



DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Robotics, Cognition, Intelligence

Evaluation and Generalization of Capsule Networks in Neurorobotics

Jean Amadeus Elsner





DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Robotics, Cognition, Intelligence

Evaluation and Generalization of Capsule Networks in Neurorobotics

Evaluation und Generalisierung von Kapsel Netzwerken in der Neurorobotik

Author:	Jean Amadeus Elsner
Supervisor:	Prof. Dr. Alois C. Knoll
Advisor:	Alexander Kuhn
Submission Date:	TBD



I confirm that this master's thesis in robotics, cognition, intelligence is my own work and I have documented all sources and material used.

Munich, TBD

Jean Amadeus Elsner

Acknowledgments

Abstract

The last few years have seen great strides being made in the fields of artificial intelligence and robotics. Many of the advancements are powered by the versatility of artificial neural networks. Especially in computer vision, deep learning architectures are particularly successful. For robotic systems, visual sensors are often their primary means to perceive their environment and negotiate it successfully. Thus computer vision systems form a vital part of the artificial intelligence necessary to allow autonomous operation of robots. There are however several factors limiting the widespread use of such systems in robotics. Namely, the need for massive data sets, the lack of generalization to novel configurations and the high computational cost associated with the training of neural networks. Neurorobotics tries to avoid these limitation by mimicking the behavior of biological systems, e.g. by using biologically inspired spiking neural networks. Recently, capsule networks have been suggested as an alternative biologically plausible computer vision system, albeit on a higher level of abstraction, based on second generation neural networks. Capsules are conceived as groups of neurons, that represent the presence of an entity and its instantiation parameters. It is proposed that by dynamically constructing a parse tree from a network of capsules, viewpoint invariance and better generalization can be achieved. In this thesis, object recognition tasks generated by the Neurorobotics Platform are used to evaluate the performance of capsule networks and quantify their ability to generalize compared to established methods such as convolutional and spiking neural networks.

Contents

Acknowledgments	iii
Abstract	v
1 Introduction	1
1.1 Motivation	2
1.2 Scope	2
2 State of the Art	3
2.1 Artificial Neural Networks for Object Recognition	3
2.2 Convolutional Neural Networks	5
2.3 Neurorobotics and Spiking Neural Networks	13
2.4 Limitations of Artificial Neural Networks	18
3 Capsule Network Architectures	21
3.1 Transforming Auto-encoders	21
3.2 Dynamic Routing Between Capsules	24
3.3 Matrix Capsules With EM Routing	27
4 Experimental Setup	29
5 Results	31
6 Discussion	33
7 Conclusion	35
List of Figures	37
List of Algorithms	39
List of Tables	41
Bibliography	43

1 Introduction

The success of techniques from the field of artificial intelligence (AI) in recent years is quite evident. Not just are there ever more real world applications of AI found in search engines, security systems, industrial automation and more. But the term *artificial intelligence* itself has become a vogue word¹, even outside of academic literature. A point can be made that much of the drive behind AI comes from the successful application of deep artificial neural networks in computer vision tasks [1]. As neural networks designed for computer vision share similarities with the mammalian visual cortex [2], a similar case can be made for the contribution of neuroscience: the eye is arguably the most intensely studied of the human sensory organs and the visual cortex is one of the best understood parts of the brain. It is therefore only sensible to try and mimic the properties and behavior of brains when developing new systems. While cameras already surpass the spatial and temporal resolution of biological eyes, the human visual system remains unmatched in visual-cognitive tasks. Moreover, computer vision solutions currently available often have a high latency, precluding them from being used in real time experiments. In order to be able to validate biologically inspired AI models it is therefore necessary to simulate a realistic environment [3]. The *Neurorobotics Platform*² was conceived to provide just such a simulation: a closed loop between a brain simulation based on artificial neural networks and a physics based simulation of robotic bodies embedded in a dynamic environment [4]. In this thesis, sensor data as generated by an experiment within the Neurorobotics Platform was used to create challenging computer vision tasks. These tasks evaluate general object recognition performance and specific generalization and robustness, e.g. generalization to new viewpoints or robustness against occlusion. The use of a high fidelity simulation platform allows for the generation of big labelled datasets that can be fine-tuned towards the strengths and weaknesses of the architecture at hand. Something that would otherwise require many hours of manual labor.

¹Often it is not immediately apparent what kind of technology is being referred to when the term AI is used. This is true to the point that opposing paradigms claim the mantle of AI. Classical AIs were conceived as systems performing logical inference on a knowledge database in the form of search algorithms (*cognitivist* paradigm) as opposed to the current use of neural networks (an aspect of the *emergent* paradigm).

²The Neurorobotics Platform is developed as a subproject of the *Human Brain Project* (HBP). The HBP is a European Commission Future and Emerging Technologies Flagship Program intended to advance knowledge in the fields of neuroscience, computing and brain-related medicine over a period of 10 years.

1.1 Motivation

The most potent of the AI technologies, artificial neural networks, has long since diverged into two mostly independent branches [5]. On one side, there are networks that are used as tools to study biological nerve cells by closely modeling the physical properties of neurons and synapses. On the other, there are neural networks optimized for machine learning tasks such as object detection in images or speech recognition. The former type of architecture, known as spiking neural network, has inherently great potential, as it is directly derived from biological examples. Still they are almost always outperformed by the latter: neural networks based on multiple layers of perceptrons [6, 7]. Multi layer perceptrons (MLP) however bring many drawbacks that ultimately limit their application (cf. section 2.4). Capsule networks have been proposed as a possible alternative to overcome some of these limitations while increasing biological plausibility over MLPs. The routing of signals between capsule layers, based on grouping votes from the lower layer, allows capsules to encode the image structure into more sophisticated internal representations. In order to address the open question of whether this leads to better generalization, capsule network architectures were explicitly evaluated in this thesis using custom-made datasets from the Neurorobotics Platform and compared to baseline networks representing both ends of the biological plausibility spectrum.

1.2 Scope

The thesis starts with a derivation of the state of the art of the artificial neural networks used as baseline in this thesis in chapter 2. The chapter also includes details about the limitations of the discussed architectures and explains classification by example of object recognition. This is followed by an introduction of the motivation and theory behind capsule networks as well as a brief review of their development in chapter 3. How the Neurorobotics Platform was used to create the datasets as well as how exactly each architecture was trained and tested and what metrics were applied is explained in chapter 4. The numerical results from the experiments described in the previous chapter are presented in chapter 5 and interpreted in a discussion found in chapter 6. Finally, the thesis concludes with an outlook on possible further developments and open questions in chapter 7.

2 State of the Art

This chapter presents an overview of state of the art approaches to object recognition, while focusing on two families of artificial neural network (ANN) architectures, which are motivated quite differently. Object recognition techniques based on *convolutional neural networks* (CNNs) currently dominate the field, achieving state of the art performance by a large margin over classical methods on many datasets [8, 9]. Even though CNNs were originally inspired by findings from neuroscience, by now they have diverged quite a bit from the computational models used to study the brain. The neurorobotics approach to cognitive systems, based on *spiking neural networks* (SNNs) on the other hand, attempts to mimic the animal brain by more closely modelling the physical properties and behavior of neurons and therefore results in biologically more plausible models [10]. Generally speaking, CNNs may be regarded as a more engineering-based approach (or top-down), while SNNs are motivated by results from neuroscience and biology (bottom-up approach).

2.1 Artificial Neural Networks for Object Recognition

Recent years have seen a surge of interest in artificial neural networks and deep learning methods, especially in the field of computer vision. While the theory for training such networks has been around for many years, their recent success is mainly due to the availability of large labelled data sets (so called big data) and the proliferation of highly parallel computing powered by GPUs. One of the specific tasks, deep learning based methods excel at, is object recognition¹: the identification of objects in images or videos (cf. figure 2.1). The significantly better performance of deep neural networks over traditional machine learning methods can be explained by: (i) their hierarchical topology of parameterized non-linear processing units is a fundamentally better probabilistic model and prior for real world data as captured in images leading to better generalization and (ii) they autonomously find good features to extract based on the training data. Interest is further fueled by the myriad potential applications of robust object recognition systems: from automated driving and image-based diagnosis in medicine to robot vision and many more. As deep learning is currently the best candidate for such a system, it is well worth exploring.

¹In machine learning, tasks such as object recognition are referred to as classification problems. The resulting classifiers are part of a broader class of methods called *supervised learning*. In supervised learning, the desired output (i.e. the label) has to be provided for all the samples used for training the classifier.

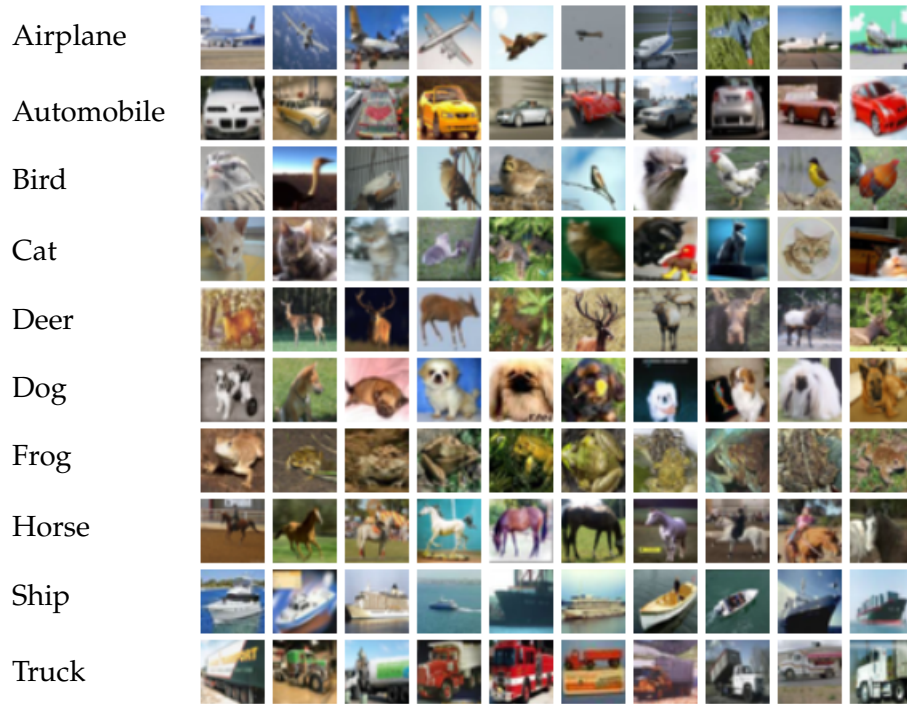


Figure 2.1: Sample images from the CIFAR-10 [11] dataset and their corresponding classes. CIFAR-10 consists of 6000 images at 32 by 32 pixels for each of the 10 classes. Datasets such as this are often used as a benchmark to evaluate the performance of novel ANN architectures for image recognition. This is done by using a subset of the dataset (often referred to as the training set) to train the neural network. The remainder of the images (accordingly called the test set), which the network has not seen before, are used to evaluate the classification accuracy.

