

# EMB-Lab, Fak. N

## Activations & Applications

DLM – Deep Learning Methoden

Benjamin Kraus

Kevin Höfle, Prof. Dr. Marcus Vetter

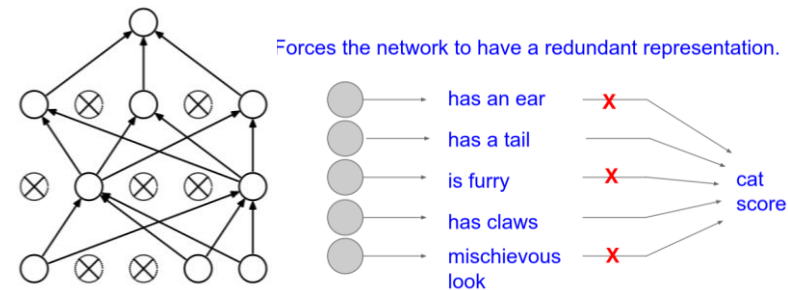
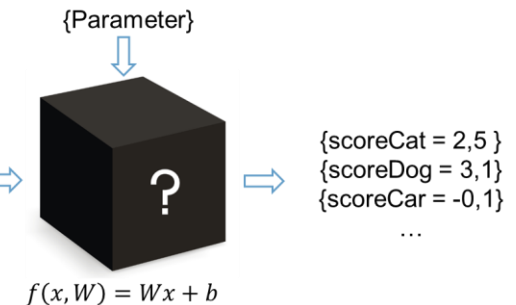
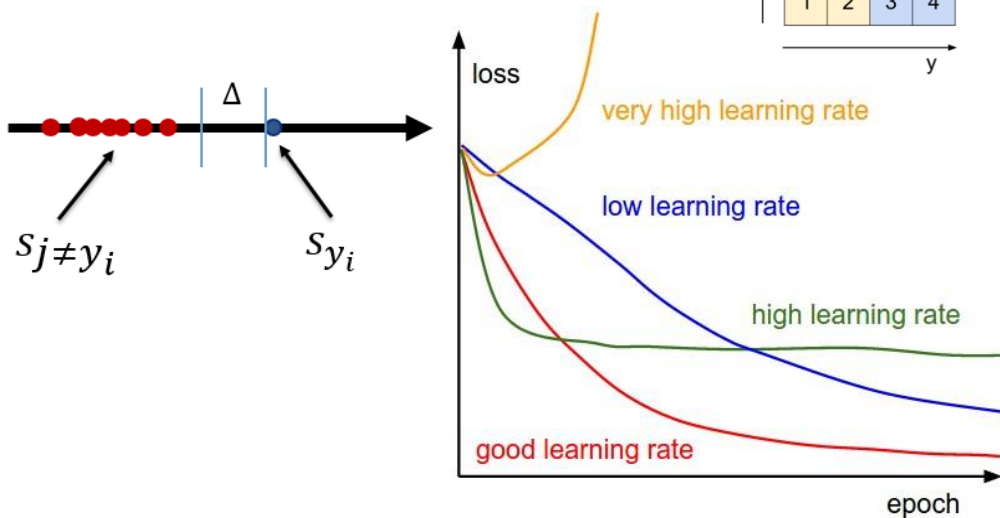
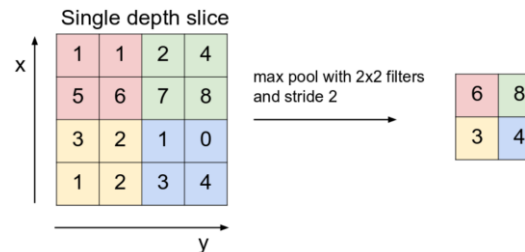
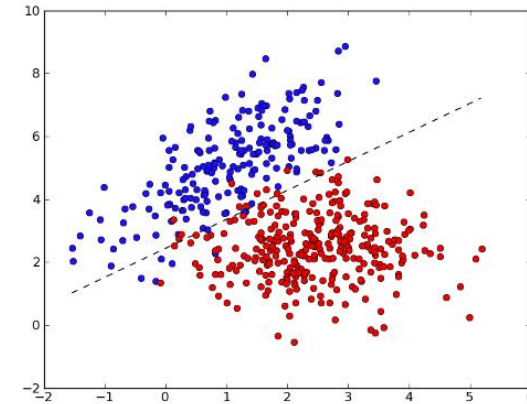
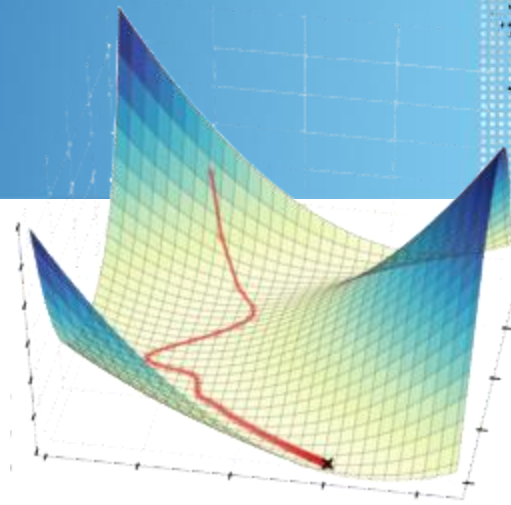
Mannheim, 11.11.2019





- Wiederholung: Fully Connected Layer, ConvLayer
- Activations
- Batchnorm
- Klassifikation, Regression:
  - Vgg16, GoogLeNet, ResNet
- Segmentierung:
  - U-Net, RCNN, YOLO
- Transfer Learning

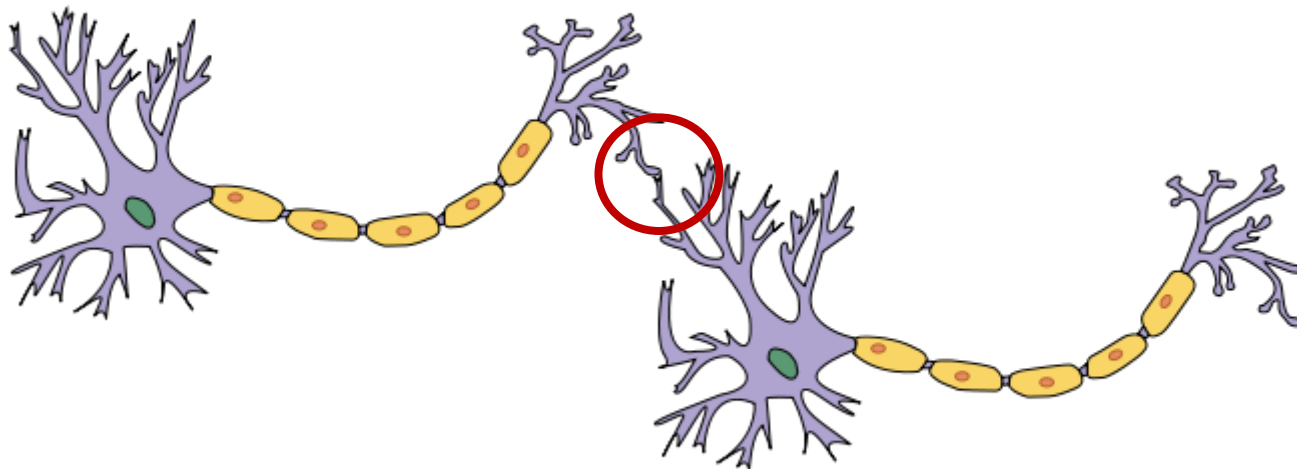
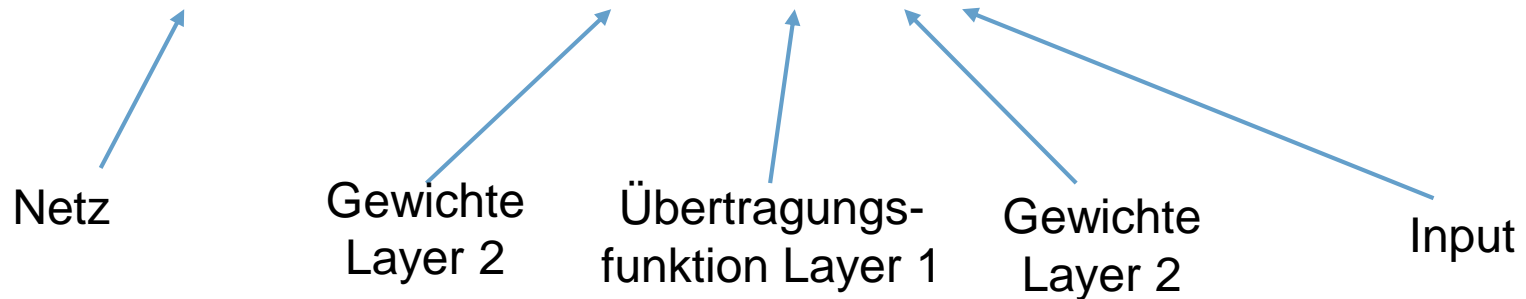
- Lineare Klassifikation
- Loss Function
- Gradienten Abstieg
- Back Propagation





- Zwei Layer NN:

