling, we can also use convolutional layers, with a stride of 2 (or bigger)
- Often used in GANs and variational autoencoders

Convolution with a 3x3 filter and stride 2



Convolution with a 3x3 filter,
padding, and stride 2



Springenberg et al. Striving for Simplicity: The All Convolutional Net (2014)

To output a classification, convolutional neural networks transition from the convolutional part (feature detection) into the classifier part, consisting of fully-connected layers

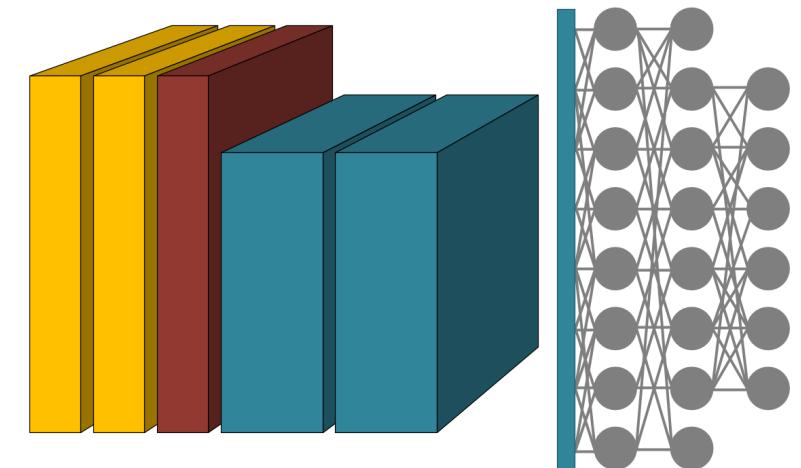
There are two ways to connect both layer types:

- Flatten
 - Use the flattened last convolutional layer as input for the first fully-connected layer
 $(\text{batch_dim}, \text{H}, \text{W}, \text{num_channels}) \rightarrow (\text{batch_dim}, \text{H} * \text{W} * \text{num_channels})$
- Global average pooling
 - Computes the mean of each feature map
 $(\text{batch_dim}, \text{H}, \text{W}, \text{num_channels}) \rightarrow (\text{batch_dim}, \text{num_channels})$
 - Also available as max operation

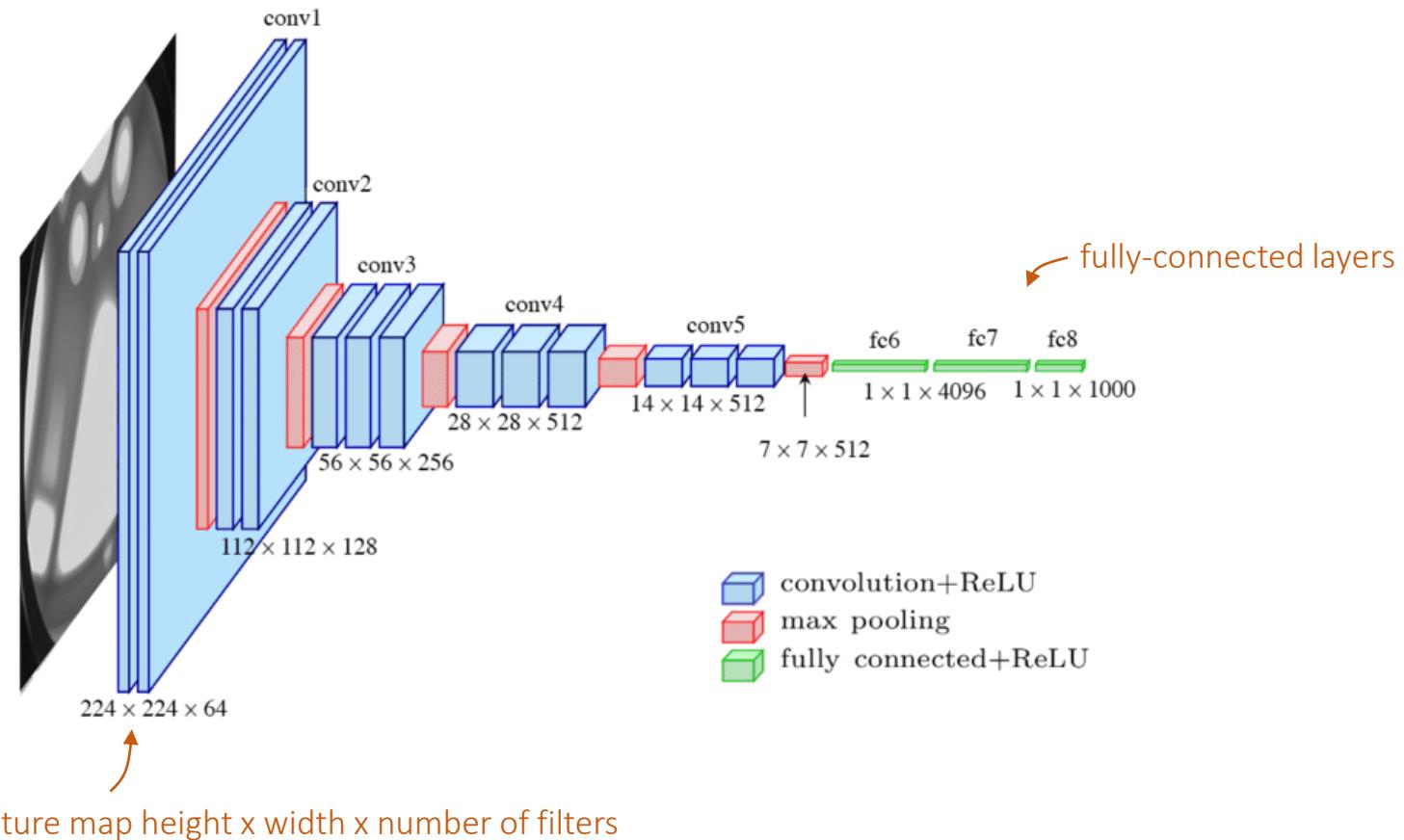
the resulting number of units depends number of feature maps in the last convolutional layer



the resulting number of units depends on the input image size



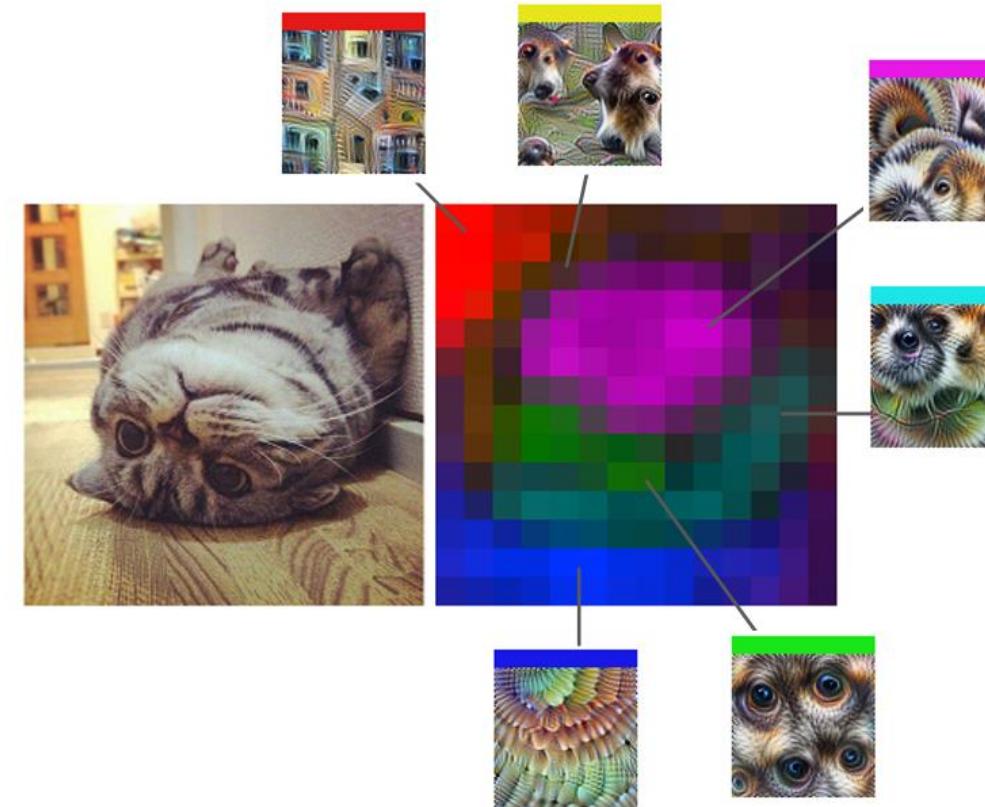
The VGG-16 network architecture



What a network learns to see

What a convolutional network learns to see

- With self-learned filters, the network learns to distinguish objects that were present in the training data
- Visualising the filters with the highest activation shows how the network analysis an image



<https://distill.pub/2018/building-blocks/>

What a convolutional network learns to see

- Visualising filters corresponding to a certain class category, here „labrador retriever“ (orange) and „tiger cat“ (blue), shows how the network collects evidence for each hypothesis
- Deeper layers further evaluate previous representations and thus become more confident



<https://distill.pub/2018/building-blocks/>
<https://distill.pub/2019/activation-atlas/>

What a convolutional network learns to see



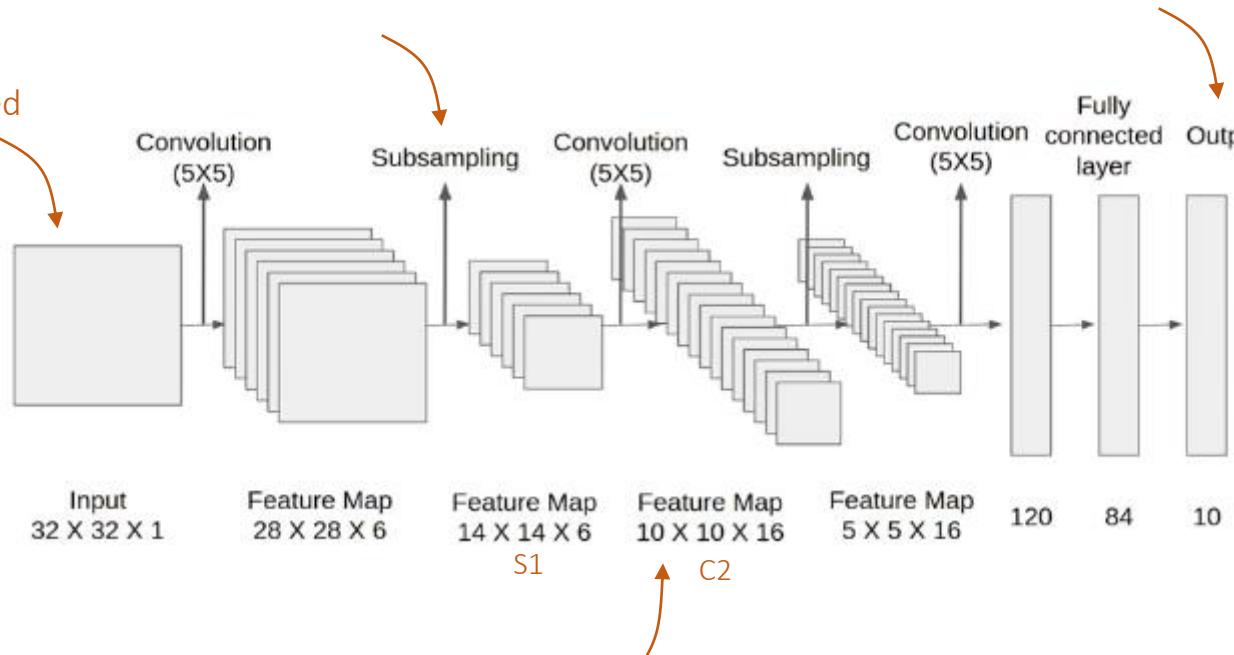
Convolutional Neural Network Architectures

Target dataset: MNIST

Input images are padded
(from 28x28) and normalised

Average pooling with 2x2 kernels and stride 2

Computes the square of the Euclidean distance between the output's input vector and the weight vector and measures how much the image belongs to a certain class



Feature maps are not padded,
thus feature map sizes shrink
with every convolution

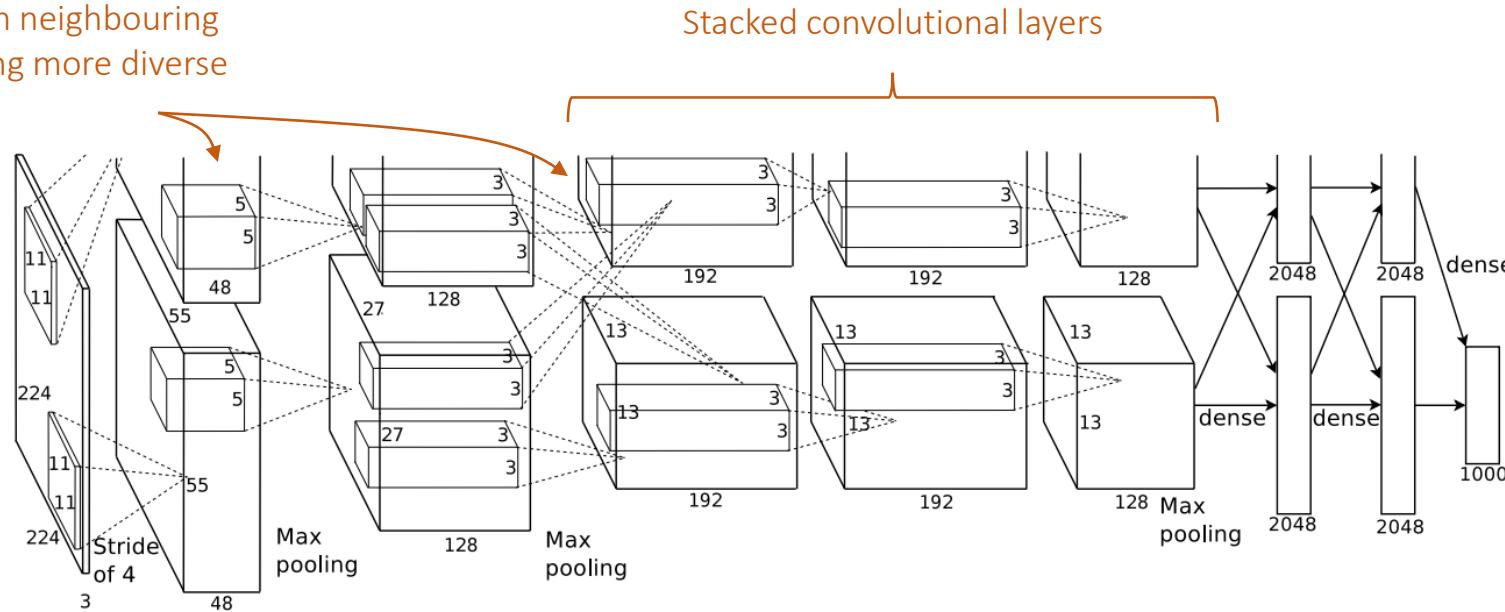
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	X			X	X	X		X	X	X	X	X			X	X
1	X	X			X	X	X		X	X	X	X	X			X
2	X	X	X			X	X	X		X	X	X			X	X
3	X	X	X		X	X	X	X		X	X			X	X	X
4	X	X	X		X	X	X	X		X	X	X	X	X	X	X
5	X	X	X		X	X	X	X		X	X	X	X	X	X	X

Feature maps of S1 combined by each C2 feature map

- Reduces the number of connections
- Forces a break of symmetry → more diversity in the learned features

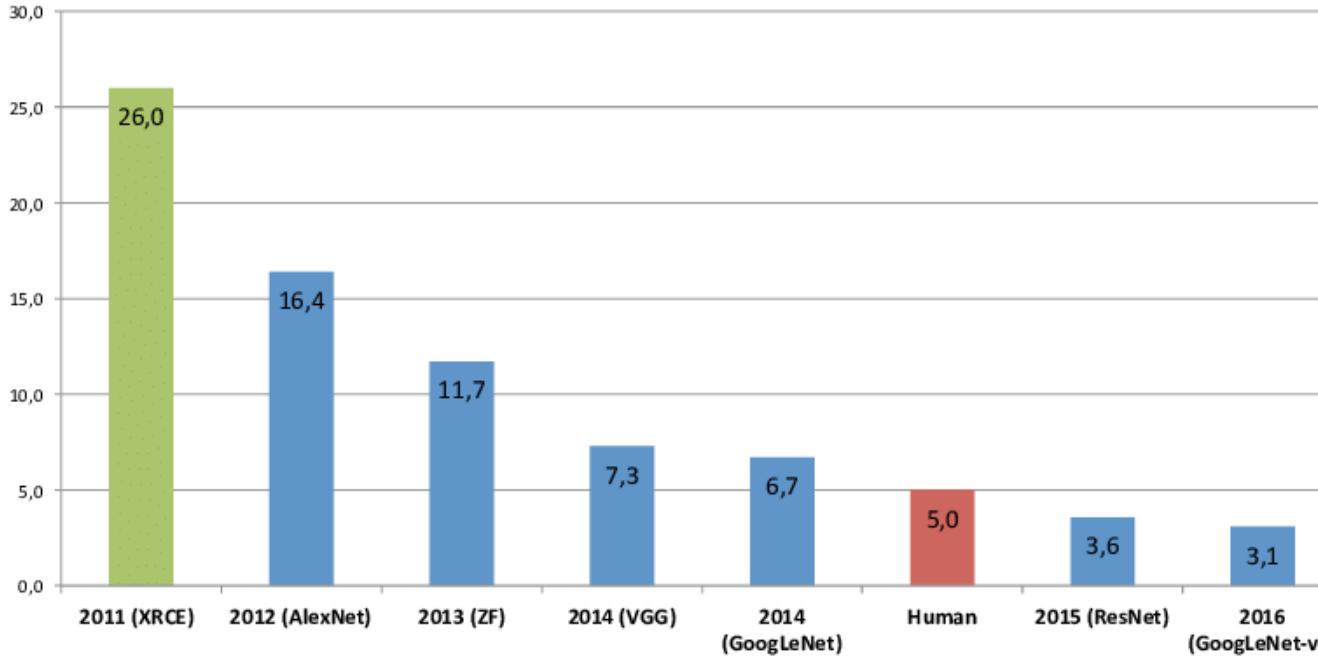
- Target dataset: ImageNet challenge dataset
- 50% dropout rate in the fully-connected layers
- Data augmentation

Local response normalisation (LRN): the most strongly activated neurons inhibit other neurons located at the same position in neighbouring feature maps, thus encouraging more diverse feature maps



- <https://www.image-net.org/index.php>

ImageNet Classification Error (Top 5)



IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.



