



Jonas Glombitza  
jonas.glombitza@fau.de

CTEQ Summer School 2024  
Hotel Idingshof, Bramsche

<https://github.com/DeepLearningForPhysicsResearchBook/deep-learning-physics>



ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



# Convolutional Neural Networks

- I. Processing image-like data
- II. Incorporating symmetries into DNNs



# Time schedule for the next days



ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



## Tutorial: Introduction to deep learning

Tuesday 1h

- Training of deep neural networks

## Hands-on

Tuesday 3:30h

- Interactive training of neural networks
- machine learning frameworks: Keras / TensorFlow
- Implementation of linear regression and fully-connected networks

## Lecture

Wednesday 1h

- Questions + Convolutional neural networks

## Advances in deep learning

Wednesday 3:30h

- Convolutional neural networks
- Transformer networks

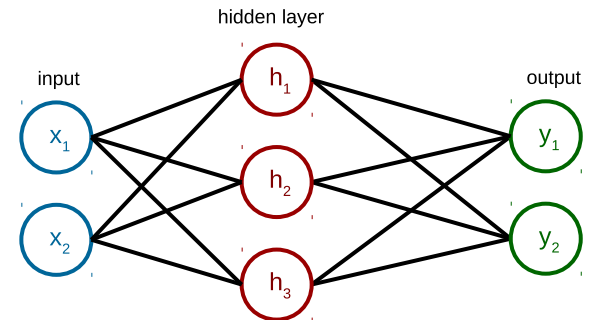
### Set up & Requirements:

<https://bitly.cx/iHcxS> & <https://bit.ly/3pyXRii>

we will use **Jupyter Notebooks** and Keras / TensorFlow  
we will use **Google Colab** → Google Account required

# Recap: Neural Networks

- Typical Machine Learning task
  - Labeled Data  $(x, y)$
  - Model with adaptive weights  $y_m(x, \theta)$
  - Objective  $J(\theta)$
  - Optimization procedure
- Neural Network  $y = \sigma(Wx + b)$ 
  - Matrix multiplication (adaptive superposition of features)
  - Add of bias
  - Nonlinear activation
- Deep Learning is form of representation learning
  - Stack multiple layers for increasing feature hierarchy



# Recap: Optimization



Why Momentum Really Works, Distill

**Mini-batch training:** Use **stochastic gradient descent** algorithm (SGD)

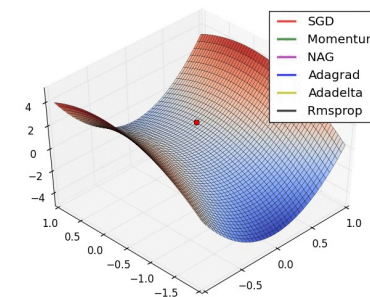
**Momentum:** Use past gradients (velocity)

- Faster convergence by **damping oscillations** and increasing the step size for more informative gradients

**Adaptive learning rate:** Scaling using past gradients

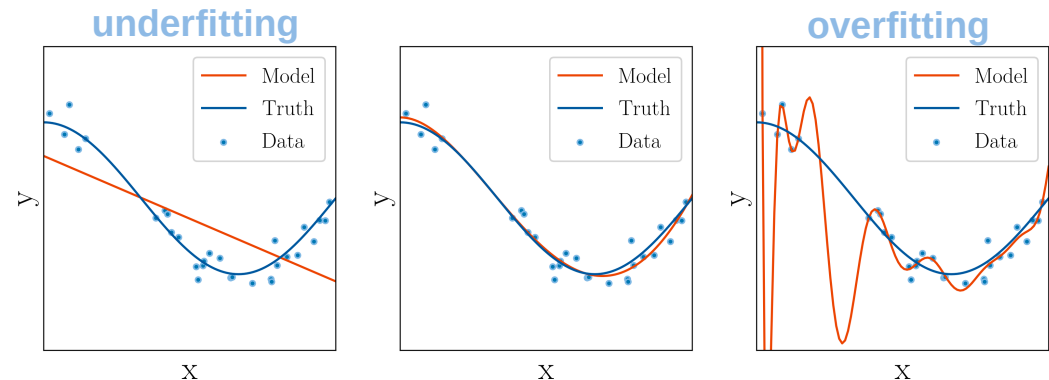
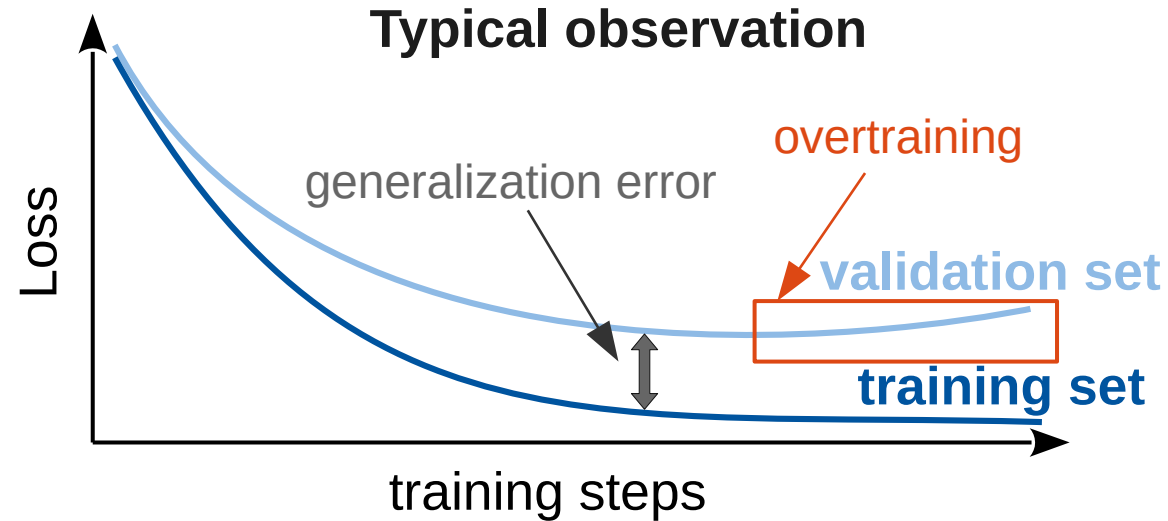
- Use adaptive learning rate for each parameter

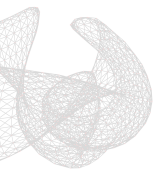
*“Friends don’t let friends use minibatches larger than 32” - Yann LeCun*



# Recap: Under- and Overtraining

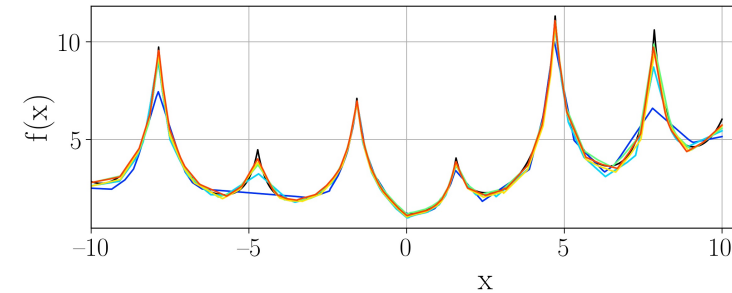
- Split data into 3 sets
  - training
  - test
  - validation
- What is a reasonably split?
- During training monitor the loss separately for training and validation set
  - check for overtraining!
  - if loss stops decreasing
    - reduce learning rate
    - stop training





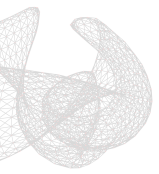
*“A feed-forward network with a linear with a finite number of nodes can approximate any reasonable function to arbitrary precision.”*

## Is the simple feed-forward network the ultimate AI building block?



- (a) Yes, as proven by the universal approximation theorem
- (b) Practically not; You know theory and theorists...
- (c) Please don't talk the next 4 hours about feed-forward networks only...





ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



## Dedicated time for **Your Questions**





Automate task for humans, very challenging for machine learning models:

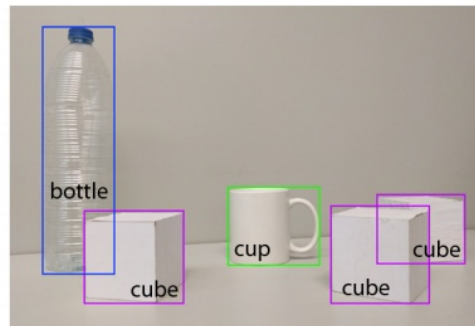
- High dimensional input (up to millions of pixels)
- Many possible classes depending on task
- Multiple variations
  - ♦ Viewing angle, light conditions, deformation, object variations, occlusions....



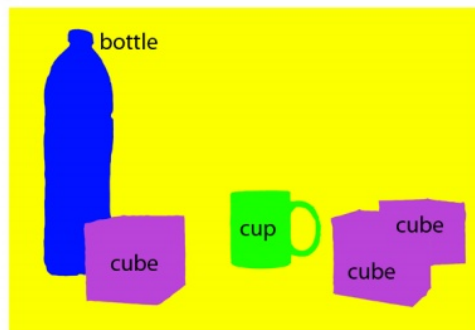
# Computer Vision Tasks



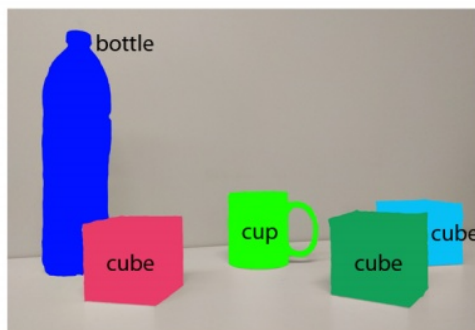
(a) Image classification



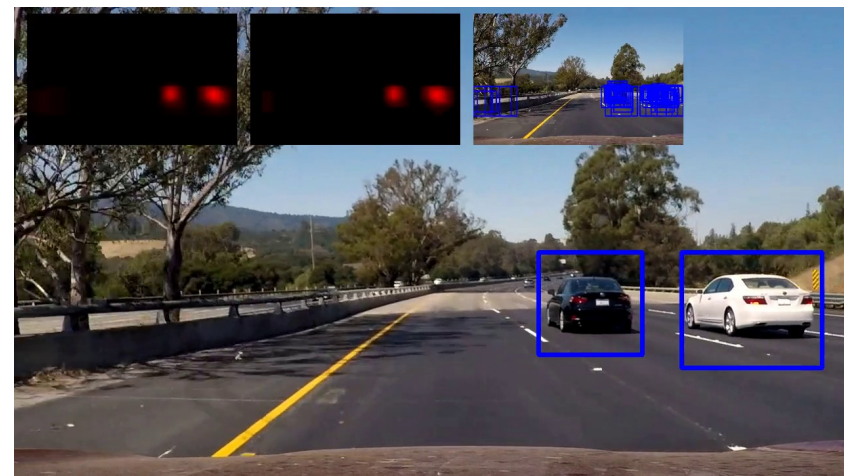
(b) Object localization



(c) Semantic segmentation

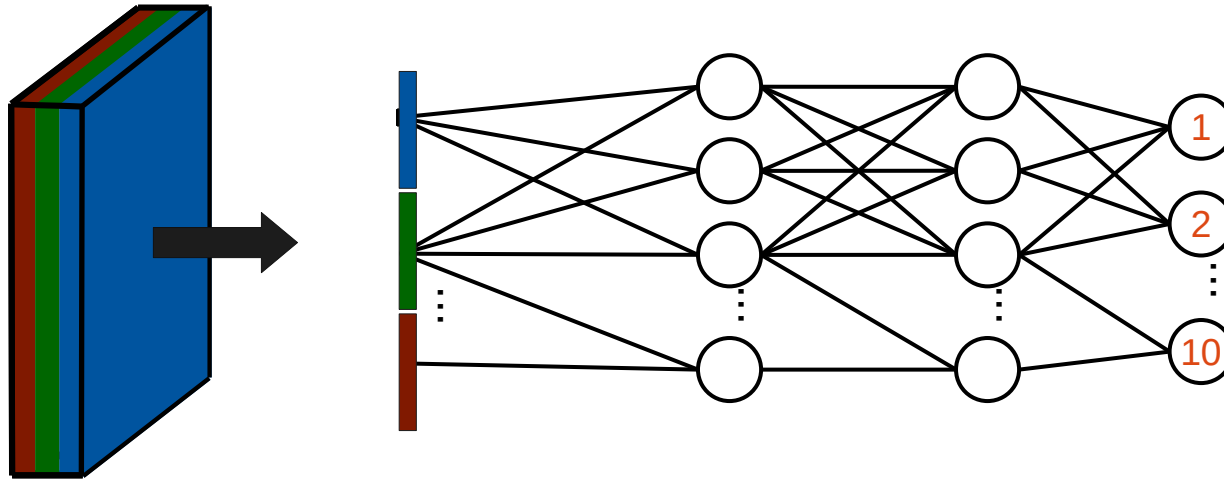


(d) Instance segmentation

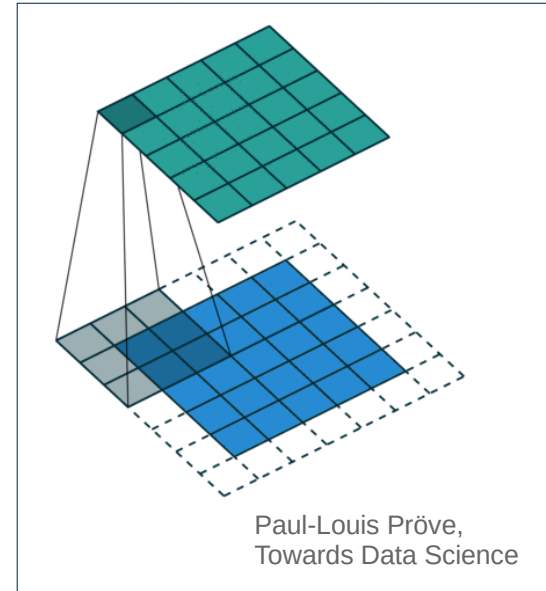
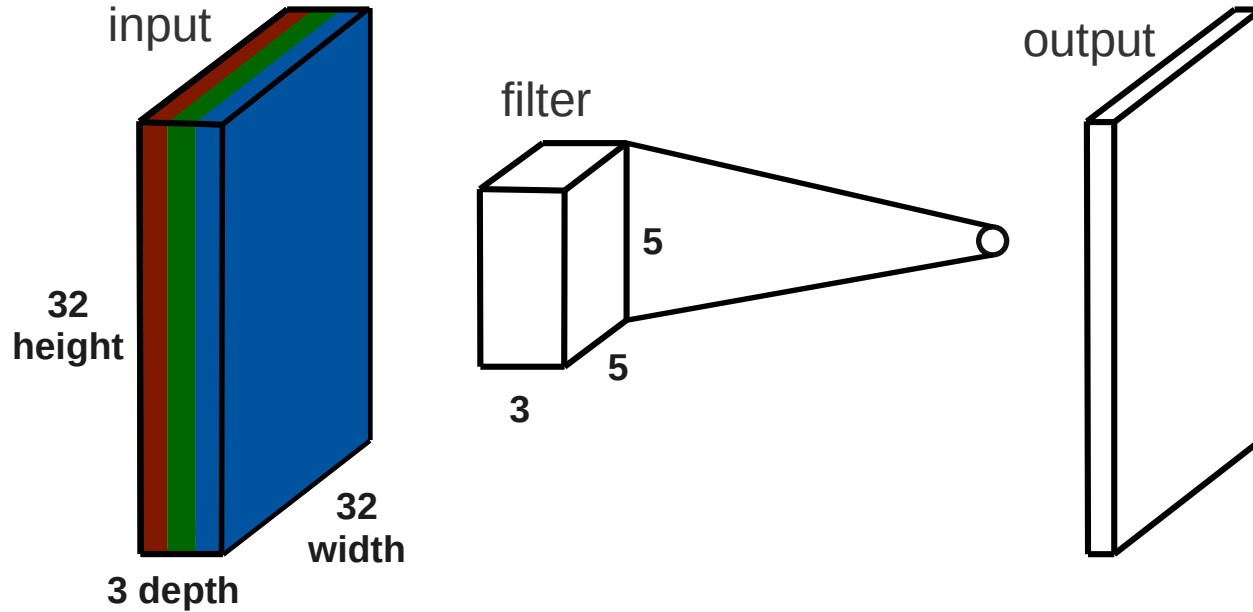


# Fully Connected Network

- Input layer: Flatten image to  $32 \times 32 \times 3 = 3072$  vector
- Fully connected: every pixel connected with each other
- × Huge number of adaptive parameters per layer
- × No use of translational variance
- × No prior on local correlations



# 2D Convolutional Neural Networks



- Consider input volume (width x height x depth), e.g., 3 color channels
- Use convolutional filter with smaller width and height but same depth
- Slide filter over the entire volume and calculate linear transformation to get one output value for each position

# Convolutional Operation

3	0	1				
4	2	0				
2	4	-3				

$$3 \cdot -1 + 0 \cdot 2 + 1 \cdot 0 + 4 \cdot 0 + 2 \cdot 3 + 0 \cdot 0 + 2 \cdot 0 + 4 \cdot 2 + -3 \cdot -5 = 26$$