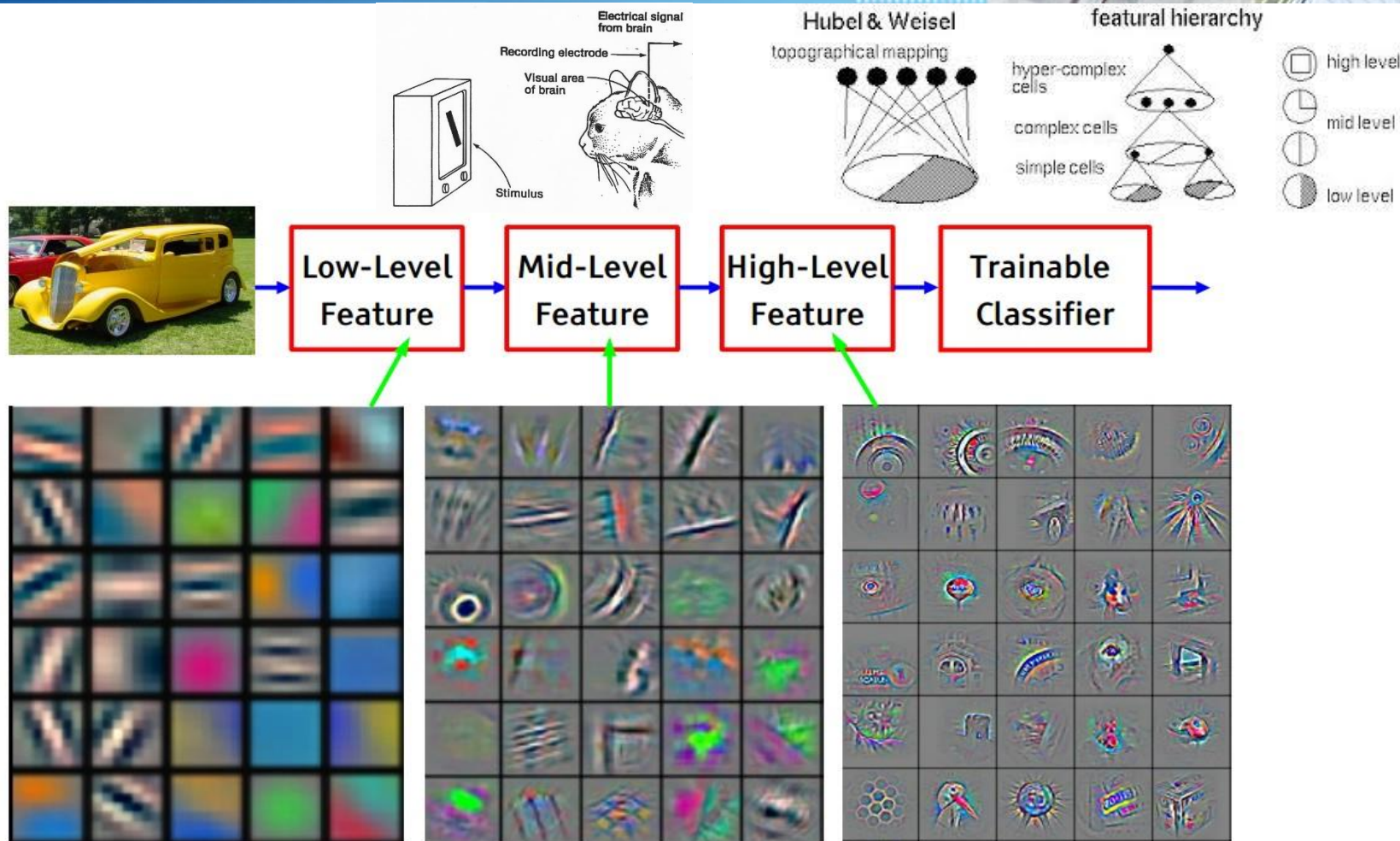
- $f(x, W_1, W_2) = W_2 * \tanh(W_1 x)$



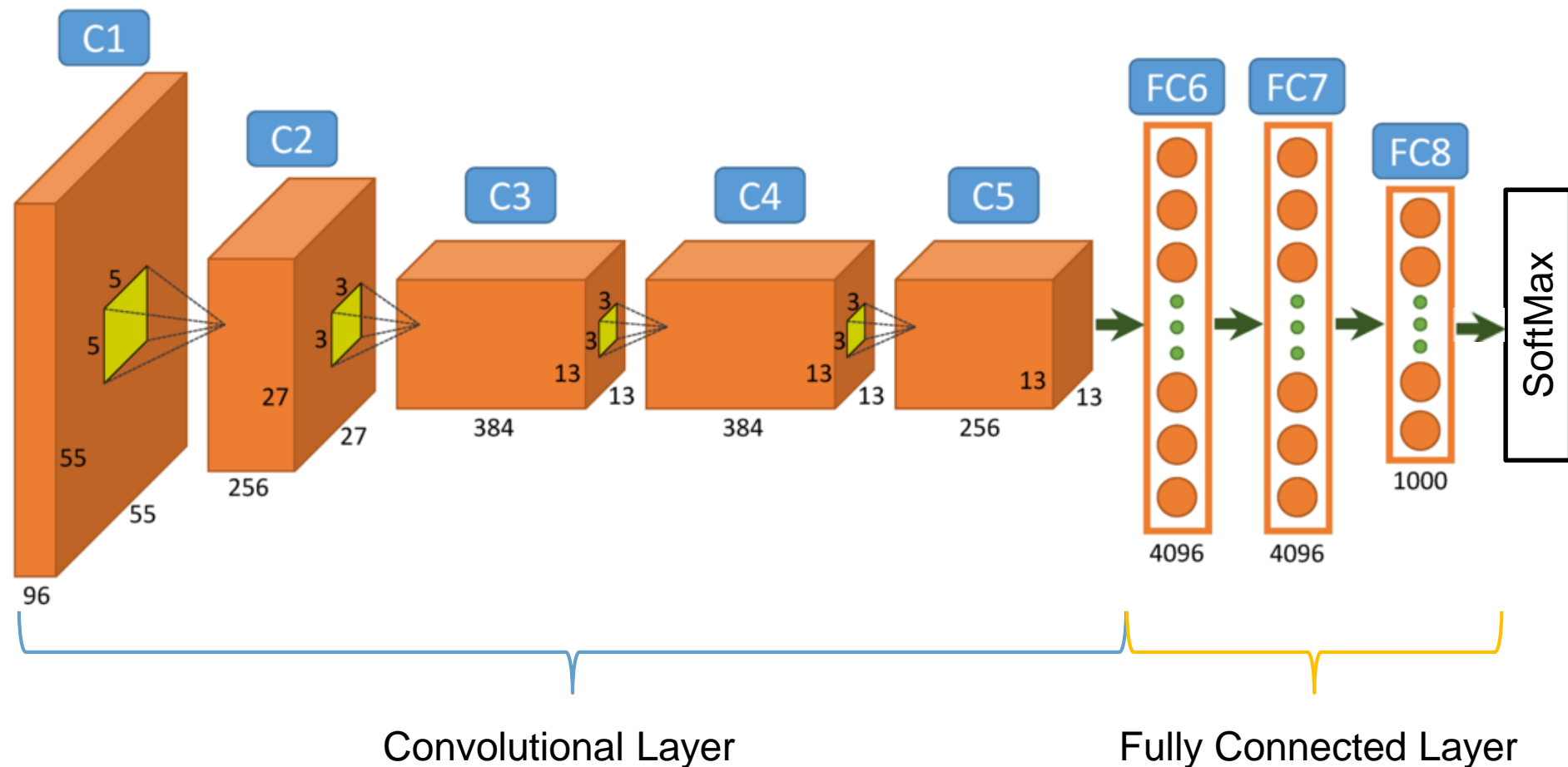




# Neuronale Netze Recap

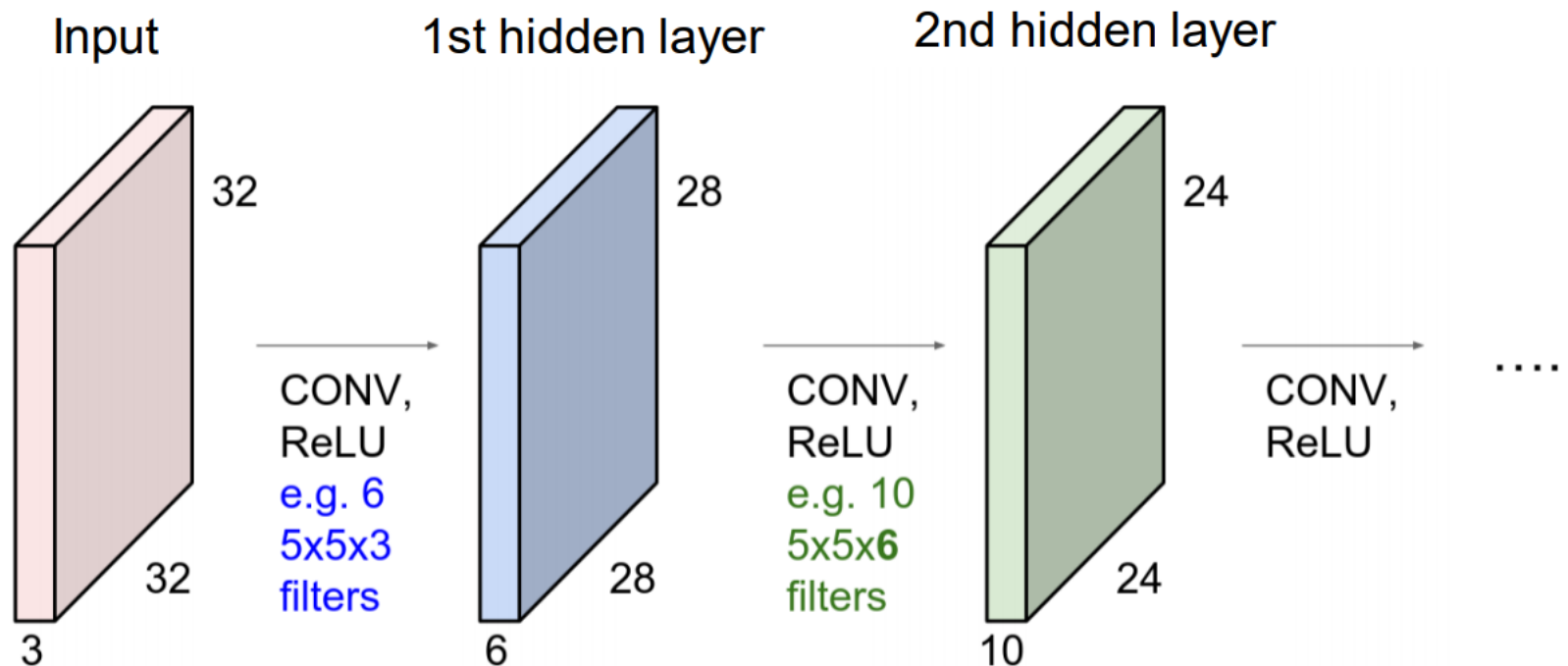


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]





Ein ConvNet besteht aus einer Folge von ConvLayers verbunden mit Aktivierungsfunktionen

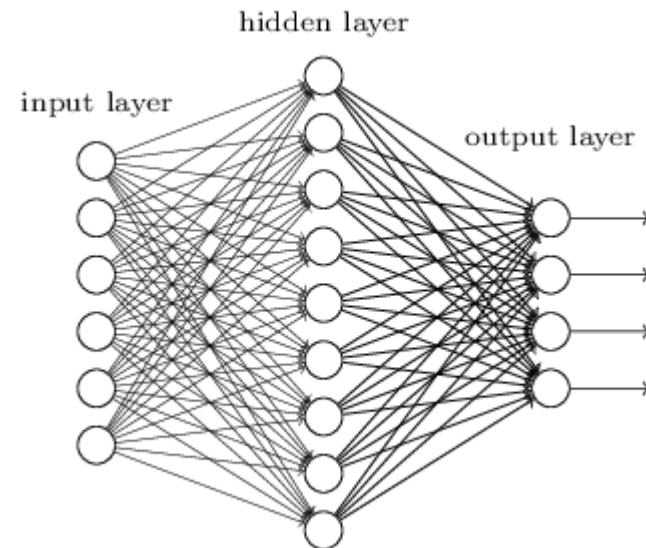




# Layer: Fully Connected

- Auch Dense Layer
- Gewichtete Verknüpfung jedes Eingangs mit jedem Ausgang
  - Sehr viele Gewichte
  - $\dim(\text{Input}) * \dim(\text{Output})$

$$\sum_{i=1}^n xW + b$$

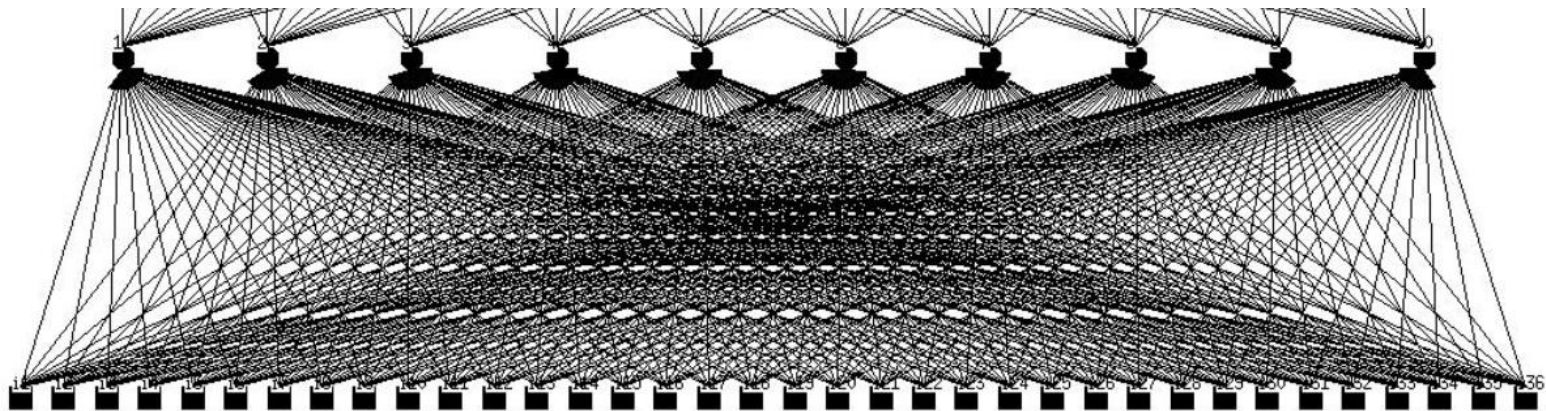






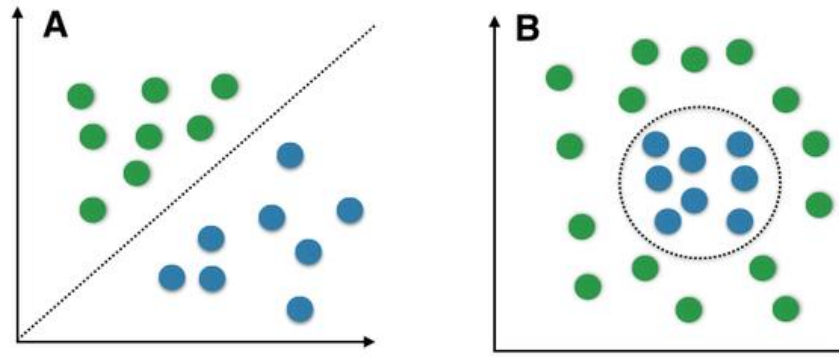
# Layer: Fully Connected

- $W = w_{ij} \in \mathbb{R}^{m \times n}$
- Bsp. Bild 32x32x3 auf 10 Neuronen
  - $W$  hat 30.720 Elemente
- Bei einem UHD-Bild
  - 248.832.000 Gewichte
  - Ca. 1GB bei Float32



[Grafik: Yoshi Komiri]

- Lineare und nicht Lineare Probleme



- Wir müssen Nichtlinearität in unseren Klassifikator einführen
- Zwei Möglichkeiten:
  - Kombinationen aus den Eingangsdaten
  - Aktivierungsfunktionen nach den Layern



- Activation Functions:

- Sigmoid

$$f(x) = 1/(1 + e^x)$$

- Tanh

$$f(x) = \tanh(x)$$

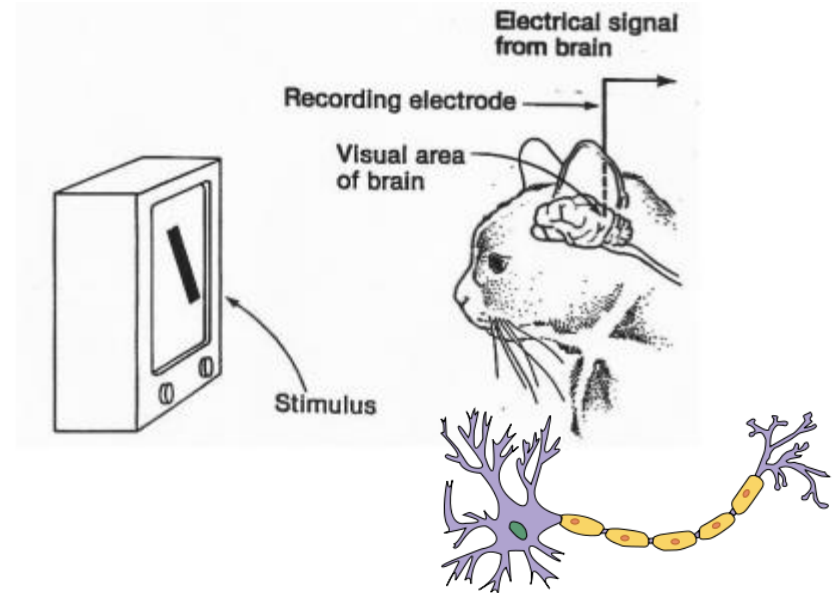
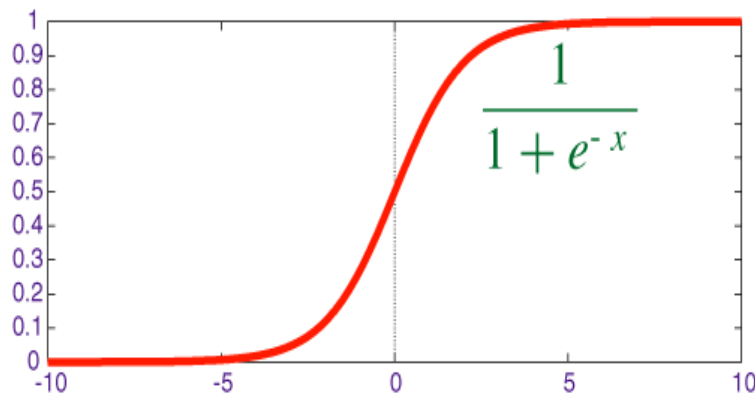
- ReLU

$$f(x) = \max(0, x)$$

- Leaky ReLU

$$f(x) = \max(ax, x)$$

- Sigmoid



“Feuerendes Neuron”

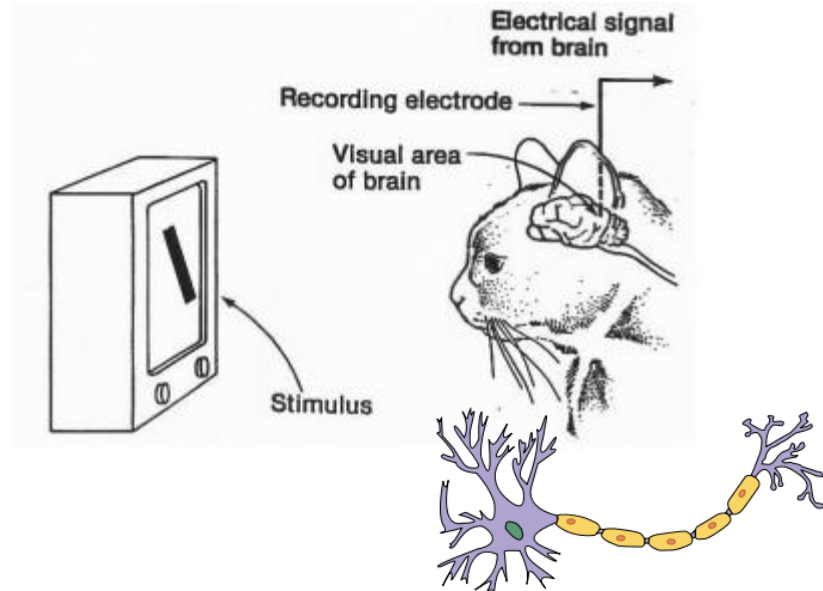
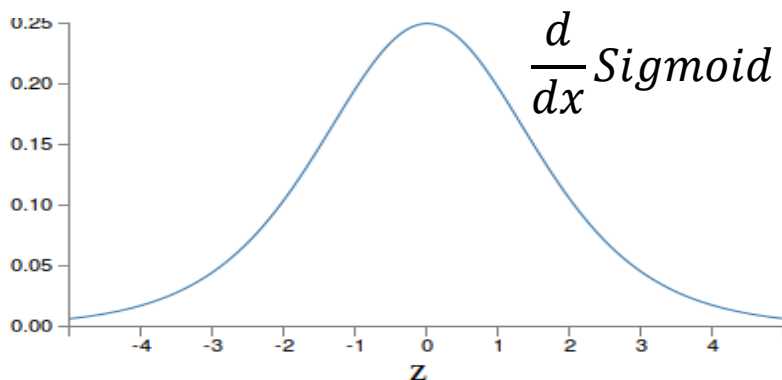
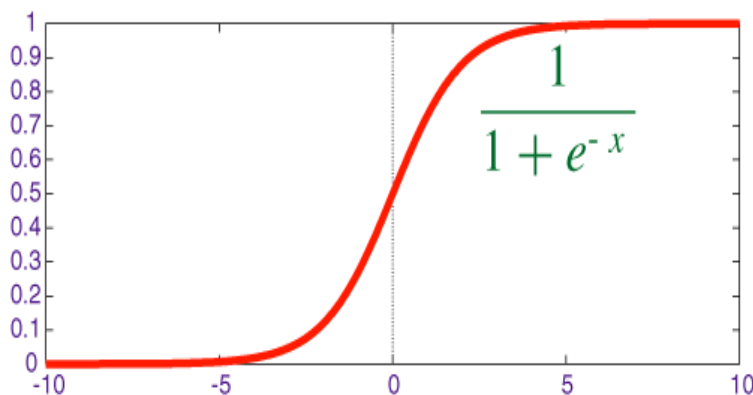
Geht in Sättigung

