



# Deep Learning

## Chapter 5: Application of CNNs

**Mina Rezaei**

Department of Statistics – LMU Munich  
Winter Semester 2020



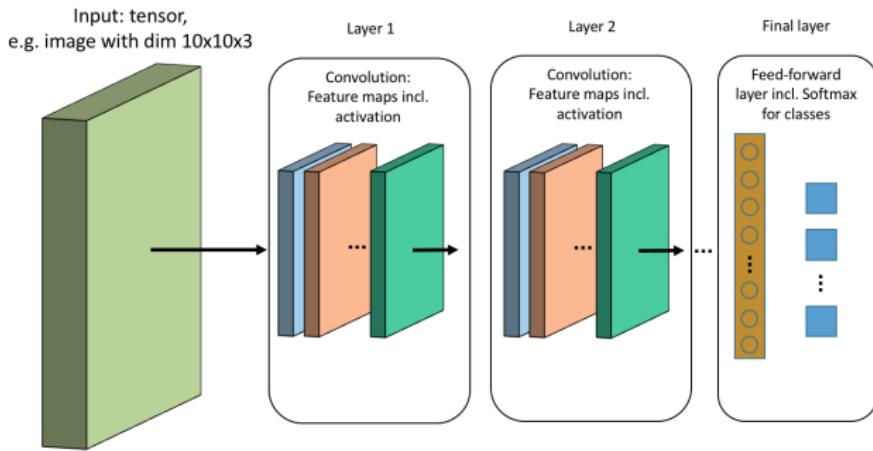
# **LECTURE OUTLINE**

## **Different Perspectives of CNNs**

## **Applications**

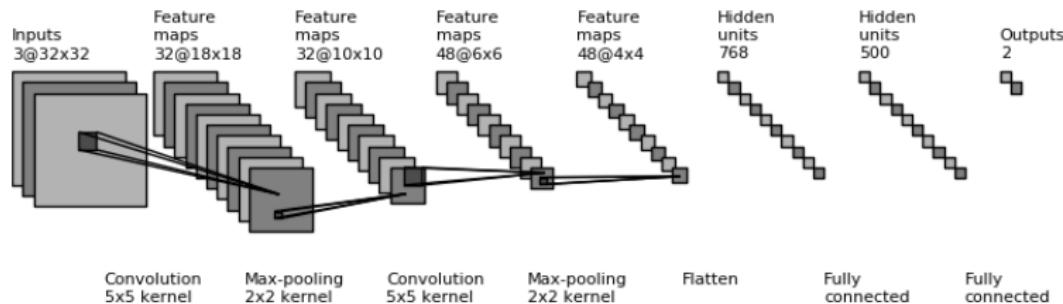
# Different Perspectives of CNNs

# CNNs - PERSPECTIVE I



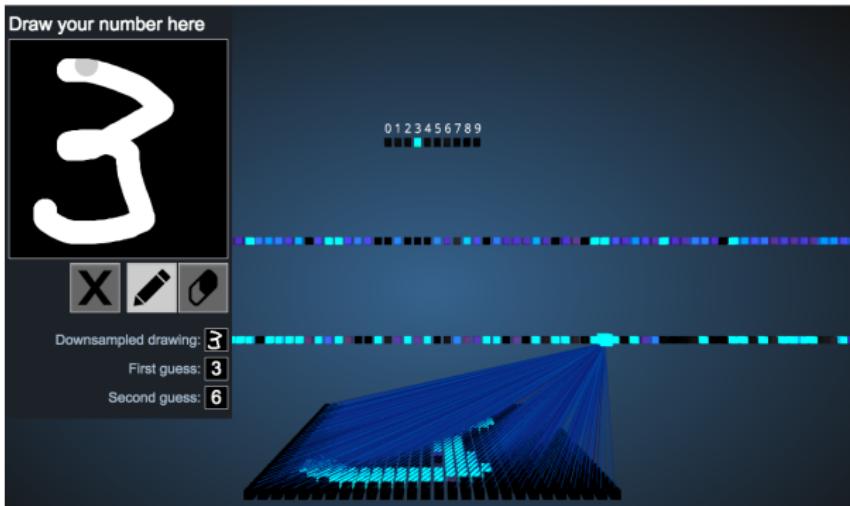
- Schematic architecture of a CNN.
- The input tensor is convolved by different filters yielding different feature maps (coloured) in subsequent layers.
- A dense layer connects the final feature maps with the softmax-activated output neurons.