-1	2	0
0	3	0
0	2	-5

$W_1$	$W_2$	$W_3$
$W_4$	$W_5$	$W_6$
$W_7$	$W_8$	$W_9$

26				

# Convolutional filters

hand-designed filters

Edge	-1	-1	-1
	-1	8	-1
	-1	-1	-1

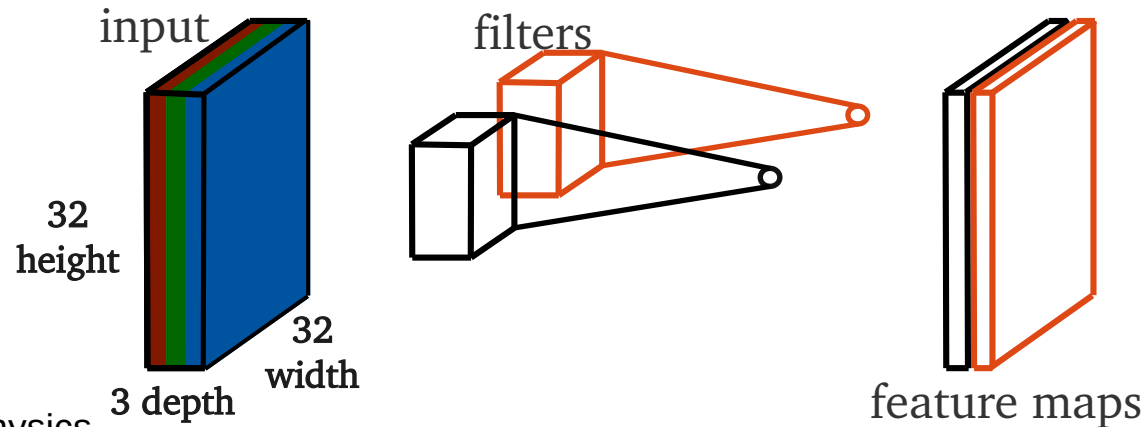
Diagonal edge	-1	-1	2
	-1	2	-1
	2	-1	-1

deep learning

Convolutional Networks	$w_1$	$w_2$	$w_3$
	$w_4$	$w_5$	$w_6$
	$w_7$	$w_8$	$w_9$

adaptive parameters

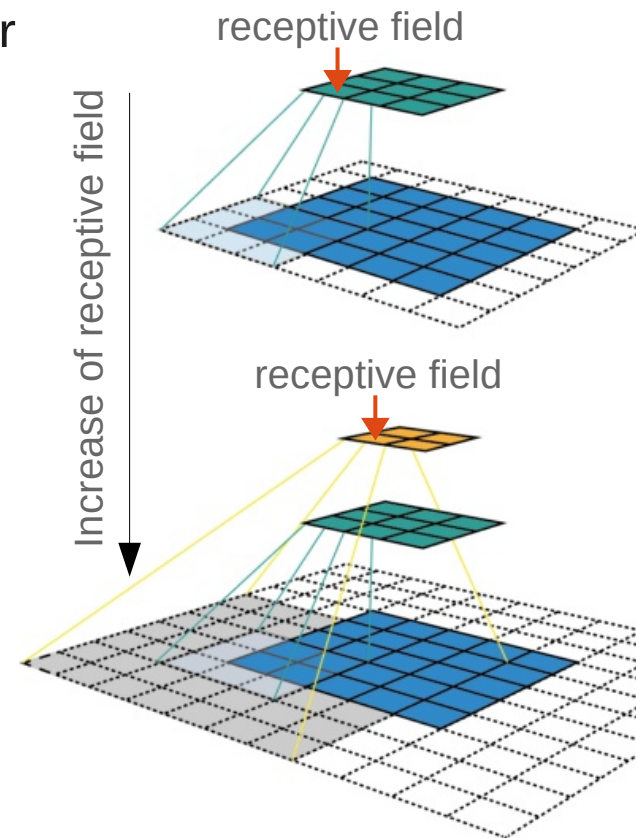
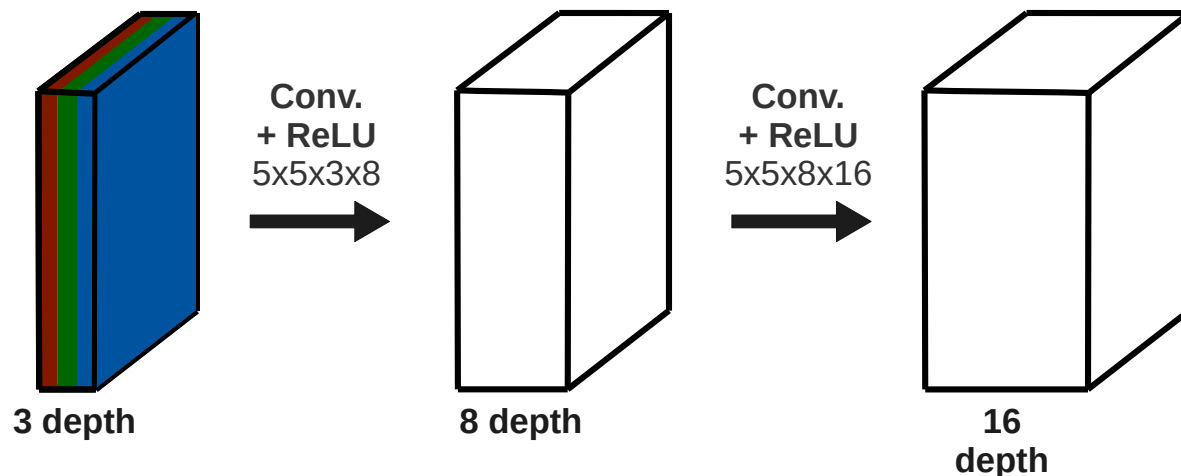
- scan input image for the presence of specific feature using **filters**
- use multiple filters and stack the results as **feature maps** (depth-wise stacking)



# 2D Convolutional Operation

Stack multiple convolutional layers + activations

- Each convolution acts on feature map of previous layer
- Increasing feature hierarchy
- Increasing of receptive field

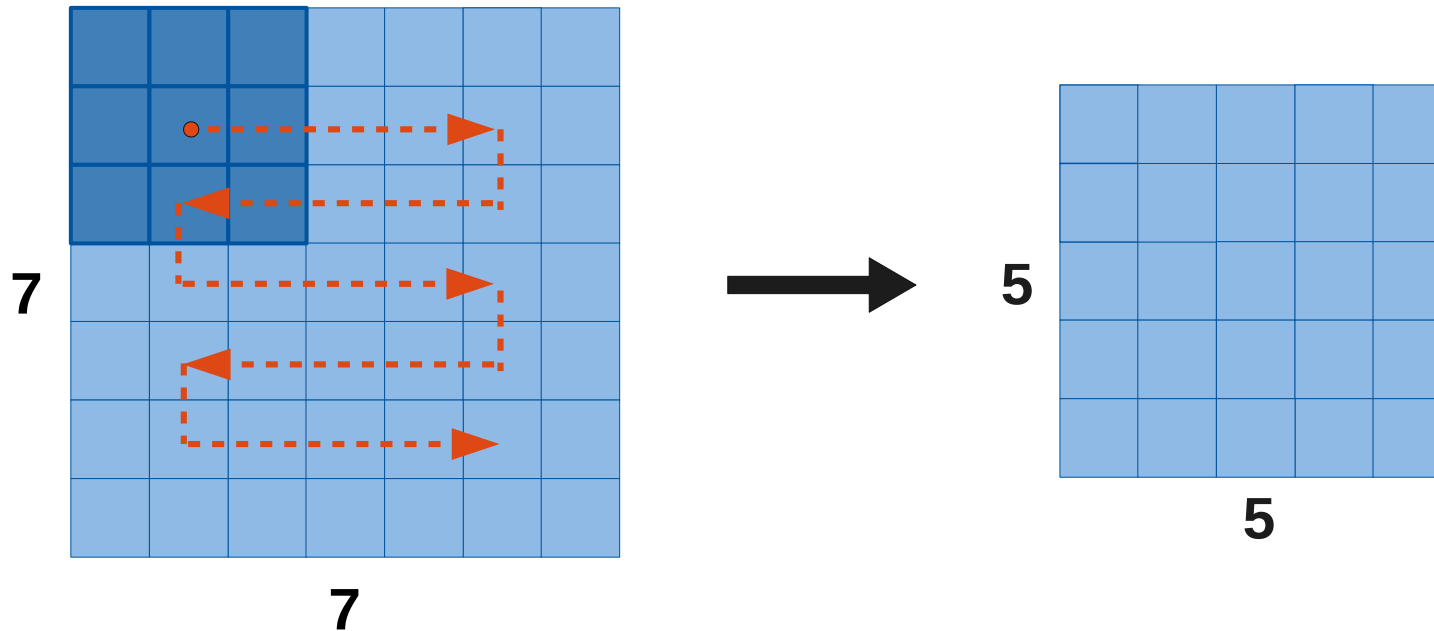


# Spatial Output Size

Standard convolution reduces the output size due to extent of the filter

- Sets upper bound to the number of convolutional layers

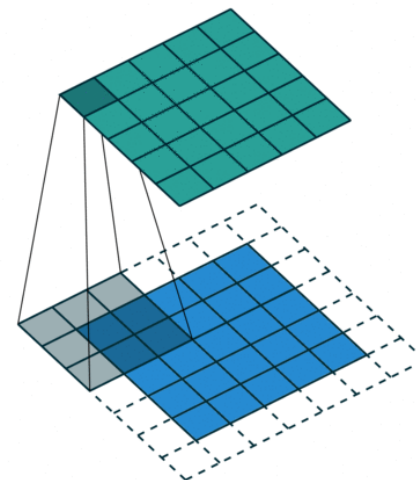
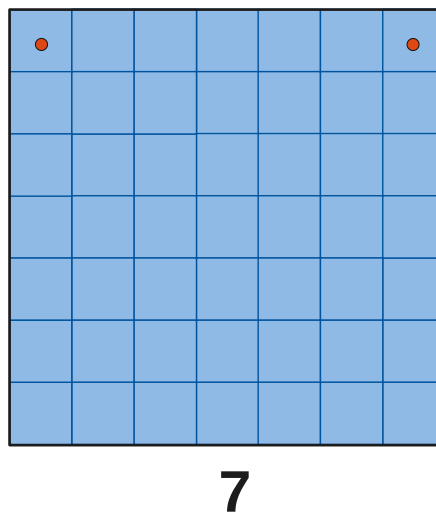
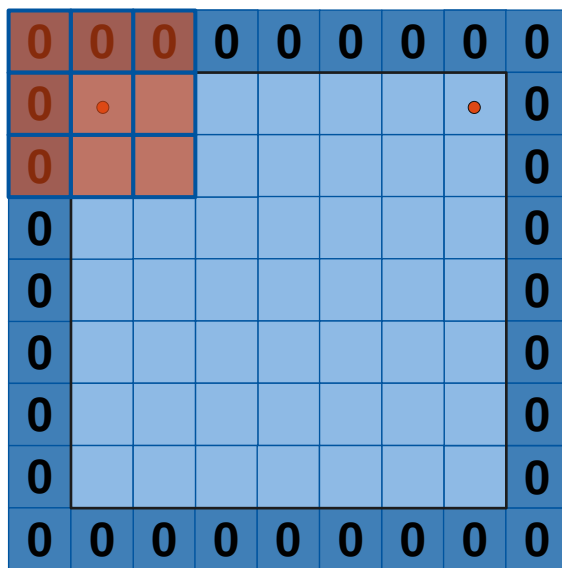
- **Example:** Convolution with 3 x 3 filter



# Padding

Add zeros around image borders to conserve the spatial extent of the input

- Prevents fast shrinking of the network input
- **Example:** Convolution with 3 x 3 filter and padding



Paul-Louis Pröve,  
Towards Data Science

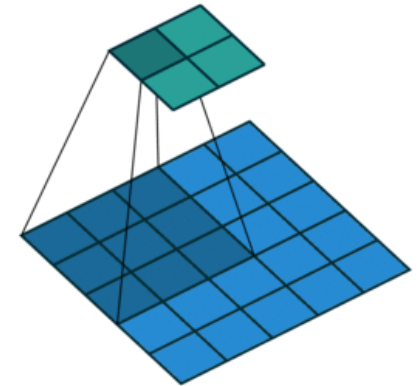
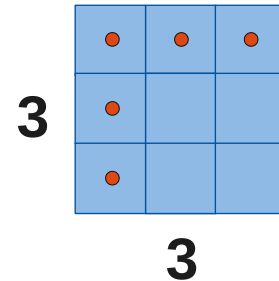
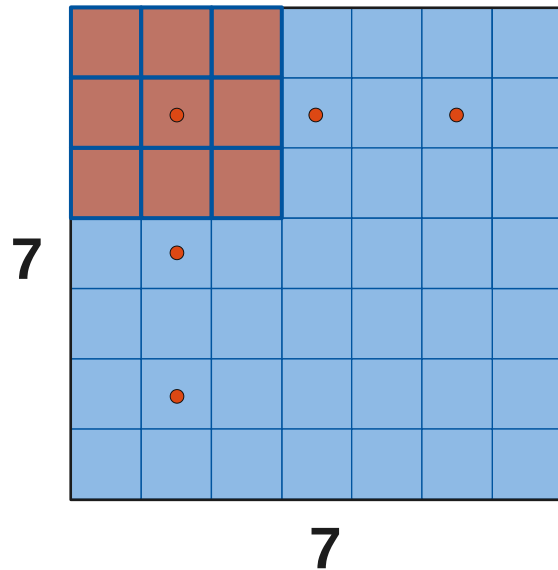


# Striding

Using a larger stride when sliding over the input, reduces the output size

➤ Useful for switching to smaller image sizes / larger scales

• **Example:** Convolution with 3 x 3 filter and stride of 2

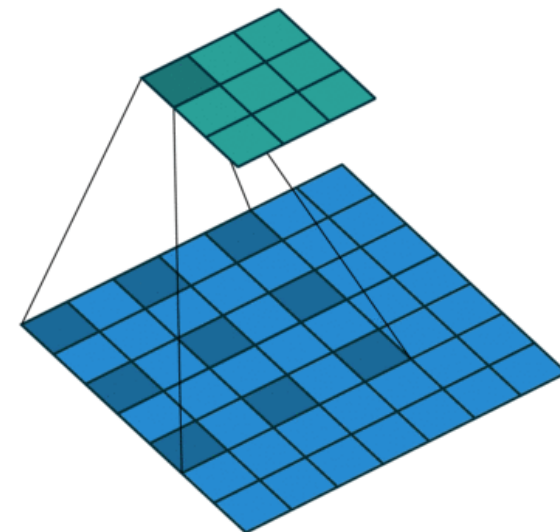
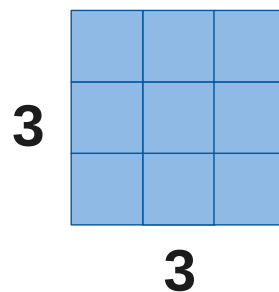
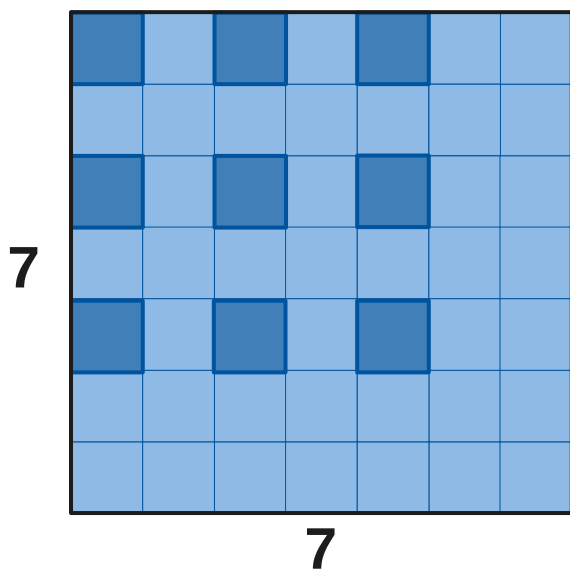


Paul-Louis Pröve,  
Towards Data Science

# Dilating

Dilation leaves holes in where the filter is applied (also called **atrous convolution**)

- Useful for aggressively merging spatial information in large images
- Allows for a large field of view
- **Example:** Convolution with 3 x 3 filter and dilation 1



# Pooling

Sub-sample the input to reduce the output size

- Used to merge semantically similar features
- Make network invariant to small translations or perturbations

**Average pooling:** Take the mean of each patch → for some regressions preferable

**Max pooling:** Take the maximum of each patch

→ in practice often better performance, applies stronger constraint

- **Typical Pooling:**

Pooling using 2 x 2 patches  
and a stride of 2

- **Overlapping Pooling:**

3 x 3 patches with stride of 2

3	2	1	1
0	5	3	-1
9	4	3	2
2	1	3	2

max pooling



5	3
9	3

average pooling



2.5	1
4	2.5

# Global Pooling Operation

- Take maximum/average over complete image → usually second last layer
- Replace fully connected layers
  - ◊ Saves parameters in later layers of the models → prevent overfitting
- Can be seen as regularizer
  - ◊ Fully connected transformation matrix with diagonal shape
- Enforcing correspondences between feature maps and categories
- Allows object detection in the input space

*“The pooling operation used in convolutional neural networks is a big mistake, and the fact that it works so well is a disaster”*