Non zero-centered

Computationally Expensive



- Sigmoid

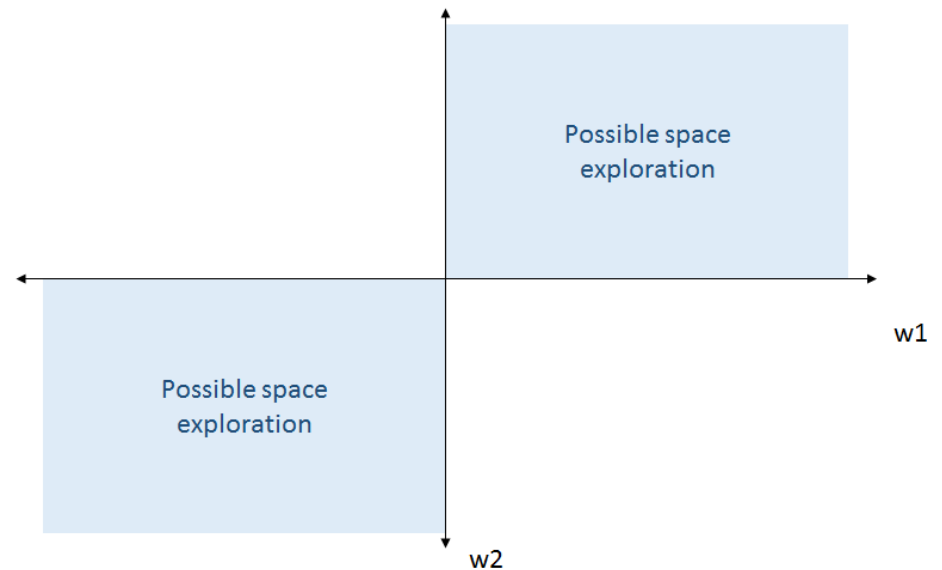
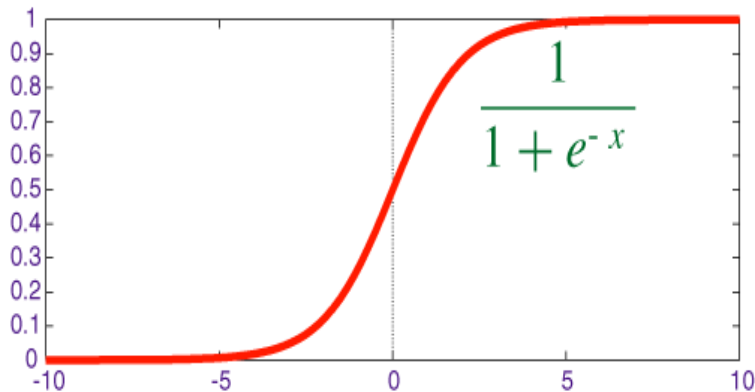


“Feuerendes Neuron”  
Geht in Sättigung  
Non zero-centered  
Computationally Expensive



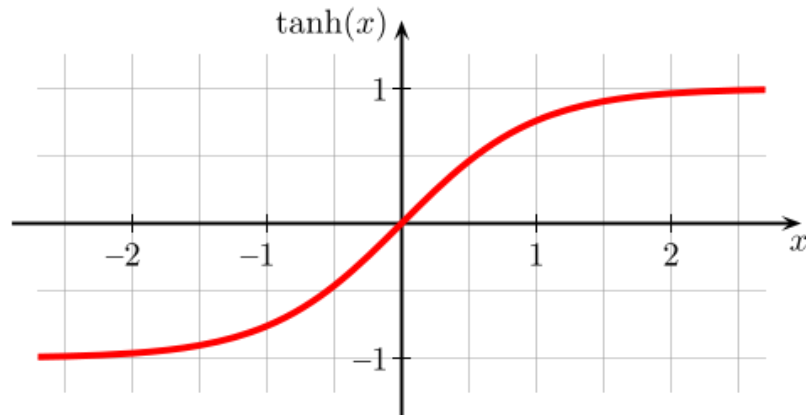


- Sigmoid
  - Non zero-centered
  - Wertebereich nur Positiv
  - Beschränkt den Freiheitsgrad beim Lernen





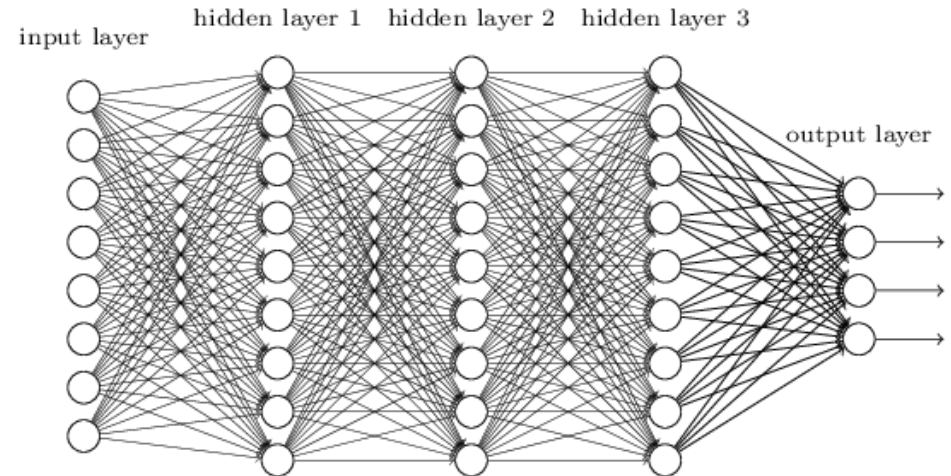
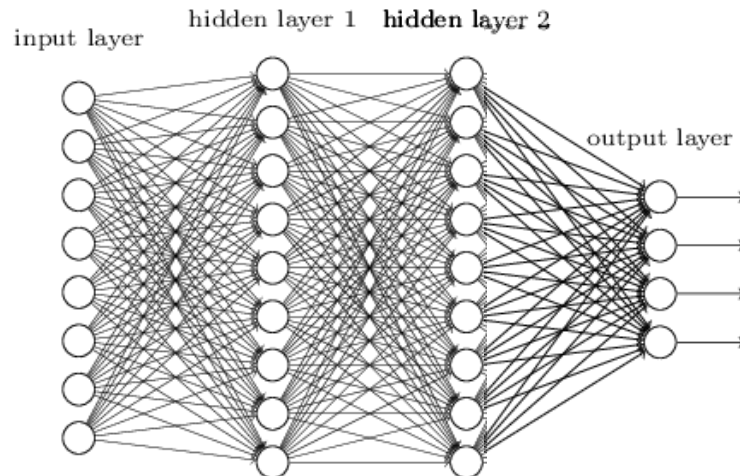
- Tangens hyperbolicus ( $\tanh$ )



- Sehr ähnlich zur Sigmoid Funktion
- Wertebereich  $[-1, 1]$
- Immernoch Sättigung



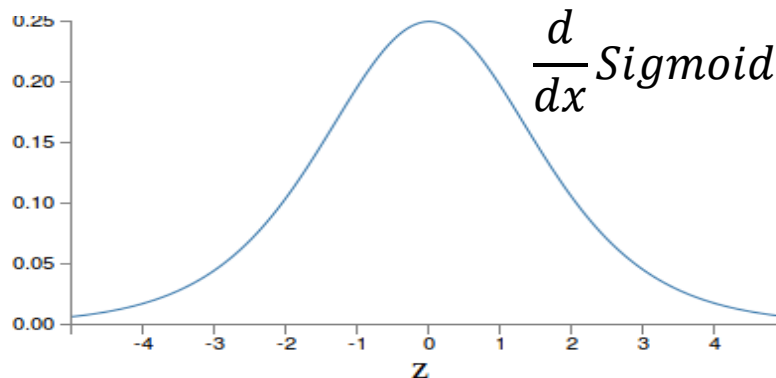
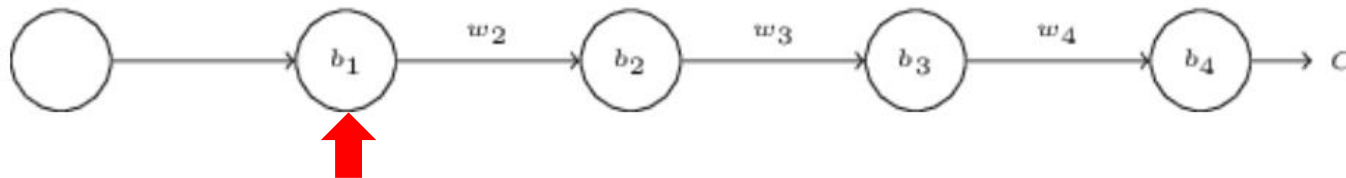
- Sigmoid / tanh
- Was Passiert bei tieferen Netzen?





- Sigmoid / tanh

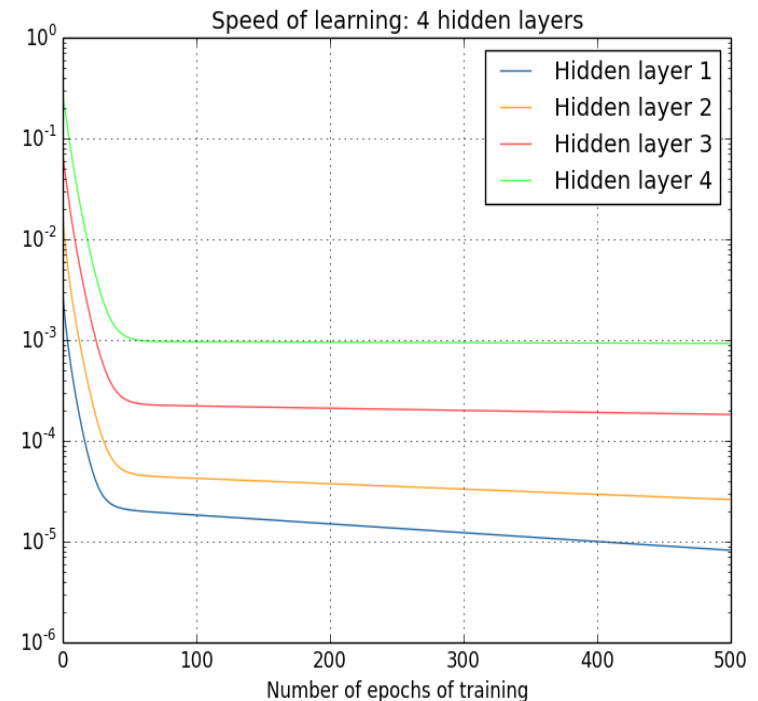
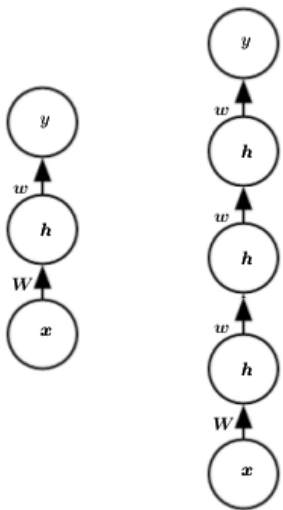
$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Geht in Sättigung?  
Was bedeutet das für den  
Gradient?



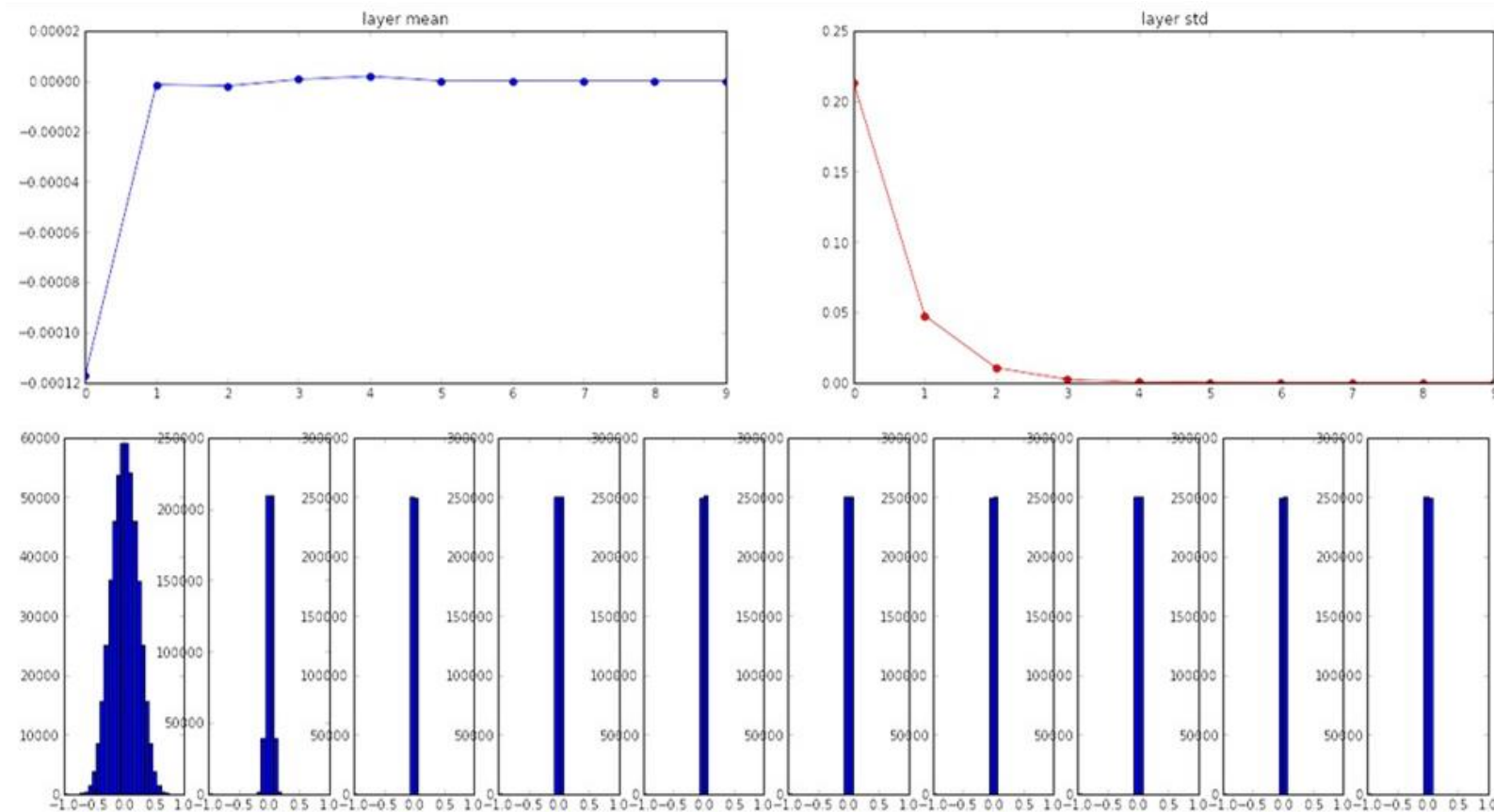
- Sigmoid / tanh
- Was Passiert bei tieferen Netzen?
- Sättigung heißt Gradient  $\sim 0$
- Die Gradienten werden pro Layer mit einer Zahl  $>1$  multipliziert
- Der weitergegebene Gradient wird immer kleiner







- Statistik: 10 Layer mit 500
- Mittelwert und Standardabweichung der Aktivierungen:

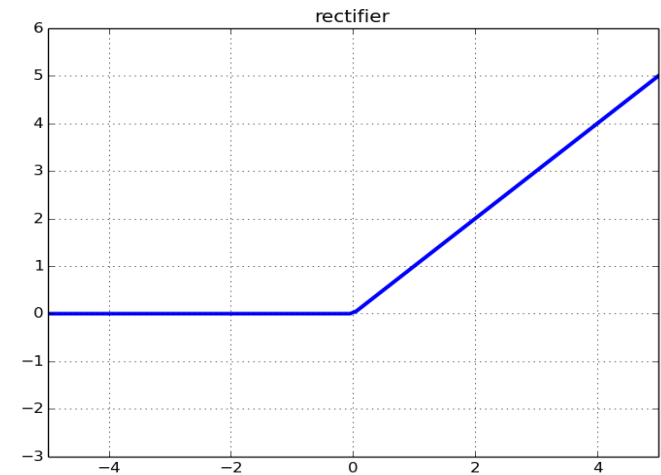




- Dieses Problem wird als „Vanishing Gradient“ bezeichnet
- Warum nicht einfach die Gradienten skalieren?
  - Das Gegenteil „Exploding Gradient“ kann entstehen



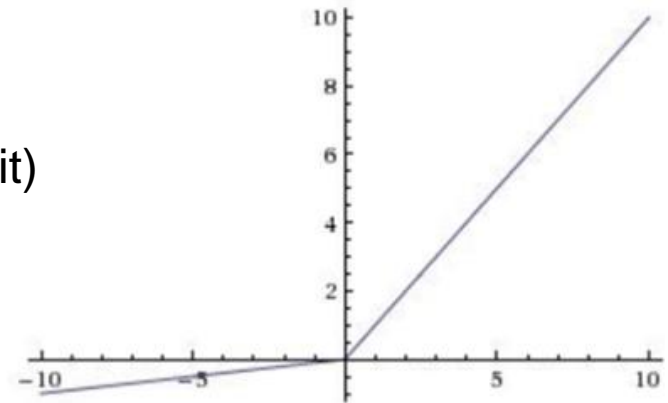
- ReLu
  - Rectified Linear Unit
    - Gradient Null (negativer Bereich)
    - Keine Sättigung (positiver Bereich)
    - Effiziente Berechnung
    - Non zero-centered
    - Konvergiert schnell