# CNNs - PERSPECTIVE II



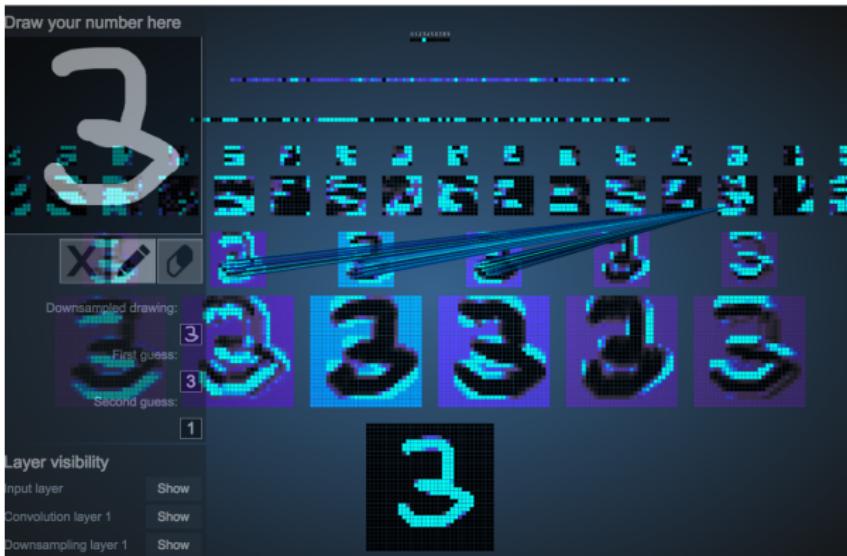
- Flat view of a CNN architecture for a classification problem.
- Consists of 2 CNN layers that are each followed by max-pooling, then flattened and connected with the final output neurons via a dense layer.

# CNNs - PERSPECTIVE III



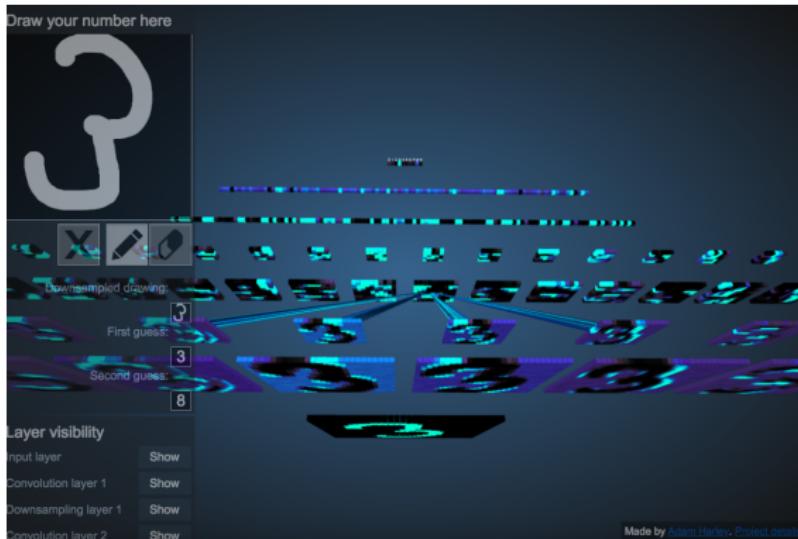
- Awesome interactive visualization (by Adam Harley) [▶ Click here](#).
- Vanilla 2-layer densely-connected net on MNIST data for input digit 3.
- Each neuron in layer 1 is connected to each of the input neurons.

# CNNs - PERSPECTIVE III



- Front view on 2-layer CNN with Pooling and final dense layer on MNIST data for input digit 3.
- Each neuron in the second CNN layer is connected to a patch of neurons from each of the previous feature maps via the convolutional kernel.

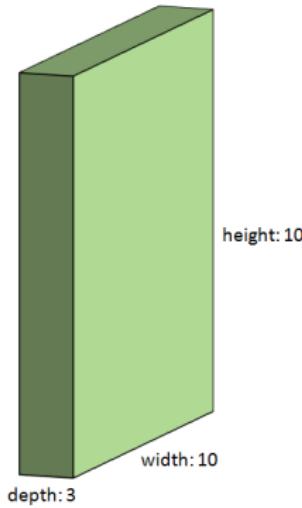
# CNNs - PERSPECTIVE III



- Bottom view on 2-layer CNN with Pooling and final dense layer on MNIST data for input digit 3.
- Each neuron in the second CNN layer is connected to a patch of neurons from each of the previous feature maps via the convolutional kernel.