*- Geoffrey Hinton*

3	2	1	1
0	5	3	-1
9	4	3	2
2	1	3	2

max pooling



9

average pooling



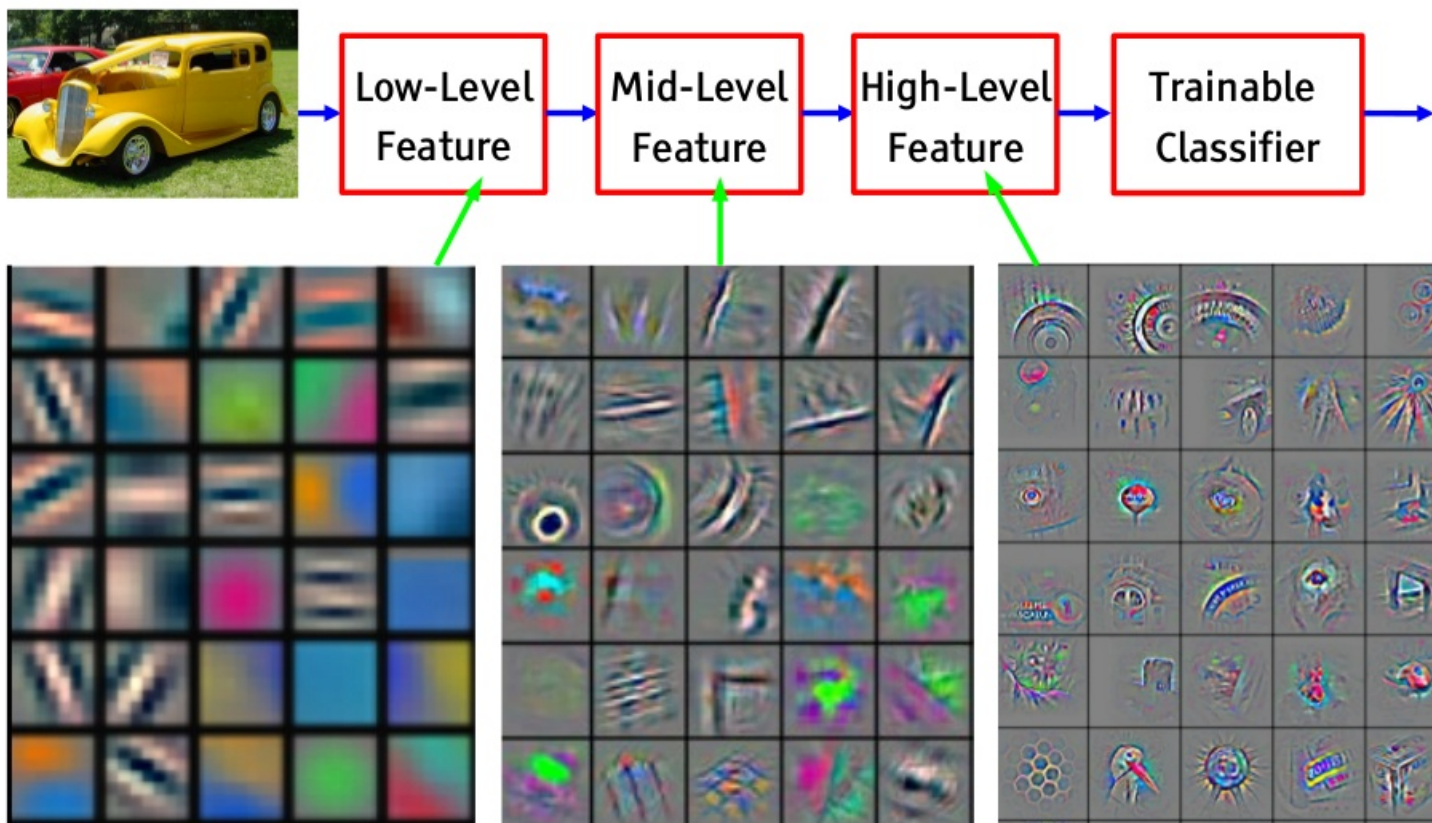
2.5

# Visualization of CNNs – 5 mins



- Inspect the CNN → is the network translational invariant or equivariant
  - ♦ invariant: translation of the image yield exactly the same output?
  - ♦ or only equivariant? (DNN can reconstruct position of a pattern)

# Feature Hierarchy



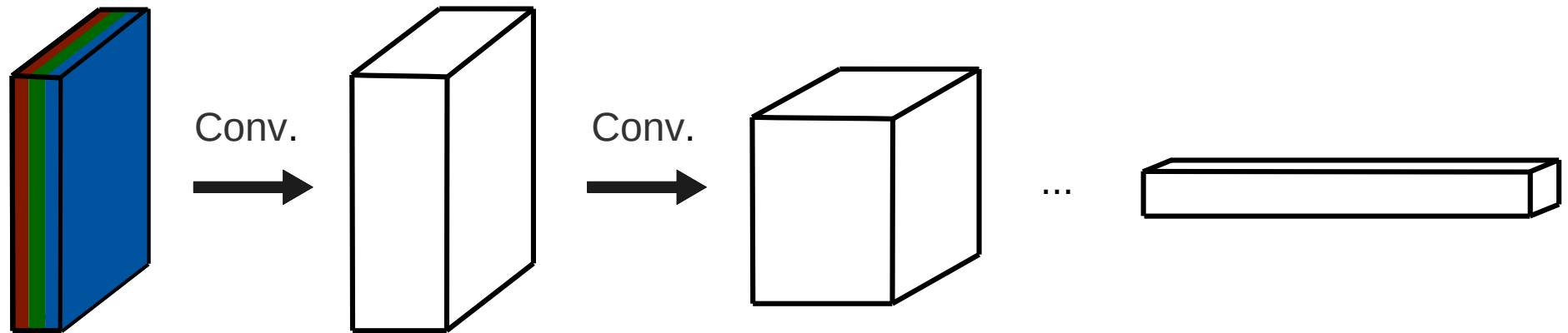
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

<https://arxiv.org/abs/1311.2901>

# Convolutional Pyramid

ConvNet architectures usually have a pyramidal shape. For deeper layers:

- Increasing of feature space
- Decreasing of spatial extent

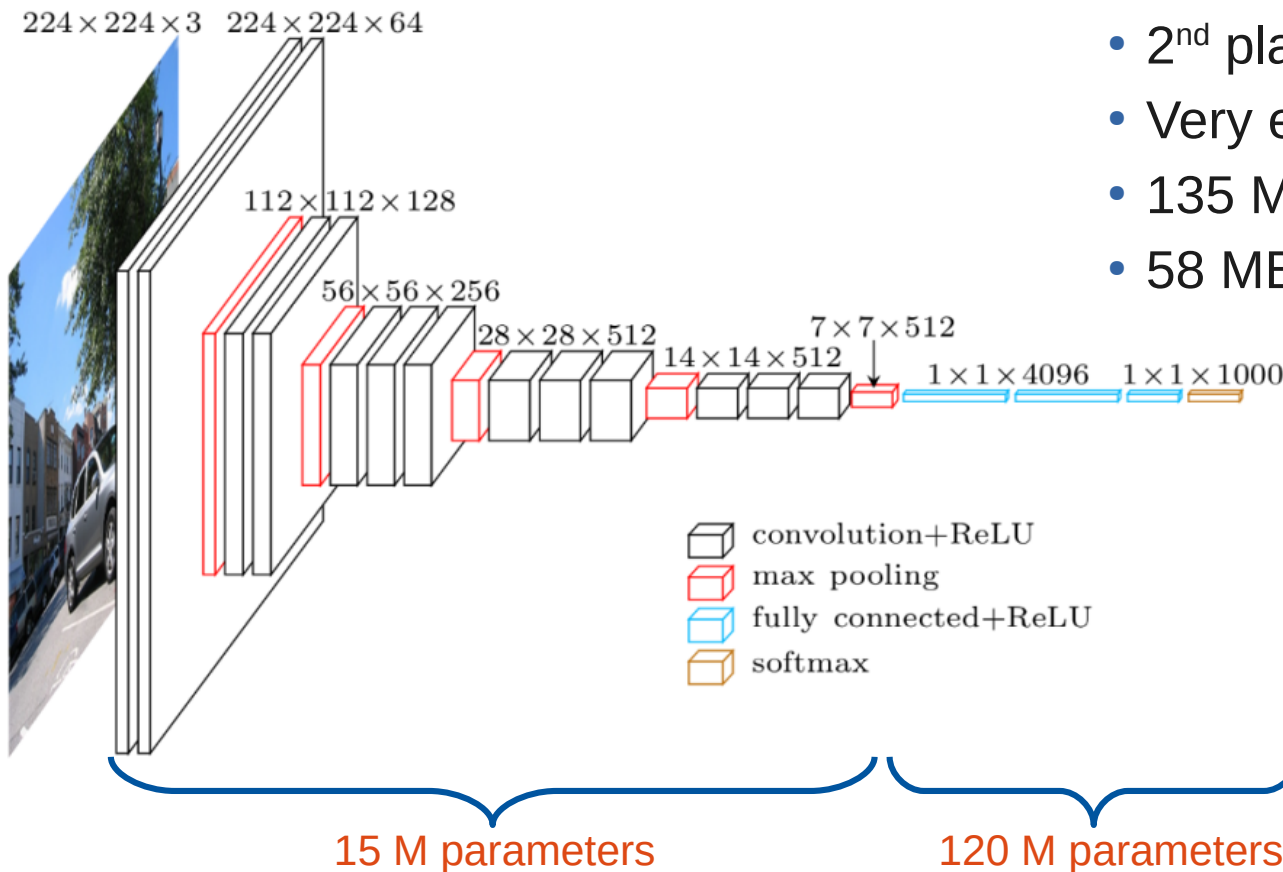


- Spatial information is converted to representational features with increasing hierarchy

# Example Architecture

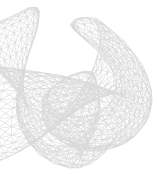
## VGGNet 16

- 2<sup>nd</sup> place ILSVRC2014
- Very easy structure
- 135 M parameters → 530 MB
- 58 MB memory for activations



Simonyan, Zissermann  
<https://arxiv.org/abs/1409.1556>





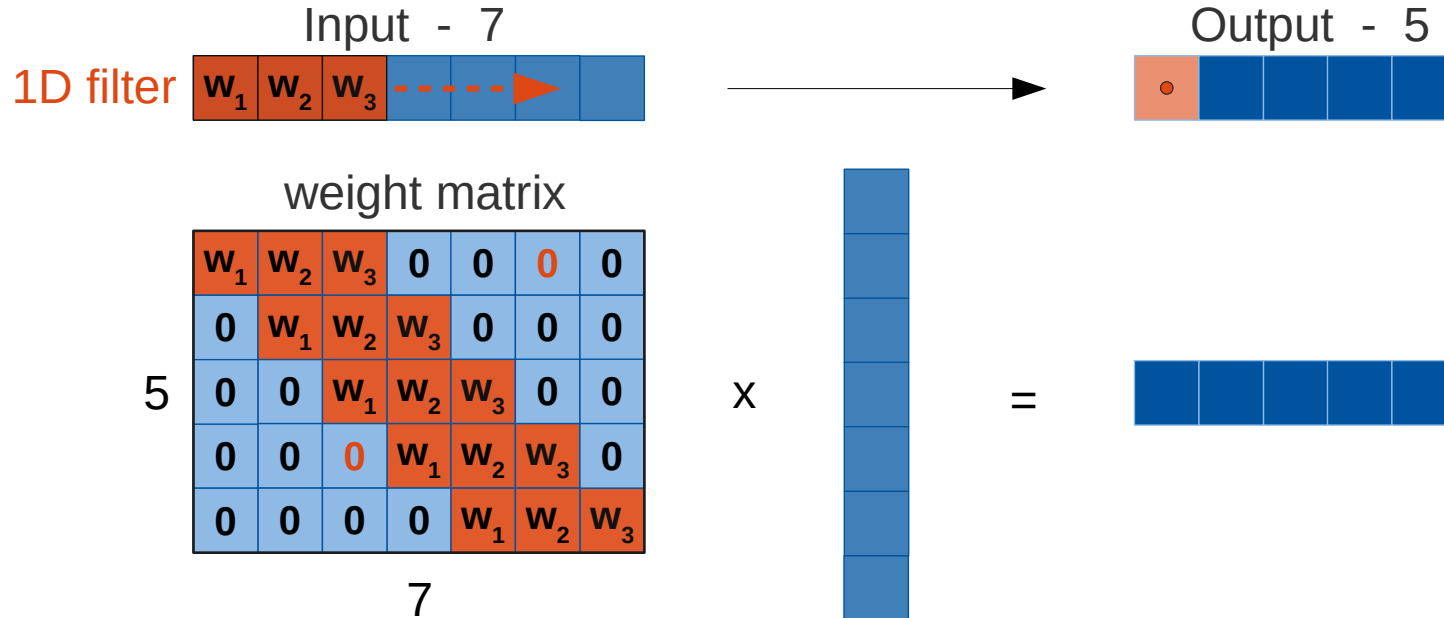
## Question: What is a good filter size?

- (a) Small (3 x 3)
- (b) Large (50 x 50)
- (c) Medium (20 x 20)



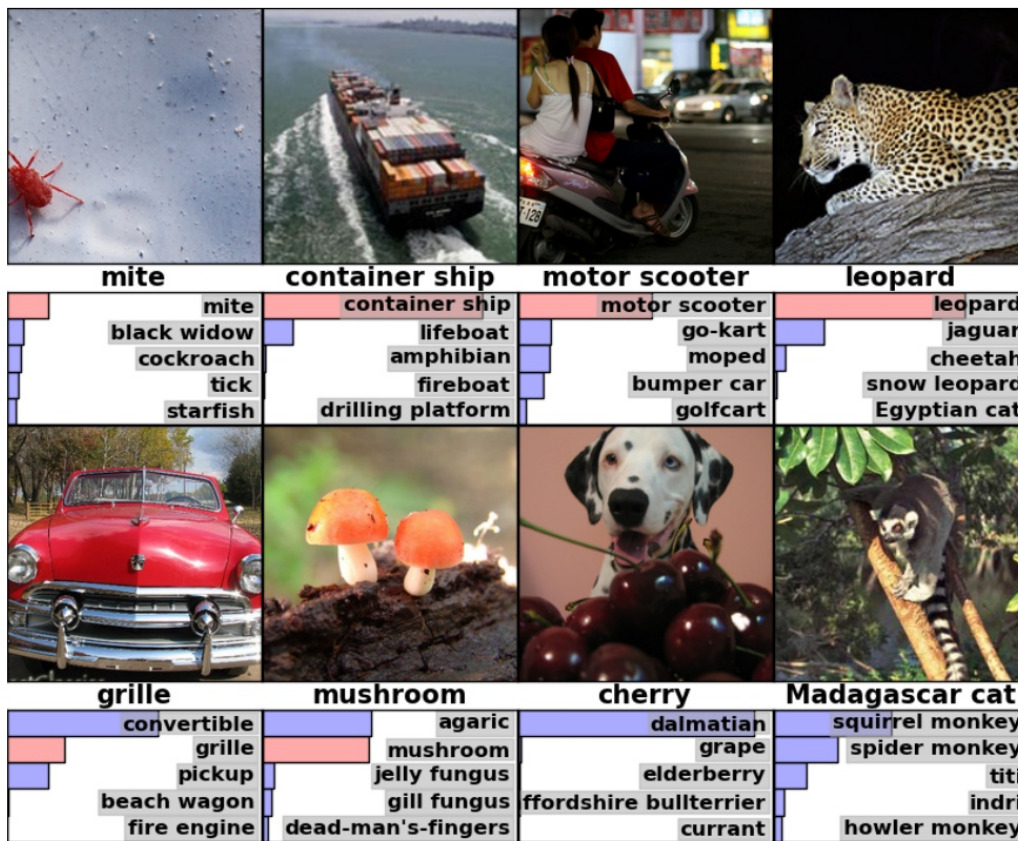
# Convolutional Operation

- Fully connected layers are special case of convolutional layers

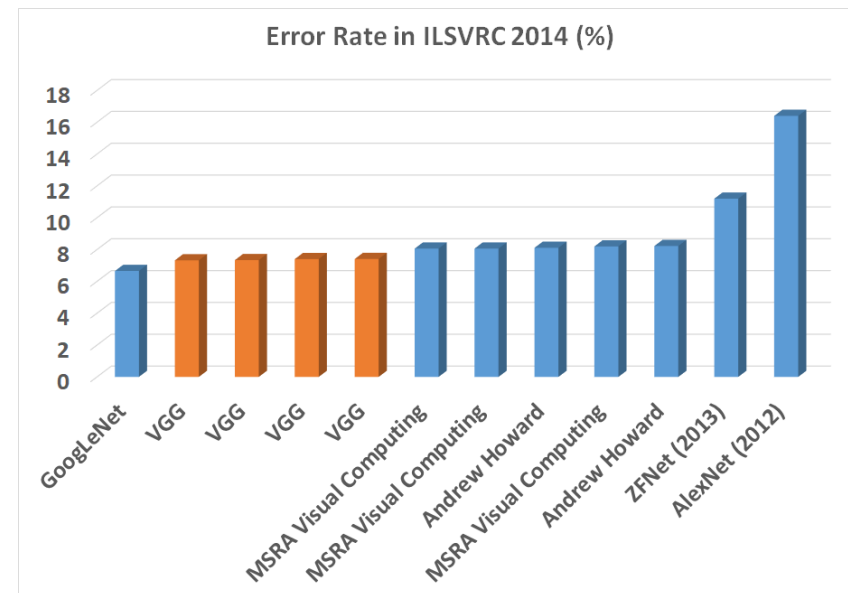


- Parameters greatly reduced due to **sparsity** and **weight sharing**
- Strong prior on **local correlation** and **translational invariance**

