

Neue Einsatzmöglichkeit von Hardwarebeschleunigern für nachhaltigere KI-Modelle: Entwicklung und Evaluation der Boltzmann Maschinen auf einem physikinspirierten Hardwarebeschleuniger

Bachelorarbeit

submitted on February 22, 2024

Fakultät Wirtschaft und Gesundheit

Studiengang Wirtschaftsinformatik

Kurs WWI2021F

von

SIMON SPITZER

Betreuer in der Ausbildungsstätte:

DHBW Stuttgart:

⟨ Hewlett Packard GmbH ⟩
⟨ Dr. Fabian Böhm ⟩
⟨ Research Scientist at Hewlett Packard Labs ⟩

⟨ Prof. Dr., Kai Holzweißig ⟩
⟨ der/des wissenschaftlichen Betreuerin/Prüferin ⟩

Unterschrift der Betreuerin/des Betreuers

Contents

List of abbreviations	IV
List of figures	V
List of tables	VI
1 Einleitung	1
1.1 Motivation	1
1.2 Problemstellung	1
1.3 Zielsetzung(ohne gneau Metriken nennen, generell halten)	3
1.4 Forschungsmethodik	3
1.5 Aufbau der Arbeit	4
2 Aktueller Stand der Forschung und Praxis (generell auch wiedergeben von aktuell existierenden Lösungsmustern)	5
2.1 Ressourcenverbrauch bei KI-Modellen	5
2.1.1 Ressourcenverbrauch bei KI-Modellen	5
2.2 Deep Neural Network - Boltzmann Maschinen (Erstmal DNN erklären generell)	5
2.2.1 Energy-based models	5
2.2.2 concept of Boltzmann Maschines	7
2.2.3 Training of Restriced Boltzmann Machines	9
2.2.4 Aktuelle Probleme mit RBM/BM	12
2.3 Hardwarebeschleuniger	13
2.3.1 Aktuelle Ansätze im Bereich KI und weitere Lösungen	13
2.3.2 ISING Maschine/ Physikinspirierter Hardwarebeschleuniger	13
2.4 Memristor Hopfield Network	13
2.4.1 Memristor	13
2.4.2 Hopfield Network	13
2.4.3 Crossbar	13
2.4.4 Output Hopfield Networtk	13
2.4.5 Noisy HNN	13
3 Zielspezifikation und Darlegung der Forschungsmethodik	14
3.1 Zielspezifikation (genauer als in Einleitung, Metriken erwähnen, Erfolg meiner Methode bewerten, Welcher Teil der Forschungsfrage wird beantwortet?)	14
3.2 Design Science Research	14
3.3 Zielsetzung(ohne gneau Metriken nennen, generell halten)	14
3.4 Laborexperiment für die Umsetzung	14
4 Implementierung/Laborexperiment der Simulator Pipeline	15
4.1 Zielsetzung und Forschungsmethodik	15
4.2 Aufbau der Simulator Pipeline	15
4.3 KI-Bibliothek Scikit-Learn	15
5 Evaluation der BM auf dem physikinspiriertem Hardwarebeschleuniger	16
5.1 Zielsetzung und Forschungsmethodik	16

5.1.1	Prediction Accuracy	16
5.1.2	Troughput (Samples/Sec)	16
5.1.3	Energieverbrauch (Energy/Operation)	16
5.2	Vergleichen mit anderen Hardwarebeschleuniger, FPGA, GPU oder CPU aus der Literatur	16
6	Kritische Reflexion und Ausblick	17
6.1	Evaluation der Erkenntnisse in Bezug auf die Zielsetzung der Arbeit	17
6.2	Kritische Reflexion der Ergebnisse und Methodik	17
6.3	Zielsetzung(ohne gneaue Metriken nennen, generell halten)	17
6.4	Ergebnisextraction für Theorie und Praxis (evtl. mit 6.4 Zusammenlegen)	17
6.5	Ausblick	17
	Appendix	18
	List of references	20

List of abbreviations

Ein Abkürzungsverzeichnis ist optional. Das Paket `acronym` kann weit mehr, als hier gezeigt.¹ Beachten Sie allerdings, dass Sie die Einträge selbst in sortierter Reihenfolge angeben müssen.

BM	Boltzmann Maschine
RBM	Restricted Boltzmann Maschine
DNN	Deep Neural Network
EBM	Energy Based Model
MCMC	Markov chain Monte Carlo

Ergänzende Bemerkung: Eine im Text verwendete Abkürzung sollte bei ihrer ersten Verwendung erklärt werden. Falls Sie sich nicht selbst darum kümmern möchten, kann das das Paket `acronym` übernehmen und auch automatisch Links zum Abkürzungsverzeichnis hinzufügen. Dazu ist an allen Stellen, an denen die Abkürzung vorkommt, `\ac{ITIL}` zu schreiben.

Das Ergebnis sieht wie folgt aus:

- erstmalige Verwendung von `\ac{ITIL}` ergibt: **ITIL!** (**ITIL!**),
- weitere Verwendung von `\ac{ITIL}` ergibt: **ITIL!**

Wo benötigt, kann man mit dem Befehl `\acl{ITIL}` wieder die Langfassung ausgeben lassen: **ITIL!**.

Falls man die Abkürzungen durchgängig so handhabt, kann man durch Paket-Optionen (in `_dhbw_preamble.tex`) erreichen, dass im Abkürzungsverzeichnis nur die tatsächlich verwendeten Quellen aufgeführt werden (Option: `printonlyused`) und zu jedem Eintrag die Seite der ersten Verwendung angegeben wird (Option: `withpage`).