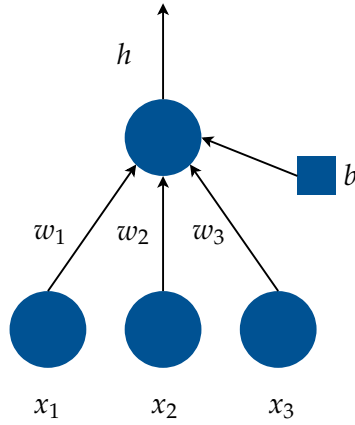


Figure 2.2: Illustration of an artificial neuron with three input connections. Artificial neurons constitute the basic non-linear computational units of neural networks. The neuron receives the activations of three lower level neurons weighted by the learnable w_i as well as a learnable bias b . After applying the nonlinearity (also known as activation function), the neuron outputs its activation h as in equation 2.1

2.2 Convolutional Neural Networks

CNN architectures are generally distinguished by their use of specific types of neuron-layers, namely, convolutional, pooling and fully connected layers. While wildly different network topologies may be found in literature, characterized by their use of skip connections, number of layers, number of paths etc., CNNs can always be reduced to these three basic layer types.

Fully Connected Layer

Each neuron in a fully connected layer is connected to all the activations in the previous layer. The activation of a single neuron (cf. figure 2.2) is calculated by applying a nonlinearity to the weighted sum of its inputs and a bias.

$$h = g\left(\sum_i w_i x_i + b\right) \quad (2.1)$$

With the nonlinear function g , the learnable weights w_i , the input activations x_i and the learnable bias b . In the case of a fully connected layer, the activations can be computed using matrix multiplication. In tensor notation this may be written as:

$$\mathbf{h}_l = g_l(\mathbf{W}_l^T \mathbf{h}_{l-1} + \mathbf{b}_l). \quad (2.2)$$

With N_l denoting the number of neurons in layer l , \mathbf{W}_l is an $N_{l-1} \times N_l$ dimensional weight matrix, \mathbf{b}_l an N_l dimensional vector and g_l the N_l dimensional vectorized activa-

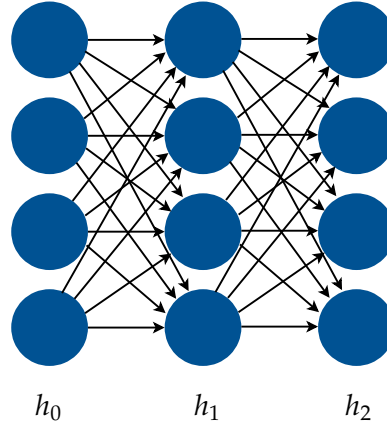


Figure 2.3: Illustration of a fully connected neural network as a directed graph with input layer h_0 and output layer h_2 . Layers between the input and output layer are usually referred to as *hidden* layers. Note that the biases are not explicitly shown.

tion function of layer l .

$$g_l(\mathbf{x}) = (g_l(x_1), \dots, g_l(x_{N_l}))^T \quad (2.3)$$

The computational power of neural networks is shown by the universal approximation theorem. The theorem states that networks with a single hidden layer (cf. figure 2.3) containing a finite number of neurons can approximate arbitrary continuous functions on compact subsets of \mathbb{R}^n [12, 13].

Convolutional Layer

In a convolutional layer, the activities of the input layer are convolved with a number of trainable kernels so as to create the same number of feature maps. In computer vision it is common to view the layer's neurons as two-dimensional grids of neurons arranged in channels (cf. figure 2.5). These grids correspond to pixels and color channels in the case of the input layer or activities (feature maps) resulting from convolution with different kernels in the case of intermediate convolutional layers. For a feature kernel $F_{m,n}^l$ of size $M_l \times M_l$ and an input layer $l-1$ with $N_{l-1} \times N_{l-1}$ neurons, the corresponding feature map activities of layer l are computed as

$$h_{m,n}^l = g_l \left(\sum_{m'} \sum_{n'} F_{m',n'}^l h_{m+m',n+n'}^l + b_l \right). \quad (2.4)$$

This is the same as a two-dimensional discrete cross correlation.

$$(f \star g)[n] \stackrel{\text{def}}{=} \sum_{m=-\infty}^{\infty} f^*[m] g[m+n] \quad (2.5)$$

Where f^* is the complex conjugate of the discrete function f . Strictly speaking, the use of the term *convolution* in neural network literature is therefore a misnomer, as either the filter or the image (or feature map) would have to be flipped before the operation. However, as the weights of the kernel are actually learned by the network, the result will be the same.

For a 2D single channel input layer of size $N_0 \times N_0$, the number of trainable parameters for a convolutional layer with a single $M \times M$ filter is $M^2 + 1$ for the kernel weights and bias respectively, compared to $N_0^2 N_1^2 + N_1^2$ for a fully connected layer of size $N_1 \times N_1$. The sparsity in learnable parameters in convolutional layers compared to fully connected layers (cf. figure 2.4) is often referred to as *weight sharing* (an entire channel *shares* the weights of a single kernel).

In the case of multiple input channels, the kernels actually extend through the whole depth of the input layer's volume. Equation 2.4 can be extended to accommodate multiple input channels and feature kernels. The activities of feature map k in layer l are then computed as

$$h_{k,m,n}^l = g_l \left(\sum_{k'}^{K_{l-1}} \sum_{m'}^{M_l} \sum_{n'}^{M_l} F_{k,k',m',n'}^l h_{k',m+m',n+n'}^l + b_{l,k} \right). \quad (2.6)$$

With K_l the number of channels in layer l . Most deep learning frameworks offer additional hyperparameters for convolutional layers, like stride and zero padding. Zero padding adds additional rows and columns of zero-activities around the input layer channels, effectively allowing the output feature maps to have the same size (i.e. height and width, depth is determined by the number of kernels) as the unpadded input layer. Stride on the other hand defines the integer value by which the kernels are moved during convolution. This can be used to keep the receptive fields from overlapping too much. Equation 2.6 is easily extended to take the stride S_l of layer l into consideration.

$$h_{k,m,n}^l = g_l \left(\sum_{k'}^{K_{l-1}} \sum_{m'}^{M_l} \sum_{n'}^{M_l} F_{k,k',m',n'}^l h_{k',S_l m+m',S_l n+n'}^l + b_{l,k} \right) \quad (2.7)$$

The spatial dimension of this feature map can be computed as a function of the input layer size $N_{l-1} \times N_{l-1}$ and the amount of zero padding P_l .

$$N_l = \frac{N_{l-1} - M_l + 2P_l}{S_l} + 1 \quad (2.8)$$

Pooling Layer

Pooling is a form of sub- or downsampling using a sliding window similar to convolution. Most often the stride is set equal to the window size - that way the windows don't overlap. An operator is applied to each window, which selects a single neuron. Common operators are maximum, average or L_2 -Norm. The pooling is applied to each channel

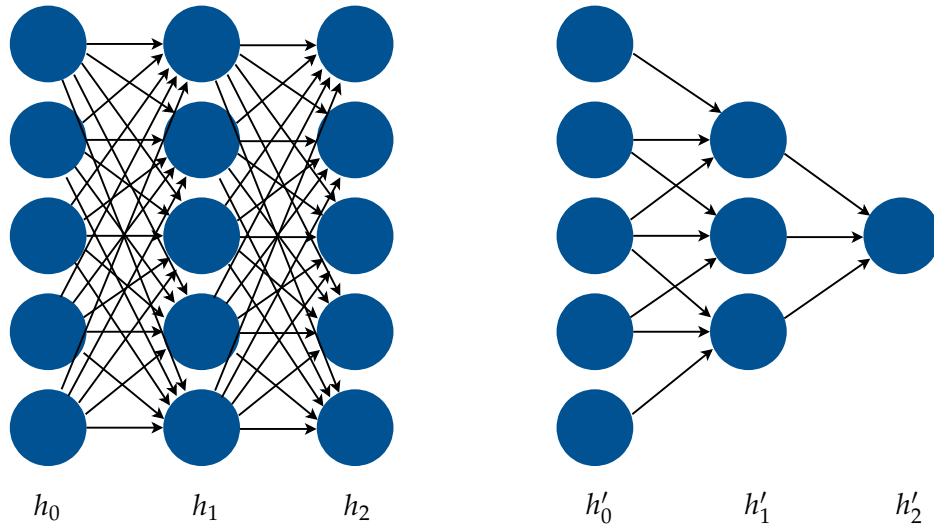


Figure 2.4: Illustration of fully connected layers compared to convolutional layers, making apparent the sparsity of the latter relative to the former. The convolutional layers h'_1 and h'_2 perform a 1D convolution with a filter of size 3. Note that due to weight sharing the number of weight parameters between h'_0 and h'_1 is actually only 3. Also note that convolution reduces the number of neurons (or pixel resolution in the case of an image). For an input layer with N neurons and a filter kernel of size M the number of neurons in the convolutional layer computes as $N - M + 1$.

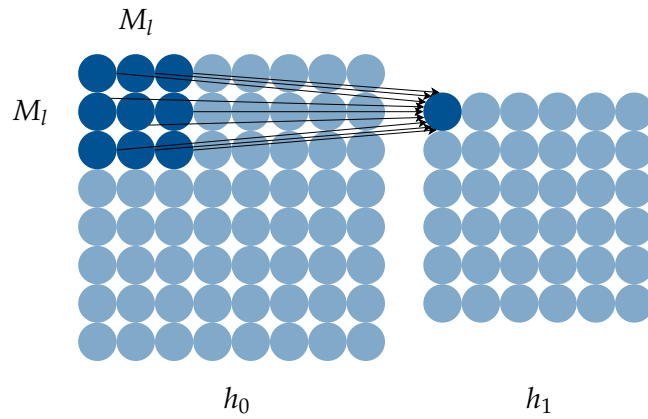


Figure 2.5: Illustration of a 2D convolution with a kernel of size $M_l \times M_l$. As the resulting feature map h_1 has a lower resolution than h_0 , applying a convolutional layer is sometimes also referred to as downsampling. For an input layer with $N \times N$ neurons, the feature map's size is computed as $(N - M_l + 1) \times (N - M_l + 1)$. The area comprised of the input neurons that are used to calculate the activity of a feature map neuron (highlighted neurons in layer h_0) are known as that neuron's *receptive field*.