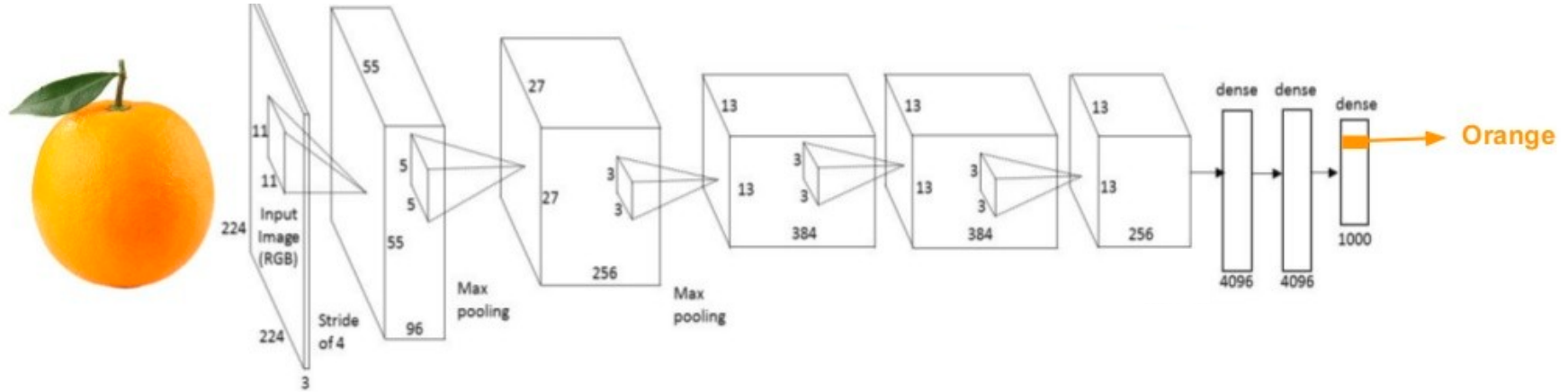
# Results on ImageNet



AlexNet (2010) - <http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>

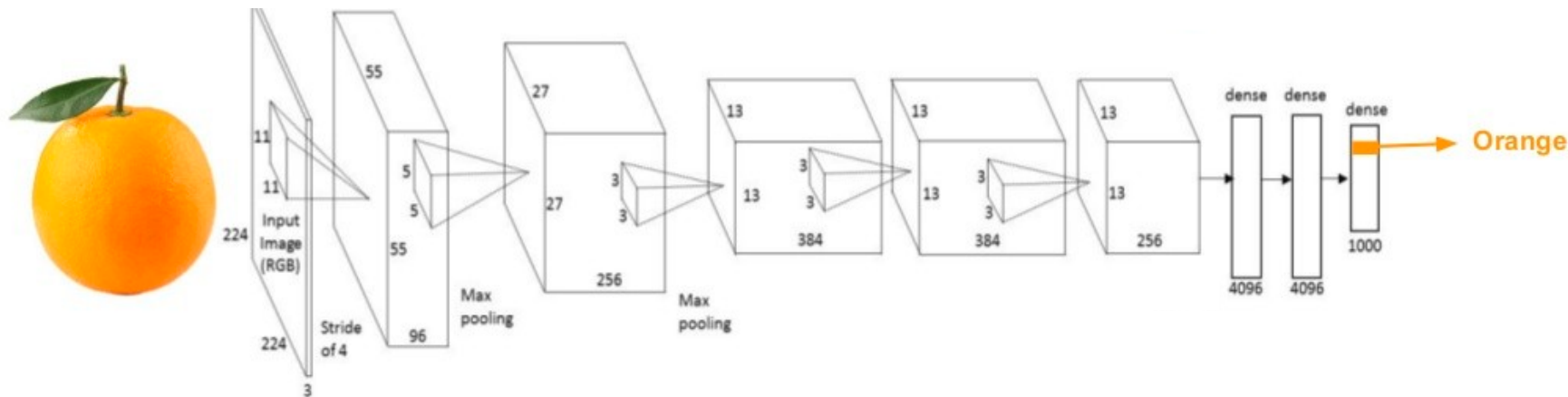


# Dropout in CNNs



## Where we have to apply dropout?

# Dropout in CNNs



## Where we have to apply dropout?

- Convolutional networks are less sensitive to over-fitting due to weight sharing
- Spatial Dropout: drop entire feature maps
- Use Dropout before MaxPooling layer
  - Make use of subdominant features

# Clarifying frequent misunderstandings



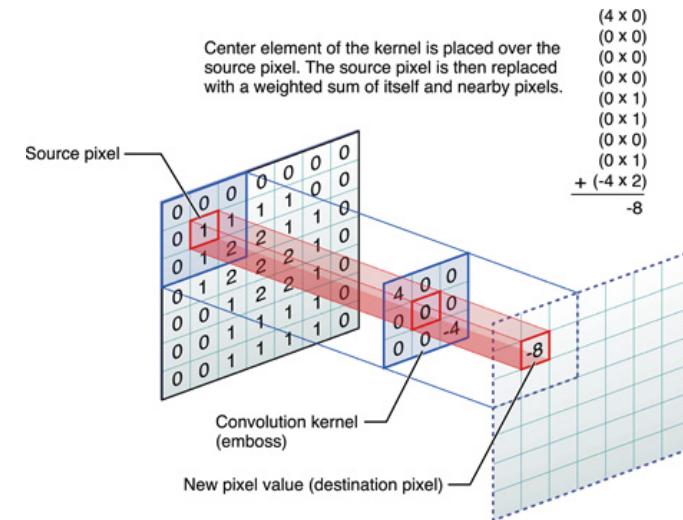
ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



- The **filters are no pre-defined** by the user → just width and depth and number
  - ♦ filters are adapted / learned by the CNN during training
- **Number of filters define number of new feature maps**
  - ♦ ten 3x3 filter applied to RGB image → 10 feature maps
- **Filter has the depth of the input image** (e.g. depth 3 for RGB images)
  - ♦ two 3x3 filter applied to RGB image → 2 feature maps, i.e. 2 channels  
→ number of adaptive parameters =  $3 \times 3 \times 3 * 2 + 2 = 56$
- **After each convolutional operation an activation is applied!** (usually)
- **CNN part is followed by a fully-connected part** (in most cases)
  - output is reshaped (flattened) to a vector → apply vanilla NN layer

# Summary

- 2D Convolution acts on 3D input (width x height x depth)
- Slide small filter over input and make linear transformation (dot product + bias)
- Hyperparameter:
  - Size of filter, typically (1 x 1), (3 x 3), (5 x 5) or (7 x 7)
  - Number of filters (feature maps)
  - **Padding** (maintain spatial extent)
  - **Striding** or **pooling** (reduce spatial extent)
- Reduction of parameters using symmetry in data:
  - Prior on **local correlations** (use small filters)
  - **Translational invariance** (weight sharing)



# References & Further Reading

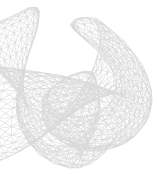


ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS

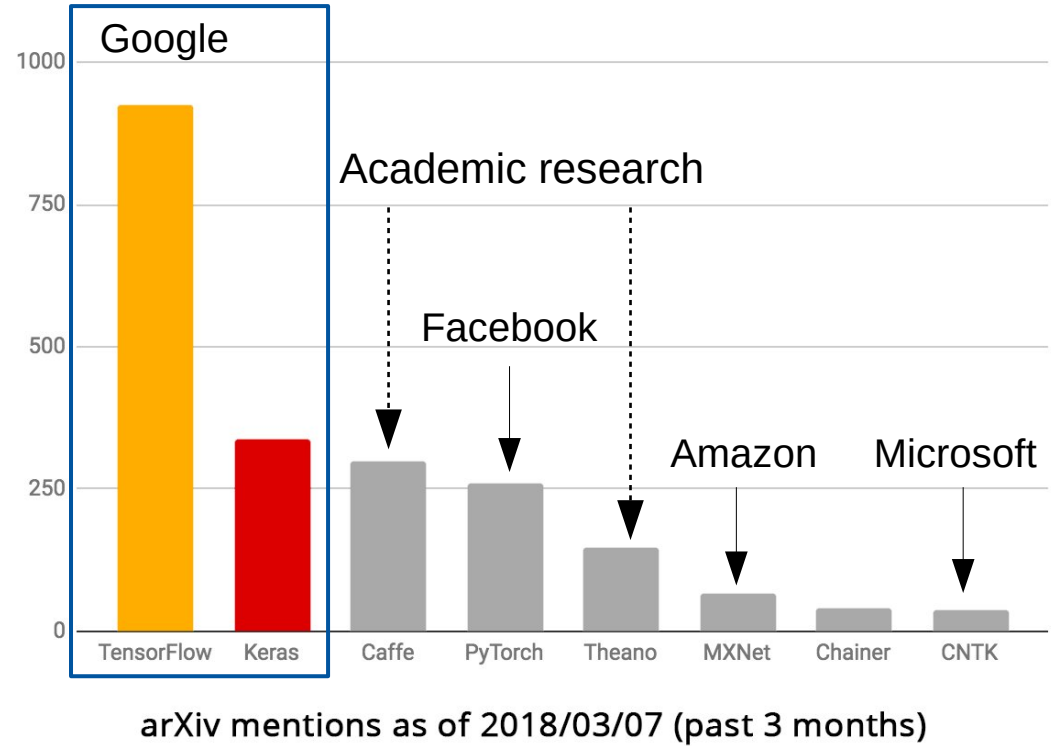


- M. Erdmann, J. Glombitza, G. Kasieczka, U. Klemradt, Deep Learning for Physics Research, World Scientific, 2021, [www.deeplearningphysics.org/](http://www.deeplearningphysics.org/)
- I. Goodfellow, Y. Bengio, A. Courville, Deep Learning, Chapter 7 / 8 / 9, MIT Press, 2016, [www.deeplearningbook.org](http://www.deeplearningbook.org)
- Xu et al. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, arXiv:1502.03044
- Y. LeCun, Y. Bengio, G. Hinton: Deep Learning, Nature 521, pages 436–444
- K. Simonyan, A. Zissermann: Very Deep Convolutional Networks for Large-Scale Image Recognition - ArXiv 1409.1556
- Toy Simulation: M. Erdmann, J. Glombitza, D. Walz, Astroparticle Physics 97, 46-53





# Practice II



# Convolutional Layers - Keras

- Same Syntax as for fully connected layers

```
layers.Convolution2D(32, kernel_size=(5, 5), padding='same', activation='relu', strides=(2, 2))
```

- layer with 32 filters, size of filter 5x5 pixels, stride of 2 in both directions, and ReLU
- Use padding='same' to keep spatial dimension (else padding='valid' )

## Pooling and transition to fully-connected networks

- Pooling layer with pooling size of 2x2 pixels and a stride of 2 in both dimensions

```
layers.MaxPooling2D((2,2), strides=(2, 2)) // layers.AveragePooling2D((2,2), strides=(2, 2))
```

- Layer flattens output to vector → allows use of Dense layers after Convolutions

```
layers.Flatten()
```

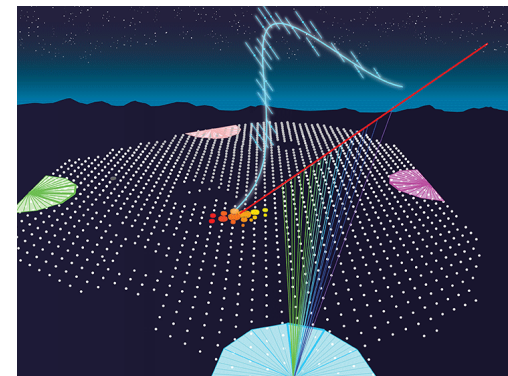
- Pooling operation on complete feature map → (remove all pixel dimensions + Flatten)

```
layers.GlobalMaxPooling2D() // layers.GlobalAveragePooling2D()
```

# Air Shower Reconstruction - CNN

Now: **OPEN** tutorial at:

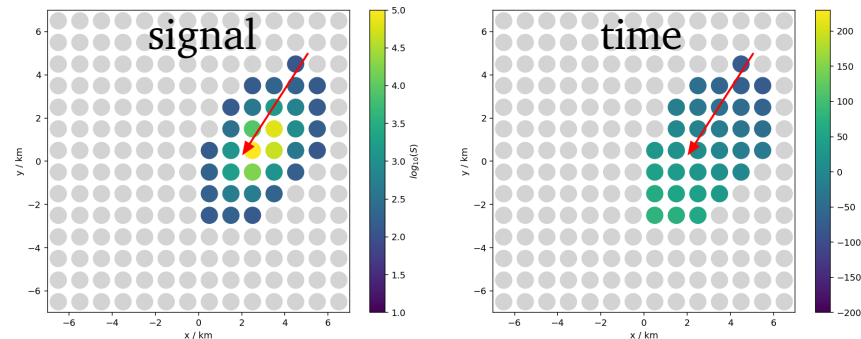
- [https://github.com/jglombitza/tutorial\\_nn\\_airshowers](https://github.com/jglombitza/tutorial_nn_airshowers)
- or click



$E = 11.4 \text{ EeV}$ ,  $\theta = 56.1^\circ$ ,  $\phi = 57.2^\circ$

## Task

- Reconstruct energy of the shower
  - ♦ Footprint is 2D image
  - ♦ **can** directly be used as input
    - input shape:  $14 \times 14 \times 2$
  - ♦ **Try to reach a resolution better than 2 EeV! (try, e.g., CNN pyramid!)**



# CIFAR 10 – Convolutional Networks



ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



Model – add:

- Conv. layers and filters
- Pooling, Dense (FC) layers
- Regularization (after Flatten)

Model – modify:

- Batch size, epochs
- Kernel size, strides
- Optimizer, learning rate

➤ **Can you achieve >75% validation accuracy?**

```
model = models.Sequential([
    layers.Convolution2D(32, kernel_size=(5, 5), padding='same',
                        activation='relu', input_shape=(32, 32, 3)),
    layers.MaxPooling2D((2,2)),
    layers.Convolution2D(64, kernel_size=(3, 3), padding='same',
                        strides=(2, 2), activation='relu'),
    layers.Flatten(),
    layers.Dropout(0.3),
    layers.Dense(10, activation='softmax')
])
```





Jonas Glombitza  
jonas.glombitza@fau.de

CTEQ Summer School 2024  
Hotel Idingshof, Bramsche

<https://github.com/DeepLearningForPhysicsResearchBook/deep-learning-physics>

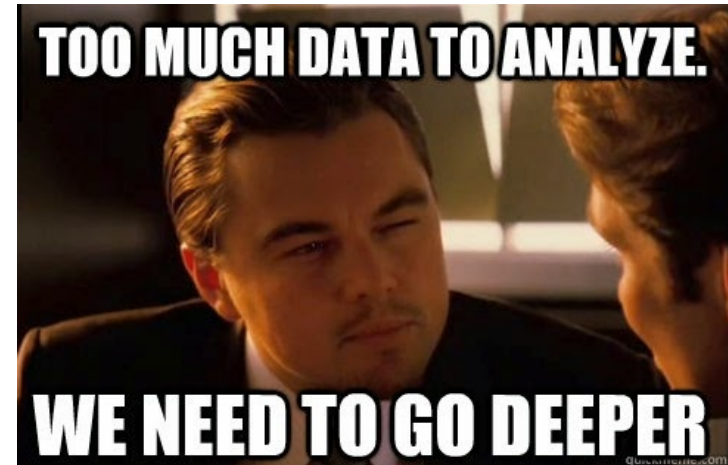


ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS

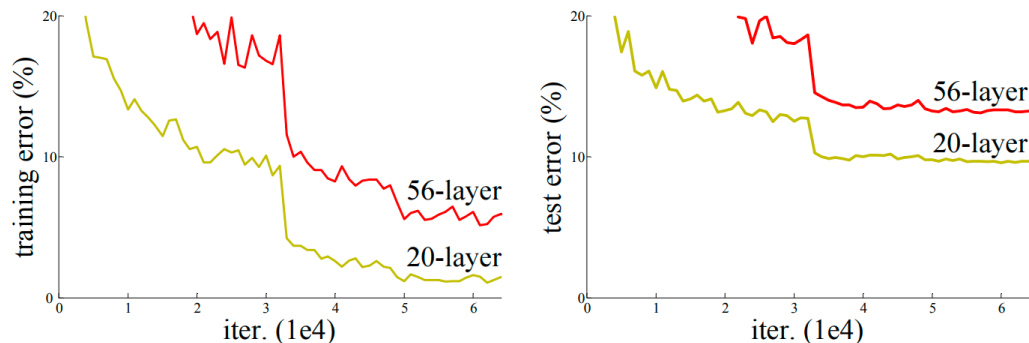


# Advanced CNN / DNN techniques

- I. Vanishing gradients
- II. Normalization layers
- III. Short cuts

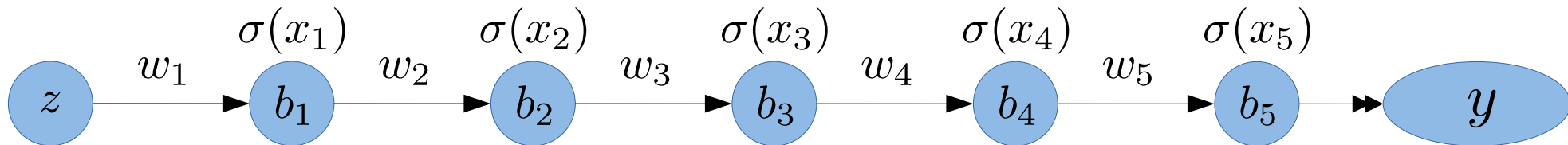


# Obstacles when Going Deeper



- Neural networks should get monotonously better when adding more layers
- **Problems**
  - ♦ Bad initialization
  - ♦ Vanishing gradients – gradients become too small
  - ♦ Shattered gradients – gradients become white noise
  - ♦ Internal covariate shift – need to constantly adapt changes in earlier layer
  - ♦ Convolutional filter show redundant behavior → advanced CNN operations