¹siehe <http://ctan.org/pkg/acronym>

List of Figures

1	Figure of a logistic sigmoid function Restricted Boltzmann Maschine (RBM)	6
2	figure of a general Boltzmann Machine	8
3	Figure of a RBM	9
4	Figure of a logistic sigmoid function RBM	10
5	Mal wieder das DHBW-Logo.	19

List of Tables

1 Einleitung

1.1 Motivation

1.2 Problemstellung

In der Forschung und Entwicklung von Generativen KI-Modellen rückt die Rechengeschwindigkeit und Energieeffizienz zunehmend in den Fokus² Die Autor*innen von Open AI bestätigen, dass die Wachstumsrate von Machine-Learning-Modellen die Effizienzrate von Computerchips schon längst übertroffen hat. So verdoppeln sich jede 3-4 Monate der Rechenbedarf dieser Modelle jedoch verdoppeln sich nach Moore's Law die Leistung der Computerchips nur jede 2 Jahre.³ Angesichts der Probleme des steigenden Energieverbrauchs von Rechenzentren und den damit verbundenen Treibhausgasemissionen dieser, ist die Suche nach effizienteren Lösungen essenziell für die Zukunft. Weltweit steigern Datenzentren ihren Energieverbrauch jährlich um 20-40%, wodurch sie 2022 etwa 1,3% des globalen Energieverbrauchs und 1% der energiebedingten globalen Treibhausgasemissionen verursacht haben.⁴ Jedoch ist hier nicht zu erkennen, wie groß dabei der KI-Anteil zur Grundgesamtheit beiträgt.

Ein bereits bekannter Ansatz ist die Benutzung von KI-Beschleunigern basierend auf ASICs (Application-specific Integrated Circuits) - also Schaltungen, die anwendungsspezifisch verwendet werden, wie zum Beispiel Google TPUs (Tensor Processing Unit).⁵ Dies ist auch sinnvoll, da die Verwendung von Mehrzweckmodellen für diskriminierende Aufgaben im Vergleich zu aufgabenspezifischen Modellen energieintensiver ist.⁶ Ein alternatives vielversprechendes Konzept in der Forschung ist die Verwendung von physikinspirierten Hardwarebeschleunigern, die primär bei Optimierungsalgorithmen eingesetzt werden aufgrund ihrer Fähigkeit Probleme schneller und effizienter als GPUs lösen zu können.⁷ Ein skalierbarer physikinspirierter Hardwarebeschleuniger (auch Ising-Maschine genannt), der die Leistung bestehender Standard-Digitalrechner übertrifft, könnte einen großen Einfluss auf praktische Anwendungen für eine Vielzahl von Optimierungsproblemen haben.⁸

Solche physikinspirierten Hardwarebeschleuniger bieten durch ihre besondere Berechnungsweise Potenzial für eine effizientere Verarbeitung von rechenintensiven Aufgaben. Konkret wird die Beschleunigung, anders als es bei digitalen Computern der Fall ist, durch die Berechnung rechenintensiver Aufgaben mit analogen Signalen erreicht. Die Implementierung auf dedizierter Hard-

²Vgl. Luccioni/Jernite/Strubell 2023, p. 1

³Vgl. Dario Amodei/Danny Hernandez 2024, p. 1

⁴Vgl. Hintemann/Hinterholzer 2022, p. 1

⁵Vgl. Wittpahl 2019, p. 39

⁶Vgl. Luccioni/Jernite/Strubell 2023, p. 5

⁷Vgl. Mohseni/McMahon/Byrnes 2022, p. 1

⁸Vgl. Mohseni/McMahon/Byrnes 2022, p. 1

ware bietet darüber hinaus die Möglichkeit, die Parallelisierung von digitalen Hardwarebeschleunigern und analogem Rechnen auszunutzen.⁹

Interessanterweise zeigen die Energiefunktionen von Hardwarebeschleunigern, die in Ising-Maschinen verwendet werden, große Parallelen zu denen in Boltzmann Maschinen, trotz ihrer unterschiedlichen Anwendungen, daher liegt es nahe, dass Ising Maschinen auch für KI gut funktionieren.¹⁰ Ising-Maschinen zielen darauf ab, ihre Energie zu minimieren, wobei sie Energie als eine paarweise Interaktion von binären Variablen „Spins“ definieren.¹¹ Boltzmann Maschinen hingegen sind energiebasierte neuronale Netzwerke, die Klassifizierungen durchführen, indem sie jeder Konfiguration der Variablen eine skalare Energie zuordnen. Die Netzwerkenergie zu minimieren ist hierbei vergleichbar mit der Lösung des Optimierungsproblems.¹² Aktuelle Probleme mit Boltzmann-Maschinen umfassen die hohe Komplexität und Anforderungen an die All-to-All-Kommunikation zwischen Verarbeitungseinheiten, was ihre Implementierung auf herkömmlichen digitalen Computern ineffizient macht, sowie eine inhärent langsame Konvergenz in bestimmten Prozessen wie Simulated Annealing.¹³ Diese Herausforderungen erschweren das Training und die Anwendung von Boltzmann-Maschinen insbesondere für große Datenmengen und komplexe Optimierungsaufgaben.¹⁴ Nichtsdestotrotz impliziert die Ähnlichkeit der beiden, dass Ising-Maschinen in der Lage sein könnten, dieses spezielle KI-Modell, energieeffizienter und mit höherer Rechengeschwindigkeit auszuführen. Aktuell existieren nur wenige Konzepte eine Implementierung von Boltzmann Maschinen auf Ising-Maschinen zu erreichen. Das Paper der Autoren Mahdi Nazm BojnordiEngin und Engin Ipek ist hier ein vielversprechender Ansatz, jedoch konnte nicht gezeigt werden, wie es auf einem richtigen Beschleunigerchip funktionieren würde.

Vor diesem Hintergrund ergeben sich folgende zentrale Forschungsfragen:

1. Können Boltzmann Maschinen auf physikinspiertenHardwarebeschleunigern durch analoge Rauschinjektion effizient implementiert werden?
 - Wie ist die Genauigkeit des KI-Modells im Hardwarebeschleuniger? Metrik: Prediction Accuracy
 - ergleichen mit anderen Hardwarebeschleuniger, FPGA, GPU oderCPU aus der Literatur (gute und schlechte) in Bezug auf Energieeffizienz und Rechengeschwindigkeit – Metriken: Troughput(Samples/Sec), Energieverbrauch (Energy/Operation)

Daher gilt es zu testen, ob dieses generative KI-Modell mit Ising Maschinen kompatibel ist und ob diese Lösung effizient ist oder nicht.

