

Universität Hamburg  
Department Informatik  
Knowledge Technology, WTM

# Deep Learning: Neural Networks for Object Detection and Tracking Tasks

Seminar Paper

Brain Modelling

Daniel Speck

Matr.Nr. 632 13 17

[2speck@informatik.uni-hamburg.de](mailto:2speck@informatik.uni-hamburg.de)

29.05.2015



## **Abstract**

Deep learning and the interrelated deep neural networks are one of the most successful learning strategies at the moment as the computing power for creating such structures rose in the past years via GPU computing. Image and video classification, object detection and tracking tasks can be fulfilled with these architectures which perform far better than more classical ones. One focus will be convolutional neural networks and the latest research on (mostly supervised) deep learning architectures.

## **Contents**

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Background information: Artificial neural networks</b>	<b>2</b>
<b>3</b>	<b>Deep neural networks</b>	<b>3</b>
<b>4</b>	<b>Convolutional neural networks and image processing</b>	<b>4</b>
<b>5</b>	<b>Research and field of application</b>	<b>6</b>
5.1	Image classification . . . . .	6
5.2	Video classification . . . . .	8
5.3	Object tracking . . . . .	8
<b>6</b>	<b>Conclusion</b>	<b>9</b>
	<b>Bibliography</b>	<b>10</b>

# 1 Introduction

Big data is a popular keyword in the last couple of years. With more and more devices being able to capture high-resolution images and videos, the data and information available on the internet increases rapidly, rather exponentially, but all this data is useless if this flood of information can not be searched, categorized and made easily accessible for humans. As a consequence of this situation classic, mostly static algorithms are “out-of-date” and new, intelligent systems have to be developed. Biological entities are, up to now, by far better when it comes to recognize, classify, sort etc. images or videos compared to computers. This circumstance led to researching intelligent systems that are able to adapt to their input, search for patterns and recognize, track and classify them.

Inspired by biological systems, deep learning and especially convolutional neural networks have become popular because their results in classification and tracking tasks are state-of-the-art. In this paper a basic information view about artificial neural networks will be given for understanding the underlying concepts which are important to solve complex problems with machine learning architectures. In addition to that a comprehension of deep neural networks and deep learning will be made to point out the abilities but also the limits of “classic” neural networks. Also some information about the role of GPU-computing in such disciplines will be provided. For solving complex problems in classification, recognition, tracking etc. domains, convolutional neural networks have shown good performance in combination with GPU-computing, so they will be focused on. Especially the latest research on these topics and fields of application will be viewed in detail. An important contest for classification tasks is measuring the performance of such architectures with the ImageNet LSVRC databases and for videos with the YouTube 1M database by Google. Both topics will be introduced. One interesting fact is that pre-trained networks on videos showing sports did better on non-sports related video categorization than “fresh” networks. This is indicating an abstract, deep correlation between different domains in recognition tasks.

# 2 Background information: Artificial neural networks

Artificial neural networks are intended to approximate certain functions for machine learning purposes. They are intended to model learning processes. Fix/static algorithms can be calculated fast by modern computers but fail most times at disciplines requiring intelligent behavior. A classic example is recognizing handwritten digits because the shape, color and contrast highly vary in dependence of the used pencil and, of course, of the writing style. As a result, classic, static algorithms would fail to identify those varying patterns and therefore fail to appropriately classify handwritten digits. Such tasks can be fulfilled by artificial neural networks. A common architecture is the multilayer-perceptron as shown in figure 1.

This artificial neural network, for example, could be used to classify hand-written digits. A database of hundreds of images would be needed for training the network accurately, using just a few training samples normally causes overfitting. The images information of input size 5x5, which results in 25 pixels, would be encoded by the 25 neurons in the input layer. The hidden layers would try to extract the features which are characteristic for each digit and the output layer with 10 neurons models the digits 0 to 9 for representing the calculated digit. Therefore each neuron of the output layer would give the probability for each digit (0-9) that it was this digit showing up in the input image. However, such architectures are insufficient for more complex tasks because they need a vast amount of resources, are not optimized and do not make use of any prior-knowledge of the problem. Non-deep structures where every layer is fully-connected turned out to be expensive in terms of computational resources and not optimal considering the results in disciplines of image and video classification, object recognition and tracking.