**$\max(0, x)$**



- Leaky ReLU
  - $\max(0.01x, x)$
- PReLU (Parametric Rectified Linear Unit)
  - $\max(\alpha x, x)$
- Gradient  $\neq 0$  (negativer Bereich)
- Keine Sättigung
- Effiziente Berechnung
- Non zero-centered
- Konvergiert schnell

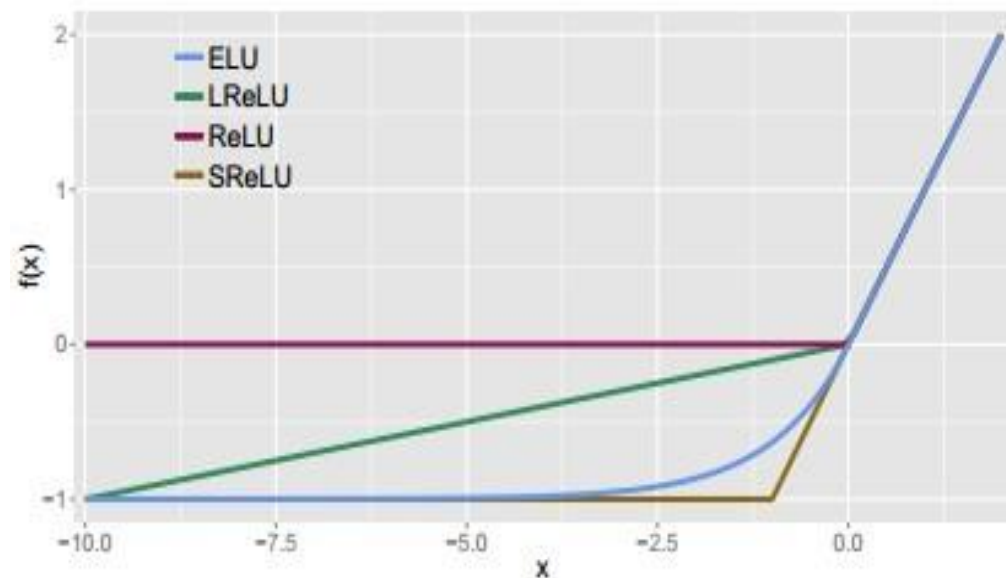




- Exponential Linear Unit (ELU)

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- Gradient  $\neq 0$  (negativer Bereich)
- Keine Sättigung
- ~~Effiziente Berechnung~~
- Non zero-centered
- Konvergiert schnell







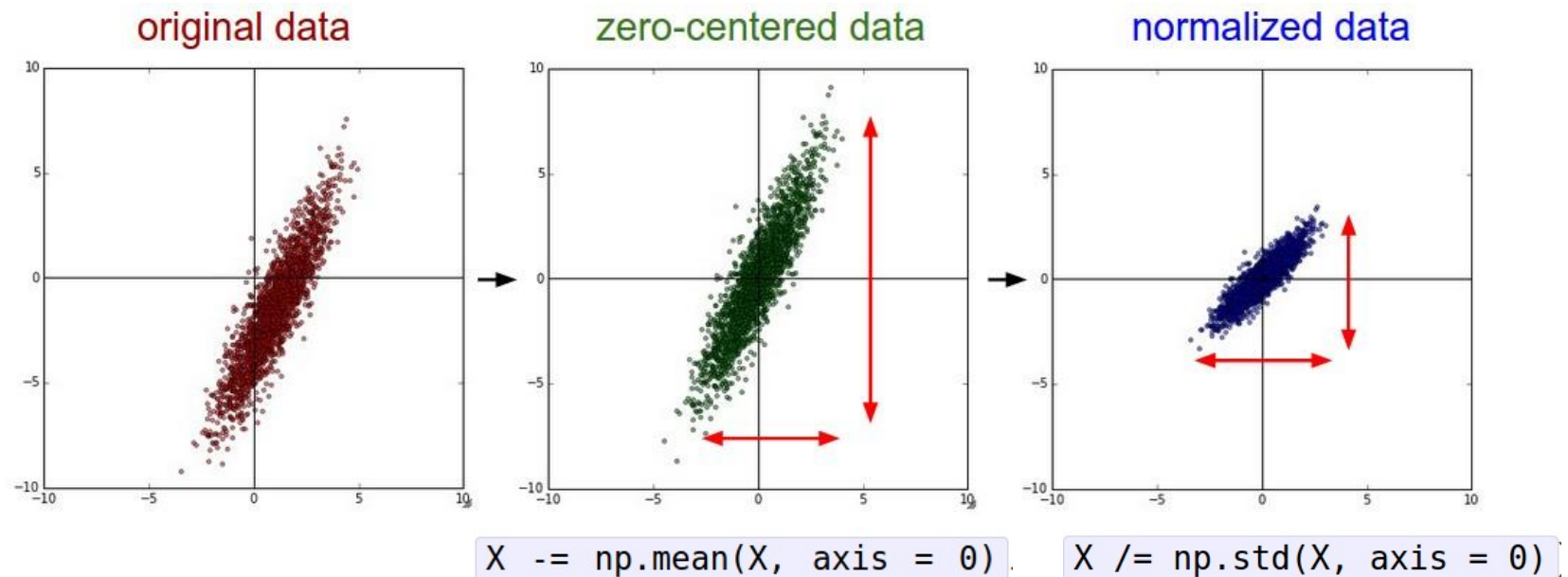
- Aus der Praxis:
  - Im Netzwerk
    - ReLU funktioniert meistens.
    - Leaky ReLU, ELU usw. sind einen Versuch wert, wenn alles andere stimmt.
    - Tanh funktioniert manchmal. Aber nicht oft.
    - Sigmoid noch weniger.
  - Am Ende des Netzwerks (nach dem Letzten Layer):
    - Was ist Ihre Problemstellung? Klassifikation? Regression?
    - Was möchten Sie? Wahrscheinlichkeiten? Koordinaten?
    - Dementsprechend wählen.



- Eingangsdaten:
- Bilder haben oft einen Wertebereich von  $[0 \dots 255]$
- Exploding Gradient Problem
- Wir hätten gern Daten von  $[-1 \dots 1]$  mit  $\sigma = 1$  und  $\mu = 0$

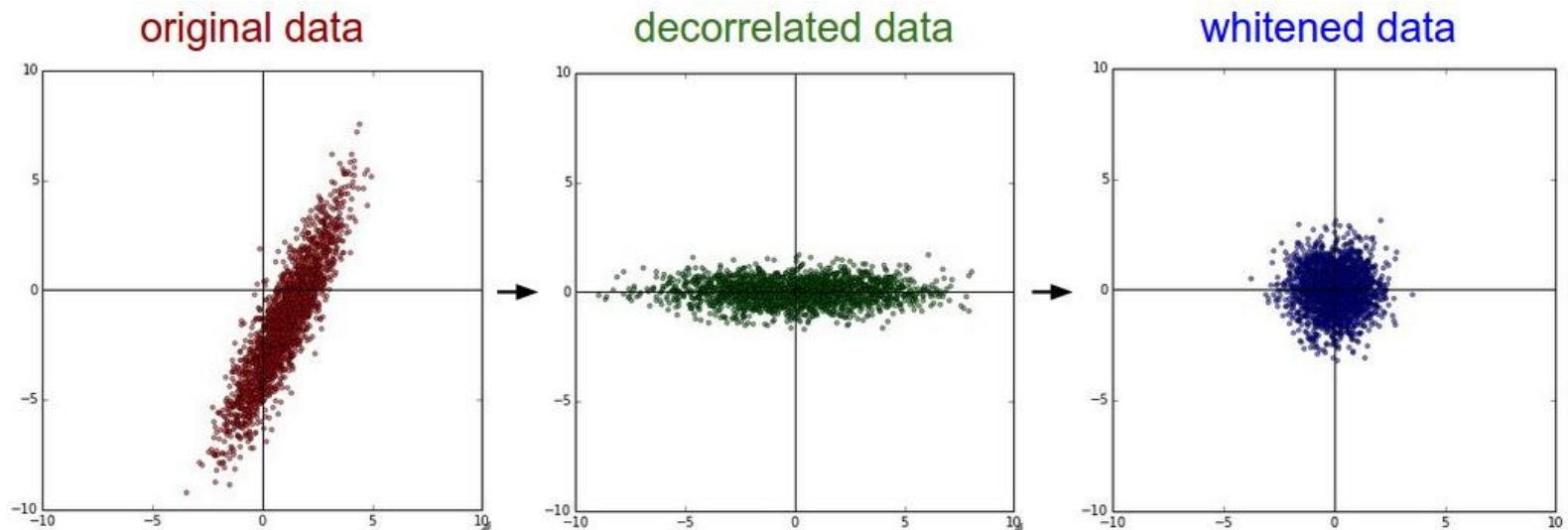


- Wir hätten gern Daten von  $[-1 \dots 1]$  mit  $\sigma = 1$  und  $\mu = 0$
- Also machen wir das einfach!





- Wir können auch noch weitergehen
- Dekorrelieren und „weiß machen“
  - z.B. mit einer PCA
  - Nicht sehr verbreitet



Kovarianz Matrix ist  
nur auf der  
Hauptachse besetzt

Kovarianz Matrix ist  
eine Einheitsmatrix

[Grafik: Karpathy]

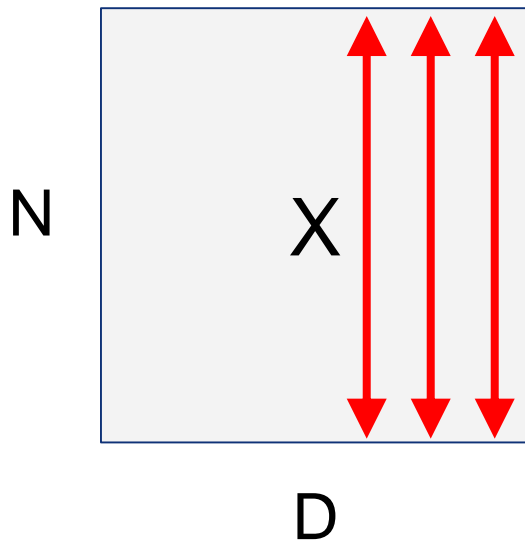


- Probleme:
  - Was ist im inneren des Netzes?
  - Was ist zur Test Zeit?





- Was ist im inneren des Netzes?
- Wir können auch die Aktivierungen im Netz normieren
  - Batch Norm Layer



1. Erwartungswert und Varianz für jede Dimension der Aktivierung berechnen

2. Normalisieren

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

[Grafik: Karpathy; Quelle: Ioffe and Szegedy]



- Batch Norm Layer

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

- Ist ableitbar -> BackProp funktioniert
- Wir erlauben es dem Netz  $\hat{x}$  zu skalieren

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

- Mit  $\gamma$  und  $\beta$  als lernbaren Parametern

Das Netzwerk kann lernen:

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbb{E}[x^{(k)}]$$

Um eine Einheitsabbildung  
zu erreichen

[Grafik: Karpathy; Quelle: Ioffe and Szegedy]



- Batch Norm Layer:
- Macht unsere Initialisierung stabiler
- Werte zwischen -1 und 1 (fast)
- Hilft gegen Vanishing- und Exploding-Gradient

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1...m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

[Grafik: Karpathy; Quelle: Ioffe and Szegedy]



- Was ist zur Test Zeit?
- Wir nutzen eine Schätzung von  $\sigma$  und  $\mu$  (z.B. das mittel aller Trainingsbatches)
- Wir können das mit einem laufenden Mittelwert anpassen

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1...m}\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

[Grafik: Karpathy; Quelle: Ioffe and Szegedy]



- Klassifikation
- Lokalisation
- Segmentierung

**Classification**



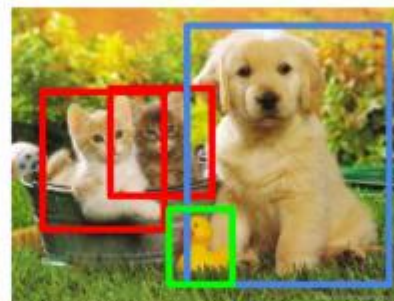
CAT

**Classification  
+ Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance  
Segmentation**



CAT, DOG, DUCK

Single object

Multiple objects

[Grafik: Jibin Mathew]

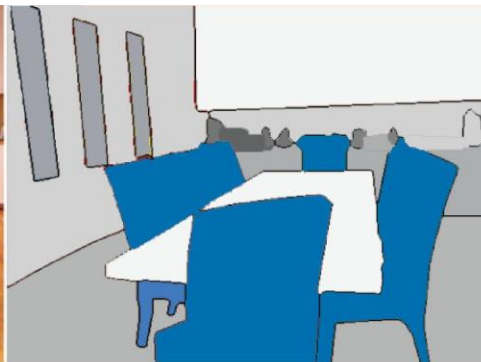




- Semantic Segmentation
- Instance Segmentation



Input Image



Semantic Segmentation

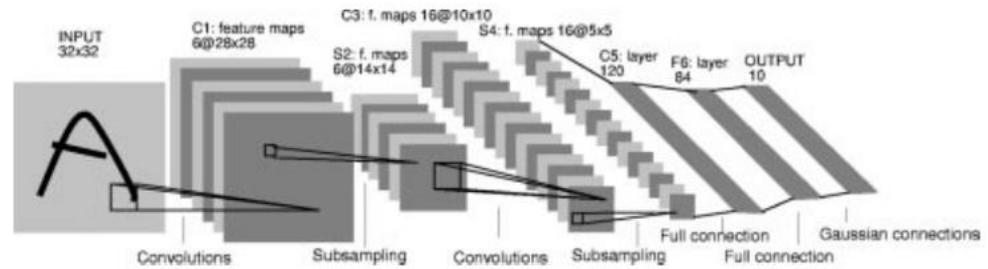


Boundary Segmentation

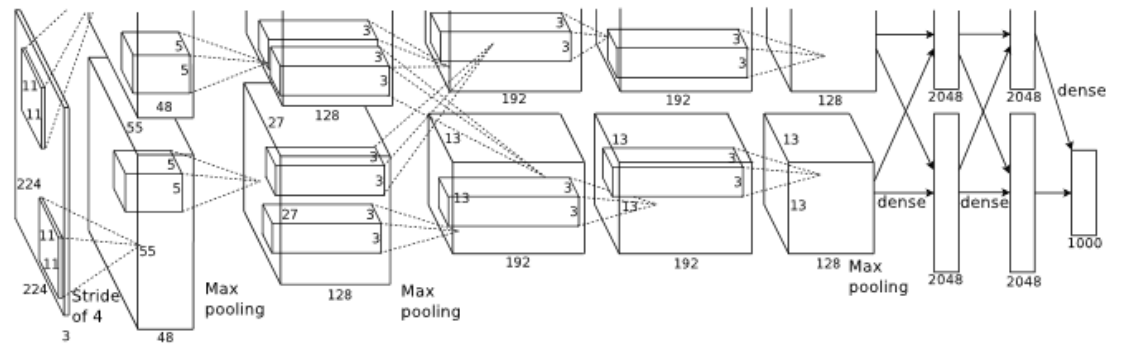


Semantic Instance Segmentation

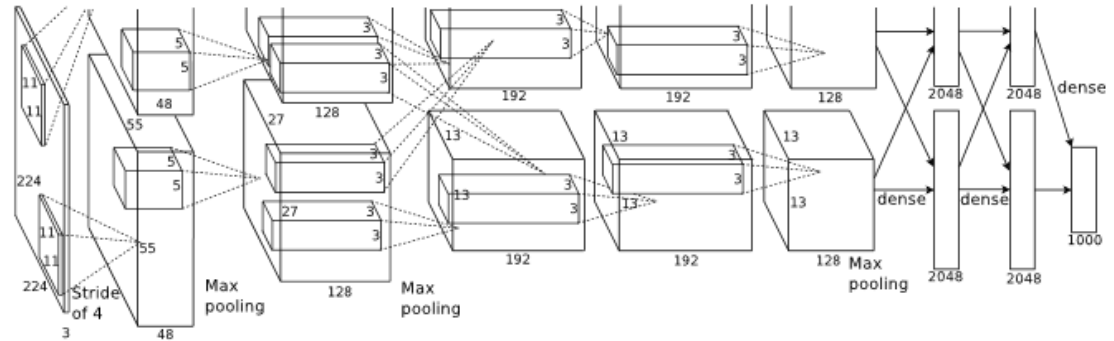
- LeNet



- AlexNet



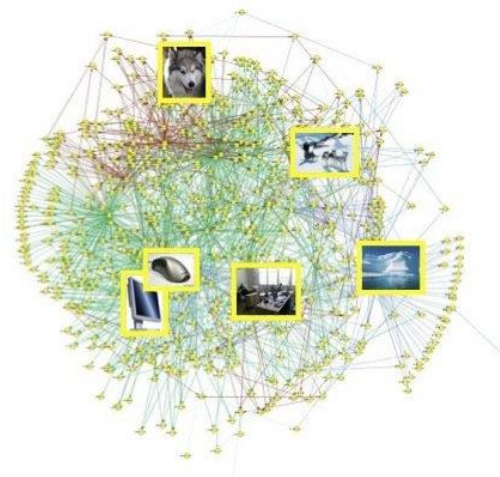
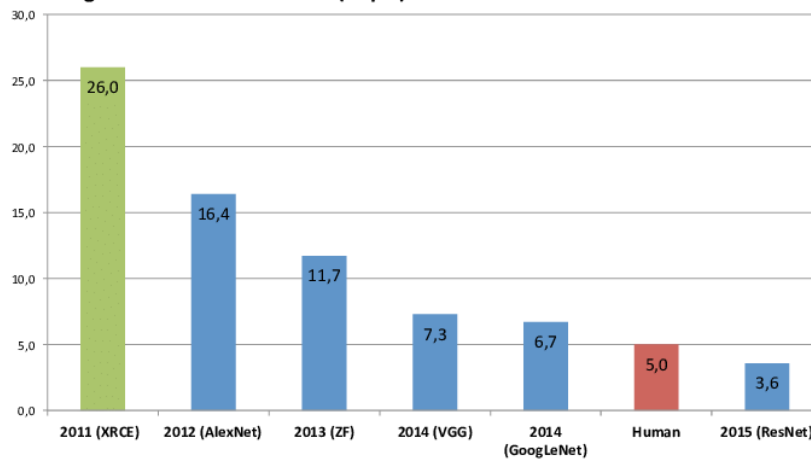
- AlexNet



Wie vergleicht man Klassifikatoren?