# CNNs - ARCHITECTURE

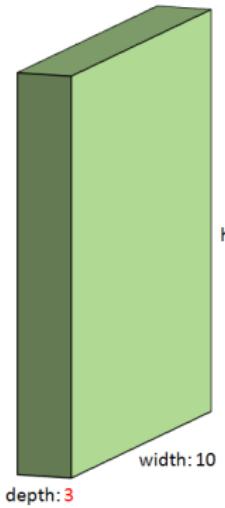
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



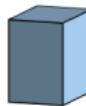
- Suppose we have the following input tensor with dimensions  $10 \times 10 \times 3$ .

# CNNs - ARCHITECTURE

Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



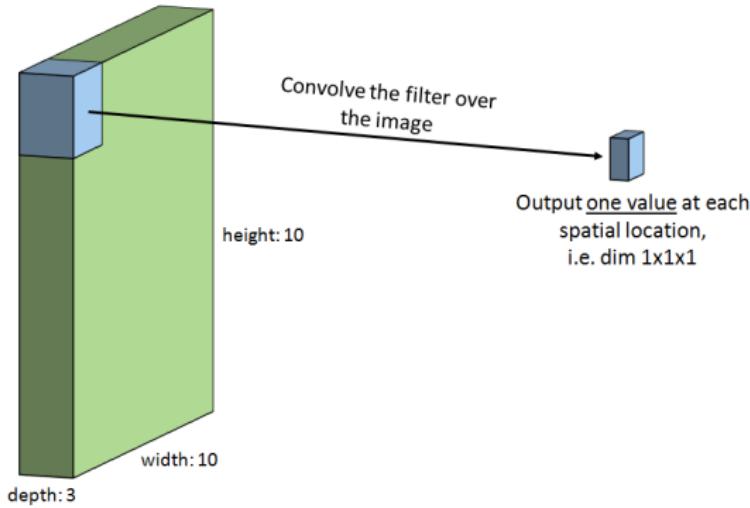
Filter/Kernel  
e.g. with dim  $2 \times 2 \times 3$



- We use a filter of size 2.

# CNNs - ARCHITECTURE

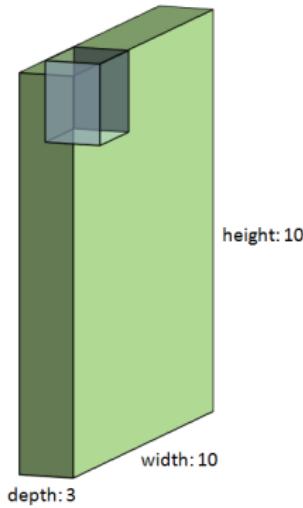
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



- Applying it to the first spatial location, yields one scalar value.

# CNNs - ARCHITECTURE

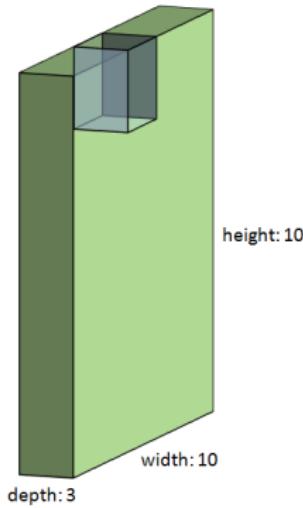
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



- The second spatial location yields another one..

# CNNs - ARCHITECTURE

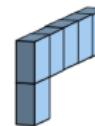
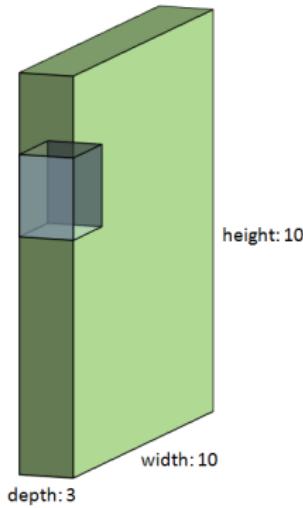
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



- ...and another one...

# CNNs - ARCHITECTURE

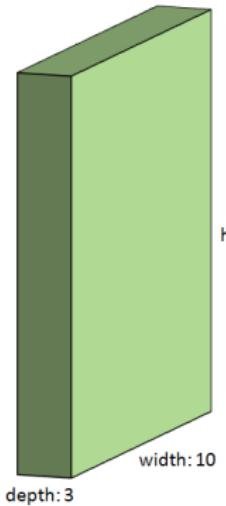
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



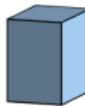
- ...and another one...

# CNNs - ARCHITECTURE

Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



Filter with  
dim  $2 \times 2 \times 3$



Output: feature map,  
here with dim  $5 \times 5 \times 1$



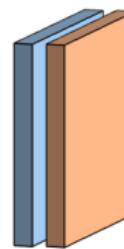
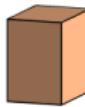
- Finally we obtain an output which is called feature map.