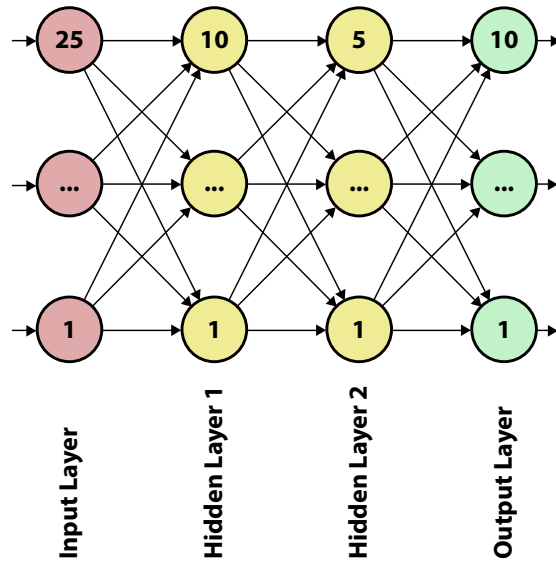


Figure 1: An example of a classic, fully-connected multilayer-perceptron. The input layer holds 20 neurons, hidden layer 1 has 10 neurons, hidden layer 2 has 5 neurons and the output layer 2 neurons. All arrows but the short ones represent randomly initialized weights.

### 3 Deep neural networks

Classic, small neural networks cannot succeed on complex tasks such as object detection in high-resolution, real world images. The variety of features, blurry backgrounds, one or several objects in an image and other factors render detection, tracking and classification tasks challenging. The ability to detect, track and distinguish real world objects demands complex structures with a much larger capability of solving complex problems. As computing power rises and GPU-computing became popular, neural networks are no longer limited to a few layers containing just some neurons. Modern solutions (in the year 2014) can carry tens of layers with thousands of neurons and millions of connections [5]. Complex structures like this enable computing complex problems which arise with tasks like image classification in real world images. To fulfill those assignments deep neural networks are used in combination with GPU-optimized code which is considerably faster than CPU-optimized code and hence allows larger, more accurate architectures [1].

## 4 Convolutional neural networks and image processing

For deep learning purposes (classic/fully-connected) multilayer perceptrons consume a sizable amount of resources for proper training when they are designed to solve complex tasks because the amount of neurons and especially weights increases rapidly with the network's size. For example, a MLP like that showed in figure 1 with four layers, an input layer with 25 neurons, two hidden layers with 10 and 5 neurons and an output layer with 10 neurons for classifying images with a size of 5x5 pixels into 10 different classes would have  $25 * 10 + 10 * 5 + 5 * 10 = 350$  weights/connections. Training this net would already consume some time and space complexity. Compared to real world images such 5x5 images are tiny, even scaling this example up to images with a dimension of 50x50 is far less than real world images but would already result in an architecture like 2500 (input), 1000 (hidden), 500 (hidden), 10 (output) neurons<sup>1</sup>. This architecture would have  $2500 * 1000 + 1000 * 500 + 500 * 10 = 3,005,000$  weights/connections. The human brain holds about 150 trillion ( $1.5 * 10^{14}$ ) synapses [9]. Moreover, as features in images capturing real world scenes are distributed in certain patterns (they cover spatially local correlation, such as shapes), there is no need to have every pixel's information being processed by one neuron. Actually in most cases results would be even better, if the pixel's information is pre-processed, for instance by edge detection filters. However a fully connected layer of neurons is not an optimal solution for this task.

Convolutional neural networks (CNNs) are inspired by biology. Instead of connecting every pixels information directly with a neuron to process its information it filters the information in the first layers [7]. This procedure is similar to the on processes happening when an biological eye receives stimuli. The receptive field<sup>2</sup> has a vast amount of photoreceptor cells<sup>3</sup> gathering information and converging the received information on to distinctly fewer retinal ganglion cells<sup>4</sup>. This process maps several features and reduces the input dimensionality as well as distinguishes the information to separate "channels" which are then transfered to the corresponding neurons to process features such as color, motion, shapes and so on separately [8]. The idea of CNNs is based on this biological process: the information of an input image is convolved by several filters which try to extract interesting features in the first layer and in following layers this information is pooled and subsampled [7]. Convolution itself is the repeated application of a function of the output of another function. In the context of CNNs it is applying different "filters" over an image to extract the already mentioned features. A convolution layer extracts the information of a pixel out of an image with kernels<sup>5</sup> [7].

