# Neuromorphic hardware for sustainable AI data centers

Bernhard Vogginger*ˣⁱ, Amirhossein Rostami*‖, Vaibhav Jain§, Sirine Arfa*, Andreas Hantsch**††, David Kappel¶,
Michael Schäfer†‡, Ulrike Faltings†, Hector A. Gonzalez*‖‡‡, Chen Liu*, Christian Mayr*‖, Wolfgang Maaß§ˣ

*Chair of Highly-Parallel VLSI-Systems and Neuro-Microelectronics, Technische Universität Dresden, Germany
†SHS - Stahl-Holding-Saar GmbH & Co. KGaA, Germany
‡KTH Royal Institute of Technology, Department of Materials Science and Engineering, Sweden
§German Research Center for Artificial Intelligence (DFKI), Saarbrücken, Germany
¶Institute for Neural Computation, Ruhr University Bochum, Germany
‖ScaDS.AI Dresden/Leipzig, Germany
**eco2050 Institut für Nachhaltigkeit – Institute for Sustainability GmbH, Nürnberg, Germany
††Hantsch Sustainability Consulting, Dresden, Germany
‡‡SpiNNcloud Systems GmbH, Dresden, Germany
ˣSaarland University, Saarbrücken, Germany
ˣⁱEmail: bernhard.vogginger@tu-dresden.de

*Abstract*—As humans advance toward a higher level of artificial intelligence, it is always at the cost of escalating computational resource consumption, which requires developing novel solutions to meet the exponential growth of AI computing demand. Neuromorphic hardware takes inspiration from how the brain processes information and promises energy-efficient computing of AI workloads. Despite its potential, neuromorphic hardware has not found its way into commercial AI data centers. In this article, we try to analyze the underlying reasons for this and derive requirements and guidelines to promote neuromorphic systems for efficient and sustainable cloud computing: We first review currently available neuromorphic hardware systems and collect examples where neuromorphic solutions excel conventional AI processing on CPUs and GPUs. Next, we identify applications, models and algorithms which are commonly deployed in AI data centers as further directions for neuromorphic algorithms research. Last, we derive requirements and best practices for the hardware and software integration of neuromorphic systems into data centers. With this article, we hope to increase awareness of the challenges of integrating neuromorphic hardware into data centers and to guide the community to enable sustainable and energy-efficient AI at scale.

*Index Terms*—neuromorphic hardware, cloud computing, artificial intelligence, data center, sustainable computing

## I. INTRODUCTION

### A. Motivation

Data centers, serving as the hub for computers and equipment necessary to manage and store vast amounts of data, play a crucial role in deploying and maintaining AI systems, especially given the exponential growth in demand for AI computing models [1]. However, advances in traditional computers (CPU, GPU, TPU) have not kept pace with this growing demand. There is thus an urgent need for innovative solutions across hardware, software, and algorithms to ensure efficient, high-throughput, and sustainable AI in data centers [2].

Among potential candidates, neuromorphic computing inspired by the human brain is emerging as a promising approach to address these challenges with the feature of energy-efficient parallel processing [3], [4]. Neuromorphic computing aims to design and build computer systems including hardware and software that can perform cognitive tasks more efficiently by emulating how neurons and synapses work in the brain. Neuromorphic systems incorporate the concept of synaptic plasticity, allowing synapses between spiking neurons to change and adapt based on the input patterns. This work addresses the integration of neuromorphic hardware into data centers for sustainable AI.

### B. Sustainable AI

Data centers have a tremendous energy demand. Whilst the estimates vary over the order of a magnitude, their median global electricity demand was 300 TWh/a in 2020 and almost tripled by 2030 [5].

For a long period, just the energy demand as part of the environmental impact was in focus [1], [6]. However, the term *"sustainable AI"* describes the creation and application of AI technologies that prioritize long-term viability, social responsibility, and reducing environmental impact. As these applications spread across various industries, there's a rising

awareness of the necessity to address ethical and environmental issues related to AI development and implementation.

Following the Environment, Social, and Governance (ESG) structure of financial ratings, we highlight the following impact points for sustainable AI:

- **Environment:** 1) Materials, water, land, refrigerants, and energy (through greenhouse gases) needed during the server and data center life cycles - including embodied carbon dioxide; 2) Computing demand for AI training and inference computations
- **Social:** 1) Gathering (and labeling) of unbiased, racist-free training data embracing human rights; 2) Transparency and explainability of algorithms; 3) Potential positive impact on society by well-designed AI applications
- **Governance:** 1) Data security; 2) Data and service availability with shifting context, data sets, and model features; 3) Virtualisation and load balancing for operation optimization; 4) Economic feasibility

We would like to highlight that the embodied carbon dioxide in IT hardware can be more than half of the total carbon footprint of this hardware's life cycle [7]. Yet, in this article, we focus on reducing the energy demand for AI processing in data centers by deploying neuromorphic computing hardware and algorithms.

### C. Related work

This work addresses a closely related issue to the task explored in [8], which involves integrating diverse emerging hardware systems, including neuromorphic systems (NC), into a unified computational environment. The authors of [9] emphasize the significance of computational environments, covering conventional digital computing (DC) systems based on the von Neumann architecture and synchronous logical processing, alongside traditional distributed computing. They define NC as event-based systems with a distinct interface, wherein the structure and function either emulate or simulate the neuronal dynamics of brains, particularly somas, and occasionally synapses, dendrites, and axons, typically represented in the form of Spiking Neural Networks (SNNs). The study highlights the interconnection between the two types of hardware, DC and NC, establishing a microservice-based conceptual framework for integrating neuromorphic systems, featuring a neuromorphic-system proxy.

In contrast to the approach in [9], we are currently placing a greater emphasis on promoting the adoption of neuromorphic systems for typical AI tasks within data centers, with the aim of optimizing efficiency in cloud computing. This shift in focus directs our attention away from the generation and decoding aspects and principles in hardware associated with event processing.

Recently, significant progress has been made in neuromorphic computing, particularly in the realm of SNNs. In [4] the authors highlighted the challenges that must be overcome within this field to fully leverage the potential of efficient AI computing. The demands placed on SNN accelerators have witnessed notable transformations, especially in the design

of large-scale systems capable of effectively leveraging the essential features of SNN algorithms. The authors emphasized the design principles of neuromorphic hardware architectures, drawing inspiration from two core tenets of SNNs: (i) event-driven sparse computations and (ii) efficient and parallel matrix operations. A comparison between these neuromorphic hardware architectures and the standard architectures for Artificial Neural Networks (ANNs) revealed the proficiency of the standard architectures in matrix operations but their shortfall in exploiting the temporal sparsity inherent in SNNs.

Similar to this work, [10] discusses the potential of neuromorphic hardware for energy-efficient and green AI computing. The approach is showcased with the Novena chip from Singapore. The complete eco-system for deploying the chip is discussed, including algorithms, software, middleware, and system integration. Instead, our work analyses the current neuromorphic computing landscape as well as the potential and challenges for the integration of NC into mainstream AI data centers.

### D. Outline

This paper considers key participants in the neuromorphic systems field, including both commercial entities and major academic contributors, in the context of AI data centers and broader solutions. We present a detailed analysis of these systems and highlight the importance of hardware and software integration. The paper showcases the need for a holistic perspective by comparing neuromorphic to conventional approaches for AI processing.

Furthermore, it identifies common AI tasks in data centers and serves as a valuable reference for researchers investigating the development of neuromorphic algorithms. Fundamentally, the overall goal of this initiative is to improve the sustainability of data centers.

## II. NEUROMORPHIC HARDWARE PLAYERS

The neuromorphic field currently has a wide diversity of hardware platforms implemented and deployed at different scales aiming to cover applications that range from low-power embedded sensing intelligence [11], [12] to green cloud services [13], [14]. Such diversity can be observed via the different board levels in Fig. 1, which presents neuromorphic chips assembled in application boards with PCIe or Ethernet, and in server boards including high-speed links to assemble complex infrastructures to maintain real-time requirements. This section provides a broad overview of some of those neuromorphic platforms with special emphasis on those reaching data center levels.