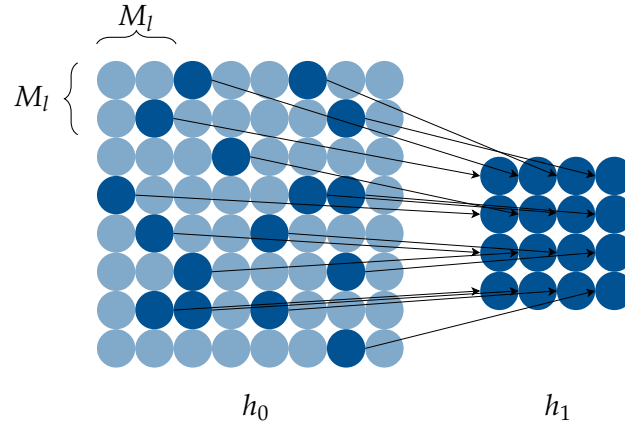


Figure 2.6: Illustration of a pooling layer. The pooling uses a 2×2 window with stride 2 to segment the input layer into non-overlapping tiles. An operator (e.g. maximum) is applied to each 2×2 segment to select a single neuron (the highlighted neurons in layer h_0). The activities of the selected neurons are arranged in a new layer h_1 of size $\frac{N_{l-1}}{M_l} \times \frac{N_{l-1}}{M_l}$ with $M_l \times M_l$ and $N_{l-1} \times N_{l-1}$ the size of the window and the input layer respectively.

and leads to a reduction in the spatial dimensions. This reduction corresponds to a loss of information or reduction in degrees of freedom and therefore reduces the amount of computation required as well as the possibility of overfitting. The choice of pooling operator has a profound effect on the generalization and speed of convergence of CNNs and recently maximum pooling has proven to be the most successful method [14, 15].

CNN Architecture

The key idea behind the combination of these three types of layers is that essential visual features such as edges and corners within a convolutional layer's receptive field are combined to form higher level features such as shapes by subsequent convolutional layers. In between these convolutions, pooling operations select the most salient features and reduce the computational size of the network. The convolutional kernels are effectively trainable feature detectors and due to weight sharing and pooling naturally incorporate a measure of translational invariance. This hierarchical organization of receptive fields is similar in structure to the mammalian visual cortex [16] and sometimes CNNs are seen to be directly derived from it [17, 18]. More often than not however, the use of convolutional layers with weight sharing is motivated as a means to achieve translational equivariance and faster computation compared to fully connected layers [19, 20, 21]. Finally, fully connected layers are placed on top of the network in order to perform high level inference (cf. figure 2.7). A CNN's architecture is therefore largely defined by its topology: the number and types of layers as well their neurons and the connections between them. The number of stacked layers is referred to as the network's

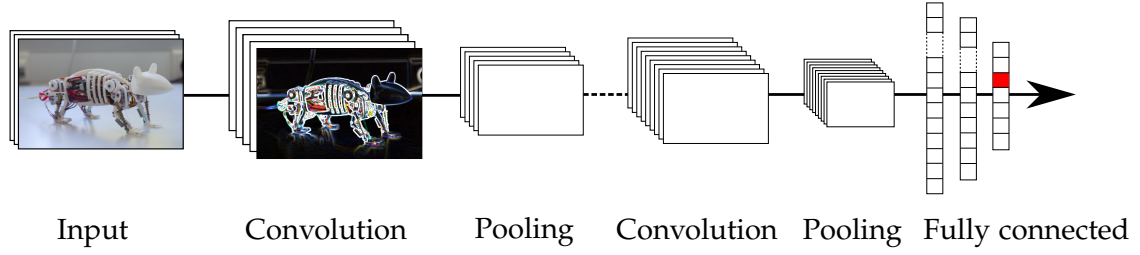


Figure 2.7: Illustration of a typical CNN architecture for object recognition. Note that the feature maps of the first convolutional layer extracted elementary features like edges. The last layer's neurons correspond to the trained classes, while the output neuron with the highest probability is picked as the network's prediction during inference (highlighted in red).

depth. Current networks often employ multiple paths and skip connections (i.e. a layer's output not only flows to its direct descendant, but skips several layers) allowing for topologies several hundred layers deep, hence *deep learning* [22, 23, 24, 25].

Training

For object recognition, the output of the final layers are the class probabilities. This is achieved by applying a normalizing activation function to the the output of the last layer, such as the *softmax* function. In case of the softmax function, the network's output $y_k(\mathbf{x})$ for the class k given input sample \mathbf{x} thus becomes

$$y_k(\mathbf{x}) = \frac{\exp(a_k)}{\sum_{k'=1}^K \exp(a_{k'})}. \quad (2.9)$$

With $\mathbf{a} = (a_1, \dots, a_K)^T$ the linear activities of the last layer's neurons (in machine learning literature often referred to as *logits*) and K the number of classes. Using one-hot coding for the target class

$$\mathbf{t} = (t_1, \dots, t_{k'}, \dots, t_K)^T = (0, \dots, 1, \dots, 0)^T, \quad (2.10)$$

the probabilistic model can be defined as a function of the neural network.

$$P(\mathbf{t} \mid \mathbf{x}, \boldsymbol{\theta}) = \prod_{k=1}^K y_k(\mathbf{x})^{t_k} \quad (2.11)$$

With $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{b})$ the set of the network's trainable parameters. The model's likelihood is then given by

$$P(\mathbf{T} \mid \mathbf{X}, \boldsymbol{\theta}) = \prod_n P(\mathbf{t}_n \mid \mathbf{x}_n, \boldsymbol{\theta}). \quad (2.12)$$

Where the index n denotes n th training sample and target class (also known as label). The negative logarithm of the likelihood, called the negative *log-likelihood* defines the classifier's error function E_{ML} .

$$E_{\text{ML}}(\boldsymbol{\theta}) = -\log(P(\mathbf{T} | \mathbf{X}, \boldsymbol{\theta})) = -\sum_n \log P(\mathbf{t}_n | \mathbf{x}_n, \boldsymbol{\theta}) \quad (2.13)$$

$$= -\sum_n \sum_{k=1}^K t_{n,k} \log(y_k(\mathbf{x}_n)) \quad (2.14)$$

The error function resulting from the softmax activation function specifically is referred to as the *cross entropy* (equation 2.14). Now the network's parameters may be trained by minimizing its classification error w.r.t. $\boldsymbol{\theta}$. As the logarithm is a convex function and the error function is the *negative* log-likelihood, this is equivalent to a maximum likelihood estimate. The gradient of the error function can be computed using the derivative chain rule (cf. figure 2.8), this method is known as *backpropagation* and forms the computational backbone of CNN architectures.

$$\frac{\partial E_{\text{ML}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_{l-1}} = \frac{\partial E_{\text{ML}}(\boldsymbol{\theta})}{\partial h_{l-1}} \frac{\partial h_{l-1}}{\partial \boldsymbol{\theta}_{l-1}} = \frac{\partial E_{\text{ML}}(\boldsymbol{\theta})}{\partial h_l} \frac{\partial h_l}{\partial h_{l-1}} \frac{\partial h_{l-1}}{\partial \boldsymbol{\theta}_{l-1}} \quad (2.15)$$

The gradient calculated based on the training samples is then used to update the network's parameter.

$$\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^{(s)} - \eta \nabla E_{\text{ML}}(\boldsymbol{\theta}^{(s)}) \quad (2.16)$$

Where s counts the number of training epochs (a run through all the samples) and η is a hyperparameter known as the *learning rate*, that effectively describes the size of the step taken in direction of the gradient (cf. figure 2.9). In practise, the gradient is often calculated as an average based on a subset of randomly selected training samples.

$$\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^{(s)} - \eta \frac{1}{M} \sum_{n=1}^M \nabla E_n(\boldsymbol{\theta}^{(s)}) \quad (2.17)$$

The $\frac{N}{M}$ sets of training samples are called (mini-) *batches* and are trained sequentially until all samples have been seen, completing an epoch. Training algorithms based on this optimization scheme are known as *stochastic gradient descent* (SGD). It is at this point that vectorization libraries are used to leverage the power of modern GPUs by processing an entire mini batch in parallel.

State of the art CNN architectures expand upon the building blocks introduced in this chapter in various ways depending on the application. Often it is necessary to regularize the network's weights in order to prevent overfitting. This can be achieved by adding penalties for large weights to the loss function (known as weight decay) or randomly setting some weights to zero during training (so-called dropout) [26, 27, 28, 29]. This

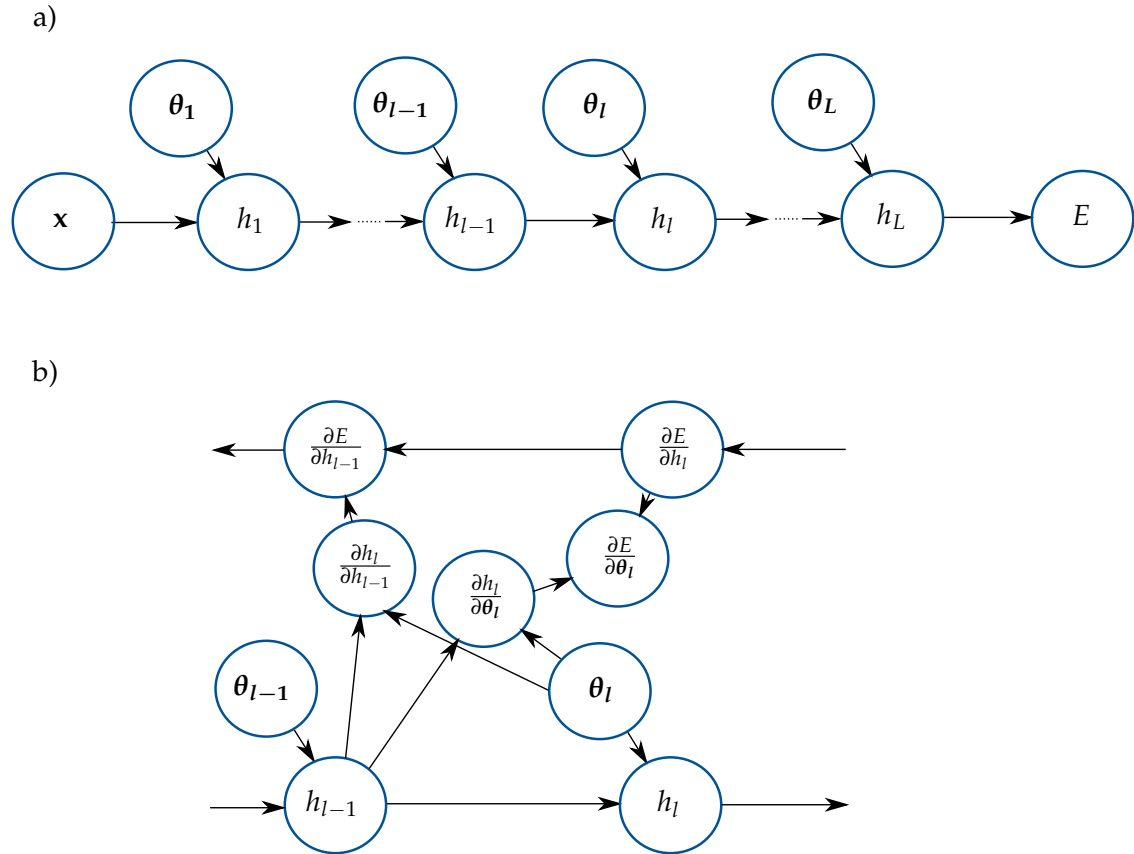


Figure 2.8: Computational graph for error propagation within a neural network. a) shows the dependency of the layers on both the lower layer's output as well as the parameters θ_l during inference (also known as forward pass). In b) the computational nodes for the error gradient as needed for error backpropagation were added (known as the backward pass).