---

<sup>1</sup>Assuming, for example, the task is classifying real hand-written digits

<sup>2</sup>[http://en.wikipedia.org/wiki/Receptive\\_field](http://en.wikipedia.org/wiki/Receptive_field)

<sup>3</sup>[http://en.wikipedia.org/wiki/Photoreceptor\\_cell](http://en.wikipedia.org/wiki/Photoreceptor_cell)

<sup>4</sup>[http://en.wikipedia.org/wiki/Retinal\\_ganglion\\_cell](http://en.wikipedia.org/wiki/Retinal_ganglion_cell)

<sup>5</sup>[http://en.wikipedia.org/wiki/Kernel\\_\(image\\_processing\)](http://en.wikipedia.org/wiki/Kernel_(image_processing))



Figure 2: Left side: original image (greyscaled), right side: edge-detection kernel processed image. Original image by Michael Plotke, 28th of January, 2013. Open creative commons license. <http://upload.wikimedia.org/wikipedia/commons/5/50/Vd-Orig.png> and <http://upload.wikimedia.org/wikipedia/commons/6/6d/Vd-Edge3.png>

Example for one possible convolution layer: figure 2 shows an image of an animal on the left side (original image, greyscaled) and a kernel-processed one on the right side. The used kernel matrix for filtering the left image is shown in equation (1). So basically each pixel's information in the right, processed image is the result of applying the kernel matrix  $K_M$  (1) on the same pixel (and the neighboring pixels) in the left image. With a wide variety of different kernels several different features can be extracted from an image. Typical CNNs use tens to hundreds of different kernel filters gathering as many features from an image as possible. The resulting pixel processed by a kernel filter is calculated via the formula:

$$K_M = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \quad (1)$$

$$I_{out}(x, y) = \left| \sum_{a=1}^3 \sum_{b=1}^3 I_{in}(x + a - c_x, y + b - c_y) * K_M(a, b) \right| \quad (2)$$

Where  $c_x$  is the coordinate of the x-center and  $c_y$  the coordinate of the y-center of the input image. After processing this filter to the whole image, edges would be highlighted and the rest becomes nearly black, like in figure 2 (right side), so that shape features are extracted out of the original image.

$$I_{in} = \begin{bmatrix} 46 & 42 & 50 \\ 44 & 65 & 56 \\ 41 & 52 & 58 \end{bmatrix} \quad (3)$$

$$I_{out}(2, 2) = 131$$

Equation (3) shows an example input and output of the image of figure 2.  $I_{in}$  is the input with a dimension of 3x3 pixels extracted out of the original image.  $I_{out}$  is the output for the center pixel ( $x = 2, y = 2$ ) of the kernel-processed output image using the kernel filter  $K_M$  of equation (1).

Another idea of CNNs are subsampling layers, which follow convolutional layers. Subsampling reduces the overall amount of information and therefore not only saves resources for the later classification tasks but also strengthens the detected features in an image. Often max-pooling is used as a subsampling strategy in CNNs [11, 7] which determines the most distinctive pixel in a given area. A max-pooling algorithm splits the input in grids and selects the maximum value out of each grid, thus non-maximal values are deleted so that only that information continues to exist which represents the current feature best [10, 7]. Before applying max-pooling the input represents the presence of a feature in one or only some pixels. After this kind of dimensionality reduction the assertion is enlarged to a greater area, corresponding to several pixels of the original image. Szegedy et al. have made heavy use of this approach in “GoogLeNet” [10]. Additionally, this technique provides robustness to the position of a feature in an image, as the position becomes less important. Thereupon the extracted features are representing invariances, which are needed for classifying the collected information regardless of its position, rotation or scaling [6]. The last layers of a CNN regularly consist of fully-connected layers. In comparison to feature extraction, strengthening and dimensionality reduction, those layers are supposed to take all gathered features/information and make predictions [11, 2, 1]. E.g. if the supplied input contains big shapes of ears, legs etc., filtered trunks, grey-ish color maps and so on, the fully-connected layers have to collect and mix all those gathered features/information and recognize it for finally classifying it as an image of elephants.

## 5 Research and field of application

### 5.1 Image classification

For measuring the performance of a neural network top-1 and top-5 error rates have become popular. While letting a neural network classify some input with test data the probabilities for each class are recorded. The error rates are arithmetic means of all test samples:

$$\text{top-1 error rate} = \sum_{S_i} x \begin{cases} 1, & \text{if guessed class matches actual class of } S_i \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

$$\text{top-5 error rate} = \sum_{S_i} x \begin{cases} 1, & \text{if the actual class of } S_i \text{ is one of the top 5} \\ & \text{probabilities guessed by the neural network} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

$S_i$  is the current test sample while  $S$  contains all test samples. Alex Krizhevsky et al. have shown in “ImageNet Classification with Deep Convolutional Neural Networks” [7] the good performance of CNNs in image classification tasks. Their



architecture achieved top-1 and top-5 error rates of 37.5% and 17.0% on the ImageNet LSVRC-2010 contest which was the best result in 2012. The CNN consists of about 650,000 neurons with 60 million parameters and nine layers total. The only preprocessing made was resizing the images to a resolution of 256x256 because the images in ImageNet vary in resolution. Out of those images 224x224x3 dimensional patches are extracted. 224x224 in horizontal/vertical dimension and 3 different colors (RGB images). The first layer is a convolution layer which filters the input of 224x224x3 to a kernel map of 96 kernels with the size of 11x11x3 each. This slightly reduces the dimensionality of the input image and extracts 96 different features out of the input.