# CNNs - ARCHITECTURE

Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



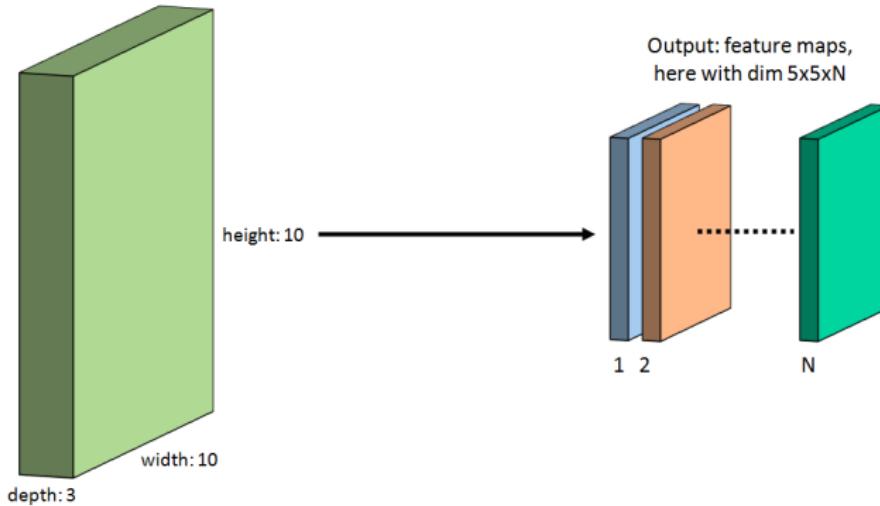
Filter with  
dim  $2 \times 2 \times 3$



- We initialize another filter to obtain a second feature map.

# CNNs - ARCHITECTURE

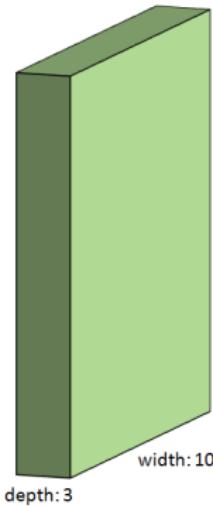
Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



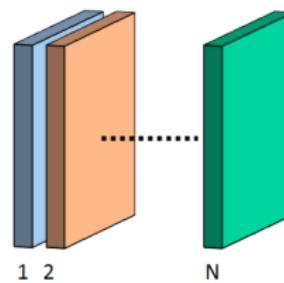
- All feature maps yield us a “new image” with dim  $h \times w \times N$ .

# CNNs - ARCHITECTURE

Input: tensor,  
e.g. image with dim  $10 \times 10 \times 3$



Output: feature maps,  
here with dim  $5 \times 5 \times N$

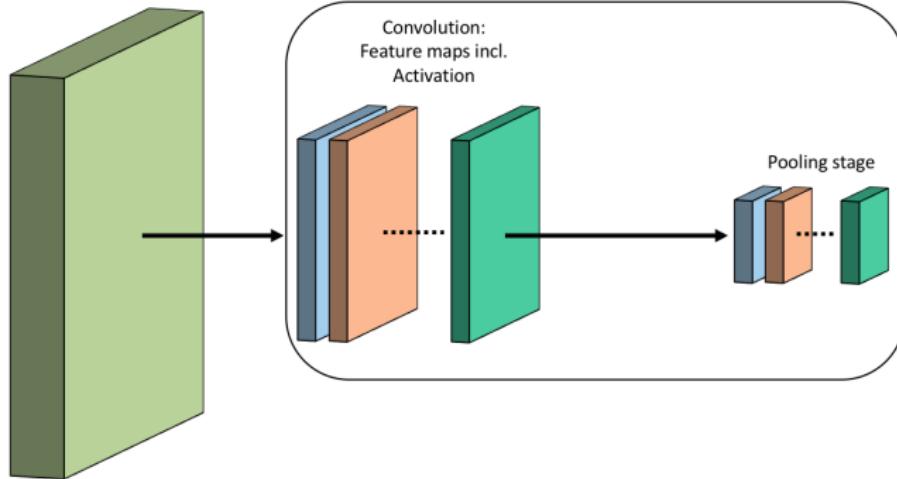


- We actually append them to a new tensor with depth = # filters.