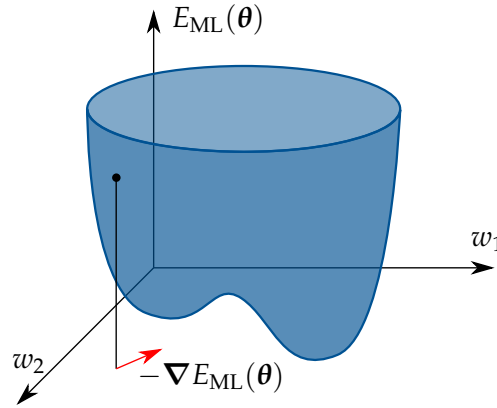


Figure 2.9: Illustration of the surface of a neural network’s error function in parameter space. The function is non-convex and contains a local minimum. Note that the resulting minimum depends on both the learning rate and the initialization of the parameters. In this case, following the negative gradient (illustrated as a red arrow) using small steps will converge on a non-optimal local minimum.

can be complemented with adaptive learning schemes² in order to reliably train neural networks [30, 31, 32, 33, 34, 35, 36]. To make the process of training a neural network even more robust and prevent the gradients from becoming either too large or too small, it has proven useful to normalize a layer’s input. In addition to stabilizing training, batch normalization can also speed up the convergence [37, 38, 39]. Finally, both the network’s topology and the loss function may be chosen in such a way as to encourage the learning of internal representations to fit the task (e.g. learn the pose of a 3D object) [40, 41, 42, 43]. Even though CNNs have been applied successfully in many fields, there is not yet a theory to derive the topology and hyperparameters best suited for a given problem or predict the network’s performance.

2.3 Neurorobotics and Spiking Neural Networks

Neurorobotics is an exciting new field that aims to study intelligence by way of an interdisciplinary approach combining computational neuroscience, robotics and artificial intelligence. Taking inspiration from biology is clearly a worthwhile effort as traditional engineering approaches like model based control systems for robotics or physical symbol systems for artificial intelligence have so far failed to replicate the intelligent behavior of even the simplest of animals. The mathematical models of the nervous system used in neurorobotics are based on results from neuroscience, which is concerned with the study

²Adaptive extensions to SGD try to adjust learning hyperparameters like the learning rate based on performance on training data. It is for example quite common to decrease the learning rate after each training epoch that didn’t increase accuracy or use such a scheme with independent learning rates for each weight.

of intelligent biological systems on all levels, from the fast signal processing of single neurons to emergent long-term behavior like memory, perception and consciousness [44, 45, 46, 47, 48]. Artificial implementations or simulations of these models are then used to control a robotic body. This idea of an *embodied* intelligence controlled by brain inspired algorithms, which is in turn *embedded* in an environment it can perceive and interact with is central to neurorobotics. Brain, body and environment in a neurorobotic experiment form a closed perception-action loop where the brain receives percepts from the body's sensors, based on which it will produce signals to move the body, causing a change in the perception of the environment etc. Observing the interactions between the robot and its environment as well as its response to specific stimuli will in turn allow scientists to draw conclusions about the biological plausibility of the used brain model and further their understanding of how the brain works in conjunction with the body. A deeper understanding again leads to better models and more realistic simulations – this way neurorobotics inherently amplifies the transfer of knowledge and feedback between the involved disciplines [49].

Spiking neural networks (SNNs) are the method of choice for brain simulations in neurorobotics. They use more elaborate neuronal models than those discussed in the previous section and therefore are able to mimic the behavior of biological networks of neurons more plausibly. Beyond biological plausibility, there are also potential engineering benefits to be gained by mimicking brains. Due to sparse encoding of information, the human brain maintains high versatility and throughput while only consuming about 20 W of power on average [50]. The ability to adapt to a dynamic environment and react in real-time under constrained resources is essential for the deployment of mobile intelligent agents like robots. In SNNs dynamics are added to the states of the neurons and synapses (in CNNs these are the activities and weights respectively). Biological neurons are surrounded by cell membranes that act as insulators, which can be charged by other neurons in the network with short electrical impulses over the synapses. As the membrane is not a perfect insulator, this charge degrades over time. If a neuron model represents the incoming signals as discrete events (this is a reasonable simplification, as the shape is roughly the same for all of a neuron's impulses) it is referred to as *integrate-and-fire*. Models that additionally include a leakage (i.e. the neuron's membrane potential degrades over time) are by far the most popular and known as *leaky integrate-and-fire* (LIF) [51, 52]. A LIF neuron can be modelled as an RC circuit (cf. figure 2.10) that fires off an impulse, once the potential reaches a threshold. Using Kirchhoff's current law, the incoming current $I(t)$ can be written as the sum of the resistive current and the current charging the capacitor.

$$I(t) = I_R + I_C = \frac{u(t) - u_{\text{rest}}}{R} + C \frac{du}{dt} \quad (2.18)$$

Where in the second step Ohm's law was used for the linear resistor R as well as the definition of capacity $C = \frac{q}{u}$ and current $I_C = \frac{dq}{dt} = C \frac{du}{dt}$ for the capacitor C . Multiplying

by the resistance R this can be written as

$$RI(t) = u(t) - u_{\text{rest}} + \underbrace{RC}_{\tau} \frac{du}{dt} \quad (2.19)$$

$$\tau \frac{du}{dt} = -(u(t) - u_{\text{rest}}) + RI(t) \quad (2.20)$$

$$\circ \quad (2.21)$$

$$\tau s U(s) - u(0^+) = -U(s) + \frac{u_{\text{rest}}}{s} + RI(s) \quad (2.22)$$

$$U(s) = u_{\text{rest}} \frac{\frac{1}{\tau}}{s(s + \frac{1}{\tau})} + u(0^+) \frac{1}{s + \frac{1}{\tau}} + \frac{R}{\tau} I(s) \frac{1}{s + \frac{1}{\tau}} \quad (2.23)$$

$$\circ \quad (2.24)$$

$$u(t) = u_{\text{rest}} \left(1 - e^{-\frac{t}{\tau}}\right) + u(0^+) e^{-\frac{t}{\tau}} + \frac{R}{\tau} \int_0^\infty e^{-\frac{t-t'}{\tau}} I(t-t') dt' \quad (2.25)$$

$$u(t) = u_{\text{rest}} + \frac{R}{\tau} \int_0^\infty e^{-\frac{t-t'}{\tau}} I(t-t') dt'. \quad (2.26)$$

Where the Laplace transform was used to find the solution of the linear differential equation and the initial membrane potential was assumed to be at resting potential $u(0^+) = u_{\text{rest}}$. In the case of LIF models, the input current will consist of discrete events represented by delta distributions.

$$I_{\text{pre}}(t) = q \sum_i \delta(t - t_{i,\text{pre}}^{(f)}) \quad (2.27)$$

This is known as a *spike train*, with $t_{i,\text{pre}}^{(f)}$ the firing times of the presynaptic neurons and q the charge delivered with each spike. In addition to the presynaptic spike train, each neuron generates its own postsynaptic spike train. The firing times $t_{k,\text{post}}^{(f)}$ are determined by the threshold ϑ .

$$u(t_{k,\text{post}}^{(f)}) = \vartheta \quad (2.28)$$

Once the threshold is reached, the potential is reset to $u_{\text{reset}} < \vartheta$, firing the postsynaptic charge $C(\vartheta - u_{\text{reset}})$. Putting both the pre- and postsynaptic activations $I(t) = I_{\text{pre}}(t) + I_{\text{post}}(t)$ into equation 2.26 leads to the momentary voltage of a LIF neuron.

$$u(t) - u_{\text{rest}} = \frac{R}{\tau} \left(q \sum_i \exp\left(-\frac{t - t_{i,\text{pre}}^{(f)}}{\tau}\right) - C(\vartheta - u_{\text{reset}}) \sum_k \exp\left(-\frac{t - t_{k,\text{post}}^{(f)}}{\tau}\right) \right) \quad (2.29)$$

As the impulses are extremely sparse in time, LIF neurons generally allow for much more efficient computation compared to the neurons found in CNNs. Alternatively this can be thought of as higher information density: the frequency of rate-based codes corresponds to averaging a temporal window and therefore a loss of information over

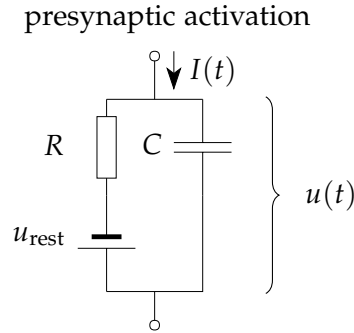


Figure 2.10: Electrical circuit of a leaky integrate and fire neuron. The unit saves incoming charges $I(t)$ called the presynaptic activation into the capacitor C , essentially summing them up or *integrating* them. Once the charge reaches a threshold value, the neuron discharges or *fires* (postsynaptic activation). The dissipation of ions through the cell’s membrane is modelled by allowing the charge to *leak* through the resistor R .

the pulse-based code. While it is possible to use a code based on the firing rate, the full potential of spiking neurons is only levelled, when information is also encoded in the timing of the firing [53, 54, 55]. A single spiking neuron, encoding information in the rate and temporal patterns of its spike train, can replace hundreds of second generation neurons³ [56, 57, 58]. The synapses interconnecting the spiking neurons in a SNN are scaled by learnable weights, the same as CNNs. Many of the ideas developed for CNNs, like perceptive fields, different types of layers, batch normalization etc., can be directly applied to SNNs. There are however two major differences: (i) the input data must be encoded into spike trains and (ii) the synaptic weights must be learned based on the dynamics of the spiking neurons. In computer vision experiments it is quite common to use rate based stochastic encoding schemes. The spike trains are generated by sampling from a stochastic distribution (e.g. Poisson) with firing rates proportional to the intensity of the pixels [59]. SNNs may be trained using biologically plausible algorithms. Learning rules which strengthen synapses based on the timing of spikes are referred to as *Hebbian* learning and often colloquially summarized by the phrase “Cells that fire together, wire together” [60]. Mathematically, the change in synaptic weight Δw_{ij} between two neurons i and j may be expressed as

$$\Delta w_{ij} \propto v_i v_j. \quad (2.30)$$

with v_i the activity of neuron i . Because there is no labelled data or any training signal involved, Hebbian learning rules are inherently unsupervised. Methods that consider

³The first generation of neurons or neural networks generally refers to networks consisting of linear neurons (called *perceptrons*), while the second generation is understood to refer to architectures using artificial neurons with nonlinear activation functions, such as CNNs as discussed in section 2.2. Networks based on the dynamics of spiking neurons are known as third generation neural networks.