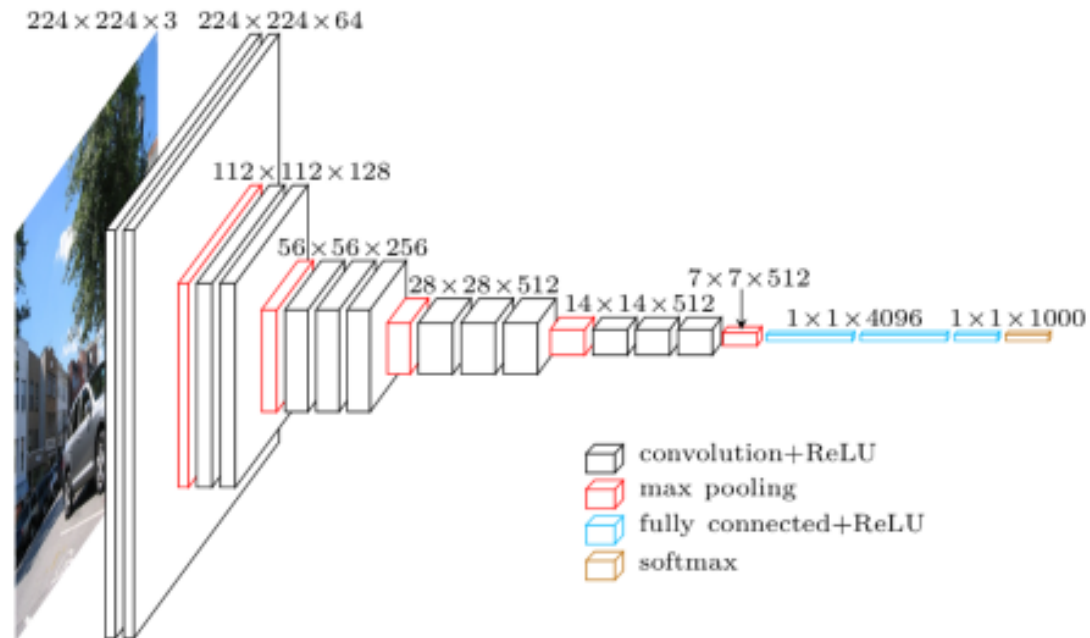
- Gewinner der ILSVRC (~1,2M Bilder, 1000 Klassen)

ImageNet Classification Error (Top 5)





- VGG16
- Sehr Einfach
- Conv -> ReLu -> Pool

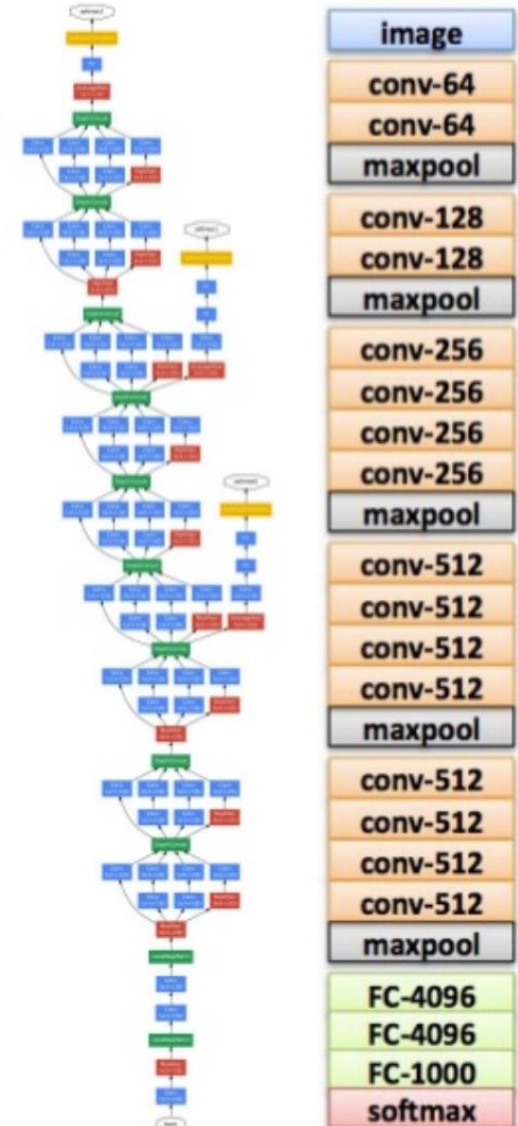
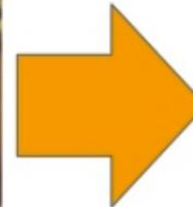
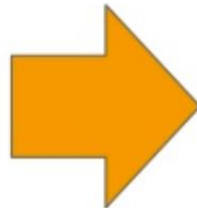


[Grafik: Davi Frossard]

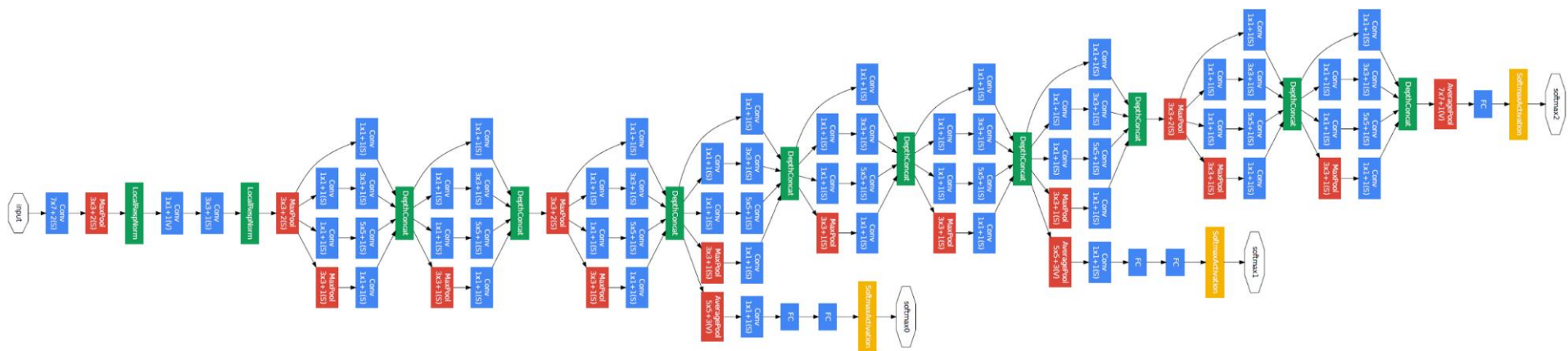


# ImageNet Challenge: 2014

## AlexNet



- GoogLeNet (auch: Inception v1)
- Going Deeper with Convolution [2014]
- Christian Szegedy et.al.

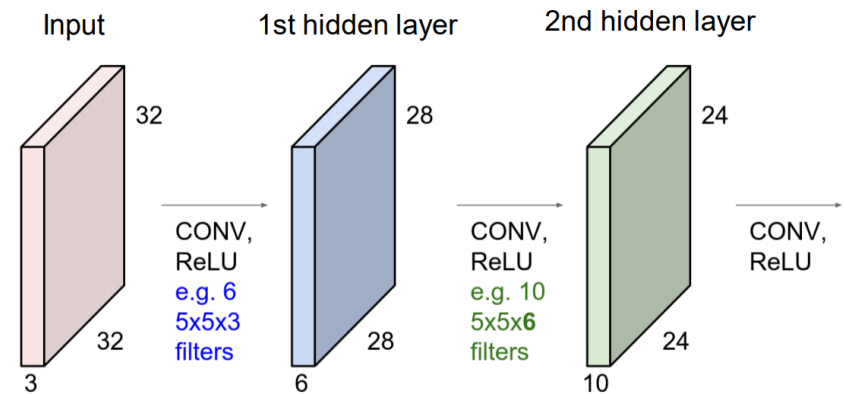




- Architektur wird trainiert
- 
- The diagram illustrates the Inception module architecture. The main part shows a detailed view of the module, which takes an input and branches into four parallel paths: a 1x1 convolution, a 3x3 convolution, a 5x5 convolution, and a 3x3 max pooling. These paths are then concatenated and passed through a DepthConcat layer. The module is shown to be part of a larger network, with a zoomed-in view of the module's internal structure and a zoomed-out view of the entire network architecture.
- Inception module
- [Grafik: Davi Frossard]

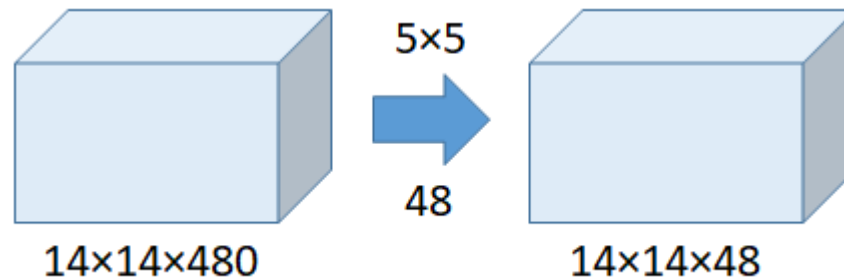


- 1x1-Kernel Faltungen?
- Was macht sowas?

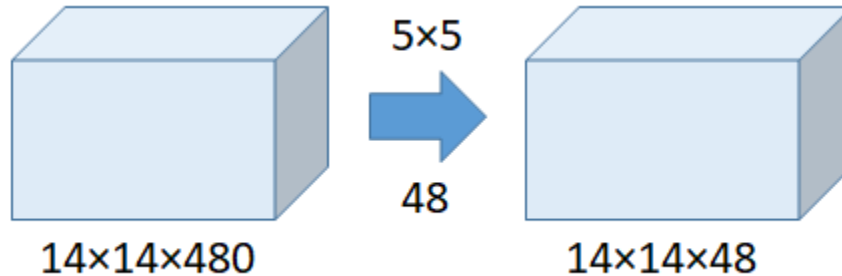




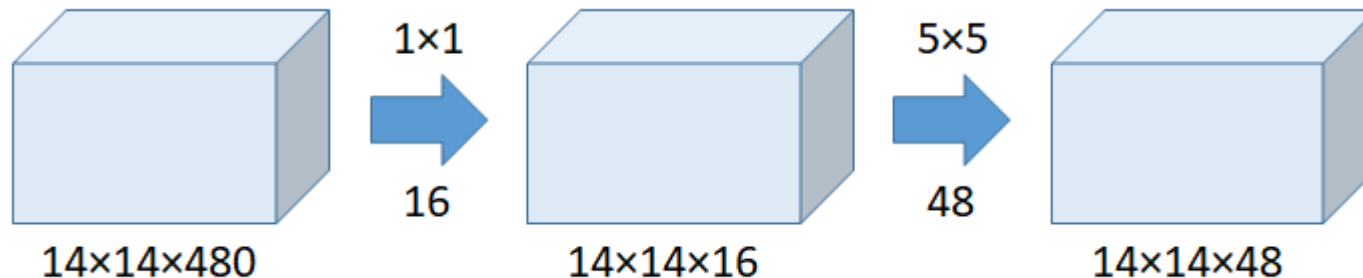
- 1x1-Kernel Faltungen?
- Was macht sowas?
- Einfaches Beispiel:



- Anzahl der Operationen  $5 \times 5 = (14 \times 14 \times 48) \times (5 \times 5 \times 16) = 112.9\text{M}$

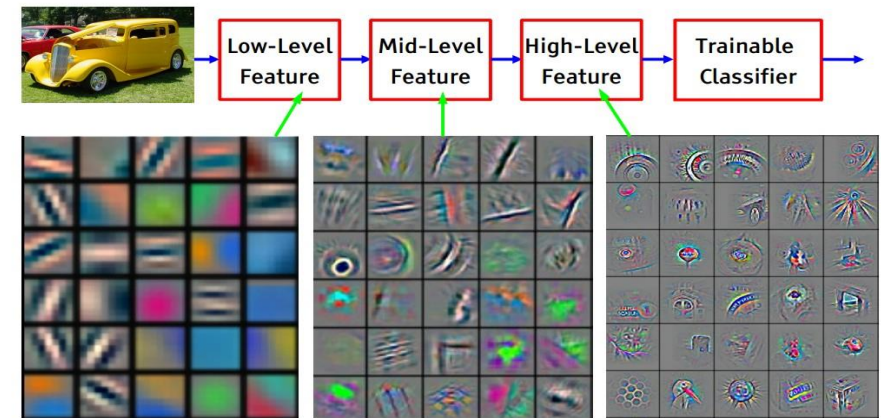
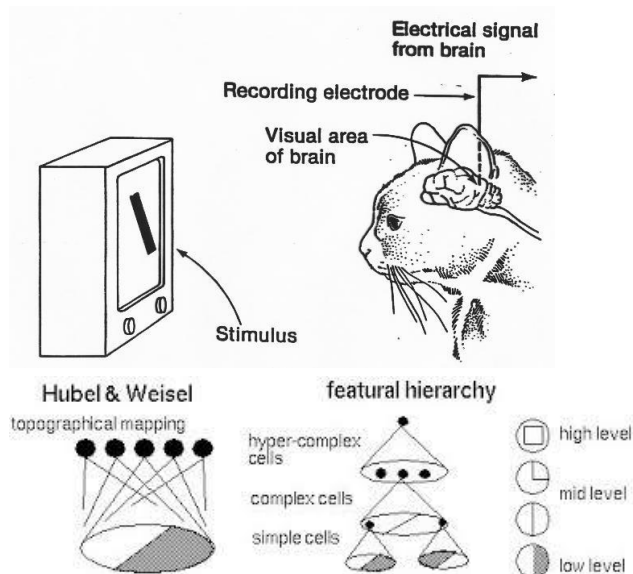


- Anzahl der Operationen =  $(14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9\text{M}$