### A. Methodology

As the term "neuromorphic" is used for various kinds of hardware methodologies, we briefly provide a taxonomy of "neuromorphic systems" in this paper. We consider *spike-based and event-based* hardware systems in all kinds of implementation approaches (analog, mixed-signal, digital, or

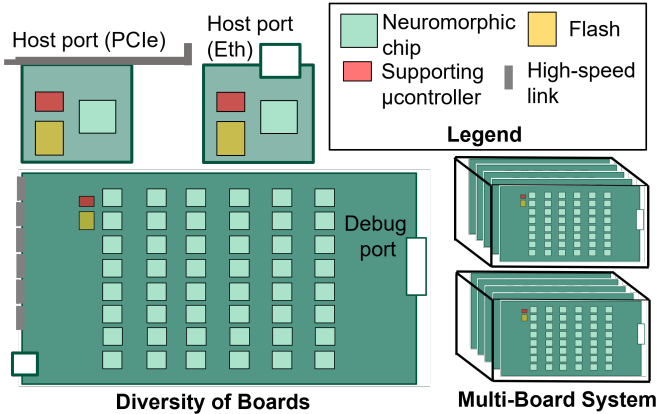| System | Developer | Technology | System size | Neurons/chip | Synapses/chip | Training | Has CPU | Framework |
|---|---|---|---|---|---|---|---|---|
| SNP T1 | Innatera | Mixed (28nm) | Edge | 1K | - | × | ✓ | Talamo |
| Speck | SynSense | Digital (65nm) | Edge | 320K | - | × | × | Sinabs |
| Xylo | SynSense | Digital (40nm) | Edge | 1K | 64K | × | × | Rockpool |
| DYNAP-SE2 | SynSense | Mixed (180nm) | Edge | 1K | 65K | ✓ | × | Rockpool |
| Akida | BrainChip | Digital (28nm) | Edge/Cloud | - | - | ✓ | ✓ | metaTF |
| GrAI VIP | GrAI Matter Labs | Digital (28nm) | Edge | 200K | - | ✓ | ✓ | GrAIFlow |
| Loihi 1 | Intel Labs | Digital (14nm) | Edge/Cloud | 128K | 128M | ✓ | ✓ | Lava/NxSDK |
| Loihi 2 | Intel Labs | Digital (Intel 4) | Edge/Cloud | 1M | 120M | ✓ | × | Lava |
| Tianjic | Tsinghua University | Digital (28nm) | Edge/Cloud | 40K | 10M | × | × | TJSim |
| BrainScaleS-1 | U Heidelberg | Mixed (180nm) | Cloud | 197K | 43M | ✓ | × | PyNN |
| BrainScaleS-2 | U Heidelberg | Mixed (65nm) | Edge/Cloud | 2K | 131K | ✓ | × | PyNN, hxTorch |
| NorthPole | IBM | Digital (12nm) | Cloud | 1M | 256M | × | × | NorthPole software toolchain |
| TrueNorth | IBM | Digital (28nm) | Edge/Cloud | 1M | 256M | × | × | CoreLet (Matlab) |
| SpiNNaker | U Manchester | Digital (130nm) | Edge/Cloud | 16K | 16M | ✓ | ✓ | PyNN |
| SpiNNaker2 | TU Dresden, U Manchester | Digital (22nm) | Edge/Cloud | 152K | 152M | ✓ | ✓ | Py-spinnaker2 |



Fig. 1. Different scales of neuromorphic systems

processor-based). Most often, those are multi-core architectures where each core combines synapse and neuron processing with a colocation of memory and computation in contrast to von Neumann architectures. We focus on neuromorphic ASICs and multi-chip systems from both industry and academia with high accessibility and mature software support. We think that these criteria are prerequisites for a near-term integration of such systems into data centers. Neuromorphic architectures such as SENECA [15] or FPGA-based systems like DeepSouth [16] are currently not covered but should be evaluated in the future.

We neither include neuromorphic compute-in-memory (CIM) architectures that perform vector-matrix multiplication in SRAM, analog crossbars, memristive or other nonvolatile memories such as [17], [18] nor neuromorphic photonics [19]. While those technologies promise very energy-efficient computing of AI models, they may face other challenges than spiking neuromorphic processors and are thus not part of this review. We also do not consider specialized digital DNN accelerator systems that aim to compete with the GPUs and TPUs. Such DNN accelerator systems are typically much closer integrated into powerful multiprocessor CPUs. But of course, neuromorphic chips should always be compared to the state-of-the-art, which includes digital DNN accelerators.

We also note that our list of hardware platforms is non-exhaustive. We refer to [20], [21], and [22] for overviews of large-scale neuromorphic systems and projects, and to [23] for a recent overview of trends in SNN processors. [24] compares digital SNN and DNN accelerators regarding efficiency and accuracy for DNN workloads.

### B. Neuromorphic Systems

We introduce mature neuromorphic systems sorted by the underlying companies or research groups. A comparison table is available at table I.

*1) Innatera:* The Innatera chip is an analog mixed-signal spiking neural processor for low-power edge applications. The chip comprises a low-power CPU, encoders, LIF neurons, and programmable synapses. There are around 1000 analog spiking neurons in each chip, connected through a multi-level crossbar structure [12]. Due to the limited number of neurons and constrained CPU power, the chip is particularly suitable for processing one-dimensional signals, such as those emerging in audio and healthcare applications.

*2) SynSense:* There are several commercial neuromorphic chips released by SynSense, such as Dynap-CNN, Speck [25],

Xylo [26], and DYNAP-SE2 [27]. The Speck and Dynap-CNN are digital chips designed for real-time vision processing applications such as gesture control, fall detection, and object tracking. The Speck chip is a system-on-chip combining a dynamic vision sensor and SNN cores containing 320K spiking neurons. It can run large-scale spiking convolutional neural networks (sCNNs) while consuming less than 1 mW of power [25]. Dynap-CNN is the processor-only variant of Speck supporting more complex sCNN models with 1 million neurons. DYNAP-SE2 [27], an analog mixed-signal chip with a built-in biosignal amplifier, is designed for wearable health devices and robotic applications and is more research-oriented.

*3) BrainChip:* Akida is an event-based processor developed by BrainChip [28]. Akida supports a wide variety of neural networks and can execute complex networks. It also supports the AXI bus for connection to CPUs, allowing custom networks not supported by Akida to be executed on the CPU. It comprises a data processing unit to preprocess input data, converting it into events, and uses an LPDDR4 interface for storing programs and parameters. Additionally, the PCIe interface can be used to connect to other Akida chips. Akida aims to support a broad range of applications, including robotics and automation in industry, real-time sensing in automotive, vital-signs prediction in on-device health monitoring, and intelligent automation in homes.

*4) GrAI Matter Labs:* The GrAI VIP chip consists of a CPU and a GrAICore with 196 NeuronFlow cores connected through an event-based network-on-chip (NoC), which is equipped with high-speed interfaces for cameras, microphones, speakers, and the host system. It is also optimized for both recurrent and feedforward models by supporting 16-bit floating-point data format and is appropriate for edge AI applications, such as audio and video processing [29].

*5) Intel Loihi 1 & 2:* Intel Loihi 1 is a programmable digital many-core neuromorphic system that approximates the behavior of biological neurons [30]. Loihi 2, the 2nd generation chip, comprises 128 neuron cores, each containing 8,192 neurons and 192 kB of memory that can be flexibly allocated between neurons and synapses. Therefore, each chip includes 1 million fully programmable neurons and 120 million synaptic connections. The neuron cores are interconnected by a NoC and support spike-based communications. The chip includes an inter-chip communication interface to facilitate the creation of large 3D chip clusters. [31]. Due to the promising scale-up ability, Loihi chips exhibit a big potential for integration into data centers.

*6) Tianjic:* The Tianjic chip is based on a 156-core architecture with localized memory and streamlined dataflow, which can be used to simulate 40,000 neurons and 10 million synapses [32]. The chip, supports both artificial neurons and spiking neurons, enabling emulation of various neural networks such as MLP, CNN, and RNN [33]. In contrast to other hybrid chips, such as SpiNNaker2, which have the flexibility to build state machines using non-neural code, the Tianjic chip uses a neural state machine to assemble its applications, trading off flexibility by high integration.

*7) BrainScaleS-1 & 2:* BrainScaleS-1 and BrainScaleS-2 (BSS) are analog mixed-signal neuromorphic chips developed by the University of Heidelberg. BSS-1 is not a programmable chip, while BSS-2 contains programmable synaptic connections. BSS-2 contains four analog neuron cores and digital synaptic arrays connected with a spike router. Each analog core includes 512 spiking neurons and 32,768 synapses. A BSS-2 system comprises multiple single-chip setups, which are interconnected to a computing cluster via Ethernet and suitable for robotic applications [34]–[37].

*8) IBM TrueNorth and NorthPole:* As the pioneer of the brain-inspired AI chips, TrueNorth [38] incorporates programmable digital neurons, whereas NorthPole [39] comprises computation units to simulate biological neurons. NorthPole contains 256 cores interconnected by two dense NoCs. Inspired by the brain's structure, one is designed for short-distance communication between nearby cores, and the other facilitates long-distance neuron activation communication across all cores. It contains 1 million programmable neurons and 256 million programmable synaptic connections. In total, 224 MB of on-chip memory is distributed across the 256 cores. The vector-matrix multiplier (VMM) can execute computations in 8-bit, 4-bit, and 2-bit fixed-point data formats. It is suitable for image classification, detection, segmentation, natural language processing, and speech recognition [38], [39].