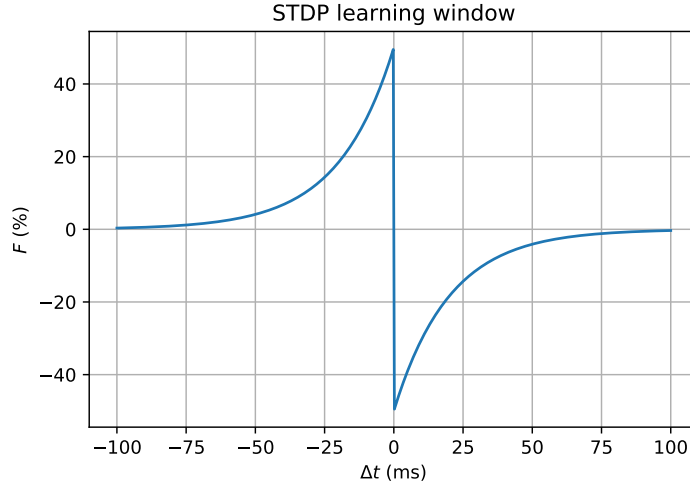


Figure 2.11: Typical STDP window function. The change in conductance for a synapse due to a single pre- and postsynaptic spike pair in STDP is proportional to $F(\Delta t)$ with Δt the time of the presynaptic spike minus the time of the postsynaptic spike. In this figure, F was expressed as a percentage. Note that in this case, action potentials that are more than 100 ms apart have virtually no influence, while rapid spikes induce up to 50 % LTP/LTD.

the precise timing and order of pre- and postsynaptic spikes are known as *spike-timing-dependent plasticity* (STDP). In STDP the firing of a postsynaptic spike increases the strength of presynaptic weights that fired shortly before. This is referred to as *long term potentiation* (LTP). The reverse order leads to a decrease in synaptic strength and is conversely known as *long term depression* (LTD). The change in weight Δw_i for a synapse from a presynaptic neuron i according to STDP can also be expressed mathematically.

$$\Delta w_i = \sum_{f=1}^N \sum_{n=1}^N F(t_n^{(j)} - t_f^{(i)}) \quad (2.31)$$

Where $t_n^{(j)}$ and $t_f^{(i)}$ denote the times of the firing of spikes by the postsynaptic neuron j and the presynaptic neuron i respectively. The STDP function F or *learning window* largely determines the specific behavior of the STDP model. A popular choice, that also corresponds well with experimental findings [61] is given by

$$F(\Delta t) = \begin{cases} A_+ \exp\left(\frac{\Delta t}{\tau_+}\right) & \text{for } \Delta t < 0 \\ -A_- \exp\left(-\frac{\Delta t}{\tau_-}\right) & \text{for } \Delta t \geq 0. \end{cases} \quad (2.32)$$

With $A_{+/-}$ linear coefficients and $\tau_{+/-}$ a kind of decay defining the size of the window (cf. figure 2.11). While research into biologically plausible learning algorithms is ongoing

[62, 63], they are generally outperformed by engineering methods such as gradient descent⁴ in tasks that allow for supervised training (e.g. classification and regression) [59]. Training a deep network of spiking neurons discriminatively using gradient descent with backpropagation is difficult, because the LIF neuron's activity is non-differentiable at the time of spikes (cf. figure 2.12). For the supervised training of feed-forward⁵ networks using backpropagation, there are generally two established approaches. The first trains a second generation network and afterwards converts it to a SNN for inference. The second type of method directly optimizes an objective function using an approximation of spatio-temporal gradient descent similar to backpropagation in CNNs. Recently, deep spiking convolutional networks have been trained successfully using backpropagation by treating the discontinuities at spike times as noise i.e. approximating a smooth signal [64, 65]. While all approaches to SNNs currently do not perform quite as well as CNNs on computer vision tasks, they are inherently computationally less expensive and could at the very least lead to power efficient hardware implementation in the form of neuromorphic systems⁶ [66].

2.4 Limitations of Artificial Neural Networks

Convolutional layers in CNNs use translated replicas (via weight sharing) of learned feature detectors. The rationale behind this is that knowledge about good features in one part of the image may be useful in other regions as well. While this has proved quite a powerful semantic and grants CNNs a measure of invariance to translations, they do not generalize well to other transformations. This requires the networks to be trained using a large number of training images, covering all desired variations in viewpoint, scale, lighting, color etc. It is common to artificially extend the training set by applying image transformations. This process is known as *data augmentation*. [67, 68]. Nevertheless the process of creating datasets for (deep) neural networks is quite time consuming and expensive. Even if labelled data is available, training is computationally expensive and can take several days [69, 70]. The same can be said for inference on trained networks, which often requires great computational resources and still suffers from relatively high latency [71]. These drawbacks along with the fact that high computational costs come with large energy requirements, make real time and especially mobile applications under constrained resources difficult. Another shortcoming of CNNs is due to their use of pooling layers. The bottle-neck architecture of CNNs allows them to use a growing

⁴Gradient descent is generally regarded as not a realistic option for how the brain learns. This is because it requires a teaching signal in the form of labelled data and calculates a gradient specific to each neuron based on that signal.

⁵Feed-forward networks allow propagation of signals in only one direction and do not contain any cycles. While there are other successful architectures like recurrent neural networks (RNNs), restricted Boltzmann machines (RBMs) or deep belief networks (DBNs), they do not play an important role in computer vision tasks compared to CNNs.

⁶A neuromorphic system is a hardware implementation of neural circuitry in silico that allows for fast and direct computation on neural networks.

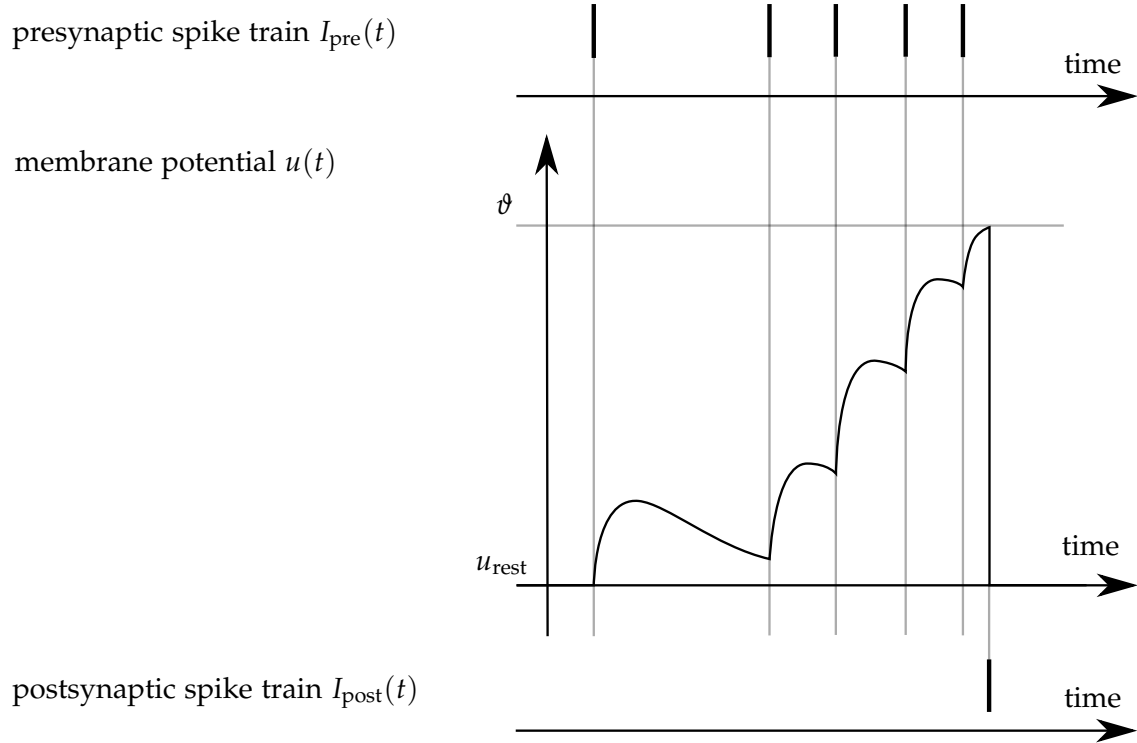


Figure 2.12: Presynaptic spike train, membrane potential and postsynaptic spike train of a LIF neuron. The three plots are synchronous in time. A presynaptic spike adds a charge to the membrane potential that decays over time. Once the membrane potential reaches the threshold ϑ , the neuron discharges, generating a postsynaptic spike and resetting to resting potential u_{rest} . Note the discontinuity whenever a pre- or postsynaptic neuron fires.

number of feature maps with large receptive fields, representing high-level features and shapes. With every pooling operation however, information on the location of features is lost. This means that CNNs cannot learn the precise spatial relations between a whole and its parts. For tasks like facial recognition however, part-whole relations play an important role [72].

SNNs on the other hand try to avoid the limitations imposed by the high computational costs and energy requirements of CNNs by mimicking the sparse encoding of information in spikes as observed in brains. SNNs running on neuromorphic hardware are extremely energy efficient and are able to operate in real time by design, independent of the size and topology of the network. Still, the computation on current neuromorphic circuitry is only a qualitative approximation of digitally simulated spiking neurons and even those do currently not approach state of the art CNN accuracy in computer vision tasks [73, 59]. Biological systems, especially the visual cortex of the primate brain, are not subject to these limitations [74, 75]. It therefore stands to reason that all of these challenges can be met and that findings from biology and neuroscience may hold a clue as to how to engineer solutions.

Recently an ANN architecture consisting of groups of neurons called *capsules* has been proposed as a possibly biologically plausible alternative for CNNs in computer vision systems. Capsule networks can be thought of as a visual attention mechanism, that encodes poses between features. Additionally, the activation functions and pooling operations found in CNNs are replaced by a sophisticated routing organism, that retains positional information. In theory this grants capsule networks a measure of viewpoint invariance as well as access to precise part-whole relations throughout the network. As capsule networks are still rate-coded and trained by supervised SGD with backpropagation, they fall somewhere inbetween CNNs and SNNs in regards to biological plausibility, but are possibly more expressive than CNNs. The theory behind capsules is explored in detail in the next chapter.

3 Capsule Network Architectures

A capsule is conceived as a group of neurons whose outputs encode different properties of an entity such as an object or part of an object in an image. The properties could represent instantiation parameters on the object’s appearance manifold like the precise pose, lighting and deformation or the probability that the entity is present. Layers within a capsule network consist of several capsules. Routing between the layers is determined by the votes from the capsules in the lower layer. Capsules can make a guess as to what they expect the higher layer’s capsules instantiation parameters to be and cast a vote. Votes are calculated by applying discriminatively learned transformations to a capsule’s neuron’s outputs. The routing algorithm will then determine, based on the lower capsules’ votes, which higher layer capsules to activate. The goal of capsules is to achieve invariance in the presence probability of the object they represent regarding the object’s instantiation parameters. A specific capsule architecture is largely defined by the nature of its routing algorithm and the dimension of its output as well as the order of the tensors used to compute the votes.