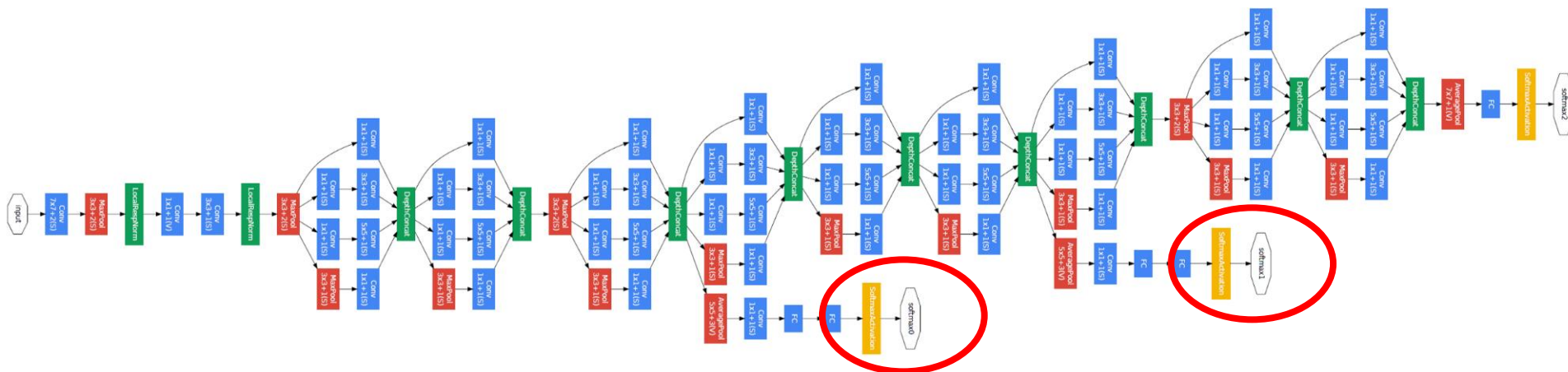
- Anzahl der Operationen für  $1 \times 1 = (14 \times 14 \times 16) \times (1 \times 1 \times 480) = 1.5\text{M}$
- Anzahl der Operationen für  $5 \times 5 = (14 \times 14 \times 48) \times (5 \times 5 \times 16) = 3.8\text{M}$
- Anzahl der Operationen Total =  $1.5\text{M} + 3.8\text{M} = 5.3\text{M}$

- 1x1 Faltung:
  - Reduziert Rechenaufwand enorm
  - Dimensions Reduktion (oder?)
    - Eine zusätzliche Verknüpfung der Features findet statt!
    - Nicht-Linearitäten werden hinzugefügt



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

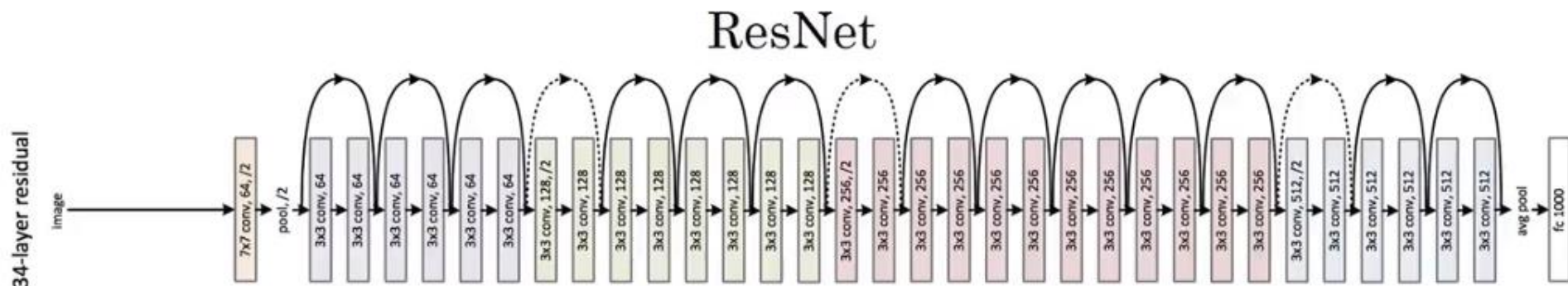
- GoogLeNet
- Auxiliary Classifiers mit Softmax
- Nur zum Training genutzt
  - Loss wird in Zwischenschichten eingebracht
  - Hilft gegen Vanishing Gradient



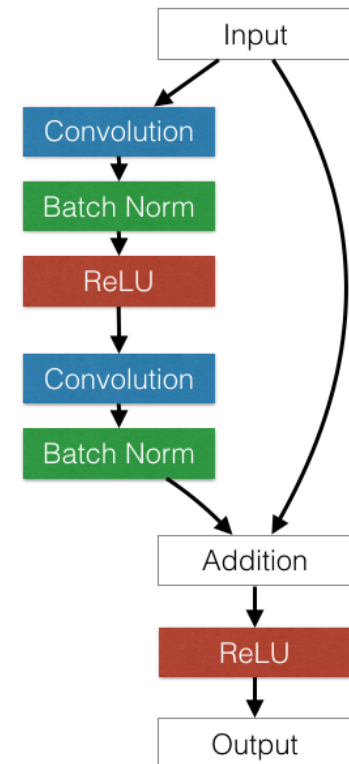




- ResNet
  - Residual Network
  - Deep Residual Learning for Image Recognition (Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, 2015)
  - ResNet-16, ResNet-34, ResNet-50...

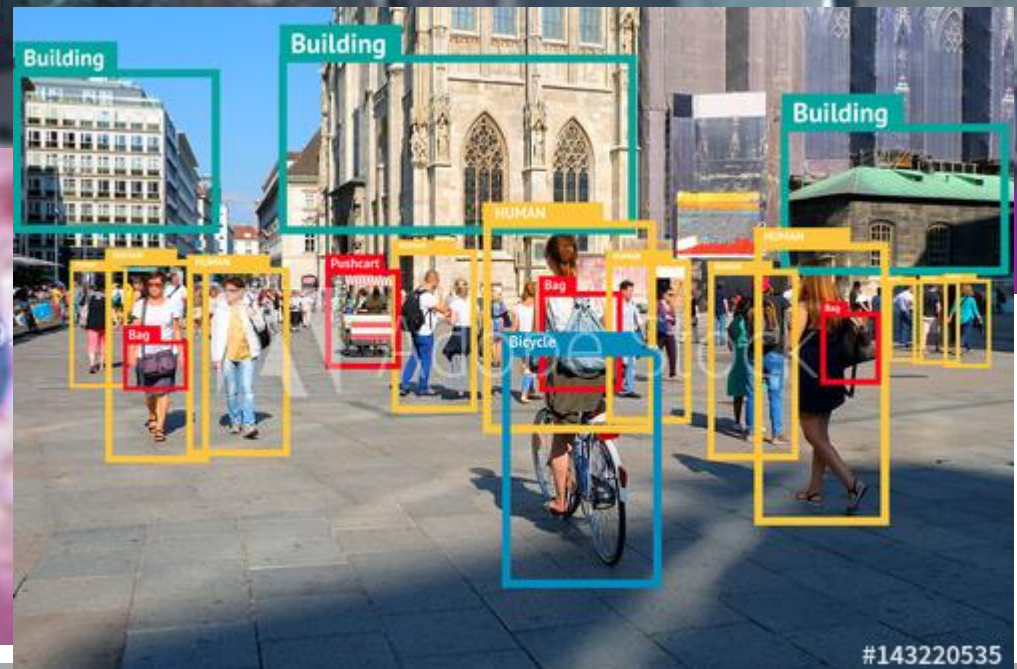
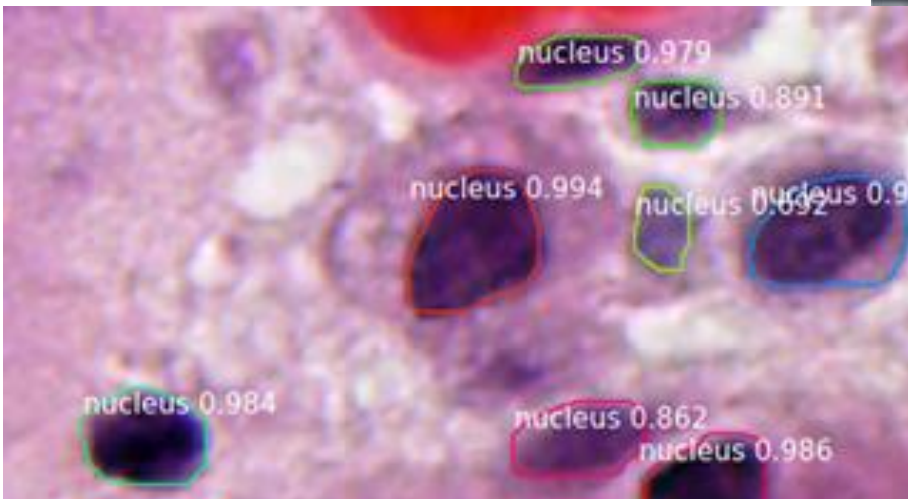
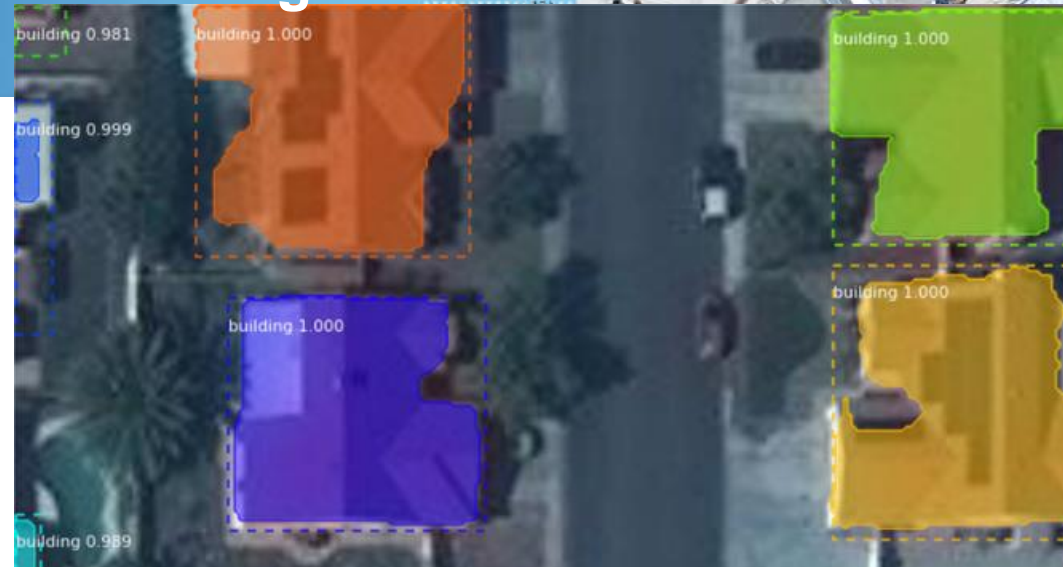


- ResNet
  - Information aus dem Input bleibt erhalten
  - Gewichtete Addition am Ende jedes Blocks
- Verringert Vanishing Gradient Problem





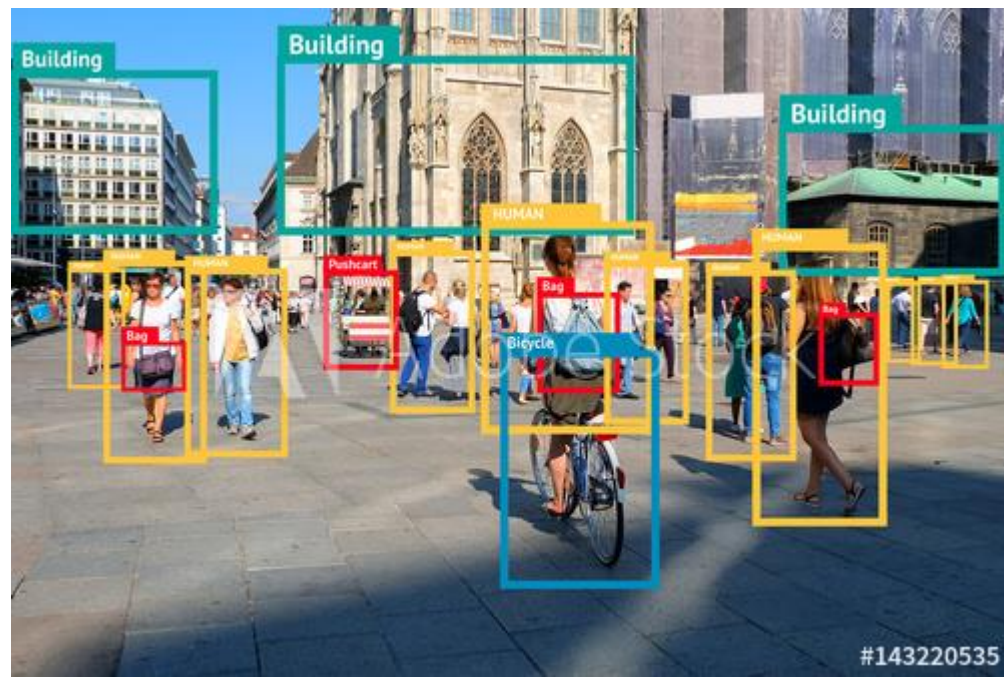
- Die R-CNN Familie
- YOLO (v1, v2...)
- U-Net





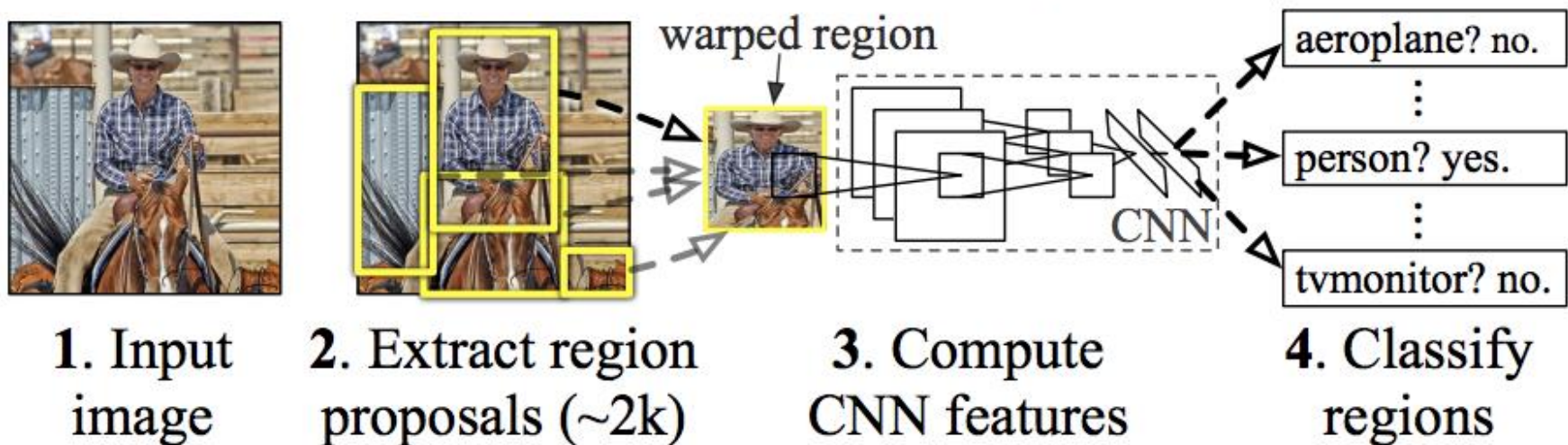


- Bounding boxes
- Vier Eckpunkte als Regression (bzw.  $x$ ,  $y$ ,  $w$ ,  $h$ )



- R-CNN (Regions with CNN features, 2014, Girshick et al.)
- Zwei Schritte:
- Region Proposals & CNN Classify