*9) SpiNNaker 1 & 2:* SpiNNaker1 is a custom ARM-based 18-core chip developed in 130nm process technology by the University of Manchester, featuring a massively parallel architecture designed for large-scale real-time brain simulations with spiking neural networks. It currently holds the record for the world's largest neuromorphic supercomputer, including a total of 1'036.800 million cores, arranged in 1,200 48-node boards highly interconnected in a toroidal mesh. Such a supercomputer has the potential to emulate roughly 1 billion neurons and 1000 billion synapses, which might vary depending on the neuron models used [40], [41].

The successor SpiNNaker2 chip was developed in 22nm FDSOI by TU Dresden and the University of Manchester within the Human Brain Project. SpiNNaker2 features 152 ARM-based processing elements (PEs) for flexible software-based execution of neural networks. SpiNNaker2 deploys the same event packet routing as SpiNNaker1 and is designed to be scaled up to 10 million cores [14]. In addition to scalable brain simulations, SpiNNaker2 also targets efficient real-time AI processing with event-based DNN and generic computation [42]. SpiNNaker2 has custom accelerators to speed up the processing of DNN layers and for compute-heavy operations in neuromorphic computing [43].

*C. Summary of neuromorphic systems*

Multifarious neuromorphic systems cover a wide range of AI applications, from ultra-low power tinyML tasks on the edge to large-scale brain simulations on the cloud. In this paper, we focus on the neuromorphic systems that can be loaded into data centers for cloud computing. Those cloud neuromorphic systems are dominated by digital technology

and tend to use a more advanced process node for higher power efficiency. Only Innatera, SynSense Dynap-SE2, and BrainScaleS adopt mixed-signal solutions that are restrained in large-scale distributed AI applications due to varied analog process errors. Tuning AI models is regarded as one of the significant tasks in data centers. However, some cloud neuromorphic platforms such as Tianjic, Truenorth, and NorthPole don't provide the training functionality and focus on inference only. Moreover, the current state of software frameworks for neuromorphic computing exhibits fragmentation, with each player developing its own software stack to be maximally adapted to its hardware.

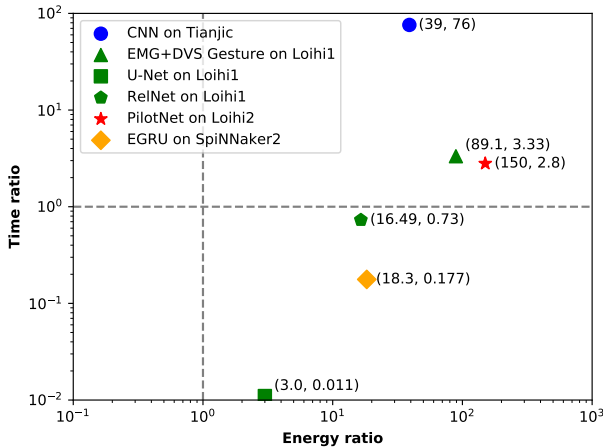## III. Comparing neuromorphic to conventional solutions for AI processing



Fig. 2. Comparison of energy and solution time ratios from table II

### A. Comparison of energy and speed

Targeting the integration of neuromorphic hardware into data centers for efficient AI processing, we next want to evaluate *for which applications* there are neuromorphic solutions and *by how much* these solutions are better in terms of energy and latency compared to conventional solutions. To do so, we have reviewed articles where the hardware systems from section V-A are benchmarked against conventional solutions (CPUs, GPUs, or accelerators like TPUs) in terms of energy, latency, and accuracy on the AI task. We neglect results on rather simple tasks like handwritten digit recognition or keyword spotting and instead look at more complex tasks as our focus is on AI deployment in data centers.

Six exemplary results on four hardware platforms are shown in table II and fig. 2: For each example, we report how much more energy-efficient the neuromorphic solution is (energy ratio) and how much faster it is to the compared solution (time ratio). In all cases, the correctness of the solution is comparable on both neuromorphic and conventional hardware. Regarding energy per inference, Loihi 1 defeats GPUs for relational reasoning with spiking LSTM [47], image segmentation

[46], and multi-sensor spatio-temporal gesture classification [45] by factor of 3 to 100. Loihi 2 achieves more than $100\times$ better energy on the automotive PilotNet benchmark with sigma-delta networks [48]. Tianjic requires on average $39\times$ less energy than a GPU for image classification using non-spiking CNNs [44]. SpiNNaker2 achieves an $18\times$ energy improvement for language modeling with event-based GRU [49] versus a data center GPU.

Concerning speed, there is no clear winner between neuromorphic systems and GPUs. Depending on the task and system, the neuromorphic solution can be up to $100\times$ faster or $100\times$ slower than on GPUs. Note that all time results in table II and fig. 2 are given for a batch size of 1. For larger batch sizes, GPUs achieve a lower average time per inference at improved energy efficiency compared to batch size 1, while on neuromorphic systems, the average time and energy per inference remain constant [46]–[49]. Hence, at larger batch sizes, the advantage of neuromorphic systems decreases.

### B. Interim conclusions

From the above results, we draw some interim conclusions:
- Tasks and models: Neuromorphic solutions are available for image processing with CNN, natural language processing with recurrent neural networks (spiking LSTM or event-based GRU), and spatiotemporal pattern recognition.
- Energy: Neuromorphic hardware is between 3 to 100 times more energy efficient per inference at batch size 1.
- Speed: For some tasks, neuromorphic hardware shows faster inference compared to conventional systems. This advantage diminishes for larger batch sizes.

### C. Limitations of analysis

In our analysis, we consider the energy ratio and solution time ratio compared to benchmarked GPU in the corresponding publications. It is clear that new conventional AI hardware platforms are available today that likely show a lower energy and lower inference time. Accordingly, an alternate approach would be to compare neuromorphic solutions to current state-of-the-art benchmark results. However, it would be unfair to benchmark a 5-year-old neuromorphic chip to a brand-new edge GPU using the latest software technology. Hence, we here defend our approach as both neuromorphic and conventional solutions represented the state-of-the-art at the time of publication.

We further note that the list of presented neuromorphic solutions is very limited and non-exhaustive. Due to a lack of public benchmarking results fulfilling our criteria, we cannot highlight any results from IBM TrueNorth or commercial neuromorphic system providers (Innatera, SynSense, BrainChip, and GrAI Matter Labs).

For a better and more transparent comparison between neuromorphic and conventional AI platforms, we recommend submitting results for standard machine learning benchmarks to MLPerf [50] and to the newly established neuromorphic benchmark NeuroBench [51].

| Reference | Hardware | Application | Task performance neurom. \| convent. | Energy ratio | Time ratio | Ref. hardware |
|---|---|---|---|---|---|---|
| Deng 2020 [44] | Tianjic | Image Classification (CNN) | 70.83%\| n.a. (Top-1 Acc. ↑) | 39 | 76 | NVIDIA V100 |
| Ceolini 2020 [45] | Loihi | EMG+DVS gesture recognition (CNN) | 96.0% \| 95.4% (Top-1 Acc. ↑) | 29.1 | 0.89 | Jetson Nano |
| Patel 2021 [46] | Loihi | Image segmentation (U-Net) | 92.13% \| 94.98% (Pixel Acc. ↑) | 3.00 | 0.011 | GeForce RTX 2080 |
| Rao 2022 [47] | Loihi | Time-series classification (RelNET) | 16/17 \| 16/17 (tasks solved ↑) | 4.36 | 0.73 | GeForce RTX 2070 |
| Shrestha 2023 [48] | Loihi 2 | Video processing (PilotNet) | 0.035 \| 0.025 (MSE ↓) | 150 | 2.8 | Jetson Orin Nano |
| Nazeer 2023 [49] | SpiNNaker2 | Language modelling (EGRU) | 97.3 \| 97.3 (Test PPL ↓) | 18.3 | 0.117 | NVIDIA A100 |

## IV. Applications, models, and algorithms

This section explores relevant AI workloads in diverse industrial applications, followed by an examination of the current state of spiking neural networks and their applications in machine learning. Finally, we outline future research directions in this evolving field.

### A. AI Workloads in Industries and Research

*1) Production Industry Applications:* The industry will be under enormous pressure to transform itself to a carbon-neutral future in the coming years. A good example of this is the steel industry, which wants to significantly reduce its $CO_2$ emissions as compared to 1990 [52]. To meet these challenges, technologies, and tools must be created at TRL8 (Technical Readiness Level) [53]. Of course, it is important to ensure that methods based on new technologies are also energy-efficient to avoid rebound effects.