3.1 Transforming Auto-encoders

The first publication on capsules came from Hinton et al. in the form of transforming auto-encoders in 2011 [76]. These capsules encode the instantiation parameters of their internal representation together with the probability that the learned entity is present into a small output vector. Each capsule would learn to recognize a single visual entity over a limited subdomain of its appearance manifold. Ideally, the probability would be locally invariant as long as the instantiation parameters remain within the trained submanifold. This is in contrast to the instantiation parameters, which are equivariant. Meaning as the viewing conditions change, the parameters encoding the appearance change by an equal amount.

The transforming auto-encoders introduced by Hinton et al. in 2011 were already motivated by a desire to retain precise part-whole relationships, which are lost by the pooling operations in convolutional neural networks (cf. section 2.4). If a capsule can be trained to encode the pose of the visual entity it represents in its output vector, it is a simple matter to test whether two capsules A and B are in agreement about a higher level capsule C . To see this, suppose that the matrix $\mathbf{T}_{A/B}$ is the pose between the canonical entity as learned by the capsule A or B respectively and the current instantiation. Multiplying this with the part-whole transformation matrix $\mathbf{T}_{A/B \rightarrow C}$ will

result in a prediction for the pose of capsule C 's entity. Now, if

$$\mathbf{T}_A \mathbf{T}_{A \rightarrow C} \approx \mathbf{T}_B \mathbf{T}_{B \rightarrow C}, \quad (3.1)$$

then the capsules are a close match and activate capsule C . Hinton et al. argue that this naturally leads to precise part-whole relationships. If the entities of capsules A and B e.g. correspond to a mouth and a nose respectively, then they have to be in right spatial relationship to form a face if they activate the higher capsule C representing the whole face. The pose of capsule C can then be determined by the average of the lower capsules' votes.

$$\mathbf{T}_C = \sum_{i \in \{A, B\}} \mathbf{T}_i \mathbf{T}_{i \rightarrow C} \quad (3.2)$$

This means that except for the first layer of capsules, only the viewpoint invariant part-whole transformations have to be learned, while the entity poses are encoded in the neural activities at run time. The rest of the paper focuses on how to extract the poses from the pixel intensities for the first layer. This is done by training capsules on transformed image pairs and providing the capsules with non-visual access to the transformation. Hinton et al. justified this by pointing out that the human visual system is also provided with information about the eye-movement while saccades correspond to a translational transformation of the retinal image.

These first capsule networks actually consist of only one layer of capsules and do not include the routing mechanism of later implementations. Instead, they include two internal layers of logistic neurons, termed recognition and generation units respectively (cf. figure 3.1). The recognition units are trained to extract pose parameters and the probability that the visual entity is present directly from the provided image. A capsule's input can be constrained to a fixed receptive field with the capsules aligned in a grid across the image. The capsules then apply a transformation to the pose and provide the generation units with this information. The generation units on the other hand are trained to output the transformed part of the input image centered at the capsule's receptive field. All the neurons within a capsule are trained discriminatively using gradient descent with backpropagation. The training signal is provided by a form of reconstruction loss where the capsules have access to the transformation and the input and reconstruction target are the original and transformed images respectively. Transforming auto-encoders are able to learn poses in the form of 3×3 matrices and can be used to successfully apply affine transformations to input images once trained. However, they are limited to one layer of capsules and cannot be scaled effectively, as the capsules are fixed to a single position, compared to the shared weights of moving filters in convolutional layers. Also they are only able to learn properties, that can be controlled and provided in explicit non-visual form to the network.

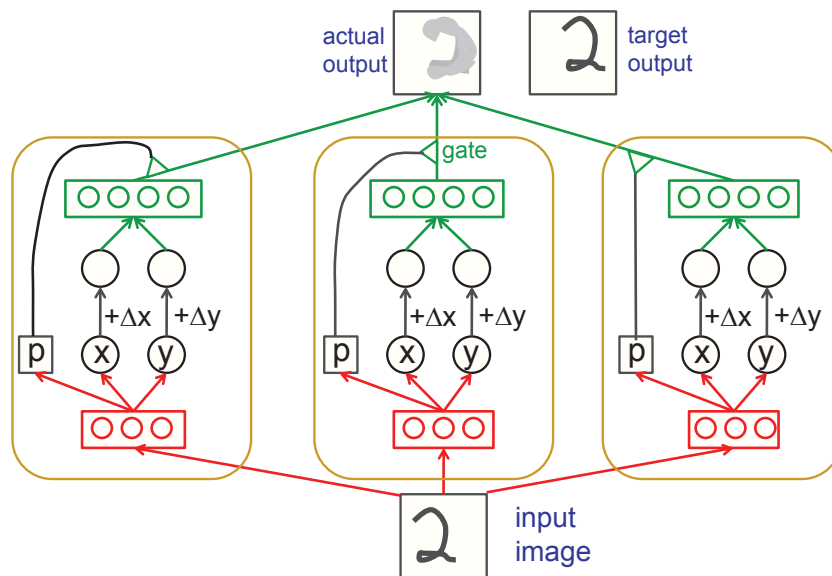


Figure 3.1: Capsules in a transforming auto-encoder [76]. The network consists of a single layer of capsules. Each capsule includes hidden recognition (in red) and generation units (in green). During training transformation parameters (in this case the translations Δx and Δy) as well as the transformed image are provided to the network. A capsule's output vector is further multiplied with its probability, thus muting capsules with a low probability.

3.2 Dynamic Routing Between Capsules

In 2017 Sabour et al. extended upon the transforming auto-encoders by suggesting a dynamic routing algorithm for capsules [77]. While the benefit of dynamic routing between capsule layers in regard to part-whole relationships was already elaborated upon by Hinton et al., Sabour et al. also motivated the mechanism biologically. They pointed out that human vision ignores irrelevant details by focusing attention on salient features in a sequence of fixations and only ever processes a fraction of the optic array at the highest resolution. They assumed that with a single fixation, groups of activated neurons build a tree structure within the fixed multi-layer visual system. The paper suggests to implement the nodes of this image parse tree as active capsules and have each capsule choose their parent node based on the routing algorithm. As with the capsules of transforming auto encoders, the capsule's neurons should learn to represent various instantiation parameters such as pose, deformation, velocity, albedo, hue, texture, etc. However, the probability of the instance's presence is encoded in the length of a vector of instantiation parameters. A nonlinearity is applied to the output vector in order to keep the length from exceeding 1, while keeping the orientation unchanged.

$$\mathbf{u}_j = \frac{\|\mathbf{s}_j\|^2}{1 + \|\mathbf{s}_j\|^2} \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|} \quad (3.3)$$

With \mathbf{u}_j the output vector of capsule j and \mathbf{s}_j the sum of its inputs. This *squashing* function will reduce short vectors to almost length zero and squash long vectors to a length slightly below 1. For every capsule j that is not a first layer capsule, the input \mathbf{s}_j is computed from the votes of the lower level capsules.

$$\mathbf{s}_j = \sum_i c_{ij} \hat{\mathbf{u}}_{j|i} \quad (3.4)$$

Where $\hat{\mathbf{u}}_{j|i}$ is the lower capsule i 's guess for the higher capsule j 's instantiation vector. Votes are calculated by applying a transformation matrix \mathbf{W}_{ij} on a capsule's output vector, as was suggested by Hinton et al. for transforming auto-encoders.

$$\hat{\mathbf{u}}_{j|i} = \mathbf{W}_{ij} \mathbf{u}_i \quad (3.5)$$

In this scheme, so far only the part-whole transformation matrices \mathbf{W}_{ij} have to be learned. The coupling coefficients c_{ij} for a lower capsule i sum to 1. They describe the influence a vote from capsule i has over capsule j and are iteratively fine-tuned by the routing algorithm during the network's forward pass. For the first iteration, the coupling coefficients are computed as a softmax function over the log prior coupling logits b_{ij} .

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})} \quad (3.6)$$

The b_{ij} can be trained discriminatively along with the other parameters. But it was found that the priors have little influence over the routing algorithm.

Using a vector as capsule output allows the use of a powerful dynamic routing-by-agreement algorithm. The agreement between capsule i and each possible parent j is calculated using the scalar product of the lower capsule's vote and the higher capsule's current output.

$$a_{ij} = \mathbf{u}_j^T \hat{\mathbf{u}}_{j|i} \quad (3.7)$$

If the scalar product is large, it induces top-down feedback, increasing the coupling coefficient between the two capsules by adding the agreement to the coupling logit b_{ij} (cf. algorithm 1). As the lower capsule's coupling coefficients are computed as a softmax of the logits, this will simultaneously decrease all the other couplings.

Algorithm 1 Dynamic routing-by-agreement

```

1: procedure ROUTING( $\hat{\mathbf{u}}_{j|i}, r, l$ )
2:   for capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $l + 1$  do
3:      $b_{ij} \leftarrow 0$ 
4:   for  $r$  iterations do
5:     for capsule  $i$  in layer  $l$  do
6:        $\mathbf{c}_i \leftarrow \text{softmax}(\mathbf{b}_i)$  ▷ computes equation 3.6
7:     for capsule  $j$  in layer  $l+1$  do
8:        $\mathbf{s}_j \leftarrow \sum_i c_{ij} \hat{\mathbf{u}}_{j|i}$ 
9:     for capsule  $j$  in layer  $l+1$  do
10:       $\mathbf{u}_j \leftarrow \text{squash}(\mathbf{s}_j)$  ▷ computes equation 3.3
11:    for capsule  $i$  in layer  $l$  and capsule  $j$  in layer  $l + 1$  do
12:       $b_{ij} \leftarrow b_{ij} + \mathbf{u}_j^T \hat{\mathbf{u}}_{j|i}$ 
  return  $\mathbf{u}_j$ 

```

Sabour et al. further improved upon transforming auto-encoders by developing convolutional capsule layers. This allows them to use learned features all across the image without losing whole-part relationships. They achieve this by having each capsule of a convolutional layer put out a local grid of vectors to each type of capsule¹ in the higher level, where every position on the grid and type of layer is assigned a different part-whole transformation matrix. While this is not strictly convolution, the net effect is the same. Capsules in a convolutional layer possess a receptive field that gets wider the higher up the layer is in the network's topology. For low level capsules, positional information is encoded in the capsule's position within the feature map. The higher up a convolutional capsule is, the smaller the feature map becomes and information on position gets encoded more and more in the instantiation parameter vector.

¹The use of different types of capsules within a convolutional layer is analogous to the use of different kernels or filters in CNNs (cf. section 2.2). The resulting map of capsule activities i.e. the lengths of the output vectors can be thought of as a feature map.