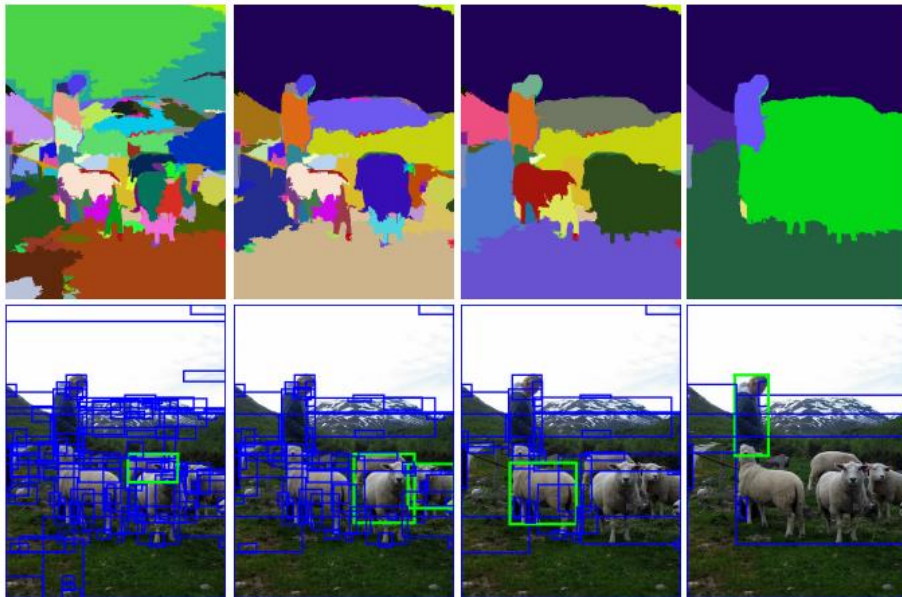
## R-CNN: *Regions with CNN features*



[Grafik: Rohith Gandhi]

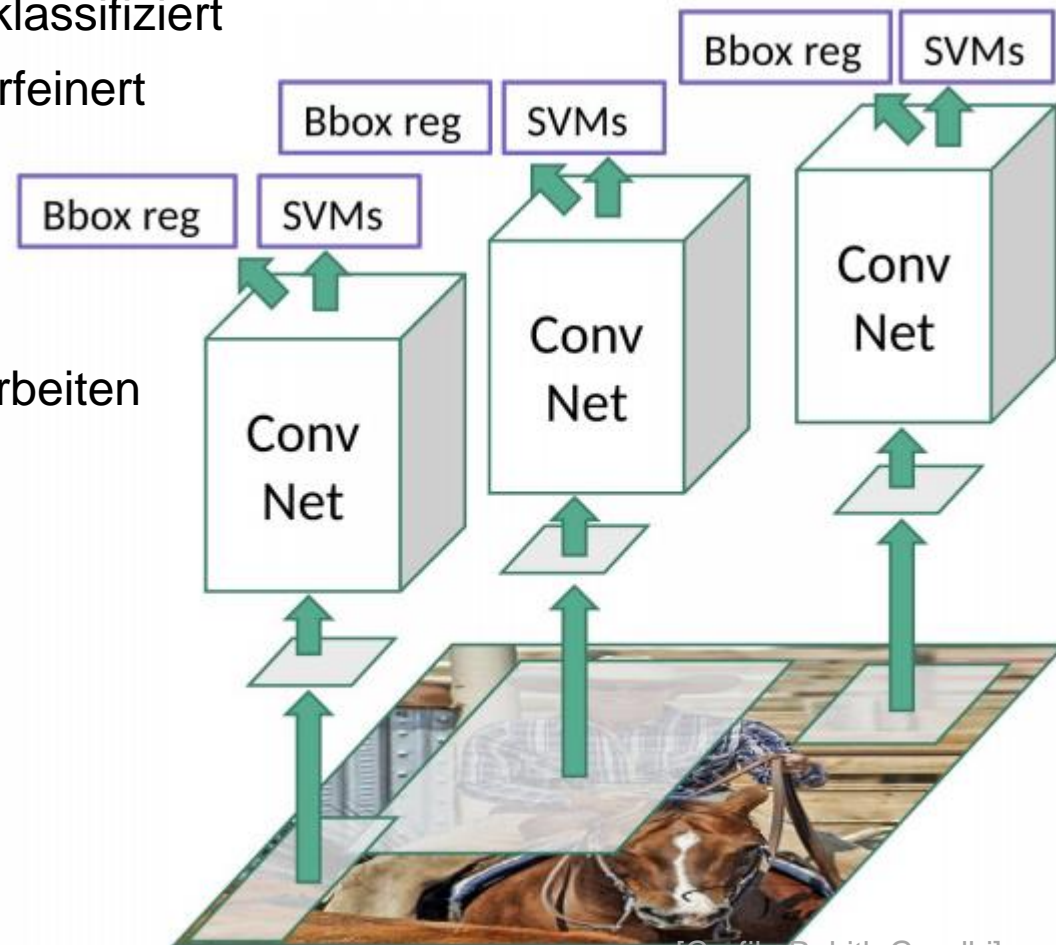


- Region Proposals:
  - Selective Search for Object Recognition [J.R.R. Uijlings et.al, 2012]
  - Statischer Algorithmus



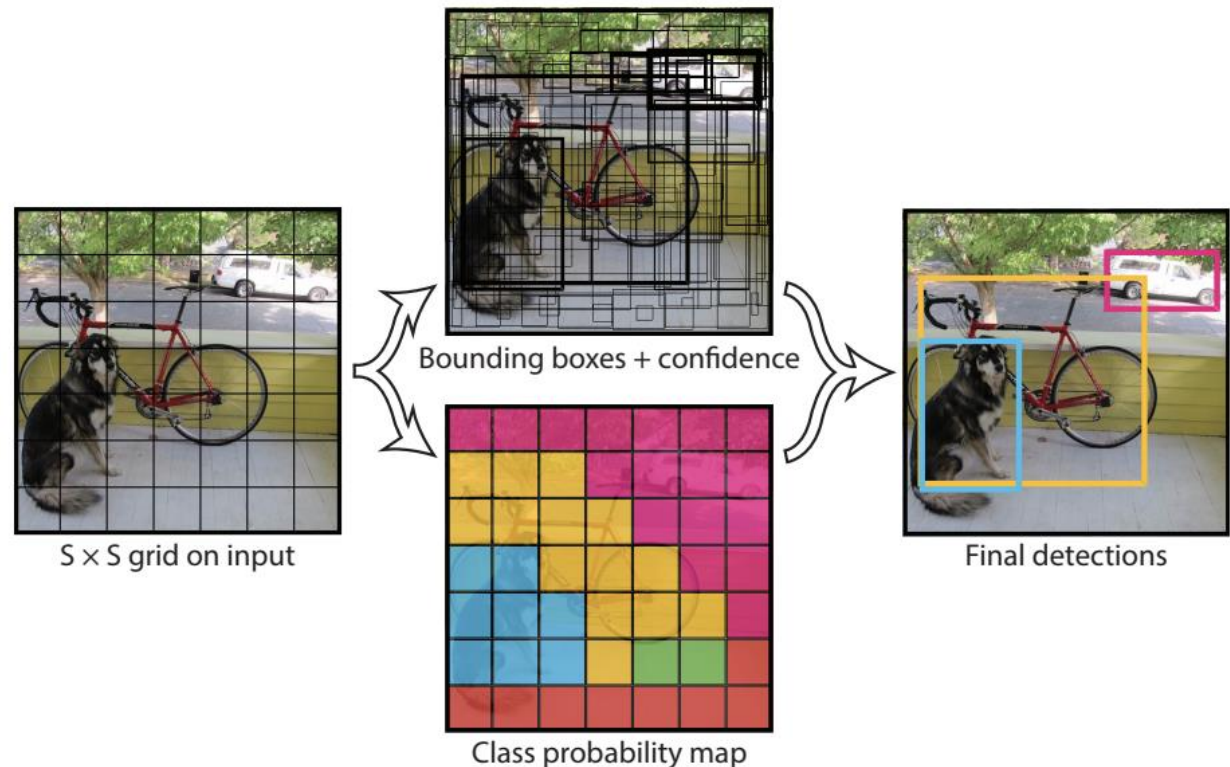


- CNN Classify und Bounding Box Regression
  - Region Proposals werden klassifiziert
  - BoundingBoxes werden verfeinert
- Sehr langsam (ca. 0.07fps)
- Mittlerweile viele Nachfolgearbeiten
  - Fast-R-CNN
  - Faster-R-CNN
  - Mask-R-CNN (ca. 5fps)

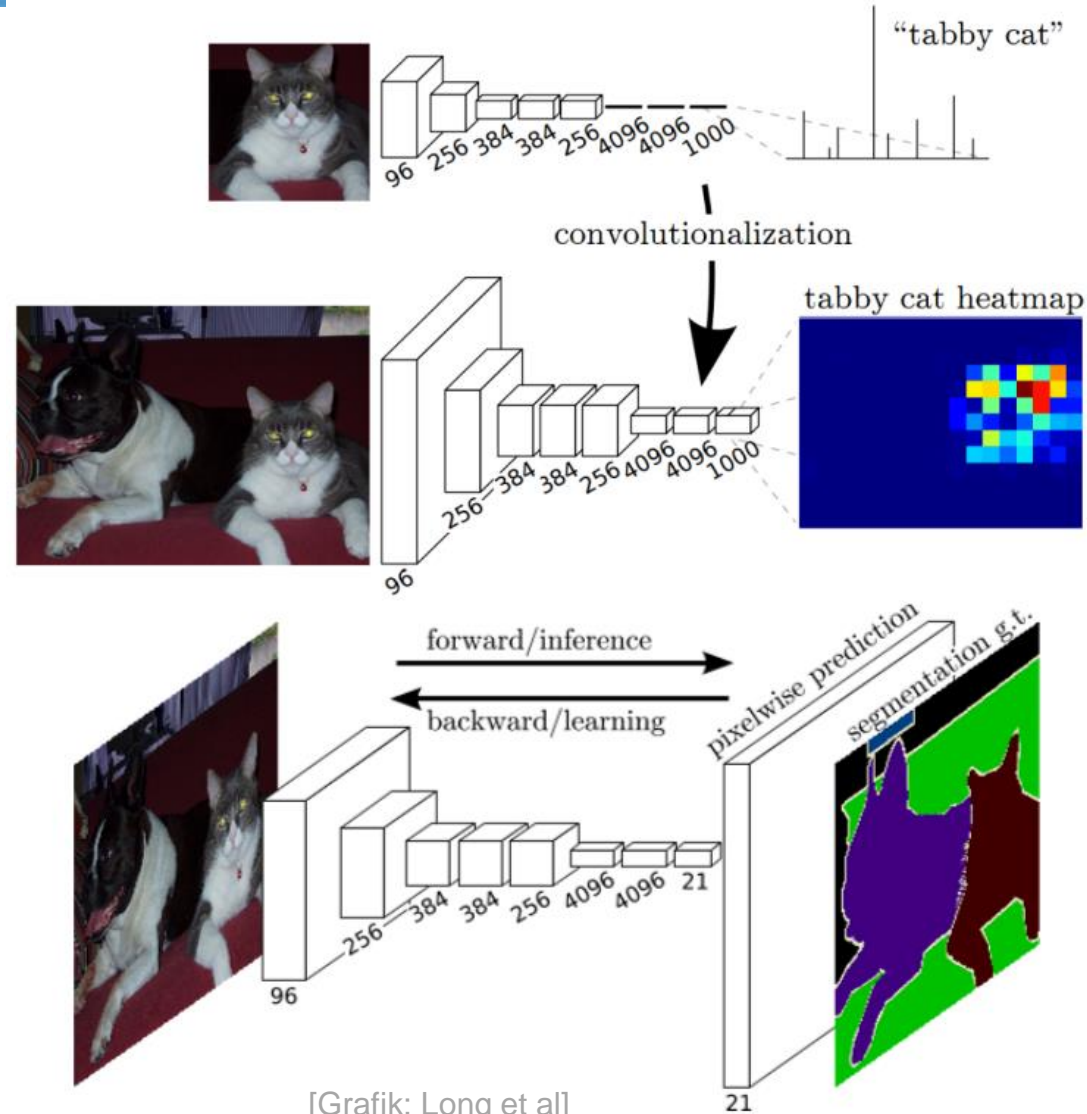


[Grafik: Rohith Gandhi]

- YOLO
  - You Only Look Once: Unified, Real-Time Object Detection
  - Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi [2016]
  - Ca. 25 fps

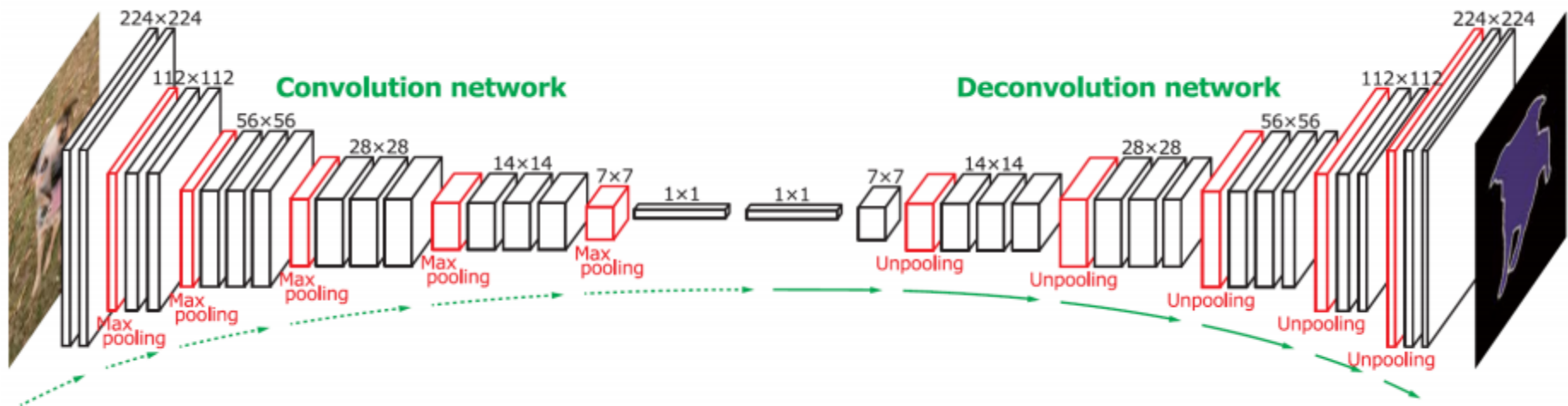


- Segmentierung
- Auch hier kommen Faltungen zum Einsatz
- Fully Convolutional
- Der Output muss wieder ein Bild sein





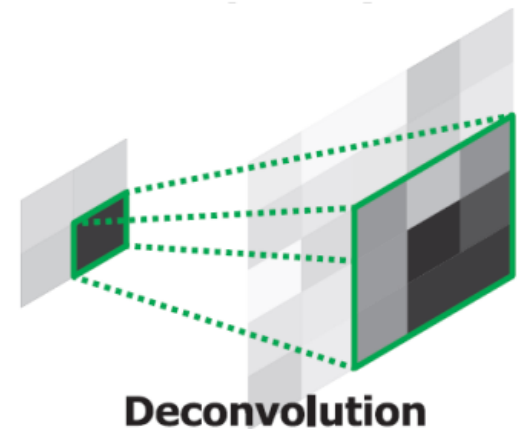
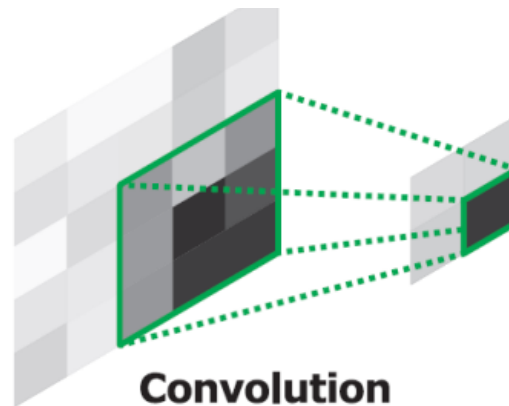
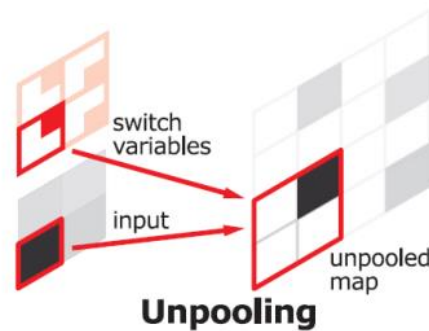
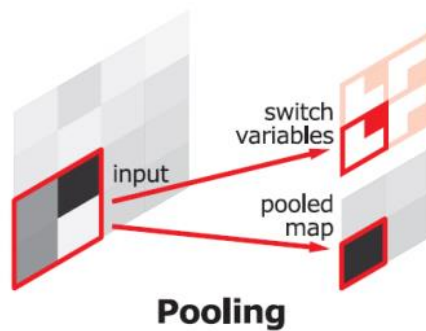
- Deconvolutional Networks