From an industrial point of view, typical applications for AI models are computer vision-related tasks. For tasks like defect detection, tracking, ID recognition, or anomaly detection, camera-based solutions are fairly common. These types of tasks tend to be operated 24/7 with high throughput (production lines are often operated at speeds > 1m/s) and high workload, as defects can be fairly small and fairly local, making rigorous requirements on the high-resolution input images or video streams. So these real-time inputs are often evaluated 24/7 on several GPUs by multiple AI models at various production stations, making the energy consumption of the productive system non-negligible.

*2) Digital Industry Applications:* AI has shown super-linear growth trends in its share of computing usage in data centers [54]. The most common applications for AI in the digital industry include recommendation models, language models, vision models, etc., and it will most likely become even more prevalent in the future. Moreover, in recent years, generative models such as Large Language Models (LLM) and Large Multimodal Models (LMM) have gained explosive growth [55]. In these applications, a significant portion of computing resources is dedicated to inference, given their

role as generative services for end-users. As a result, energy consumption becomes a crucial consideration for the models.

*3) AI workloads in Research:* Looking at deep learning research, we see that the following AI model types are still relevant: convolutional neural networks, transformers, graph neural networks, generative adversarial networks, variational autoencoders, normalizing flows, diffusion models and deep reinforcement learning [56], [57]. From the models above, neuromorphic computing has mainly looked into convolutional networks (see section III) and into reinforcement learning [58]. Recurrent neural networks, for which SNN and neuromorphic computing have shown very efficient solutions [31], have moved out of focus a bit.

*4) Need for AI workload statistics:* From the academic literature and from general media it is possible to extract which AI applications and which AI models are trending. Also, it is now common to report the cost for training AI models in the machine learning literature as the GPU hours used [59], or even to provide the estimated $CO_2e$ emissions [60]. Unfortunately, we miss the public information about which AI models are run, how often in commercial data centers, and what is their share of the total compute resources. The big tech companies only share rough information, e.g., about the relative increase of AI tasks or the share of inference, training, and network architecture search [1]. Detailed AI workload statistics would help determine which AI tasks and models to focus on for developing energy-efficient neuromorphic solutions. Finally, this could help to apply the greatest leverage for reducing the operational cost in AI data centers.

### B. Spiking Neural Networks in Machine Learning: Current Landscape

Spiking neural networks (SNN) [61] are well-suited for efficient event-based implementations and have been scaled to very large sizes [62]. However, for many years, they have been avoided in machine learning because of their non-differentiable dynamics, which at first glance made them ill-suited for gradient-based learning. However, several methods have been proposed to tackle this problem, enabling gradient-based end-to-end learning for SNNs [63]–[65]. This development has led to several successful implementations of SNN

learning on neuromorphic hardware. Algorithms like Spiking-based Backpropagation [66] [67] and other hybrid methods have been shown to achieve comparable performance as their ANNs counterparts at image classification tasks [68] [69]. Additionally, training techniques like Spiking Generative Networks in Lifelong Learning Environment [70] demonstrated their effectiveness in image classification and generation.

The most powerful SNN models that reach close to the state-of-the-art performance of conventional machine learning models have so far only been demonstrated in software. Recently, there has been increasing interest in porting large transformer-based models [71] to SNNs [72]–[76]. Several recent studies have highlighted the high levels of sparsity in transformers [77]–[83], and spiking neurons are an efficient method to make use of this feature. Specifically, Spikformer [72] is a variant of the transformer network architecture based on spiking neural networks. They introduce a spiking variant of self-attention to efficiently implement transformers with SNNs and demonstrate compelling performance on a range of benchmark tasks. Recently, a newer version of this model has been published, which reaches up to the performance of its non-spiking counterparts at a significantly reduced compute budget [76]. In addition, other research groups have introduced spiking variants of specific popular transformer architectures for large language models, such as GPT [73] or BERT [74], [75], demonstrating promising results.

### C. Research Directions for SNNs in Machine Learning

As we approach the forefront of the current research, we suggest that the future lies in incorporating SNN-based models into mainstream applications, focusing on extending large-scale applications beyond image classification, refining training techniques, and optimizing mapping strategies for diverse hardware platforms. Efforts in standardization and compatibility with existing frameworks will enhance SNN adoption across industries. Additionally, delving into novel applications, particularly in real-time scenarios, and further integration with transformer-based models could unlock new frontiers for SNN research. Emphasizing energy efficiency, economic viability, and scalability will be pivotal for solidifying SNNs as a key player in the evolving landscape of machine learning and neuromorphic computing.

## V. Hardware and software integration

In this section, we discuss the hardware and software challenges for integrating neuromorphic hardware platforms into AI data centers.

### A. Hardware integration

*1) Status quo:* So far, the following neuromorphic systems have been integrated at large scale into data centers: SpiNNaker 1 [20], TrueNorth (NS16-4e) [84], Loihi (Pohoiki Springs) [85], BrainScaleS-1 [86] and Tianjic [33]. All systems use slide-in modules with custom printed circuit boards (PCBs) for integration into standard 19" server racks. Typically, the neuromorphic chips are accessed via Ethernet,

only the TrueNorth NS16e-4 uses PCIe for communication with the host chip. Baseboard Management Controller (BMC) or similar controllers are used for booting and monitoring the boards. All platforms also include field-programmable gate arrays (FPGAs) or system-on-chips (SoCs), most often as middleware between host computers and neuromorphic systems. Some systems have already integrated a host CPU, e.g., Pohoiki Springs or the Tianjic server, while the other systems require external host CPU servers for the configuration and control of the neuromorphic systems and for preprocessing.

As an exception, BrainChip offers PCIe boards for integrating their Akida chips with CPU servers. This represents another option for integrating neuromorphic computing systems into data centers, similar to normal GPUs. Note, however, that this might limit the size of neuromorphic models that can be implemented compared to the larger systems discussed above.

*2) Conclusion:* The above examples show that a variety of neuromorphic systems have been successfully integrated into standard data center server racks. Thus, technically, the hardware integration does not pose a problem. Yet, we observe a diversity in how large neuromorphic systems are assembled into server boards, e.g., many of them leverage FPGAs or SoCs as middleware. These extra devices and the host CPU add a power overhead to the very energy-efficient neuromorphic systems. Optimizing for system-level efficiency of AI compute servers, these components need to be included when performing benchmarking on AI workloads. Another requirement for the industry-level deployment of neuromorphic chips is high reliability and robustness. The chips and boards need to be designed for a 24/7 operation, e.g., the server board should keep working if a single chip or processor fails. Replacement parts should be available for a long period.

### B. Software and operation

*1) Operation principles:* The way how neuromorphic systems are operated significantly differs from other AI accelerators or GPUs. There is no operating system or runtime that schedules compute tasks sequentially on neuromorphic cores. Instead, the synaptic weights and neuron parameters are configured first on all cores and chips of the neuromorphic systems. Then, input data, such as spikes or scalar events, are streamed into the NC system and processed by the neurons and synapses. Most of the systems operate in real-time, which means the individual chips run asynchronously, and spike events are processed as they arrive. In mixed-signal systems, the decay time constants of the silicon neurons define the speed of operation, which may be equal to the speed of biological neurons like in DYNAP-SE2 [27], or accelerated by a factor of 1000 to 10000 in the BrainScaleS systems [37]. Digital systems such as TrueNorth, SpiNNaker, or Loihi typically split the neuron updates into timesteps. How long it takes to process one timestep then defines how fast a spiking network can be executed. Loihi offers a barrier synchronization to continue with the next time step once all cores in all jobs are done with the current step. In contrast, in SpiNNaker, neuron updates are triggered in regular intervals (e.g., 1 ms),

effectively yielding a real-time system. Typically, for AI inference neuromorphic systems process with batch size 1. Multiple inputs are processed sequentially or need to be distributed spatially onto different system cores. Note, however, that some SNN architectures allow for pipelining input data when each layer requires the same fixed number of timesteps [87], [88].

The real-time operation of neuromorphic systems poses a challenge for integration into cloud-based digital computing systems. Nilsson et al. [9] frame this challenge in detail and suggest a conceptual framework based on micro-services to solve it. For data centers, we consider several scenarios on how AI compute load arrives at the host CPU managing the neuromorphic system:

- Streaming application in edge cloud, e.g., for processing data from a video camera at 30 FPS in real-time. In this case, a neuromorphic system is pre-configured with a neural network model. Each image needs to be converted to data suitable for the hardware (e.g., spikes) and then provided into a neuromorphic system. Results are recorded and processed further on the host CPU. Depending on the SNN model, the states need to be reset after each input image. Another option might be to use a network model for video streaming that does not require an external state reset.
- Irregular requests from the web. The frequency of AI compute requests arriving in the data center may vary across the time of the day, the week, or even seasons. This creates very irregular workload statistics. Data centers with neuromorphic hardware will need to handle such scenarios. Again, for a specific task, the neural network weights can be pre-loaded. Then, inference is triggered by arriving requests. For neural network architectures supporting pipelining, multiple requests can be buffered and forwarded sequentially into the NC hardware to achieve the highest possible throughput. Maximum response latencies need to be adhered to. If the incoming request rate becomes too high, the processing of the AI tasks needs to be redirected to other neuromorphic systems or conventional compute units. In longer idle intervals, the neuromorphic chips may even be switched off or sent into sleep or retention mode.
- Regular, low-priority AI tasks that can be scheduled to free neuromorphic resources. For such tasks, the execution typically also includes the generation and loading of the hardware configuration. There is no specific challenge except for achieving a high utilization of the neuromorphic chip resources.