- <https://www.image-net.org/index.php>
- 2019: Cleaning of the dataset
Yang et al. Towards Fairer Datasets: Filtering and Balancing the Distribution of the People Subtree in the ImageNet Hierarchy <https://arxiv.org/pdf/1912.07726.pdf>
- 2020: New Labels
Beyer et al. Are we done with ImageNet? (2020) <https://arxiv.org/pdf/2006.07159.pdf>

IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.



Old label: pier
ReaL: dock; pier;
speedboat; sandbar;
seashore



Old label: hammer
ReaL: screwdriver;
hammer; power drill;
carpenter's kit

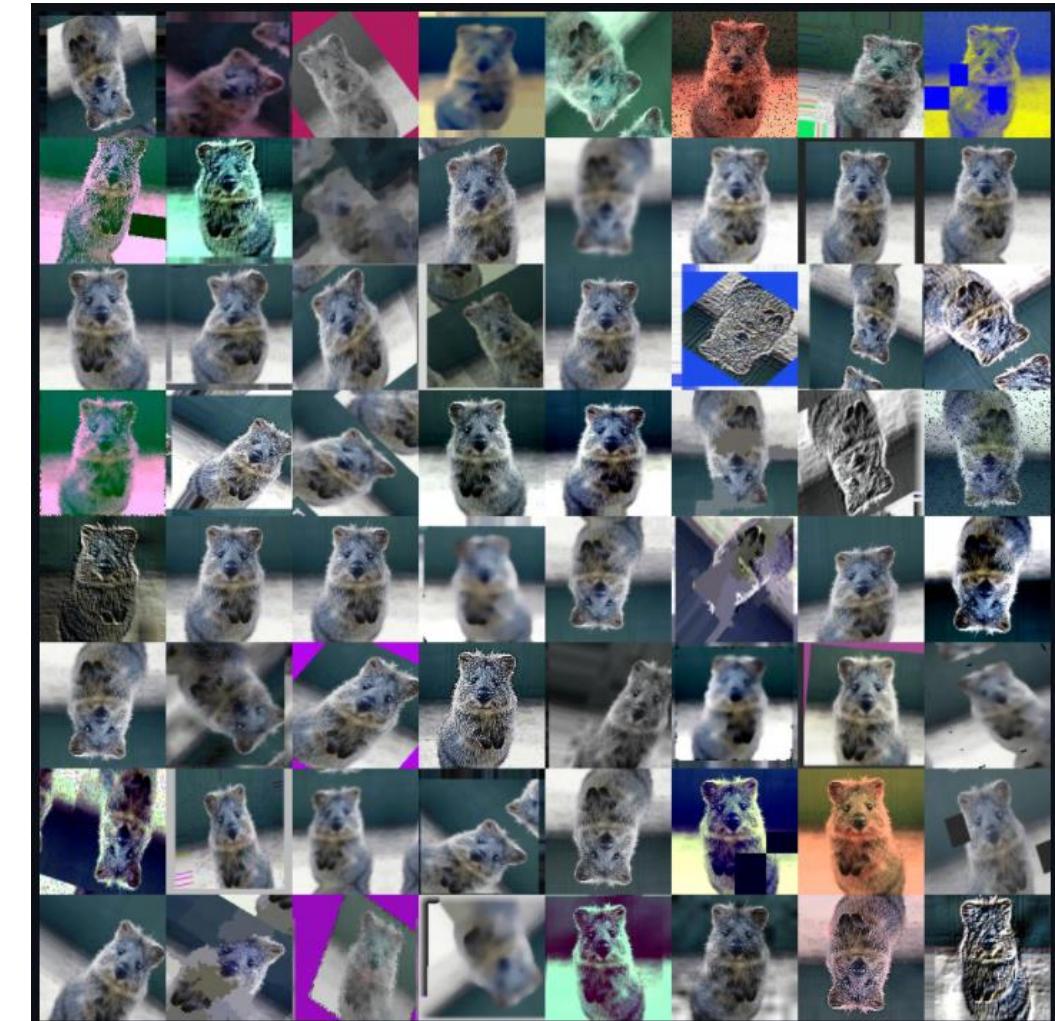


Increase the training dataset size by generating many **realistic** variants of the original training samples

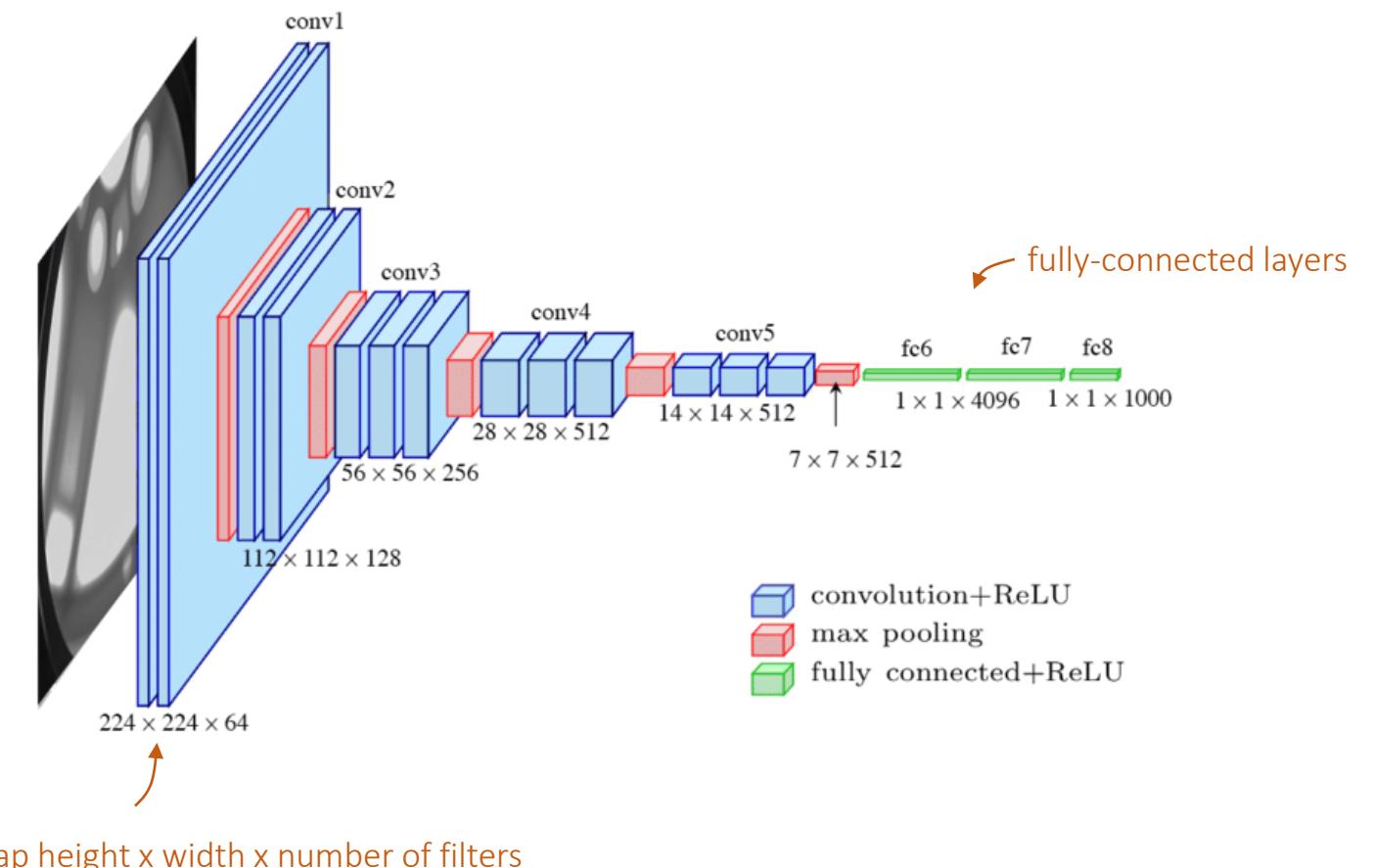
- Reduces overfitting by increasing the model's robustness

Augmentation libraries:

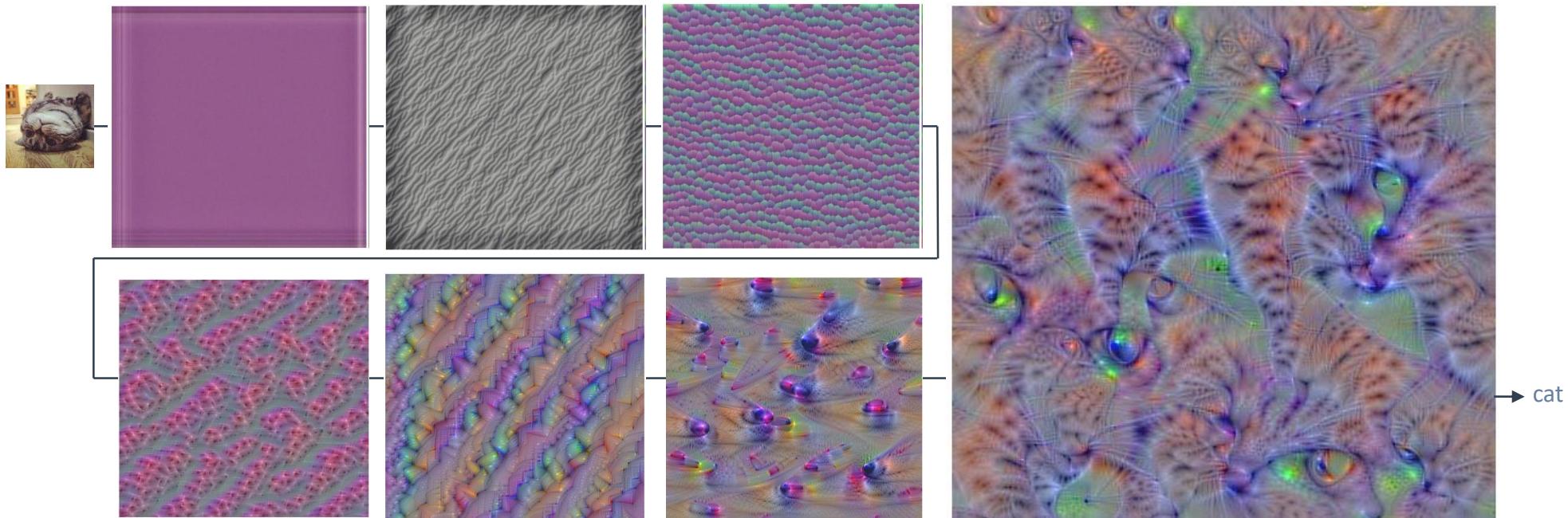
- https://www.tensorflow.org/tutorials/images/data_augmentation
- <https://github.com/aleju/imgaug>
- <https://ai.facebook.com/blog/ugly-a-new-data-augmentation-library-to-help-build-more-robust-ai-models/>
- <https://albumentations.ai/>



- Target dataset: ImageNet
- Second place, ImageNet challenge 2014
- Two networks: VGG-16 and VGG-19
- 3x3 filters throughout the network
- Transfer learning



- The shallow layers of a CNN learn basic patterns, such as colours, edges, etc.
 - With increasing network depth, the learned filters get more specific as they operate on more and more abstracted features of the original input and have a larger receptive field
- Transfer parameters from a pre-trained network and use them to initialise (parts of) your network



In case of similar tasks

- Choose a suitable network architecture

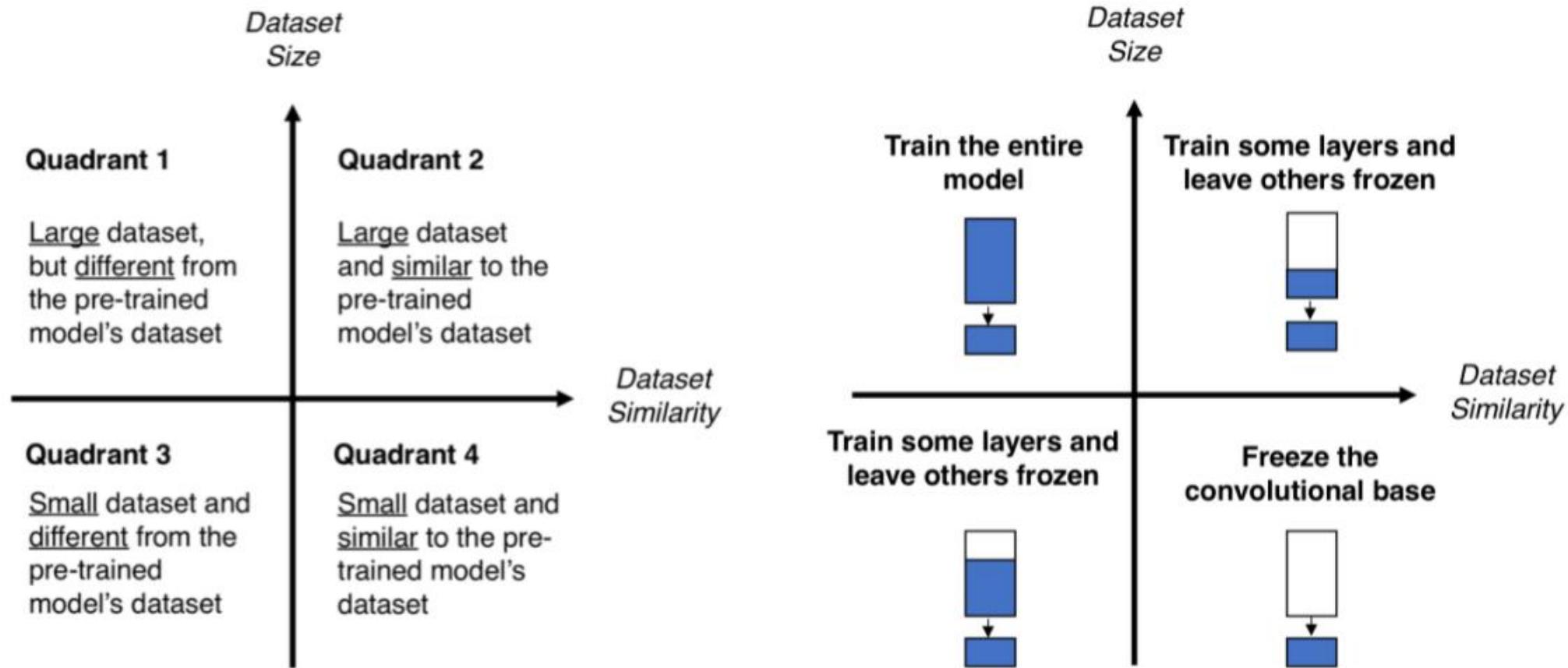
„False“ removes the
fully-connected layers

```
tf.keras.applications.vgg16.VGG16(  
    include_top=True, weights='imagenet', input_tensor=None,  
    input_shape=None, pooling=None, classes=1000,  
    classifier_activation='softmax'  
)
```