- Deconvolutional Networks
- Unpooling & Deconv



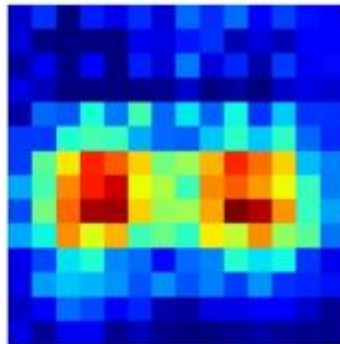
Credit: Noh et al



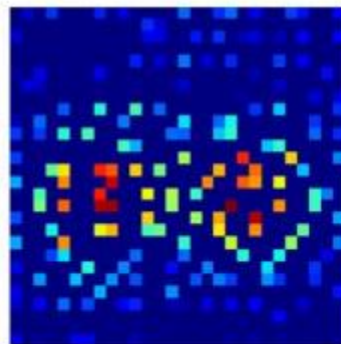
- Deconvolutional Networks
- Unpooling & Deconv



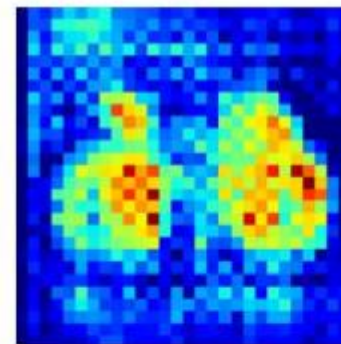
(a)



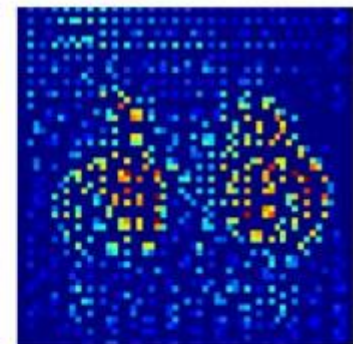
(b)



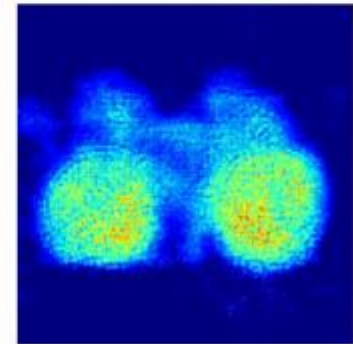
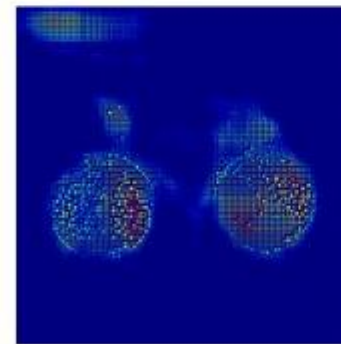
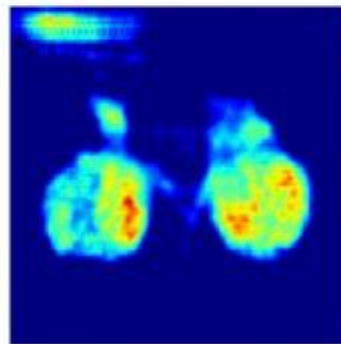
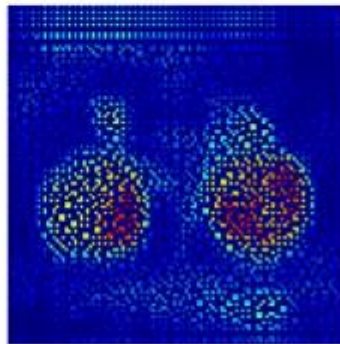
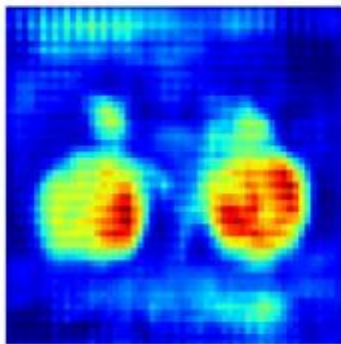
(c)



(d)



(e)

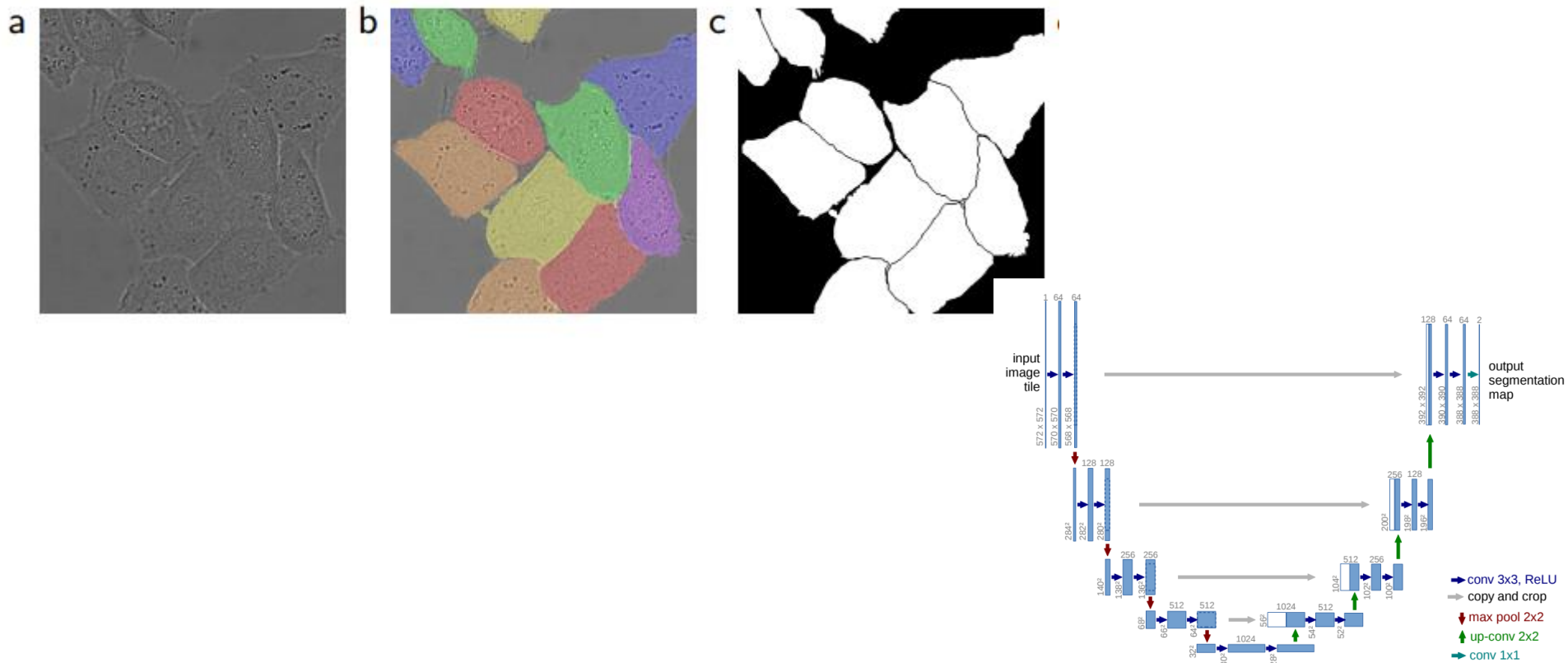


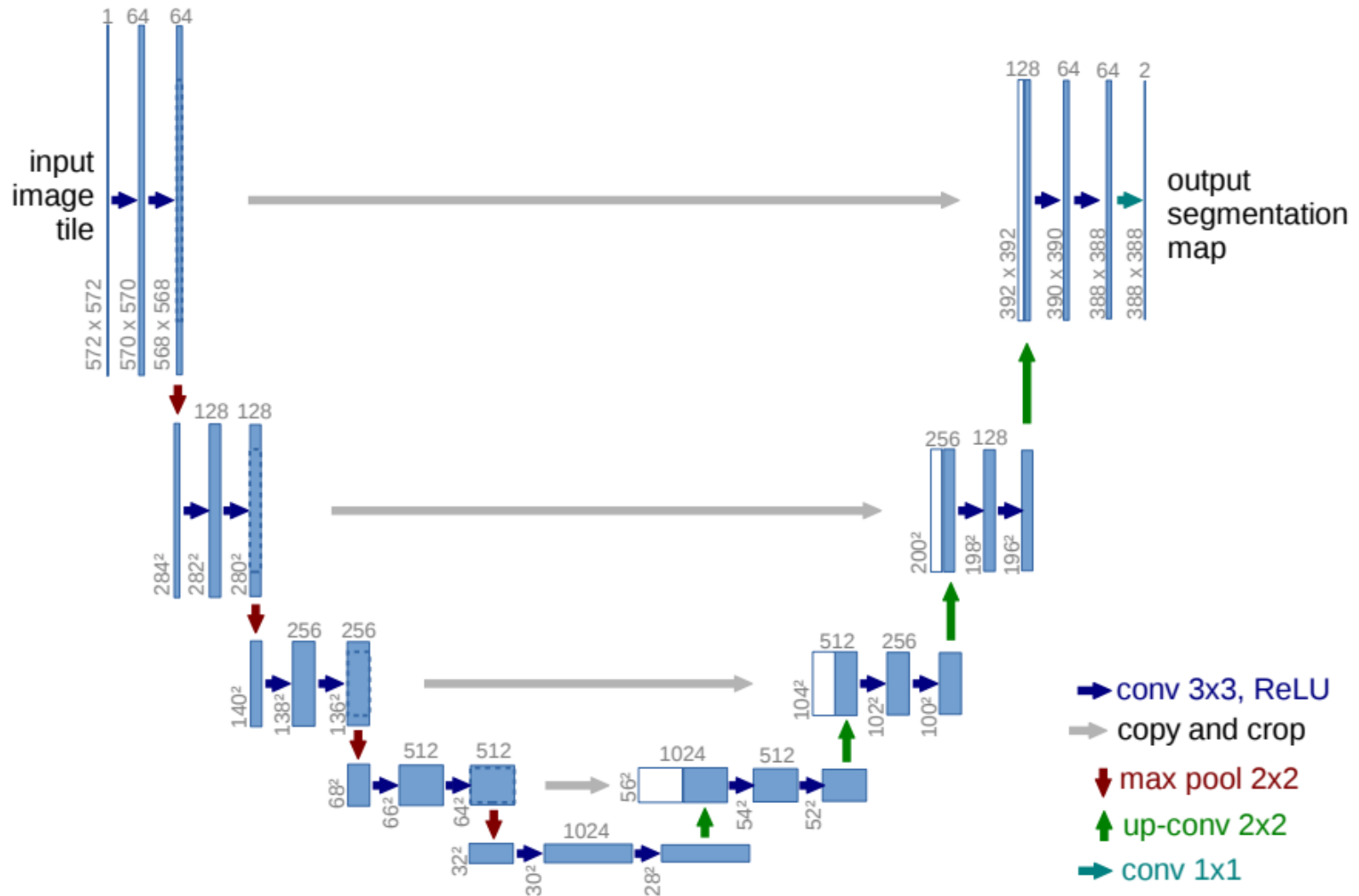
Credit: Noh et al





- U-Net: Convolutional Networks for Biomedical Image Segmentation
  - Olaf Ronneberger, Philipp Fischer, Thomas Brox
  - Uni Freiburg [2015]







- Manchmal reicht es nicht eine passende Architektur zu haben
- ConvNets brauchen viele Daten zum Training
- Oftmals sucht man Lösungen für einen Spezialfall
  - sehr kleiner Datensatz
- Lösung:
  - mit einem anderen (größeren) Datensatz vortrainieren
  - Die vorbereiteten Gewichte auf das eigene Problem anpassen

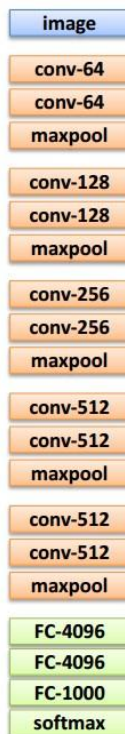


ImageNet data

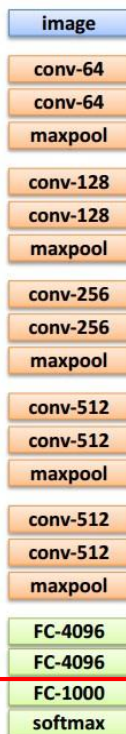


your data

[Grafik: Karpathy]

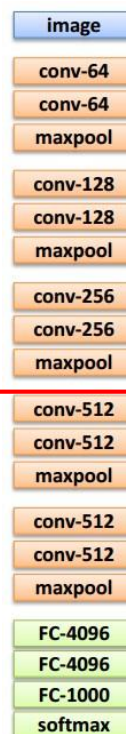


1. Training  
mit  
ImageNet



2. Bei kleinem  
Dataset: fixieren  
aller Weights  
(Das CNN als  
feature extractor),  
Nachtrainieren  
des Klassifikators

FC und Softmax  
layer am Ende  
des Netzes  
austauschen



3. Bei mittelgroßem  
Dataset: “**finetuning**”  
Die Vortrainierten  
Weights als  
Initialisierung  
beibehalten.  
Das ganze Netz oder  
mehrere Layer  
trainieren

Größeren Teil des  
Netzes modifizieren



- Es existieren sehr viele vortrainierte Netzwerke
- Keras bietet diese als „Applications“:

## Available models

### Models for image classification with weights trained on ImageNet:

- Xception
- VGG16
- VGG19
- ResNet50
- InceptionV3
- InceptionResNetV2
- MobileNet
- DenseNet
- NASNet
- MobileNetV2



- Generell ist die Keras Docu sehr gut!
- FAQs:

## What does "sample", "batch", "epoch" mean?

Below are some common definitions that are necessary to know and understand to correctly utilize Keras:

- **Sample:** one element of a dataset.
- *Example:* one image is a **sample** in a convolutional network
- *Example:* one audio file is a **sample** for a speech recognition model
- **Batch:** a set of  $N$  samples. The samples in a **batch** are processed independently, in parallel. If training, a batch results in only one update to the model.
- A **batch** generally approximates the distribution of the input data better than a single input. The larger the batch, the better the approximation; however, it is also true that the batch will take longer to process and will still result in only one update. For inference (evaluate/predict), it is recommended to pick a batch size that is as large as you can afford without going out of memory (since larger batches will usually result in faster evaluating/prediction).
- **Epoch:** an arbitrary cutoff, generally defined as "one pass over the entire dataset", used to separate training into distinct phases, which is useful for logging and periodic evaluation.
- When using `evaluation_data` or `evaluation_split` with the `fit` method of Keras models, evaluation will be run at the end of every **epoch**.
- Within Keras, there is the ability to add **callbacks** specifically designed to be run at the end of an **epoch**. Examples of these are learning rate changes and model checkpointing (saving).





- Paper Referenzen:

## ResNet50

```
keras.applications.resnet50.ResNet50(include_top=True, weights='imagenet', input_tensor=None, input_shape=None,
```

### Arguments

- include\_top: whether to include the fully-connected layer at the top of the network.
- weights: one of `None` (random initialization) or `'imagenet'` (pre-training on ImageNet).
- input\_tensor: optional Keras tensor (i.e. output of `layers.Input()`) to use as image input for the model.
- input\_shape: optional shape tuple, only to be specified if `include_top` is `False` (otherwise the input shape has to be `(224, 224, 3)` (with `'channels_last'` data format) or `(3, 224, 224)` (with `'channels_first'` data format). It should have exactly 3 inputs channels, and width and height should be no smaller than 32. E.g. `(200, 200, 3)` would be one valid value.
- pooling: Optional pooling mode for feature extraction when `include_top` is `False`.
  - `None` means that the output of the model will be the 4D tensor output of the last convolutional layer.
  - `'avg'` means that global average pooling will be applied to the output of the last convolutional layer, and thus the output of the model will be a 2D tensor.
  - `'max'` means that global max pooling will be applied.
- classes: optional number of classes to classify images into, only to be specified if `include_top` is `True`, and if no `weights` argument is specified.

### References

- Deep Residual Learning for Image Recognition



- Wenn keine Application existiert:
  - Paper lesen
  - GitHub
  - Model Zoo



- Wie geht es Weiter?
- Was ist wenn meine Daten keine Bilder sind?
  - Audio und Video?
- Wie Trainiere ich meine Netze richtig?
  - Überwachung des Trainingsvorgangs
  - Hyperparameter
  - Initialisierungen
  - Versuchsplanung