We note that for most of the tasks, a pre-processing of data (e.g., encoding into spikes) and the post-processing of data (e.g., decoding of spike rates into class probabilities) are required

Regarding the dynamic scheduling of AI compute jobs onto neuromorphic hardware, data center standards need to be adopted. Current cloud-scale neuromorphic systems either use SLURM [89] or custom schedulers [90], while commercial data centers use container orchestration platforms like Kubernetes or alternatives.

*2) High-level software:* To enable neuromorphic computing for efficient AI processing, advanced software tools are needed for the training and deployment of SNN or other brain-inspired models [91]. While there exists a multitude of software frameworks for training SNN in PyTorch [92]–[94] or Jax [95], deep SNN needs to be further optimized for each hardware platform considering details of the neuron model or quantization of weights. Because of this, each neuromorphic system provides its own software framework, as is shown in the last column of Table I. While there are some approaches for unifying the programming model of NC [96]–[98], none of them have evolved to a standard yet.

However, the standardization of tools and ease of use are key factors for a new technology to obtain acceptance in the industry. Thus, to proliferate SNNs further in an industrial setting, standard APIs for developing and training SNNs are required, e.g. comparable to Keras [99], TensorFlow [100] or PyTorch [101], ideally even a frontend such as Keras with the backend, whether it be SNN-based or a conventional TensorFlow or PyTorch, freely interchangeable. This would make development on standard hardware possible and an easy transition to specialized SNNs for a productive system, regardless of whether the SNNs are located on-premise or in an external data center, aka in the cloud. This kind of portability could also help avoid the fear of vendor lock-in when transitioning to a comparably new technology stack, as well as questions regarding long-term support when switching from established suppliers to a new, comparably small hardware supplier.

In addition to that, it is essential to provide a large number of examples and a so-called model zoo of validated AI models for each hardware platform. So far, only BrainChip provides a model zoo for their Akida systems [102]. This allows us to retrain or fine-tune existing model architectures to customers' needs, which not only offers a shorter time-to-solution but also reduces the environmental footprint, as training from scratch can be avoided.

## VI. Conclusion

In this article, we reviewed neuromorphic hardware platforms and algorithms for their suitability in reducing energy consumption in AI data centers. We also discussed the current challenges that neuromorphic computing faces in becoming a mainstream technology used by the industry. In particular, we analyzed that the current AI model types supported by neuromorphic computing only partially match the AI models commonly run in AI data centers. We conclude that the neuromorphic computing community should focus on state-of-the-art ML technologies, such as transformers, and needs to establish standardized software frameworks that ensure interoperability among hardware vendors.

Data center sustainability is not only about saving energy during operations but also about saving water and materials while keeping social and governance issues in mind. These

latter issues are becoming increasingly important as AI models use vast amounts of personal data. When considering the carbon footprint, one will eventually face the question of whether or not to integrate specialized hardware: the embodied footprint of an additional device may be greater than the operational footprint savings due to specialized solutions [103]. Neuromorphic engineers should therefore focus on the high utilization of their platforms.

## REFERENCES

[1] C.-J. Wu, R. Raghavendra, U. Gupta, B. Acun, N. Ardalani, K. Maeng, G. Chang, F. Aga, J. Huang, C. Bai *et al.*, "Sustainable ai: Environmental implications, challenges and opportunities," *Proceedings of Machine Learning and Systems*, vol. 4, pp. 795–813, 2022.

[2] J. Guitart, "Toward sustainable data centers: a comprehensive energy management strategy," *Computing*, vol. 99, no. 6, pp. 597–615, 2017.

[3] K. Roy, A. Jaiswal, and P. Panda, "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, vol. 575, no. 7784, p. 607–617, Nov 2019. [Online]. Available: https://doi.org/10.1038/s41586-019-1677-2

[4] N. Rathi, I. Chakraborty, A. Kosta, A. Sengupta, A. Ankit, P. Panda, and K. Roy, "Exploring neuromorphic computing based on spiking neural networks: Algorithms to hardware," *ACM Comput. Surv.*, vol. 55, no. 12, mar 2023. [Online]. Available: https://doi.org/10.1145/3571155

[5] D. Mytton and M. Ashtine, "Sources of data center energy estimates: A comprehensive review," *Joule*, vol. 6, no. 9, pp. 2032–2056, 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2542435122003580

[6] A. Hantsch, L. M. Banzer, C. Wächter, A. Weisemann, and J. Struckmeier, "Pushing the boundary conditions of data centers facilitates innovative circular economy approaches," in *Open Compute Project Future Technology Symposium, San José, CA, USA*, 2021.

[7] G. McGovern, *World Wide Waste - How digital is killing our planet - And what we can do about it.* Gormanston: Silver Beach Publishing, 2020.

[8] J. S. Vetter, R. Brightwell, M. Gokhale, P. McCormick, R. Ross, J. Shalf, K. Antypas, D. Donofrio, T. Humble, C. Schuman, B. Van Essen, S. Yoo, A. Aiken, D. Bernholdt, S. Byna, K. Cameron, F. Cappello, B. Chapman, A. Chien, M. Hall, R. Hartman-Baker, Z. Lan, M. Lang, J. Leidel, S. Li, R. Lucas, J. Mellor-Crummey, P. Peltz Jr., T. Peterka, M. Strout, and J. Wilke, "Extreme heterogeneity 2018 - productive computational science in the era of extreme heterogeneity: Report for doe ascr workshop on extreme heterogeneity," USDOE Office of Science (SC), Washington, DC (United States), Tech. Rep., 12 2018. [Online]. Available: https://www.osti.gov/biblio/1473756

[9] M. Nilsson, O. Schelén, A. Lindgren, U. Bodin, C. Paniagua, J. Delsing, and F. Sandin, "Integration of neuromorphic ai in event-driven distributed digitized systems: Concepts and research directions," *Frontiers in Neuroscience*, vol. 17, 2023. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fnins.2023.1074439

[10] T. Luo, W.-F. Wong, R. S. M. Goh, A. T. Do, Z. Chen, H. Li, W. Jiang, and W. Yau, "Achieving green ai with energy-efficient deep learning using neuromorphic computing," *Commun. ACM*, vol. 66, no. 7, p. 52–57, jun 2023. [Online]. Available: https://doi.org/10.1145/3588591

[11] "Speck," https://www.synsense.ai/products/speck/, accessed: 2024-18-01.

[12] S. Ward-Foxton, "Innatera unveils neuromorphic ai chip to accelerate spiking networks," Tech. Rep., 2021.

[13] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain *et al.*, "Loihi: A neuromorphic manycore processor with on-chip learning," *Ieee Micro*, vol. 38, no. 1, pp. 82–99, 2018.

[14] C. Mayr, S. Hoeppner, and S. Furber, "Spinnaker 2: A 10 million core processor system for brain simulation and machine learning," 2019.

[15] G. Tang, K. Vadivel, Y. Xu, R. Bilgic, K. Shidqi, P. Detterer, S. Traferro, M. Konijnenburg, M. Sifalakis, G.-J. van Schaik, and A. Yousefzadeh, "Seneca: building a fully digital neuromorphic processor, design trade-offs and challenges," *Frontiers in Neuroscience*, vol. 17, 2023. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fnins.2023.1187252

[16] R. M. Wang, C. S. Thakur, and A. van Schaik, "An fpga-based massively parallel neuromorphic cortex simulator," *Frontiers in Neuroscience*, vol. 12, 2018. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fnins.2018.00213

[17] W. Wan, R. Kubendran, C. Schaefer, S. B. Eryilmaz, W. Zhang, D. Wu, S. Deiss, P. Raina, H. Qian, B. Gao *et al.*, "A compute-in-memory chip based on resistive random-access memory," *Nature*, vol. 608, no. 7923, pp. 504–512, 2022.