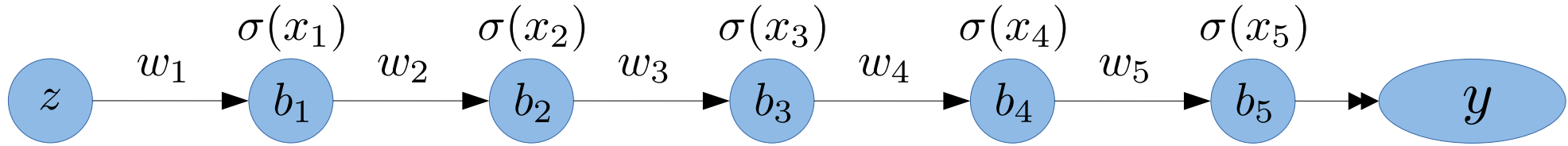
# Vanishing Gradient Problem



$$y = \sigma(x_5) = \sigma(w_5 \cdot \sigma(x_4) + b_5) = \sigma(w_5 \cdot \sigma(w_4 \cdot \sigma(x_3) + b_4) + b_5) \dots$$

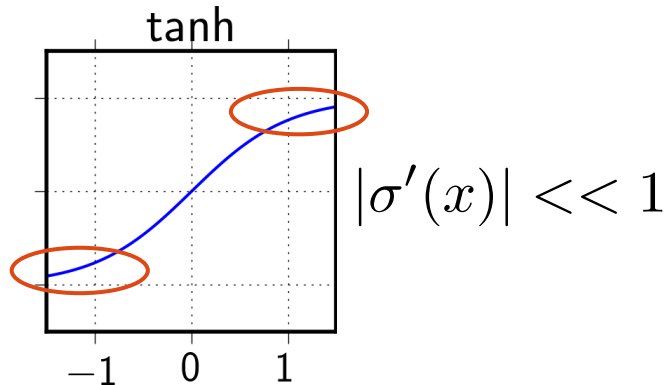
$$\frac{\partial y}{\partial w_1} = \frac{\partial \sigma(x_5)}{\partial x_5} \frac{\partial x_5}{\partial \sigma(x_4)} \frac{\partial \sigma(x_4)}{\partial x_4} \frac{\partial x_4}{\partial w_1} \dots = \underline{\sigma'(x_5)} w_5 \cdot \underline{\sigma'(x_4)} w_4 \dots \cdot \underline{\sigma'(x_1)} z$$

# Vanishing Gradient Problem



$$y = \sigma(x_5) = \sigma(w_5 \cdot \sigma(x_4) + b_5) = \sigma(w_5 \cdot \sigma(w_4 \cdot \sigma(x_3) + b_4) + b_5) \dots$$

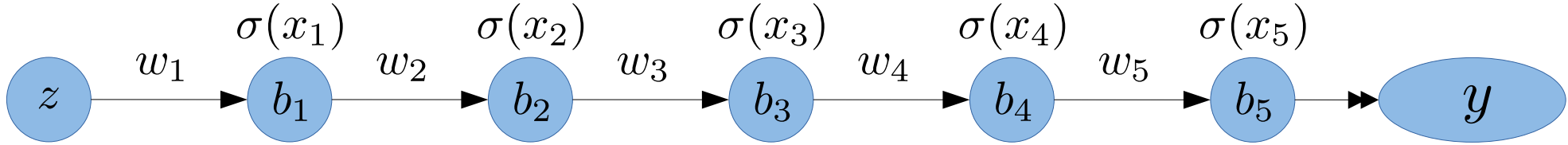
$$\frac{\partial y}{\partial w_1} = \frac{\partial \sigma(x_5)}{\partial x_5} \frac{\partial x_5}{\partial \sigma(x_4)} \frac{\partial \sigma(x_4)}{\partial x_4} \frac{\partial x_4}{\partial w_1} \dots = \underline{\sigma'(x_5)} w_5 \cdot \underline{\sigma'(x_4)} w_4 \dots \cdot \underline{\sigma'(x_1)} z$$



- Stacking many layers can lead to vanishing gradients
- Activation saturates
  - Updates in early layers become very tiny
  - **No learning**
  - Don't use sigmoids / tanh only rarely → use shortcuts
  - still useful as last layer if data ranges [0;1] or [1;-1]



# Vanishing Gradient Problem



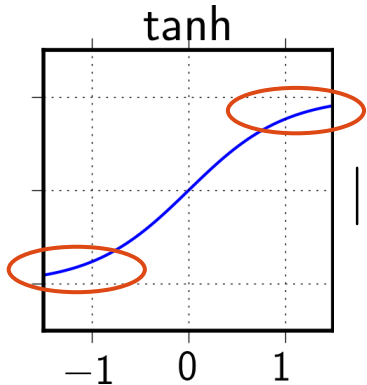
$$y = \sigma(x_5) = \sigma(w_5 \cdot \sigma(x_4) + b_5) = \sigma(w_5 \cdot \sigma(w_4 \cdot \sigma(x_3) + b_4) + b_5) \dots$$

$$\frac{\partial y}{\partial w_1} = \frac{\partial \sigma(x_5)}{\partial x_5} \frac{\partial x_5}{\partial \sigma(x_4)} \frac{\partial \sigma(x_4)}{\partial x_4} \frac{\partial x_4}{\partial w_1} \dots = \sigma'(x_5) w_5 \cdot \sigma'(x_4) w_4 \dots \cdot \sigma'(x_1) z$$

**VISUALIZATION**

[LINK](#)

- Stacking many layers can lead to vanishing gradients
- Activation saturates



$$|\sigma'(x)| \ll 1$$

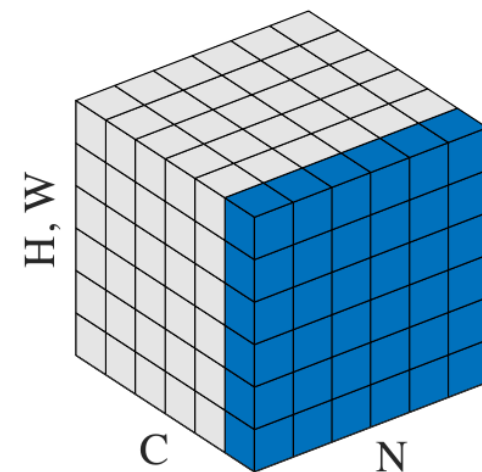
- Updates in early layers become very tiny
- **No learning**
- Don't use sigmoids / tanh only rarely → use shortcuts
- still useful as last layer if data ranges [0;1] or [1;-1]

# Batch Normalization

- Calculate batch-wise for each channel:
  - Mean:  $\mu_B$  and Variance:  $\sigma_B^2$
  - Add free parameters  $\gamma, \beta$ 
    - Easily control first moments of distribution

$$\triangleright y = \frac{x - \mu_B}{\sigma_B} \gamma + \beta$$

- Makes DNN robust against poor initializations
- Helps with vanishing gradient / less sensitive to high learning rates
- Has regularizing effect (no large weights, noise because of batch dependency)
- Reduce internal covariate shift
- **Very successful for convolutional architectures**



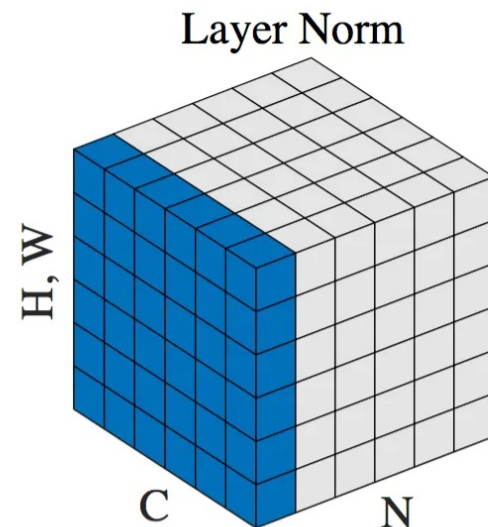
# Layer Normalization

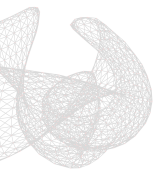
- Calculate sample-wise (across channels)
  - ♦ Mean:  $\mu_l$  and Variance:  $\sigma_l^2$
  - ♦ Add free parameters  $g$ ,  $b \rightarrow$  control first moments

$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l, \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$\triangleright y = \frac{x - \mu_B}{\sigma_B} g + b$$

- Makes DNN robust against poor initializations
- Improves training behaviour  $\rightarrow$  normalized gradients
- $\triangleright$  Very successful for recurrent / transformer architectures (small batches)**





## Does it Make sense to use a $(1 \times 1)^*$ filter

- (a) No! Number of learnable parameters just too small
- (b) Yes! It can help to reduce number of parameters!
- (c) Usually not, but sometimes in an engineering task ...

\*you can think of a  $1 \times 1$  filter as a fully-connected (linear projection) in the feature space  
[Think of one feature map as a node]

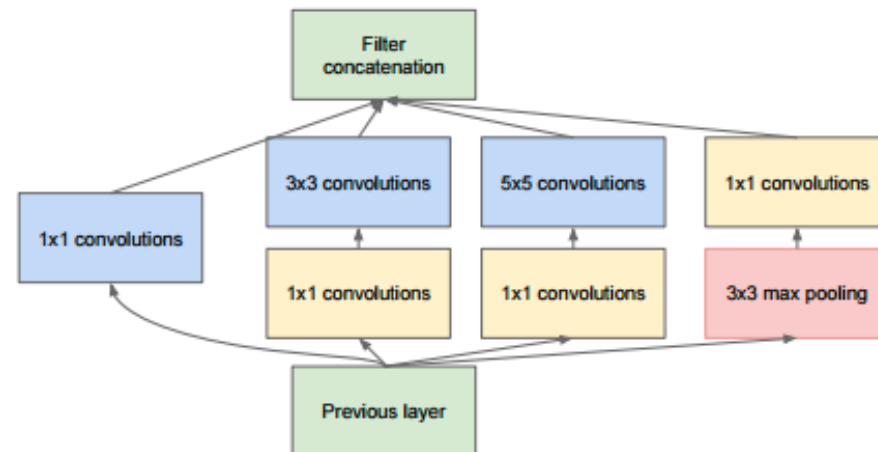
# Inception Module

**Key observation:** Convolutional filters show redundant behavior



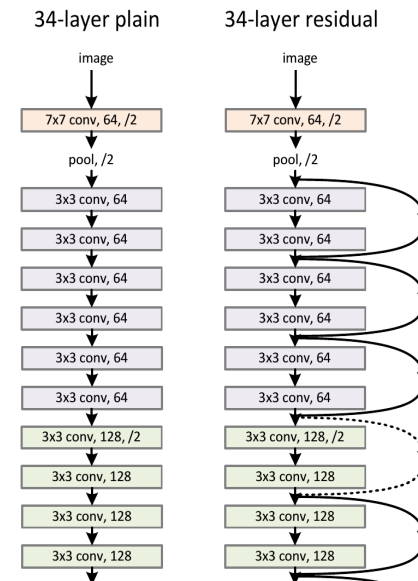
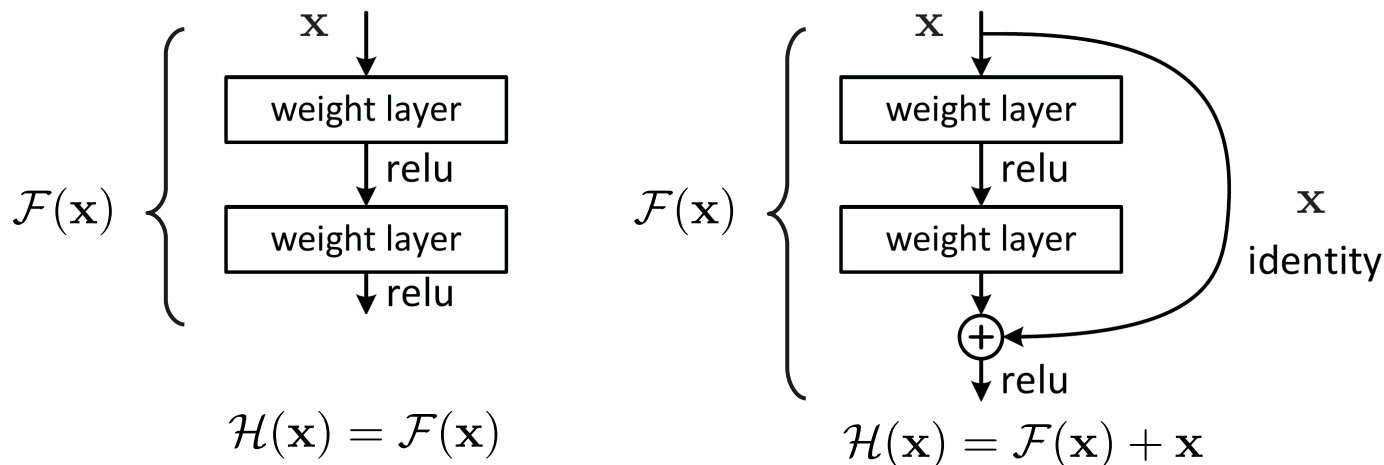
**Idea:** Factorize convolution operation

- Use different small convolutions in parallel and concatenate outputs
- Massive use of (1 x 1) convolutions
- Increase model complexity
- Make model sensitive to different scales



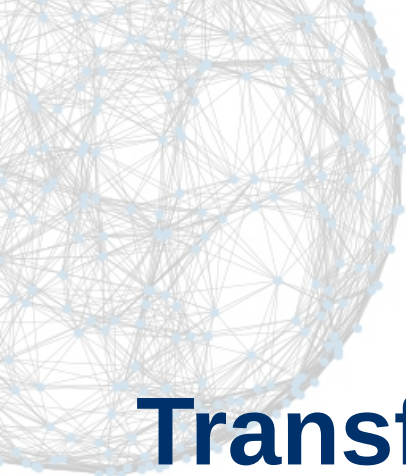
# Residual Unit

**Idea:** Residual unit consisting of small network and a shortcut (identity mapping)



*Up to several of  
hundreds layer deep!*

- Weight block learns small residual  $\mathcal{F}(\mathbf{x})$  on top of input  $\mathbf{x}$ 
  - Output of residual unit  $\mathcal{H}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x}$
- Shortcut let gradient propagate easily to earlier layers
- Later layers can easily turn weights to zero by  $\mathcal{F}(\mathbf{x}) \rightarrow 0$



Jonas Glombitza  
jonas.glombitza@fau.de

CTEQ Summer School 2024  
Hotel Idingshof, Bramsche

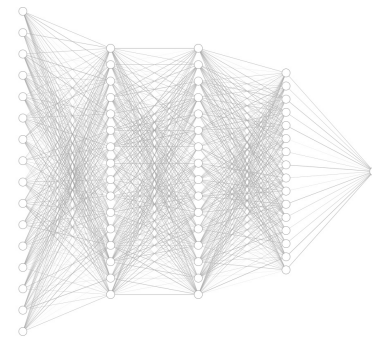
# Transformers

- I. Processing image-like data
- II. Point cloud transformers

<https://github.com/DeepLearningForPhysicsResearchBook/deep-learning-physics>



ERLANGEN CENTRE  
FOR ASTROPARTICLE  
PHYSICS



# Transformers

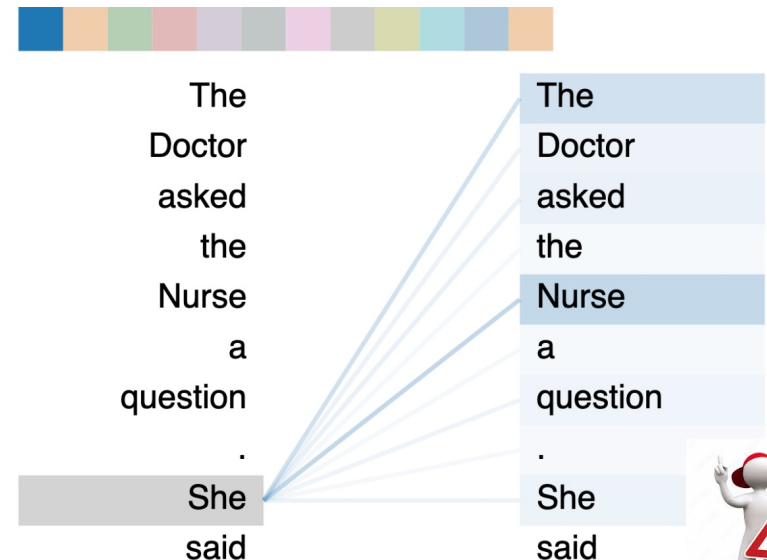
- Transformers are **backbone of latest breakthroughs**: LLMs / Stable Diffusion
- Building blocks: DNNs with attention mechanism
  - ◊ Which parts of sequence semantically correlated → analyze together

In a nutshell: extension of fully-connected DNNs

- listen to all inputs
- *but* focus on most important inputs
- increase noise robustness

**Analyze sequences (different lengths):**

- $(X_1, X_2, X_3, X_4, X_5, \dots, X_n)$ 
  - ◊ single element called *token* (e.g, word)





# Simple case: Attention of 2 vectors

$$X = \begin{pmatrix} \vec{x}_1^T \\ \vec{x}_2^T \end{pmatrix} = \begin{pmatrix} x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & x_2^{(2)} \end{pmatrix}$$

- Input: 2 vectors (arbitrary length), could, e.g., represent words / measurements