⁹Vgl. Mohseni/McMahon/Byrnes 2022, p. 4

¹⁰Vgl. Cai et al. 2019, p. 10

¹¹Vgl. Wang/Roychowdhury 2017, p. 1

¹²Vgl. Nazm Bojnordi/Ipek 2016, p. 2

¹³Vgl. Nazm Bojnordi/Ipek 2016, p. 1

¹⁴Vgl. Nazm Bojnordi/Ipek 2016, p. 2

1.3 Zielsetzung(ohne gneaue Metriken nennen, generell halten)

Das primäre Ziel dieser Bachelorarbeit ist die Erforschung und Erweiterung eines bestehenden physikinspirierten Hardwarebeschleunigers (ISING Maschine) zur Implementierung und Evaluation von Boltzmann Maschinen, einem energiebasierten KI-Modell. Dabei sollen die aufgestellten Forschungsfragen beantwortet werden.

Hierzu ist es zu Beginn nötig eine Simulator Pipeline zu konstruieren mit der Boltzmann Maschinen auf dem Hardwarebeschleuniger übersetzt werden. Die Simulator Pipeline besteht dabei aus einer bestehender KI-Bibliothek und bestehenden Hardwarebeschleuniger, die miteinander verbunden werden. Mit der Simulator Pipeline soll gezeigt werden, dass der Hardwaresimulator die Boltzmann Maschinen umsetzen kann. Aus der Simulator Pipeline heraus werden die Aktivierungswahrscheinlichkeiten der einzelnen Neuronen auf der simulierten Hardware gemessen und bei Erfolg bis zu einem vollständigen Neuronalen Netzwerk erweitert. Finaler Schritt ist, dass der Hardwarebeschleuniger für Training und Interferenz genutzt werden kann und dabei vergleichbar mit herkömmlichen ML Libraries ist. Diese Phase umfasst die sorgfältige Anpassung und möglicherweise Erweiterung des bestehenden Beschleunigers, um die spezifischen Anforderungen der Boltzmann Maschinen zu erfüllen.

Wenn die Simulator Pipeline validiert werden kann, wird ein Workload auf ein Standarddataset zur Handschrifterkennung getestet. Dabei werden die Prediction Accuracy, Troughput (Samples/Sec) und der Energieverbrauch (Energy/Operation) der Boltzmann Maschinen auf dem ISING Hardwarebeschleuniger untersucht und dadurch die aufgestellten Forschungsfragen beantwortet.

1.4 Forschungsmethodik

Design Science Research

1. **Problemorientierung:** DSR fokussiert auf die Lösung praktischer Probleme, wie die Forschung zur Steigerung der Effizienz und Rechengeschwindigkeit in KI-Modellen.
2. **Artefakt Entwicklung:** Zentral in DSR ist die Entwicklung innovativer Artefakte. Die Arbeit zielt darauf ab, ein solches Artefakt in Form des physikinspirierten Hardwarebeschleunigers weiterzuentwickeln und für KI-Modelle einzusetzen.
3. **Iterative Evaluation:** Durch die iterative Vorgehensweise in DSR kann die Ausarbeitung der Lösung fortlaufend verbessert und angepasst werden, was für die Entwicklung und Optimierung von KI-Systemen entscheidend ist (ebenfalls das Konzept).
4. **Beitrag zur Wissensbasis und Praxisrelevanz:** DSR unterstützt die Generierung neuer Erkenntnisse und stellt sicher, dass Forschungsergebnisse sowohl theoretisch fundiert

als auch praktisch anwendbar sind, was mit den Zielen Ihres Projekts im Einklang steht. Untermethodik könnte hierbei eine Simulation sein. Variabel, je nach Verlauf der Forschung.

1.5 Aufbau der Arbeit

2 Aktueller Stand der Forschung und Praxis (generell auch wiedergeben von aktuell existierenden Lösungsmustern)

2.1 Ressourcenverbrauch bei KI-Modellen

2.1.1 Ressourcenverbrauch bei KI-Modellen

Nachhaltigkeit

Stromverbrauch

Rechenleistung begrenzt, KI-Modelle wachsen schneller als verfügbare Leistung

2.2 Deep Neural Network - Boltzmann Maschinen (Erstmal DNN erklären generell)

(Erklären von Deep Neurol Network und Neurol Network) -> Anwendungsbereiche Spracherkennung, Image recognition Solche Deep Neural Networks, sind sehr ressourceneffizient und möglicher Forschungsbereich für Nutzung des Hardwarebeschleunigers Idee dabei, die repräsentationspower von energybased model höher als bei LLMs, mit weniger Neuronen besser als LLMs

Some regression tasks within computer vision in Deep Neural Network (DNN) include object detection, medical image registration, head- and body-pose estimation, age estimation and visual tracking.¹⁵

->RBMs einführen und sagen warum Training vereinfacht ist

2.2.1 Energy-based models

An Energy Based Model (EBM) is a type of statistical model where the likelihood of a particular state is determined by an energy function.¹⁶ Since 1982, those models have been continuously emerging in the machine learning field when J.J. Hopfield introduced the Hopfield Network.¹⁷ Current developments include their use in reinforcement learning, potential replacements for discriminators in generative adversarial networks and for quantum EBMs.¹⁸ The underlying idea

¹⁵Vgl. Gustafsson et al. 2020, pp. 325–326

¹⁶Vgl. Huembeli et al. 2022, p. 2

¹⁷Vgl. Hopfield 1982

¹⁸Vgl. Verdon et al. 2019, p. 1; Vgl. Du/Lin/Mordatch 2021, p. 1

behind EBMs is to establish a probabilistic physical system that is able to learn and memorize patterns but most importantly generalize it.¹⁹ Especially, it involves learning an energy function $E_\theta(x) \in \mathbb{R}$ and assigning the low energy to observed data x_i and high energy to other values x .²⁰