[18] M. L. Gallo, R. Khaddam-Aljameh, M. Stanisavljevic, A. Vasilopoulos, B. Kersting, M. Dazzi, G. Karunaratne, M. Braendli, A. Singh, S. M. Mueller, J. Büchel, X. Timoneda, V. Joshi, M. J. Rasch, U. Egger, A. Garofalo, A. Petropoulos, T. A. Antonakopoulos, K. Brew, S. Choi, I. Ok, T. M. Philip, V. Chan, C. Silvestre, I. Ahsan, N. Saulnier, V. Narayanan, P. A. Francese, E. Eleftheriou, and A. Sebastian, "A 64-core mixed-signal in-memory compute chip based on phase-change memory for deep neural network inference," *Nature Electronics*, vol. 6, pp. 680–693, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:254275062

[19] B. J. Shastri, A. N. Tait, T. Ferreira de Lima, W. H. Pernice, H. Bhaskaran, C. D. Wright, and P. R. Prucnal, "Photonics for artificial intelligence and neuromorphic computing," *Nature Photonics*, vol. 15, no. 2, pp. 102–114, 2021.

[20] S. Furber, "Large-scale neuromorphic computing systems," *Journal of neural engineering*, vol. 13, no. 5, p. 051001, 2016.

[21] C. S. Thakur, J. L. Molin, G. Cauwenberghs, G. Indiveri, K. Kumar, N. Qiao, J. Schemmel, R. Wang, E. Chicca, J. Olson Hasler, J.-s. Seo, S. Yu, Y. Cao, A. van Schaik, and R. Etienne-Cummings, "Large-scale neuromorphic spiking array processors: A quest to mimic the brain," *Frontiers in Neuroscience*, vol. 12, 2018. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fnins.2018.00891

[22] D. Ivanov, A. Chezhegov, M. Kiselev, A. Grunin, and D. Larionov, "Neuromorphic artificial intelligence systems," *Frontiers in Neuroscience*, vol. 16, p. 1513, 2022.

[23] A. Basu, L. Deng, C. Frenkel, and X. Zhang, "Spiking neural network integrated circuits: A review of trends and future directions," in *2022 IEEE Custom Integrated Circuits Conference (CICC)*, 2022, pp. 1–8.

[24] F. Ottati, C. Gao, Q. Chen, G. Brignone, M. R. Casu, J. K. Eshraghian, and L. Lavagno, "To spike or not to spike: A digital hardware perspective on deep learning acceleration," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 13, no. 4, pp. 1015–1025, 2023.

[25] M. Yao, O. Richter, G. Zhao, N. Qiao, Y. Xing, D. Wang, T. Hu, W. Fang, T. Demirci, M. De Marchi, L. Deng, T. Yan, C. Nielsen, S. Sheik, C. Wu, Y. Tian, B. Xu, and G. Li, "Spike-based dynamic computing with asynchronous sensing-computing neuromorphic chip," *Nature Communications*, vol. 15, no. 1, p. 4464, May 2024. [Online]. Available: https://www.nature.com/articles/s41467-024-47811-6

[26] H. Bos and D. R. Muir, "Sub-mw neuromorphic snn audio processing applications with rockpool and xylo," *ArXiv*, vol. abs/2208.12991, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:251903240

[27] O. Richter, C. Wu, A. M. Whatley, G. Köstinger, C. Nielsen, N. Qiao, and G. Indiveri, "Dynap-se2: a scalable multi-core dynamic neuromorphic asynchronous spiking neural network processor," *Neuromorphic Computing and Engineering*, 2024. [Online]. Available: http://iopscience.iop.org/article/10.1088/2634-4386/ad1cd7

[28] B. M. Posey, "What Is the Akida Event Domain Neural Processor?" Tech. Rep., 2020. [Online]. Available: https://brainchip.com/wp-content/uploads/2020/03/BrainChip_tech-brief_What-is-Akida_v3-1.pdf

[29] O. Moreira, A. Yousefzadeh, F. Chersi, G. Cinserin, R.-J. Zwartenkot, A. Kapoor, P. Qiao, P. Kievits, M. Khoei, L. Rouillard, A. Ferouge, J. Tapson, and A. Visweswara, "NeuronFlow: a neuromorphic processor architecture for Live AI applications," in *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, Mar. 2020, pp. 840–845, iSSN: 1558-1101. [Online]. Available: https://ieeexplore.ieee.org/abstract/document/9116352

[30] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain, Y. Liao, C.-K. Lin, A. Lines, R. Liu, D. Mathaikutty, S. McCoy, A. Paul, J. Tse, G. Venkataramanan, Y.-H. Weng, A. Wild, Y. Yang, and H. Wang, "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82–99, 2018.

[31] M. Davies, "Taking neuromorphic computing to the next level with Loihi2," *Intel Labs' Loihi*, vol. 2, pp. 1–7, 2021. [Online]. Available: https://download.intel.com/newsroom/2021/new-technologies/neuromorphic-computing-loihi-2-brief.pdf

[32] J. Pei, L. Deng, S. Song, M. Zhao, Y. Zhang, S. Wu, G. Wang, Z. Zou, Z. Wu, W. He, F. Chen, N. Deng, S. Wu, Y. Wang, Y. Wu, Z. Yang, C. Ma, G. Li, W. Han, H. Li, H. Wu, R. Zhao, Y. Xie, and L. Shi, "Towards artificial general intelligence with hybrid tianjic chip architecture," *Nature*, vol. 572, no. 7767, pp. 106–111, 2019, cited By 52.

[33] J. Pei, L. Deng, C. Ma, X. Liu, and L. Shi, "Multi-grained system integration for hybrid-paradigm brain-inspired computing," *Science China Information Sciences*, vol. 66, pp. 1–14, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:257498470

[34] S. Schmitt, J. Klähn, G. Bellec, A. Grübl, M. Güttler, A. Hartel, S. Hartmann, D. Husmann, K. Husmann, S. Jeltsch, V. Karasenko, M. Kleider, C. Koke, A. Kononov, C. Mauch, E. Müller, P. Müller, J. Partzsch, M. A. Petrovici, S. Schiefer, S. Scholze, V. Thanasoulis, B. Vogginger, R. Legenstein, W. Maass, C. Mayr, R. Schüffny, J. Schemmel, and K. Meier, "Neuromorphic hardware in the loop: Training a deep spiking network on the brainscales wafer-scale system," in *2017 International Joint Conference on Neural Networks (IJCNN)*, 2017, pp. 2227–2234.

[35] E. Müller, S. Schmitt, C. Mauch, S. Billaudelle, A. Grübl, M. Güttler, D. Husmann, J. Ilmberger, S. Jeltsch, J. Kaiser, J. Klähn, M. Kleider, C. Koke, J. Montes, P. Müller, J. Partzsch, F. Passenberg, H. Schmidt, B. Vogginger, J. Weidner, C. Mayr, and J. Schemmel, "The operating system of the neuromorphic brainscales-1 system," *Neurocomputing*, vol. 501, pp. 790–810, 2022. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0925231222006646

[36] S. Neuwirth, D. Frey, M. Nuessle, and U. Bruening, "Scalable communication architecture for network-attached accelerators," in *2015 IEEE 21st International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2015, pp. 627–638.

[37] E. Müller, E. Arnold, O. Breitwieser, M. Czierlinski, A. Emmel, J. Kaiser, C. Mauch, S. Schmitt, P. Spilger, R. Stock *et al.*, "A scalable approach to modeling on accelerated neuromorphic hardware," *Frontiers in neuroscience*, vol. 16, p. 690, 2022.

[38] P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura *et al.*, "A million spiking-neuron integrated circuit with a scalable communication network and interface," *Science*, vol. 345, no. 6197, pp. 668–673, 2014.

[39] D. S. Modha, F. Akopyan, A. Andreopoulos, R. Appuswamy, J. V. Arthur, A. S. Cassidy, P. Datta, M. V. DeBole, S. K. Esser, C. O. Otero *et al.*, "Neural inference at the frontier of energy, space, and time," *Science*, vol. 382, no. 6668, pp. 329–335, 2023.

[40] "SpiNNworld: SpiNNaker presence worldwide," 2023, accessed on October 3, 2023. [Online]. Available: https://www.google.com/maps/d/u/0/edit?mid=1jrbV2OVaBFqGlVYMxerSh0Pcexd_wznQ&ll=1.5170674512027844%2C0&z=2

[41] S. Furber and P. Bogdan, Eds., *SpiNNaker: A Spiking Neural Network Architecture*. Boston-Delft: now publishers, 2020. [Online]. Available: http://dx.doi.org/10.1561/9781680836523

[42] H. A. Gonzalez, J. Huang, F. Kelber, K. K. Nazeer, T. Langer, C. Liu, M. Lohrmann, A. Rostami, M. Schöne, B. Vogginger *et al.*, "Spinnaker2: A large-scale neuromorphic system for event-based and asynchronous machine learning," *arXiv preprint arXiv:2401.04491*, 2024.