# CNNs - ARCHITECTURE

Input: tensor,  
e.g. image with dim 10x10x3

One Layer

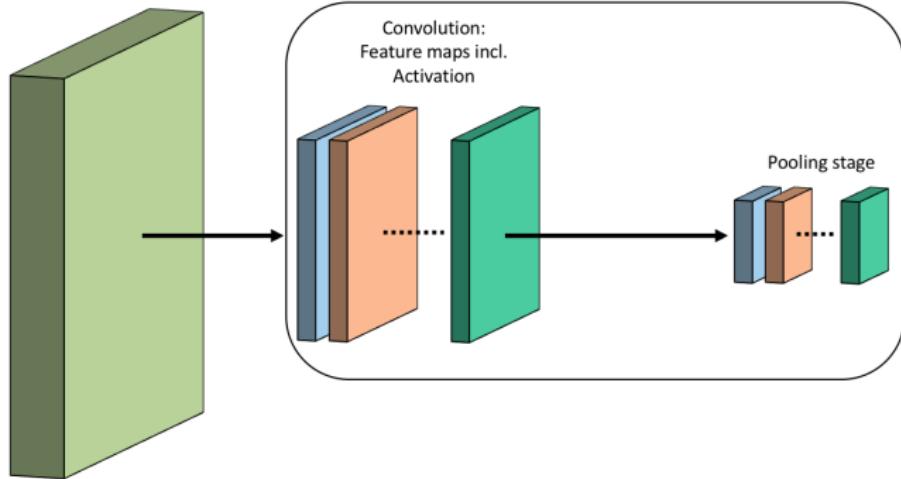


- All feature map entries will then be activated (e.g. via ReLU), just like the neurons of a standard feedforward net.

# CNNs - ARCHITECTURE

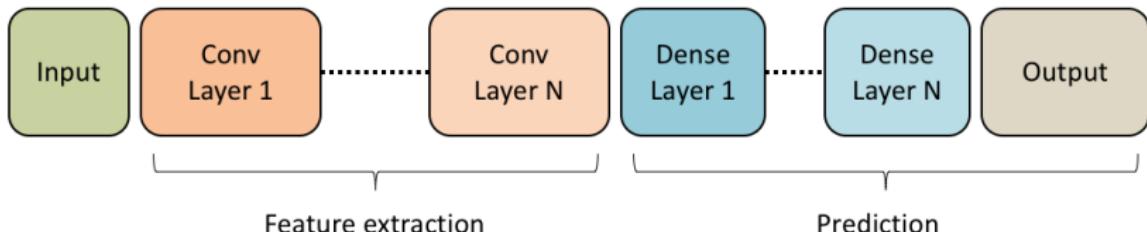
Input: tensor,  
e.g. image with dim 10x10x3

One Layer



- One may use pooling operations to downsample the dimensions of the feature maps.
- Pooling is applied on each feature map independently: the latter, blue block is the pooled version of the previous, blue feature map.

# CNNs - ARCHITECTURE



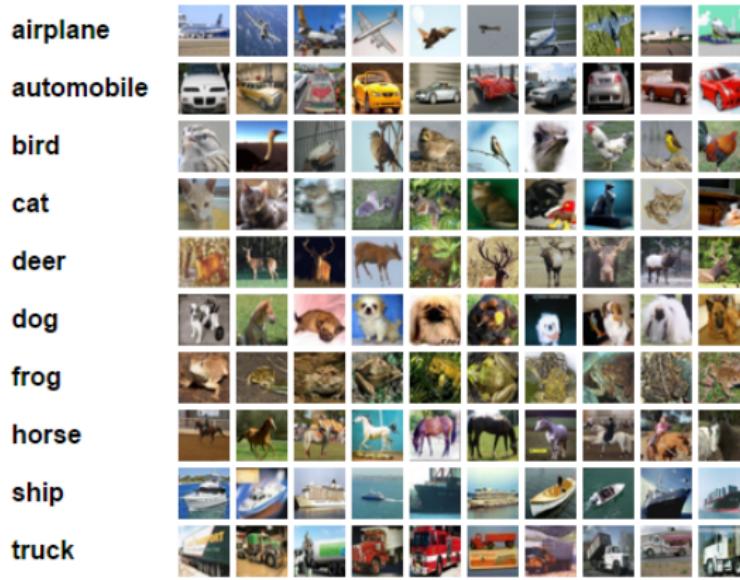
- Many of these layers can be placed successively, to extract even more complex features.
- The feature maps are fed into each other sequentially. For instance, each filter from the second conv layer gets all previous feature maps from the first conv layer as an input. Each filter from the first layer extracts information from the input image tensor.
- The feature maps of the final conv layer are flattened (into a vector) and fed into a dense layer which, in turn, is followed by more dense layers and finally, the output layer.