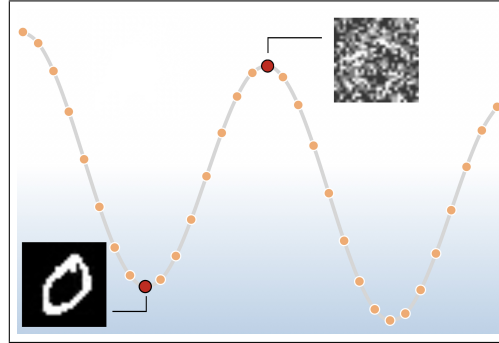


Abb. 1: Figure of a logistic sigmoid function RBM

In this figure a simplified energy landscape is shown where the local minima correspond to states that encode an MNIST digit.²¹ It is visible that observed data settles in the local minimum of the energy landscape, in this case a clear 0. On the other hand close to the local maxima of the energy landscape the 0 is only barely recognizable and therefore got a higher energy value assigned to it. The assumption of the underlying distribution function $P(x)$ is equal to the solution of the optimization problem:

$$P(x) = \frac{1}{Z} \exp \left(-\frac{E(x)}{T} \right), \quad (2.1)$$

where Z is given by the partition function to ensure that the density function normalizes to a total probability of 1 and T is interpreted as the temperature.²² As a result the behavior of a EBM is determined by 2.2. The aim of the training is to match the real data P_{data} as closely as possible with the internal model P_{model} . A practical method to achieve this goal is to use the KL divergence. KL divergence is a mathematical equation that helps to measure how close the predictions are by comparing the model's learned distribution to the true distribution of the data:

$$G = \sum_x P^+(x) \ln \left(\frac{P^+(x)}{P^-(x)} \right) \quad (2.2)$$

Here, $P^+(x)$ is the probability when the states are determined by a data input from the environment, while $P^-(x)$ represents the internal network running freely, also referred to as “dreaming”.²³ To optimise the KL divergence, in this case G , the energy is adjusted, whereby data is assigned to low energy states (according to 2.1) and the training data receives high energy and therefore

¹⁹Vgl. Huembeli et al. 2022, p. 2

²⁰Vgl. Gustafsson et al. 2020, p. 330

²¹Vgl. Huembeli et al. 2022, p. 6

²²Vgl. Huembeli et al. 2022, pp. 2–3

²³Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, pp. 154–155

high probabilities.²⁴ To complete the section the “partition function”, Z , used in 2.1 is given by summing over all possible pairs of visible and hidden vectors:

$$Z = \sum_x \exp \left(-\frac{E(x)}{T} \right) \quad (2.3)$$

As a side note that is worth mentioning is, that using the maximum likelihood estimator for Z is impractical due to the requirement of summing over all possible states, which leads to an exponential increase in the number of states for larger systems.²⁵

2.2.2 concept of Boltzmann Machines

A Boltzmann Maschine (BM) is a specific symmetrical EBM consisting of binary neurons $\{0, 1\}$.²⁶ The neurons of the network can be split into two functional groups, a set of visible neurons and a set of hidden neurons.²⁷ Therefore, the BM is a two-layer model with a visible layer (“v”) and a hidden layer (“h”).²⁸ The visible layer is the interface between the network and the environment. It receives data inputs during training and sets the state of a neuron to either $\{0, 1\}$ which represents activated or not activated. On the other hand, the hidden units are not connected to the environment and can be used to “explain” underlying constraints in the internal model of input vectors and they cannot be represented by pairwise constraints.²⁹ The connection between the individual neurons is referred to as bidirectional, as each neuron communicates with each other in both directions.³⁰

As early as 1985, one of the founding fathers of artificial intelligence, “Geoffrey Hinton”, was aware that an BM is able to learn its underlying features by looking at data from a domain and developing a generative internal model.³¹ In the next step, it is possible to generate examples with the same probability distribution as the input data examples shown. In the following figure 2, a general BM is depicted, where the upper layer embodies a vector of stochastic binary ‘hidden’ features, while the lower layer embodies a vector of stochastic binary ‘visible’ variables.³²

²⁴Vgl. Zhai et al. 2016, pp. 2–3

²⁵Vgl. Zhai et al. 2016, pp. 2–3

²⁶Vgl. Amari/Kurata/Nagaoka 1992, p. 260

²⁷Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 154

²⁸Vgl. Salakhutdinov/Hinton, G. 2009, p. 448

²⁹Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 154

³⁰Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 149

³¹Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 148

³²Vgl. Salakhutdinov/Hinton, G. 2009, p. 449

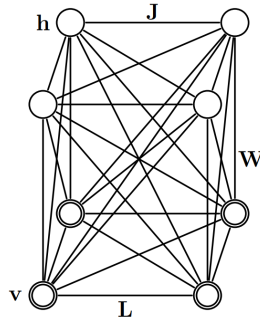


Abb. 2: figure of a general Boltzmann Machine

The model contains a set of visible units $v \in \{0, 1\}$, and a set of hidden units $h \in \{0, 1\}$ (see Fig. 1). The energy function of the BM with the states $\{v, h\}$ is defined as:

$$E(v, h; \theta) = -\frac{1}{2}v^T L v - \frac{1}{2}h^T J h - v^T W h, \quad (2.4)$$

where $\theta = \{W, L, J\}$ are the model parameters.³³ W, L, J represent visible-to-hidden, visible-to-visible and hidden-to-hidden weights. The individual neurons can be made to try to minimize the global energy by setting the right assumptions.³⁴ Entering a particular input to the machine, the system will find the minimum energy configuration that can illustrate the input.³⁵ A simple method to find a local energy minimum is to switch into whichever of the two states of a neuron hold the lower energy given the current state of the other neurons.³⁶ The exact reason for this is the following: “If all the connection strengths are symmetrical, which is typically the case for constraint satisfaction problems, each unit can compute its effect on the total energy from information that is locally available.”³⁷ By inserting the function 2.4 into the earlier introduced KL-divergence 2.2 and doing gradient descend the following learning rule to update the weights and biases results³⁸:

$$\Delta w_{ij} = \epsilon(\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}) \quad (2.5)$$