[43] S. Höppner, Y. Yan, A. Dixius, S. Scholze, J. Partzsch, M. Stolba, F. Kelber, B. Vogginger, F. Neumärker, G. Ellguth, S. Hartmann, S. Schiefer, T. Hocker, D. Walter, G. Liu, J. Garside, S. Furber, and C. Mayr, "The spinnaker 2 processing element architecture for hybrid digital neuromorphic computing," 2022.

[44] L. Deng, G. Wang, G. Li, S. Li, L. Liang, M. Zhu, Y. Wu, Z. Yang, Z. Zou, J. Pei, Z. Wu, X. Hu, Y. Ding, W. He, Y. Xie, and L. Shi, "Tianjic: A unified and scalable chip bridging spike-based and continuous neural computation," *IEEE Journal of Solid-State Circuits*, vol. 55, no. 8, pp. 2228–2246, 2020.

[45] E. Ceolini, C. Frenkel, S. B. Shrestha, G. Taverni, L. Khacef, M. Payvand, and E. Donati, "Hand-gesture recognition based on emg and event-based camera sensor fusion: A benchmark in neuromorphic computing," *Frontiers in neuroscience*, vol. 14, p. 637, 2020.

[46] K. Patel, E. Hunsberger, S. Batir, and C. Eliasmith, "A spiking neural network for image segmentation," *arXiv preprint arXiv:2106.08921*, 2021.

[47] A. Rao, P. Plank, A. Wild, and W. Maass, "A long short-term memory for ai applications in spike-based neuromorphic hardware," *Nature Machine Intelligence*, vol. 4, no. 5, pp. 467–479, 2022.

[48] S. B. Shrestha, J. Timcheck, P. Frady, L. Campos-Macias, and M. Davies, "Efficient video and audio processing with loihi 2," *arXiv preprint arXiv:2310.03251*, 2023.

[49] K. K. Nazeer, M. Schöne, R. Mukherji, C. Mayr, D. Kappel, and A. Subramoney, "Language modeling on a spinnaker 2 neuromorphic chip," *arXiv preprint arXiv:2312.09084*, 2023.

[50] V. J. Reddi, C. Cheng, D. Kanter, P. Mattson, G. Schmuelling, C.-J. Wu, B. Anderson, M. Breughe, M. Charlebois, W. Chou, R. Chukka, C. Coleman, S. Davis, P. Deng, G. Diamos, J. Duke, D. Fick, J. S. Gardner, I. Hubara, S. Idgunji, T. B. Jablin, J. Jiao, T. S. John, P. Kanwar, D. Lee, J. Liao, A. Lokhmotov, F. Massa, P. Meng, P. Micikevicius, C. Osborne, G. Pekhimenko, A. T. R. Rajan, D. Sequeira, A. Sirasao, F. Sun, H. Tang, M. Thomson, F. Wei, E. Wu, L. Xu, K. Yamada, B. Yu, G. Yuan, A. Zhong, P. Zhang, and Y. Zhou, "Mlperf inference benchmark," in *Proceedings of the ACM/IEEE 47th Annual International Symposium on Computer Architecture*, ser. ISCA '20. IEEE Press, 2020, p. 446–459. [Online]. Available: https://doi.org/10.1109/ISCA45697.2020.00045

[51] J. Yik, K. V. den Berghe, D. den Blanken, Y. Bouhadjar, M. Fabre, P. Hueber, D. Kleyko, N. Pacik-Nelson, P.-S. V. Sun, G. Tang, S. Wang, B. Zhou, S. H. Ahmed, G. V. Joseph, B. Leto, A. Micheli, A. K. Mishra, G. Lenz, T. Sun, Z. Ahmed, M. Akl, B. Anderson, A. G. Andreou, C. Bartolozzi, A. Basu, P. Bogdan, S. Bohte, S. Buckley, G. Cauwenberghs, E. Chicca, F. Corradi, G. de Croon, A. Danielescu, A. Daram, M. Davies, Y. Demirag, J. Eshraghian, T. Fischer, J. Forest, V. Fra, S. Furber, P. M. Furlong, W. Gilpin, A. Gilra, H. A. Gonzalez, G. Indiveri, S. Joshi, V. Karia, L. Khacef, J. C. Knight, L. Kriener, R. Kubendran, D. Kudithipudi, Y.-H. Liu, S.-C. Liu, H. Ma, R. Manohar, J. M. Margarit-Taulé, C. Mayr, K. Michmizos, D. Muir, E. Neftci, T. Nowotny, F. Ottati, A. Ozcelikkale, P. Panda, J. Park, M. Payvand, C. Pehle, M. A. Petrovici, A. Pierro, C. Posch, A. Renner, Y. Sandamirskaya, C. J. Schaefer, A. van Schaik, J. Schemmel, S. Schmidgall, C. Schuman, J. sun Seo, S. Sheik, S. B. Shrestha, M. Sifalakis, A. Sironi, M. Stewart, K. Stewart, T. C. Stewart, P. Stratmann, J. Timcheck, N. Tömen, G. Urgese, M. Verhelst, C. M. Vineyard, B. Vogginger, A. Yousefzadeh, F. T. Zohora, C. Frenkel, and V. J. Reddi, "Neurobench: A framework for benchmarking neuromorphic computing algorithms and systems," 2024.

[52] EUROFER, "Low carbon roadmap, pathways to a co2-neutral european steel industry," https://www.eurofer.eu/assets/Uploads/EUROFER-Low-Carbon-Roadmap-Pathways-to-a-CO2-neutral-European-Steel-Industry.pdf, 2019.

[53] "Proposal for clean steel partnership under the horizon europe programme," https://www.estep.eu/assets/Uploads/ec-rtd-he-partnerships-for-clean-steel-low-carbon-steelmaking.pdf, 2020.

[54] A. S. Andrae and T. Edler, "On global electricity usage of communication technology: trends to 2030," *Challenges*, vol. 6, no. 1, pp. 117–157, 2015.

[55] J. Sevilla, L. Heim, A. Ho, T. Besiroglu, M. Hobbhahn, and P. Villalobos, "Compute trends across three eras of machine learning," in *2022 International Joint Conference on Neural Networks (IJCNN)*, 2022, pp. 1–8.

[56] S. J. Prince, *Understanding Deep Learning*. MIT press, 2023.

[57] C. M. Bishop and H. Bishop. Springer, 2023.

[58] T. Wunderlich, A. F. Kungl, E. Müller, A. Hartel, Y. Stradmann, S. A. Aamir, A. Grübl, A. Heimbrecht, K. Schreiber, D. Stöckel *et al.*, "Demonstrating advantages of neuromorphic computation: a pilot study," *Frontiers in neuroscience*, vol. 13, p. 260, 2019.

[59] "Neurips code of ethics." [Online]. Available: https://neurips.cc/public/EthicsGuidelines

[60] L. F. W. Anthony, B. Kanding, and R. Selvan, "Carbontracker: Tracking and predicting the carbon footprint of training deep learning models," *arXiv preprint arXiv:2007.03051*, 2020.

[61] W. Maass, "Networks of spiking neurons: the third generation of neural network models," *Neural networks*, vol. 10, no. 9, pp. 1659–1671, 1997.

[62] J. Jordan, T. Ippen, M. Helias, I. Kitayama, M. Sato, J. Igarashi, M. Diesmann, and S. Kunkel, "Extremely scalable spiking neuronal

network simulation code: from laptops to exascale computers," *Frontiers in neuroinformatics*, vol. 12, p. 2, 2018.

[63] G. Bellec, F. Scherr, A. Subramoney, E. Hajek, D. Salaj, R. Legenstein, and W. Maass, "A solution to the learning dilemma for recurrent networks of spiking neurons," *Nature Communications*, vol. 11, no. 1, p. 3625, Jul. 2020.

[64] T. C. Wunderlich and C. Pehle, "Event-based backpropagation can compute exact gradients for spiking neural networks," *Scientific Reports*, vol. 11, no. 1, p. 12829, Jun 2021.

[65] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud, "Neural ordinary differential equations," *Advances in neural information processing systems*, vol. 31, 2018.

[66] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, 2019.

[67] Y. Zhu, Z. Yu, W. Fang, X. Xie, T. Huang, and T. Masquelier, "Training spiking neural networks with event-driven backpropagation," *Advances in Neural Information Processing Systems*, vol. 35, pp. 30 528–30 541, 2022.

[68] W. Fang, Z. Yu, Y. Chen, T. Masquelier, T. Huang, and Y. Tian, "Incorporating learnable membrane time constant to enhance learning of spiking neural networks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 2661–2671.

[69] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, and Y. Tian, "Deep residual learning in spiking neural networks," *Advances in Neural Information Processing Systems*, vol. 34, pp. 21 056–21 069, 2021.

[70] J. Zhang, W. Fan, and X. Liu, "Spiking generative networks in life-long learning environment," in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*. Springer, 2023, pp. 353–364.