Use weights pre-trained
on the ImageNet dataset

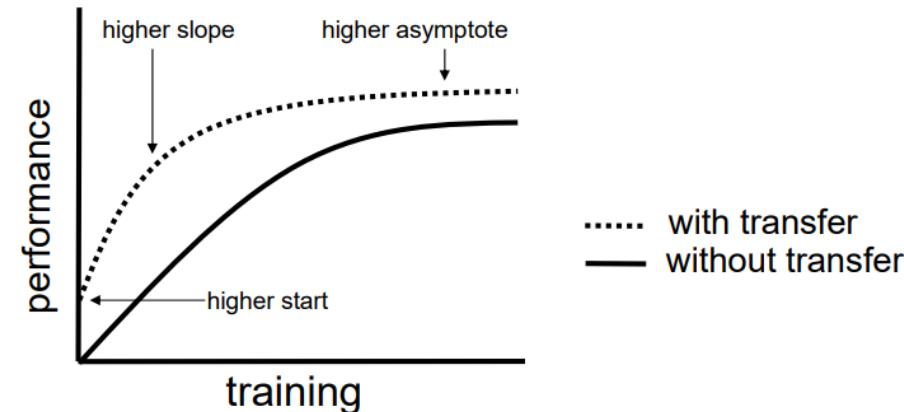
- Add your own classifier (fully-connected layers, randomly initialised)
- Freeze the pre-trained layer parameters
- Train the new layers on your dataset
- If your data is quite different to the dataset the network was trained on:
 - Unfreeze the convolutional layers and fine-tune the whole network using a very small learning rate
 - For very dissimilar data: Transfer only parameters of the shallow convolutional layers and randomly initialise the deep conv layers & fully-connected layers

- ▼ applications
 - [Overview](#)
 - [MobileNetV3Large](#)
 - [MobileNetV3Small](#)
 - [densenet](#)
 - [efficientnet](#)
 - [efficientnet_v2](#) ↗
 - [imagenet_utils](#)
 - [inception_resnet_v2](#)
 - [inception_v3](#)
 - [mobilenet](#)
 - [mobilenet_v2](#)
 - [mobilenet_v3](#)
 - [nasnet](#)
 - [resnet](#)
 - [resnet50](#)
 - [resnet_v2](#)
 - [vgg16](#)
 - [vgg19](#)
 - [xception](#)

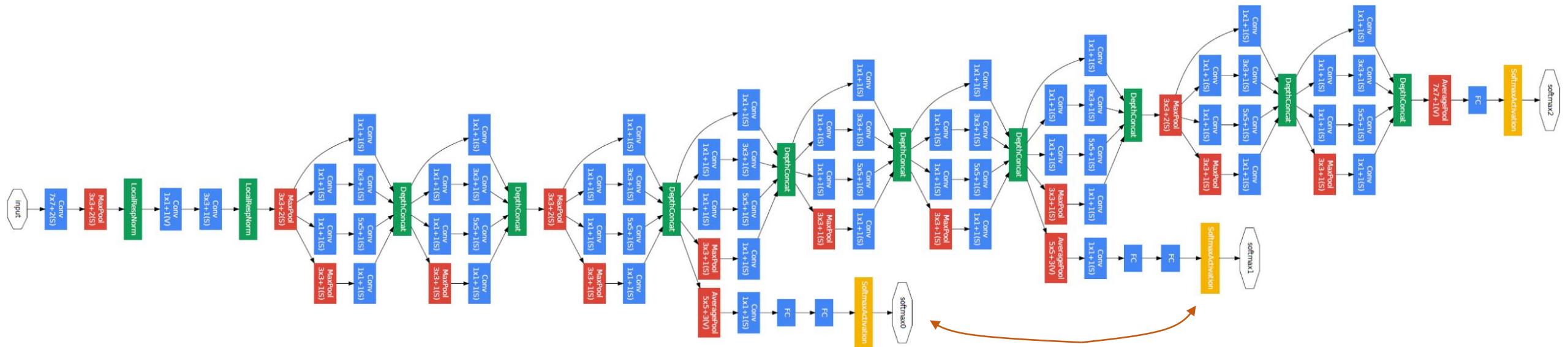


Transfer learning

- Exploits the fact that images typically feature the same patterns up to a certain degree (e.g. edges, light-dark transitions, colours, ...)
- Hence, transfer learning speeds up the training and requires less data
- Works best on similar data, but can also improve the performance of dissimilar images
- Popular not only in computer vision but also natural language processing (NLP), where large transformer models (ELMO, BERT, GPT-3, etc.) are trained on very, very large datasets A single training of GPT-3 costs >4.5 Million \$

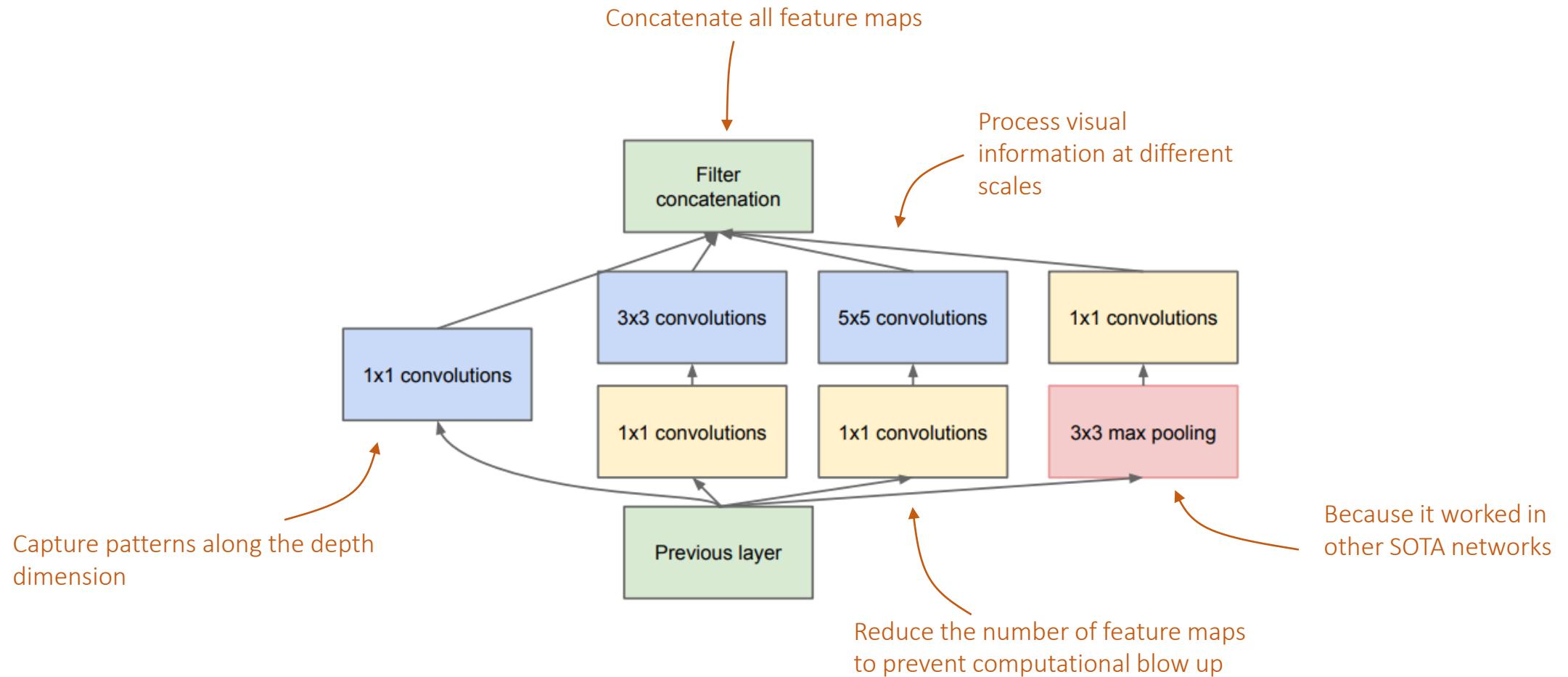


- Target dataset: ImageNet
- First place ImageNet challenge 2014 (ensemble of 7 networks)
- 22-layer deep network, but only ≈ 6 million parameters (AlexNet ≈ 60 million)

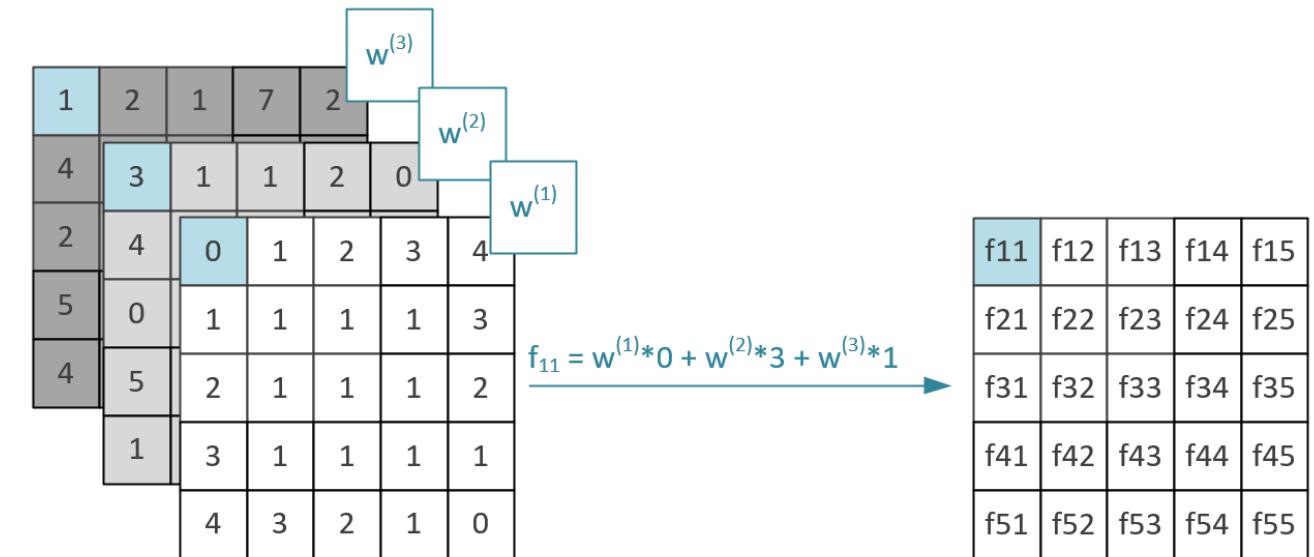


Auxiliary classifiers:

- Improve gradients
- Encourage feature discrimination in shallow layers
- Loss is added to the final loss

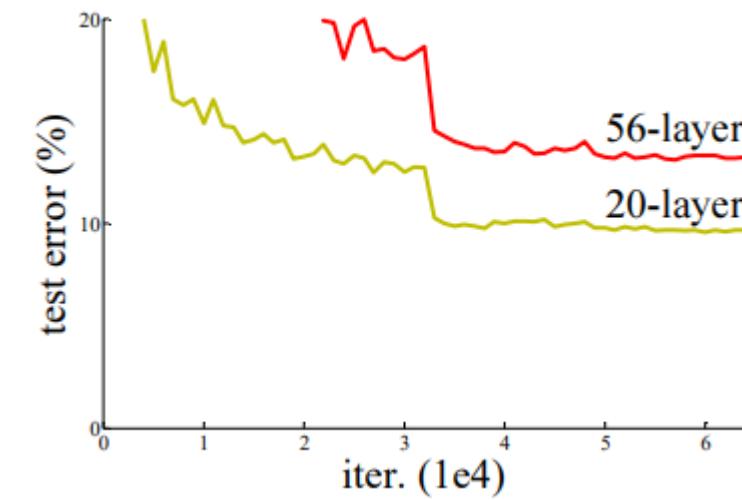
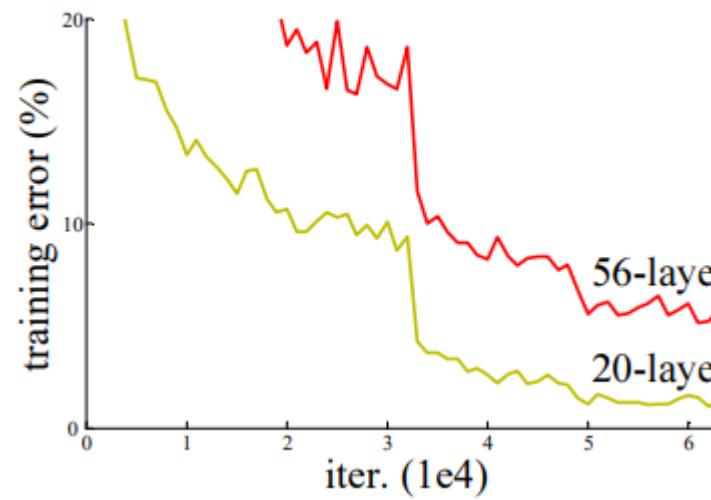


- Also called feature map pooling or projection layer
- Used to reduce the number of feature maps (depth)
- Kind of pixel-wise summary of the input feature maps



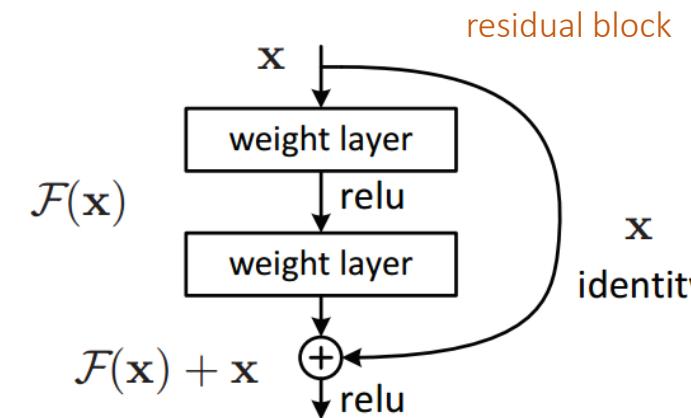
ResNets address the so-called degradation problem: with increasing network depth, accuracy saturates and then degrades rapidly

Not overfitting, because both, training error
and test error are higher in the deeper network