# **Applications**

# APPLICATION - IMAGE CLASSIFICATION

- One of use case for CNNs is image classification.
- There exist a broad variety of battle-proven image classification architecture such as the AlexNet , the Inception Net or the ResNet which will be discussed in detail in the next lecture.
- All those architectures rely on a set of subsequent convolutional filters and aim to learn the mapping from an image to a probability score over a set of classes.

# APPLICATION - IMAGE CLASSIFICATION



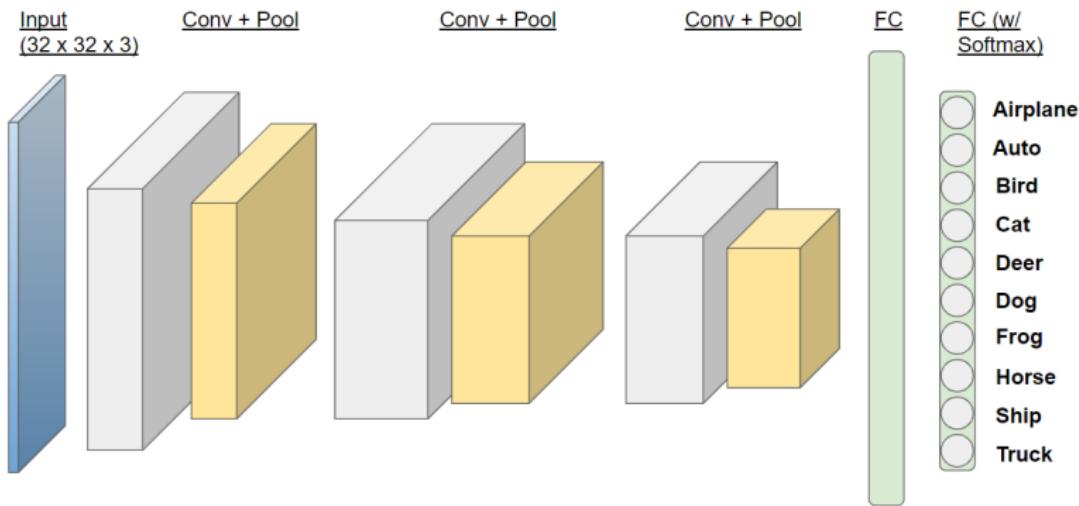
**Figure:** Image classification with Cifar 10: famous benchmark dataset with 60000 images and 10 classes (Alex Krizhevsky (2009)). There is also a much more difficult version with 60000 images and 100 classes.

# APPLICATION - IMAGE CLASSIFICATION



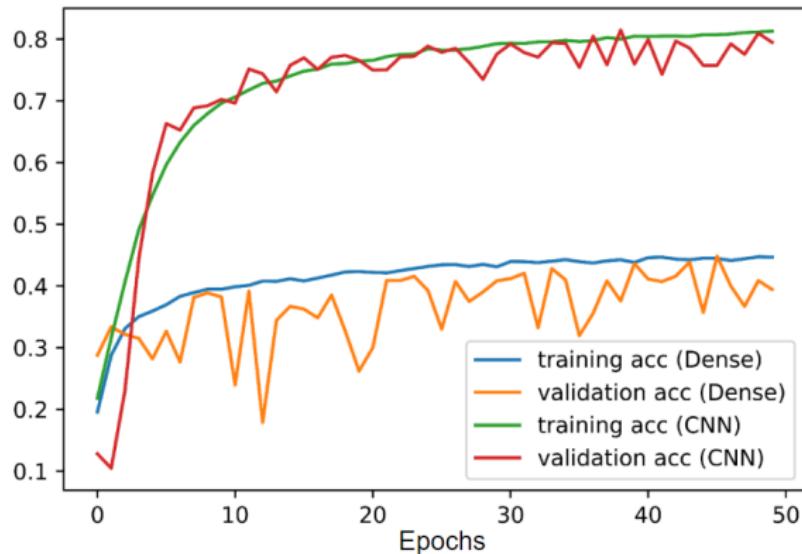
**Figure:** One example of the Cifar10 data: A highly pixelated, coloured image of a frog with dimension [32, 32, 3].

# APPLICATION - IMAGE CLASSIFICATION



**Figure:** An example of a CNN architecture for classification on the Cifar10 dataset (FC = Fully Connected).

# CNN VS A FULLY CONNECTED NET ON CIFAR10



**Figure:** Performance of a CNN and a fully connected neural net ("Dense") on Cifar-10. Both networks have roughly the same number of layers. They were trained using the same learning rate, weight decay and dropout rate. Of course, the CNN performs better because it has less learnable parameters and the right inductive bias for the task.