Dot product measures similarity: parallel = 1, orthogonal = 0

$$A = \begin{pmatrix} \vec{x}_1 \cdot \vec{x}_1 & \vec{x}_1 \cdot \vec{x}_2 \\ \vec{x}_2 \cdot \vec{x}_1 & \vec{x}_2 \cdot \vec{x}_2 \end{pmatrix} = X X^T$$

**Attention matrix:**  $A$

- Estimate attention (similarity of vectors)
- Calculate dot product
- Add learnable parameters to compare similarities in feature space (after filtering), use  $W^Q$
- ✗ Attention matrix symmetric (not useful)

$$A = \begin{pmatrix} W^Q \vec{x}_1 \cdot W^Q \vec{x}_1 & W^Q \vec{x}_1 \cdot W^Q \vec{x}_2 \\ W^Q \vec{x}_2 \cdot W^Q \vec{x}_1 & W^Q \vec{x}_2 \cdot W^Q \vec{x}_2 \end{pmatrix} = X W^Q (X W^Q)^T$$

=1 → boring  
← same! →  
=1 → boring

# Simple case: Attention of 2 vectors

$$X = \begin{pmatrix} \vec{x}_1^T \\ \vec{x}_2^T \end{pmatrix} = \begin{pmatrix} x_1^{(1)} & x_1^{(2)} \\ x_2^{(1)} & x_2^{(2)} \end{pmatrix}$$

- Input: 2 vectors (arbitrary length), could, e.g., represent words / measurements

$$A = \begin{pmatrix} \vec{x}_1 \cdot \vec{x}_1 & \vec{x}_1 \cdot \vec{x}_2 \\ \vec{x}_2 \cdot \vec{x}_1 & \vec{x}_2 \cdot \vec{x}_2 \end{pmatrix} = X X^T$$

**Attention matrix:**  $A$

- Estimate attention (similarity of vectors)
- Calculate dot product
- Add 2 learnable sets of parameters to compare in feature space, project using  $W^Q$ ,  $W^K$ 
  - ✓ Attention matrix not symmetric & diagonal entries not unity

$$A = \begin{pmatrix} W^Q \vec{x}_1 \cdot W^K \vec{x}_1 & W^Q \vec{x}_1 \cdot W^K \vec{x}_2 \\ W^Q \vec{x}_2 \cdot W^K \vec{x}_1 & W^Q \vec{x}_2 \cdot W^K \vec{x}_2 \end{pmatrix} = X W^Q (X W^K)^T$$

# Scaled dot product attention

- Create three projections of input,  $W^i \in \mathbb{R}^{f \times d}$ 
  - perform scaled dot product attention

$$X = \begin{pmatrix} \vec{x}_i^T \\ \vdots \end{pmatrix}$$

Input: arbitrary number of vectors, of dimension  $f$

queries

$$Q = X \cdot W^Q \rightarrow \text{projection 1 for attention matrix}$$

keys

$$K = X \cdot W^K \rightarrow \text{projection 2 for attention matrix}$$

values

$$V = X \cdot W^V \rightarrow \text{project input to multiply with attention matrix}$$

Output shape:  $n_{\text{inputs}} \times d$

$$\text{Attention}(Q, K, V) \hat{=} A V = \text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V = Z$$

$$\frac{1}{\sqrt{d}} \rightarrow \text{ensure } \text{Var}(QK^T) = 1 \text{ after initialization}$$

$$\text{Attention matrix: } A_{ij} = \frac{e^{q_i \cdot k_j}}{\sum_j e^{q_i \cdot k_j}}$$

# Scaled dot product attention

- Create three projection of inputs  $W^i \in \mathbb{R}^{4 \times 3}$ 
  - perform scaled dot product attention

$$X = \begin{pmatrix} \vec{x}_1^T \\ \vec{x}_2^T \end{pmatrix} = \begin{array}{|c|} \hline \vec{x}_1 \\ \hline \vec{x}_2 \\ \hline \end{array}$$

Input: 2 vectors, of dimension  $f = 4$

queries  $Q = X \cdot W^Q$

keys  $K = X \cdot W^K$

values  $V = X \cdot W^V$

Output shape:  $n_{\text{inputs}} \times d = 2 \times 3$

$$\text{Attention}(Q, K, V) \hat{=} A V = \underbrace{\text{softmax} \left( \frac{Q \times K^T}{\sqrt{3}} \right)}_{\text{Attention matrix}} \cdot V = Z$$

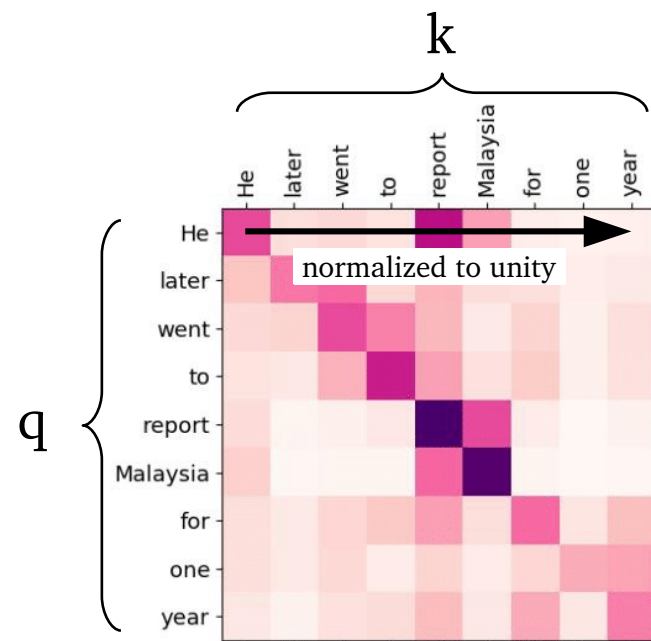
Attention matrix:  $A_{ij} = \frac{e^{q_i \cdot k_j}}{\sum_j e^{q_i \cdot k_j}}$

# Attention

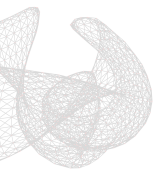
- Enables to learn long range correlations
  - ♦ CNNs: prior on local correlations (small filters)
  - ♦ Correlates all inputs with each other (like in fully-connected networks)
  - ♦ Attention: clever way to tell DNN where to focus on
- Very memory expensive
  - ♦ Matrix scales with sequence length squared:  $n_{inp}^2$

Attention mechanism very flexible

- Due to dot product: independent of sequence length



Visualized attention matrix



## Why don't we multiply with the inputs but values (V)?

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V$$

- (a) In transformers, it doesn't matter to what we apply the attention
- (b) We are still multiplying with input (but with increased capacity)
- (c) Projecting stuff, always helps ...



# Values: Attention is applied to input

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right) V, \quad V \hat{=} X \cdot W^V$$

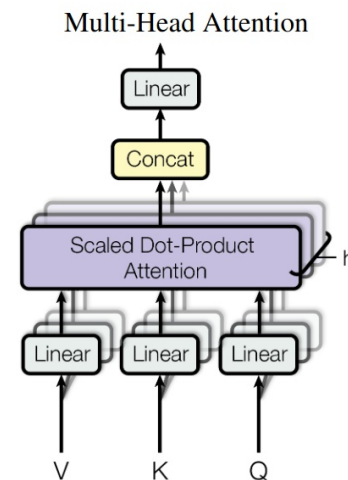
- Attention matrix multiplied by values (projected input)
- Value projection is corresponding linear transformation as in fully-connected layer
- Matrix multiplication is associative

$$\underbrace{\text{softmax} \left( \frac{QK^T}{\sqrt{d}} \right)}_{\text{Attention matrix}} V = A \underset{\text{I.}}{(X W^V)} = \underset{\text{II.}}{(A X)} W^V = Z$$

- Two interpretations possible:
  - I. Multiply projected input by attention matrix
  - II. “input multiplied by attention” is projected

# Multi-head attention

- Scaled dot-product attention performed in parallel  $\rightarrow$  called *heads*
  - Analyze  $h$  subspaces with different representation
- Keep complexity
  - dimensionality per head reduced:  $d' = d/h$   
(Note:  $d$  needs to be dividable by  $h$ )
- Add final projection to combined heads
  - $V = X \cdot W^O$



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) \cdot W^O$$

$$\text{head} = \text{Attention} \left( QW_i^Q, KW_i^K, VW_i^V \right)$$

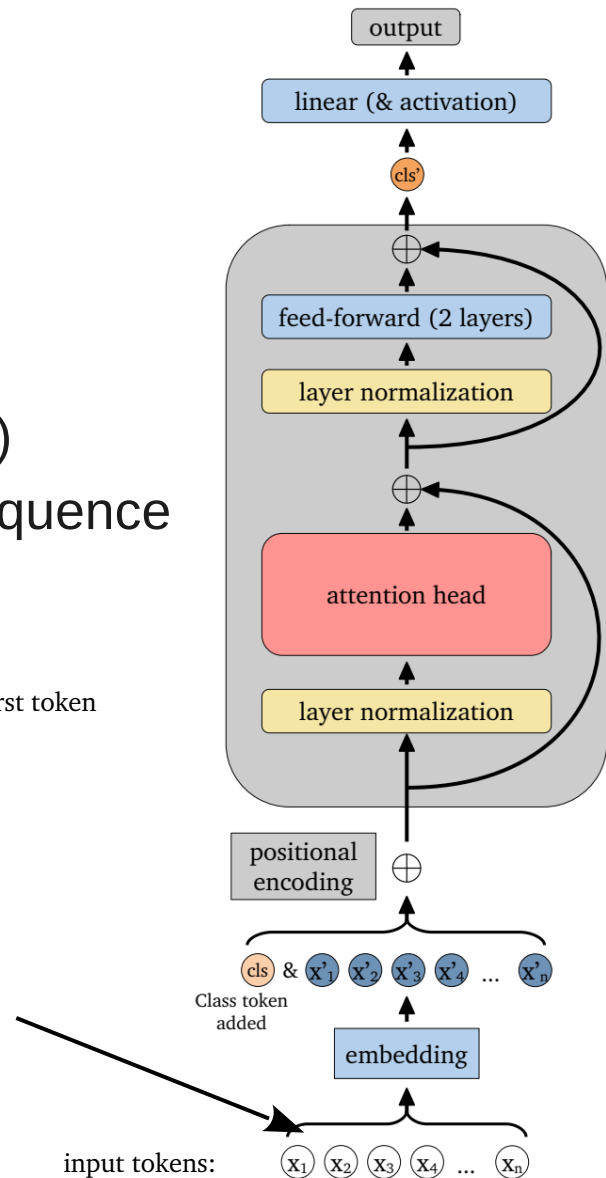
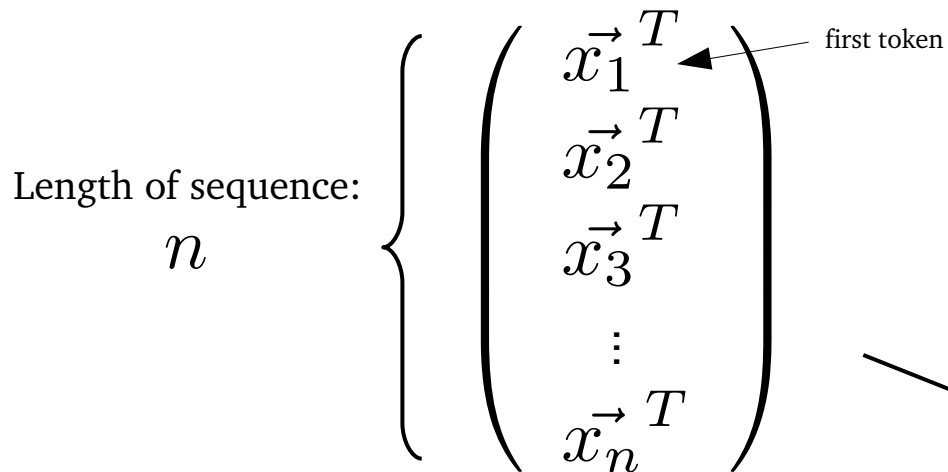


# Inputs to the transformer

\*For simplicity we consider batch size of 1  
usually training is performed in batches

## Inputs\*

- Input is a **single** sequence (one sentence)
- Arbitrary number of tokens (words) per sequence
  - Each word / vector has dimension  $f$



# Embedding

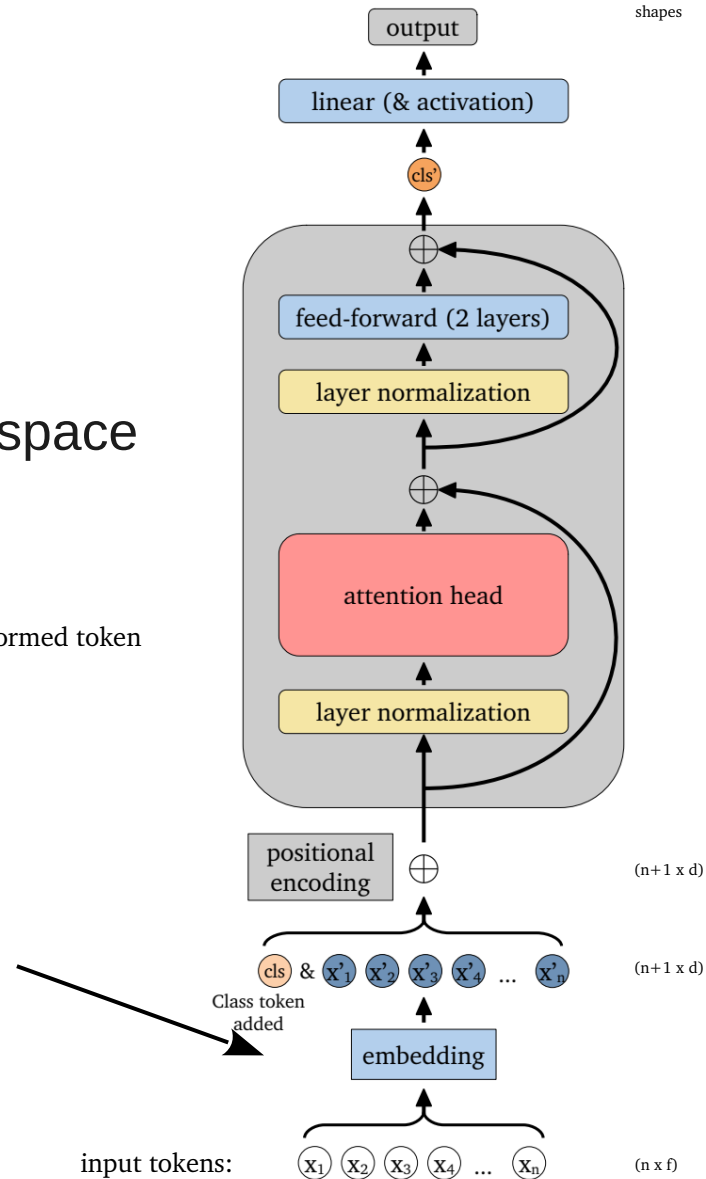
## Embedding layer

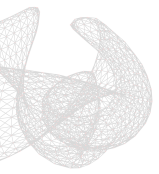
- Project inputs vectors into high-dimensional space
- Find suitable representation

Embedding matrix

$$X' = W_{\text{embed}} X = \begin{pmatrix} \vec{x}'_1 \\ \vec{x}'_2 \\ \vec{x}'_3 \\ \vdots \\ \vec{x}'_n \end{pmatrix} \begin{matrix} T \\ T \\ T \\ T \\ T \end{matrix}$$

transformed token





## Does the attention matrix considers token ordering?

- (a) Yes! How should large language models (LLMs) work if not?!
- (b) No, the decomposition is permutational invariant

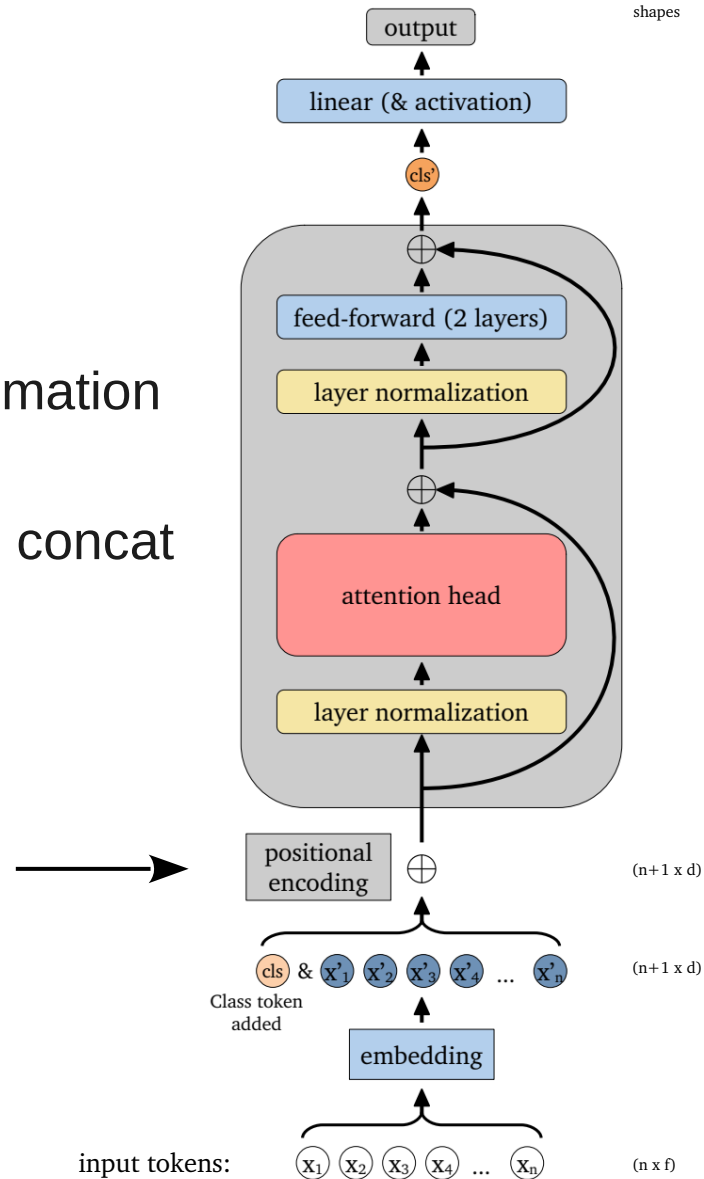


# Positional encoding

## Positional encoding

- Transformers are permutational invariant
- If positions / ordering contain important information
  - ◊ add positions as additional information
- Many options possible: encode & add / pure concat