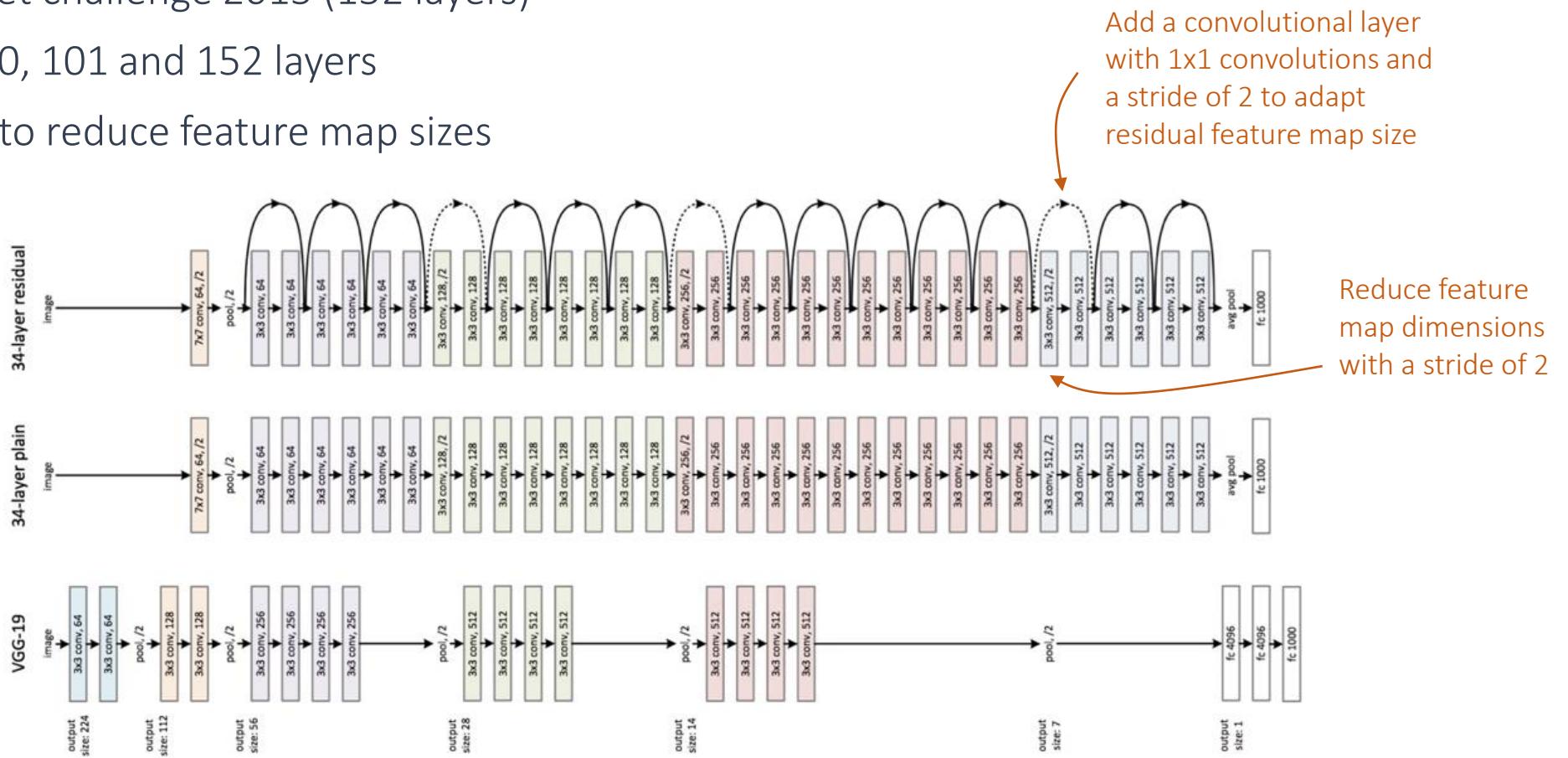


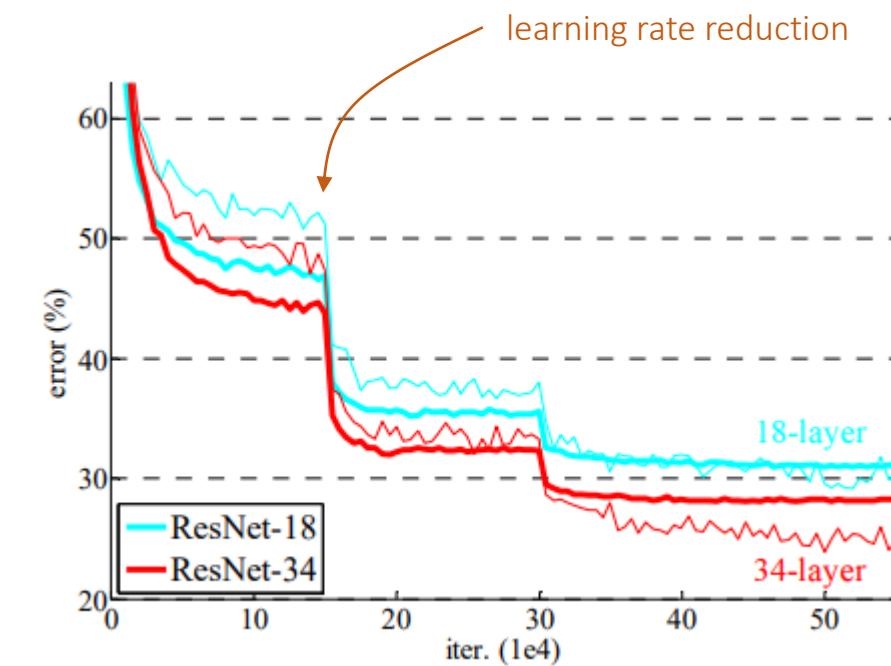
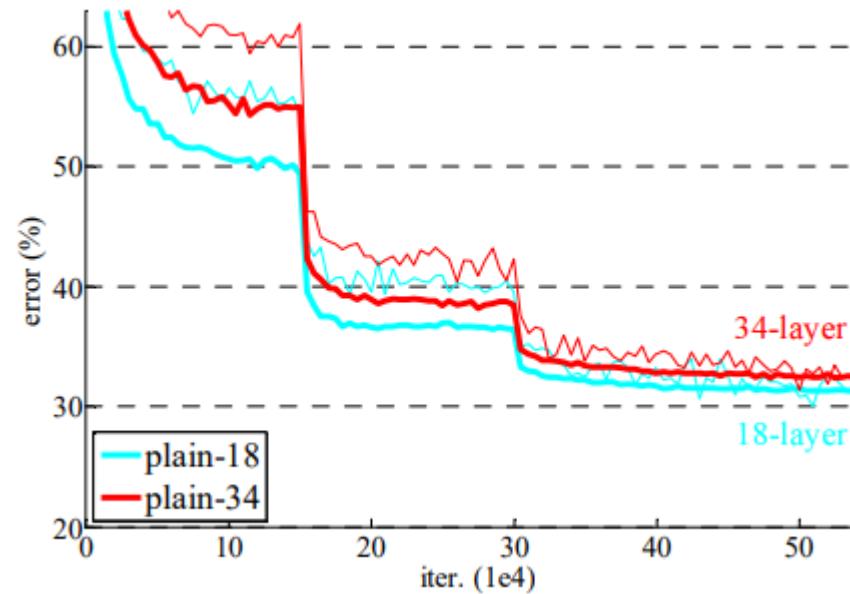
To address the degradation problem, ResNets introduce residual blocks

- Randomly initialised layers (with standard-normally distributed weights close to zero) output standard-normally distributed values close to zero
- Hence, untrained residual blocks output the identity function, that is they output a copy of their input
- Often, the target function is close to the identity function, hence residual blocks speed up training



- Target dataset: ImageNet
- First place ImageNet challenge 2015 (152 layers)
- Variants with 34, 50, 101 and 152 layers
- Uses convolutions to reduce feature map sizes





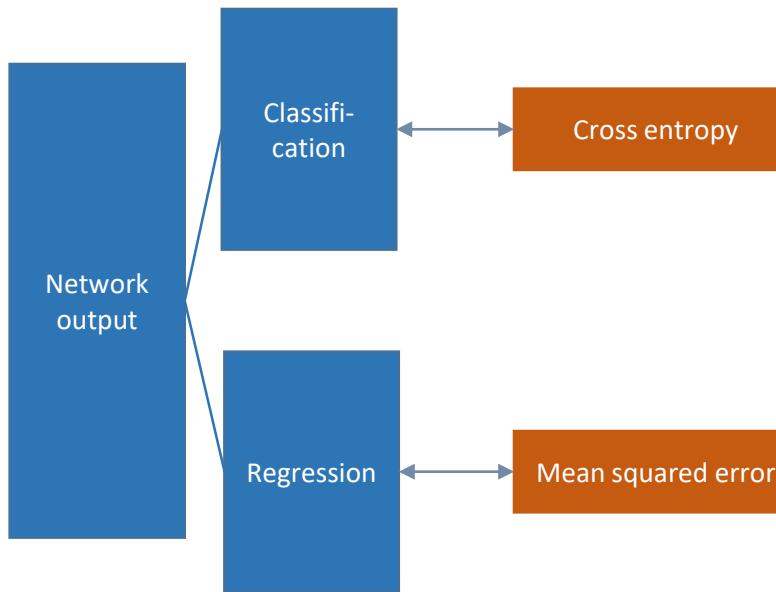
Object Localisation & Detection

- Determine the object position in addition to the object classification
- So-called bounding boxes visualise the object position in the image
- Bounding boxes are comparably cheap to label (in contrast to pixel-wise labels)
- Research dataset: <https://cocodataset.org/#home> (object detection, segmentation, and captioning)

bounding box



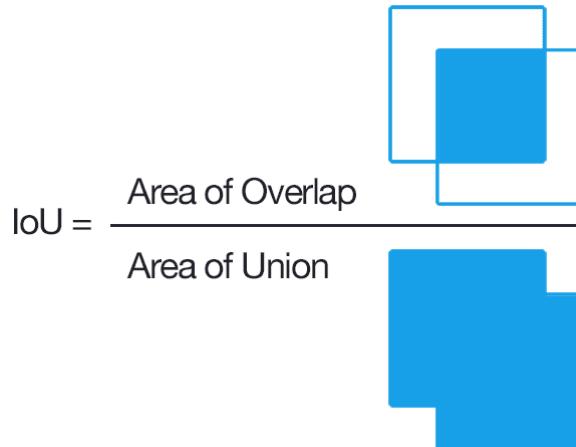
- Determine the object position in addition to the object classification
- So-called bounding boxes visualise the object position in the image
- Add an additional regression output to the network, predicting bounding box coordinates (center (x,y), height, width)



bounding box
center coordinates
bounding box
height & width
coordinates



- Determine the object position in addition to the object classification
- So-called bounding boxes visualise the object position in the image
- Add an additional regression output to the network, predicting bounding box coordinates (center (x,y), height, width)
- Performance is typically determined using the intersection over union metric

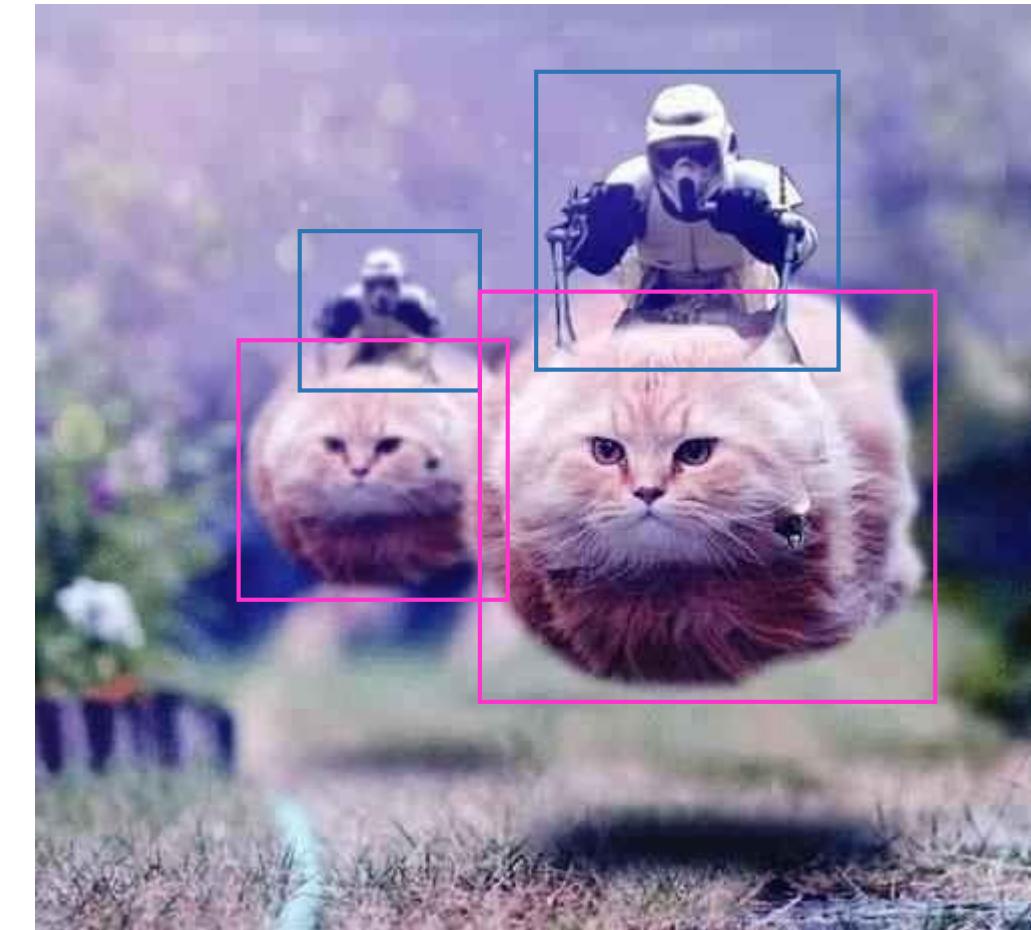


IoU values over 0.5 are usually considered a good prediction

bounding box
 center coordinates
 bounding box
 height & width
 coordinates



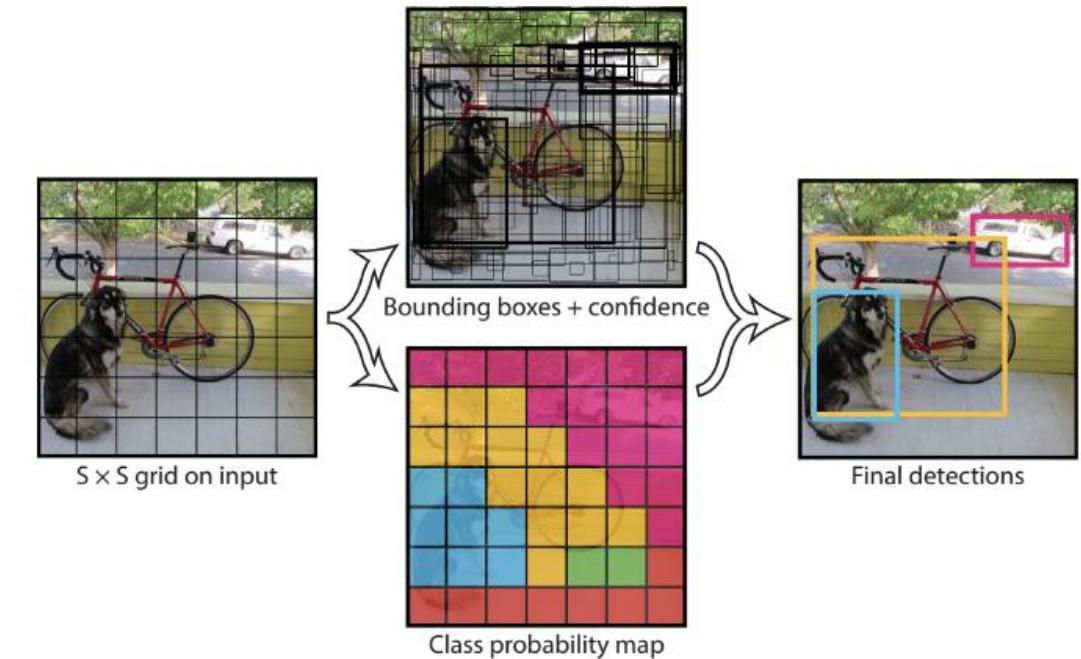
- Classify and localise different objects in an image
- The problem: the length of the output layer depends on the variable number of objects in the image



- One approach (very simplified) is to split the image into $s \times s$ grids and predict m differently sized bounding boxes per grid
- For each of the bounding boxes, an objectness score is computed (prediction confidence) as well as the class probabilities

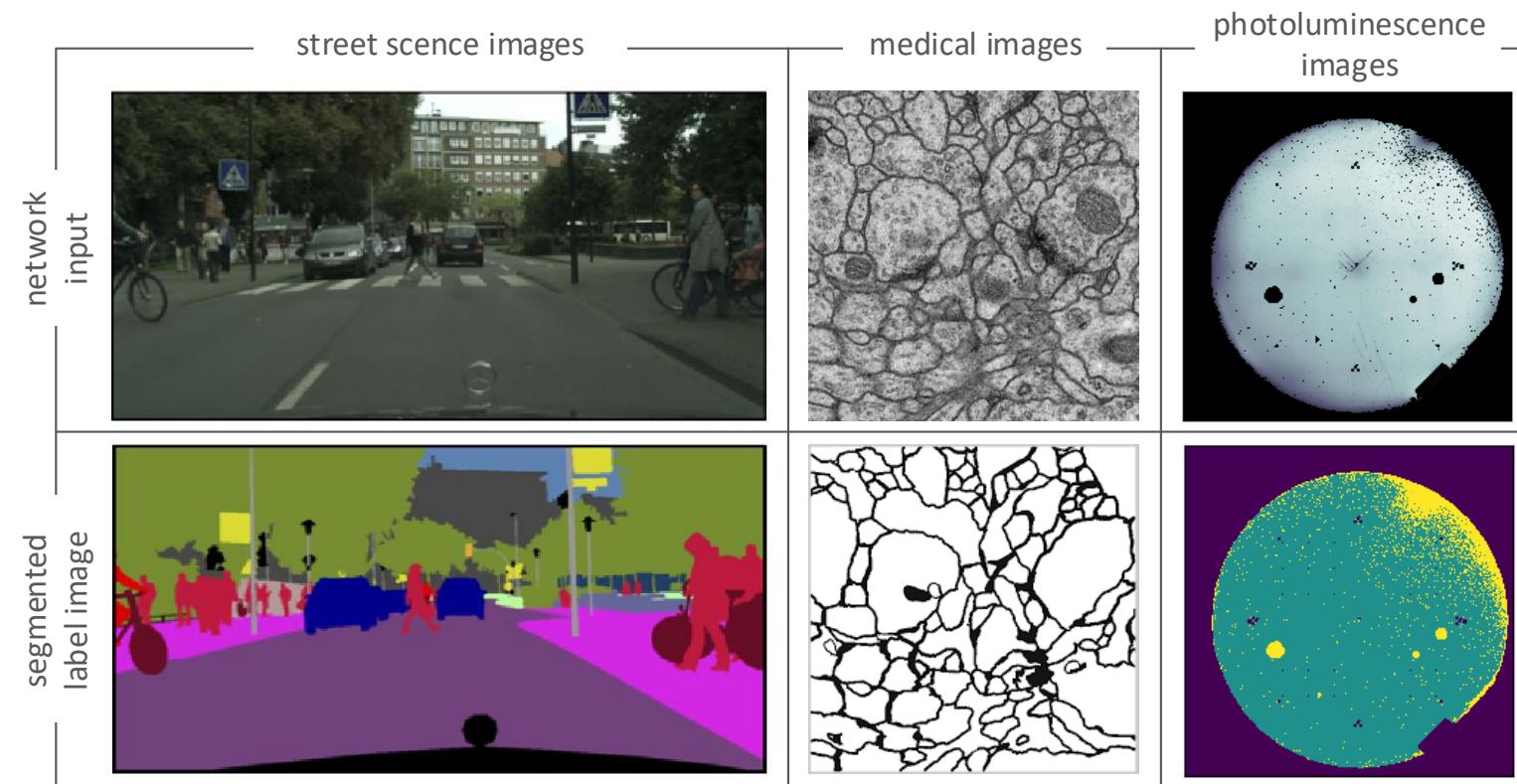
State-of-the-art architectures

- R-CNN, Fast R-CNN and Faster R-CNN
(<https://paperswithcode.com/method/r-cnn>)
- YOLO, YOLO9000, YOLOv3(-6)
(<https://paperswithcode.com/paper/yolov3-an-incremental-improvement>)
- SSD
(<https://paperswithcode.com/method/ssd>)
- RetinaNet
(Keras: <https://keras.io/examples/vision/retinanet/>)