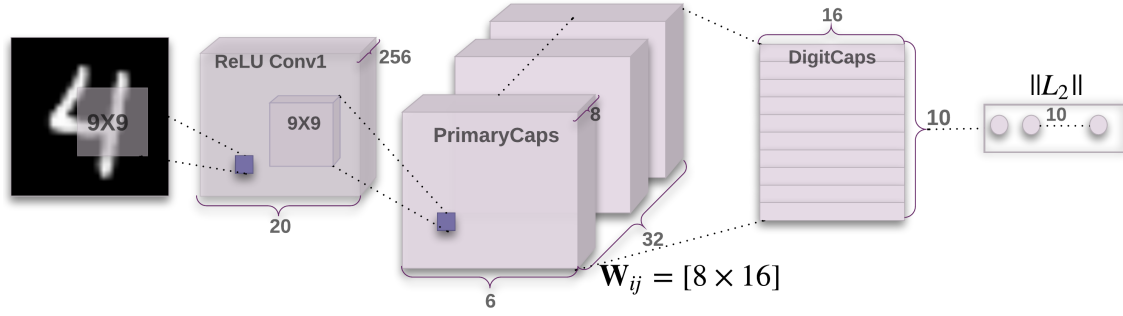


Figure 3.2: Vector capsule network with 3 layers [77]. Conv1 is a standard convolutional layer consisting of 256 9×9 kernels that learn basic features from pixel intensities, that are converted to 8D vector representations by the primary capsule layer. The primary capsule layer is the only convolutional capsule layer in the network and has 32 types of capsules of which each has a receptive field of the size 9×9 and a stride of 2. The class capsule layer (here referred to as DigitCaps) uses 16D vectors and 10 fully connected capsules. Note that the routing algorithm is applied only between the primary and the class capsule layer.

The primary capsules i.e. the first capsule layer computes its input vectors from the reshaped output of a standard convolutional layer. For the next capsule layers, the routing algorithm is applied. This is yet another improvement over transforming auto-encoders, as the network can be trained on images only without needing access to the applied transformations. The final capsule layer is fully connected and includes one capsule per class (cf. figure 3.2). *Margin loss* is used to compute the error signal based on the class capsules' output vector length.

$$L_k = T_k \max(0, m^+ - \|\mathbf{v}_k\|)^2 + \lambda(1 - T_k) \max(0, \|\mathbf{v}_k\| - m^-)^2 \quad (3.8)$$

Where $\mathbf{T} = (0, \dots, 1, \dots, 0)^T$ is a one-hot encoded vector of the true class labels and λ is used to down-weight the loss for absent classes so that the vectors don't become too short during initial training. In addition Sabour et al. used a reconstruction loss not unlike with transforming auto-encoders in order to encourage the class capsules to encode useful instantiation parameters for each class. This is done by masking out all but the correct class capsule during training and feeding its output vector to a 3-layer decoder (cf figure 3.3). The reconstruction loss is calculated as the euclidean distance between the decoder's output and the input image and scaled down so as not to dominate the margin loss. Sabour et al. were able to demonstrate that a shallow network with only two capsule layers can already achieve good results in object recognition tasks and that the class capsules encoded sensible instantiation parameters like the thickness of strokes on the MNIST dataset.

As the internal representation, the votes and the agreement are all vectors, this type of architecture will be referred to as *vector capsules* for the remainder of this thesis.

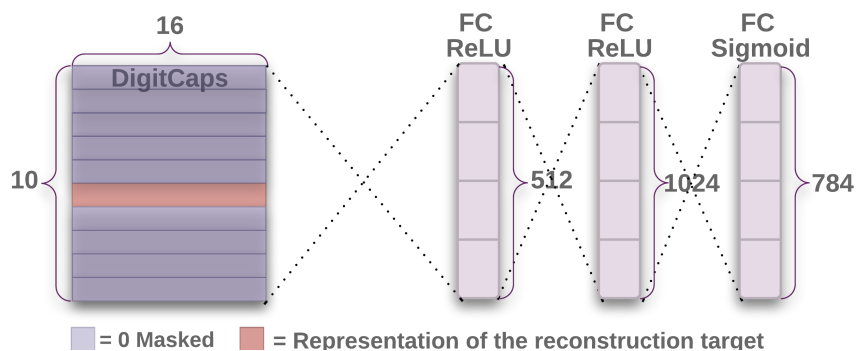


Figure 3.3: Decoder on a class vector capsule layer [77]. The output of the correct class capsule only is used as input for a 3-layer fully connected decoder in order to reconstruct the input image and encourage the class capsule to encode useful instantiation parameters.

3.3 Matrix Capsules With EM Routing

4 Experimental Setup

5 Results

6 Discussion

7 Conclusion

List of Figures

2.1	CIFAR-10 classes and sample images	4
2.2	Illustration of an artificial neuron with three input connections	5
2.3	Illustration of fully connected layers	6
2.4	Illustration of convolutional layers	8
2.5	Illustration of receptive field in 2D convolutional layer	8
2.6	Illustration of a pooling layer	9
2.7	Illustration of a typical CNN architecture	10
2.8	Computational graph for error propagation	12
2.9	Illustration of a neural network's error function	13
2.10	Electrical circuit of a leaky integrate and fire neuron	16
2.11	Typical STDP window function	17
2.12	Presynaptic spike train, membrane potential and postsynaptic spike train of a LIF neuron	19
3.1	Capsules in a transforming auto-encoders	23
3.2	Vector capsule network with 3 layers	26
3.3	Decoder on a class vector capsule layer	27

List of Algorithms

1 Dynamic routing-by-agreement 25

List of Tables

Bibliography

- [1] L. Paulun, A. Wendt, and N. K. Kasabov. “A retinotopic spiking neural network system for accurate recognition of moving objects using NeuCube and dynamic vision sensors.” In: *Frontiers in Computational Neuroscience* 12 (2018), p. 42.
- [2] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. “Gradient-based learning applied to document recognition.” In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [3] E. Falotico, L. Vannucci, A. Ambrosano, U. Albanese, S. Ulbrich, J. C. Vasquez Tieck, G. Hinkel, J. Kaiser, I. Peric, O. Denninger, et al. “Connecting artificial brains to robots in a comprehensive simulation framework: The neurorobotics platform.” In: *Frontiers in neurorobotics* 11 (2017), p. 2.
- [4] A. Knoll, F. Röhrbein, A. Kuhn, M. Akl, and K. Sharma. “Neurorobotics.” In: *Informatik-Spektrum* 40.2 (2017), pp. 161–164.
- [5] J. L. Krichmar. “Neurorobotics—A Thriving Community and a Promising Pathway Toward Intelligent Cognitive Robots.” In: *Frontiers in Neurorobotics* 12 (2018).
- [6] F. Rosenblatt. “The perceptron: a probabilistic model for information storage and organization in the brain.” In: *Psychological review* 65.6 (1958), p. 386.
- [7] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. *Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [8] A. Diba, V. Sharma, A. M. Pazandeh, H. Pirsiavash, and L. V. Gool. “Weakly Supervised Cascaded Convolutional Networks.” In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), pp. 5131–5139.
- [9] W. Ouyang, X. Zeng, X. Wang, S. Qiu, P. Luo, Y. Tian, H. Li, S. Yang, Z. Wang, H. Li, K. Wang, J. Yan, C. C. Loy, and X. Tang. “DeepID-Net: Object Detection with Deformable Part Based Convolutional Neural Networks.” In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39.7 (July 2017), pp. 1320–1334. issn: 0162-8828. doi: 10.1109/TPAMI.2016.2587642.
- [10] A. J. Schofield, I. D. Gilchrist, M. Bloj, A. Leonardis, and N. Bellotto. “Understanding images in biological and computer vision.” In: *Interface Focus* 8.4 (2018). issn: 2042-8898. doi: 10.1098/rsfs.2018.0027. eprint: <http://rsfs.royalsocietypublishing.org/content/8/4/20180027.full.pdf>.
- [11] A. Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. Tech. rep. 2009.

- [12] G. Cybenko. "Approximation by superpositions of a sigmoidal function." In: *Mathematics of Control, Signals and Systems* 2.4 (Dec. 1989), pp. 303–314. doi: <https://doi.org/10.1007/BF02551274>.
- [13] K. Hornik. "Approximation capabilities of multilayer feedforward networks." In: *Neural Networks* 4.2 (1991), pp. 251–257. issn: 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T).
- [14] Y.-L. Boureau, J. Ponce, and Y. Lecun. "A Theoretical Analysis of Feature Pooling in Visual Recognition." In: *ICML 2010 - Proceedings, 27th International Conference on Machine Learning*. Nov. 2010, pp. 111–118.
- [15] D. Scherer, A. Müller, and S. Behnke. "Evaluation of Pooling Operations in Convolutional Architectures for Object Recognition." In: *Artificial Neural Networks – ICANN 2010*. Ed. by K. Diamantaras, W. Duch, and L. S. Iliadis. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 92–101. isbn: 978-3-642-15825-4.
- [16] H. D. H. and W. T. N. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex." In: *The Journal of Physiology* 160.1 (), pp. 106–154. doi: 10.1113/jphysiol.1962.sp006837. eprint: <https://physoc.onlinelibrary.wiley.com/doi/pdf/10.1113/jphysiol.1962.sp006837>.
- [17] K. Fukushima. "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position." In: *Biological Cybernetics* 36.4 (Apr. 1980), pp. 193–202. issn: 1432-0770. doi: 10.1007/BF00344251.
- [18] Y. LeCun, Y. Bengio, and G. Hinton. "Deep learning." In: *nature* 521.7553 (2015), p. 436.
- [19] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. Chap. 9. eprint: <http://www.deeplearningbook.org/contents/convnets.html>.
- [20] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. "Is object localization for free? - Weakly-supervised learning with convolutional neural networks." In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2015, pp. 685–694. doi: 10.1109/CVPR.2015.7298668.
- [21] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. "Going deeper with convolutions." In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2015, pp. 1–9. doi: 10.1109/CVPR.2015.7298594.
- [22] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. "Densely Connected Convolutional Networks." In: *CVPR*. Vol. 1. 2. July 2017, p. 3.
- [23] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. "Inception-v4, inception-resnet and the impact of residual connections on learning." In: *AAAI*. Vol. 4. 2017, p. 12.