A permutational invariant transformer would confuse  
 “I like physics more than ice cream”  
 with  
 “I like ice cream more than physics”



# Positional encoding

## Positional encoding

- Transformers are permutational invariant
- If positions / ordering contain important information
  - add positions as additional information
- Many options possible: encode & add / pure concat

### Example positions

Coordinate of point:

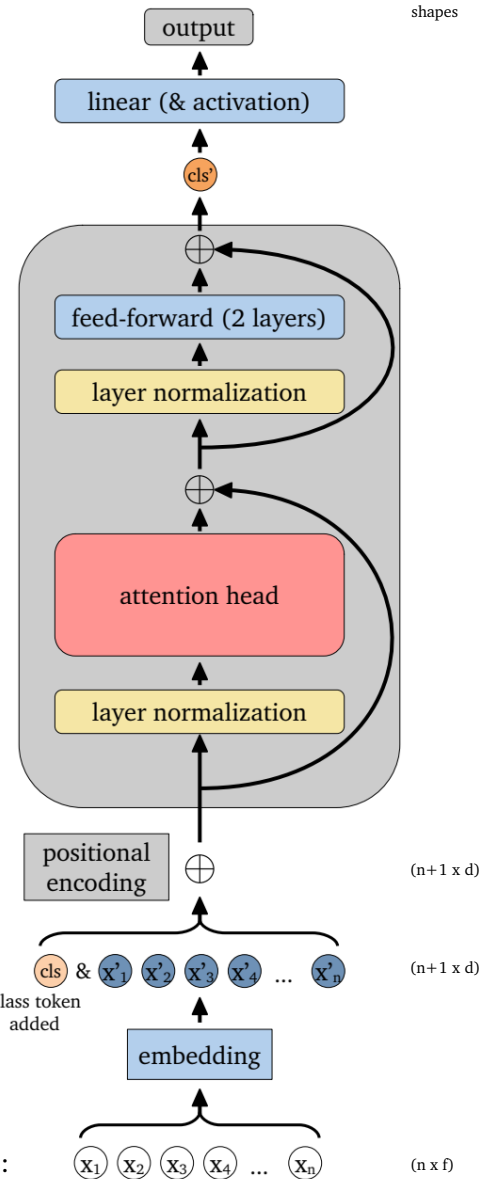
$$\vec{x}_{pos} = \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Index of sequence:

$$\vec{x}_{pos} = ( i )$$

Encoding matrix

$$X' = P_{enc} X_{pos} = \begin{pmatrix} \vec{x}'_{1,pos} \\ \vec{x}'_{2,pos} \\ \vec{x}'_{3,pos} \\ \vdots \\ \vec{x}'_{n,pos} \end{pmatrix}$$



# Positional encoding

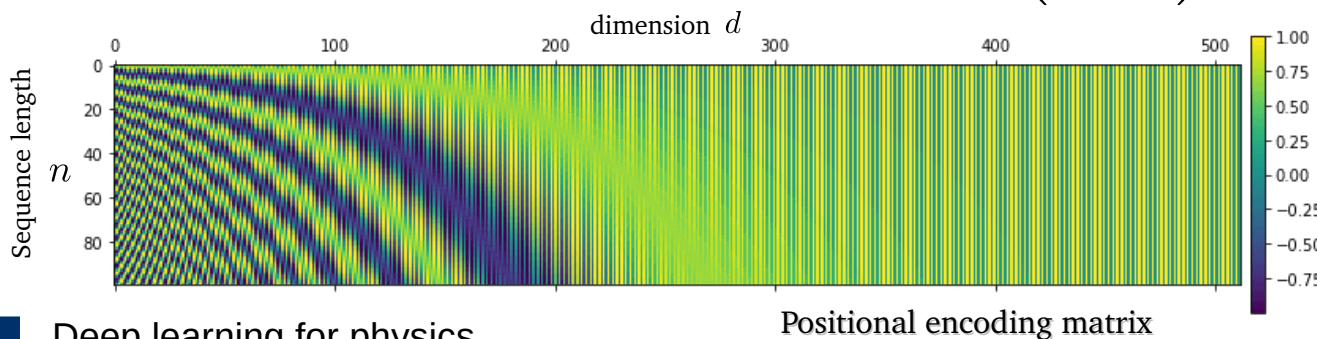
## Positional encoding

- Transformers are permutational invariant
- If positions / ordering contain important information
  - add positions as additional information

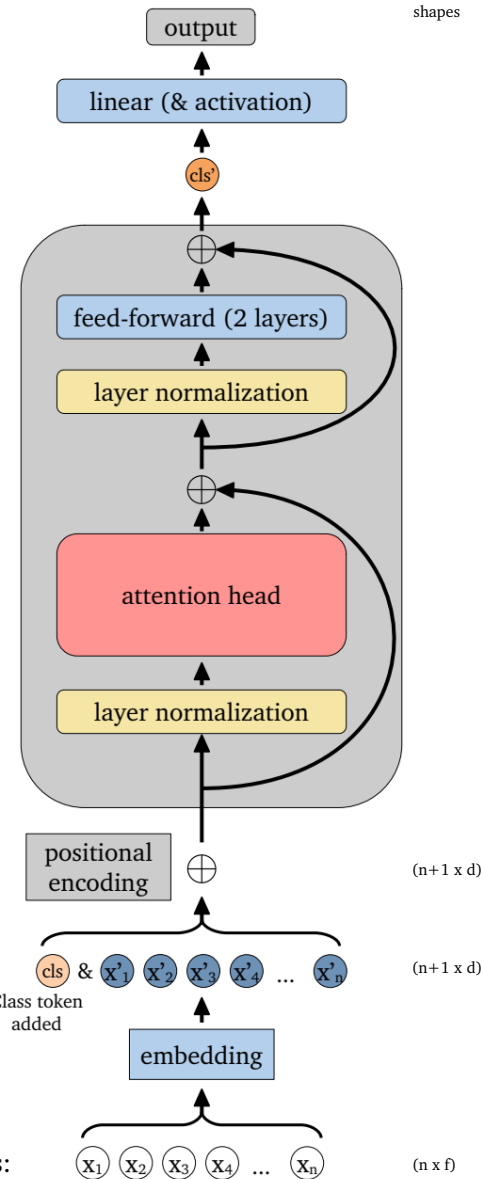
Positional encoding  
for vision transformer

$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P(k, 2i + 1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$



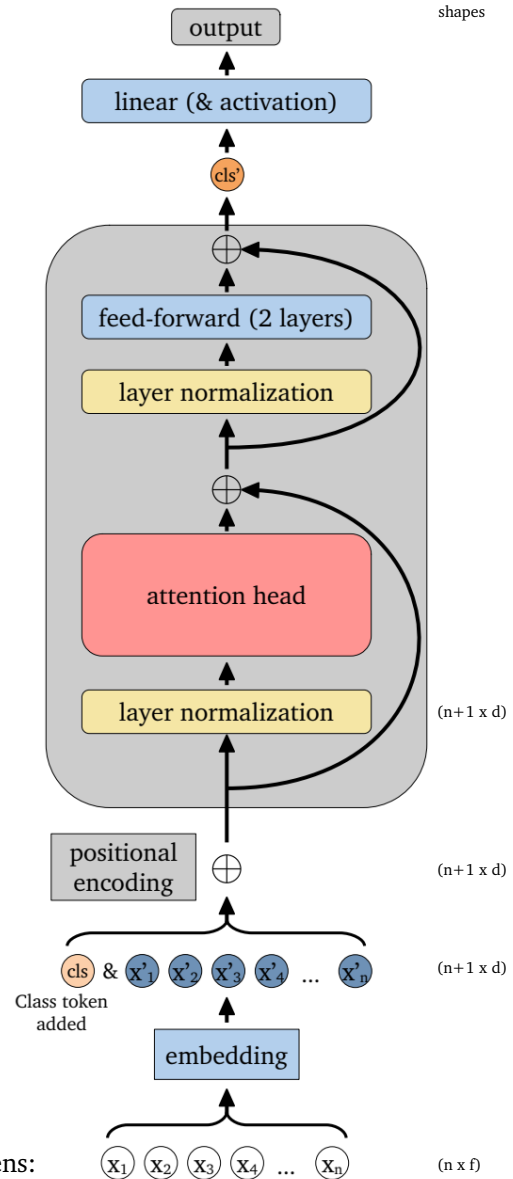
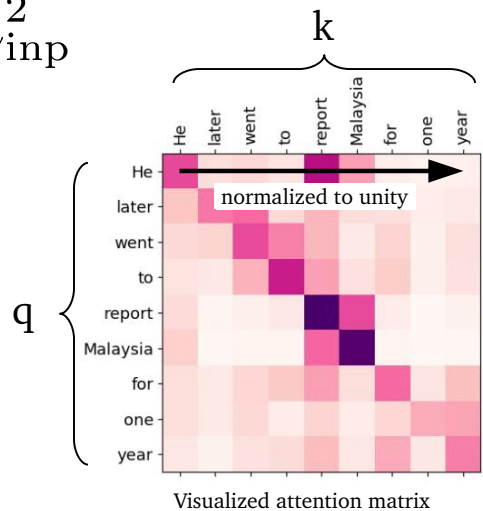
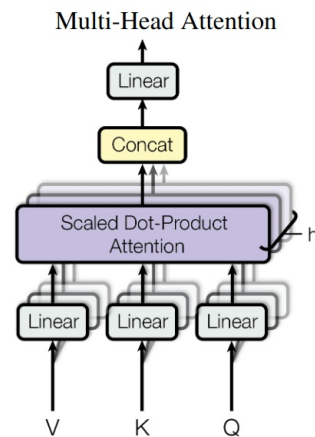
Positional encoding matrix



# Attention block

## Attention block

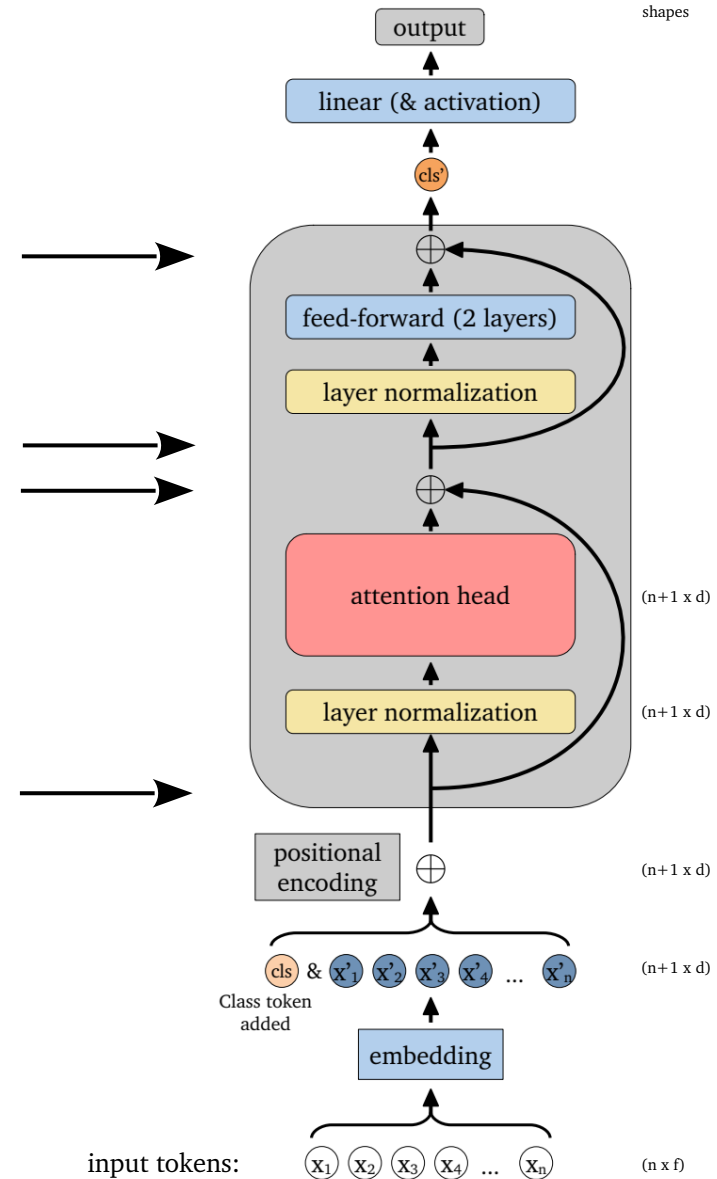
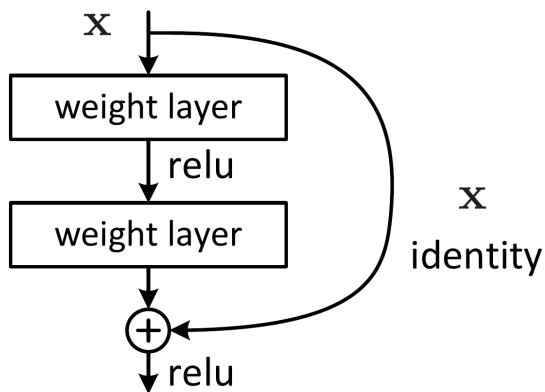
- Multi-head attention performed
  - $h$  heads run in parallel  $d' = d/h$
- Learn long-range relation
- Expensive, scales  $n_{inp}^2$



# Residual connection

## Residual connections

- Learn residual:  $\mathcal{H}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{x}$
- Improved gradient propagation
- Previous information easily “kept”

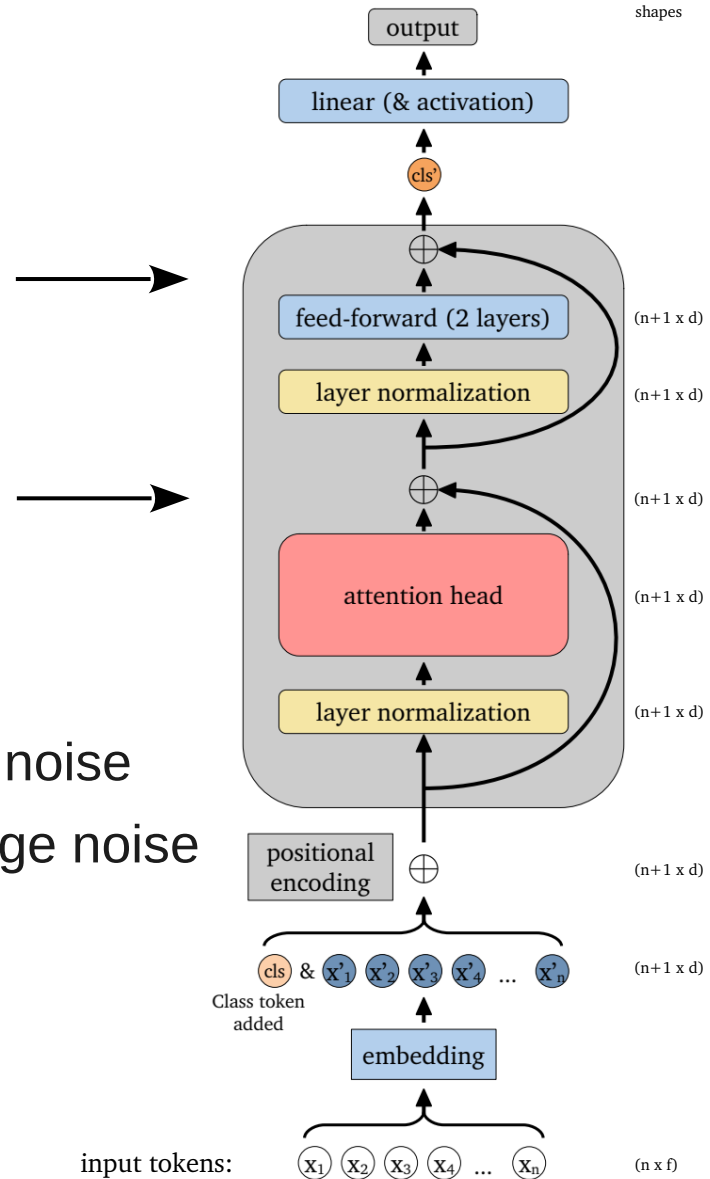
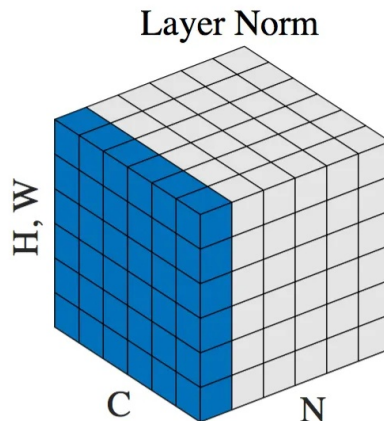