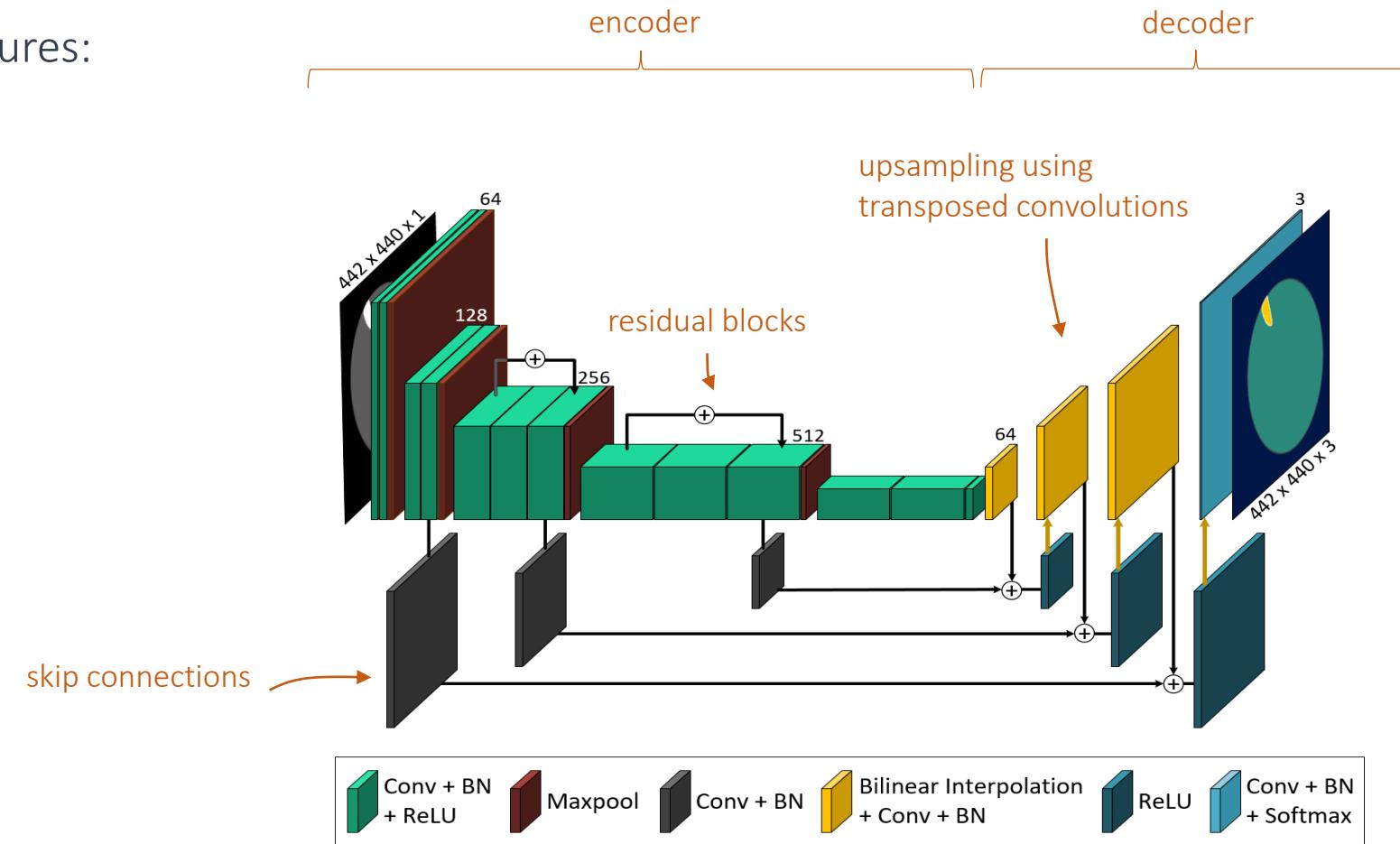
Semantic Segmentation

Pixel-wise classification of images

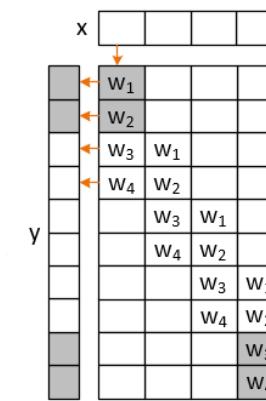
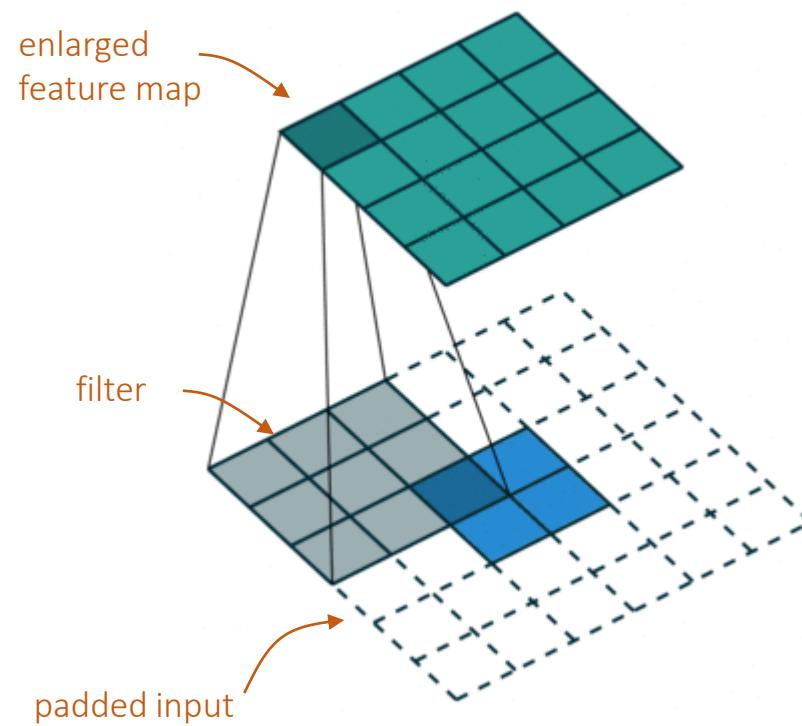


State-of-the-art architectures:

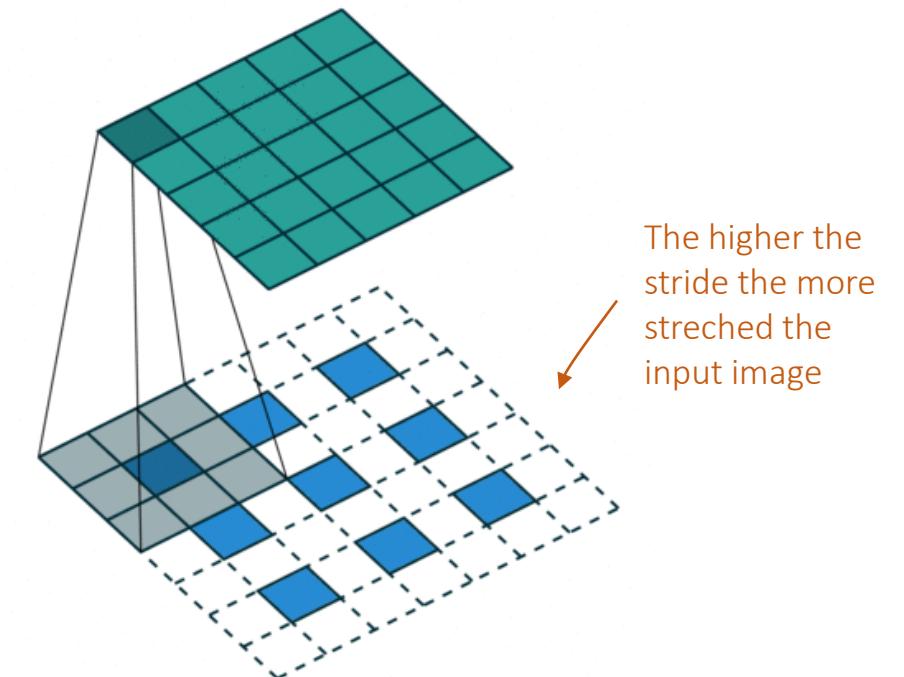
- U-Net
- DeepLab1-3
- Segnet



Transposed convolution
with a 3x3 filter

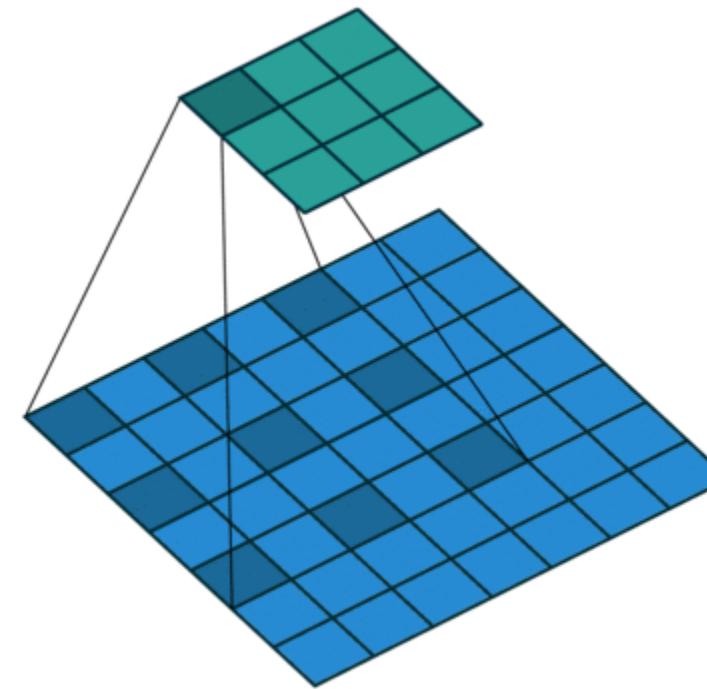


Transposed convolution with a
3x3 filter and a stride of 2



Dumoulin & Visin. A guide to convolution arithmetic for deep learning (2018) https://github.com/vdumoulin/conv_arithmetic

Dilated convolution with a
3x3 filter



Dumoulin & Visin. A guide to convolution arithmetic for deep learning (2018) https://github.com/vdumoulin/conv_arithmetic

Generative Adversarial Networks

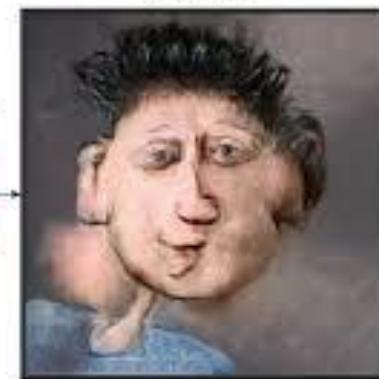
GANs are used to generate new data, e.g. images of non-existent people, new music pieces, or manipulate existing data

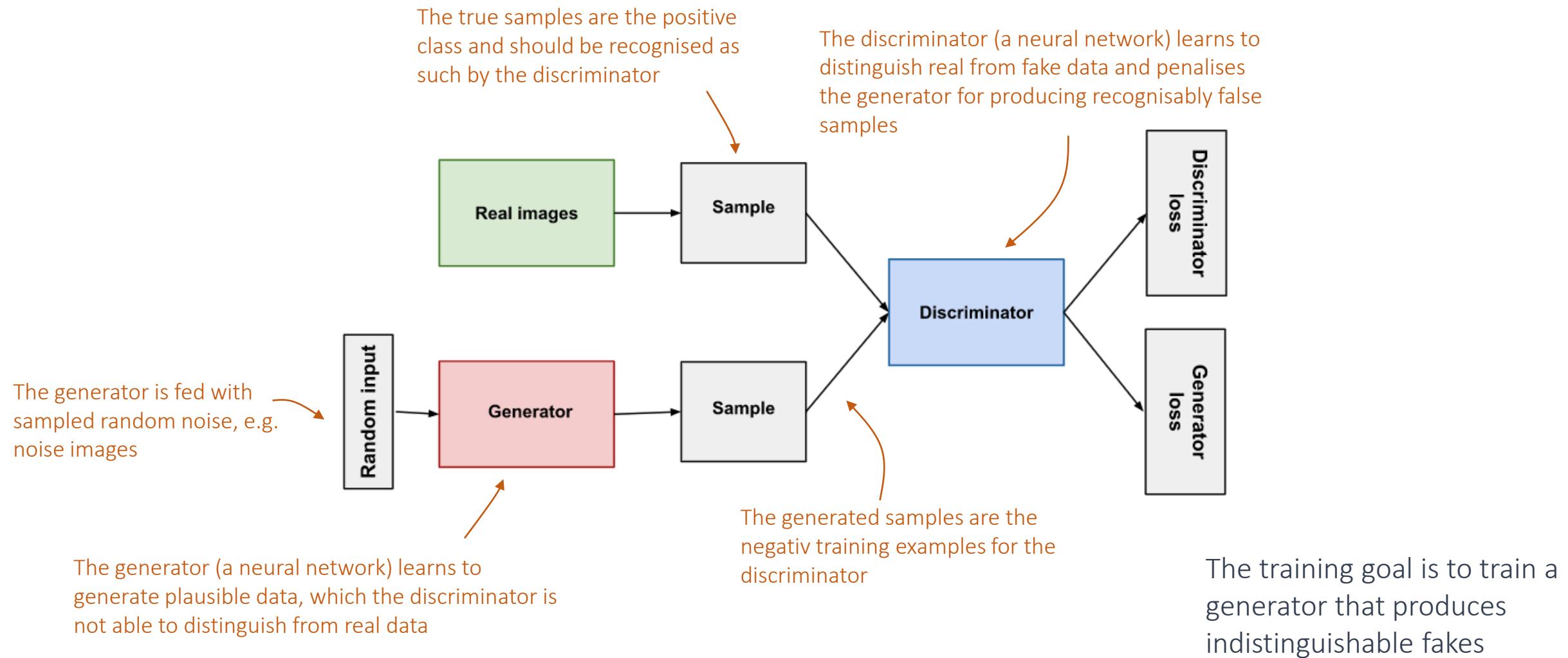


INPUT



OUTPUT





Network Training

In general

- If you are building your own architecture, start with a simple design / design that is known to work
- Keep in mind, that CNNs learn hierarchical feature representations, that is the shallow layers learn simple, general patterns/features (e.g. edges, colours, etc.) while the deep layers learn data-specific features → the less information images contain (e.g. measurement images) the fewer layers your network needs
- Run the sanity checks (next slide)
- Monitor the loss / metrics and tune the hyperparameters (slides after the next slide)
 - For tuning (e.g. with Keras tuner) start with a coarse net of the most important hyperparameters and run optimisation for a few epochs only (e.g. 10), which should be sufficient to discover a trend. Then, focus on the promising areas and run a fine-tuning search (with early stopping!)
 - e.g. start with learning rate values: 0.1, 0.01, 0.001, 0.0001, 0.00001
 - If 0.001 achieves the most promising performance, run a tuning search between 0.01 and 0.0001 (in log space!)
 - Tune sensible combinations, e.g. learning rate and regularisation, number of units and dropout, etc.
- Analyse the results, especially wrong predictions to evaluate certain decisions, e.g. the used loss function
- Take a baseline and add more fancy methods to address the shortcomings