[71] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[72] Z. Zhou, Y. Zhu, C. He, Y. Wang, S. Yan, Y. Tian, and L. Yuan, "Spikformer: When spiking neural network meets transformer," *arXiv preprint arXiv:2209.15425*, 2022.

[73] R.-J. Zhu, Q. Zhao, and J. K. Eshraghian, "Spikegpt: Generative pre-trained language model with spiking neural networks," *arXiv preprint arXiv:2302.13939*, 2023.

[74] C. Lv, T. Li, J. Xu, C. Gu, Z. Ling, C. Zhang, X. Zheng, and X. Huang, "Spikebert: A language spikformer trained with two-stage knowledge distillation from bert," *arXiv preprint arXiv:2308.15122*, 2023.

[75] M. Bal and A. Sengupta, "Spikingbert: Distilling bert to train spiking language models using implicit differentiation," *arXiv preprint arXiv:2308.10873*, 2023.

[76] Z. Zhou, K. Che, W. Fang, K. Tian, Y. Zhu, S. Yan, Y. Tian, and L. Yuan, "Spikformer v2: Join the high accuracy club on imagenet with an snn ticket," *arXiv preprint arXiv:2401.02020*, 2024.

[77] Z. Li, M. S. Asif, and Z. Ma, "Event transformer," 2022.

[78] Z. Liu, J. Wang, T. Dao, T. Zhou, B. Yuan, Z. Song, A. Shrivastava, C. Zhang, Y. Tian, C. Re *et al.*, "Deja vu: Contextual sparsity for efficient llms at inference time," in *International Conference on Machine Learning*. PMLR, 2023, pp. 22 137–22 176.

[79] S. Jaszczur, A. Chowdhery, A. Mohiuddin, L. Kaiser, W. Gajewski, H. Michalewski, and J. Kanerva, "Sparse is enough in scaling transformers," *Advances in Neural Information Processing Systems*, vol. 34, pp. 9895–9907, 2021.

[80] I. Han, R. Jayaram, A. Karbasi, V. Mirrokni, D. P. Woodruff, and A. Zandieh, "Hyperattention: Long-context attention in near-linear time," *arXiv preprint arXiv:2310.05869*, 2023.

[81] Z. Li, C. You, S. Bhojanapalli, D. Li, A. S. Rawat, S. Reddi, K. Ye, F. Chern, F. Yu, R. Guo *et al.*, "The lazy neuron phenomenon: On emergence of activation sparsity in transformers," in *Conference on Parsimony and Learning (Recent Spotlight Track)*, 2023.

[82] G. M. Correia, V. Niculae, and A. F. Martins, "Adaptively sparse transformers," *arXiv preprint arXiv:1909.00015*, 2019.

[83] W. Fedus, B. Zoph, and N. Shazeer, "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity," *The Journal of Machine Learning Research*, vol. 23, no. 1, pp. 5232–5270, 2022.

[84] M. V. DeBole, B. Taba, A. Amir, F. Akopyan, A. Andreopoulos, W. P. Risk, J. Kusnitz, C. Ortega Otero, T. K. Nayak, R. Appuswamy, P. J. Carlson, A. S. Cassidy, P. Datta, S. K. Esser, G. J. Garreau,

K. L. Holland, S. Lekuch, M. Mastro, J. McKinstry, C. di Nolfo, B. Paulovicks, J. Sawada, K. Schleupen, B. G. Shaw, J. L. Klamo, M. D. Flickner, J. V. Arthur, and D. S. Modha, "Truenorth: Accelerating from zero to 64 million neurons in 10 years," *Computer*, vol. 52, no. 5, pp. 20–29, 2019.

[85] E. P. Frady, G. Orchard, D. Florey, N. Imam, R. Liu, J. Mishra, J. Tse, A. Wild, F. T. Sommer, and M. Davies, "Neuromorphic nearest neighbor search using intel's pohoiki springs," in *Proceedings of the 2020 Annual Neuro-Inspired Computational Elements Workshop*, 2020, pp. 1–10.

[86] H. Schmidt, J. Montes, A. Grübl, M. Güttler, D. Husmann, J. Ilmberger, J. Kaiser, C. Mauch, E. Müller, L. Sterzenbach, J. Schemmel, and S. Schmitt, "From clean room to machine room: commissioning of the first-generation brainscales wafer-scale neuromorphic system," *Neuromorphic Computing and Engineering*, vol. 3, no. 3, p. 034013, sep 2023. [Online]. Available: https://dx.doi.org/10.1088/2634-4386/acf7e4

[87] S. K. Esser, P. A. Merolla, J. V. Arthur, A. S. Cassidy, R. Appuswamy, A. Andreopoulos, D. J. Berg, J. L. McKinstry, T. Melano, D. R. Barch, C. di Nolfo, P. Datta, A. Amir, B. Taba, M. D. Flickner, and D. S. Modha, "Convolutional networks for fast, energy-efficient neuromorphic computing," *Proceedings of the National Academy of Sciences*, vol. 113, no. 41, pp. 11 441–11 446, 2016. [Online]. Available: https://www.pnas.org/doi/abs/10.1073/pnas.1604850113

[88] J. Lopez-Randulfe, N. Reeb, and A. Knoll, "Conversion of convnets to spiking neural networks with less than one spike per neuron," in *2022 Conference on Cognitive Computational Neuroscience*, 2022, pp. 553–555.

[89] E. Müller, S. Schmitt, C. Mauch, H. Schmidt, J. Montes, J. Ilmberger, J. Klähn, F. Passenberg, C. Koke, M. Kleider, S. Jeltsch, M. Güttler, D. Husmann, S. Billaudelle, P. Müller, A. Grübl, J. Kaiser, J. Weidner, B. Vogginger, J. Partzsch, C. Mayr, and J. Schemmel, *Neurocomputing*, may 2022, in typesetting.

[90] A. G. Rowley, C. Brenninkmeijer, S. Davidson, D. Fellows, A. Gait, D. R. Lester, L. A. Plana, O. Rhodes, A. B. Stokes, and S. B. Furber, "Spinntools: the execution engine for the spinnaker platform," *Frontiers in neuroscience*, vol. 13, p. 231, 2019.

[91] P. Qu, L. Yang, W. Zheng, and Y. Zhang, "A review of basic software for brain-inspired computing," *CCF Transactions on High Performance Computing*, vol. 4, no. 1, pp. 34–42, 2022.

[92] J. K. Eshraghian, M. Ward, E. Neftci, X. Wang, G. Lenz, G. Dwivedi, M. Bennamoun, D. S. Jeong, and W. D. Lu, "Training spiking neural networks using lessons from deep learning," *Proceedings of the IEEE*, vol. 111, no. 9, pp. 1016–1054, 2023.

[93] C. Pehle and J. E. Pedersen, "Norse - A deep learning library for spiking neural networks," Jan. 2021, documentation: https://norse.ai/docs/. [Online]. Available: https://doi.org/10.5281/zenodo.4422025

[94] W. Fang, Y. Chen, J. Ding, Z. Yu, T. Masquelier, D. Chen, L. Huang, H. Zhou, G. Li, and Y. Tian, "Spikingjelly: An open-source machine learning infrastructure platform for spike-based intelligence," *Science Advances*, vol. 9, no. 40, p. eadi1480, 2023.

[95] K. Heckel, "kmheckel/spyx: v0.1.0-beta," Aug. 2023, https://doi.org/10.5281/zenodo.8241588.

[96] J. B. Aimone, W. Severa, and C. M. Vineyard, "Composing neural algorithms with fugu," in *Proceedings of the International Conference on Neuromorphic Systems*, ser. ICONS '19. New York, NY, USA: Association for Computing Machinery, Jul 2019, p. 1–8. [Online]. Available: https://doi.org/10.1145/3354265.3354268

[97] M. G. K. Williams, P. Plank, and S. B. Shrestha, "Lava - a software framework for neuromorphic computing," Oct. 2023. [Online]. Available: https://github.com/lava-nc/lava

[98] J. E. Pedersen, S. Abreu, M. Jobst, G. Lenz, V. Fra, F. C. Bauer, D. R. Muir, P. Zhou, B. Vogginger, K. Heckel, G. Urgese, S. Shankar, T. C. Stewart, J. K. Eshraghian, and S. Sheik, "Neuromorphic intermediate representation: A unified instruction set for interoperable brain-inspired computing," 2023.

[99] "Keras," https://keras.io/.

[100] "Tensorflow," https://www.tensorflow.org/.

[101] "Pytorch," https://pytorch.org/.

[102] "Model zoo performance." [Online]. Available: https://doc.brainchipinc.com/model_zoo_performance.html#akida-2-0-models

[103] L. Eeckhout, "A first-order model to assess computer architecture sustainability," *IEEE Computer Architecture Letters*, vol. 21, no. 2, pp. 137–140, 2022.