The second layer is a max-pooling layer which filters the information via 256 5x5x48 kernels. Therefore the dimensionality is further reduced and already extracted features strengthened.

Afterwards convolutional and max-pooling layers follow the the fifth layer until three fully-connected layers with 4096 neurons each try to collect and mix the information and then classify the images.

One key problem with deep learning structures and especially CNNs is overfitting. Because of the size of those networks a huge amount of training data is needed. An interesting method used by Krizhevsky is to design first layer (224x224x3 for Krizhevsky's architecture) slightly smaller than input data (256x256x3) and then not just apply kernel filters on those input data but first extract random 224x224x3 patches out of the original data. Additionally they the horizontal reflections of those patches and altered the intensities of the RGB channels.

This technique can boost the training data set by hundreds or even thousands of samples [7].

Another procedure to reduce overfitting and test errors is called "dropout" [4]. The idea is to set the output of hidden neurons to zero with a certain probability, mostly 0.5 [7]. This prevents neurons of co-adapting each other which means that neurons can no longer distort its output because they "familiarize" with some neurons of their parent layer and ignore others. With deactivating random neurons in each learning step the presence of an other neuron is not guaranteed so neurons are more likely to learn and focus on robust features rather than on single attractive neurons.

This has shown especially effective when applied in fully-connected layers [7].

Normally the transfer function for a neurons output is modeled a hyperbolic tangent or sigmoid function but those functions can increase the training time [7]. Using Rectified Linear Units (ReLUs) [3] with a function like  $f(x) = \max(0, x)$  not only makes CNNs train a lot faster but also slightly reduces the training error, because low training error rates are reached faster (in less epochs) the whole network concentrates on robust features faster and this has an additional effect of convergence.

In Krizhevsky's architecture the 25% error rate on the CIFAR-10 training dataset is reached six times faster with ReLUs than with classic hyperbolic tangent transfer functions [7].

## 5.2 Video classification

A more challenging and likewise interesting task is the classification of videos because instead of images where the input remains static the input of videos changes over time.

CNNs already showed state-of-the-art results in image recognition and classification tasks so experiments with videos were just a matter of time. A huge problem was the changing input of videos over time in combination with the long periods of training for CNNs. One solution figured out by a group of researchers of Stanford and Google was to improve the runtime performance. They split up the input data in two different streams, the *context* and the *fovea stream* and created a multiresolutional architecture [5].

The original input dimension was 178x178. The fovea stream keeps track of the centered 89x89 pixels of a frame while the context stream analyzes the whole frame but downsampled to a size of 89x89 total. Therefore the total input dimensionality is halved [5]. An illustration of this architecture is shown in figure 3.

Four different solutions have been introduced by Karpathy et al. to create CNNs capable of classifying videos properly: Single Frame, Early Fusion, Late Fusion and Slow Fusion.

Single Frame classification is an adoption already used for classifying images: only one image is being classified at a time.

Early Fusion combines the pixel input information for a specific time window.

Late Fusion calculates two different single frames which are combined at the last fully-connected layers.

Slow Fusion is a mix of Early and Late Fusion: Several input images pixel information is combined two several streams and those streams are additionally combined every few layers.

According to Karpathy et al. [5] Slow Fusion had the best results with top-1 and top-5 error rates of 60.9% and 80.2%.

## 5.3 Object tracking

Object tracking in videos is an even more complex scenario because some object in a video has to be recognized and followed by its position over time. Varying background, changing in the lighting, unexpected motions etc. render this a challenging task.

One supposed solution is the use of stacked denoising autoencoders to learn common features of a dataset and then have an additional deep neural network for tracking those learned patterns over time [12].

Also some kind of preprocessing with particle filters showed good results for tracking purposes [12].