<https://cs231n.github.io/neural-networks-3/>

Run sanity checks

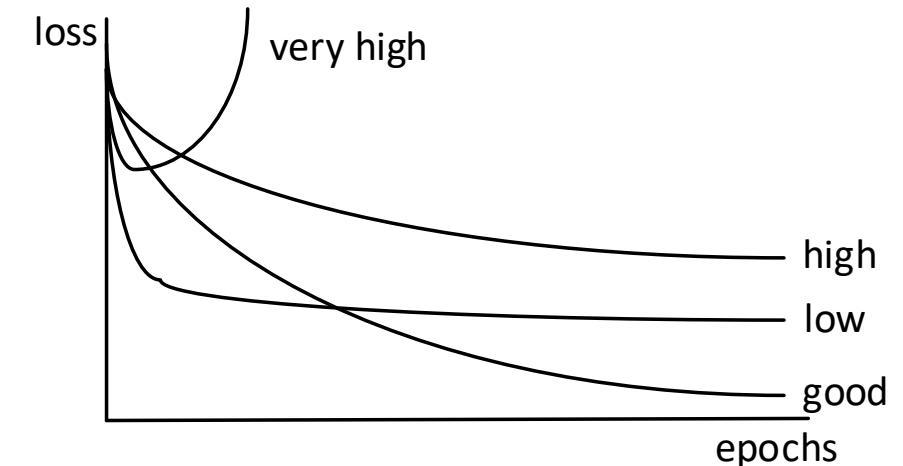
- Compute the expected loss for chance prediction and check if the initial loss is in that range, e.g. classification with 10 class categories and a Softmax classifier:
 - Set the regularisation strength to zero
 - For a randomly initialised network, we expect a diffuse probability of 0.1 for each of the ten class categories
 - The expected cross-entropy loss is then $-\ln(0.1) \approx 2.3$
- Increasing the regularisation strength should increase the loss
- Overfit on a tiny subset of the data, e.g. 20 samples, and try to achieve zero **training loss**
 - Set the regularisation strength to zero

<https://cs231n.github.io/neural-networks-3/>

Monitor the training loss

- Monitoring the **training** loss indicates how well the network architecture learns in general
- A wiggly (noisy) loss function is a symptom of a small batch size
 - Some noise is normal, but if the loss is very noisy you might want to increase the batch size (if possible)
- The loss also indicates if the learning rate is chosen well
- Also, if the initial **training** loss looks good but the **validation** loss is very high → might indicate a too high learning rate as well
- If network training is driving you nuts: <https://lossfunctions.tumblr.com/>

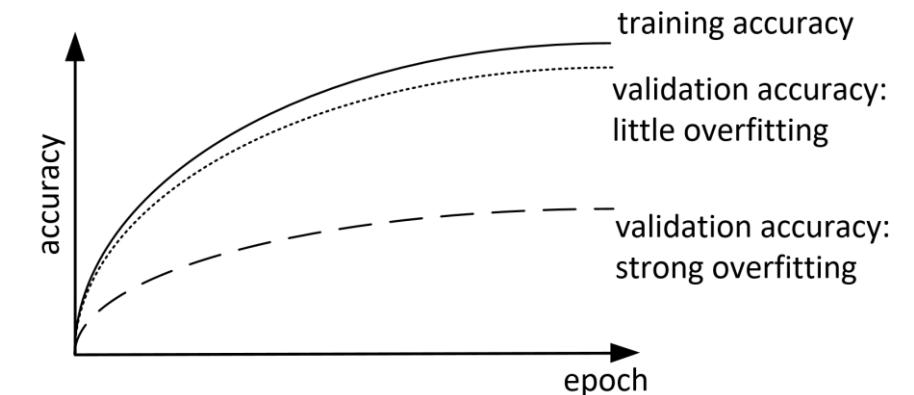
Monitoring the loss may indicate if the learning rate is too high or too low



Babysitting the training process | Training & validation accuracy

Compare training and validation accuracy

- Monitoring the **training** accuracy with respect to the **validation** accuracy indicates whether the model starts to overfit
- Overfitting: the model learns the training dataset by heart and thus generalises badly to new samples. Thus, we are interested in a network with a high validation / test accuracy
- Overfitting is indicated by a high **training** accuracy in combination with a distinctively lower **validation** accuracy
- Overfitting with a very noisy validation accuracy / loss might also indicate that the training dataset is too small (the training data does not represent the validation data well)
- Reducing overfitting:
 - Early stopping (typically, overfitting increases with increasing training time)
 - Regularisation (Dropout, batch norm, L1/L2, ...)
 - More training data
 - Use transfer learning to initialise the shallow layers (or all conv layers, if the dataset is similar)
 - Decrease the number of layers



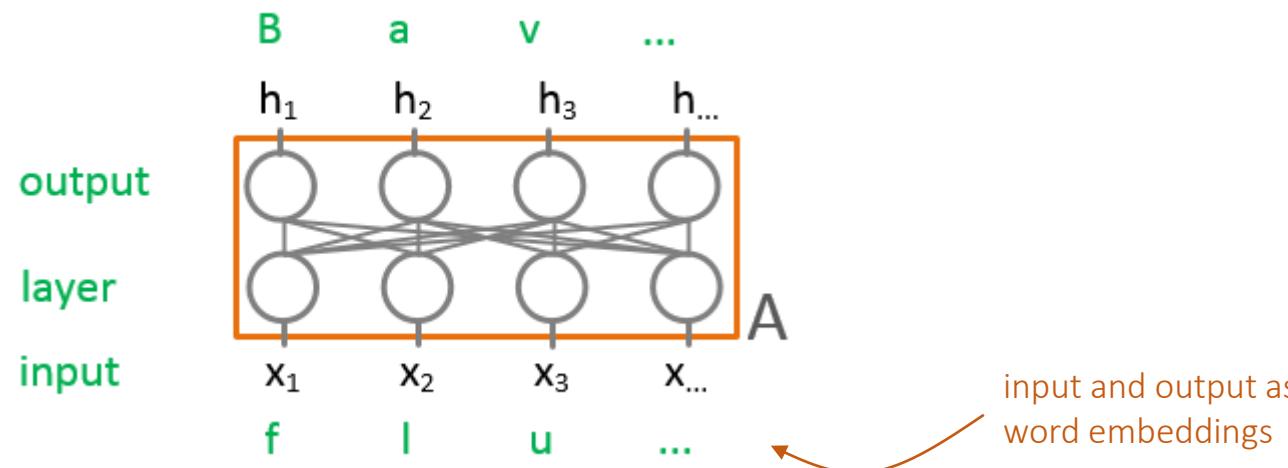
Network Architectures for Natural Language Processing (NLP)



„I grew up in **Bavaria**, in the beautiful city of Aschaffenburg with its castle and Pompejanum and many parks. Of course, I speak fluent **Bavarian**“

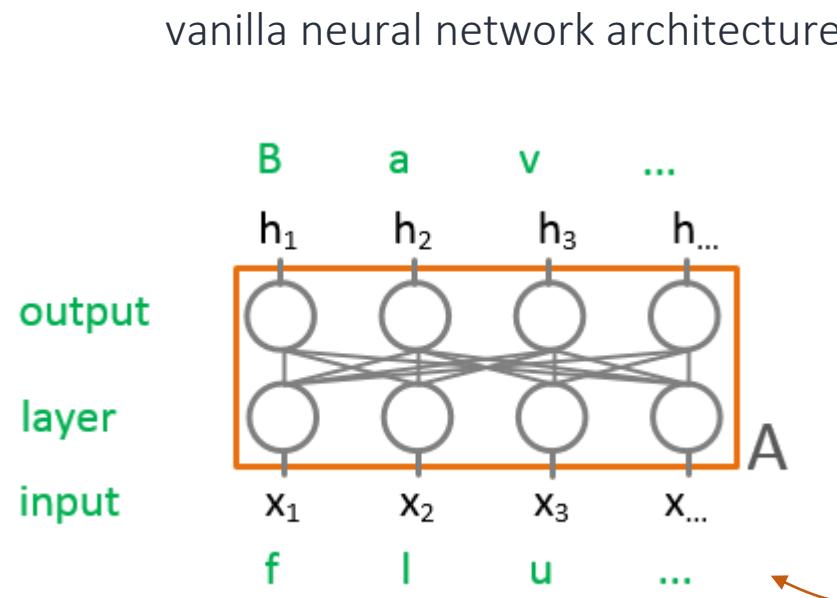
Task: Predict the next word in the sentence

vanilla neural network architecture



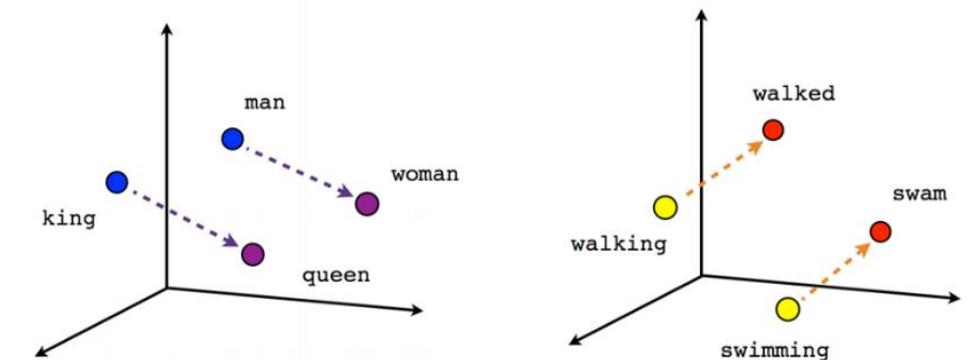
„I grew up in Bavaria, in the beautiful city of Aschaffenburg with its castle and Pompejanum and many parks. Of course, I speak fluent ?“

Task: Predict the next word in the sentence



input and output as
word embeddings

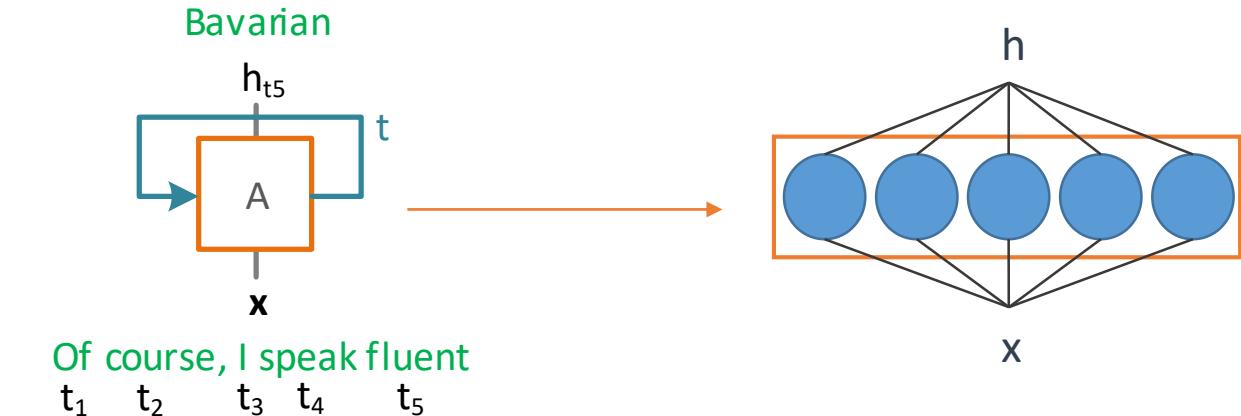
word embeddings (GloVe, Word2Vec)



„I grew up in Bavaria, in the beautiful city of Aschaffenburg with its castle and Pompejanum and many parks. Of course, I speak fluent ?,,

- Task: Predict the next word in the sentence
- Recurrent neural networks have connections which point backwards in time

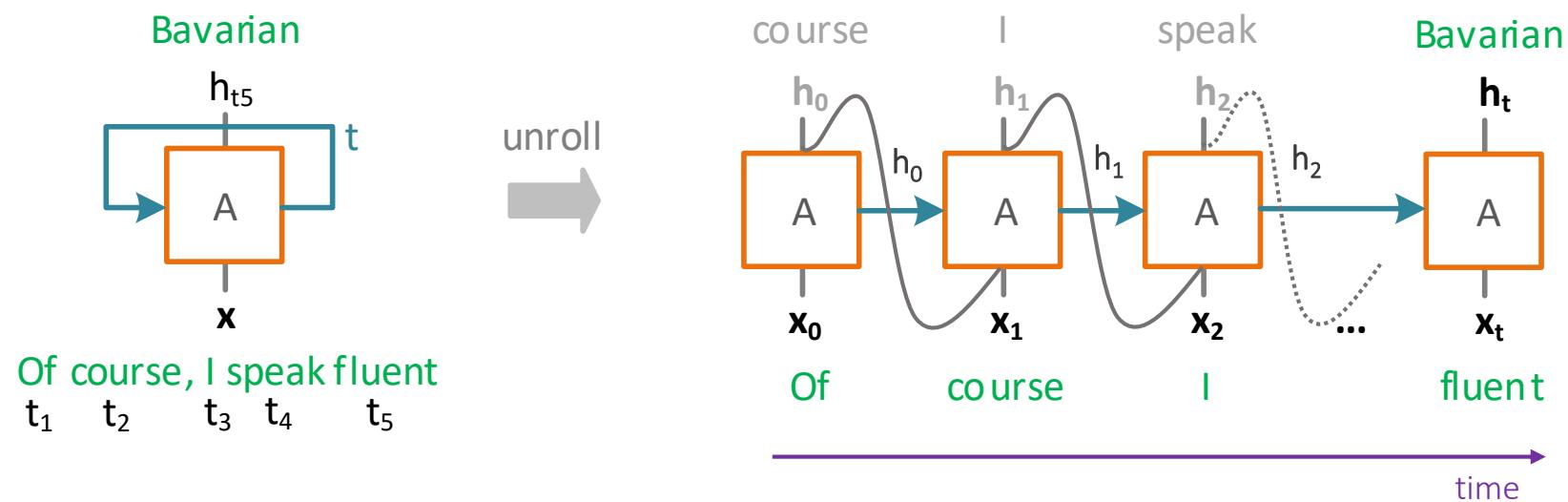
recurrent neural network architecture



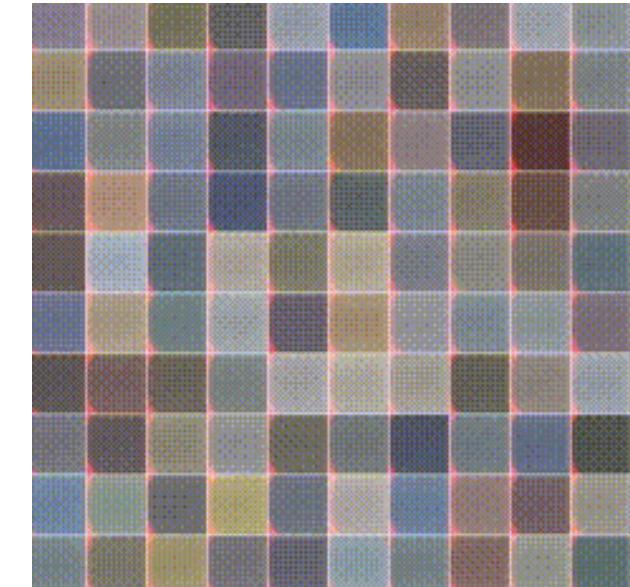
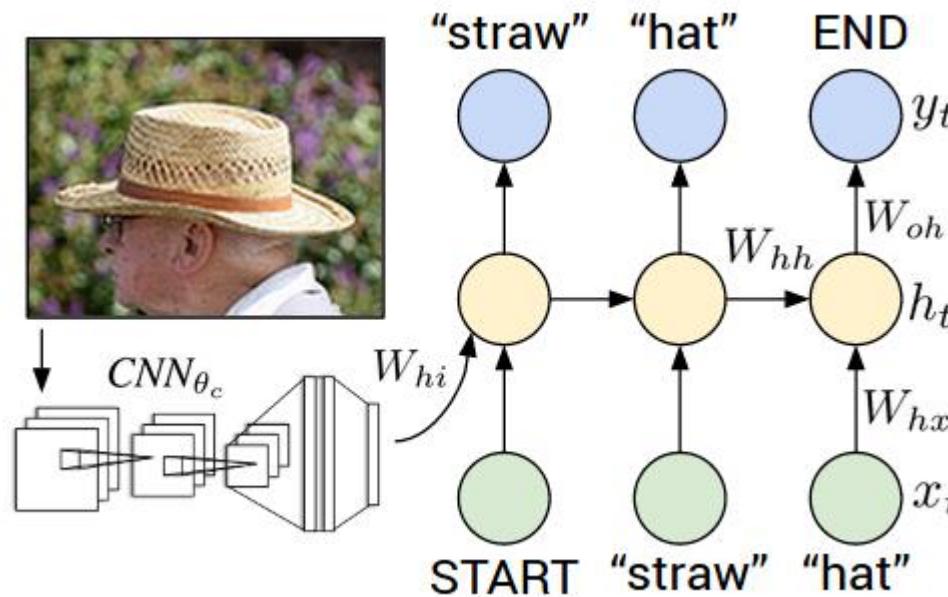
„I grew up in Bavaria, in the beautiful city of Aschaffenburg with its castle and Pompejanum and many parks. Of course, I speak fluent ?,“

- Task: Predict the next word in the sentence
- Recurrent neural networks have connections which point backwards in time
- At each time step t , a recurrent neuron receives its input \mathbf{x}_t and its own output from the previous step, \mathbf{h}_{t-1}
- Accordingly, each recurrent neuron has two sets of weights, one for the inputs \mathbf{x}_t and one for the outputs of the previous step \mathbf{h}_{t-1} : $\mathbf{h}_t = \phi(\mathbf{W}_x^T \mathbf{x}_t + \mathbf{W}_h^T \mathbf{h}_{t-1} + b)$, where is ϕ the tanh or ReLU activation function

\mathbf{h}_t is a function of all the previous inputs since time $t = 0$



Recurrent neural networks are the go-to network architecture for speech recognition, language modeling, translation, image captioning, timeseries, ...