The network can now update the weights “W” that exist between the neurons through the training rule based on the observations that served as input and modified by the learning rate ϵ .³⁹

Performing exact maximum likelihood learning in this model is intractable because exact computation of the data predictions and the model predictions takes a time that is exponential in

³³Vgl. Salakhutdinov/Hinton, G. 2009, p. 448

³⁴Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 150

³⁵Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 150

³⁶Vgl. Fahlman/Hinton, G./Sejnowski, T. 1983, p. 110

³⁷Fahlman/Hinton, G./Sejnowski, T. 1983, p. 110

³⁸Vgl. Hinton, G. E. 2012b, p. 5

³⁹Vgl. Barra et al. 2012, pp. 1–2

the number of hidden units.⁴⁰ When the number of hidden units is large compared to the number of visible units it is impossible to achieve a perfect model because of the totally connected network and the resulting 2^n possibilities.⁴¹ This leads back to the briefly mentioned constraint of equation 2.3, that is needed to calculate an activation probability of a neuron, which is required to update a weight in the training process shown in 2.5.

A specific example to demonstrate why it is intractable to calculate a activation of a BM is the following. The BM has 80 visible nodes and 120 hidden nodes and therefore the possibilities of states of neurons are 2^{200} , which is 1.61×10^{60} . To put this in perspective the total atoms that exist on earth are only estimated to be around 1.33×10^{50} .⁴² That means even if it would be possible to store one information per atom it would just not be enough.

2.2.3 Training of Restricted Boltzmann Machines

As a solution for the training problem Hinton and Sejnowski proposed Gibbs sampling as an algorithm to approximate both expectations.⁴³ Furthermore, the intralayer connections of the model got removed and the result is the so called RBM. To transform an BM into a RBM the diagonal elements L and J introduced earlier, are set to 0 and as a result the well-known model of a RBM establishes shown in fig.2.⁴⁴

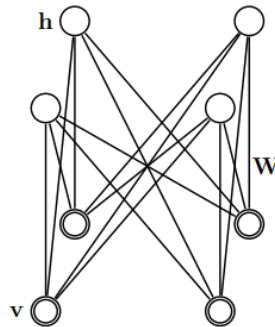


Abb. 3: Figure of a RBM

What can be recognized that no more visible-to-visible and hidden-to-hidden connections can be found in the model. The configuration of the visible and hidden units (v, h) therefore has also an updated energy function (Hopfield, 1982) given by:

$$E(v, h) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij}, \quad (2.6)$$

⁴⁰Vgl. Salakhutdinov/Hinton, G. 2009, p. 449

⁴¹Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, p. 154

⁴²Vgl. Helmenstine 2022, p. 478-480; Vgl. Schlamming 2014, p. 1

⁴³Vgl. Ackley/Hinton, G. E./Sejnowski, T. J. 1985, pp. 158-165

⁴⁴Vgl. Salakhutdinov/Hinton, G. 2009, p. 449

where v_i, h_j are the binary states of a visible unit i and hidden unit j , a_i, b_j are their biases and w_{ij} is the weight between them.⁴⁵ Despite, compared to the fully connected BM the RBM is less complex the advantages of training surpasses the loss in expressivity.⁴⁶ The RBM has recently been drawing attention in the machine learning community because of its adaption and extension for various tasks such as representational learning, document modeling, image recognition and for serving as foundational components for deep networks including Deep Boltzmann Machine, Deep Belief Networks and hybrid models with CNNs.⁴⁷ The training of the model can be split up into the following steps:

1. Forward Pass (positive phase)

During the forward pass using the Gibbs Sampling method, the visible units are set to a completely random state. Next up the hidden units are computed. The computation of the hidden units involves calculating their activation probabilities and performing an actual sampling with their calculated activation probabilities. With the RBM it is now easy to get an analytical calculated unbiased sample of $(\mathbf{v}_i \mathbf{h}_j)_{data}$.⁴⁸ Given a input data out of the training images, v , the binary state, j , of each hidden unit, h_j , is set to 1 with following probability:

$$p(h_j = 1|\mathbf{v}) = \sigma(b_j + \sum_i v_i w_{ij}), \quad (2.7)$$

where $\sigma(x)$ is the logistic sigmoid function is then an unbiased sample. The sigmoid function is defined as $\sigma(x) = \frac{1}{1+\exp(-x)}$ and shown in figure 4:

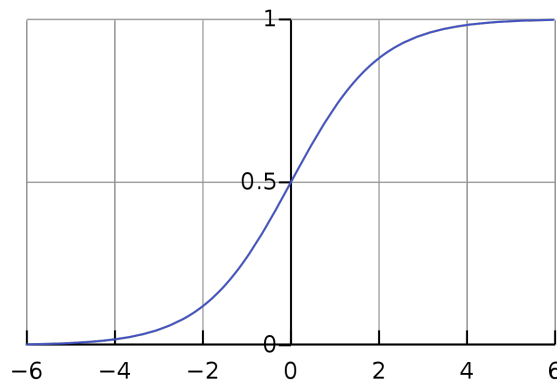


Abb. 4: Figure of a logistic sigmoid function RBM

The result is a set of probabilities that reflects how likely it is for each hidden unit to be on or off given the input data.⁴⁹ The sampling part of the positive phase uses the just calculated activation probability of each hidden unit and performs a random experiment with it. Afterwards,

⁴⁵Vgl. Hinton, G. E. 2012a, pp. 3–4

⁴⁶Vgl. Huembeli et al. 2022, p. 4

⁴⁷Vgl. Zhang et al. 2018, p. 1186

⁴⁸Vgl. Hinton, G. E. 2012b, p. 5

⁴⁹Vgl. Huembeli et al. 2022, p. 6

the hidden unit is either activated or not activated and the training process continues with the new state of the hidden units (activated or not activated).