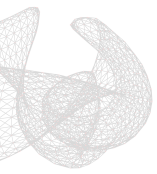


# Layer Normalization

## Layer normalization

- Stabilize training
- Normalize across features
- Useful for transformers
  - ◆ Usually small batches → batch norm: large noise
  - ◆ Varying sequence length → batch norm: large noise





## Why don't we use batch normalization?

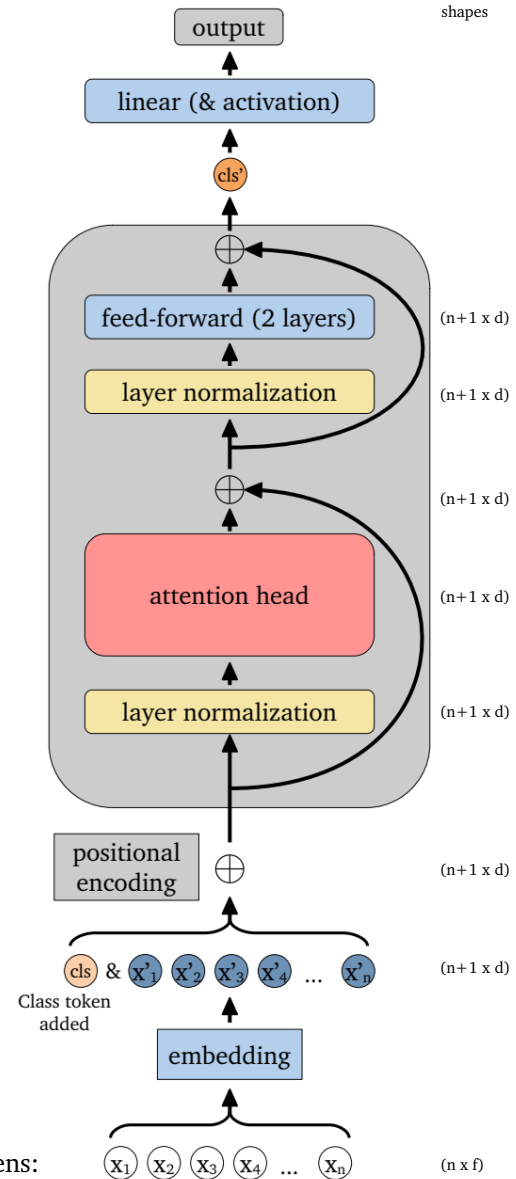
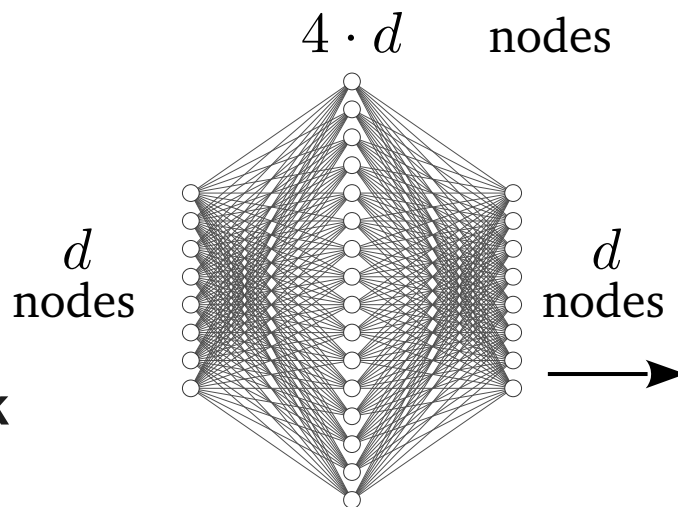
- (a) Layer normalization causes less noise for varying sequence lengths
- (b) We can also use batch normalization
- (c) In transformer training batches are usually too small



# Feed-forward part

## Fully-connected network

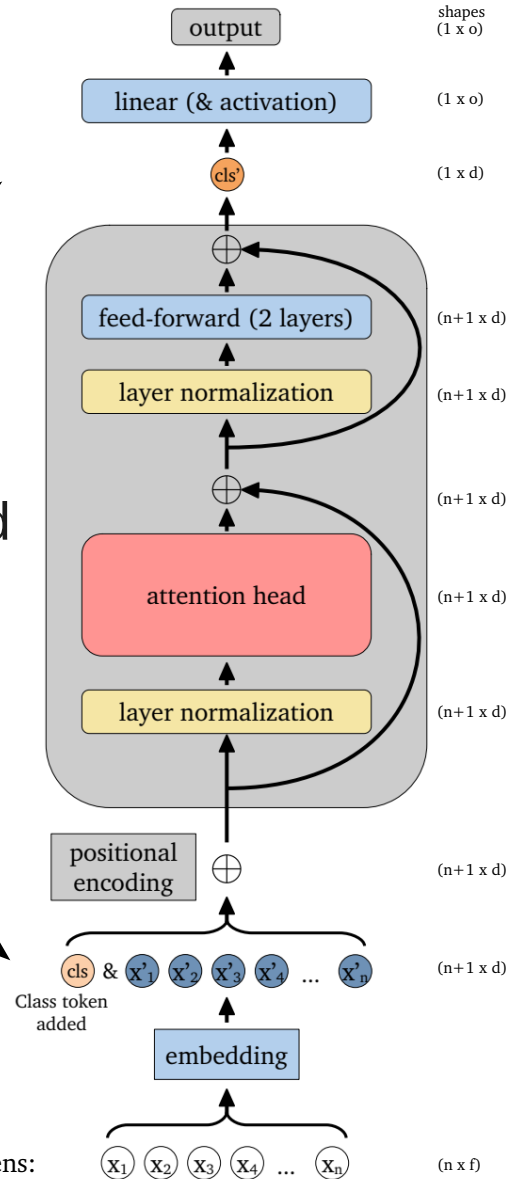
- Affine transformation
  - ◆ Includes bias (not only projections)
- Operates on each token in parallel (**weight sharing**)
- Often 2 layers (hidden layer holds larger capacity)



# CLS (class) token

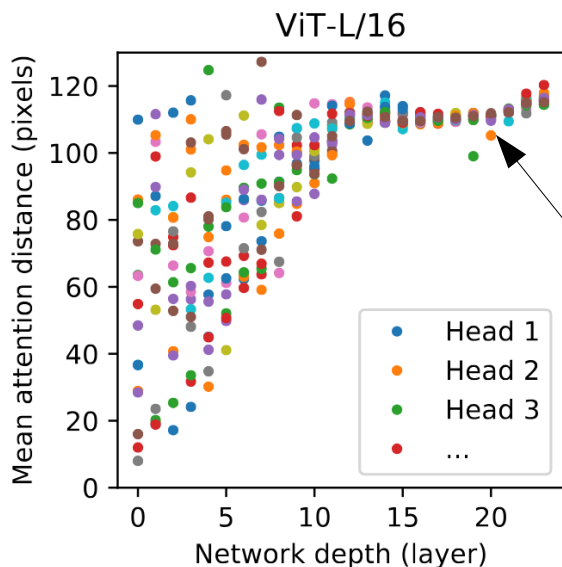
## Class token

- Used for the final prediction
  - ◊ Different to RNNs where usually last time step used
- CLS token can be interpreted as memory register
  - ◊ Store and read from that token

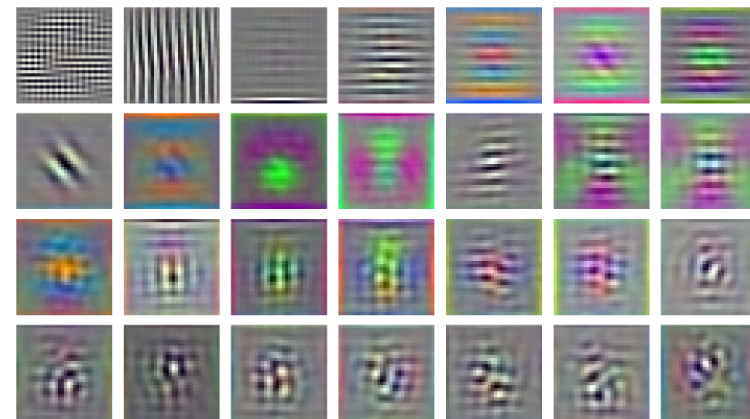


# Transformer Networks

- Interestingly transformers show similar behavior as CNNs
  - Visualized filters in first layers very similar
- Deeper layers focus on long-range correlations



RGB embedding filters  
(first 28 principal components)

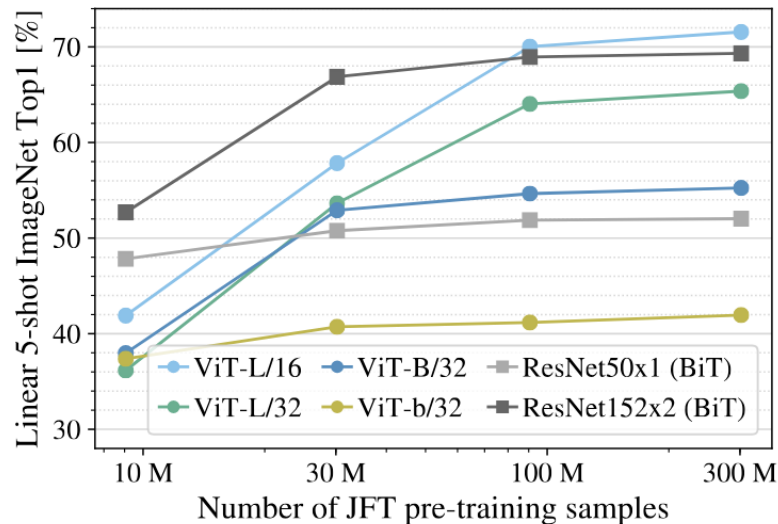


Similar to:  
Gabor filters, edge detection, color detection

This is almost impossible to  
be extracted by CNNs

# Scaling laws of transformer networks

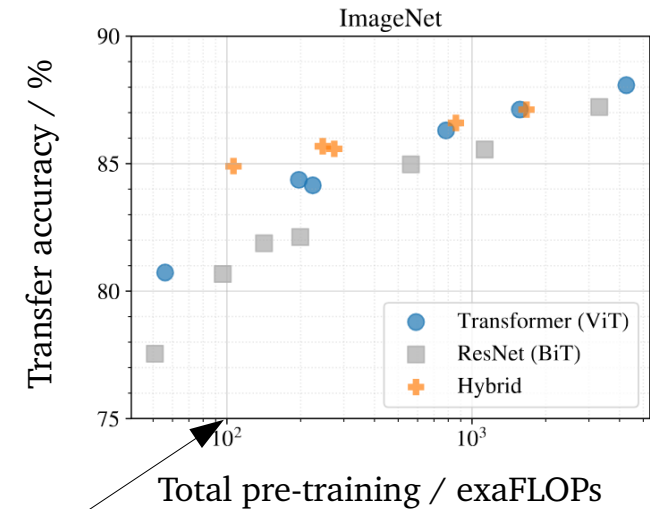
- Transformers outperform CNNs on large data sets
- CNN prior (locality, translational invariance) powerful, if less information available
- Hybrid approach: first CNN layers followed by ViT backbone → performs best



## JFT-300M

Internal Google data set with  
300M images, 1B labels

*Transformers usually  
always pre-trained!*



Consumer GPU: (1k€ - 5k€)

Single GPU: ~ 50 terraFLOPs

→ 6h for 1 exaFLOPs

Second bin: 1 month training time!

# CIFAR 10 – Transformer Networks

Model – add:

- Implement MLP
- Add MultiHeadAttention layer
- Add MLP to transformer block
- Add DNN output

Model – modify:

- Numbers of transformer blocks
- Hidden layer size
- Optimizer, learning rate, epochs



➤ **Can you achieve >70% validation accuracy?**

# Bonus: Try PointCloud Transformer

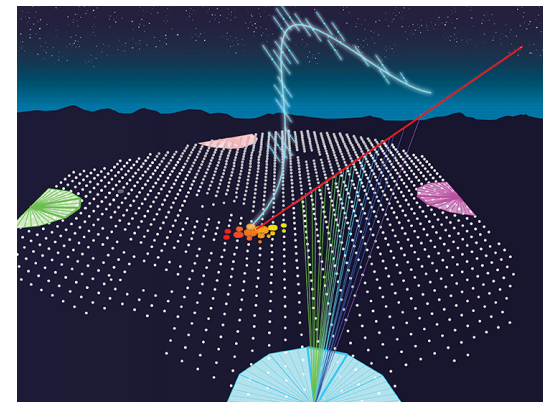
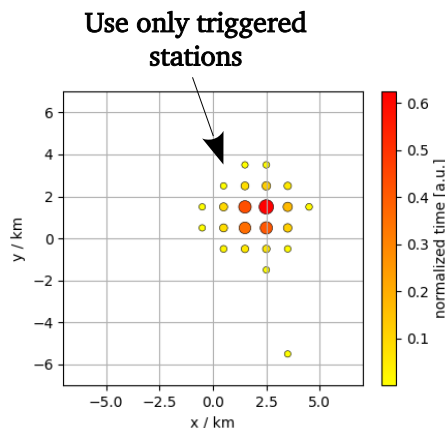
Model – add:

- Implement MLP
- Add MultiHeadAttention layer
- Add MLP to transformer block
- Add DNN output

Model – modify:

- Numbers of transformer blocks
- Hidden layer size
- Optimizer, learning rate, epochs

➤ **Can you achieve a good energy resolution?**





# Summary

- We only touched surface of deep learning in the last days
  - Field is rapidly evolving
- What's newer not always better for your task: tree learning still powerful (*if you have low dimensional data / small datasets*)
- You're now prepared for exploring advanced neural network stuff

Check out more advanced stuff

## Covered basics

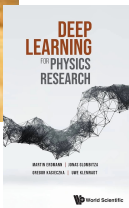
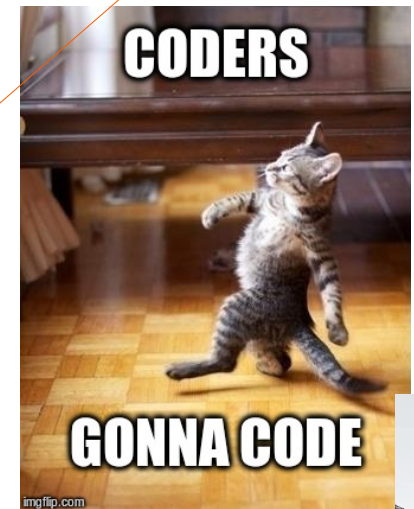
- DNN training
- Generalization, Regularization
- FCN, CNN, Transformers
- Keras / tensorflow basics
- First training experience

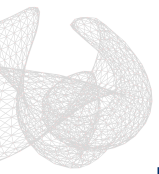
## Tryout Deep Learning Yourself!

Find many physics examples at:  
<http://www.deeplearningphysics.org/>

For example:

- CNNs, RNNs, GCNs
- GANs and WGANs
- Anomaly detection, Denosing AEs
- Visualization & introspection and more





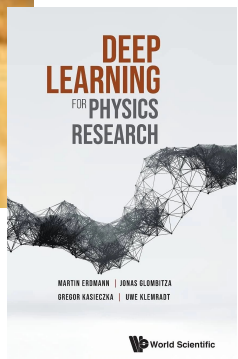
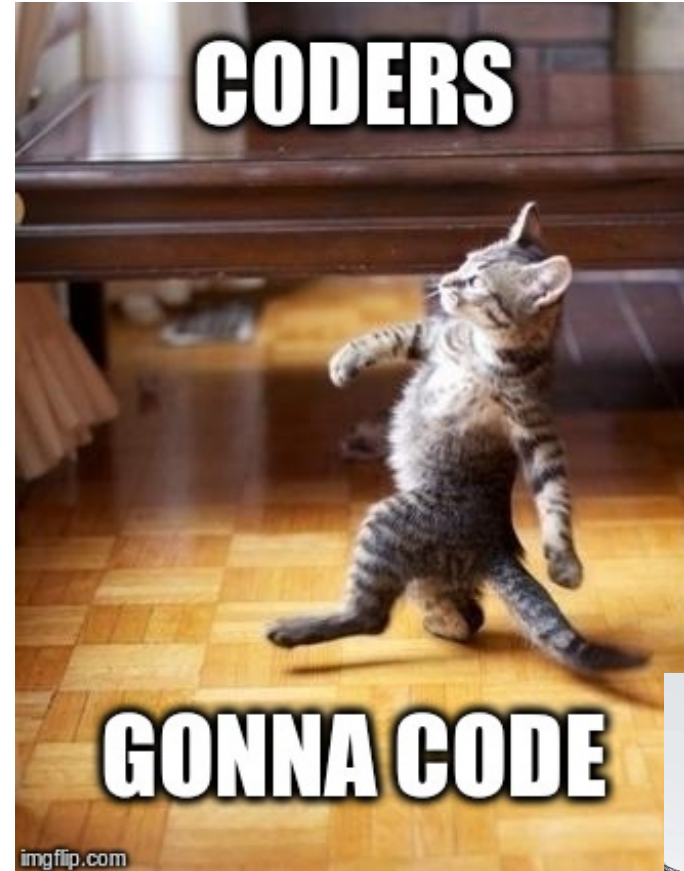
# Tryout Deep Learning Yourself!

Find many physics examples at:

<http://www.deeplearningphysics.org/>

For example:

- CNNs, RNNs, GCNs
- GANs and WGANs
- Anomaly detection, Denoising AEs
- Visualization & introspection and more



# Shapes during transformer operations