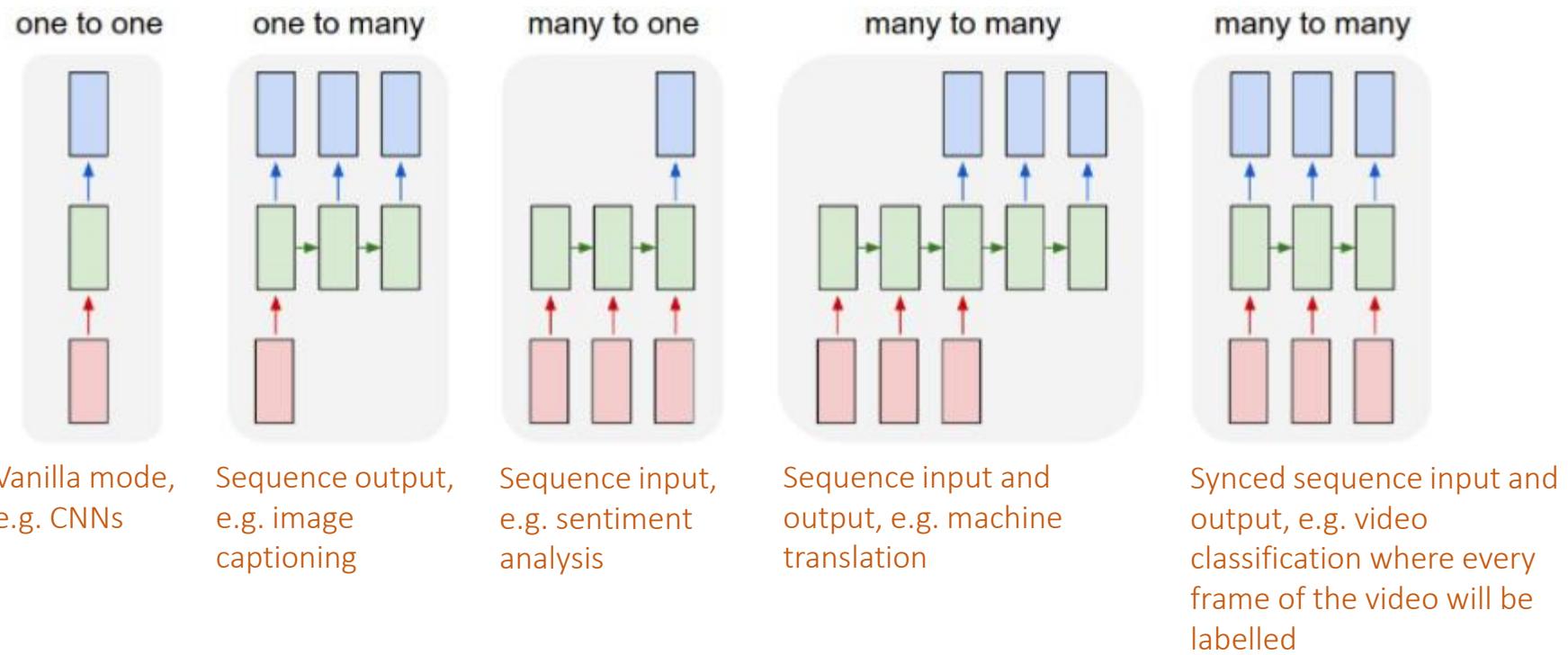


<https://cs.stanford.edu/people/karpathy/>

<https://karpathy.github.io/2015/05/21/rnn-effectiveness/>

<https://arxiv.org/abs/1502.04623>

... because they are flexible regarding input and output dimensions



Training examples: generate new text character-wise, based on Leo Tolstoy's War and Peace

100 iterations

```
tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e  
plia tkldgd t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng
```

700 iterations

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

1200 iterations

"Kite vouch!" he repeated by her
door. "But I would be done and quarts, feeling, then, son is people...."

2000 iterations

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

<https://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Training examples: generate new text character-wise, based on Latex source code for an algebraic stacks book

Vanilla RNNs lack long-term memory

Proof. Omitted. \square

Lemma 0.1. Let \mathcal{C} be a set of the construction.

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. \square

Lemma 0.2. This is an integer \mathcal{Z} is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. \square

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram

$$\begin{array}{ccccc}
 S & \xrightarrow{\quad} & & & \\
 \downarrow & & & & \\
 \xi & \xrightarrow{\quad} & \mathcal{O}_{X'} & \xleftarrow{\quad} & \\
 \text{gor}_s & & \uparrow & & \\
 & & =\alpha' \longrightarrow & & \\
 & & \downarrow & & \\
 & & =\alpha' \longrightarrow & & X \\
 & & \downarrow & & \downarrow \\
 \text{Spec}(K_\psi) & & \text{Mor}_{\text{Sets}} & & \text{d}(\mathcal{O}_{X_{/\mathbb{A}}}, \mathcal{G})
 \end{array}$$

is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

Proof. We have see that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . \square

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}}^{-1}(\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X_\ell}^{-1}\mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\eta}^{\overline{v}})$$

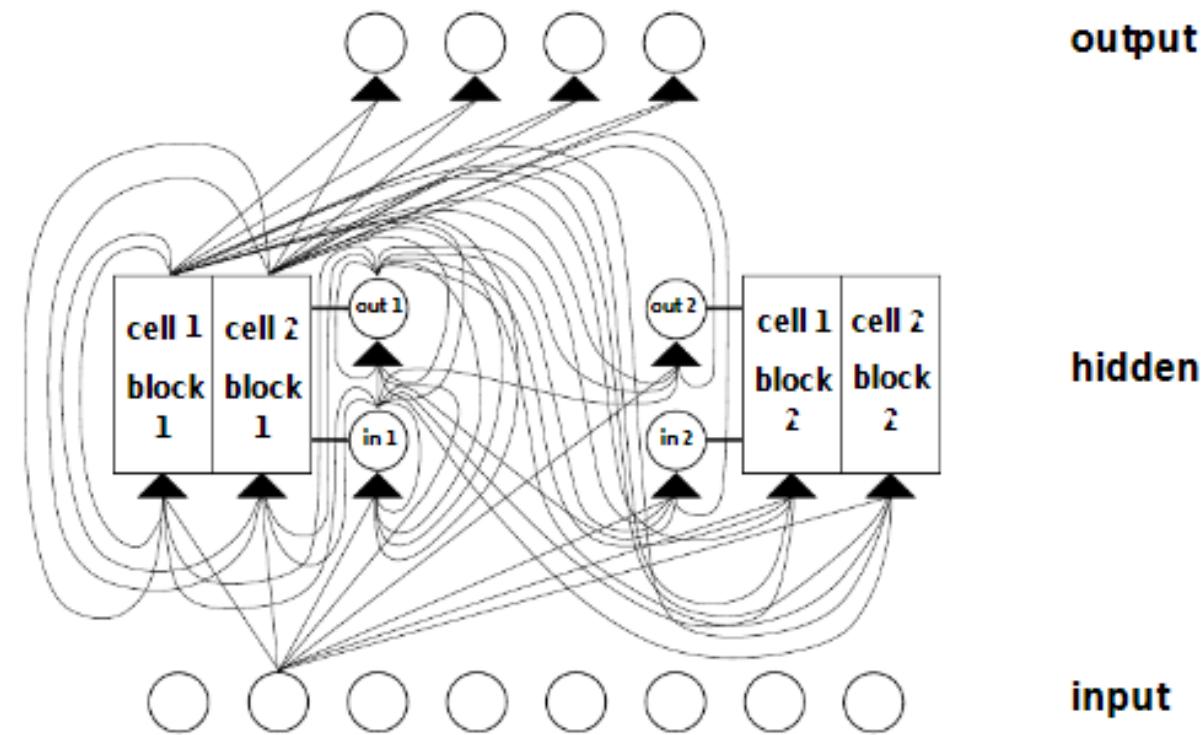
is an isomorphism of covering of \mathcal{O}_{X_ℓ} . If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

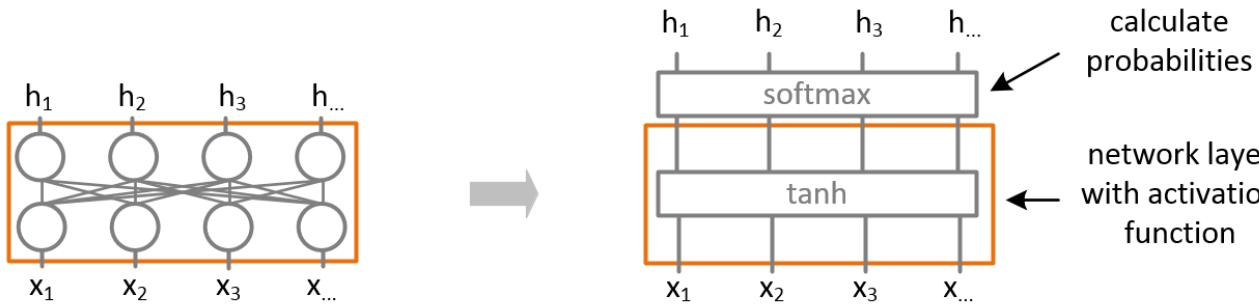
The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S . If \mathcal{F} is a scheme theoretic image points. \square

If \mathcal{F} is a finite direct sum \mathcal{O}_{X_λ} is a closed immersion, see Lemma ???. This is a sequence of \mathcal{F} is a similar morphism.

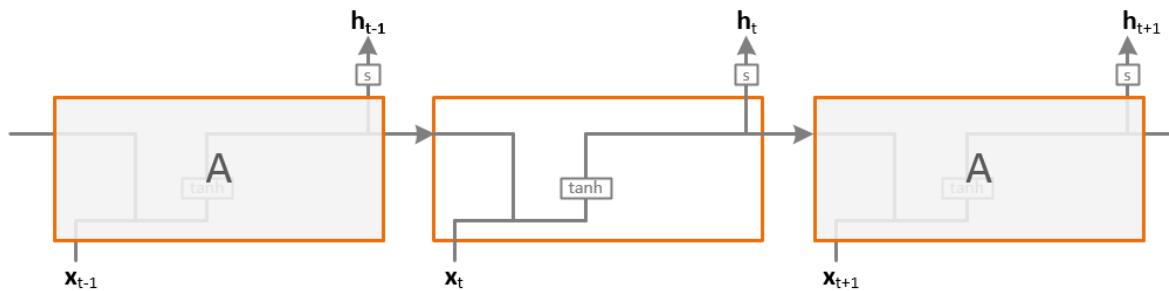
<https://karpathy.github.io/2015/05/21/rnn-effectiveness/>

- Long short-term memory networks improve the memory of recurrent networks

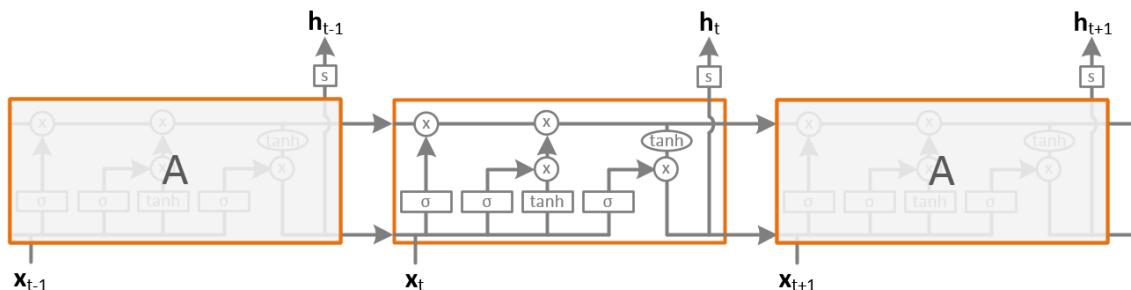




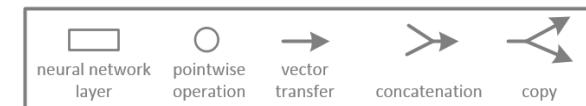
vanilla neural network architecture



vanilla recurrent network architecture

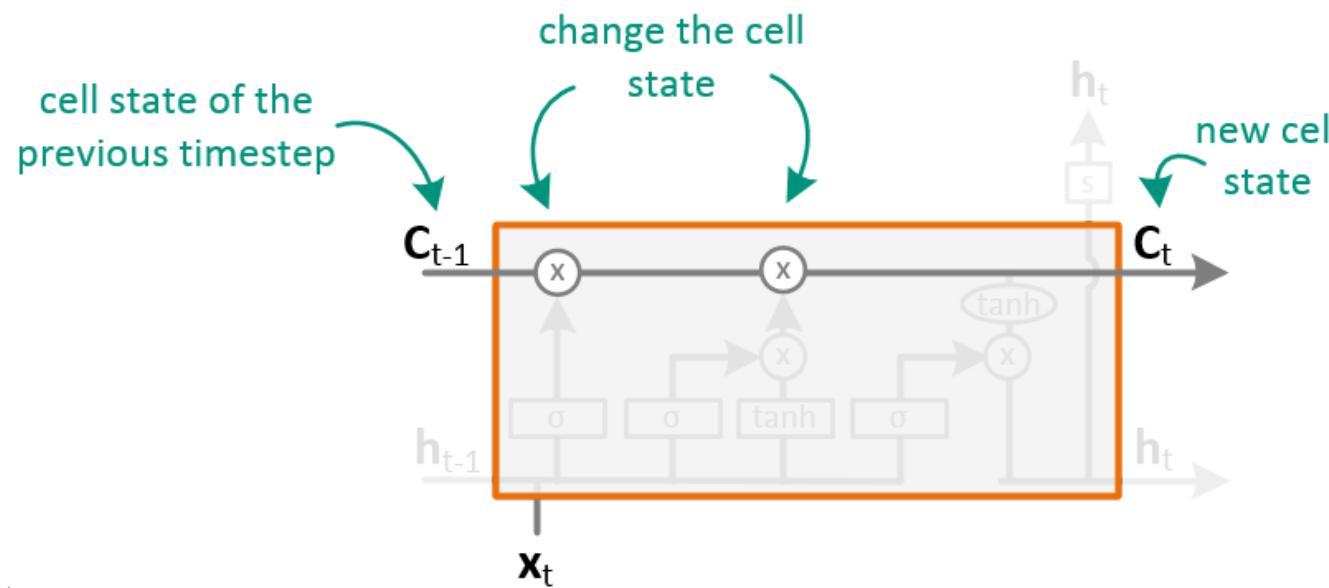


LSTM network architecture



Cell state

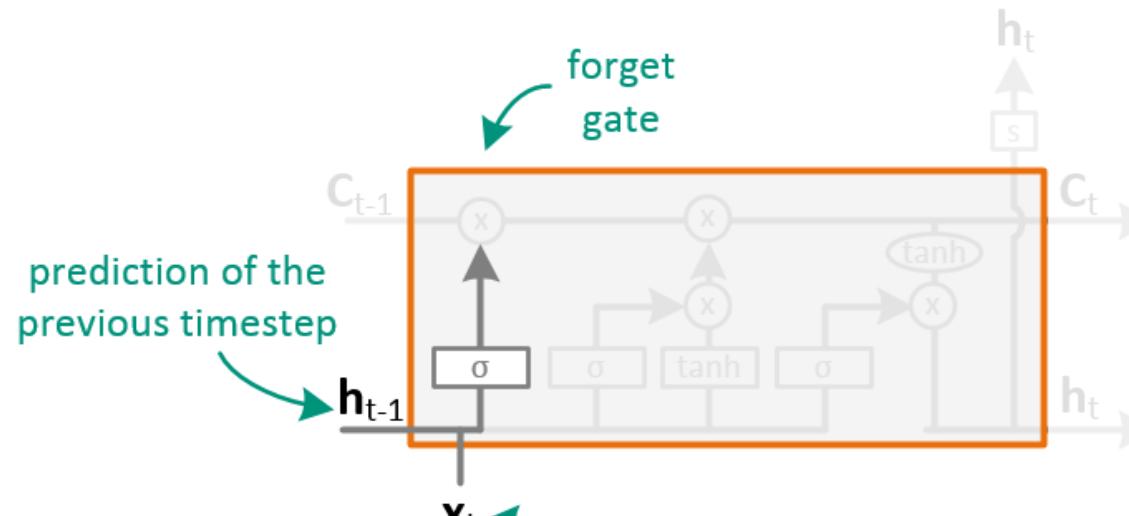
- The cell state stores the information and passes it to the next timestep
- Can be changed by so-called gates