2. Reconstruction (negative phase)

In this phase, the sampled hidden states are used to reconstruct the visible units. This is essentially a prediction of the input, which is calculated with following probability:⁵⁰

$$p(v_i = 1|\mathbf{h}) = \sigma(a_i + \sum_j h_j w_{ij}) \quad (2.8)$$

The sampling part of the negative phase uses the just calculated activation probability of each visible unit and performs a random experiment, like in the positive phase. Now, the result is a prediction of the input in the visible nodes. Afterwards, a half forward pass is made to calculate the activation probability of a hidden unit again based on the activated or not activated visible units.

3. Updating the weights

Meanwhile, all the requirements to update the weights are satisfied and can be used within the equation 2.5. The delta that results is summed to the current weight and therefore the internal model gets closer to predicting the observed data. Therefore, one training iteration consisting of 1 Forward Pass, 1 Reconstruction and 0.5 Forward Pass again is accomplished. Repeating this training steps N times for a suitable chosen N the model learns better, since more steps of alternating Gibbs sampling were performed.⁵¹

Markov-Chain-Monte-Carlo-Verfahren

Contrastive Divergence: Contrastive divergence is a special Gibbs Sampling training method developed by Geoffrey Hinton for the efficient training of Boltzmann machines especially RBMs.⁵² In traditional Gibbs sampling would have to generate a long chain of samples, until independent samples are obtained from the observed data distribution of the model.⁵³ The samples are needed for each iteration of the gradient ascent on the log-likelihood resulting in large computational costs.⁵⁴ To solve this issue contrastive divergence minimizes an approximation of the Kullback-Leibler divergence between the empirical distribution of the training data and the distribution generated by the model.⁵⁵ The way to achieve this is by initializing the Markov chain with the samples from the data distribution.⁵⁶ The outcome has been shown to heavily increase the training time while only adding a small bias.⁵⁷ What this means is initializing the visible units

⁵⁰Vgl. Hinton, G. E. 2012b, p. 6

⁵¹Vgl. Huembeli et al. 2022, p. 6

⁵²Vgl. Hinton, G. E. 2012b, pp. 4–5

⁵³Vgl. Huembeli et al. 2022, pp. 5–6

⁵⁴Vgl. Upadhyaya/Sastry 2019, pp. 7–8

⁵⁵Vgl. Mocanu et al. 2016, p. 246

⁵⁶Vgl. Upadhyaya/Sastry 2019, pp. 7–8

⁵⁷Vgl. Larochelle/Bengio 2008, p. 537

with a real data input for example a MNIST sample and starting the proposed steps with the underlying states. Often the process can be stopped after only sampling a very small number of steps (often only one).⁵⁸

Metropolis Hastings:

2.2.4 Current Problems with BMs and RBMs

One general problem that occurs in the learning process of a BM is that it is both time-consuming and difficult.⁵⁹ This is because sampling from an undirected graphical model is not straightforward and therefore RBMs can make use of Markov chain Monte Carlo (MCMC) proposed methods like Contrastive Divergence and Metropolis Hastings.⁶⁰ In addition to that the selection of hyperparameters can be difficult since for the training a practical model the a large hyper-parameter space needs to be explored.⁶¹ Especially finding the right size of the hidden layer, the learning rate and number of training iteration but also the method for calculating activation probabilities(Contrastive Divergence, Metropolis Hastings, etc.) can be seen as art. Furthermore, training can become unstable due to the system's low temperature, which impacts the training negatively.⁶² A lower temperature reduces the system's possibility to explore the energy landscape thoroughly, leading to the false selection of local minima instead of finding the global minimum.

⁵⁸Vgl. Larochelle/Mandel, et al. 2012, p. 646

⁵⁹Vgl. Fischer/Igel 2012, pp. 1-2

⁶⁰Vgl. Fischer/Igel 2012, p. 2

⁶¹Vgl. Larochelle/Bengio 2008, p. 536

⁶²Vgl. Huembeli et al. 2022, pp. 3-4

2.3 Hardwarebeschleuniger

2.3.1 Aktuelle Ansätze im Bereich KI und weitere Lösungen

Asics

Quantencomputing

2.3.2 ISING Maschine/ Physikinspirierter Hardwarebeschleuniger

Konzept (mit Energiefunktion), Probleme der Digitalrechner bzw. Unterschied zu Digitalrechner

Aktuelle Anwendung

Potentielle Einsatzgebiete für KI-Modelle

Parallelen Energiefunktion BM und ISING Maschine

2.4 Memristor Hopfield Network

2.4.1 Memristor

2.4.2 Hopfield Network

2.4.3 Crossbar

2.4.4 Output Hopfield Network

2.4.5 Noisy HNN

3 Zielspezifikation und Darlegung der Forschungsmethodik

3.1 Zielspezifikation (genauer als in Einleitung, Metriken erwähnen, Erfolg meiner Methode bewerten, Welcher Teil der Forschungsfrage wird beantwortet?)