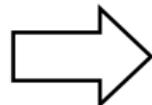
# APPLICATION - IMAGE COLORIZATION

- Basic idea (introduced by Zhang et al., 2016):
  - Train the net on pairs of grayscale and coloured images.
  - Force it to make a prediction on the colour-value **for each pixel** in the grayscale input image.
  - Combine the grayscale-input with the colour-output to yield a colorized image.
- Very comprehensive material on the method is provided on the author's website. [▶ Click here](#)

# APPLICATION - IMAGE COLORIZATION



**Input:** Grayscale image  
 $L$  channel



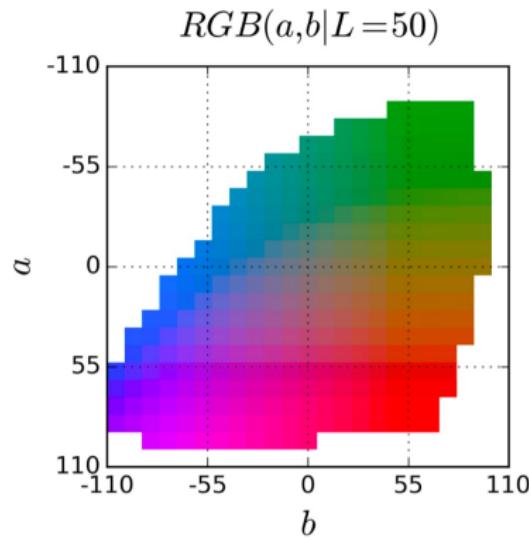
**Output:** Color information  
 $ab$  channels



Concatenate ( $L, ab$ )  
for plausible colorization

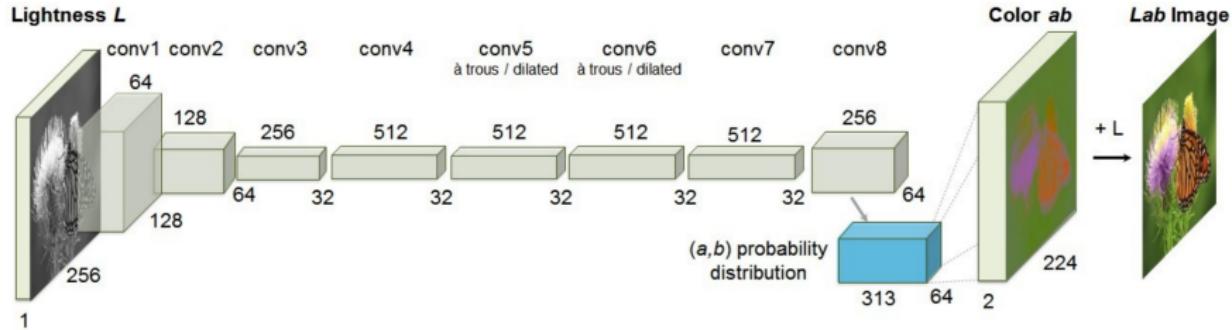
**Figure:** The CNN learns the mapping from grayscale ( $L$ ) to color ( $ab$ ) for each pixel in the image. The  $L$  and  $ab$  maps are then concatenated to yield the colorized image. The authors use the LAB color space for the image representation.

# APPLICATION - IMAGE COLORIZATION



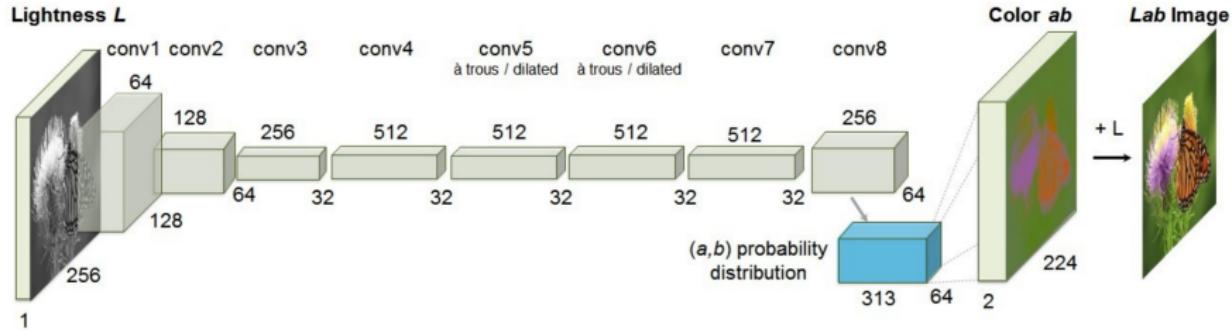
**Figure:** The colour space (ab) is quantized in a total of 313 bins. This allows to treat the color prediction as a classification problem where each pixel is assigned a probability distribution over the 313 bins and that with the highest softmax-score is taken as predicted color value. The bin is then mapped back to the corresponding, numeric (a,b) values. The network is optimized using a multinomial cross entropy loss over the 313 quantized (a,b) bins.