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Forget Gate

- decides which values will be updated based on the current input and previous prediction
- outputs values between 0 (forget everything) and 1 (keep information)



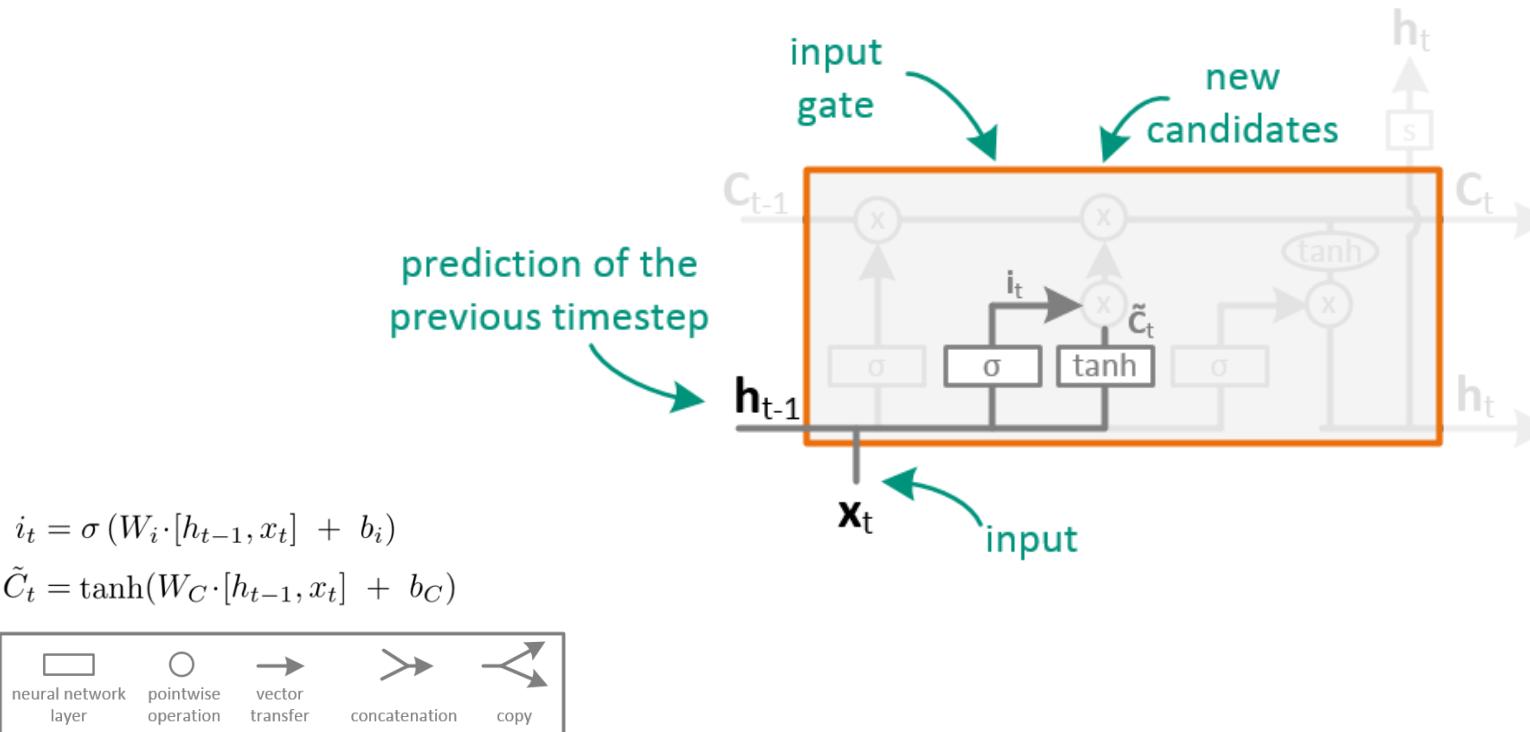
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Update Gate, composed of two parts:

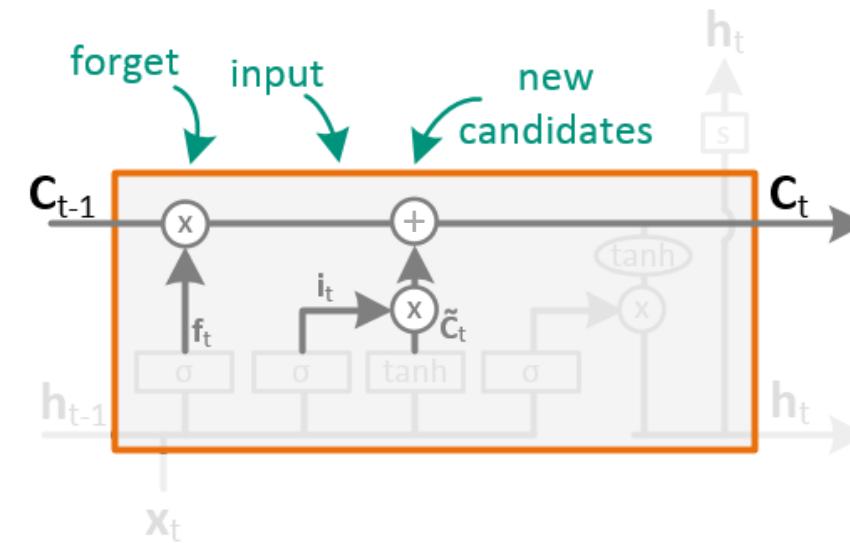
- input gate: decides, which information to update and to what scale
- tanh layer: calculates new candidates



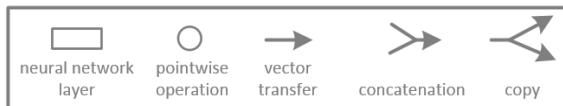
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Update the Cell State

- delete obsolete information from the cell state
- add the selection of new candidates



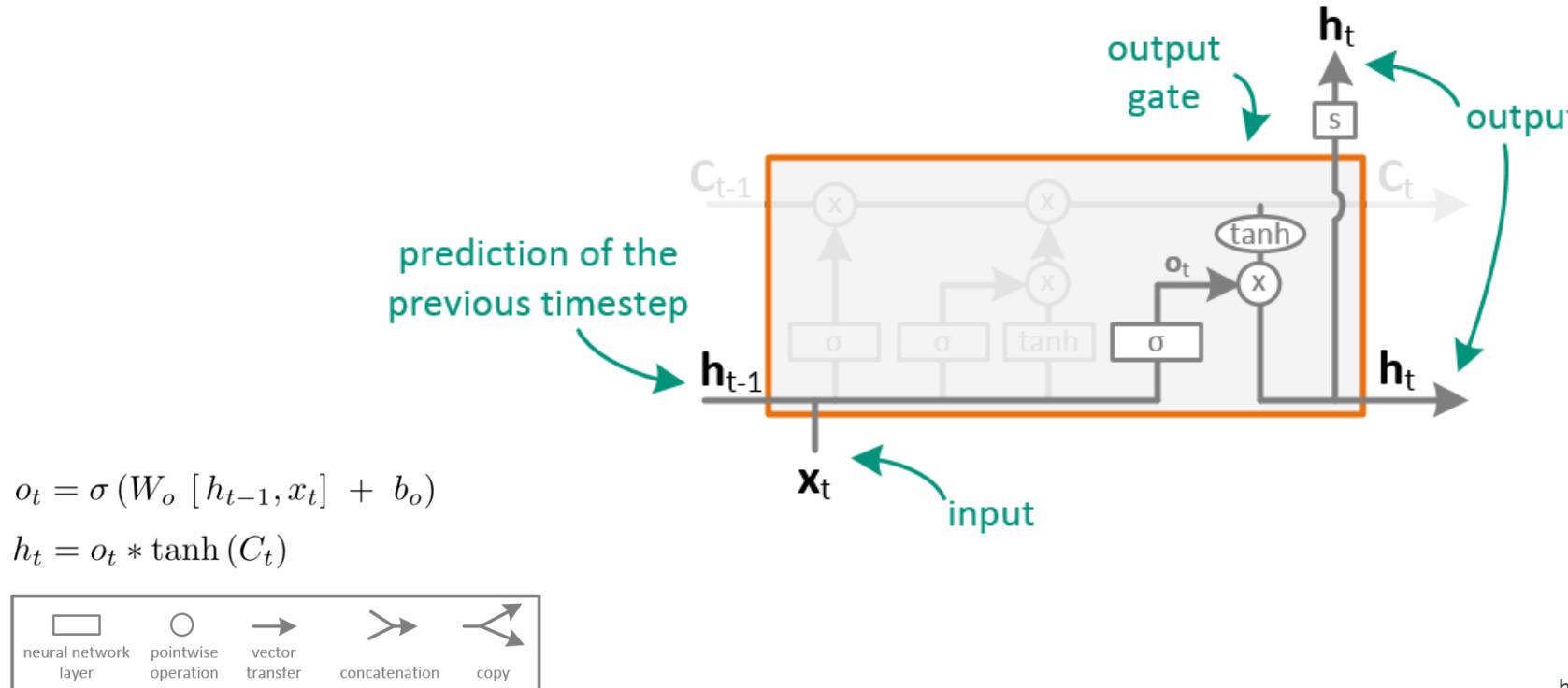
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

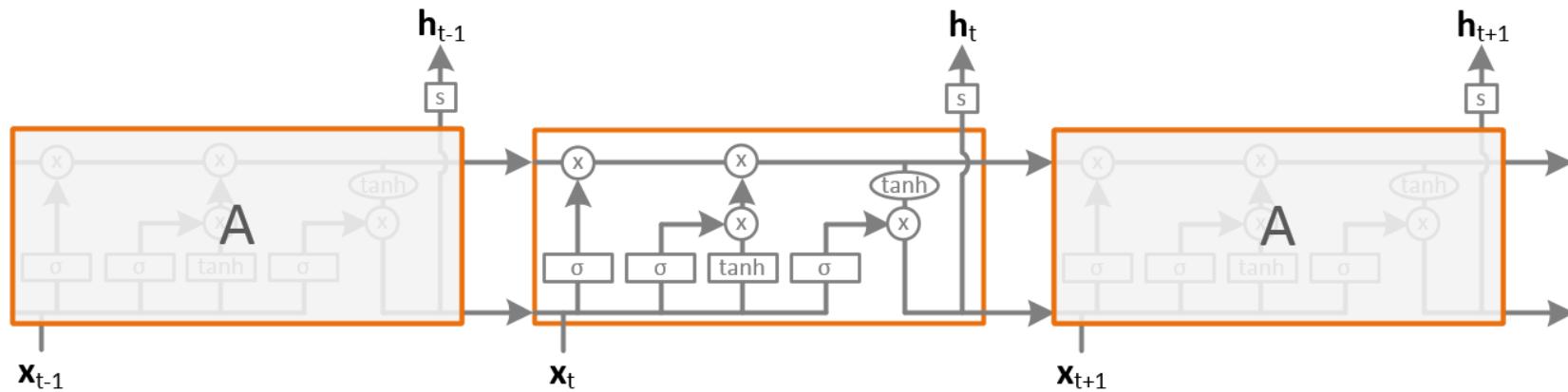
Output Gate

- determines which information to output, based of the previous timestep's prediction and the input
- tanh squeezes the cell state information between -1 and 1



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- There are various adaptations of LSTM cells, such as GRU (Gated Recurrent Unit), which simplify LSTMs
- Which works best depends on the dataset



Current state-of-the-art networks for natural language processing, e.g. GPT-3, T5, BERT, use neither convolutional layers nor recurrent layers but the Transformer architecture, which improves upon the long-term memory

- GPT-NeoX by Eleuther
 - www.eleuther.ai/
- Google BERT and T5
 - https://colab.research.google.com/github/tensorflow/tpu/blob/master/tools/colab/bert_finetuning_with_cloud_tpus.ipynb
 - <https://github.com/google-research/text-to-text-transfer-transformer>
- DistilBERT by Huggingface
 - <https://medium.com/huggingface/distilbert-8cf3380435b5>
- GPT-3 by OpenAI
 - <https://openai.com/blog/openai-api/>



- Transformers are encoder-decoder architectures for sequence-to-sequence (seq2seq) applications, e.g. language translation
- The encoder maps the input (source sentence) into a higher dimensional space (representation)
- This representation is then fed into the decoder, which turns it into an output sentence (target sentence)

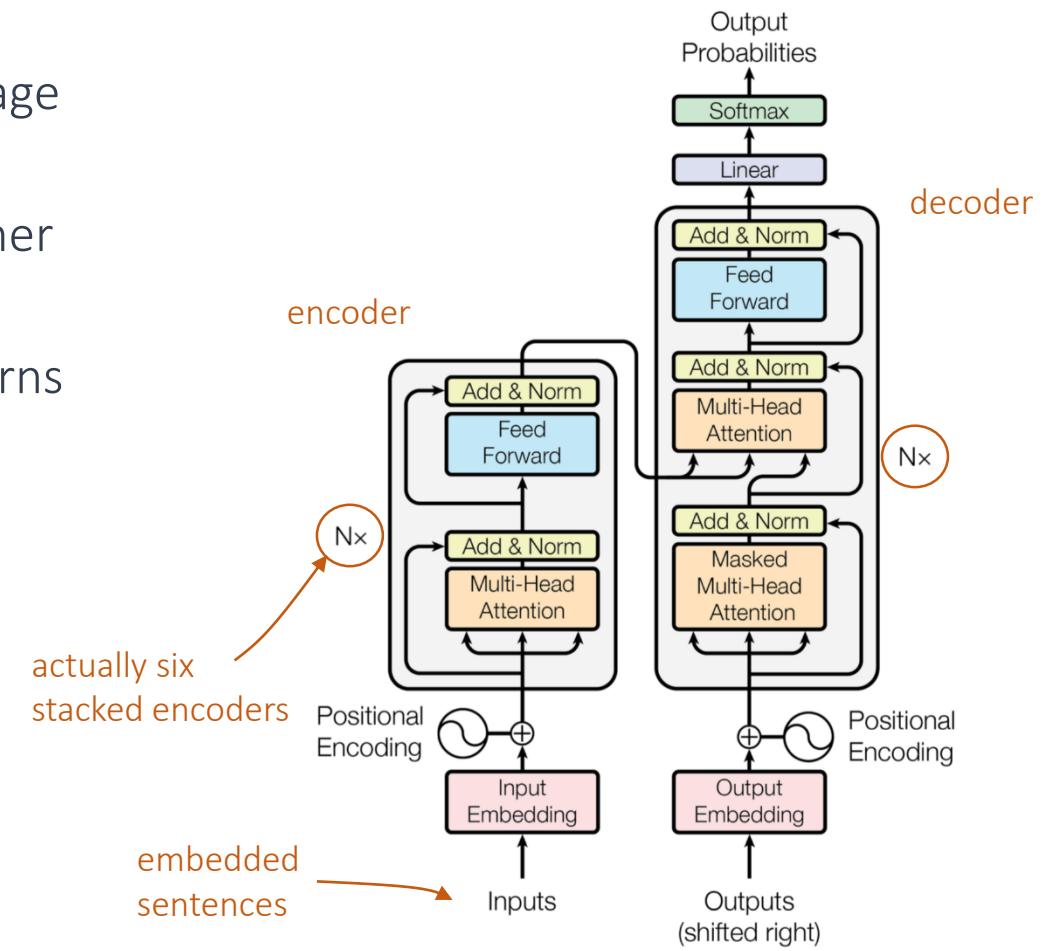
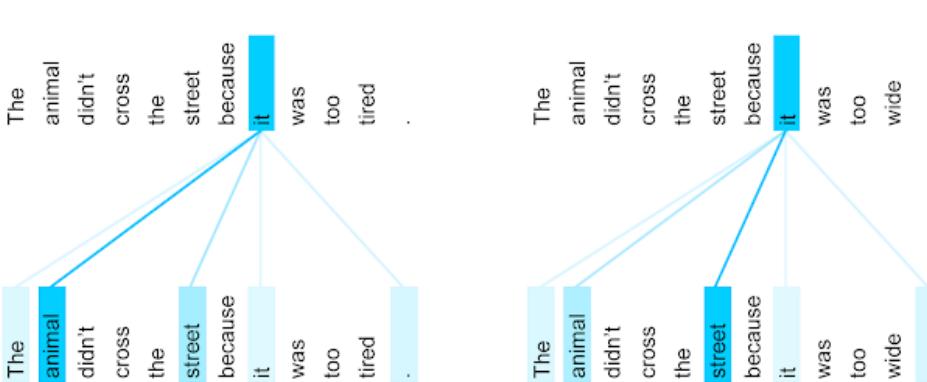


Figure 1: The Transformer - model architecture.

CNNs, Explainable AI, RNNs:

- <https://karpathy.github.io/>
- <https://distill.pub/>
- <https://colah.github.io/>

Computer vision:

- <http://cs231n.stanford.edu/>

Natural Language Processing

- <https://web.stanford.edu/class/cs224n/>
- <https://huggingface.co/course/chapter1>