3.2 Design Science Research

3.3 Zielsetzung(ohne genaue Metriken nennen, generell halten)

3.4 Laborexperiment für die Umsetzung

4 Implementierung/Laborexperiment der Simulator Pipeline

4.1 Zielsetzung und Forschungsmethodik

4.2 Aufbau der Simulator Pipeline

4.3 KI-Bibliothek Scikit-Learn

5 Evaluation der BM auf dem physikinspiriertem Hardwarebeschleuniger

5.1 Zielsetzung und Forschungsmethodik

5.1.1 Prediction Accuracy

5.1.2 Troughput (Samples/Sec)

5.1.3 Energieverbrauch (Energy/Operation)

5.2 Vergleichen mit anderen Hardwarebeschleuniger, FPGA, GPU oder CPU aus der Literatur

6 Kritische Reflexion und Ausblick

6.1 Evaluation der Erkenntnisse in Bezug auf die Zielsetzung der Arbeit

6.2 Kritische Reflexion der Ergebnisse und Methodik

6.3 Zielsetzung(ohne genaue Metriken nennen, generell halten)

6.4 Ergebnisextraktion für Theorie und Praxis (evtl. mit 6.4 Zusammenlegen)

6.5 Ausblick

Appendix

List of appendices

Anhang 1	So funktioniert's	19
Anhang 1/1	Wieder mal eine Abbildung	19

Appendix 1: So funktioniert's

Um den Anforderungen der Zitierrichtlinien nachzukommen, wird das Paket `tocloft` verwendet. Jeder Anhang wird mit dem (neu definierten) Befehl `\anhang{Bezeichnung}` begonnen, der insbesondere dafür sorgt, dass ein Eintrag im Anhangsverzeichnis erzeugt wird. Manchmal ist es wünschenswert, auch einen Anhang noch weiter zu unterteilen. Hierfür wurde der Befehl `\anhangteil{Bezeichnung}` definiert.

In Anhang 1/1 finden Sie eine bekannte Abbildung und etwas Source Code in ??.

Anhang 1/1: Wieder mal eine Abbildung



Abb. 5: Mal wieder das DHBW-Logo.

List of references

- Ackley, D. H./Hinton, G. E./Sejnowski, T. J. (1985):** A Learning Algorithm for Boltzmann Machines. In: *Cognitive Science* 9.1, pp. 147–169. ISSN: 0364-0213. DOI: 10.1016/S0364-0213(85)80012-4. URL: <https://www.sciencedirect.com/science/article/pii/S0364021385800124> (retrieval: 02/16/2024).
- Amari, S./Kurata, K./Nagaoka, H. (1992):** Information Geometry of Boltzmann Machines. In: *IEEE Transactions on Neural Networks* 3.2, pp. 260–271. ISSN: 1941-0093. DOI: 10.1109/72.125867. URL: <https://ieeexplore.ieee.org/abstract/document/125867> (retrieval: 02/16/2024).
- Barra, A./Bernacchia, A./Santucci, E./Contucci, P. (2012):** On the Equivalence of Hopfield Networks and Boltzmann Machines. In: *Neural Networks* 34, pp. 1–9. ISSN: 0893-6080. DOI: 10.1016/j.neunet.2012.06.003. URL: <https://www.sciencedirect.com/science/article/pii/S0893608012001608> (retrieval: 02/16/2024).
- Cai, F./Kumar, S./Van Vaerenbergh, T./Liu, R./Li, C./Yu, S./Xia, Q./Yang, J. J./Beusoleil, R./Lu, W./Strachan, J. P. (2019):** Harnessing Intrinsic Noise in Memristor Hopfield Neural Networks for Combinatorial Optimization. DOI: 10.48550/arXiv.1903.11194. arXiv: 1903.11194 [cs]. URL: <http://arxiv.org/abs/1903.11194> (retrieval: 02/15/2024). preprint.
- Dario Amodei/Danny Hernandez (2024):** AI and Compute. URL: <https://openai.com/research/ai-and-compute> (retrieval: 02/15/2024).
- Du, Y./Lin, T./Mordatch, I. (2021):** Model Based Planning with Energy Based Models. DOI: 10.48550/arXiv.1909.06878. arXiv: 1909.06878 [cs, stat]. URL: <http://arxiv.org/abs/1909.06878> (retrieval: 02/19/2024). preprint.
- Fahlman, S./Hinton, G./Sejnowski, T. (1983):** Massively Parallel Architectures for AI: NETL, Thistle, and Boltzmann Machines., p. 113. 109 pp.
- Fischer, A./Igel, C. (2012):** An Introduction to Restricted Boltzmann Machines. In: *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*. Ed. by Luis Alvarez/Marta Mejail/Luis Gomez/Julio Jacobo. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, pp. 14–36. ISBN: 978-3-642-33275-3. DOI: 10.1007/978-3-642-33275-3_2.
- Gustafsson, F. K./Danelljan, M./Bhat, G./Schön, T. B. (2020):** Energy-Based Models for Deep Probabilistic Regression. In: *Computer Vision – ECCV 2020*. Ed. by Andrea Vedaldi/Horst Bischof/Thomas Brox/Jan-Michael Frahm. Lecture Notes in Computer Science. Cham: Springer International Publishing, pp. 325–343. ISBN: 978-3-030-58565-5. DOI: 10.1007/978-3-030-58565-5_20.
- Helmenstine, A. (2022):** How Many Atoms Are in the World? Science Notes and Projects. URL: <https://sciencenotes.org/how-many-atoms-are-in-the-world/> (retrieval: 02/21/2024).
- Hintemann, R./Hinterholzer, S. (2022):** Data Centers 2021: Data Center Boom in Germany Continues - Cloud Computing Drives the Growth of the Data Center Industry and Its Energy Consumption. DOI: 10.13140/RG.2.2.31826.43207.