(More to come - complex stuff has not be shortened and summarized enough until now)

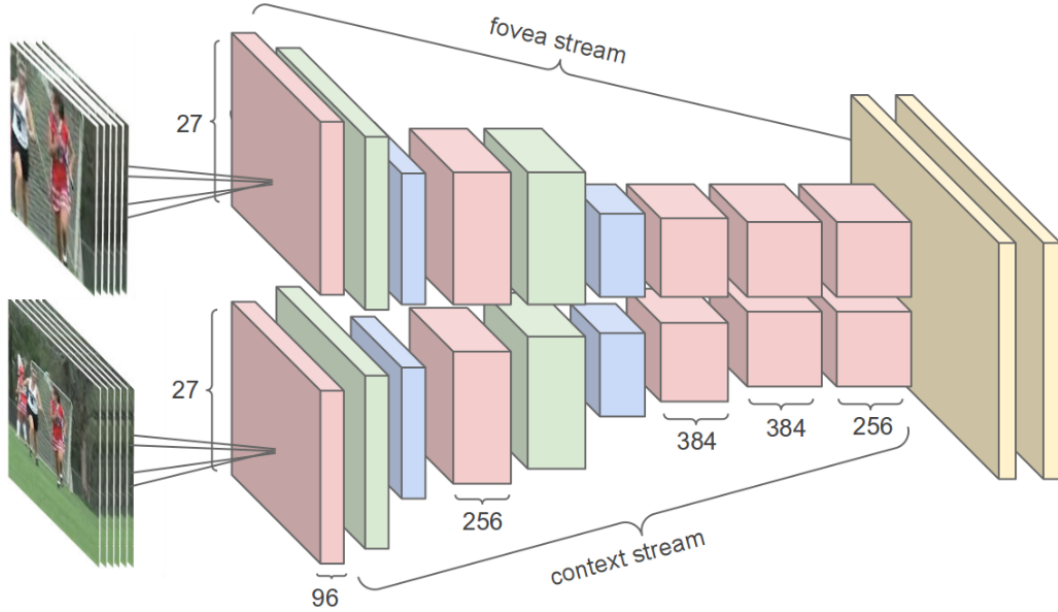


Figure 3: CNN architecture by Karpathy et al. The input information is processed in parallel by the two streams. Graphics retrieved from “Large-scale Video Classification with Convolutional Neural Networks” by Andrej Karpathy et al. [5]

## 6 Conclusion

In the beginning of 2012 the best top-5 error rate on ImageNet was 16.4% by “Team SuperVision” and in 2014 the best top-5 error rate was 6.67% by Team GoogLeNet [10]. This shows the rapid development of deep learning and CNNs but also that there is enough potential for even further improvement. At least for the next years the performance of humans in classifying images will be unmatched which indicates that today’s models have to be improved.

In many papers [10, 11, 7, 12, 1] researchers say improving the computational resources also improves machine learning architectures because CNNs and other neural network architectures could deliver better results if they would be enlarged (which is not possible with current resources) but also new concepts, filters and architectures can improve machine learning.

On the one hand CNNs are one of maybe the most promising architecture at the moment [10, 7] with astounding results in classification tasks. Larger CNNs, more training data, enhanced precision and perhaps some architectural / mathematical improvements are likely to boost their performance even further. On the other hand (although they are less resource demanding than pure, fully-connected neural networks for example) their resource consumption is relatively high and as shown in “Deep Hierarchies In The Primate Visual Cortex” [8] they are not equivalent nor equally effective as biological neural architectures. This displays that CNNs can only be a temporal model but not the ultimate solution and new architectures should be researched.

## References

- [1] D. Cireşan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3642 – 3649, June 2012.
- [2] Dan Cireşan, Ueli Meier, Jonathan Masci, and Jürgen Schmidhuber. Multi-column deep neural network for traffic sign classification. *Neural Networks*, 32:333–338, 2012.
- [3] Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines vinod nair. In *Proc. 27th International Conference on Machine Learning*, 2010.
- [4] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.
- [5] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [6] Alireza Khotanzad and Jiin-Her Lu. Classification of invariant image representations using a neural network. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 38(6):1028–1038, 1990.
- [7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [8] Norbert Kruger, Peter Janssen, Sinan Kalkan, Markus Lappe, Ales Leonardis, Justus Piater, Antonio Jose Rodriguez-Sanchez, and Laurenz Wiskott. Deep hierarchies in the primate visual cortex: What can we learn for computer vision? *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(8):1847–1871, 2013.
- [9] Bente Pakkenberg, Dorte Pelvig, Lisbeth Marner, Mads J. Bundgaard, Hans Jørgen G. Gundersen, Jens R. Nyengaard, and Lisbeth Regeur. Aging and the human neocortex. *Experimental gerontology*, 38(1):95–99, 2003.
- [10] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- [11] Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2553–2561. Curran Associates, Inc., 2013.

- [12] Naiyan Wang and Dit-Yan Yeung. Learning a deep compact image representation for visual tracking. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 809–817. Curran Associates, Inc., 2013.