# APPLICATION - IMAGE COLORIZATION



**Figure:** The architecture consists of stacked CNN layers which are upsampled towards the end of the net. It makes use of *dilated convolutions* and *upsampling layers* which are explained in the next lecture. The output is a tensor of dimension [64, 64, 313] that stores the 313 probabilities for each element of the final, downsampled 64x64 feature maps.

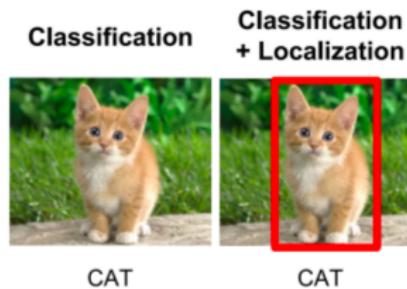
# APPLICATION - IMAGE COLORIZATION



**Figure:** This block is then upsampled to a dimension of 224x224 and the predicted color bins are mapped back to the (a,b) values yielding a depth of 2. Finally, the L and the ab maps are concatenated to yield a colored image.

# APPLICATION - OBJECT LOCALIZATION

- Until now, we used CNNs for *single-class* classification of images - **which object is on the image?**
- Now we extend this framework - **is there an object in the image and if yes, where and which?**



**Figure:** Classify and detect the location of the cat.

# APPLICATION - OBJECT LOCALIZATION

- Bounding boxes can be defined by the location of the left lower corner as well as the height and width of the box:  $[b_x, b_y, b_h, b_w]$ .
- We now combine three tasks (detection, classification and localization) in one architecture.
- This can be done by adjusting the label output of the net.
- Imagine a task with three classes (cat, car, frog).
- In standard classification we would have:

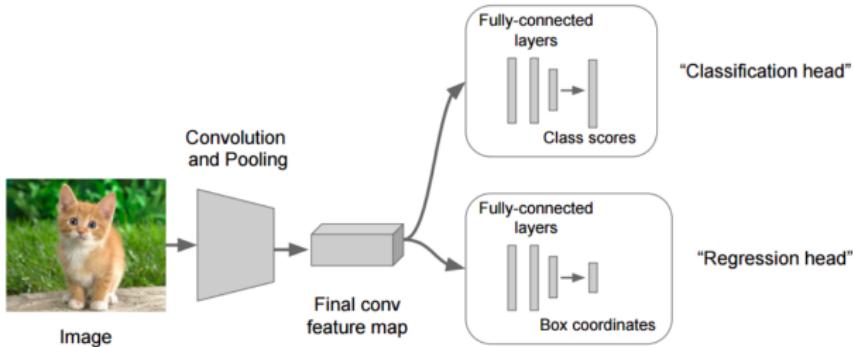
label vector  $\begin{bmatrix} c_{cat} \\ c_{car} \\ c_{frog} \end{bmatrix}$  and softmax output  $\begin{bmatrix} P(y = cat|X, \theta) \\ P(y = car|X, \theta) \\ P(y = frog|X, \theta) \end{bmatrix}$

# APPLICATION - OBJECT LOCALIZATION

- We include the information, if there is an object as well as the bounding box parametrization in the label vector.
- This gives us the following label vector:

$$\begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ c_o \\ c_{cat} \\ c_{car} \\ c_{frog} \end{bmatrix} = \begin{bmatrix} \text{x coordinate box} \\ \text{y coordinate box} \\ \text{height box} \\ \text{width box} \\ \text{presence of object, binary} \\ \text{class cat, one-hot} \\ \text{class car, one-hot} \\ \text{class frog, one-hot} \end{bmatrix}$$

# APPLICATION - OBJECT LOCALIZATION



- Naive approach: use a CNN with two heads, one for the class classification and one for the bounding box regression.
- But: What happens, if there are two cats in the image?
- Different approaches: "Region-based" CNNs (R-CNN, Fast R-CNN and Faster R-CNN) and "single-shot" CNNs (SSD and YOLO).

# REFERENCES

-  Ian Goodfellow, Yoshua Bengio and Aaron Courville (2016)  
Deep Learning  
*http://www.deeplearningbook.org/*
-  Alex Krizhevsky (2009)  
Learning Multiple Layers of Features from Tiny Images  
*https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf*
-  Adam W Harley (2015)  
An Interactive Node-Link Visualization of Convolutional Neural Networks  
*http://scs.ryerson.ca/~aharley/vis/*