-
- [24] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. "Aggregated residual transformations for deep neural networks." In: *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*. IEEE. 2017, pp. 5987–5995.
 - [25] L. Wang, Y. Xiong, Z. Wang, and Y. Qiao. "Towards good practices for very deep two-stream convnets." In: *arXiv preprint arXiv:1507.02159* (2015).
 - [26] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting." In: *The Journal of Machine Learning Research* 15.1 (2014), pp. 1929–1958.
 - [27] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, and R. Fergus. "Regularization of Neural Networks using DropConnect." In: *Proceedings of the 30th International Conference on Machine Learning*. Ed. by S. Dasgupta and D. McAllester. Vol. 28. Proceedings of Machine Learning Research 3. Atlanta, Georgia, USA: PMLR, June 2013, pp. 1058–1066.
 - [28] A. Krogh and J. A. Hertz. "A simple weight decay can improve generalization." In: *Advances in neural information processing systems*. 1992, pp. 950–957.
 - [29] N. K. Treadgold and T. D. Gedeon. "Simulated annealing and weight decay in adaptive learning: The SARPROP algorithm." In: *IEEE Transactions on Neural Networks* 9.4 (1998), pp. 662–668.
 - [30] M. D. Zeiler. "ADADELTA: an adaptive learning rate method." In: *arXiv preprint arXiv:1212.5701* (2012).
 - [31] J. Duchi, E. Hazan, and Y. Singer. "Adaptive subgradient methods for online learning and stochastic optimization." In: *Journal of Machine Learning Research* 12.Jul (2011), pp. 2121–2159.
 - [32] D. P. Kingma and J. Ba. "Adam: A method for stochastic optimization." In: *arXiv preprint arXiv:1412.6980* (2014).
 - [33] B. T. Polyak and A. B. Juditsky. "Acceleration of stochastic approximation by averaging." In: *SIAM Journal on Control and Optimization* 30.4 (1992), pp. 838–855.
 - [34] A. Graves. "Generating sequences with recurrent neural networks." In: *arXiv preprint arXiv:1308.0850* (2013).
 - [35] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. "On the importance of initialization and momentum in deep learning." In: *International conference on machine learning*. 2013, pp. 1139–1147.
 - [36] I. Loshchilov and F. Hutter. "Sgdr: Stochastic gradient descent with warm restarts." In: *arXiv preprint arXiv:1608.03983* (2016).
 - [37] J. L. Ba, J. R. Kiros, and G. E. Hinton. "Layer normalization." In: *arXiv preprint arXiv:1607.06450* (2016).
 - [38] S. Ioffe and C. Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." In: *arXiv preprint arXiv:1502.03167* (2015).

- [39] T. Salimans and D. P. Kingma. "Weight normalization: A simple reparameterization to accelerate training of deep neural networks." In: *Advances in Neural Information Processing Systems*. 2016, pp. 901–909.
- [40] D. Worrall and G. Brostow. "CubeNet: Equivariance to 3D Rotation and Translation." In: *arXiv preprint arXiv:1804.04458* (2018).
- [41] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow. "Harmonic networks: Deep translation and rotation equivariance." In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2. 2017.
- [42] U. Schmidt and S. Roth. "Learning rotation-aware features: From invariant priors to equivariant descriptors." In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE. 2012, pp. 2050–2057.
- [43] T. Cohen and M. Welling. "Group equivariant convolutional networks." In: *International conference on machine learning*. 2016, pp. 2990–2999.
- [44] F. E. Bloom. "Chapter 1 - Fundamentals of Neuroscience." In: *Fundamental Neuroscience (Fourth Edition)*. Ed. by L. R. Squire, D. Berg, F. E. Bloom, S. du Lac, A. Ghosh, and N. C. Spitzer. Fourth Edition. San Diego: Academic Press, 2013, pp. 3–13. ISBN: 978-0-12-385870-2. DOI: <https://doi.org/10.1016/B978-0-12-385870-2.00001-9>.
- [45] A. H. Bell, L. Pessoa, R. B. Tootell, and L. G. Ungerleider. "Chapter 44 - Visual Perception of Objects." In: *Fundamental Neuroscience (Fourth Edition)*. Ed. by L. R. Squire, D. Berg, F. E. Bloom, S. du Lac, A. Ghosh, and N. C. Spitzer. Fourth Edition. San Diego: Academic Press, 2013, pp. 947–968. ISBN: 978-0-12-385870-2. DOI: <https://doi.org/10.1016/B978-0-12-385870-2.00044-5>.
- [46] J. R. Manns and E. A. Buffalo. "Chapter 48 - Learning and Memory: Brain Systems." In: *Fundamental Neuroscience (Fourth Edition)*. Ed. by L. R. Squire, D. Berg, F. E. Bloom, S. du Lac, A. Ghosh, and N. C. Spitzer. Fourth Edition. San Diego: Academic Press, 2013, pp. 1029–1051. ISBN: 978-0-12-385870-2. DOI: <https://doi.org/10.1016/B978-0-12-385870-2.00048-2>.
- [47] C. Koch. "Chapter 51 - The Neuroscience of Consciousness." In: *Fundamental Neuroscience (Fourth Edition)*. Ed. by L. R. Squire, D. Berg, F. E. Bloom, S. du Lac, A. Ghosh, and N. C. Spitzer. Fourth Edition. San Diego: Academic Press, 2013, pp. 1091–1103. ISBN: 978-0-12-385870-2. DOI: <https://doi.org/10.1016/B978-0-12-385870-2.00051-2>.
- [48] J. H. Byrne. "Chapter 47 - Learning and Memory: Basic Mechanisms." In: *Fundamental Neuroscience (Fourth Edition)*. Ed. by L. R. Squire, D. Berg, F. E. Bloom, S. du Lac, A. Ghosh, and N. C. Spitzer. Fourth Edition. San Diego: Academic Press, 2013, pp. 1009–1027. ISBN: 978-0-12-385870-2. DOI: <https://doi.org/10.1016/B978-0-12-385870-2.00047-0>.
- [49] A. Knoll and M.-O. Gewaltig. "Neurorobotics: a strategic pillar of the Human Brain Project." In: *Science Robotics* (2016).

- [50] D. Drubach. *The Brain Explained*. Prentice-Hall, 2000.
- [51] A. L. Hodgkin and A. F. Huxley. "A quantitative description of membrane current and its application to conduction and excitation in nerve." In: *The Journal of physiology* 117.4 (1952), pp. 500–544.
- [52] A. N. Burkitt. "A review of the integrate-and-fire neuron model: I. Homogeneous synaptic input." In: *Biological cybernetics* 95.1 (2006), pp. 1–19.
- [53] S. M. Bohte, H. La Poutré, and J. N. Kok. "Unsupervised clustering with spiking neurons by sparse temporal coding and multilayer RBF networks." In: *IEEE Transactions on neural networks* 13.2 (2002), pp. 426–435.
- [54] S. M. Bohte. "The evidence for neural information processing with precise spike-times: A survey." In: *Natural Computing* 3.2 (2004), pp. 195–206.
- [55] J. J. Hopfield. "Pattern recognition computation using action potential timing for stimulus representation." In: *Nature* 376.6535 (1995), p. 33.
- [56] W. Gerstner and W. M. Kistler. *Spiking neuron models: Single neurons, populations, plasticity*. Cambridge university press, 2002.
- [57] W. Maass. "Networks of spiking neurons: the third generation of neural network models." In: *Neural networks* 10.9 (1997), pp. 1659–1671.
- [58] F. Rieke and D. Warland. *Spikes: exploring the neural code*. MIT press, 1999.
- [59] P. U. Diehl and M. Cook. "Unsupervised learning of digit recognition using spike-timing-dependent plasticity." In: *Frontiers in computational neuroscience* 9 (2015), p. 99.
- [60] D. Hebb. "The organization of behavior." In: *Wiley and Sons, New York, NY, USA* (1949).
- [61] S. Song, K. D. Miller, and L. F. Abbott. "Competitive Hebbian learning through spike-timing-dependent synaptic plasticity." In: *Nature neuroscience* 3.9 (2000), p. 919.
- [62] Z. Bing, C. Meschede, F. Röhrbein, K. Huang, and A. C. Knoll. "A Survey of Robotics Control Based on Learning-Inspired Spiking Neural Networks." In: *Frontiers in neurorobotics* 12 (2018), p. 35.
- [63] S. Bartunov, A. Santoro, B. A. Richards, G. E. Hinton, and T. Lillicrap. "Assessing the scalability of biologically-motivated deep learning algorithms and architectures." In: *arXiv preprint arXiv:1807.04587* (2018).
- [64] J. H. Lee, T. Delbruck, and M. Pfeiffer. "Training Deep Spiking Neural Networks Using Backpropagation." In: *Frontiers in Neuroscience* 10 (2016), p. 508. ISSN: 1662-453X. DOI: 10.3389/fnins.2016.00508.
- [65] P. Panda and K. Roy. "Unsupervised regenerative learning of hierarchical features in spiking deep networks for object recognition." In: *arXiv preprint arXiv:1602.01510* (2016).

- [66] M. Hopkins, G. Pineda-García, P. A. Bogdan, and S. B. Furber. "Spiking neural networks for computer vision." In: *Interface Focus* 8.4 (2018), p. 20180007.
- [67] S. C. Wong, A. Gatt, V. Stamatescu, and M. D. McDonnell. "Understanding data augmentation for classification: when to warp?" In: *arXiv preprint arXiv:1609.08764* (2016).
- [68] L. Perez and J. Wang. "The effectiveness of data augmentation in image classification using deep learning." In: *arXiv preprint arXiv:1712.04621* (2017).
- [69] C.-S. Lee, M.-H. Wang, S.-J. Yen, T.-H. Wei, I.-C. Wu, P.-C. Chou, C.-H. Chou, M.-W. Wang, and T.-H. Yan. "Human vs. Computer Go: Review and Prospect [Discussion Forum]." In: *Comp. Intell. Mag.* 11.3 (Aug. 2016), pp. 67–72. ISSN: 1556-603X. DOI: 10.1109/MCI.2016.2572559.
- [70] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. "Mastering the game of Go with deep neural networks and tree search." In: *nature* 529.7587 (2016), p. 484.
- [71] Y. Dong, Y. Wang, Z. Lin, and T. Watanabe. "High performance and low latency mapping for neural network into network on chip architecture." In: *ASIC, 2009. ASICON'09. IEEE 8th International Conference on*. IEEE. 2009, pp. 891–894.
- [72] J. W. Tanaka and D. Simonyi. "The "parts and wholes" of face recognition: A review of the literature." In: *The Quarterly Journal of Experimental Psychology* 69.10 (2016), pp. 1876–1889.
- [73] G. Indiveri, B. Linares-Barranco, T. J. Hamilton, A. Van Schaik, R. Etienne-Cummings, T. Delbruck, S.-C. Liu, P. Dudek, P. Häfliger, S. Renaud, et al. "Neuromorphic silicon neuron circuits." In: *Frontiers in neuroscience* 5 (2011), p. 73.
- [74] T. Stone, M. Mangan, A. Wystrach, and B. Webb. "Rotation invariant visual processing for spatial memory in insects." In: *Interface Focus* 8.4 (2018), p. 20180010.
- [75] J. B. Isbister, A. Eguchi, N. Ahmad, J. M. Galeazzi, M. J. Buckley, and S. Stringer. "A new approach to solving the feature-binding problem in primate vision." In: *Interface Focus* 8.4 (2018), p. 20180021.
- [76] G. E. Hinton, A. Krizhevsky, and S. D. Wang. "Transforming auto-encoders." In: *International Conference on Artificial Neural Networks*. Springer. 2011, pp. 44–51.
- [77] S. Sabour, N. Frosst, and G. E. Hinton. "Dynamic routing between capsules." In: *Advances in Neural Information Processing Systems*. 2017, pp. 3856–3866.