- Hinton, G. E. (2012a):** ‘A Practical Guide to Training Restricted Boltzmann Machines’. In: *Neural Networks: Tricks of the Trade: Second Edition*. Ed. by Grégoire Montavon/Geneviève B. Orr/Klaus-Robert Müller. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, pp. 599–619. ISBN: 978-3-642-35289-8. DOI: 10.1007/978-3-642-35289-8_32. URL: https://doi.org/10.1007/978-3-642-35289-8_32 (retrieval: 02/15/2024).
- **(2012b):** ‘A Practical Guide to Training Restricted Boltzmann Machines’. In: *Neural Networks: Tricks of the Trade*. Ed. by Grégoire Montavon/Geneviève B. Orr/Klaus-Robert Müller. Vol. 7700. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 599–619. ISBN: 978-3-642-35288-1 978-3-642-35289-8. DOI: 10.1007/978-3-642-35289-8_32. URL: http://link.springer.com/10.1007/978-3-642-35289-8_32 (retrieval: 02/15/2024).
- Hopfield, J. J. (1982):** Neural Networks and Physical Systems with Emergent Collective Computational Abilities. In: *Proceedings of the National Academy of Sciences* 79.8, pp. 2554–2558. DOI: 10.1073/pnas.79.8.2554. URL: <https://www.pnas.org/doi/10.1073/pnas.79.8.2554> (retrieval: 02/19/2024).
- Huembeli, P./Arrazola, J. M./Killoran, N./Mohseni, M./Witteck, P. (2022):** The Physics of Energy-Based Models. In: *Quantum Machine Intelligence* 4.1, p. 1. ISSN: 2524-4914. DOI: 10.1007/s42484-021-00057-7. URL: <https://doi.org/10.1007/s42484-021-00057-7> (retrieval: 02/19/2024).
- Larochelle, H./Bengio, Y. (2008):** Classification Using Discriminative Restricted Boltzmann Machines. In: *Proceedings of the 25th International Conference on Machine Learning*. ICML ’08. New York, NY, USA: Association for Computing Machinery, pp. 536–543. ISBN: 978-1-60558-205-4. DOI: 10.1145/1390156.1390224. URL: <https://dl.acm.org/doi/10.1145/1390156.1390224> (retrieval: 02/22/2024).
- Larochelle, H./Mandel, M./Pascanu, R./Bengio, Y. (2012):** Learning Algorithms for the Classification Restricted Boltzmann Machine. In: *The Journal of Machine Learning Research* 13, pp. 643–669.
- Luccioni, A. S./Jernite, Y./Strubell, E. (2023):** Power Hungry Processing: Watts Driving the Cost of AI Deployment? DOI: 10.48550/arXiv.2311.16863. arXiv: 2311.16863 [cs]. URL: <http://arxiv.org/abs/2311.16863> (retrieval: 02/15/2024). preprint.
- Mocanu, D. C./Mocanu, E./Nguyen, P. H./Gibescu, M./Liotta, A. (2016):** A Topological Insight into Restricted Boltzmann Machines. In: *Machine Learning* 104.2, pp. 243–270. ISSN: 1573-0565. DOI: 10.1007/s10994-016-5570-z. URL: <https://doi.org/10.1007/s10994-016-5570-z> (retrieval: 02/22/2024).
- Mohseni, N./McMahon, P. L./Byrnes, T. (2022):** Ising Machines as Hardware Solvers of Combinatorial Optimization Problems. DOI: 10.48550/arXiv.2204.00276. arXiv: 2204.00276 [physics, physics:quant-ph]. URL: <http://arxiv.org/abs/2204.00276> (retrieval: 02/15/2024). preprint.
- Nazm Bojnordi, M./Ipek, E. (2016):** Memristive Boltzmann Machine: A Hardware Accelerator for Combinatorial Optimization and Deep Learning, p. 13. 1 p. DOI: 10.1109/HPCA.2016.7446049.

- Salakhutdinov, R./Hinton, G. (2009):** Deep Boltzmann Machines. In: *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*. Artificial Intelligence and Statistics. PMLR, pp. 448–455. URL: <https://proceedings.mlr.press/v5/salakhutdinov09a.html> (retrieval: 02/16/2024).
- Schlamming, S. (2014):** A Cool Way to Measure Big G. In: *Nature* 510.7506 (7506), pp. 478–480. ISSN: 1476-4687. DOI: 10.1038/nature13507. URL: <https://www.nature.com/articles/nature13507> (retrieval: 02/21/2024).
- Upadhy, V./Sastry, P. (2019):** An Overview of Restricted Boltzmann Machines. In: *Journal of the Indian Institute of Science* 99. DOI: 10.1007/s41745-019-0102-z.
- Verdon, G./Marks, J./Nanda, S./Leichenauer, S./Hidary, J. (2019):** Quantum Hamiltonian-Based Models and the Variational Quantum Thermalizer Algorithm. arXiv: 1910.02071 [quant-ph]. URL: <http://arxiv.org/abs/1910.02071> (retrieval: 02/19/2024). preprint.
- Wang, T./Roychowdhury, J. (2017):** Oscillator-Based Ising Machine. DOI: 10.48550/arXiv.1709.08102. arXiv: 1709.08102 [physics]. URL: <http://arxiv.org/abs/1709.08102> (retrieval: 02/15/2024). preprint.
- Wittpahl, V., ed. (2019):** Künstliche Intelligenz: Technologie | Anwendung | Gesellschaft. Berlin, Heidelberg: Springer. ISBN: 978-3-662-58041-7 978-3-662-58042-4. DOI: 10.1007/978-3-662-58042-4. URL: <http://link.springer.com/10.1007/978-3-662-58042-4> (retrieval: 02/15/2024).
- Zhai, S./Cheng, Y./Lu, W./Zhang, Z. (2016):** Deep Structured Energy Based Models for Anomaly Detection. In: *Proceedings of The 33rd International Conference on Machine Learning*. International Conference on Machine Learning. PMLR, pp. 1100–1109. URL: <https://proceedings.mlr.press/v48/zhai16.html> (retrieval: 02/19/2024).
- Zhang, N./Ding, S./Zhang, J./Xue, Y. (2018):** An Overview on Restricted Boltzmann Machines. In: *Neurocomputing* 275, pp. 1186–1199. ISSN: 0925-2312. DOI: 10.1016/j.neucom.2017.09.065. URL: <https://www.sciencedirect.com/science/article/pii/S0925231217315849> (retrieval: 02/15/2024).

Erklärung

Ich versichere hiermit, dass ich die vorliegende Arbeit mit dem Thema: *Mein Titel* selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

(Ort, Datum)

(Unterschrift)